

# Unsupervised Doodling and Painting with Improved SPIRAL

John F. J. Mellor   Eunbyung Park   Yaroslav Ganin   Igor Babuschkin   Tejas Kulkarni  
Dan Rosenbaum   Andy Ballard   Theophane Weber   Oriol Vinyals   S. M. Ali Eslami

## Abstract

We investigate using reinforcement learning agents as generative models of images (Ganin et al., 2018). A *generative agent* controls a simulated painting environment, and is trained with rewards provided by a discriminator network simultaneously trained to assess the realism of the agent’s samples, either unconditional or reconstructions. Compared to prior work, we make a number of improvements to the architectures of the agents and discriminators that lead to intriguing and at times surprising results. We find that when sufficiently constrained, generative agents can learn to produce images with a degree of visual abstraction, despite having only ever seen real photographs (no human brush strokes). And given enough time with the painting environment, they can produce images with considerable realism. These results show that, under the right circumstances, some aspects of human drawing can emerge from simulated embodiment, without the need for external supervision, imitation or social cues. Finally, we note the framework’s potential for use in creative applications.

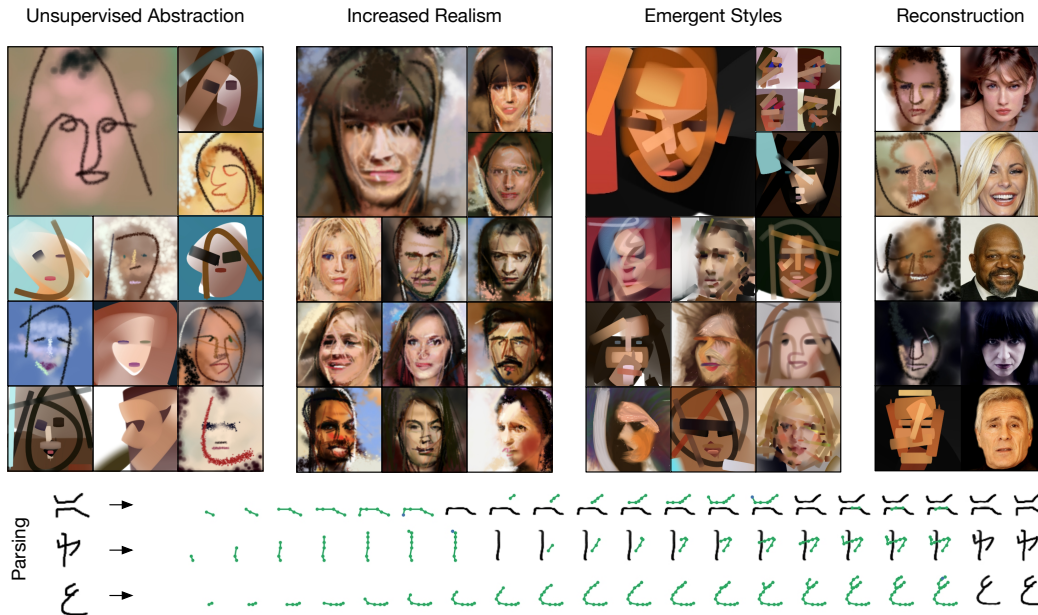


Figure 1: Despite not seeing examples of human drawings, generative agents can produce images with a degree of visual abstraction, in a diverse array of styles, and they can scale to approach realistic results. Agents can also learn to reconstruct, for instance parsing Omniglot characters into strokes.

# 1 Introduction

Through a human’s eyes, the world is much more than just the patterns of photons reflected in our retinas. We have the ability to form high-level interpretations of objects and scenes, and we use these to reason, to build memories, to communicate, to plan and to take actions.

An important open problem in artificial intelligence is how these interpretations are produced in the absence of labelled datasets to train from. The idea of generative modelling has offered a possible solution, whereby high-level representations are obtained by first training models to generate outputs that resemble real images (Kingma and Welling, 2013; Gregor et al., 2015; Eslami et al., 2018). The underlying hypothesis is that if a concise representation is sufficient for the model to reproduce an image, then that representation has, in some form, captured the essence of the contents of that image.

Modern implementations of the generative modeling approach use large, flexible, mostly unstructured neural networks to model the generative process, with state-of-the-art models now achieving photo-realistic results (Karras et al., 2018; Brock et al., 2019) even at high resolutions. Whilst immensely impressive, such results raise two intriguing questions: 1. To what extent is photo-realism necessary for learning of high-level descriptions of images? 2. Is it ever beneficial to incorporate grounded *knowledge* and *structure* in the generative model?

Humans have been representing and reconstructing their visual sensations, using physical tools to draw and sculpt depictions of those sensations, for at least 60,000 years (Hoffmann et al., 2018), well before the development of symbolic writing systems. And it has been hypothesized that abstraction away from raw sensations and use of physical tools are essential components of certain human cognitive abilities (Clark and Chalmers, 1998; Lake et al., 2015, 2017). Within the field of artificial intelligence itself, human figurative drawings have long provided inspiration for the study of meaningful representation of object concepts (Minsky and Papert, 1972).

We would like to be able to create generative models that similarly use physical grounding. In this work, we equip artificial agents with the same tools that we use to create reconstructions of images (namely digital brushes, pens and spray cans). We train these agents with reinforcement learning to interact with digital painting environments (Renold, 2004; Li, 2017). The agents act by placing strokes on a simulated canvas and changing the brush size, pressure and colour as they do so. Building on the work by Ganin et al. (2018), we consider a setting where the agents’ rewards are specified by jointly-trained adversarial discriminator networks. An important aspect of this approach is that the structure is not in the agent itself, but mostly in the environment that it interacts with.

In summary, we make the following contributions: **(a) Quality:** We scale up and tune the work of Ganin et al. (2018). We use several simple but general tricks that improve the performance of the framework and show that these changes positively impact the fidelity and realism of generated images. **(b) Abstraction:** We demonstrate that generative agents can learn to draw recognizable objects with a small number of strokes of the simulated brush and without supervision. These drawings appear to capture the essential elements of objects at a surprisingly high level: for instance generating faces by drawing two eyes, a nose and a mouth each with a single stroke. We also show how the discriminators learn distance metrics that can prioritize semantic likeness over pixel similarity. Our results in this work are primarily qualitative, in part due to the subjective nature of the phenomena being studied, however we provide systematic ablations of the results in the supplementary material.

## 2 SPIRAL++

In this work, we build closely on the SPIRAL architecture (Ganin et al., 2018), see Figure 2. The goal of SPIRAL is to train a policy  $\pi$  that controls a black-box rendering simulator  $\mathcal{R}$  using adversarial learning (Goodfellow et al., 2014). At every time step,  $\pi$  observes the current state of the drawing and produces a rendering command  $a_t$  used by  $\mathcal{R}$  to update the canvas. Each command specifies the appearance of a single brush stroke (*e.g.*, its shape and colour, the amount of pressure applied to the brush and so on – the action space is discussed in greater detail in appendix A). In SPIRAL, the learning signal for the policy comes in the form of the reward computed by a discriminator network  $D$  simultaneously optimized to tell apart examples from some target distribution of images  $p_d$  and final renders produced by the agent (denoted as  $p_g$ ). The outcome of the training procedure is a policy that generates a sequences of commands such that  $p_g$  is close to  $p_d$ . Although the SPIRAL framework was shown to perform well on a number of tasks, in its original form it is limited in terms

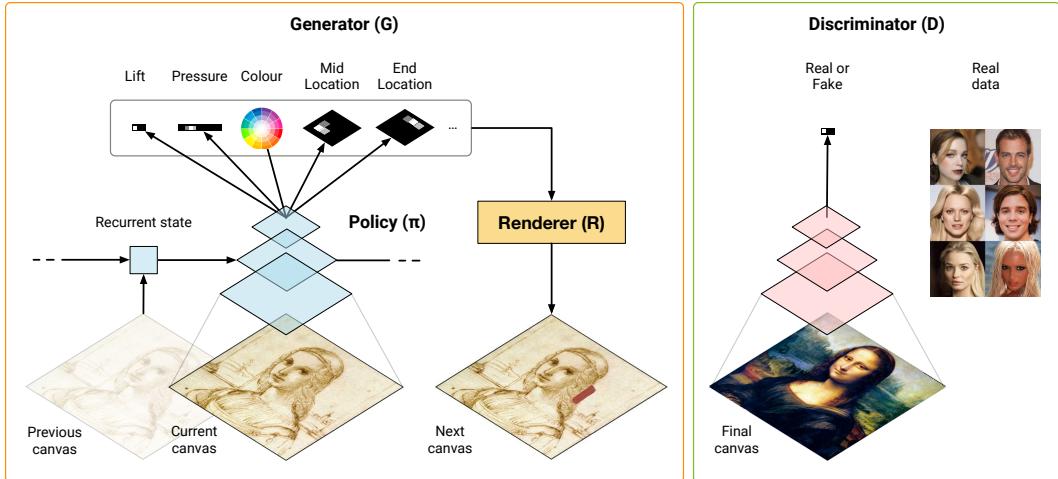


Figure 2: **An architecture for generative agents.** There are two learnable components in the SPIRAL architecture: the policy network or agent  $\pi$  and the discriminator network  $D$ . The policy network takes as input a partially completed canvas and produces the parameters of the agent’s action. This action is sent to the renderer  $\mathcal{R}$  as a command to update the canvas. The discriminator network takes as input the final canvas and attempts to classify it as real or fake. This network is trained using both images from the real dataset and rendered samples from the agent.

of image fidelity and scalability. We improve these aspects of the framework by introducing the simple but effective modifications detailed in the following subsections. We refer to the resulting model as SPIRAL++ in order to distinguish it from prior work.

**2.1 Spectral Normalization.** SPIRAL (Ganin et al., 2018) followed the setup of Gulrajani et al. (2017), namely that of a WGAN-GP discriminator architecture. In SPIRAL++, we instead revert to the standard GAN loss (Goodfellow et al., 2014), but with spectral normalization as described in Miyato et al. (2018). Spectral normalization controls the Lipschitz constant of the discriminator by literally setting the spectral norm of each of its layers to equal 1. The motivation for this is to stabilize the training of the discriminator. In practice we find that spectral normalization significantly improves the fidelity of the generated images, as ablated in subsection F.1 in the appendix.

**2.2 Temporal Credit Assignment.** SPIRAL’s reward structure (namely that of rewarding the agent only on the last step of each episode) provides agents with the flexibility to learn non-greedy image generation policies, however in practice, the sparsity of the learning signal limits the length of the episodes for which reinforcement learning techniques yield positive results. Ganin et al. (2018) only report results on episodes of length 20. It would be desirable to maintain this flexibility to achieve non-greedy policies whilst relaxing the difficulty of the learning problem. In SPIRAL++, instead of only providing the agent a reward at the end of the episode based on the final discriminator loss, we calculate the discriminator loss at every step based on the partially-drawn canvas, and at every timestep the agent gets the 1-step improvement in loss as its reward:  $r_t = D(\mathcal{R}(a_{1:t})) - D(\mathcal{R}(a_{1:t-1}))$ <sup>1</sup>. Similarly to Ng et al. (1999), this reward redistribution leads to the same optimal policies if the reinforcement learning discount factor  $\gamma$  is set to equal 1. Intuitively, when  $\gamma = 1$ , all terms cancel out in the calculation of returns apart from the final step’s reward. However we obtained best results with  $\gamma$  between 0 and 0.99, which in this reward redistribution scheme causes the agent to be more greedy by introducing some bias (Guez et al., 2018). Huang et al. (2019) similarly use this technique, concurrent with this work. We provide an ablation of the effectiveness of temporal credit assignment in subsection F.2 in the appendix.

**2.3 Complement Discriminator.** When in conditional training mode, the discriminator is used to train agents such that when given a target image  $x_{\text{target}}$ , they produce reconstructions  $x$  that are as similar as possible to the target. Note that conditional agents are trained differently than their unconditional counterparts. Ganin et al. (2018) achieve this by conditioning the discriminator on the target. Specifically, in each forward pass, the input to the discriminator is either: 1. a *fake* image pair

<sup>1</sup>Note lack of  $\gamma$  before the first term, unlike in Ng et al. (1999).

( $\mathbf{x}_{\text{target}}, \mathbf{x}$ ), or 2. a *real* image pair ( $\mathbf{x}_{\text{target}}, \mathbf{x}_{\text{target}}$ ). Agents are encouraged to move in a direction that makes renderings  $\mathbf{x}$  more similar to the targets  $\mathbf{x}_{\text{target}}$  in order to fool the discriminator. A potential downfall of this method is that the discriminator can in theory find a simple shortcut to achieve perfect performance: it can compare the two images in each input pair, and if they are not identical, pixel to pixel, then it will know that it is a fake image pair. Whilst this strategy allows the discriminator to achieve its training objective, it can lead to a sub-optimal loss surface for the agent.

To overcome this we use what we call a *complement discriminator*. In each forward pass, a binary image mask  $\mathbf{m}$  is sampled. We mask the rendered image by  $\mathbf{m}$  to obtain  $\mathbf{x}^m = \mathbf{m} \cdot \mathbf{x}$ , and mask the target both by  $\mathbf{m}$  to obtain  $\mathbf{x}_{\text{target}}^m = \mathbf{m} \cdot \mathbf{x}_{\text{target}}$  and by the complement of  $\mathbf{m}$  to obtain  $\mathbf{x}_{\text{target}}^{1-m} = (1 - \mathbf{m}) \cdot \mathbf{x}_{\text{target}}$ . We then define the fake input to the discriminator to be ( $\mathbf{x}_{\text{target}}^{1-m}, \mathbf{x}^m$ ) and the real input to the discriminator to be ( $\mathbf{x}_{\text{target}}^{1-m}, \mathbf{x}_{\text{target}}^m$ ). In our experiments we mask the left/right half (see *e.g.* Figure 3) or top/bottom half of the image, however random masks are likely to work just as well (Pathak et al., 2016). This trick makes it impossible for the discriminator to detect the real pair by simply comparing pixels, therefore creating a more suitable loss surface for the agent. We ablate the effect of the complement discriminator in subsection F.3 in the appendix.



Figure 3: Fake input ( $\mathbf{x}_{\text{target}}^{1-m}, \mathbf{x}^m$ )

**2.4 Population discriminator.** For best results we dynamically tune learning rate and entropy cost hyper-parameters using population based training (PBT, Jaderberg et al., 2017). However, unlike Ganin et al. (2018) which also employs this technique, we typically train with a periodically-evolving population of 10 generative agents sharing a *single* discriminator. The effect of this modification is three-fold. Since the discriminator receives fake samples from all the policy learners, this means that it gets updated more frequently than individual generators allowing us to remove the replay buffer from the original SPIRAL distributed setup. PBT also automatically resurrects individual generators should they collapse during training. Additionally each generator can now specialize in a subset of the modes of the full distribution, such that the population collectively covers the full distribution. Put another way, this technique allows us to increase the relative representational capacity of the generators as compared to the discriminator. We show in the experiments that this setup sometimes leads to the different agents in a single population learning different painting styles. Note, however, that we do not employ any explicit techniques to encourage or guarantee population diversity.

### 3 Experimental Results

Owing to the diverse range of settings under which the generative agents framework can be examined (with different action sets, rendering environments, brush types, episode lengths, agent and discriminator hyper-parameters, and so on), we first provide highlights of our main findings in section 3. In figures, we show selected agents to illustrate phenomena that are *possible* to observe with the framework. This is not to imply that the settings those agents were trained with are necessary or sufficient for the observed results. In the appendix we provide a systematic analysis of the major components of SPIRAL++ via controlled ablations.

We show highlights of our experiments on Celeba-HQ (Karras et al., 2017) in Figure 4. Amongst all the framework’s settings, the one we observed to have the most profound impact on the agents’ behaviour was the number of brush strokes (steps) they were allowed in each episode to generate an image. Agents that were constrained with short episodes learned qualitatively different policies than those that could afford numerous interactions, producing images with a degree of visual abstraction, not unlike how humans do when similarly constrained (Selim, 2018; Fan et al., 2019). For this reason we structure the results into two sections: short episodes in subsection 3.1 and long episodes in subsection 3.2.

**3.1 Short episodes.** We encourage the reader to take a few moments to inspect the samples in Figure 4(a-c). Each row was generated by a different agent, differing in the settings of their environments (*e.g.* brush type, action space, episode length).

It was surprising for us to see the aesthetically pleasing way in which the agent draws faces: *e.g.* using a large circle to delineate the outline of the face, dots for each of the eyes, a line for the nose and a line for the mouth. Note that the agent has never been exposed to human *drawings* of human faces, but only to realistic *photographs* of human faces. Also note that the agents choose to use bright



Figure 4: **Artificial portraiture with generative agents.** Random unconditional samples generated by agents trained on CelebA-HQ. Separate agents with varying hyper-parameters were trained to generate samples (a, b, c) in 17 steps with various brushes and action spaces, (d) in 19 steps using an oil paint simulator (Li, 2017), (e) in 31 steps, (f) in 400 steps, (g) in 1000 steps. As a baseline, (h) shows unconditional samples generated in 19 steps by our best reimplement of SPIRAL (Ganin et al., 2018, WGAN-GP + single discriminator per population + hyperparameter tweaks). Our improved method scales with episode length from relatively abstract to approach realistic-looking results. See <https://learning-to-paint.github.io> for videos.

colours and thin strokes to depict salient elements of faces despite no element of the framework encouraging such behaviour. In all cases, the architecture of the agent is constant, and it is the variation in the characteristics of the agent’s environment that creates the diversity of styles. These results serve as an existence proof for the conjecture that some aspects of human drawings of faces can emerge automatically from a learning framework as simple as that of generative agents (namely agents equipped with brushes working against a discriminator), without the need for supervision, imitation or social cues.

Unlike in most modern GANs which directly output pixels, in this setting there is a large discrepancy between what the generator *can* produce, and what it *should* produce. The brushes are too constrained and there is simply not enough time in the episodes for the generator to be able to produce a completely photo-realistic image. And in our experiments the discriminators can always distinguish between generated and real images with high confidence. Nevertheless we observe that when sufficiently regularised, discriminators can provide sufficient signal for meaningful learning to take



(a) gray with gray (diverse)      (b) color with gray (mode-collapse)      (c) gray with black (mode-collapse)

Figure 5: **Tasked with the impossible, these agents make do.** Samples are selected from three 20 step models trained on CelebA-HQ with modified environments. As a baseline, (a) was tasked with generating grayscale photos using a black brush with variable opacity. (b) was tasked with generating color photos using the same variable opacity black brush. (c) was tasked with generating grayscale photos using an opaque black pen. Models (b) and (c) often manage to draw recognizable faces despite the huge gap between what they *can* and *should* produce, however both experience severe mode collapse: each agent in these populations generates minor variations on a single image rather than a full distribution of images.

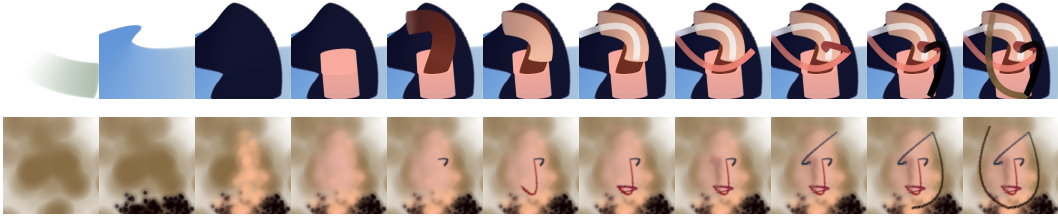


Figure 6: **10 step canvas sequences show precise control.** Agents interact with the canvas to produce the final image, manipulating the location, colour and thickness of the brush with precision. See <https://learning-to-paint.github.io> for videos.

place, suggesting that they are still capable of ranking generated images in a useful way. We explore this further in Figure 5, by training agents on color photos but only with various grayscale brushes.

It is informative to examine how agents interact with the simulated canvas to produce these images. In Figure 6 we show a number of episodes as sequences. We see that agents can learn to manipulate the location, colour and thickness of the brush with precision, constructing the final images stroke by stroke. It is worth noting that due to RL’s objective of maximizing potentially delayed rewards, agent policies often deviate from being purely greedy, and they often take actions that appear to reduce the quality of the image, especially early in the episode, only for it to be revealed to have been important for the final drawing. We revisit this point in subsection 3.2 and Figure 17 in the appendix.

As described in subsection 2.4, we use population based training (PBT, Jaderberg et al., 2017). In Figure 7 (and Figure 11 in the appendix) we show samples from three different agents that were trained as part of the same population. We see that the different generators each specialize in a subset of the modes of the full distribution, each producing images with a perceptibly different style. Note that in all three cases the architecture of the agents, their action spaces, and the settings of their environments were identical, and they differ only in their weights and evolved learning rate and entropy cost. Finally, we show conditional samples in Figure 8. The agent is able to match the higher level statistics of the target images, and even appears to capture the faces’ smiles. In Figure 12 in the appendix we visualise the agent’s stochasticity by producing multiple samples for the same target image.

**3.2 Long episodes.** We found temporal credit assignment (TCA, subsection 2.2) to be crucial for training agents on long episodes. Without TCA, agents only learn to perform meaningful actions in the last 15 to 60 steps or so of the episode, regardless of the episode length. However, with TCA, agents make full use of episodes of up to 1000 steps, leading in turn to qualitatively different generation policies. We train with discounts of 0.9 or 0.99 and unroll lengths of 20-50 in the experiments that follow. See Figure 4(f,g) for samples (as well as Figure 11 and Figure 18 in the appendix).

In Figure 9 we show a visualisation of how images are constructed by TCA agents. It can be seen that agents are successful in making good use of hundreds of timesteps, controlling the brush with precision to form comparatively realistic images. Due to the smaller-than-1 discount factor, agents are incentivised to have somewhat believable intermediate canvases, leading to the formation of the most critical elements of faces (*e.g.* eyes, nose and mouth) even in the earliest steps of the episode.



Figure 7: **Within-population and within-agent diversity (short episodes).** We show three sets of samples for three different agents in the same population, showing both the diversity of samples for each agent, and the diversity of samples across agents in the same population. The first agent attempts more figurative line drawings, whilst the agent on the right uses more realistic shading. We do occasionally observe ‘mode collapse’ where certain samples repeat themselves, for instance in the middle agent though there is still slight variation in the way each image is drawn.



Figure 8: **Conditional image generation.** The agent is generally able to match the higher level statistics of the targets (top), and even appears to capture the faces’ smiles.

In unconditional generation, much of the form of the final is determined by the randomness of the first few strokes on the canvas, with the agent doing its best towards the end of the episode to complete the result. In this light, the agent can be seen to be similar in spirit to auto-regressive models such as Pixel-RNNs and Pixel-CNNs (van den Oord et al., 2016), but in stroke-space as opposed to pixel-space.

**3.3 Other datasets and environments.** We conduct a set of additional experiments to assess the generality of the proposed approach. appendix C presents results for several ImageNet (Russakovsky et al., 2015) classes. In appendix D, we discuss how a slight modification to the painting environment can lead to meaningful parses of the Omniglot (Lake et al., 2015) characters. Finally, appendix G introduces a new oil paint simulator and shows that our agent is able to learn to control it despite the added complexity of the setting.

## 4 Related Work

Non-photorealistic rendering emerged as an area of digital art that focuses on enabling a wide variety of artistic styles (in contrast to standard computer graphics which has a focus on photorealism). Stroke based rendering is a particular approach in non-photorealistic rendering that creates images by placing discrete elements such as paint strokes or stipples (see Hertzmann, 2003 for a comprehensive review). A common goal in stroke based rendering is to make the painting look like some other image, whilst limiting the total number of strokes (Hertzmann, 2003). Techniques from computer vision have in some cases been used to control the positions of the strokes (Zeng et al., 2009). One attractive property of stroke based rendering is that its interpretations of images can, in theory, be executed in the physical world using real robots (Kotani and Tellex, 2019; El Helou et al., 2019). Closely connected are notions of ‘low-complexity’ art (Schmidhuber, 1997), which describe an algorithmic theory of beauty and aesthetics based on the principles of information theory. It postulates that amongst several comparable images, the most pleasing one is the one with the shortest description. In almost all cases, the aforementioned techniques operate by optimising, either explicitly or implicitly, a fixed objective function without any learnable parameters. New large-scale sources of data have made it possible to model human behaviour more directly in simple drawing domains (Ha and Eck,



Figure 9: **Longer sequences demonstrate higher levels of realism.** SPIRAL++ agents make good use of hundreds of timesteps, controlling the brush with precision. Here we show logarithmically spaced frames from a 1000 and 500 step agent. See <https://learning-to-paint.github.io> for videos.

2017; Zhou et al., 2018; Li et al., 2019), potentially making it possible to capture our internal object function in a data-driven manner.

Superficially, SPIRAL (Ganin et al., 2018) and SPIRAL++ can be seen as a technique for stroke based, non-photorealistic rendering. Similarly to some modern stroke based rendering techniques, the positions of strokes are determined by a learned system (namely, the agent). Unlike traditional methods however, the objective function that determines the goodness of the output image is learned unsupervised via the adversarial objective. This adversarial objective allows us to train without access to ground-truth strokes, in contrast to *e.g.* Ha and Eck (2017); Zhou et al. (2018); Li et al. (2019),



enabling the method to be applied to domains where labels are prohibitively expensive to collect, and for it to discover surprising and unexpected styles.

There are a number of works that use constructions similar to SPIRAL to tune the parameters of non-differentiable simulators (Louppe et al., 2019; Ruiz et al., 2018; Kar et al., 2019). Frans and Cheng (2018); Nakano (2019); Zheng et al. (2019); Huang et al. (2019) achieve remarkable learning speed by relying on back-propagation through a learned model of the renderer, however they rely on being able to train an accurate model of the environment in the first instance. We further elaborate on this difference in appendix G. On the applications side, El-Nouby et al. (2018); Lázaro-Gredilla et al. (2019); Agrawal et al. (2018) build systems that generate outcomes iteratively, conditioned on ongoing linguistic input or feedback. Ellis et al. (2018); Liu et al. (2019) use inference techniques to convert simple scenes into programs expressing abstract regularities. Technical advances such as these (for instance the use of neural simulators and model-based RL) are largely orthogonal to what we introduce in this paper, and are likely to improve the performance of SPIRAL++ further when used in combination. Nonetheless we note that to the best of our knowledge, very few models have the capability to paint images without conditioning on a target image, and even fewer do so at a similar level of abstraction, and lead to as many surprising, emergent styles as those obtained with SPIRAL++.

## 5 Discussion and Conclusions

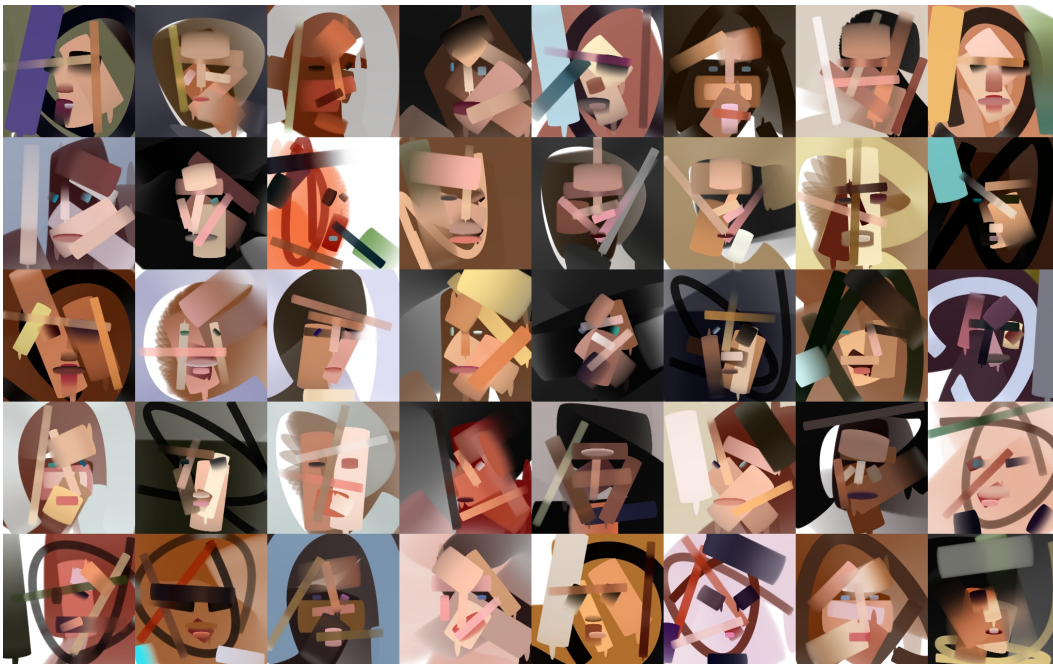


Figure 10: Selected samples from a single population of agents. The agents find a great variety of ways of working within the constraints of the brush they were given to evoke faces in just 32 steps.

We showed that when sufficiently constrained, generative agents can learn to produce images with a degree of visual abstraction, despite having never been shown human drawings. To the best of our knowledge, this is the first time that an unsupervised artificial system has discovered visual abstractions of this kind, abstractions that resemble those made by children and novice illustrators. We stress that our aim has been to show that certain results are at least possible to achieve with the generative agents framework, not that the specifics of our setting are necessarily the best way of achieving those results. It is important to acknowledge that our results are mostly qualitative, in part due to the subjective nature of the phenomena being studied (*e.g.* improved realism and abstraction as compared to SPIRAL). Nevertheless an important line of work will be to quantify these effects via controlled human studies, or even by exposing the generated samples to the markets (Cohn, 2018)!

The abstraction results arise due to the constraints imposed on the generator by the renderer and the episode’s finite horizon. The generator cannot manipulate pixels directly in the way that generative models typically do; instead it has to obey the physics of the rendering simulation. In other words, the generator is *embodied*, albeit in a simulation. The discriminator must also be well regularized for this result. It is worth noting that the framework does utilize structure to obtain its results, but this structure is built into the environment, not the agent itself. This stands in contrast with efforts that achieve abstraction by structuring the model or agent, and provides a potentially useful data point for recent debates around the utility of structure (see *e.g.* Sutton, 2019; Welling, 2019).

We also showed that such generative agents can produce outputs with considerable realism and complexity when given enough time with the environment. This is a considerable step up from the original SPIRAL, providing encouraging evidence that the approach can be scaled to other sequence generation tasks where expert sequences are difficult to obtain for imitation, for instance in music generation, program synthesis, chemical synthesis and protein folding.

Finally, we note the framework’s potential for creative and artistic applications. Machine learning techniques such as GANs (Goodfellow et al., 2014) and neural style transfer (Gatys et al., 2016) are increasingly being used for creative investigation in the artistic community. Generative agents have an uncommon ability to generate novel aesthetic styles with no human input except for a choice of brush (see *e.g.* Figure 10). The extent to which any of these systems can be considered *creative* in and of themselves is open for discussion (Hertzmann, 2018), however we look forward to seeing whether the proposed framework can provide new points for consideration in that debate.

**Acknowledgements** We would like to thank Aishwarya Agrawal, Nicolas Heess, Simon Kohl, Adam Kosiorek and David Ha for insightful discussions and comments on this manuscript.

## References

- Agrawal, A., M. Malinowski, F. Hill, S. M. A. Eslami, O. Vinyals, and T. Kulkarni  
2018. Generating diverse programs with instruction conditioned reinforced adversarial learning. [arXiv:1812.00898](https://arxiv.org/abs/1812.00898) [cs.LG].
- Brock, A., J. Donahue, and K. Simonyan  
2019. Large scale GAN training for high fidelity natural image synthesis. In *International Conference on Learning Representations*.
- Clark, A. and D. Chalmers  
1998. The extended mind. *analysis*, 58(1):7–19.
- Cohn, G.  
2018. AI art at christie’s sells for \$432,500. <https://www.nytimes.com/2018/10/25/arts/design/ai-art-sold-christies.html>.
- El Helou, M., S. Mandt, A. Krause, and P. Beardsley  
2019. Mobile robotic painting of texture. *2019 International Conference on Robotics and Automation (ICRA)*, Pp. 640–647.
- El-Nouby, A., S. Sharma, H. Schulz, D. Hjelm, L. El Asri, S. E. Kahou, Y. Bengio, and G. W. Taylor  
2018. Tell, draw, and repeat: Generating and modifying images based on continual linguistic instruction. [arXiv:1811.09845](https://arxiv.org/abs/1811.09845) [cs.CV].
- Ellis, K., D. Ritchie, A. Solar-Lezama, and J. Tenenbaum  
2018. Learning to infer graphics programs from hand-drawn images. In *Advances in Neural Information Processing Systems*, Pp. 6059–6068.
- Eslami, S. M. A., D. J. Rezende, F. Besse, F. Viola, A. S. Morcos, M. Garnelo, A. Ruderman, A. A. Rusu, I. Danihelka, K. Gregor, and others  
2018. Neural scene representation and rendering. *Science*, 360(6394):1204–1210.
- Fan, J. E., R. D. Hawkins, M. Wu, and N. D. Goodman  
2019. Pragmatic inference and visual abstraction enable contextual flexibility during visual communication. *Computational Brain & Behavior*.

- Fernando, R. and others  
2004. *GPU gems: programming techniques, tips, and tricks for real-time graphics*, volume 590. Addison-Wesley Reading.
- Frans, K. and C.-Y. Cheng  
2018. Unsupervised image to sequence translation with canvas-drawer networks. [arXiv:1809.08340](https://arxiv.org/abs/1809.08340) [cs.CV].
- Ganin, Y., T. Kulkarni, I. Babuschkin, S. M. A. Eslami, and O. Vinyals  
2018. Synthesizing programs for images using reinforced adversarial learning. In *Proceedings of the 35th International Conference on Machine Learning*.
- Gatys, L., A. Ecker, and M. Bethge  
2016. A neural algorithm of artistic style. *Journal of Vision*, 16(12):326.
- Goodfellow, I., J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio  
2014. Generative adversarial nets. In *Advances in neural information processing systems*, Pp. 2672–2680.
- Gregor, K., I. Danihelka, A. Graves, D. J. Rezende, and D. Wierstra  
2015. Draw: A recurrent neural network for image generation. In *ICML*.
- Guez, A., T. Weber, I. Antonoglou, K. Simonyan, O. Vinyals, D. Wierstra, R. Munos, and D. Silver  
2018. Learning to search with MCTSnets. In *Proceedings of the International Conference on Machine Learning (ICML)*.
- Gulrajani, I., F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville  
2017. Improved training of Wasserstein GANs. In *Advances in Neural Information Processing Systems 30*.
- Ha, D. and D. Eck  
2017. A neural representation of sketch drawings. [arXiv:1704.03477](https://arxiv.org/abs/1704.03477) [cs.NE].
- Hertzmann, A.  
2003. A survey of stroke-based rendering. *IEEE Computer Graphics and Applications*, 23:70–81.
- Hertzmann, A.  
2018. Can computers create art? *Arts*, 7(2):18.
- Hoffmann, D. L., C. Standish, M. García-Diez, P. B. Pettitt, J. Milton, J. Zilhão, J. J. Alcolea-González, P. Cantalejo-Duarte, H. Collado, R. De Balbín, and others  
2018. U-Th dating of carbonate crusts reveals Neandertal origin of Iberian cave art. *Science*, 359(6378):912–915.
- Huang, Z., W. Heng, and S. Zhou  
2019. Learning to paint with model-based deep reinforcement learning. [arXiv:1903.04411](https://arxiv.org/abs/1903.04411) [cs.CV].
- Jaderberg, M., V. Dalibard, S. Osindero, W. M. Czarnecki, J. Donahue, A. Razavi, O. Vinyals, T. Green, I. Dunning, K. Simonyan, and others  
2017. Population based training of neural networks. [arXiv:1711.09846](https://arxiv.org/abs/1711.09846) [cs.LG].
- Kar, A., A. Prakash, M.-Y. Liu, E. Cameracci, J. Yuan, M. Rusiniak, D. Acuna, A. Torralba, and S. Fidler  
2019. Meta-sim: Learning to generate synthetic datasets. [arXiv:1904.11621](https://arxiv.org/abs/1904.11621) [cs.CV].
- Karras, T., T. Aila, S. Laine, and J. Lehtinen  
2017. Progressive growing of gans for improved quality, stability, and variation. [arXiv:1710.10196](https://arxiv.org/abs/1710.10196) [cs.NE].
- Karras, T., S. Laine, and T. Aila  
2018. A style-based generator architecture for generative adversarial networks. [arXiv:1812.04948](https://arxiv.org/abs/1812.04948) [cs.NE].

- Kingma, D. P. and J. Ba  
2014. Adam: A method for stochastic optimization. [arXiv:1412.6980](https://arxiv.org/abs/1412.6980) [cs.LG].
- Kingma, D. P. and M. Welling  
2013. Auto-encoding variational bayes. [arXiv:1312.6114](https://arxiv.org/abs/1312.6114) [stat.ML].
- Kotani, A. and S. Tellex  
2019. Teaching Robots to Draw. In *IEEE International Conference on Robotics and Automation (ICRA)*.
- Lake, B. M., R. Salakhutdinov, and J. B. Tenenbaum  
2015. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338.
- Lake, B. M., T. D. Ullman, J. B. Tenenbaum, and S. J. Gershman  
2017. Building machines that learn and think like people. *Behavioral and Brain Sciences*, 40.
- Lázaro-Gredilla, M., D. Lin, J. S. Guntupalli, and D. George  
2019. Beyond imitation: Zero-shot task transfer on robots by learning concepts as cognitive programs. *Science Robotics*, 4(26):eaav3150.
- Li, D.  
2017. Fluid paint. <http://david.li/paint/>.
- Li, M., Z. Lin, R. Mech, E. Yumer, and D. Ramanan  
2019. Photo-sketching: Inferring contour drawings from images. In *WACV 2019: IEEE Winter Conf. on Applications of Computer Vision*.
- Liu, Y., Z. D. Wu, D. A. Ritchie, W. T. Freeman, J. B. Tenenbaum, and J. Wu  
2019. Learning to describe scenes with programs. In *ICLR*.
- Louppe, G., J. Hermans, and K. Cranmer  
2019. Adversarial variational optimization of non-differentiable simulators. In *The 22nd International Conference on Artificial Intelligence and Statistics*, Pp. 1438–1447.
- Minsky, M. and S. A. Papert  
1972. Artificial intelligence progress report.
- Miyato, T., T. Kataoka, M. Koyama, and Y. Yoshida  
2018. Spectral normalization for generative adversarial networks. In *International Conference on Learning Representations*.
- Nakano, R.  
2019. Neural painters: A learned differentiable constraint for generating brushstroke paintings. [arXiv:1904.08410](https://arxiv.org/abs/1904.08410) [cs.CV].
- Ng, A. Y., D. Harada, and S. J. Russell  
1999. Policy invariance under reward transformations: Theory and application to reward shaping. In *Proceedings of the Sixteenth International Conference on Machine Learning (ICML 1999)*.
- Pathak, D., P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros  
2016. Context encoders: Feature learning by inpainting. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Renold, M.  
2004. Mypaint. <http://mypaint.org/>.
- Ruiz, N., S. Schuler, and M. Chandraker  
2018. Learning to simulate. [arXiv:1810.02513](https://arxiv.org/abs/1810.02513) [cs.LG].
- Russakovsky, O., J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei  
2015. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252.

- Schmidhuber, J.  
1997. Low-complexity art. *Leonardo*, 30(2):97–103.
- Selim, M.  
2018. Spiderman 10 min 1 min 10 sec speed challenge. [https://www.youtube.com/watch?v=x9wn633v1\\_c](https://www.youtube.com/watch?v=x9wn633v1_c).
- Sutton, R.  
2019. The bitter lesson. <http://www.incompleteideas.net/IncIdeas/BitterLesson.html>.
- van den Oord, A., N. Kalchbrenner, and K. Kavukcuoglu  
2016. Pixel recurrent neural networks. In *International Conference on Machine Learning*, Pp. 1747–1756.
- Welling, M.  
2019. Do we still need models or just more data and compute? <https://staff.fnwi.uva.nl/m.welling/wp-content/uploads/Model-versus-Data-AI-1.pdf>.
- Zeng, K., M. Zhao, C. Xiong, and S.-C. Zhu  
2009. From image parsing to painterly rendering. *ACM Trans. Graph.*, 29:2:1–2:11.
- Zheng, N., Y. Jiang, and D. jiang Huang  
2019. Strokenet: A neural painting environment. In *ICLR*.
- Zhou, T., C. Fang, Z. Wang, J. Yang, B. Kim, Z. Chen, J. Brandt, and D. Terzopoulos  
2018. Learning to sketch with deep Q networks and demonstrated strokes. [arXiv:1810.05977](https://arxiv.org/abs/1810.05977) [cs.CV].

## Appendix A Compound action space

In Ganin et al. (2018), a new stroke is deposited onto the canvas at every step of the environment. Consequently, each stroke has a fixed quadratic Bézier form with start coordinates, mid-point coordinates and end coordinates (in addition to thickness and colour). The agent produces new values for the mid-point and end coordinates at every step and selects a thickness and colour for the stroke. The previous step’s end coordinate is used for the current step’s start coordinate. We refer to this interface as the *simple* interface.

In this work, we also investigate a *compound* interface, where each stroke is built up over the course of several environment steps. In each step, the agent produces the coordinates of a new control point in the current stroke. In order to terminate a stroke, the agent can execute a discrete action, at which point the stroke is deposited onto the canvas as a cubic spline going through the control points. The thickness and colour of the stroke are determined at the beginning of a sequence of strokes. The compound interface allows the agents to produce smoother, more precise curves, which can sometimes lead to more aesthetically pleasing results. The compound interface is also instrumental in parsing Omniglot characters, as we show in appendix D.

## Appendix B Sample diversity

Figure 7 and Figure 11 show diversity of unconditional samples. Figure 12 shows that there is also considerable diversity in reconstructions of the same target.



Figure 11: **Within-population and within-agent diversity (long episodes)**. We show three sets of samples for three different agents in the same population, showing both the diversity of samples for each agent, and the diversity of samples across agents in the same population. Mode collapse occurs less frequently in long episodes.

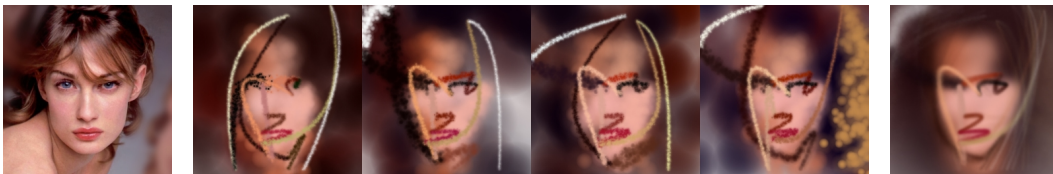


Figure 12: **Conditional generation stochasticity**. We visualise an agent’s stochasticity by producing multiple samples for the same target image. Agent stochasticity is due to each action being sampled auto-regressively. Left: Target image. Middle: Multiple samples from the same agent. Right: Mean image produced by averaging 100 samples.

## Appendix C Evaluation on ImageNet

We also try the generative agents framework on the ImageNet dataset (Russakovsky et al., 2015) to see if it generalizes to more varied images that have not been aligned/cropped. However to keep the task simpler we train on individual classes from ImageNet, so only about 1000 images are used for training each agent, again downscaled to  $64 \times 64$ .

We sometimes instead train using an infinite dataset of samples for single ImageNet classes generated by a pre-trained BigGAN (Brock et al., 2019) with a truncation threshold of 0.4. This is similar to training on ImageNet directly, but due to the truncation the samples generated are less diverse: objects tend to be centered, face on, with normal appearances and less cluttered backgrounds.

In both cases generative agents are able to reproduce the rough colors, and sometimes as in Figure 13 can extract high-level structure (catamarans have masts, daisies have radial lines, and so on), but results are less reliable than for faces.



Figure 13: **Cherry-picked samples of reconstruction with complement discriminator, trained for 20 steps on single classes of ImageNet.** Left: Trained on ImageNet catamaran class. Middle: Trained on BigGAN samples of daisies. Right: Trained on BigGAN samples of phones.

## Appendix D Evaluation on Omniglot

We also experiment with applying the generative agents framework to the Omniglot dataset (Lake et al., 2015). In this setting, we train conditional generative agents to reconstruct Omniglot characters, thereby learning a mapping from bitmap representations into strokes. Note that we do not use the human stroke data for training. We primarily use agents with compound action spaces (appendix A) to allow them to express more complex, intricately curved strokes when parsing. We additionally introduce a number of hyperparameters to encourage the agent to learn more natural behaviours. In particular, the environment rewards the agent with a small penalty for each time a new stroke is created, as well as a penalty on the total length of each stroke. Without these penalties, agents achieve almost perfect reconstructions but with un-natural movements. We set the values of these parameters via trial and error and visual inspection. A sample of these results can be seen in Figure 14 and Figure 15.

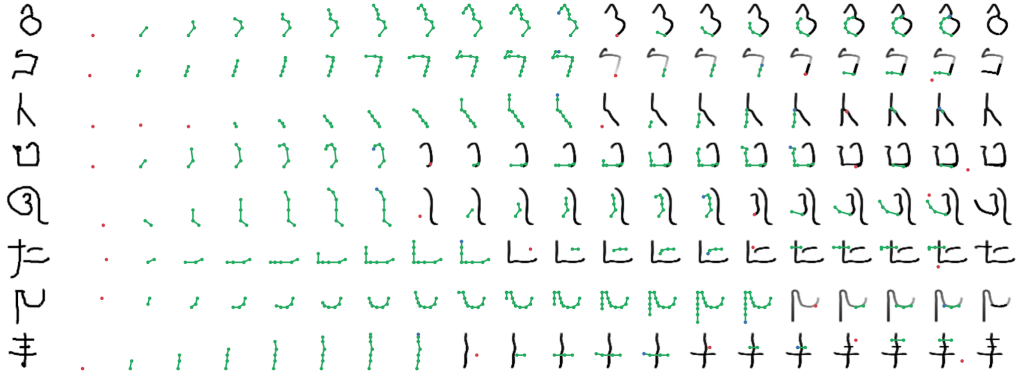


Figure 14: **Omniglot parsing.** Left: Target images. Right: Sequence of actions taken by a single agent to recreate the target images. Red dots signify a request from the agent to commit the current compound stroke to the canvas and to start a new stroke. Green dots indicate locations of commands that build up a compound stroke. The agent is mostly successful in parsing characters into plausible strokes. At times, it fails to complete the character in the allocated episode length.

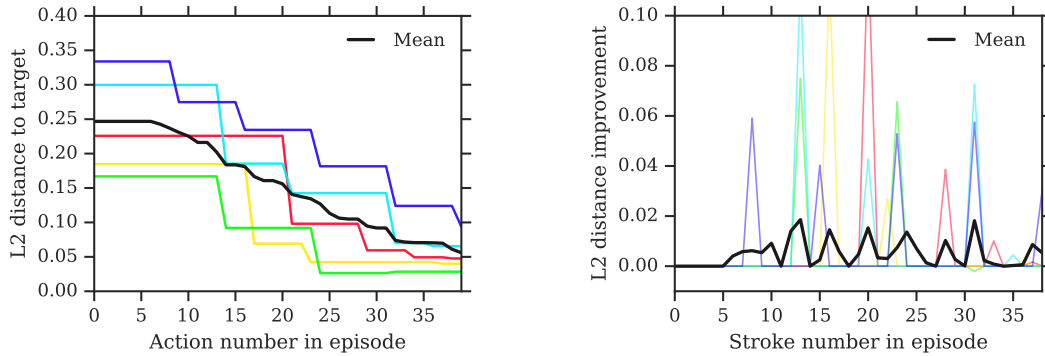


Figure 15: **Within episode 'improvement' of image quality (Omniglot).** Left: We plot the L2 distance of the canvas to its target image for 5 different episodes of a 40 step agent (each curve a different episode). Right: The amount of improvement in L2 distance in each step.



## Appendix E Training curves

Figure 16, Figure 17 show training curves for representative agents on the CelebA-HQ dataset.

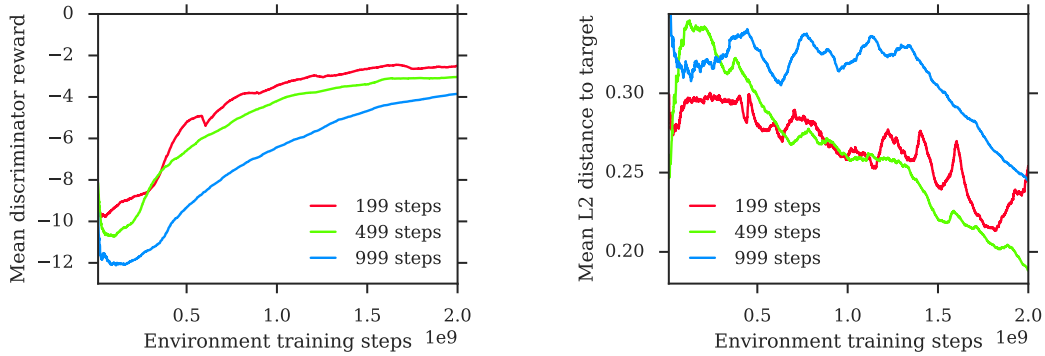


Figure 16: **Conditional training curves with long episodes.** Left: Mean discriminator reward achieved by the generators on the final step of the episode as training progresses. Note that the absolute value of the curves cannot be compared across different agents, due to each competing with their own discriminators. Right: Mean L2 distance of the final generated image and the target image. Note that the agents are trained to maximise the discriminator reward, but we observe some correlation between that objective and minimisation of L2 error.

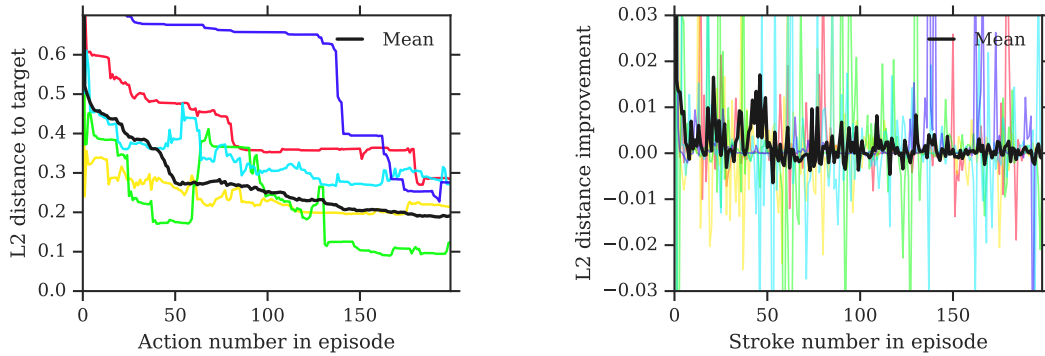


Figure 17: **Within episode 'improvement' of image quality (CelebA).** Left: We plot the L2 distance of the canvas to its target image for 5 different episodes of a conditional 200 step agent trained to reconstruct specific target images with discount factor 0.99 (each curve a different episode). Right: The amount of improvement in L2 distance in each step. We see that the agent often takes actions that increase the distance between the canvas and the target.

## Appendix F Ablations

We provide ablations to show the impact of our design choices. Figure 18 shows the combined improvement from all our modifications detailed under section 2. The following sections separately ablate important modifications.

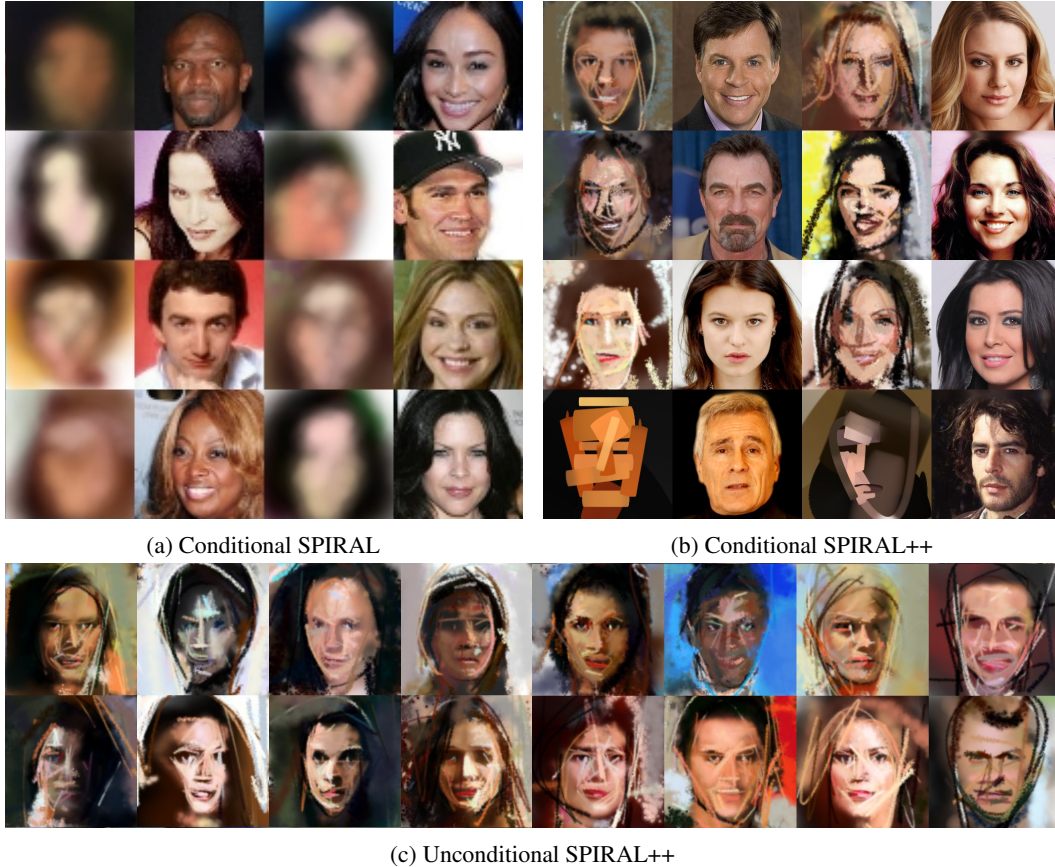


Figure 18: **Quality improvements over SPIRAL.** (a) Reconstructions taken from (Ganin et al., 2018). Note that these results are for CelebA (not CelebA-HQ used in the present paper). Unfortunately (Ganin et al., 2018) did not include generation results. (b) Selected reconstructions with our improvements. (c) Selected unconditional generation with our improvements.

**F.1 Spectral Normalization ablation.** Figure 19 shows the impact of the discriminator regularization modifications in subsection 2.1.

**F.2 Temporal Credit Assignment ablation.** As mentioned in subsection 2.2, agents without TCA often struggle to make good use of long episodes. This can be seen in Figure 20. It happens even if episode length is gradually increased as a curriculum. The impact on sample quality of adding TCA can be seen in Figure 21.

Remarkably, we were also able to train good agents with a discount factor of precisely zero, meaning they completely greedily paint one stroke at a time all the way from a blank white canvas to a final image. This suggests that in expectation the discriminator loss monotonically decreases from a blank canvas to a completed image. This works both when training the discriminator to reject intermediate canvases as fakes (in which case there is greater diversity of intermediate canvases the further away you get from a blank canvas) and when only training the discriminator on the final canvas.

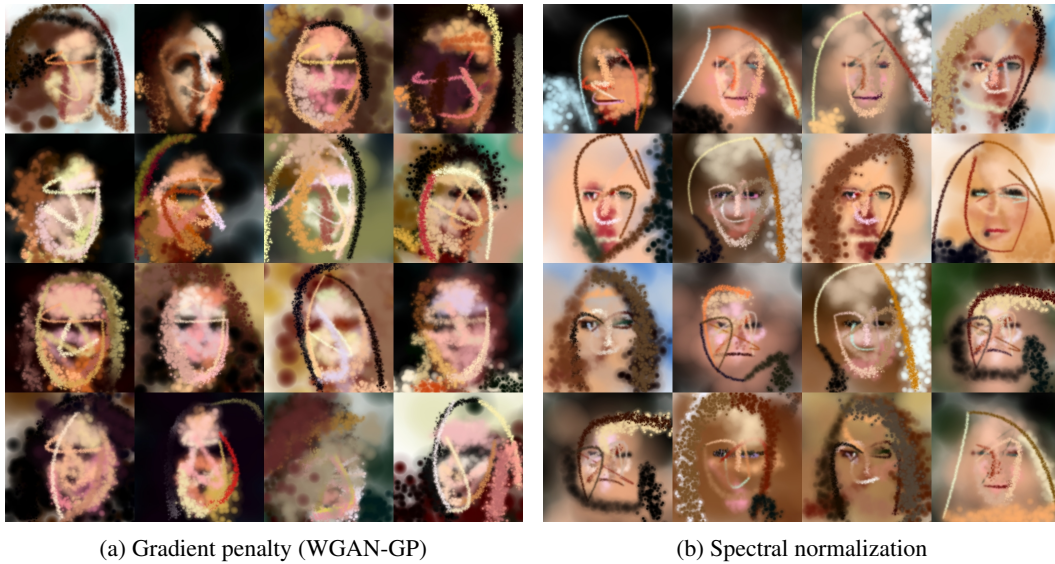


Figure 19: **Comparison of discriminator regularization on 19 step episodes.** Spectral Normalization lets the discriminator guide the agent into producing finer details like eyes, noses and mouths. These samples are from a single agent in the population; Figure 7 shows how different agents in the same population focus on reproducing different aspects of the target image distribution.

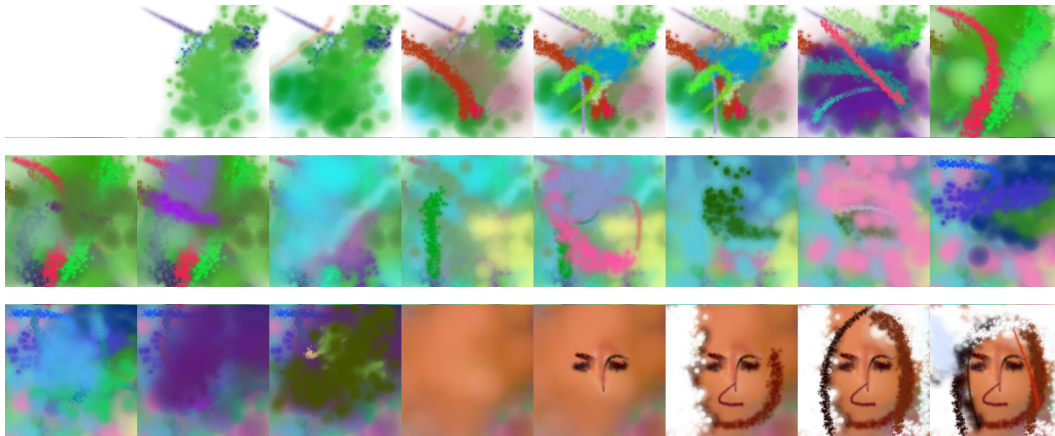


Figure 20: **Agents without TCA often waste actions.** These are 24 equally spaced frames taken from a single episode with 199 steps. The agent wastes its first 156 actions drawing brightly colored strokes that it eventually paints over, finally drawing a rough face in the final 43 steps. By comparison, agents with TCA start drawing immediately and gradually refine the image, as in Figure 9.

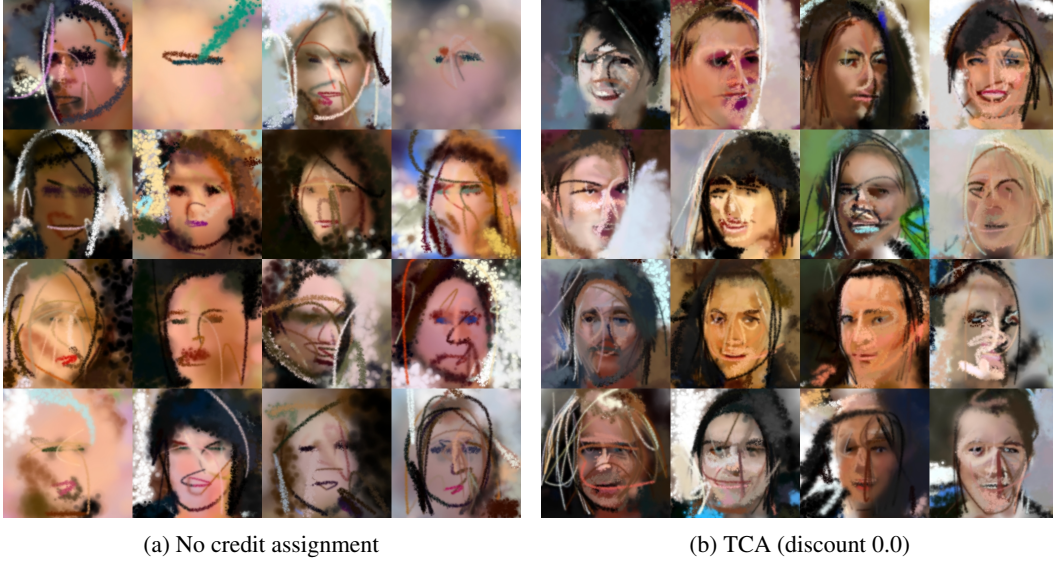


Figure 21: TCA lets agents make better use of 199 steps. Random unconditional samples.

**F.3 Complement discriminator ablation.** Figure 22 and Figure 23 contrast reconstruction with L2 loss or fully-conditioned discriminator against reconstruction with the Complement Discriminator from subsection 2.3.

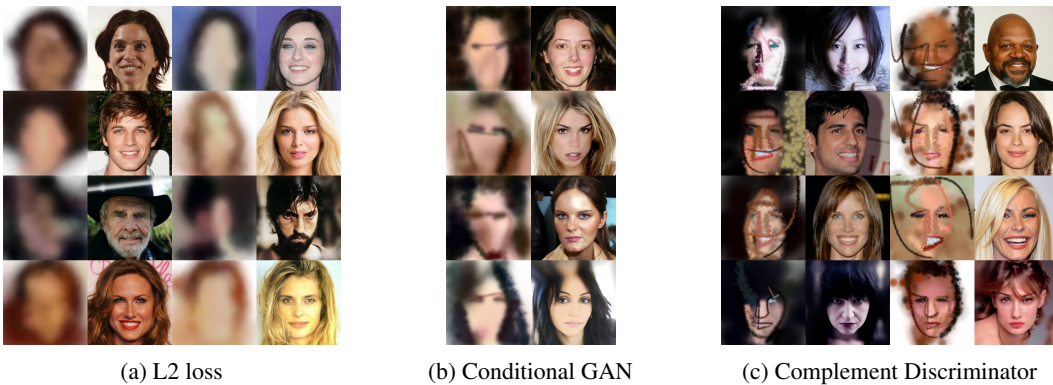


Figure 22: **Types of conditioning 20 step.** The Complement Discriminator gives rise to more interesting reconstructions that are semantically similar rather than similar in pixel space, for example mouths are open wider if the person was smiling widely.



Figure 23: **Types of conditioning 200 step.** Four 200-step agents reconstruct the target images in the top row. The top two agents paint bezier curves and were trained using L2 loss with discount  $\gamma$  of 0 and 0.9 respectively. The bottom two agents were both trained with Complement Discriminator, and paint bezier and spline curves respectively (see appendix A).

**F.4 Action space ablation.** When generative agents learn to operate human software like painting programs, they have a human-interpretable action space, and you can control the model by changing the actions offered by the environment. Figure 24 and Figure 25 show how retraining the agent with different brushes leads to very different styles.



Figure 24: **Training with different brushes leads to very different styles.** Agents with identical architectures perform reconstruction of the same target photos (top row) for 20 steps; only the brush they were trained with varies (except the first row which uses the Fluid Paint renderer, appendix G). Continued on next page.



Figure 25: **Training with different brushes leads to very different styles.** Agents with identical architectures perform reconstruction of the same target photos (top row) for 20 steps; only the brush they were trained with varies. Continued from previous page.

## Appendix G Comparison with recent model-based approaches to painting

Recently, several papers (Frans and Cheng, 2018; Nakano, 2019; Zheng et al., 2019; Huang et al., 2019) proposed to improve sample complexity and stability of generative agents by replacing the actual rendering simulator with a differentiable neural surrogate. The latter is trained offline to predict how a certain action (usually random) affects the state of the canvas. Although this is a promising avenue for research, we would like to mention several scenarios in which model-free approaches like SPIRAL and SPIRAL++ would be more suitable alternatives.

First, one of the most appealing features of the neural environments is that they allow to train agents by directly backpropagating the gradient coming from the objective of interest (*e.g.*, reconstruction loss or adversarial generator loss). Unfortunately, this means that one has to stick to continuous actions in order to avoid usage of brittle gradient approximators. Continuous actions may be fine for defining locations (*e.g.*, the end point in a Bézier curve) but are not appropriate for the situations requiring inherently discrete decisions like whether the agent should lift the brush or add one more control point to a spline. In appendix D, we show that SPIRAL++ performs reasonably well in this latter case.

Secondly, the success of the agent training largely depends on the quality of the neural model of the environment. The simulator used in Ganin et al. (2018) and in the experiments discussed so far is arguably easy to learn since new stokes interact with the existing drawing in a fairly straightforward manner. Environments which are highly stochastic or require handling of object occlusions and lighting might pose a challenge for neural network based environment models.

Lastly, while in principle possible, designing a model for a simulator with complex dynamics may be a non-trivial task. The majority of the recent works relying on neural renderers assume that the simulator state is fully represented by the appearance of the canvas and therefore only consider non-recurrent state transition models. There is no need for such an assumption in the SPIRAL framework. We demonstrate this advantage by training our agent in a new painting environment based on Fluid Paint (Li, 2017). A distinctive feature of this renderer is that under the hood it performs fluid simulation on a grid of cells. The simulation is governed by the Navier-Stokes equations and requires the access to the velocity field of the fluid (Fernando and others, 2004) which is not directly observable from the drawing. On top of that, unlike MyPaint, Fluid Paint models the behaviour of the brush bristles moving against the canvas surface. Despite having to deal with an arguably more complex setting, SPIRAL++ managed to not only learn the basic control of tool at hand but also exploit some of its peculiarities in an interesting way (*e.g.*, use of the brush bristles to imitate hair). Please refer to Figure 26 and the second row of Figure 24 for qualitative results.





Figure 26: **Selected unconditional 20-step samples.** Notice how these agent can paint hair in just one or two brush strokes, relying on the paint physics to draw many individual strands of hair.

## Appendix H Additional samples

Figure 27 and Figure 28 show further samples of agents that learned interesting painting styles.

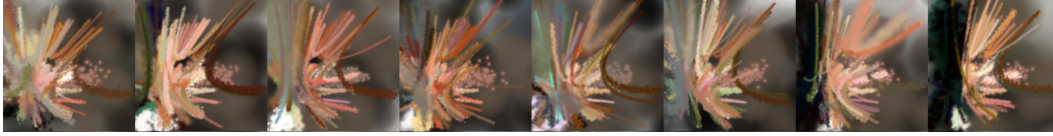


Figure 27: **Radial faces.** This agent learned a rather unusual style, building up faces with arcs originating from a single corner of the image.



Figure 28: **Focus on color.** Agents are occasionally Warholesque in their use of color. Perhaps these are compensating for other agents in the population using drab color palettes?

## Appendix I Hyperparameters

The generator learning rates and RL entropy costs were evolved using population based training (subsection 2.4), but we found that bad initial values could cause training to fail to converge. For all our experiments each generator sampled initial values of the learning rate from the range ( $1 \times 10^{-5}$ ,  $3 \times 10^{-4}$ ) and entropy cost from the range ( $2 \times 10^{-3}$ ,  $1 \times 10^{-1}$ ). As in Ganin et al. (2018) we trained the discriminator using Adam (Kingma and Ba, 2014) with a learning rate of  $1 \times 10^{-4}$ ,  $\beta_1$  set to 0.5 and  $\beta_2$  set to 0.999. For short episodes we used an RL discount factor  $\gamma$  of 1.0, but for long episodes using temporal credit assignment (subsection 2.2) we found that using values anywhere between 0 (completely greedy) and 0.99 gave good results, on episodes of length up to even 1000 steps. Most of our experiments use unroll lengths of 20 steps; others use unroll lengths of 50 steps with no noticeable performance differences (in particular, longer unroll lengths do not seem to reduce the need for temporal credit assignment (subsection 2.2)). We used batch sizes of 64 for both generators and discriminators. Brush strokes were rendered on a 256x256 canvas, but to reduce compute costs the canvas and images from the dataset were downsampled to 64x64 before being fed to the generator or discriminator.