
DiffInfinite: Large Mask-Image Synthesis via Parallel Random Patch Diffusion in Histopathology

Marco Aversa

University of Glasgow and Dotphoton
Glasgow, United Kingdom
marco.aversa@glasgow.ac.uk

Gabriel Nobis

Fraunhofer HHI
Berlin, Germany
gabriel.nobis@hhi.fraunhofer.de

Miriam Hägele

Aignostics
Berlin, Germany
miriam.haegele@aignostics.com

Kai Standvoss

Aignostics
Berlin, Germany
kai.standvoss@aignostics.com

Mihaela Chirica

Institute of Pathology, LMU Munich
Munich, Germany
mihaela.chirica@med.uni-muenchen.de

Roderick Murray-Smith

University of Glasgow
Glasgow, United Kingdom
roderick.murray-smith@glasgow.ac.uk

Ahmed Alaa

UC Berkeley
Berkeley, California
amalaa@berkeley.edu

Lukas Ruff

Aignostics
Berlin, Germany
lukas.ruff@aignostics.com

Daniela Ivanova

University of Glasgow
Glasgow, United Kingdom
daniela.ivanova@glasgow.ac.uk

Wojciech Samek

Fraunhofer HHI and TU Berlin
Berlin, Germany
wojciech.samek@hhi.fraunhofer.de

Frederick Klauschen

Institute of Pathology, LMU Munich
Munich, Germany
f.klauschen@lmu.de

Bruno Sanguinetti

Dotphoton
Zug, Switzerland
bruno.sanguinetti@dotphoton.com

Luis Oala

Dotphoton
Zug, Switzerland
luis.oala@dotphoton.com

Abstract

We present DiffInfinite, a hierarchical diffusion model that generates arbitrarily large histological images while preserving long-range correlation structural information. Our approach first generates synthetic segmentation masks, subsequently used as conditions for the high-fidelity generative diffusion process. The proposed sampling method can be scaled up to any desired image size while only requiring

small patches for fast training. Moreover, it can be parallelized more efficiently than previous large-content generation methods while avoiding tiling artifacts. The training leverages classifier-free guidance to augment a small, sparsely annotated dataset with unlabelled data. Our method alleviates unique challenges in histopathological imaging practice: large-scale information, costly manual annotation, and protective data handling. The biological plausibility of DiffInfinite data is evaluated in a survey by ten experienced pathologists as well as a downstream classification and segmentation task. Samples from the model score strongly on anti-copying metrics which is relevant for the protection of patient data.

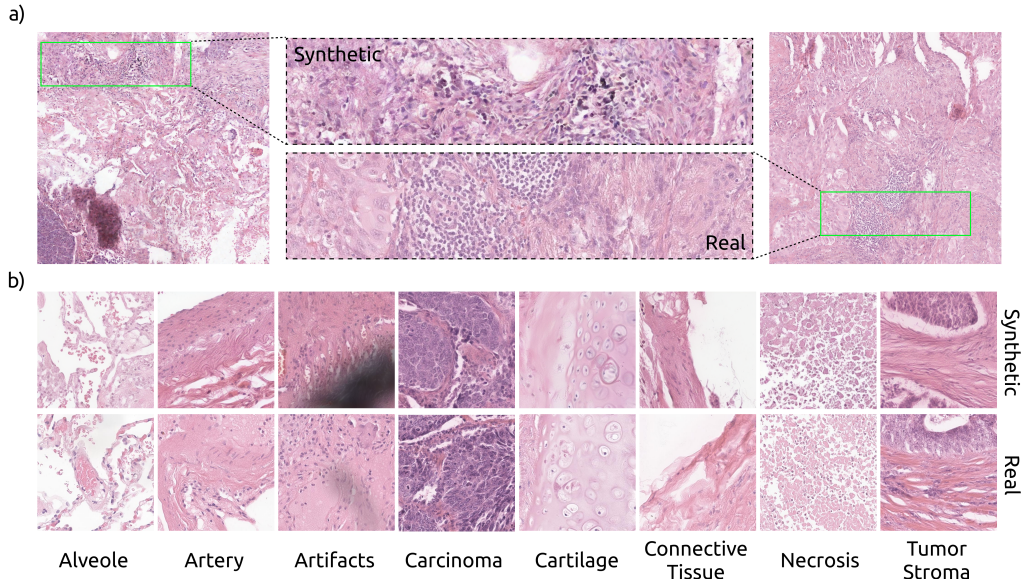


Figure 1: a) Examples of synthetic and real 2048×2048 images. b) Pairs of 512×512 synthetic tiles (top) with the closest real images found with Inception-v3 near-neighbour (bottom).

1 Introduction

Deep learning (DL) models are promising auxiliary tools for medical diagnosis [1–3]. Applications like segmentation and classification have been refined and pushed to the limit on natural images [4]. However, these models trained on rich datasets still have limited applications in medical data. While segmentation models rely on sharp object contours when applied to natural data, in medical imaging, the model struggles to detect a specific feature because it has a “limited ability to handle objects with missed boundaries” and often “miss tiny and low-contrast objects” [5, 6]. Therefore, task-specific medical applications require their own specialised and fine-grained annotation. Data labelling is arguably one of the most critical bottlenecks in healthcare machine learning (ML) applications. In histopathology, pathologists examine the histological slide at multiple levels, usually starting with a lower magnification to analyse the tissue architecture and cellular arrangement and gradually proceeding to a higher magnification to examine cell morphology and subcellular features, such as the appearance and number of nucleoli, chromatin density and cytoplasm appearance. Annotating features within gigapixel whole slide images (WSIs) with this level of detail demands effort and time, often leading to sparse, limited annotated data. In addition, due to privacy regulations and ethics [7, 8], having access to medical data can be challenging since it has been shown that it is possible to extract patients’ sensitive information [9] from this data.

In histopathology, state-of-the-art ML models require the context of the entire WSIs, with features at different scales, in order to distinguish between different tumor sub-types, grades and stages [10]. Despite the demonstrated effectiveness of diffusion models (DMs) in generating natural images compared to other approaches, they still have rarely been applied in medical imaging. Existing

generative models in histopathology can generate images of relatively small resolution compared to WSIs. To give a few examples, the application of Generative Adversarial Networks (GANs) in cervical dysplasia detection [11], glioma classification [12], and generating images of breast and colorectal cancer [13], generate images with 256×128 px, 384×384 px and 224×224 px, respectively. In spite of their current limitations in generating images at scales necessary to fully address all medical concerns, the use of synthetic data in medical imaging can provide a valuable solution to the persistent issue of data scarcity [14–17]. Models generally improve after data augmentation and synthetic images are equally informative as real images when added to the training set [18, 19]. Data augmentation could also help with the underrepresentation in data sets of rare cancer subtypes. By adding synthetic images to the training set, Chen et al. [20] demonstrated that their model had better accuracy in detecting chromophobe renal cell carcinoma, which is a rare subtype of renal cell carcinoma. Furthermore, Doleful et al. [21] showed how synthetic histological images could be used for educational purposes for pathology residents. Regarding the challenges highlighted before, we present a novel sampling method to generate large histological images with long-range pixel correlation (see Fig. 1), aiming to extend up to the resolution of the WSI.

Our contributions are as follows: 1) We introduce DiffInfinite, a hierarchical generative framework that generates arbitrarily large images, paired with their segmentation masks. 2) We introduce a fast outpainting method that can be efficiently parallelized. 3) The quality of DiffInfinite data is evaluated by ten experienced pathologists as well as downstream machine learnings tasks (classification and segmentation) and anti-duplication metrics to assess the leakage of patient data from the training set.

2 Related Work

Large-content image generation can be reduced to inpainting/outpainting tasks. Image inpainting is the problem of reconstructing unknown or unwanted areas within an image. A closely related task is image outpainting, which aims to predict visual content beyond the boundaries of an image. In both cases, the newly in- or outpainted image regions have to be visually indistinguishable with respect to the rest of the image. Such image completion approaches can help utilise models trained on smaller patches for the purpose of generating large images, by initially generating the first patch, followed by its extension outward in the desired direction.

Traditional approaches Traditional methods for image region completion rely on repurposing known image features, necessitating costly nearest neighbour searches for suitable pixels or patches [22–26]. Such methods often falter with complex or large regions [24]. In contrast, DL enables novel, realistic image synthesis for inpainting and outpainting. Some methods like Deep Image Prior [27] condition new image areas on the existing image, while others aim to learn natural image priors for realistic generation [28, 29].

Generative modelling for conditional image synthesis GANs have dominated image-to-image translation tasks like inpainting and outpainting for years [28–42]. Recently, DMs have surpassed GANs in various image generation tasks [43]. Palette [44] was the first to apply DMs to tasks like inpainting and outpainting. RePaint [45] and ControlNet [46] demonstrate resampling and masking techniques for conditioning using a pre-trained diffusion model. SinDiffusion [47] and DiffCollage [48] offer state-of-the-art outpainting solutions using DMs trained with overlapping patches. In parallel to our work, Bond-Taylor and Willcocks [49] developed a related approach called ∞ -Diff which trains on random coordinates, allowing the generation of infinite-resolution images during sampling. However, in contrast to our approach the method does not involve image compression in a latent space.

Synthetic data assessment The authenticity of synthetic data produced by DMs, trained on vast paired labelled datasets [50], remains contentious. Ethical implications necessitate distinguishing if generated images are replicas of training data [51, 52]. The task is complicated due to subjective visual similarities and diverse dataset ambiguities. Various metrics have been proposed for quantifying data replication, including information theory distances from real data [53], consistency measurements using downstream models [54, 55], comparison with inpainted areas [52], and detection of “forgotten” examples [56].

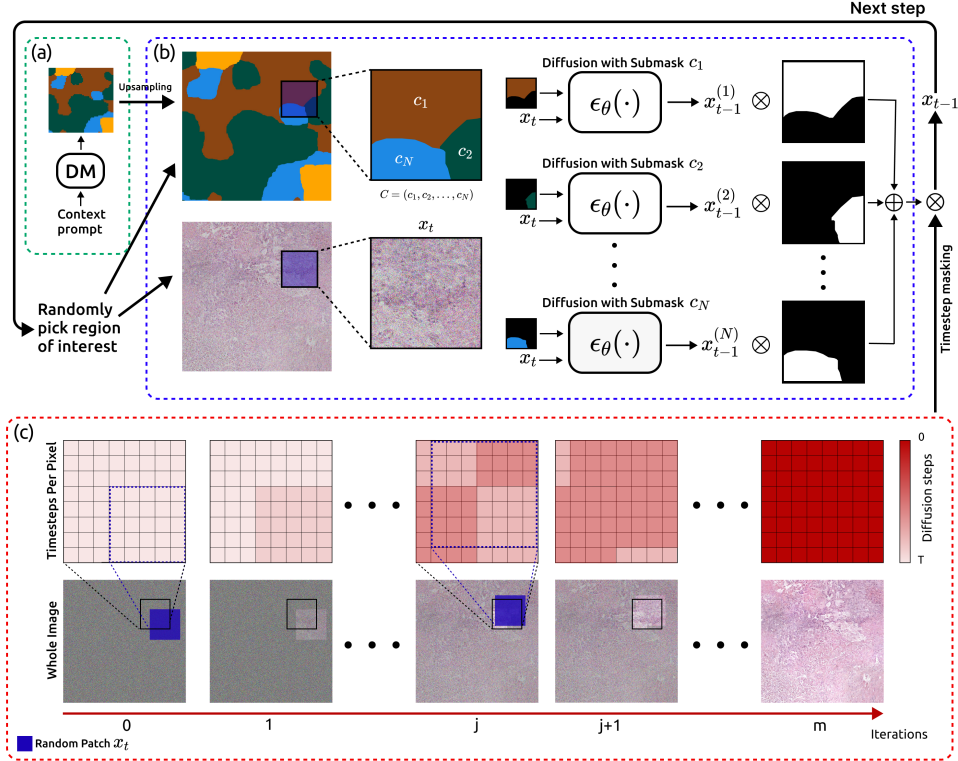


Figure 2: DiffInfinite generation method. a) Large-scale context mask generation. A diffusion model conditioned on a large-scale conditional prompt (e.g. Adenocarcinoma subtype) generates a low-resolution mask. The mask is upsampled via linear interpolation to the desired image size. b) Diffusion steps on large images. Given a random position, we select a sub-tile with its segmentation mask. A diffusion model generates in parallel the next step conditioned on each conditional label, or prompt, found in the mask. The outputs are masked individually with the corresponding label. The next step is the union of all the sub-patches. c) Tracking time steps pixel-wise. We keep track of the time step of each pixel in the large image. The model evolves only the pixels with the higher time step on each iteration.

3 Preliminaries

Diffusion Models DMs [57–59] represent a class of parameterized Markov chains that effectively optimize the lower variational bound associated with the likelihood function of the unknown data distribution. By iteratively adding small amounts of noise until the image signal is destroyed and then learning to reverse this process, DMs can approximate complex distributions much more faithfully than GANs [60]. The increased diversity of samples while preserving sample fidelity comes at the cost of training and sampling speed, with DMs being much slower than GANs [43]. The universally adopted solution to this problem is to encode the images from pixel space into a lower dimensional latent space via a Vector Quantised-Variational AutoEncoder (VQ-VAE), and perform the diffusion process over the latents, before decoding back to pixel space [61]. Pairing this with the Denoising Diffusion Implicit Models (DDIMs) sampling method [62] leads to faster sampling while preserving the DM objective

$$z_{t-1} = \sqrt{\alpha_{t-1}} \left(\frac{z_t - \sqrt{1 - \alpha_t} \epsilon_\theta(z_t, t)}{\sqrt{\alpha_t}} \right) + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \epsilon_\theta(z_t, t) + \sigma_t \epsilon_t, \quad (1)$$

where z_t is the latent variable at time step t in the VQ-VAE latent space, α_t is the noise scheduler, ϵ_θ is the noise learned by the model and ϵ_t is random noise. Conditioning can be achieved either by specifically feeding the condition with the noised data [44, 63], by guiding an unconditional model using an external classifier [64, 65] or by classifier-free guidance [66] used in this work, where the

convex combination

$$\tilde{\epsilon}_\theta(z_t, c) = (1 + \omega)\epsilon_\theta(z_t, c) - \omega\epsilon_\theta(z_t, \emptyset), \quad (2)$$

of a conditional diffusion model $\epsilon_\theta(z_t, c)$ and an unconditional model $\epsilon_\theta(z_t, \emptyset)$ is used for noise estimation. The parameter ω controls the tradeoff between conditioning and diversity, since $\omega > 0$ introduces more diversity in the generated data by considering the unconditional model while $\omega = 0$ uses only the conditional model.

4 Infinite Diffusion

The DiffInfinite approach we present here¹, is a generative algorithm to generate arbitrarily large images without imposing conditional independence, allowing for long-range correlation structural information. The method overcomes this limitation of DMs for large-content generation by deploying multiple realizations of a DM on smaller patches. In this section, we first define a mathematical description of this hierarchical generation model and then describe the sampling method paired with a masked conditioned generation process.

4.1 The Method

Let $X \sim \mathcal{X}$ be a large-content generating random variable taking values in \mathbb{R}^{KD} . Using the approach of latent diffusion models [61], the high-dimensional content is first mapped to the latent space \mathbb{R}^D by $\Phi(X) = Y \sim \mathcal{Y}_\Phi$. For simplicity, we assume throughout this work the existence of an ideal encoder-decoder pair (Φ, Ψ) such that $\Psi(\Phi(X)) = X$ is the identity on \mathbb{R}^{KD} . Assume further, to have a reverse time model $(SM_\theta, \epsilon_\theta)$ at hand consisting of a sampling method SM_θ and a learned model ϵ_θ trained on small patches $Z \sim \mathcal{Z}_\Phi$ taking values in \mathbb{R}^d . The reverse time model transforms $z_T \sim \mathcal{N}(0, I_d)$ over the time steps $t \in \{T, T-1, \dots, 1\}$ recursively by

$$z_{t-1} = SM_\theta(z_t) \quad (3)$$

to an approximate instance of \mathcal{Z}_Φ . We aim to sample instances from \mathcal{Y}_Φ by deploying multiple realizations of the reverse time model $(SM_\theta, \epsilon_\theta)$. Towards that goal, define the set of projections

$$\mathcal{C} := \{proj_I : \mathbb{R}^D \rightarrow \mathbb{R}^d \mid I \subset \mathbb{N} \text{ correspond to } d \text{ indices of connected pixels in } \mathbb{R}^D\}, \quad (4)$$

where $proj \in \mathcal{C}$ models a crop $proj(Y) \in \mathbb{R}^d$ of d connected pixels from the latent image Y . Since the model ϵ_θ is trained on images taking values in \mathbb{R}^d the standing assumption is

Assumption 1 Any projection $proj \in \mathcal{C}$ maps Y to the same distribution $proj(Y) \sim \mathcal{Z}_\Phi$ in \mathbb{R}^d .

Since the goal is to approximate an instance of \mathcal{Y}_Φ , we initialize the sampling method by $y_T \sim \mathcal{N}(0, I_D)$ and proceed in the following way: Given y_t , randomly choose $proj_{I_1}, \dots, proj_{I_m} \in \mathcal{C}$ independent of the state y_t such that $proj_{I_1}, \dots, proj_{I_m}$ are non equal crops that cover all latent pixels in \mathbb{R}^D . To be more precise, for every $i \in \{1, \dots, D\}$ we find at least one $j \in \{1, \dots, m\}$ with $i \in I_j$. For every projection $proj_{I_1}, \dots, proj_{I_m}$ we calculate the crop $z_t^j = proj_{I_j}(y_t)$ of the current state y_t and perform one step of the reverse time model following the sampling scheme

$$z_{t-1}^j = SM_\theta(z_t^j), \quad j \in \{1, \dots, m\}. \quad (5)$$

This results in overlapping estimates $z_{t-1}^1, \dots, z_{t-1}^m$ of the subsequent state $t-1$ and we simply assign to every pixel in the latent space the first value computed for this pixel such that

$$[y_{t-1}]_i = [z_{t-1}^j]_l, \quad \text{where } j = \min \{j' \mid i \in I_{j'}\} \quad (6)$$

and l refers to the entry in z_{t-1}^j corresponding to i with $[proj_{I_j}(y_{t-1})]_l = [y_{t-1}]_i$. Hence, starting from $y_T \sim \mathcal{N}(0, I_D)$ we sample in the first step from a distribution

$$y_{T-1} \sim p_{T-1, \theta}(y \mid y_T, proj_{I_1}, \dots, proj_{I_m}). \quad (7)$$

Using Bayes' theorem, this distribution simplifies to

$$p_{T-1, \theta}(y \mid y_T, proj_{I_1}, \dots, proj_{I_m}) = p_{T-1, \theta}(y \mid y_T), \quad (8)$$

¹Code available at <https://github.com/marcoaversa/diffinfinite>

since we sample the projections independently from y_T . Repeating the argument, we sample in every step from a distribution $y_{t-1} \sim p_{t-1,\theta}(y|y_t, \dots, y_T)$ over \mathbb{R}^D instead of sampling from $z_{t-1} \sim q_{t-1,\theta}(z|z_t, \dots, z_T)$ over \mathbb{R}^d . Hence, we approximate the true latent distribution \mathcal{Y}_Φ by the approximate distribution with density $p_{0,\theta}(y|y_1, \dots, y_T)$. In contrast to [48], our method does not use the assumption of conditional independence and the method can be applied to a wide range of DMs, without an adjustment of the training method. As the authors of [48] point out in their section on limitations, the assumption of conditional independence is not well-suited in cases of a data distribution with long-range dependence. For image generation in the medical context, we aim to circumvent this assumption as we do not want to claim that the density of a given region depends only on one neighboring region. The drawback of dropping the assumption is that we only approximate the reverse time model of the latent image distribution \mathcal{Y}_Φ indirectly, by multiple realizations of a reverse time model that approximates \mathcal{Z}_Φ .

4.2 Semi-supervised Guidance

In order to generate diverse high-fidelity data, DMs require lots of training data. Perhaps, training on a few samples still extracts significant features but it lacks variability, resulting in simple replicas. Here, we show how to enhance synthetic data diversity using classifier-free guidance as a semi-supervised learning method. In the classifier-free guidance [66], a single model is trained conditionally and unconditionally on the same dataset. We adapt the training scheme using two separate datasets. The model is guided by a small and sparsely annotated dataset q_1 , used for the conditional training step, while extracts features by the large unlabelled dataset q_0 , used on the unconditional training step (see Alg.1)

$$(z_0, c) = \begin{cases} (z_0, \emptyset) \sim q_0(z_0) & \text{if } u \geq p_{unc} \\ (z_0, c) \sim q_1(z_0, c) & \text{otherwise} \end{cases}, \quad (9)$$

where u is sampled from a uniform distribution in $[0,1]$, p_{unc} is the probability of switching from the conditional to the unconditional setting and \emptyset is a null label. During the sampling, a tradeoff between conditioning and diversity is controlled via the parameter ω in eq.2.

4.3 Sampling

High-level content generation The outputs of DMs have pixel consistency within the training image size. Outpainting an area with a generative model might lead to unrealistic and odd artifacts due to poor long-range spatial correlations. Here, we show how to predict pixels beyond the image’s boundaries by generating a hierarchical mapping of the data. The starting point is the generation of the highest-level representation of the data. In our case, it is the sketch of the cellular arrangement in the WSI (see Figure 2a). Since higher-frequency details are unnecessary at this stage, we can downsample the masks until the clustering pattern is still recognizable. The diffusion model, conditioned on the context prompt (e.g. Adenocarcinoma subtype), learns the segmentation masks which contain the cellular macro-structures information.

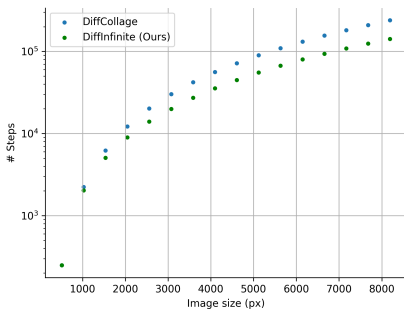


Figure 3: Comparison of sampling speed for DiffCollage and DiffInfinite, measuring diffusion steps required for image sampling. Demonstrating increased efficiency of DiffInfinite for larger images.

Random patch diffusion Given a segmentation mask M , we can proceed with the large image sampling according to Section 4.1 in the latent space \mathbb{R}^D of $Y = \Phi(X)$ (see Alg.2). Since we trained a conditional diffusion model with conditions c_1, \dots, c_N , the learned model takes the form $\epsilon_\theta(z_t, t) = (\epsilon_\theta(z_t, t|c_1), \dots, \epsilon_\theta(z_t, t|c_N))$. Given y_t , we first sample projections $proj_{I_1}, \dots, proj_{I_m} \in \mathcal{C}$, corresponding to different crops of d connected pixels up to the m -th projection with $\cup_{j=1}^m I_j = \{1, \dots, D\}$ and $\cup_{j=1}^m I_j \setminus \cup_{j=1}^{m-1} I_j \neq \emptyset$ (see the left hand-side of Figure 2b). Note that m is not fixed, but varies over the sampling steps and is upper bounded by the number of possible crops of d connected pixels. The random selection of

the projection is implemented such that regions with latent pixels of low projection coverage are more likely. Secondly, we calculate for every projection $j \in \{1, \dots, m\}$ the crop $proj_{I_j}(y_t)$ and perform one step of the DDIM sampling procedure using the classifier-free guidance model $(1 + \omega)\epsilon_\theta(z_t^j, t, c) - \omega\epsilon_\theta(z_t^j, t, \emptyset)$, where ϵ_θ is the learned model and $z_t^j = proj_{I_j}(y_t)$. This results for every pixels $i \in I_j$ in N values $DDIM_{\theta, c_1}(proj_{I_j}(y_t)), \dots, DDIM_{\theta, c_N}(proj_{I_j}(y_t))$, one for every condition c_i (see the right hand-side of Figure 2b). If $i \notin I_{j'}$ for all j' , the pixel i has not been considered yet and we assign i the value $[y_{t-1}]_i = [DDIM_{\theta, M_i}(proj_{I_j}(y_t))]_l$, where l corresponds to the pixel i under the projection I_j and M_i is the value of i in the mask M . Since we are updating random projections of the overall image, in the t -th step pixels either have the time index t or $t + 1$, resulting in a reversed diffusion process of differing time states. We initialize a tensor L_t , with the same size D as the latent variable, to keep track of the time index for each pixel. Each element is set to $L_T \equiv T$. In the j -th iteration of the t -th step we only update the pixels that have not been considered in one of the previous iterations of the t -th diffusion step, hence all the pixels in $i \in I_j$ with $proj_{I_j}(L_t)_i = t + 1$, similarly to the inpainting mask in the Repaint sampling method [45].

To restore the pixels that already received an update, i.e. every pixel $i \in I_j$ with $proj_{I_j}(L_t)_i = t$, we store a replica of the previous diffusion step for every pixel. Finally, we update all the time states in L_t that received an update in the j -th iteration to t resulting in $proj_{I_j}(L_t)_i = t$ for all $i \in I_j$. See the top row of Fig.2c for an illustration of the evolution of L_t . The random patch diffusion can also be applied to mask generation, where the only condition is the context prompt. This method can generate segmentation masks of arbitrary sizes with the correlation length bounded by two times the training mask image size.

Parallelization The sampling method proposed has several advantages. In Zhang et al. [48] each sequential patch is outpainted from the previous one with 50% of the pixels shared. Here, the randomization eventually leads to every possible overlap with the neighboring patches. This introduces a longer pixel correlation across the whole generated image, avoiding artifacts due to tiling. In Figure 3, we show that the number of steps in the whole large image generation process is drastically reduced with the random patching method with respect to the sliding window one. Moreover, in the sliding sampling method, the model can be paralleled only 2 or 4 times, depending if we are outpainting the image horizontally or on both axis. In our approach, we can parallelize the sampling up to the computational resource limit.

5 Data Assessment

To assess synthetic images for medical image analysis, we need to take various dimensions of data assessment into account. We extend traditional metrics from the natural image community with qualitative and quantitative assessments specific to the medical context. For the qualitative analysis, a team of pathologists evaluated the images for histological plausibility. The quantitative assessment entailed a proof-of-concept that a model can learn sensible features from the synthetically generated image patches for a relevant downstream task. As data protection is highly relevant regarding patient data, we performed evaluations to rule out memorization effects of the generative model.

Algorithm 1 DiffInfinite Training

Repeat
1: Randomly train on labelled or unlabelled data with probability p_{unc} .
 $u \sim Uniform[0, 1]$
 $(z_0, c) = \begin{cases} (z_0, \emptyset) \sim q_0(z_0) & \text{if } u \geq p_{unc} \\ (z_0, c) \sim q_1(z_0, c) & \text{otherwise} \end{cases}$
2: Sample random time step
 $t \sim Uniform\{1, \dots, T\}$
3: Sample noise, $\epsilon \sim \mathcal{N}(0, \mathbf{I})$
4: Corrupt data, $z_t = \gamma_t z_0 + \sigma_t \epsilon$
5: Take gradient descent step:
 $\nabla_\theta \|\epsilon - \epsilon_\theta(z_t, t, c)\|^2$
6: **until** converged

Algorithm 2 DiffInfinite Sampling

Input: High-level segmentation mask $M \in \mathbb{R}^D$ and learned model ϵ_θ
Output: Synthetic image X with the mask size
Initialization:
 $y_T \sim \mathcal{N}(0, \mathbf{I})$, index set $I_0 = \emptyset$ and time state tensor $L_T \equiv T$
Repeat
1: **for** $t \in \{T - 1, \dots, 0\}$ **do**
2: **while** $\cup_{j=0}^m I_j \neq \{1, \dots, D\}$ **do**
3: $m \leftarrow m + 1$
4: Select randomly $proj_{I_m} \in C \setminus \{proj_{I_1}, \dots, proj_{I_{m-1}}\}$
5: Crop $z_t^m = proj_{I_m}(y_t)$
6: **for** all conditions $n \in \{1, \dots, N\}$ **do**
7: DDIM sampling with classifier-free guidance
8: $z_{t-1}^m | c_n \sim p_{\theta, t}(z | z_t^m, c_n)$
9: **end for**
10: **for** all indices $i \in I_m$ **do**
11: **if** $i \notin I_j$ for all $j < m$ **then**
12: $[y_{t-1}]_i \leftarrow [z_{t-1}^m]_{M_i}$
13: $proj_{I_m}(L_t)_i \leftarrow t$
14: such that $[proj_{I_m}(y_{t-1})]_l = [z_{t-1}^m]_{M_i}$
15: **end if**
16: **end for**
17: **end while**
18: **end for**
19: $X \leftarrow \Psi(y_0)$

Table 1: Metrics to quantitatively evaluate the quality of the generated images. Left: scores for images of size 512×512 . DiffInfinite (a) first generates a mask and secondly an image following Section 4.1. Right: scores for real and generated images of size 2048×2048 resized to 512×512 . All methods use the same model trained on small patches of size 512×512 . DiffCollage corresponds to the method proposed in [48]. DiffInfinite (b) uses the real masks, while DiffInfinite (c) first generates a mask and secondly the large image. DiffInfinite (b) & (c) refers to the mixture of the generated dataset from DiffInfinite (b) and DiffInfinite (c).

	IP \uparrow	IR \uparrow	IS \uparrow	FID \downarrow		IP \uparrow	IR \uparrow
Morph-Diffusion [72]	0.26	0.85	2.1	20.1	DiffCollage	0.94	0.22
NASDM [73]	-	-	2.7	15.7	DiffInfinite (b)	0.95	0.48
DiffInfinite (a)	0.94	0.70	2.7	26.7	DiffInfinite (c)	0.98	0.44
					DiffInfinite (b) & (c)	0.98	0.33

5.1 Traditional Fidelity

We evaluate the fidelity of synthetic 512×512 images by calculating Improved Precision (IP) and Improved Recall (IR) metrics between 10240 real and synthetic images [67].² The IP evaluates synthetic data quality, while the IR measures data coverage. Despite their unsuitability for histological data [68, 69], Frechet-Inception Distance (FID) and Inception Score (IS) [70, 71] are reported for comparison with [72] and Shrivastava and Fletcher [73].³ The metrics’ explanations and formulas can be found in Appendix C.

In Table 1 (left), we report an IP of 0.94 and an IR of 0.70, indicating good quality and coverage of the generated samples. However, we note that these metrics are only somewhat comparable due to the different types of images generated by MorphDiffusion [72] and NASDM [73]. For the large images of size 2048×2048 , we rely solely on the IP and IR for quantitative evaluation due to the limited number of 200 generated large images. As shown in Figure 3(a) of [67], FID is unsuitable for evaluating such a small sample size, while IP and IR are more reliable. In Table 1 (right), we find that generating images first results in slightly higher IR, while generating the mask first achieves an IP of 0.98. For the sake of completeness we also report the scores then combining the two datasets. To compare our method to DiffCollage we generate 200 images using [48]. DiffInfinite performs better than DiffCollage wrt. to IP and IR. The drop of IR to 0.22 might be a result of the tiling artifacts observable in the LHS of Figure 11.

5.2 Domain Experts Assessment

To assess the histological plausibility of our generated images, we conducted a survey with a cohort of ten experienced pathologists, averaging 8.7 years of professional tenure. The pathologists were tasked with differentiating between our synthesized images and real image patches extracted from whole slide images. We included both small patches (512×512 px) as commonly used for downstream tasks as well as large patches (2048×2048 px). Including large patches enabled us to additionally evaluate the modelled long-range correlations in terms of transitions between tissue types as well as growth patterns which are usually not observable on the smaller patch sizes but essential in histopathology. In total the survey contained 60 images, in equal parts synthetic and real images as well as small and large patches. The overall ability of pathologists to discern between real and synthetic images was modest, with an accuracy of 63%, and an average reported confidence level of 2.22 on a 1-7 Likert scale. While we observed high inter-rater variance, there was no clear correlation between experience and accuracy ($r(8) = .069, p=.850$), nor between confidence level and accuracy ($r(8) = .446, p=.197$). Furthermore there was no significant correlation between the participants’ completion time of the survey and the number of correct responses ($r(8) = -.08, p=.826$).

Surprisingly, we found a similar performance for both, real and synthetic images. This indicates that, while clinical practice is mostly based on visual assessment, it is not a common task for pathologists to be restricted to parts of the whole slide image only. More detailed visualizations of the individual scores can be found in Appendix B. Besides this satisfactory result, we additionally wanted to

²<https://github.com/blandocs/improved-precision-and-recall-metric-pytorch>

³<https://github.com/toshast/torch-fidelity>

Table 2: Zero-shot evaluation results of the downstream tasks, encompassing both classification and segmentation scenarios. We employed three distinct models for each scenario: The first, "Trained Real," was trained using real data (in-house IH1), which also served as the training set for DiffInfinite. The second, "Trained Synthetic," was trained using samples generated from DiffInfinite, and the third, "Trained Augmented," utilized a combination of real and synthetic data. Our evaluation extends across separate lung cohorts (internal datasets IH2 and IH3) and additional indications (external datasets NCT, CRC, PCam), with varying degrees of data drift introduced.

	IH1	IH2	IH3	NCT-100K	CRC-7K	PCam-327K
Drift components	-	Patient change Different center		Patient change Different center Indication change		Patient change Different center Indication change Lower resolution
Trained Real	0.846 ± 0.005	0.733 ± 0.021	0.598 ± 0.049	0.857 ± 0.009	0.822 ± 0.034	0.628 ± 0.035
Trained Synthetic	0.747 ± 0.025	0.753 ± 0.005	0.699 ± 0.002	0.796 ± 0.023	0.753 ± 0.038	0.628 ± 0.012
Trained Augmented	0.852 ± 0.007	0.732 ± 0.027	0.637 ± 0.025	0.847 ± 0.044	0.811 ± 0.057	0.641 ± 0.035

(a) Classification results

	IH2
Trained Real	0.614 ± 0.009
Trained Synthetic	0.471 ± 0.039
Trained Augmented	0.710 ± 0.021

(b) Segmentation results

explore the limitations of our method by assessing the nuanced differences pathologists observed between synthetic and real images. While overall the structure and features seemed similar and hard to discern, they sometimes reported regions of inconsistent patterns, overly homogeneous chromatin in some of the synthetic nuclei, peculiarities in cellular and intercellular structures, and aesthetic elements. These seemed to be especially pronounced in tumorous regions where sometimes the tissue architecture appeared exaggerated, the transition to stroma or surrounding tissue was too abrupt and some cells lacked distinguishable nucleoli or cytoplasm. We attribute the nuanced effect of larger image size on the accuracy on this observation (cf. Fig. 5C). Overall the finding of the conducted survey demonstrates how complex the task of distinguishing between real and synthetically generated data is even for experienced pathologists while still highlighting potential areas to improve the generative model.

5.3 Synthetic Data for Downstream Tasks

A major interest in the availability of high quality labeled synthetic images is their use in downstream digital pathology applications. In this area, two primary challenges are the binary classification of images into cancerous or healthy tissues and the segmentation of distinct tissue areas in the tumor microenvironment. The unique ability of our technique to generate images of different cancer subtypes through the context prompt as well as the ability to create new segmentation masks and their corresponding H&E images specifically addresses these two challenges. Notably, expert annotations are costly and time consuming to acquire thus emphasizing the benefits of being able to train on purely synthetic datasets or augmenting annotated data in the low data regime. To showcase these two usecases we performed a series of experiments in both classification and segmentation settings. For all experiments, we trained a baseline classifier on a relatively small number of expert annotations IH1 (#patches = 3726) — the same that were used to train DiffInfinite — and additionally trained one model purely on synthetic data (IH1-S, #patches = 9974, $\omega = 0$), and one model on the real data augmented with the synthetic images. To generate target labels for the classification experiments, we simplified the segmentation challenge by categorizing patches with at least 0.05% of pixels labeled as 'Carcinoma' in the segmentation masks as 'Carcinoma'. All other patches were labeled 'Non-Carcinoma'. We evaluated all three classification models on several out-of-distribution datasets. We utilized two proprietary datasets (from the same cancer type with similar attributes but from distinct patient groups: IH1 (# patients=13, # patches=704) and IH3 (# patients=2, # patches=2817). Moreover, we assessed the models using two public datasets (NCK-CRC [74] and PatchCamelyon [75]), both representing tissue from different organs with distinct morphologies. Our findings, summarized in Table 2a, suggest that a classifier's out-of-distribution performance, trained with limited sample size and morphological diversity, can vary significantly (ranging from 0.628 to 0.857 balanced accuracy). This variability cannot be attributed solely to morphology but may also be influenced by

factors such as resolution and variations in scanning and staining techniques. Training exclusively with a larger set of synthetic images can enhance performance on some datasets (specifically IH2 and IH3), underscoring the advantages of leveraging the full training data in a semi-supervised manner within the generative model. Incorporating synthetic data as an augmentation to real data not only prevents the classifier’s performance decline, as seen on NCT-CRC and Patchcamelyon, on similar datasets but also bolsters its efficiency on more distinct ones. For the more challenging segmentation task we again trained three segmentation models to differentiate between carcinoma, stroma, necrosis, and a miscellaneous class that included all other tissue types, such as artifacts. The baseline performance of the real data model on a distinct group of lung patients (dataset IH2) of a F_1 score of 0.614 ± 0.009 (across three random seeds) highlights the difficulty of generalizing out of distribution in this tasks. While the purely synthetic model was not able to fully recover the baseline performance (0.471 ± 0.039), augmenting the small annotated dataset with synthetic data enhanced predictive performance to an F_1 score of 0.710 ± 0.021 . This boost of 10 percentage points in performance demonstrates that the synthetic data provide new, relevant information to the downstream task. In summary, our findings demonstrate the feasibility of meeting or surpassing baseline performance levels for both tasks using either entirely synthetic data or within an augmented context. Nevertheless, the advantages of employing synthetic data in downstream tasks continue to pose a challenge, not only within the medical image domain but also across various other domains [76–78], thus requiring more comprehensive assessment and thorough examination.

5.4 Considerations on Memorization

In medicine the adherence to privacy regulations is a sensitive requirement. While it is generally not possible for domain experts to infer patient identities from the image content of a histological tile or slide alone [79], developers and users of generative models are well advised to understand the risk of correspondence between the training data and the synthesized data. To this end, we evaluate the training and synthesized data against two memorization measures. The authenticity score $A \in [0, 1]$ by [54] aims to measure the rate by which a model generates new samples (higher score means more innovative samples). Similarly, [80] aims to estimate the degree of data copying C_T from the training data by the generative model. A $C_T \ll 0$ implies data copying, while a $C_T \gg 0$ implies an underfitting of the model. The closer to 0 the better. See Appendix C for a precise closed form of the measures and Table 5 for the full quantitative results, indicating that the DiffInfinite model is not prone to data copying across all resolutions and variations considered here ⁴. The A range between 0.86 and 0.98, signifying a high rate of authenticity. While other papers unfortunately do not report such detailed memorization statistics for their models, the results by [54] suggest that a score $\gg 0.8$ is not trivial to achieve. None of the models under consideration in [54] (VAE, DCGAN, WGAN-GP, ADS-GAN) achieve more than 0.82 in A on simpler data (MNIST). This interpretation is strengthened by the results of a $C_T \gg 0$ which indicates that the model might even be underfitting and is not in a data copying regime. Qualitative results on the nearest neighbour search between training and synthetic data in Figure 1 further corroborate these quantitative results.

6 Conclusions

DiffInfinite offers a novel sampling method to generate large images in digital pathology. Due to the high-level mask generation followed by the low-level image generation, synthetic images contain long-range correlations while maintaining high-quality details. Since the model trains and samples on small patches, it can be efficiently parallelized. We demonstrated that the classifier-free guidance can be extended to a semi-supervised learning method, expanding the labelled data feature space with unlabelled data. The biological plausibility of the synthetic images was assessed in a survey by 10 domain experts. Despite their training, most participants found it challenging to differentiate between real and synthetic data, reporting an average low confidence in their decisions. We found that samples from DiffInfinite can help in certain downstream machine learning tasks, on both in- as well as out-of-distribution datasets. Finally, authenticity metrics validate DiffInfinite’s capacity to generate novel data points with little similarity to the training data which is beneficial for the privacy preserving use of generative models in medicine.

⁴We use <https://github.com/marcojira/fls> from [81] to calculate both scores.

7 Acknowledgements

We would like to acknowledge our team of pathologists who provided valuable feedback in and outside of the conducted survey - special thank you to Frank Dubois, Niklas Prenissl, Cleopatra Schreiber, Vitaly Garg, Alexander Arnold, Sonia Villegas, Rosemarie Krupar and Simon Schallenberg. Furthermore, we would like to thank Marvin Sextro for his support in the analyses. This work was supported by the Federal Ministry of Education and Research (BMBF) as grants [SyReal (01IS21069B)]. RM-S is grateful for EPSRC support through grants EP/T00097X/1, EP/R018634/1 and EP/T021020/1, and DI for EP/R513222/1. MA is funded by Dotphoton, QuantIC and a UofG Ph.D. scholarship.

References

- [1] Ravi Aggarwal, Viknesh Sounderajah, Guy Martin, Daniel SW Ting, Alan Karthikesalingam, Dominic King, Hutan Ashrafian, and Ara Darzi. Diagnostic accuracy of deep learning in medical imaging: a systematic review and meta-analysis. *NPJ digital medicine*, 4(1):65, 2021.
- [2] Antonio Parziale, Monica Agrawal, Shengpu Tang, Kristen Severson, Luis Oala, Adarsh Subbaswamy, Sayantan Kumar, Elora Schoerverth, Stefan Hegselmann, Helen Zhou, Ghada Zamzmi, Purity Mugambi, Elena Sizikova, Girmaw Abebe Tadesse, Yuyin Zhou, Taylor Killian, Haoran Zhang, Fahad Kamran, Andrea Hobby, Mars Huang, Ahmed Alaa, Harvineet Singh, Irene Y. Chen, and Shalmali Joshi. Machine learning for health (ml4h) 2022. In Antonio Parziale, Monica Agrawal, Shalmali Joshi, Irene Y. Chen, Shengpu Tang, Luis Oala, and Adarsh Subbaswamy, editors, *Proceedings of the 2nd Machine Learning for Health symposium*, volume 193 of *Proceedings of Machine Learning Research*, pages 1–11. PMLR, 28 Nov 2022.
- [3] Maximilian Springenberg, Annika Frommholz, Markus Wenzel, Eva Weicken, Jackie Ma, and Nils Strodthoff. From modern cnns to vision transformers: Assessing the performance, robustness, and classification strategies of deep learning models in histopathology. *Medical Image Analysis*, 87:102809, 2023. ISSN 1361-8415.
- [4] Xinlong Wang, Xiaosong Zhang, Yue Cao, Wen Wang, Chunhua Shen, and Tiejun Huang. Seggpt: Segmenting everything in context. *arXiv preprint arXiv:2304.03284*, 2023.
- [5] Miriam Hägele, Philipp Seegerer, Sebastian Lapuschkin, Michael Bockmayr, Wojciech Samek, Frederick Klauschen, Klaus-Robert Müller, and Alexander Binder. Resolving challenges in deep learning-based analyses of histopathological images using explanation methods. *Scientific reports*, 10(1):1–12, 2020.
- [6] Jun Ma and Bo Wang. Segment anything in medical images. *arXiv preprint arXiv:2304.12306*, 2023.
- [7] Nicola Rieke, Jonny Hancox, Wenqi Li, Fausto Milletari, Holger R Roth, Shadi Albarqouni, Spyridon Bakas, Mathieu N Galtier, Bennett A Landman, Klaus Maier-Hein, et al. The future of digital health with federated learning. *NPJ digital medicine*, 3(1):119, 2020.
- [8] Luis Oala, Jana Fehr, Luca Gilli, Pradeep Balachandran, Alixandro Werneck Leite, Saul Calderon-Ramirez, Danny Xie Li, Gabriel Nobis, Erick Alejandro Muñoz Alvarado, Giovanna Jaramillo-Gutierrez, Christian Matek, Arun Shroff, Ferath Kherif, Bruno Sanguinetti, and Thomas Wiegand. Ml4h auditing: From paper to practice. In Emily Alsentzer, Matthew B. A. McDermott, Fabian Falck, Suproteem K. Sarkar, Subhrajit Roy, and Stephanie L. Hyland, editors, *Proceedings of the Machine Learning for Health NeurIPS Workshop*, volume 136 of *Proceedings of Machine Learning Research*, pages 280–317. PMLR, 11 Dec 2020.
- [9] Christopher G Schwarz, Walter K Kremers, Terry M Therneau, Richard R Sharp, Jeffrey L Gunter, Prashanthi Vemuri, Arvin Arani, Anthony J Spsychalla, Kejal Kantarci, David S Knopman, et al. Identification of anonymous mri research participants with face-recognition software. *New England Journal of Medicine*, 381(17):1684–1686, 2019.

- [10] Richard J Chen, Chengkuan Chen, Yicong Li, Tiffany Y Chen, Andrew D Trister, Rahul G Krishnan, and Faisal Mahmood. Scaling vision transformers to gigapixel images via hierarchical self-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16144–16155, 2022.
- [11] Yuan Xue, Jiarong Ye, Qianying Zhou, L Rodney Long, Sameer Antani, Zhiyun Xue, Carl Cornwell, Richard Zaino, Keith C Cheng, and Xiaolei Huang. Selective synthetic augmentation with histogan for improved histopathology image classification. *Medical image analysis*, 67:101816, 2021.
- [12] Le Hou, Ayush Agarwal, Dimitris Samaras, Tahsin M Kurc, Rajarsi R Gupta, and Joel H Saltz. Unsupervised histopathology image synthesis. *arXiv preprint arXiv:1712.05021*, 2017.
- [13] Adalberto Claudio Quiros, Roderick Murray-Smith, and Ke Yuan. PathologyGAN: learning deep representations of cancer tissue. *Journal of Machine Learning for Biomedical Imaging*, 4:1–48, 2021.
- [14] Ahmad B Qasim, Ivan Ezhov, Suprosanna Shit, Oliver Schoppe, Johannes C Paetzold, Anjany Sekuboyina, Florian Kofler, Jana Lipkova, Hongwei Li, and Bjoern Menze. Red-gan: Attacking class imbalance via conditioned generation. yet another medical imaging perspective. In *Medical Imaging with Deep Learning*, pages 655–668. PMLR, 2020.
- [15] Pierre Chambon, Christian Bluethgen, Jean-Benoit Delbrouck, Rogier Van der Sluijs, Małgorzata Połacin, Juan Manuel Zambrano Chaves, Tanishq Mathew Abraham, Shivanshu Purohit, Curtis P Langlotz, and Akshay Chaudhari. Roentgen: Vision-language foundation model for chest x-ray generation. *arXiv preprint arXiv:2211.12737*, 2022.
- [16] Pierre Chambon, Christian Bluethgen, Curtis P Langlotz, and Akshay Chaudhari. Adapting pretrained vision-language foundational models to medical imaging domains. *arXiv preprint arXiv:2210.04133*, 2022.
- [17] Luis Oala, Marco Aversa, Gabriel Nobis, Kurt Willis, Yoan Neuenschwander, Michèle Buck, Christian Matek, Jerome Extermann, Enrico Pomarico, Wojciech Samek, Roderick Murray-Smith, Christoph Clausen, and Bruno Sanguinetti. Data models for dataset drift controls in machine learning with optical images. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856.
- [18] Adrian B Levine, Jason Peng, David Farnell, Mitchell Nursey, Yiping Wang, Julia R Naso, Hezhen Ren, Hossein Farahani, Colin Chen, Derek Chiu, Aline Talhouk, Brandon Sheffield, Maziar Riazy, Philip P Ip, Carlos Parra-Herran, Anne Mills, Naveena Singh, Basile Tessier-Cloutier, Taylor Salisbury, Jonathan Lee, Tim Salcudean, Steven JM Jones, David G Huntsman, C Blake Gilks, Stephen Yip, and Ali Bashashati. Synthesis of diagnostic quality cancer pathology images by generative adversarial networks. *The Journal of Pathology*, 252(2):178–188, 2020.
- [19] Virginia Fernandez, Walter Hugo Lopez Pinaya, Pedro Borges, Petru-Daniel Tudosiu, Mark S Graham, Tom Vercauteren, and M Jorge Cardoso. Can segmentation models be trained with fully synthetically generated data? In *Simulation and Synthesis in Medical Imaging: 7th International Workshop, SASHIMI 2022, Held in Conjunction with MICCAI 2022, Singapore, September 18, 2022, Proceedings*, pages 79–90. Springer, 2022.
- [20] Richard Chen, Ming Lu, Tiffany Chen, Drew Williamson, and Faisal Mahmood. Synthetic data in machine learning for medicine and healthcare. *Nature Biomedical Engineering*, 5:1–5, 06 2021.
- [21] James M Dolezal, Rachele Wolk, Hanna M Hieromnimon, Frederick M Howard, Andrew Srisuwananukorn, Dmitry Karpeyev, Siddhi Ramesh, Sara Kochanny, Jung Woo Kwon, Meghana Agni, et al. Deep learning generates synthetic cancer histology for explainability and education. *NPJ Precision Oncology*, 7(1):49, 2023.
- [22] Alexei A Efros and Thomas K Leung. Texture synthesis by non-parametric sampling. In *Proceedings of the seventh IEEE international conference on computer vision*, volume 2, pages 1033–1038. IEEE, 1999.

- [23] Li-Yi Wei and Marc Levoy. Fast texture synthesis using tree-structured vector quantization. In *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, pages 479–488, 2000.
- [24] Marcelo Bertalmio, Guillermo Sapiro, Vincent Caselles, and Coloma Ballester. Image inpainting. In *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, pages 417–424, 2000.
- [25] Zongben Xu and Jian Sun. Image inpainting by patch propagation using patch sparsity. *IEEE transactions on image processing*, 19(5):1153–1165, 2010.
- [26] Lin Liang, Ce Liu, Ying-Qing Xu, Baining Guo, and Heung-Yeung Shum. Real-time texture synthesis by patch-based sampling. *ACM Transactions on Graphics (ToG)*, 20(3):127–150, 2001.
- [27] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Deep image prior. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9446–9454, 2018.
- [28] Avisek Lahiri, Arnav Kumar Jain, Sanskar Agrawal, Pabitra Mitra, and Prabir Kumar Biswas. Prior guided gan based semantic inpainting. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13696–13705, 2020.
- [29] Liang Liao, Jing Xiao, Zheng Wang, Chia-Wen Lin, and Shin’ichi Satoh. Guidance and evaluation: Semantic-aware image inpainting for mixed scenes. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVII 16*, pages 683–700. Springer, 2020.
- [30] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.
- [31] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2536–2544, 2016.
- [32] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Generative image inpainting with contextual attention. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5505–5514, 2018.
- [33] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Free-form image inpainting with gated convolution. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4471–4480, 2019.
- [34] Yi Wang, Ying-Cong Chen, Xin Tao, and Jiaya Jia. Vcnet: A robust approach to blind image inpainting. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXV 16*, pages 752–768. Springer, 2020.
- [35] Shengyu Zhao, Jonathan Cui, Yilun Sheng, Yue Dong, Xiao Liang, Eric I Chang, and Yan Xu. Large scale image completion via co-modulated generative adversarial networks. *arXiv preprint arXiv:2103.10428*, 2021.
- [36] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky. Resolution-robust large mask inpainting with fourier convolutions. *arXiv preprint arXiv:2109.07161*, 2021.
- [37] Mark Sabini and Gili Rusak. Painting outside the box: Image outpainting with gans. *arXiv preprint arXiv:1808.08483*, 2018.
- [38] Basile Van Hoorick. Image outpainting and harmonization using generative adversarial networks. *arXiv preprint arXiv:1912.10960*, 2019.

- [39] Chieh Hubert Lin, Yen-Chi Cheng, Hsin-Ying Lee, Sergey Tulyakov, and Ming-Hsuan Yang. InfinityGAN: Towards infinite-pixel image synthesis. In *International Conference on Learning Representations*, 2022.
- [40] Yen-Chi Cheng, Chieh Hubert Lin, Hsin-Ying Lee, Jian Ren, Sergey Tulyakov, and Ming-Hsuan Yang. Inout: diverse image outpainting via gan inversion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11431–11440, 2022.
- [41] Yaxiong Wang, Yunchao Wei, Xueming Qian, Li Zhu, and Yi Yang. Sketch-guided scenery image outpainting. *IEEE Transactions on Image Processing*, 30:2643–2655, 2021.
- [42] Yaxiong Wang, Yunchao Wei, Xueming Qian, Li Zhu, and Yi Yang. Rego: Reference-guided outpainting for scenery image. *arXiv preprint arXiv:2106.10601*, 2021.
- [43] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021.
- [44] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–10, 2022.
- [45] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11461–11471, 2022.
- [46] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*, 2023.
- [47] Weilun Wang, Jianmin Bao, Wengang Zhou, Dongdong Chen, Dong Chen, Lu Yuan, and Houqiang Li. Sindiffusion: Learning a diffusion model from a single natural image. *arXiv preprint arXiv:2211.12445*, 2022.
- [48] Qinsheng Zhang, Jiaming Song, Xun Huang, Yongxin Chen, and Ming-Yu Liu. Diffcollage: Parallel generation of large content with diffusion models. *arXiv preprint arXiv:2303.17076*, 2023.
- [49] Sam Bond-Taylor and Chris G. Willcocks. ∞ -diff: Infinite resolution diffusion with subsampled mollified states, 2023.
- [50] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *arXiv preprint arXiv:2210.08402*, 2022.
- [51] Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion art or digital forgery? investigating data replication in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6048–6058, 2023.
- [52] Nicholas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramèr, Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models. *arXiv preprint arXiv:2301.13188*, 2023.
- [53] Nikhil Vyas, Sham Kakade, and Boaz Barak. Provable copyright protection for generative models. *arXiv preprint arXiv:2302.10870*, 2023.
- [54] Ahmed Alaa, Boris Van Breugel, Evgeny S Saveliev, and Mihaela van der Schaar. How faithful is your synthetic data? sample-level metrics for evaluating and auditing generative models. In *International Conference on Machine Learning*, pages 290–306. PMLR, 2022.
- [55] Luis Oala. *Metrological machine learning (2ML)*. 1 edition, 2023. URL <https://metrological.ml>.

- [56] Matthew Jagielski, Om Thakkar, Florian Tramer, Daphne Ippolito, Katherine Lee, Nicholas Carlini, Eric Wallace, Shuang Song, Abhradeep Thakurta, Nicolas Papernot, et al. Measuring forgetting of memorized training examples. *arXiv preprint arXiv:2207.00099*, 2022.
- [57] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015.
- [58] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. *Advances in neural information processing systems*, 33:12438–12448, 2020.
- [59] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851. Curran Associates, Inc., 2020.
- [60] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021.
- [61] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021.
- [62] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2020.
- [63] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *J. Mach. Learn. Res.*, 23(47):1–33, 2022.
- [64] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- [65] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2020.
- [66] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- [67] Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved precision and recall metric for assessing generative models. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [68] Tuomas Kynkäänniemi, Tero Karras, Miika Aittala, Timo Aila, and Jaakko Lehtinen. The role of imagenet classes in fréchet inception distance. In *Proc. ICLR*, 2023.
- [69] Shane T. Barratt and Rishi Sharma. A note on the inception score. *ArXiv*, abs/1801.01973, 2018.
- [70] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [71] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, Xi Chen, and Xi Chen. Improved techniques for training gans. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- [72] Puria Azadi Moghadam, Sanne Van Dalen, Karina C. Martin, Jochen Lennerz, Stephen Yip, Hossein Farahani, and Ali Bashashati. A morphology focused diffusion probabilistic model for synthesis of histopathology images. *arXiv preprint arXiv: arXiv:2209.13167v2*, 2022.

- [73] Aman Shrivastava and P Thomas Fletcher. Nasdm: Nuclei-aware semantic histopathology image generation using diffusion models. *arXiv preprint arXiv:2303.11477*, 2023.
- [74] Jakob Nikolas Kather, Niels Halama, and Alexander Marx. 100,000 histological images of human colorectal cancer and healthy tissue, April 2018. URL <https://doi.org/10.5281/zenodo.1214456>.
- [75] Bastiaan S Veeling, Jasper Linmans, Jim Winkens, Taco Cohen, and Max Welling. Rotation equivariant cnns for digital pathology. In *Medical Image Computing and Computer Assisted Intervention—MICCAI 2018: 21st International Conference, Granada, Spain, September 16-20, 2018, Proceedings, Part II 11*, pages 210–218. Springer, 2018.
- [76] Shenghuan Sun, Gregory M Goldgof, Atul Butte, and Ahmed M Alaa. Aligning synthetic medical images with clinical knowledge using human feedback. *arXiv preprint arXiv:2306.12438*, 2023.
- [77] Xian Li, Ping Yu, Chunting Zhou, Timo Schick, Luke Zettlemoyer, Omer Levy, Jason Weston, and Mike Lewis. Self-alignment with instruction backtranslation. *arXiv preprint arXiv:2308.06259*, 2023.
- [78] Iliia Shumailov, Zakhar Shumaylov, Yiren Zhao, Yarin Gal, Nicolas Papernot, and Ross Anderson. The curse of recursion: Training on generated data makes models forget. *arXiv preprint arXiv:2305.17493*, 2023.
- [79] Petr Holub, Heimo Müller, Tomáš Bíl, Luca Pireddu, Markus Plass, Fabian Prasser, Irene Schlünder, Kurt Zatloukal, Rudolf Nenutil, and Tomáš Brázdil. Privacy risks of whole-slide image sharing in digital pathology. *Nature Communications*, 14(1):2577, 2023.
- [80] Casey Meehan, Kamalika Chaudhuri, and Sanjoy Dasgupta. A non-parametric test to detect data-copying in generative models. *CoRR*, abs/2004.05675, 2020.
- [81] Marco Jiralerspong, Avishek Joey Bose, Ian Gemp, Chongli Qin, Yoram Bachrach, and Gauthier Gidel. Feature likelihood score: Evaluating generalization of generative models using samples, 2023.
- [82] Nicolas Pielawski and Carolina Wählby. Introducing hann windows for reducing edge-effects in patch-based image segmentation. *PloS one*, 15(3):e0229839, 2020.

A Histological Dataset

Table 3: Details of the histological dataset

Histological dataset	
Number of whole slide images	41
Image type	H&E-stained whole slide images
Whole slide image size	$\sim 100,000 \times 100,000$
Magnification	40x
Image scanner	Aperio scanner
Number of annotation categories	40
Annotation distribution	37% Carcinoma, 36% Stroma, 3.5% Necrosis, 23.5% Other
Resolution	0.5 microns per pixel
Number of patches (image training)	4,781 labelled + 255,643 unlabelled
Patch size	512×512 px
Train/Test split	90/10 stratified by annotation categories
Number of patches (large mask training)	1,183
Patch size	2048×2048 px

The real-world data used for training the generative model consisted of 41 high-resolution Hematoxylin and Eosin (H&E)-stained whole slide images of lung tissue biopsies from different cancer patients. These images were evenly split between cases diagnosed with adenocarcinoma of the lung and squamous cell carcinoma, representing the two most common sub-types in lung cancer. The images were scanned on an Aperio scanner at a resolution of 0.25 microns per pixel (40x). Different classes used for conditioning were annotated digitally by a pathologist using an apple pencil with the instruction to clearly demarcate boundaries between tissue regions. The pathologist could choose from a list of 40 distinct annotation categories, aiming to cover all possible annotation requirements. 37% of the annotations belonged to the Carcinoma category, 36% to Stroma, 3.5% to Necrosis and the remaining 23.5% to other smaller categories summarized as Other. All data handling was performed in strict accordance with privacy regulations and ethical standards, ensuring the protection of patient information at all times. For training the diffusion model, we utilized a patch dataset derived from expert annotations. In total, the dataset contained 4,781 patches of size 512×512 px. The dataset was split into train/ test sets with a ratio of 90/ 10, stratified by annotation categories. This test split was used for the generative model as well as to evaluate the downstream task. We also tiled the slides with size 2048×2048 from the same annotations, extracting 1,183 patches. These masks are used for training the mask generative model.

A.1 Downstream task datasets

For the downstream tasks we utilized additional internal and external datasets to assess the predictive performance of our models.

Table 4: Details of downstream task datasets

Dataset	Indication	Patch size	mpp	Number of patches
IH1	lung	512	0.5	3.7 K
IH2	lung	512	0.5	0.7 K
IH3	lung	512	0.5	2.8 K
NCT	colorectal	224	0.5	100.0 K
CRC	colorectal	224	0.5	7.0 K
PCam	lymph nodes	96	0.972	327.0 K

B Survey

The survey was sent out to 10 pathologists with varying years of experience. Fig.4 shows the setup of the survey. The presented images were shown in randomized order.

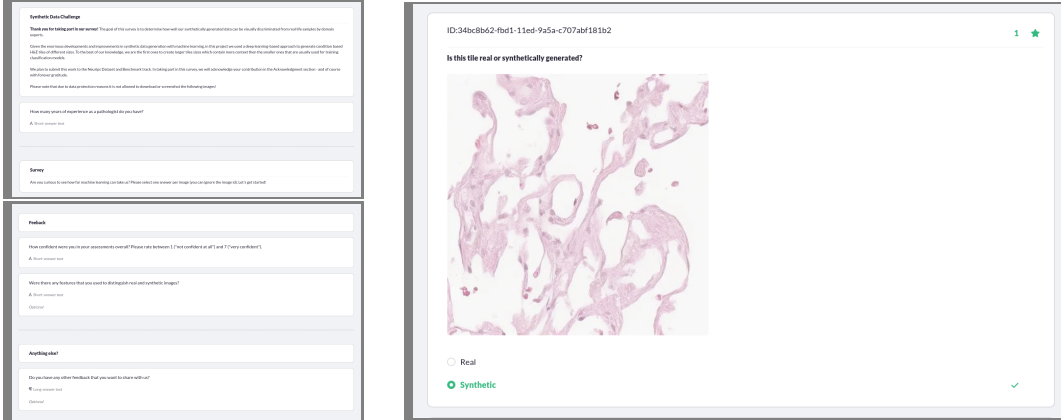


Figure 4: Survey interface for the domain expert assessment of real versus synthetic data.

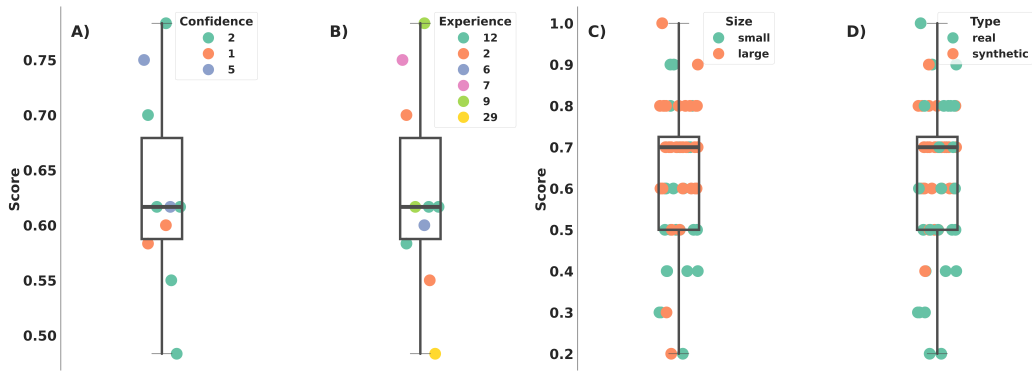


Figure 5: Results of the survey. Left: Accuracy per pathologist, color-coded by subjective confidence level (A) and years of experience (B). Right: Average accuracy across pathologists for each image patch, color-coded by path-size (C) and veracity (D).

C Metrics for Data Assessment

In this section we provide the definitions of the metrics used in Section 5 to assess the fidelity and degree of memorization of DiffInfinite. Following the notation of Section 4.1, denote by $X_r \sim \mathcal{X}_r$ the real data distribution and by $X_g = \Psi(\hat{Y}) \sim \mathcal{X}_g$ the distribution from which the generative model samples. For the quantitative evaluation of the quality and the coverage of the data generated by DiffInfinite we use

Improved recall and improved precision [67] A pre-trained classifier⁵ maps the samples into a high-dimensional feature space resulting in the feature vectors Φ_r and Φ_g . For $\Phi \in \{\Phi_r, \Phi_g\}$ denote by $NN_k(\phi', \Phi)$ the k th nearest feature vector of ϕ' from set Φ and define the binary function

$$f(\phi, \Phi) = \begin{cases} 1, & \text{if } \|\phi - \phi'\|_2 \leq \|\phi' - NN_k(\phi', \Phi)\|_2 \text{ for at least one } \phi' \in \Phi \\ 0, & \text{otherwise} \end{cases} \quad (10)$$

that identifies whether a given sample ϕ is within the estimated manifold volume of Φ corresponding to NN_k . To measure the similarity of Φ_g to the estimated manifold of the real images, define improved precision (IP) by

$$\text{precision}(\Phi_r, \Phi_g) = \frac{1}{|\Phi_g|} \sum_{\phi_g \in \Phi_g} f(\phi_g, \Phi_r) \quad (11)$$

and to measure the similarity of Φ_r to the estimated manifold of the generated images, define improved recall (IR) by

$$\text{recall}(\Phi_r, \Phi_g) = \frac{1}{|\Phi_r|} \sum_{\phi_r \in \Phi_r} f(\phi_r, \Phi_g). \quad (12)$$

The rate of DiffInfinite to innovate a new sample is approximated by the

Authenticity score [54] For the definition of the authenticity score $A \in [0, 1]$, assume that the probability measure \mathbb{P}_g corresponding to \mathcal{X}_g is a mixture of the probability measures

$$\mathbb{P}_g = A \cdot \mathbb{P}'_g + (1 - A) \cdot \delta_{g,\epsilon}, \quad (13)$$

where \mathbb{P}'_g characterizes the generative distribution, excluding synthetic samples that are duplicates of training samples and $\delta_{g,\epsilon} = \delta_g * \mathcal{N}(0, \epsilon^2)$ is the noisy distribution over training data implied by an unknown discrete probability measure δ_g placing probability mass on each data point used for training.

To test DiffInfinite for data-copying we compute the

C_T score [80] For a set of training images $\mathcal{D}_{train} = \{x_1, \dots, x_k | x_i \sim \mathcal{X}_r\}$ and $y \in \mathbb{R}^{KD}$ define the distance measure $d(y) = \min_{x \in \mathcal{D}_{train}} \|x - y\|_2^2$. Denote by $L(\mathcal{V})$ the one dimensional distribution $d(V)$ of any random variable $V \sim \mathcal{V}$ with the same instance space as \mathcal{X}_r . For the test set of the real data $\mathcal{D}_{test} = \{y_1, \dots, y_n | y_i \sim \mathcal{X}_r\}$, define the fraction $P_n(\pi) = |\{y \in \mathcal{D}_{test} | y \in \pi \in \Pi\}| / n$ of test points in cell $\pi \in \Pi$, where Π is a partition of \mathbb{R}^{KD} resulting from applying the k -means algorithm on \mathcal{D}_{train} . Similar for a set of generated images $\mathcal{D}_{gen} = \{\hat{x}_1, \dots, \hat{x}_m\}$ sampled from \mathcal{X}_g , define the fraction $Q_n(\pi)$ of generated samples in cell $\pi \in \Pi$. Denote by Z_U the z -scored Mann-Whitney U statistic from Section 3.1 of [80] with $L_\pi(\mathcal{D}) = \{d(x) | x \in \mathcal{D}, \pi \in \Pi\}$ for $\mathcal{D} \in \{\mathcal{D}_{test}, \mathcal{D}_{gen}\}$ and let Π_τ be the set of all cells in Π for which $Q_m(\pi) \geq \tau$ holds true. The C_T score is finally defined as the average

$$C_T(P_n, Q_m) = \frac{\sum_{\pi \in \Pi_\tau} P_n(\pi) Z_U(L_\pi(P_n), L_\pi(Q_m); T)}{\sum_{\pi \in \Pi_\tau} P_n(\pi)}. \quad (14)$$

across all cells represented by \mathcal{X}_g .

⁵We use the pre-trained VGG-16 classifier from <https://github.com/blandocs/improved-precision-and-recall-metric-pytorch>.

Table 5: Quantitative memorization metrics for the variants of DiffInfinite described in Section 5.1. For consistency, we consider all methods from Table 1 in our evaluation, including the comparison to DiffCollage. For the methods that output a large image of size 2048 we consider the *tiled* patches resulting in 16 patches per large image and the *resized* image resulting in 200 images of size 512×512 .

	$A \uparrow$		$C_T \downarrow \uparrow$	
	<i>tiled</i>	<i>resized</i>	<i>tiled</i>	<i>resized</i>
DiffCollage	0.89	0.97	11.02	7.00
DiffInfinite (a)	0.86	-	4.99	-
DiffInfinite (b)	0.86	0.97	3.29	8.11
DiffInfinite (c)	0.86	0.98	9.61	11.56
DiffInfinite (b) & (c)	0.87	0.95	5.31	10.96

D Data Samples

Mask-image pairs In Fig. 6, we show the control on the mask-image generation for 512×512 patches. The *Unknown* class corresponds to pixels which were not annotated due to a sparse annotation strategy. The images show that the cross-attention layer controls mask conditioning effectively. As a proof of concept, we generated images at different scales ($512 \times 512, 1024 \times 1024, 2048 \times 2048$) with a simple squares mask (see Fig. 7). In Figure 8, we see that for the small masks of size 512×512 , the frequency of labels in the real masks are reproduced well by the generated masks. For the large masks of size 2048×2048 , the labels that occur most frequently in the real masks are underrepresented in the generated masks, while all other labels are overrepresented in the generated masks.

Random patch advantages Sampling with the random patch (RP) method leads to several benefits compared to the sliding windows (SW) approach (see Fig. 11). First, the sliding window method starts from the centre of the image and outpaints in four directions. As a consequence, the model needs to condition on previously generated areas, leading to blurriness on the farther pixels. With the random patch method, every area is conditioned only on its neighbour, avoiding error propagation. Moreover, while SWs have only information on the closest neighbour, RPs consider long-range correlations. On every diffusion step, we have every possible overlap between near patches, extending correlation lengths to twice the diffusion model output size. Furthermore, this random overlap avoids any tiling effect.

Inpainting Using the segmentation images and masks of the test set, we inpainted the annotated areas with the same corresponding class (see Fig. 9). We show that the model generates new content respect to the real one. We run the same experiment by inpainting one area with different classes (see Fig. 10). Keeping the same seed, we show how the generation changes while ω increases. By increasing ω , we enhance the diversity at the cost of losing some conditioning.

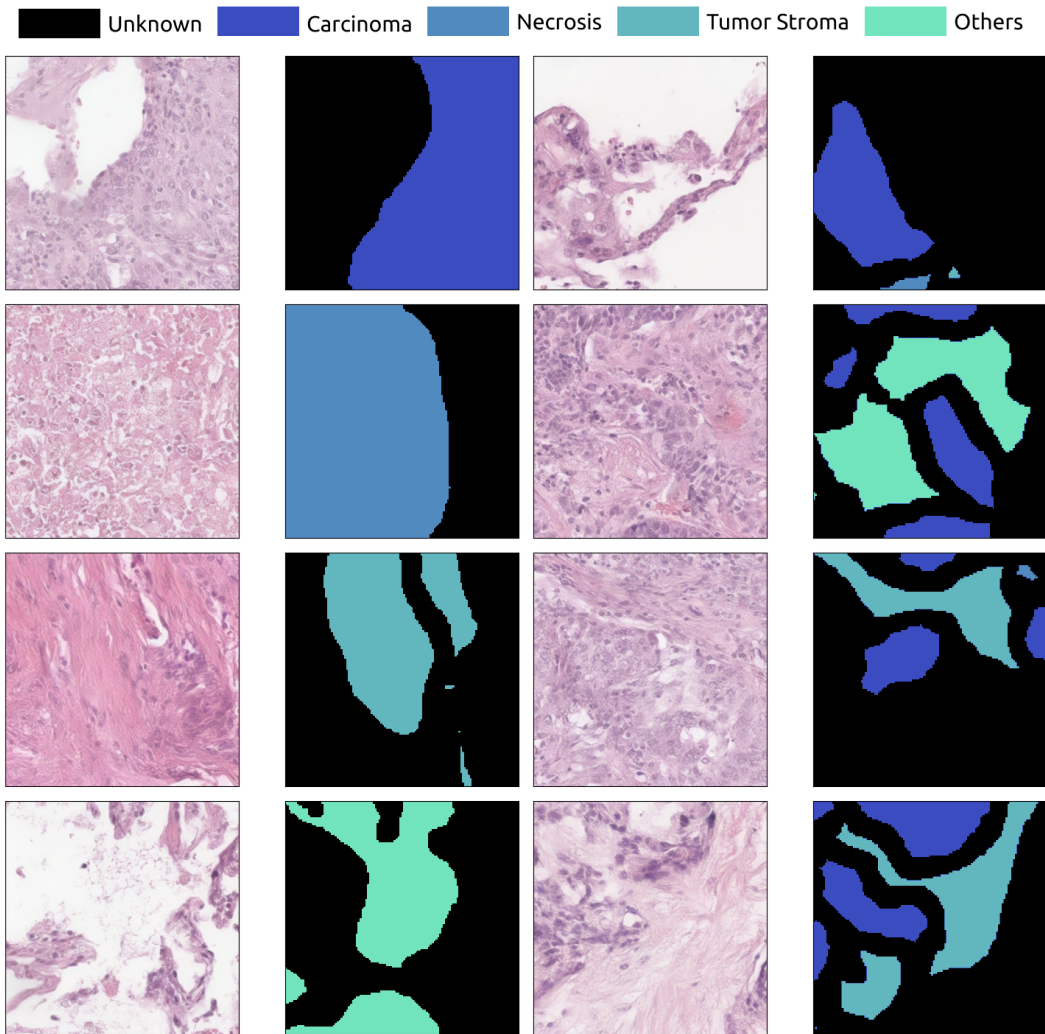


Figure 6: Generated images conditioned on the synthetic segmentation masks.

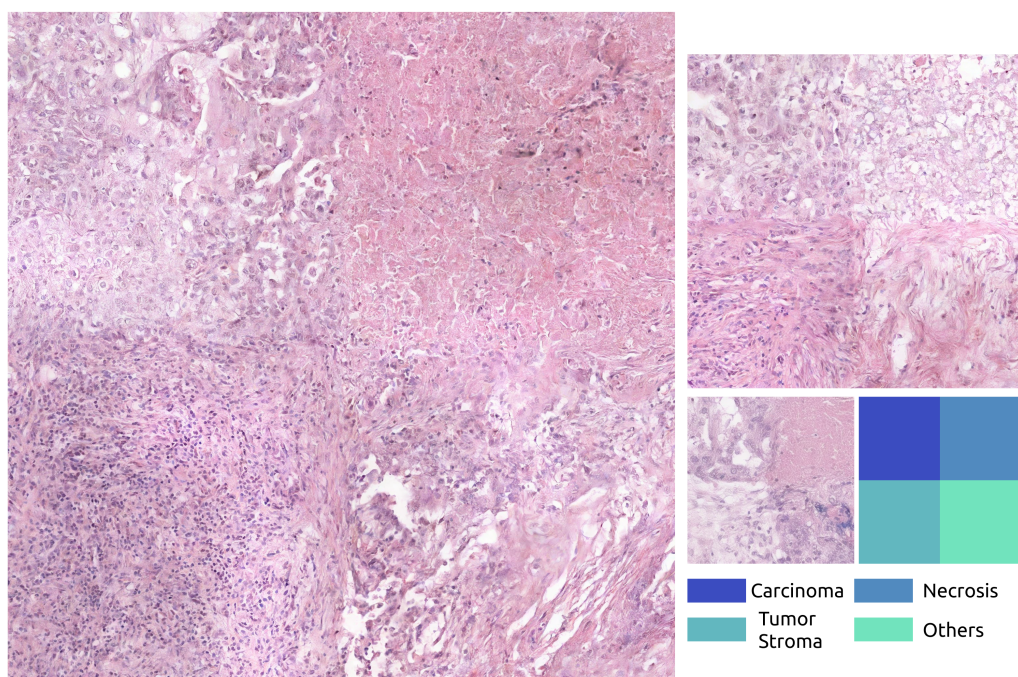


Figure 7: Conditioning visualization. All the images are conditioned with the squared mask shown. Left) 2048 × 2048 image. Top-Right) 1024 × 1024 image. Bottom-Right) 512 × 512 images.

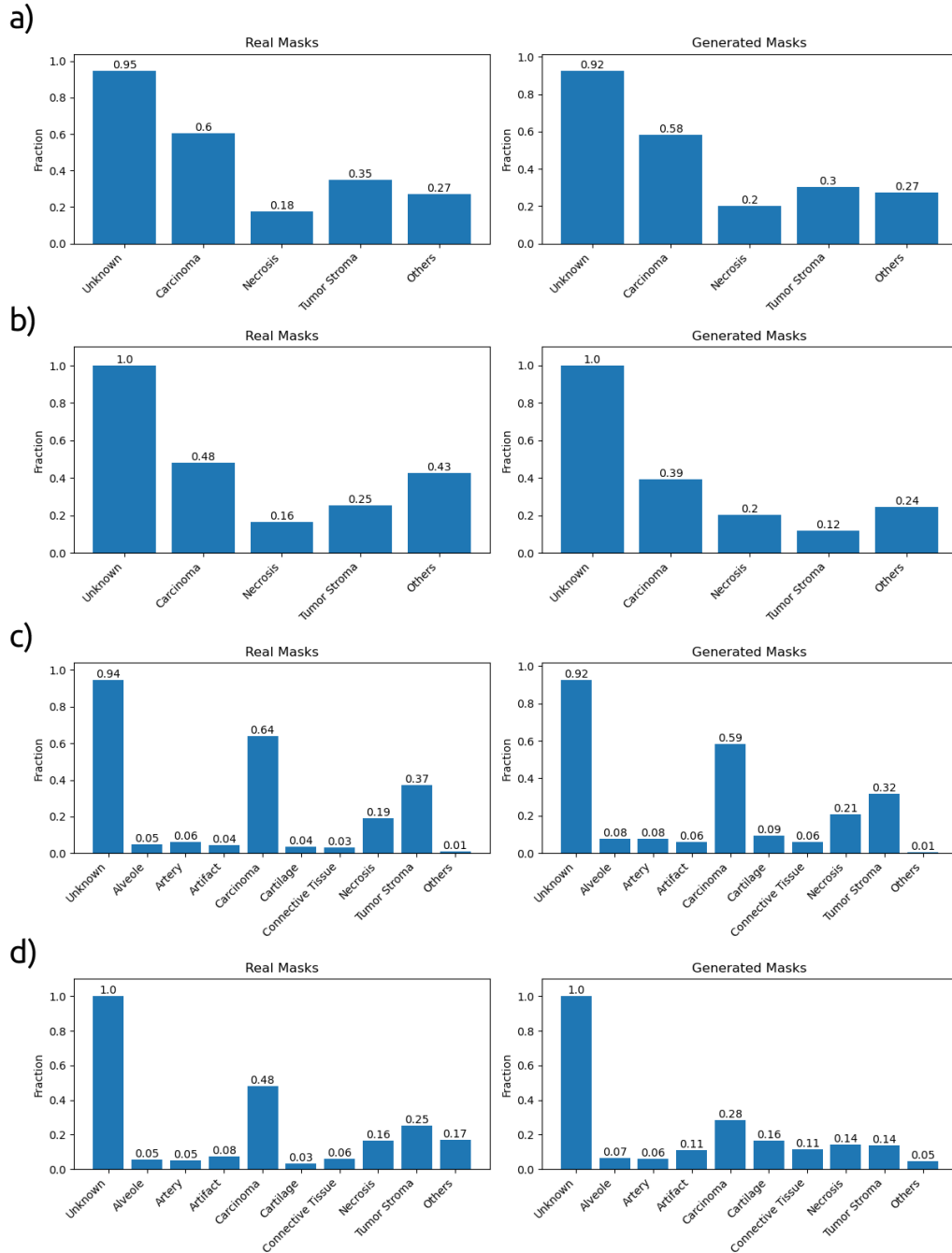


Figure 8: Fraction of label appearance in the segmentation masks with 5 classes in a,b) and 10 in c,d). Fractions estimated over a) 4205 real masks of size 512×512 and 20719 generated masks, b) 1183 real masks of size 2048×2048 and 22705 generated masks, c) 4205 real masks of size 512×512 and 22604 generated masks, d) 1183 real masks of size 2048×2048 and 22560 generated masks.

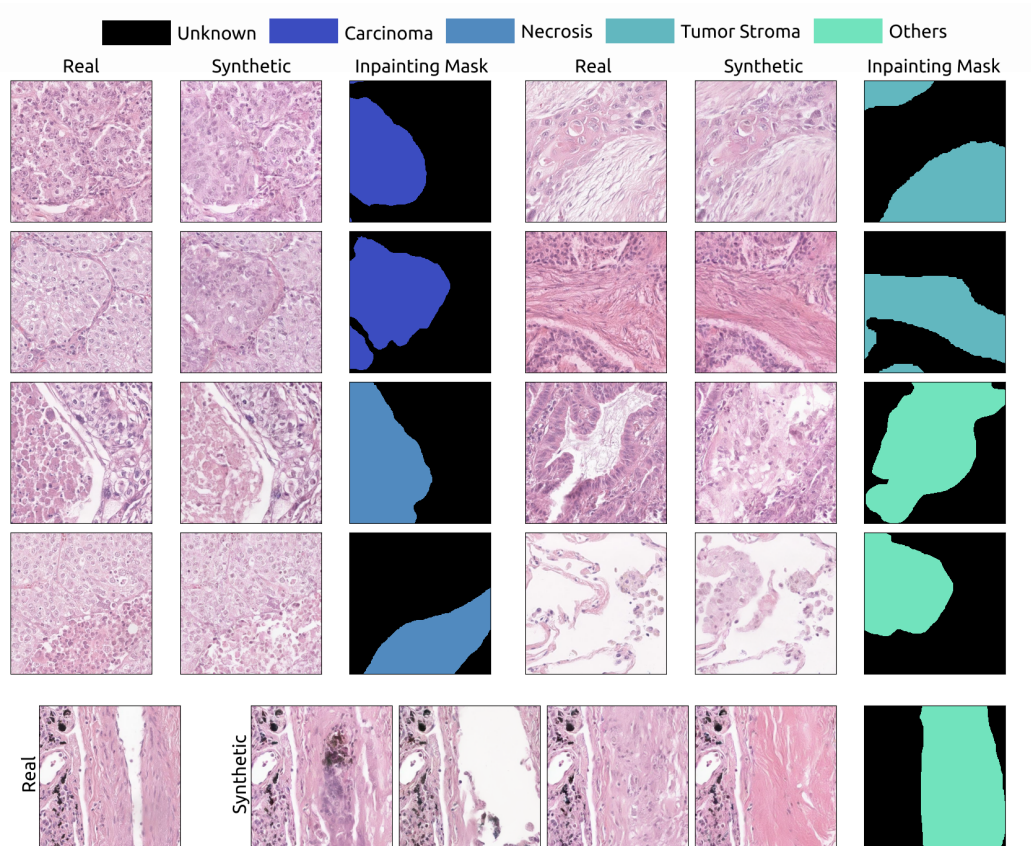


Figure 9: Inpainting test data with the corresponding label. Top) Inpainting for different labels. Bottom) Different inpainted synthetic areas for the same mask.

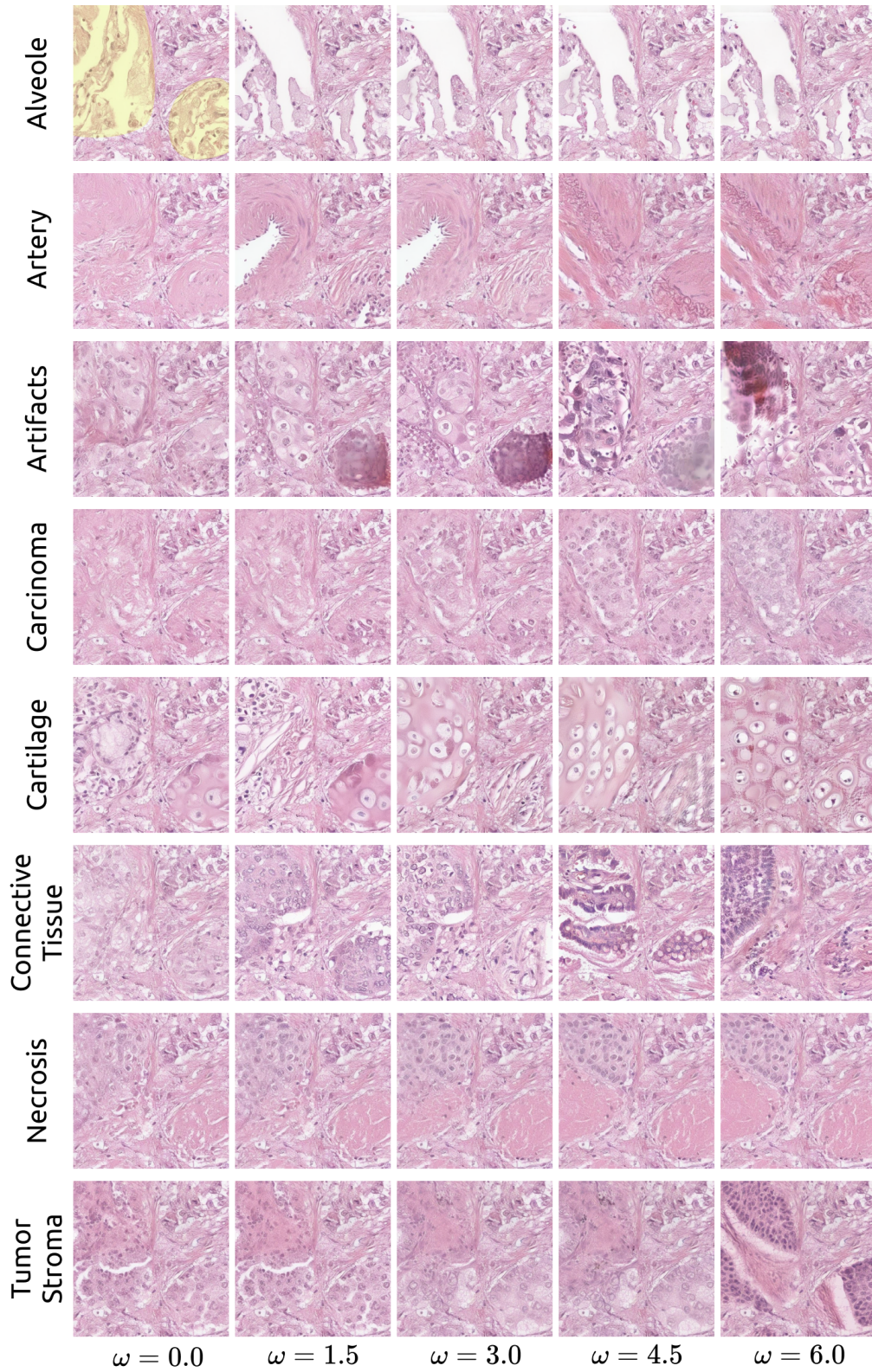


Figure 10: Proof of concept with inpainting. We inpainted the same base image with different classes and different strengths of conditioning (small ω corresponding to less diversity). The corresponding inpainting mask is displayed as an overlay on the top left patch (in yellow).

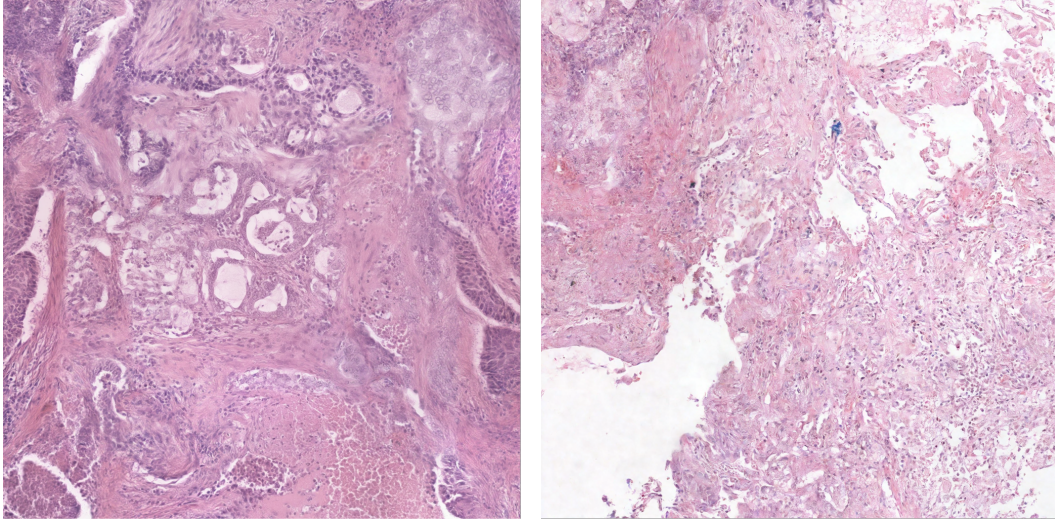


Figure 11: Comparison of different methods to generate large images (2048×2048). Left) DiffCollage image generation using the grid graph [48]. Right) DiffInfinite (ours) image generation using the proposed random patch sampling.

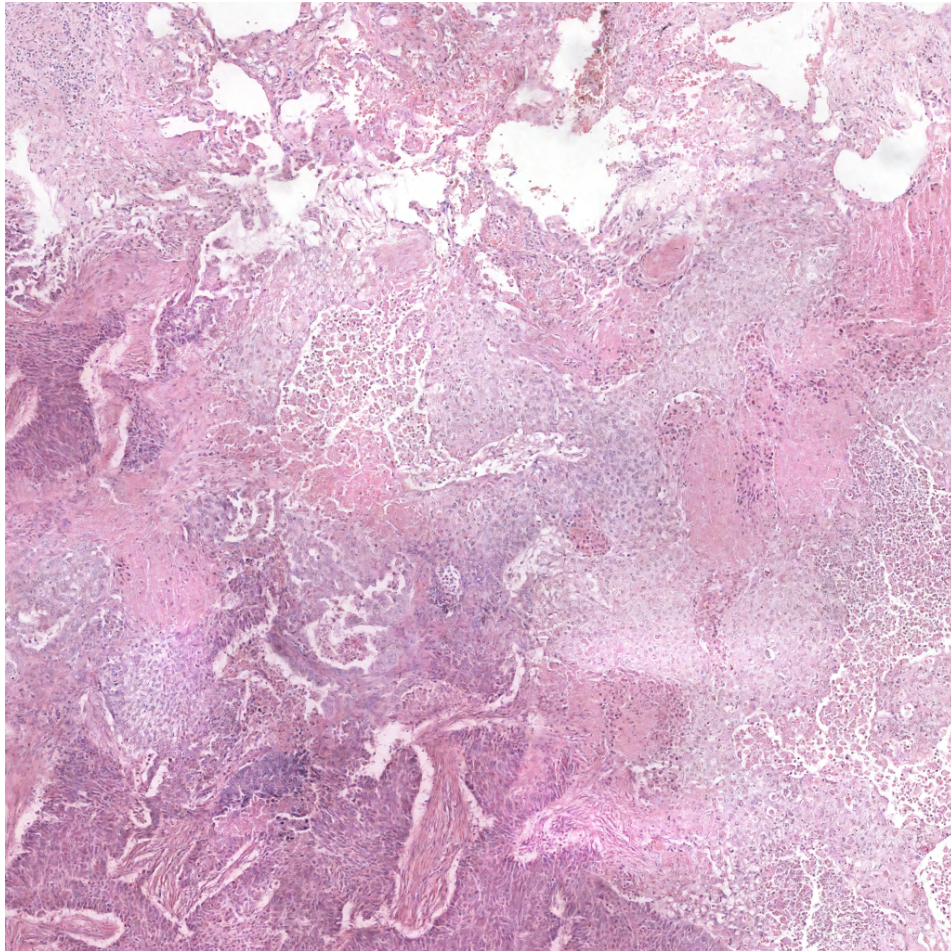


Figure 12: Large-content synthetic image with a size of 4096×4069 pixels.

E Training Details

Training on images The core model used in the diffusion process is a U-Net⁶. Every U-net’s block is composed of two ResNet blocks, a cross-attention layer and a normalization layer. On each ResNet block, we feed the output x^l of the previous block l , the time t and the label c_i . The cross-attention mask is performed using the mask corresponding to the label c_i as query and the input x^l as key and value.

Training on masks We replace the cross-attention with a linear self-attention layer for mask generation. Here, the model is conditioned with binary labels $\{0, 1\}$, where 0 corresponds to adenocarcinoma and 1 corresponds to squamous cell carcinoma. The masks of size 512×512 is first downsampled to size $1 \times 128 \times 128$. We stack the downsampled mask to the size $(3, 128, 128)$ to make it compatible with a pre-trained VAE⁷. We repeated the same training for the larger masks 2048×2048 , downsampling them to 128×128 as well.

Table 6: Details of the parameters used for training

	Model parameters image generation		Model parameters mask generation
Image X shape	(3,512,512)	Mask M shape	(3,128,128)
Latent Y shape	(4,64,64)	Latent Y shape	(4,16,16)
VAE	stabilityai/stable-diffusion-2-base	VAE (repo id)	stabilityai/stable-diffusion-2-base
Num classes	5 and 10	Num classes	2
Loss	L2	Loss	L2
Diffusion steps	1000	Diffusion steps	1000
Training steps	250000	Training steps	100000
Sampling steps	250	Sampling steps	250
Heads	4	Heads	4
Heads channels	32	Heads channels	32
Attention resolution	32,16,8	Attention resolution	32,16,8
Num Resblocks	2	Num Resblocks	2
Probability p_{unc}	0.5	Probability p_{unc}	0.5
Batch size	128	Batch size	64
Number of workers	32	Number of workers	1
GPUs Training	4 NVIDIA GeForce RTX 3090, 24Gb each	GPUs Training	2 Ampere A100, 40Gb each
GPUs Inference	1 NVIDIA GeForce RTX 3090	GPUs Inference	1 NVIDIA GeForce RTX 3090
Training time	~ 1 week	Training time	~ 4 hours
Optimizer	Adam	Optimizer	Adam
Scheduler	OneCycleLR(max lr=1e-4)	Scheduler	OneCycleLR(max lr=1e-4)

⁶Baseline, <https://github.com/lucidrains/classifier-free-guidance-pytorch>

⁷<https://huggingface.co/stabilityai/stable-diffusion-2>

F Sampling Details

Mask cleaning The diffusion model samples a latent mask in the VAE’s latent space. After mapping the latent mask back to the pixel space we average over the channels to have a mask with one channel and round the pixel values to the integers $\{0, 1, \dots, num_values\}$. Since we note some boundary artifacts between regions of different values we first apply a method from skimage⁸ to find these boundary artifacts and replace it by 0, corresponding to unknown area. Before resizing the mask to the full size, we apply a minpooling operation to erase labelled regions of small magnitude and replace it as well with unknowns.

Hann windows decoding After the diffusion model samples Z in the VAE’s latent space, the latent variable Z needs to be decoded into the pixel space. However, due to computational constraints, it is not feasible to decode Z all at once. Therefore, we tile it into smaller patches. Decoding smaller patches would introduce tiling effects. In order to reduce edge artifacts, we used an overlapping window method using Hann windows as weights [82]. In Fig. 13, we tile the image in four different configurations such that the edges and corners are overlapping, and then we perform a weighted sum over the upsampled outputs.

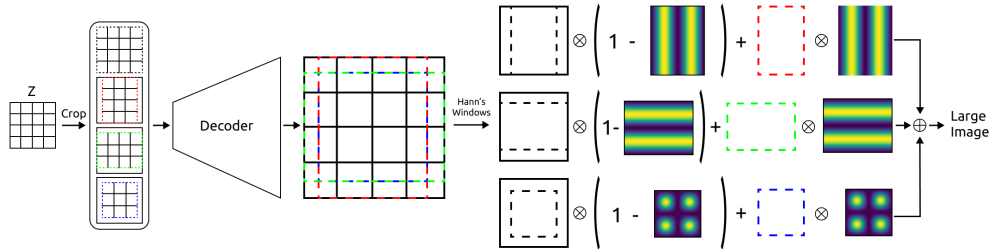


Figure 13: Hann window overlapping illustration.

⁸<https://scikit-image.org/docs/stable/api/skimage.segmentation.html>