

SegDiff: Image Segmentation with Diffusion Probabilistic Models

Tomer Amit¹, Tal Shaharbany¹, Eliya Nachmani^{1,2}, and Lior Wolf¹

¹Tel-Aviv University

²Facebook AI Research

{[tomeramit1](mailto:tomeramit1@tau.ac.il), [shaharabany](mailto:shaharabany@tau.ac.il), [eliyan](mailto:eliyan@tau.ac.il), [wolf](mailto:wolf@tau.ac.il)}@mail.tau.ac.il

Abstract

Diffusion Probabilistic Methods are employed for state-of-the-art image generation. In this work, we present a method for extending such models for performing image segmentation. The method learns end-to-end, without relying on a pre-trained backbone. The information in the input image and in the current estimation of the segmentation map is merged by summing the output of two encoders. Additional encoding layers and a decoder are then used to iteratively refine the segmentation map, using a diffusion model. Since the diffusion model is probabilistic, it is applied multiple times, and the results are merged into a final segmentation map. The new method produces state-of-the-art results on the Cityscapes validation set, the Vaihin-gen building segmentation benchmark, and the MoNuSeg dataset.

1. Introduction

Diffusion methods, which iteratively improve a given image, obtain image quality that is on par with or better than other types of generative models, including other forms of log-likelihood models and adversarial models [10, 19]. Such methods have been shown to excel in many generation tasks, both conditional and unconditional.

The vast majority of diffusion models are applied in domains in which there is no absolute ground truth result and the output is evaluated either through a user study or using several quality and diversity scores. As far as we know, with the exception of super resolution [19, 27, 41], diffusion models have not been applied to problems in which the ground truth result is unique.

In this work, we tackle the problem of image segmentation. This problem is a cornerstone of both classical computer vision and the deep learning methods of the last decade. The leading methods in the field employ

encoder-decoder networks of varied architectures [4, 31, 38, 50, 52, 53]. While adversarial methods have been attempted [12, 33, 49, 51], they do not constitute the current state of the art.

Therefore, it is uncertain whether diffusion models, which have been used primarily for GAN-like generation tasks, would be competitive in this domain. In this work, we propose applying a diffusion model to learn the image segmentation map. Unlike other recent improvements in the field of image segmentation [13, 22, 44], we train our method end-to-end, without relying on a pre-trained backbone network.

The diffusion model employs a denoising network conditioned on the input image only through a sum in which this information is aggregated with information arising from the current estimate x_t . Specifically, the input image I and the current estimate x_t of the binary segmentation map are passed through two different encoders, and the sum of these multi-channel tensors is passed through a U-Net [38] to provide the next estimate x_{t-1} .

Since the generation process is stochastic in its nature, one may obtain multiple solutions. As we show, merging these solutions, by simply averaging multiple runs, leads to an improvement in overall accuracy.

The novel method presented produces state-of-the-art results on multiple benchmarks: Cityscapes [9], building segmentation [39], and nuclei segmentation [25, 26].

Our main contributions are:

- We are the first to apply diffusion models to the image segmentation problem.
- We propose a new way to condition the model on the input image.
- We introduce the concept of multiple generations, in order to improve performance and calibration of the diffusion model.

- We obtained state-of-the-art results on multiple benchmarks. The margin is especially large for small data sets.

2. Related work

Image segmentation is a problem of assigning each pixel a label that identifies whether it belongs to a specific class or not. This problem is widely investigated using different architectures. These include fully convolutional networks [31], encoder-decoder architectures with skip-connections, such as U-Net [38], transformer-based architectures, such as the segformer [50], and even architectures that combine hypernetworks, such as [36].

Diffusion Probabilistic Models (DPM) [43] are a class of generative models based on a Markov chain, which can transform a simple distribution (e.g. Gaussian) to data that is sampled in a complex distribution. Diffusion models are capable of generating high-quality images that can compete with and even outperform the latest GAN methods [10, 18, 35, 43]. A variational framework for the likelihood estimation of diffusion models was introduced by Huang et al. [21]. Subsequently, Kingma et al. [23] proposed a Variational Diffusion Model that produces state-of-the-art results in likelihood estimation for image density.

Diffusion models were also applied to language modeling [2, 20], where a novel diffusion model for categorical data was used.

Conditional Diffusion Probabilistic Models In our work, we use diffusion models to solve the image segmentation problem as conditional generation, given the image. Conditional generation with diffusion models includes methods for class-conditioned generation, which is obtained by adding a class embedding to the timestamp embedding [35]. In [8] a method for guiding the generative process in DDPM is present. This method allows the generation of images based on a given reference image without any additional learning.

In the domain of super resolution, the lower-resolution image is upsampled and then concatenated, channelwise, to the generated image at each iteration [19, 41]. A similar approach passes the low-resolution images through a convolutional block [27] prior to the concatenation. Concurrently with our work, diffusion models were applied to image-to-image translation tasks [40]. These tasks include uncropping, inpainting, and colorization. The results obtained outperform strong GAN baselines.

Conditional diffusion models have also been used for voice generation. The mel-spectrogram is processed with a convolutional network, and is used as an additional input to the DPM denoising network ϵ [6, 24, 30]. Furthermore, in [37] a text-to-speech diffusion model is introduced, which uses text as a condition to the diffusion model.

In our work, we take a different approach to conditioning, adding (not concatenating) the input image, after it passes through an convolutional encoder, to the current estimation of the segmentation image. In other words, we learn the DPM of a residual model.

3. Background

We briefly introduce the formulation of diffusion models mentioned in [18]. Diffusion models are generative models parametrized by a Markov chain and composed of forward and backward processes. The forward process q is described by the formulation:

$$q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1}), \quad (1)$$

where T is the number of steps in the diffusion model, x_1, \dots, x_T are latent variables, and x_0 is a sample from the data. At each iteration of the forward process, Gaussian noise is added according to

$$q(x_t|x_{t-1}) = N(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t I_{n \times n}), \quad (2)$$

where β_t is a constant that defines the schedule of added noise, and $I_{n \times n}$ is the identity matrix of size n . As described in [18],

$$\alpha_t = 1 - \beta_t, \bar{\alpha}_t = \prod_{s=0}^t \alpha_s. \quad (3)$$

The forward process supports sampling at an arbitrary timestamp t , with the formula

$$q(x_t|x_0) = N(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I_{n \times n}), \quad (4)$$

which can be reparametrized to:

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{(1 - \bar{\alpha}_t)}\epsilon, \epsilon \sim N(0, I_{n \times n}). \quad (5)$$

The reverse process is parametrized by θ and defined by

$$p_\theta(x_{0:T-1}|x_T) = \prod_{t=1}^T p_\theta(x_{t-1}|x_t). \quad (6)$$

Starting from $p_\theta(x_T) = N(x_T; 0, I_{n \times n})$, the reverse process transforms the latent variable distribution $p_\theta(x_T)$ to the data distribution $p_\theta(x_0)$. The reverse process steps are performed by taking small Gaussian steps described by

$$p_\theta(x_{t-1}|x_t) = N(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)). \quad (7)$$

Calculating $q(x_{t-1}|x_t, x_0)$ using Bayes' theorem, one obtains:

$$q(x_{t-1}|x_t, x_0) = N(x_{t-1}; \tilde{\mu}(x_t, x_0), \tilde{\beta}_t I_{n \times n}), \quad (8)$$

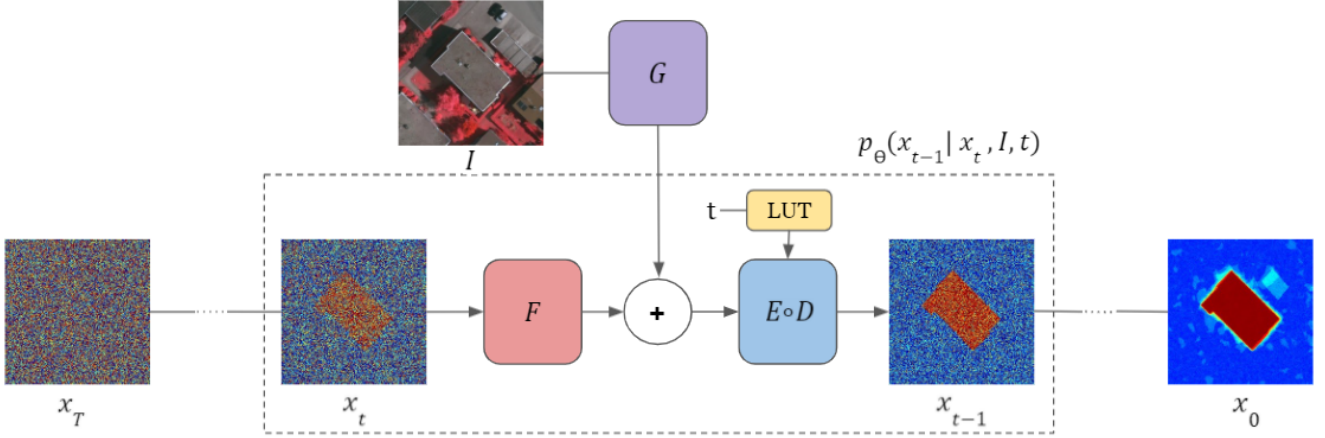


Figure 1. Our proposed diffusion method for image segmentation encodes the input signal, x_t , with F . The extracted features are summed with the feature map of the conditioned image I generated by network G . Networks E and D are a U-net encoder and decoder [35, 38], respectively, that refine the estimated segmentation map, obtaining x_{t-1} .

where

$$\tilde{\mu}_t(x_t, x_0) = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}x_0 + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}x_t, \quad (9)$$

$$\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\beta_t. \quad (10)$$

The neural network μ_θ predicts the noise ϵ , which is parametrized using Eq. 5, 9 to obtain:

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_\theta(x_t, t)\right). \quad (11)$$

Following [18] we set

$$\Sigma_\theta(x_t, t) = \sigma_t^2 I_{n \times n}, \quad (12)$$

where

$$\sigma_t^2 = \tilde{\beta}_t. \quad (13)$$

The forward process variance parameter is chosen to be a linearly increasing constant from $\beta_1 = 10^{-4}$ to $\beta_T = 2 * 10^{-2}$, formally:

$$\beta_t = \frac{10^{-4}(T - t) + 2 * 10^{-2}(t - 1)}{T - 1}. \quad (14)$$

Finally, we will minimize the term

$$E_{x_0, \epsilon, t} [|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)|^2], \quad (15)$$

where $\epsilon \sim N(0, I_{n \times n})$.

For inference, we can reparametrize the reverse process, Eq. 7, with Eq. 11, obtaining

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_\theta(x_t, t)\right) + \sigma_\theta z. \quad (16)$$

Algorithm 1 Inference Algorithm

Input total diffusion steps T , image I

$x_T \sim N(\mathbf{0}, \mathbf{I}_{n \times n})$

for $t = T, T - 1, \dots, 1$ **do**

$z \sim N(\mathbf{0}, \mathbf{I}_{n \times n})$

$\beta_t = \frac{10^{-4}(T-t) + 2 * 10^{-2}(t-1)}{T-1}$

$\alpha_t = 1 - \beta_t$

$\bar{\alpha}_t = \prod_{s=0}^t \alpha_s$

$\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\beta_t$

$x_{t-1} = \alpha_t^{-\frac{1}{2}}\left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_\theta(x_t, I, t)\right) + \mathbb{1}_{[t > 1]}\tilde{\beta}_t^{\frac{1}{2}}z$

return x_0

Algorithm 2 Training Algorithm

Input total diffusion steps T , images and segmentation masks dataset $D = \{(I_k, M_k)\}_k^K = 1$

repeat

Sample $(I_i, M_i) \sim D$, $\epsilon \sim N(\mathbf{0}, \mathbf{I}_{n \times n})$

Sample $t \sim \text{Uniform}(\{1, \dots, T\})$

$\beta_t = \frac{10^{-4}(T-t) + 2 * 10^{-2}(t-1)}{T-1}$

$\alpha_t = 1 - \beta_t$

$\bar{\alpha}_t = \prod_{s=0}^t \alpha_s$

Take gradient step on $\nabla_\theta ||\epsilon - \epsilon_\theta(x_t, I_i, t)||, x_t = \sqrt{\bar{\alpha}_t}M_i + \sqrt{1 - \bar{\alpha}_t}\epsilon$

until convergence

4. Method

Our method modifies the diffusion model by conditioning the step estimation function ϵ_θ on an input tensor that

combines information derived from both the current estimate x_t and the input image I .

In diffusion models, ϵ_θ is typically a U-Net [38]. In our work, ϵ_θ can be expressed in the following form:

$$\epsilon_\theta(x_t, I, t) = D(E(F(x_t) + G(I), t), t). \quad (17)$$

In this architecture, the U-Net’s decoder D is conventional and its encoder is broken down into three networks: E , F , and G . The last encodes the input image, while F encodes the segmentation map of the current step x_t . The two processed inputs have the same spatial dimensionality and number of channels. Based on the success of residual connections [17], we sum these signals $F(x_t) + G(I)$. This sum then passes to the rest of the U-Net encoder E .

The current step index t is passed to two different networks D and E . In each of these, it is embedded using a shared learned look-up table.

The output of ϵ_θ from Eq. 17, which is conditioned on I , is plugged into Eq. 16, replacing the unconditioned ϵ_θ network. This resulting inference time procedure is illustrated in Fig. 1 and detailed in Alg. 1.

4.1. Employing multiple generations

Since calculating x_{t-1} during inference includes the addition of $\sigma_\theta(x_t, t)z$, where z is from a standard distribution, there is significant variability between different runs of the inference method on the same inputs, see Fig. 2(b).

In order to exploit this phenomenon, we run the inference algorithm multiple times, then average the results.

This way, we stabilize the results of segmentation and improve performance, as demonstrated in Fig. 2(c). We use thirty generated instances in all experiments, except for the experiments in the ablation study, which quantifies the gain of this averaging procedure.

4.2. Training

The training procedure is depicted in Alg. 2. The total number of diffusion steps T is set by the user. For each iteration, a random sample is obtained (I_i, M_i) (an image and the associated ground truth binary segmentation map). The iteration number $1 \leq t \leq T$ is sampled from a uniform distribution, and epsilon from a standard distribution.

We then sample x_t according to Eq. 5, compute $F(x_t) + G(I_i)$, and apply networks E and D to obtain $\epsilon_\theta(x_t, I_i, t)$.

The loss being minimized is a modified version of Eq 15, namely:

$$E_{x_0, \epsilon, x_e, t} [|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, I_i, t)|^2]. \quad (18)$$

At training time, the ground truth segmentation of the input image I_i is known, and the loss is computed by setting $x_0 = M_i$.

4.3. Architecture

The input image encoder G is built from Residual in Residual Dense Blocks [47] (RRDBs), which combine multi-level residual connections without batch normalization layers. G has an input 2D-convolutional layer, an RRDB with a residual connection around it, followed by another 2D-convolutional layer, leaky RELU activation and a final 2D-convolutional output layer. F is a 2D-convolutional layer with a single-channel input and an output of C channels.

The encoder-decoder part of ϵ_θ , i.e., D and E , is based on U-Net, similarly to [35]. Each level is composed of residual blocks, and at resolution 16x16 and 8x8 each residual block is followed by an attention layer. The bottleneck contains two residual blocks with an attention layer in between. Each attention layer contains multiple attention heads.

The residual block is composed of two convolutional blocks, where each convolutional block contains group-norm, Silu activation, and a 2D-convolutional layer. The residual block receives the time embedding through a linear layer, Silu activation, and another linear layer. The result is then added to the output of the first 2D-convolutional block. Additionally, the residual block has a residual connection that passes all its content.

On the encoder side (network E), there is a downsample block after the residual blocks of the same depth, which is a 2D-convolutional layer with a stride of two. On the decoder side (network D), there is an upsample block after the residual blocks of the same depth, which is composed of the nearest interpolation that doubles the spatial size, followed by a 2D-convolutional layer. Each layer in the encoder has a skip connection to the decoder side.

5. Experiments

We present segmentation results for three datasets, as well as an ablation study.

Datasets The Cityscapes dataset [9] is an instance segmentation dataset containing 5,000 annotated images divided into 2,975 images for training, 500 for validation, and 1,525 for testing.

The experimental setting used is sometimes referred to as interactive segmentation and is motivated by the need to accelerate object annotation [1]. Under this setting, there are eight object categories, and the goal is to recover the objects’ per-pixel masks, given a cropped patch that contains the bounding box around each object.

Our per-object training and validation sets are created by taking crops from images in the original Cityscapes sets using the locations of the ground truth classes (we do not have access to the ground truth labels of the original Cityscapes test set).

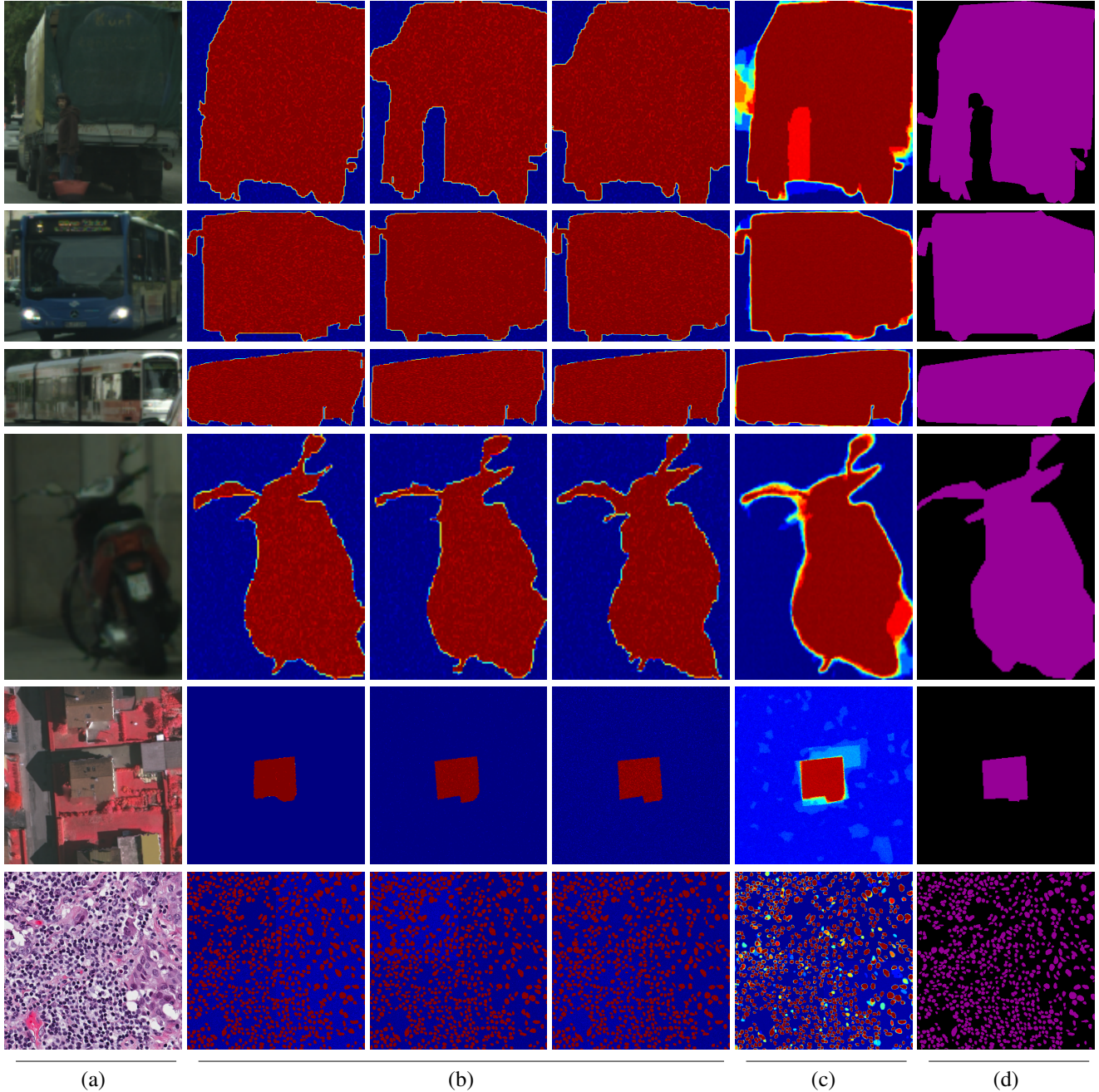


Figure 2. Obtaining multiple segmentation results for Cityscapes, Vaihingen, and MoNuSeg. (a) input image, (b) a subset of the obtained results for multiple runs on the same input, visualized by the jet color scale between 0 in blue and 1 in red, (c) average result, and (d) ground truth.

We compared our method for the Cityscapes dataset with PSPDeepLab [5], Polygon-RNN++ [1], CurveGCN [29] Deep active contours [14], Segformer-B5 [50] and Stdcl [11]. For most baselines, we report the results obtained from previous publications. For Segformer and Stdcl, we train from scratch.

We did not perform a comparison with PolyTrans-

form [28], since it uses a different protocol. Specifically, this method, which improves upon Mask R-CNN [16], utilizes the entire image (and not just the segmentation patch) as part of its inputs, and does not work on standard patches in a way that would enable a direct comparison.

The Vaihingen dataset [39] contains 168 aerial images of Vaihingen, in Germany, divided into 100 images for train-

ing and 68 for the test. The task is to segment the central building in each image. For this dataset, the leading baselines are DSAC [34], DarNet [7], TDAC [15], Deep active contours [14], FCN-UNET [38], FCN-ResNet-34, FCN-HarDNet-85 [4], Segformer-B5 [50] and Stdcl [11].

The MoNuSeg dataset [25, 26] contains a training set with 30 microscopic images from seven organs, with annotations of 21,623 individual nuclei. The test dataset contains 14 similar images. We resized the images to a resolution of 512×512 , following [45]. The relevant baseline methods are FCN [3], UNET [38], UNET++ [53], ResUnet [48], Axial attention (A.A) Unet [46] and Medical transformer [45].

Evaluation The Cityscapes dataset is evaluated using the common metrics of mean Intersection-over-Union (mIoU) per class.

$$mIoU(y_i, \hat{y}_i) = \sum_{i=1}^N \frac{TP(y_i, \hat{y}_i)}{TP(y_i, \hat{y}_i) + FN(y_i, \hat{y}_i) + FP(y_i, \hat{y}_i)} \quad (19)$$

Where N is the number of classes in the dataset, TP is the true positive between the ground truth y and output mask \hat{y} , FN is a false negative, and FP is a false positive.

The Vaihingen dataset is evaluated using several metrics: mIoU, F1-score, Weighted Coverage (WCov), and Boundary F-score (BoundF), as described in [7]. Briefly, the prediction is correct if it is within a certain distance threshold from the ground truth. The benchmarks use five thresholds, from 1px to 5px, for evaluating performance.

Following previous work, evaluation on the MoNuSeg dataset is performed using mIoU and the F1-score.

Training details The number of diffusion steps in previous works was 1000 [18] and even 4000 [35]. The literature suggests that more is better [42]. In our main experiments, we employ 100 diffusion steps to reduce inference time. An additional set of experiments investigated the influence of the number of diffusion steps on the performance and runtime of the method.

The AdamW [32] optimizer is used in all our experiments. Based on the intuition that the more RRDB blocks, the better the results, we used as many blocks as we could fit on the GPU without overly reducing batch size. The Unet used for datasets with a resolution of 256×256 has one additional layer with respect to the dataset with half that resolution, in order to account for the spatial dimensions.

On the Cityscapes dataset, the input resolution of our model is 128×128 . The test metrics are computed on the original resolution; therefore, we resized the prediction to the original image size.

Training took place with a batch size of 30 images. The network had 15 RRDB blocks and a depth of six. The number of channels was set to $[C, C, 2C, 2C, 4C, 4C]$ with $C = 128$. We followed the same augmentation scheme as in [14],

including random scaling in the range of $[0.75, 1.25]$, with up to 22 degrees rotation in each direction, and a horizontal flip with a probability of 0.5.

For the Vaihingen dataset, the size of the input image and the test image resolution was 256×256 . The experiments were performed with a batch size of eight images, six RRDB blocks, and a depth of seven. The number of channels was set to $[C, C, C, 2C, 2C, 4C, 4C]$ with $C = 128$.

The same augmentations are used as in [7]: random scaling by a factor sampled uniformly in the range $[0.75, 1.5]$, a rotation sampled uniformly between zero and 360 degrees, independent horizontal and vertical flips, applied with a probability of 0.5, and a random color jitter, with a maximum value of 0.6 brightness, 0.5 contrast, 0.4 saturation, and 0.025 hue.

For MoNuSeg, the input image resolution was 256×256 , but the test resolution was 512×512 . To address this, we applied a sliding window of 256×256 with a stride of 256, i.e., we tested each quadrant of the image separately.

The experiments were carried out with a batch size of eight images, with 12 RRDB blocks. The network depth was seven, and the number of channels in each depth was $[C, C, C, 2C, 2C, 4C, 4C]$, with $C = 128$. We used the same augmentation scheme as in [45] with random cropping of 256×256 to adjust for GPU memory.

It is worth noting that all baseline methods except Segformer and Stdcl rely on pre-trained weights obtained on the ImageNet, PASCAL or COCO datasets. Our networks are initialized with random weights.

Results Following previous work, Cityscapes is evaluated in one of two settings. *Tight*: in this setting, the samples (image and associated segmentation map) are extracted by a tight crop around the object mask. *Expansion*: samples are extracted by a crop around the object mask, which is 15% larger than the tight crop. The inputs of the model are crops 10% - 20% larger than the tight one. This setting is slightly more challenging, since there is less information on the location of the target object.

The results for the Cityscapes dataset are reported in Tab. 1. As can be seen, our method outperforms all baseline methods, across all categories and in both settings.

The gap is apparent even for the most recent baseline methods and, as can be seen in Fig. 3, the gap in performance is especially sizable for datasets with less training images.

The results for the Vaihingen dataset are presented in Tab. 2. As can be seen, our method outperforms the results reported in previous work for all four scores.

The results for the MoNuSeg dataset are presented in Tab. 3. In both segmentation metrics, our method outperforms all previous works, including very recent variants of U-Net and transformers that were developed specifically for this segmentation task.

	Method	Bicycle	Bus	Person	Train	Truck	M.cycle	Car	Rider	Mean
expansion	Polygon-RNN++ [1]	63.06	81.38	72.41	64.28	78.90	62.01	79.08	69.95	71.38
	PSP-DeepLab [5]	67.18	83.81	72.62	68.76	80.48	65.94	80.45	70.00	73.66
	Polygon-GCN [29]	66.55	85.01	72.94	60.99	79.78	63.87	81.09	71.00	72.66
	Spline-GCN [29]	67.36	85.43	73.72	64.40	80.22	64.86	81.88	71.73	73.70
	SegDiff (ours)	69.80	85.97	76.09	75.95	80.68	67.06	83.40	72.57	76.44
	Deep contour [14]	68.08	83.02	75.04	74.53	79.55	66.53	81.92	72.03	75.09
	Segformer-B5 [50]	68.02	78.78	73.53	68.46	74.54	64.06	83.20	69.12	72.46
	Std1 [11]	67.86	80.67	74.20	69.73	77.02	64.52	83.53	69.58	73.39
	Std2 [11]	68.67	81.29	74.41	71.36	75.71	63.69	83.51	69.90	73.57
	SegDiff (ours)	69.62	84.64	75.18	74.89	80.34	67.75	83.63	73.49	76.19

Table 1. Cityscapes segmentation results for two protocols: the top part refers to segmentation results with 15% expansion around the bounding box; the bottom part refers to segmentation results with a tight bounding box.

Method	F1-Score	mIoU	WCov	FBound
FCN-UNet [38]	87.40	78.60	81.80	40.20
FCN-ResNet34	91.76	87.20	88.55	75.12
FCN-HarDNet [4]	93.97	88.95	93.60	80.20
DSAC [34]	-	71.10	70.70	36.40
DarNet [7]	93.66	88.20	88.10	75.90
Deep contour [14]	94.80	90.33	93.72	78.72
TDAC [15]	94.26	89.16	90.54	78.12
Segformer-B5 [50]	93.94	88.57	91.91	77.95
Std1 [11]	94.04	88.75	92.78	78.86
Std2 [11]	93.97	88.62	92.59	77.3
SegDiff (ours)	95.14	91.12	93.83	85.09

Table 2. Segmentation results for the Vaihingen dataset.

Method	Dice	mIoU
FCN [3]	28.84	28.71
U-Net [38]	79.43	65.99
U-Net++ [53]	79.49	66.04
Res-UNet [48]	79.49	66.07
A.A U-Net [46]	76.83	62.49
MedT [45]	79.55	66.17
Ours	81.59	69.00

Table 3. Segmentation results for the MoNuSeg dataset.

The performance of the mIoU segmentation metric as a function of the number of iterations is presented for the three datasets in Fig. 4. It is interesting to note that the number of diffusion steps required to achieve the maximal score differs across the datasets.

All Cityscapes classes present similar behavior (with different levels of performance), saturating around the 60th iteration. The Vaihingen score takes longer to reach its maximal value. While one may attribute this to the larger input image size observed by the network, MoNuSeg, which has the same input image size as Vaihingen, reaches satura-

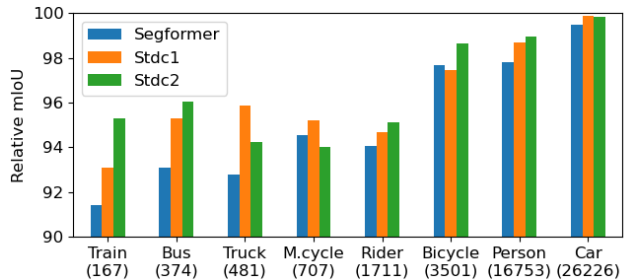


Figure 3. mIoU relative to SegDiff being 100% for each Cityscapes class, sorted by the number of training images per class.

	Method	F1-Score	IoU	WCov	FBound
Vaihingen	Variant one	91.60	85.45	88.67	71.70
	Variant two	90.92	84.00	89.67	70.00
	Variant three	93.77	88.67	91.69	80.15
	Variant four	94.77	90.27	93.82	82.64
	Variant five	93.16	87.76	91.08	79.89
	Variant six	91.97	85.57	89.83	71.04
	Full method	94.95	90.64	94.00	84.37
Cityscapes "Bus"	Variant one	90.52	84.15	90.37	62.66
	Variant two	85.21	75.92	81.15	38.81
	Variant three	90.35	83.76	88.56	58.80
	Variant four	91.30	85.17	90.34	63.85
	Variant five	89.57	82.73	88.87	58.40
	Variant six	82.97	72.66	80.85	34.38
	Full method	90.72	84.35	89.96	63.87

Table 4. Ablation study for different conditioning methods.

tion earlier, similarly to the Cityscapes classes. An alternative hypothesis can relate the number of required iterations to the ratio of pixels within the segmentation mask, which is higher for Cityscapes and MoNuSeg than for Vaihingen. This requires further validation.

We next study the effect of the number of generated in-

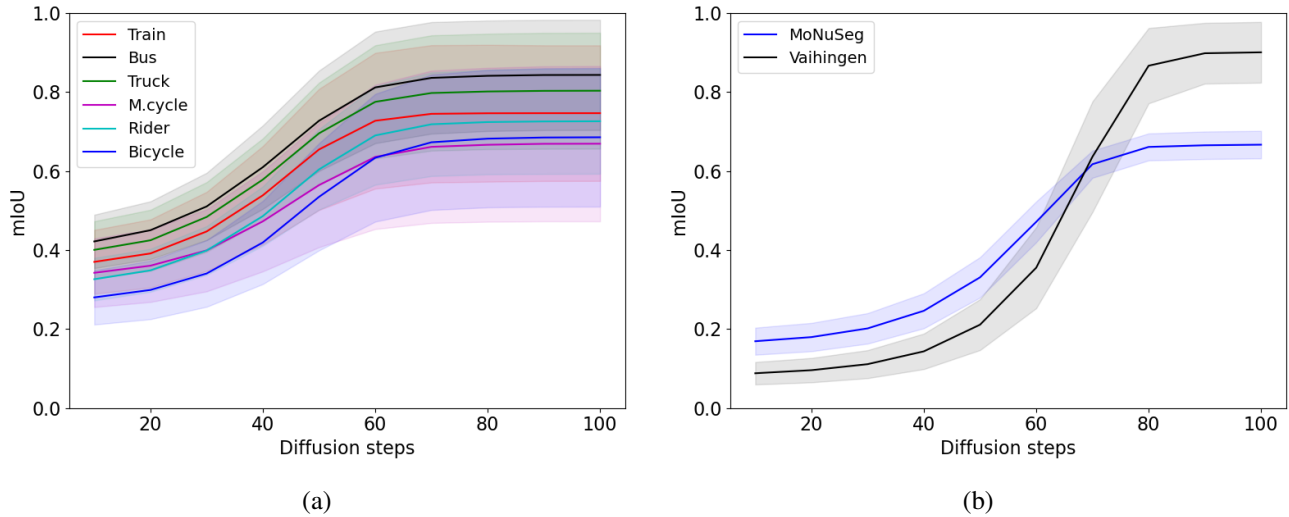


Figure 4. mIoU (mean and variance) across the test images as a function of the number of diffusion steps. (a) Results for the Cityscapes classes, with 128×128 image resolution. (b) Results for the Vaihingen and MoNuSeg datasets, with 256×256 image resolution.

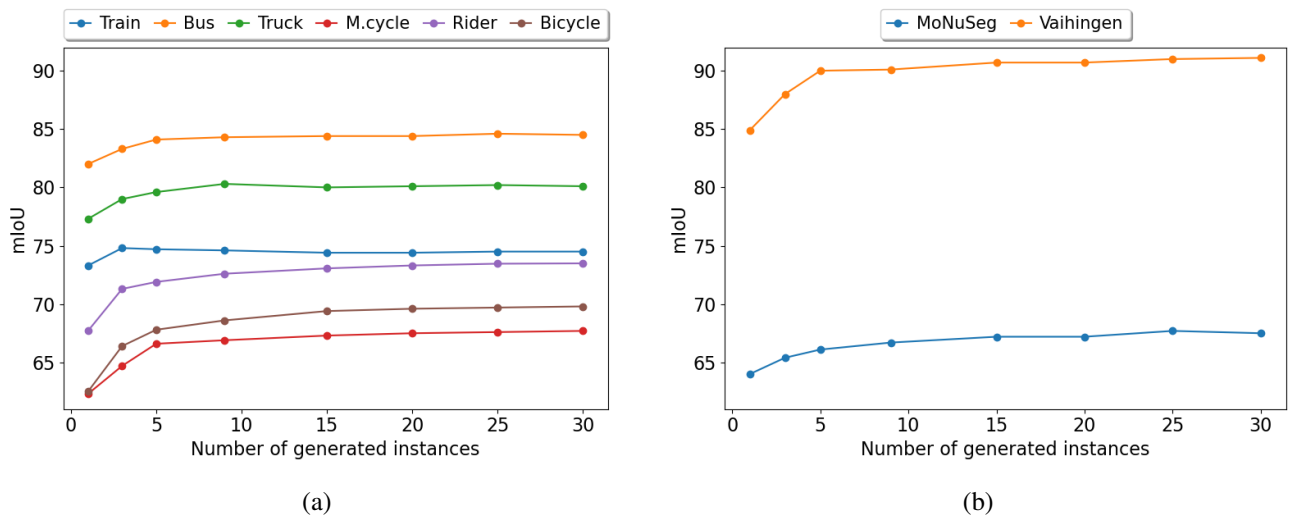


Figure 5. mIoU per number of generated inferences. (a) Results for the Cityscapes classes, with 128×128 image resolution. (b) Results for the Vaihingen and MoNuSeg datasets, with 256×256 image resolution.

stances on performance. The results can be seen in Fig. 5. In general, increasing the number of generated instances tends to increase the mIoU score. However, the number of runs required to reach optimal performance varies between classes. For example, for the “Bus” and “Train” classes of Cityscapes, the best score is achieved when using 10 and 3 generated instances, respectively. MoNuSeg, requires considerably more runs (25) for maximal performance. On the other hand, when the number of generated instances is increased, inference time also increases linearly, resulting in a slower method compared to architectures such as Segformer and Std.

Another aspect of achieving improvement by employing multiple generations is calibration. The calibration score is measured as the difference between the prediction probability and the true probability of the event. For example, a perfectly calibrated model is defined by $\mathbb{P}(\hat{Y} = Y | \hat{P} = p) = p$, which means that the prediction probability equals the true probability of the event. We estimate the calibration score by splitting the $[0, 1]$ range into ten uniform bins, then average the squared difference between each bin’s mean prediction probability and the percentage of positive samples.

The results of examining the calibration scores are pre-

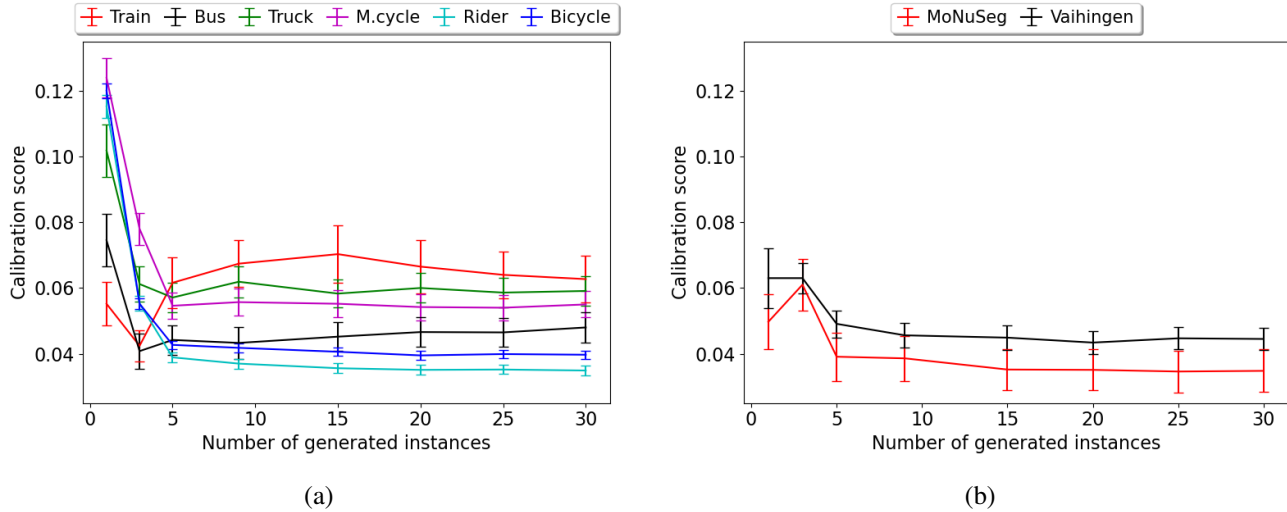


Figure 6. Mean calibration score (lower is better) per number of generated inferences. The error bars depict the standard error. (a) Results for the Cityscapes classes, with an image resolution of 128×128 . (b) Results for the Vaihingen and MoNuSeg datasets, with the 256×256 image resolution.

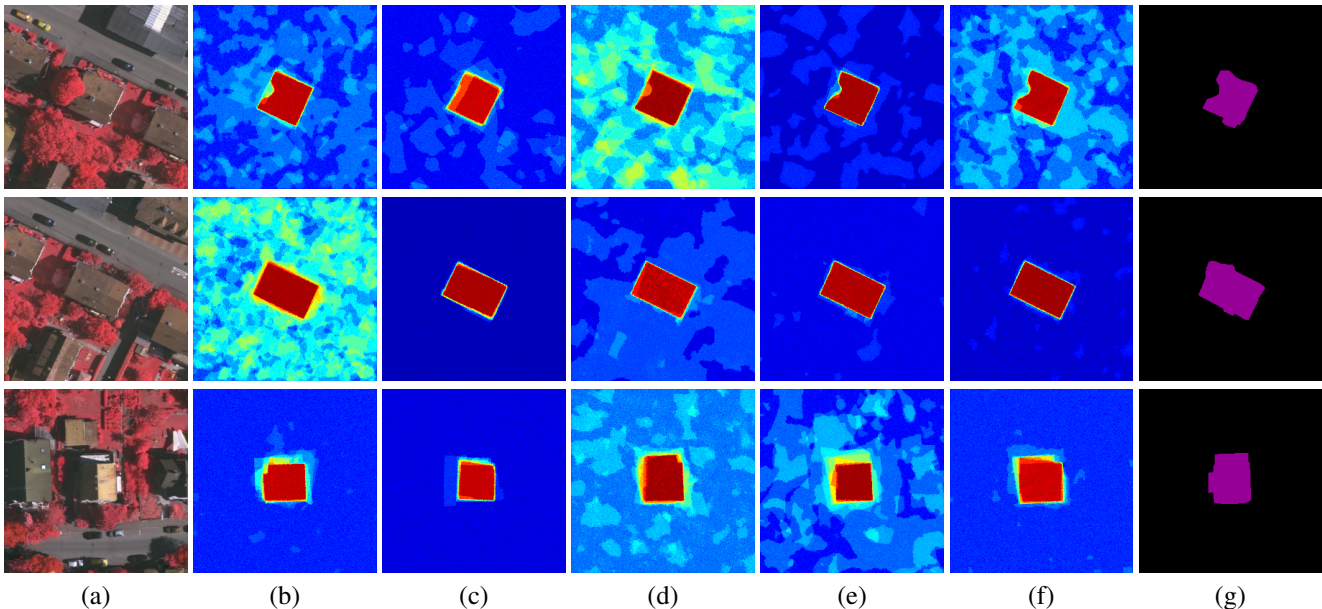
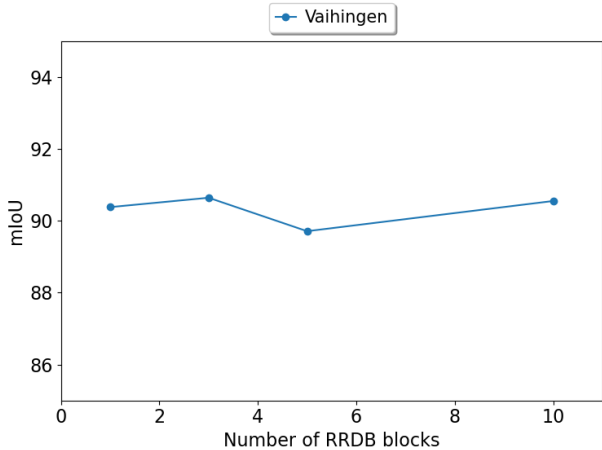


Figure 7. Results of the ablation study. (a) the input image, (b-e) results for variants one–four of our method, respectively, (f) the result of our method, and (g) ground truth. Panels (b-f) employ the jet color scale between 0 in blue and 1 in red

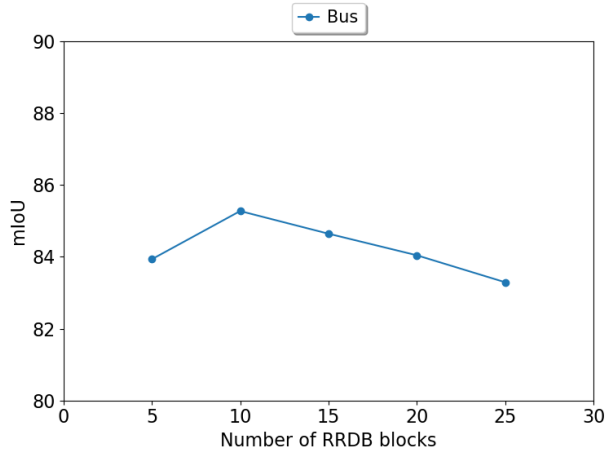
sented in Fig. 6. For most datasets, increasing the number of generated instances improves the calibration score, especially when the increase is from a single instance. In addition, for the larger classes in Cityscapes - Rider and Bicycle - and for the MoNuSeg and Vaihingen datasets, the improvement continues to increase even more compared to the other datasets. The “Train” class in Cityscapes is an exception; here, the single-instance calibration score is better

than other experiments with a larger number of generated instances. This phenomenon may be a result of the highly varied size and the small number of test images.

Ablation Study We evaluate various alternatives to our method. The first variant concatenates $[F(x_t), G(I)]$ at the channel dimension. The second variant employs FC-HarDNet-70 V2 [4] instead of RRDBs. The third variant, following [19, 41], concatenates I channelwise to x_t , with-

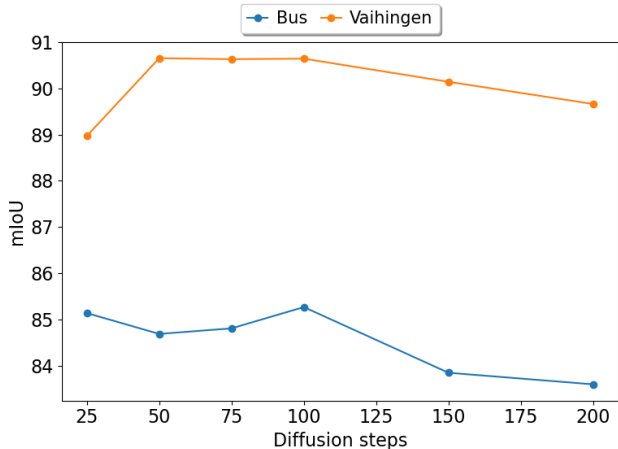


(a)

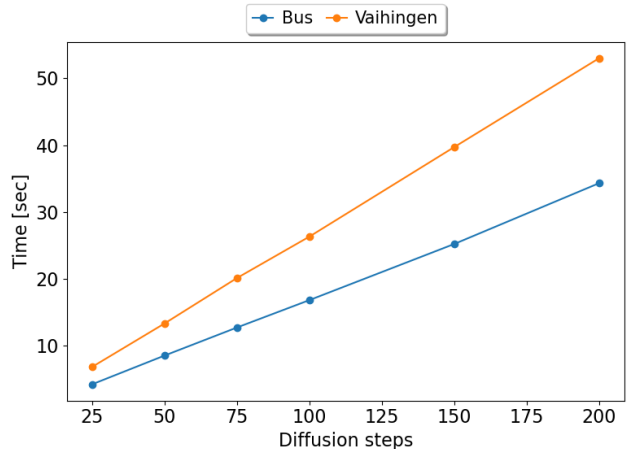


(b)

Figure 8. mIoU per number of RRDB blocks. (a) Results on Vaihingen, (b) Results on Cityscapes “Bus”.



(a)



(b)

Figure 9. Generation time in seconds and mIoU per number of diffusion steps for Vaihingen and Cityscapes “Bus”. (a) mIoU per diffusion step, (b) Time per diffusion step.

out using an encoder. The last alternative method is to propagate $F(x_t)$ through the U-Net module and add it to $G(I)$ after the first, third, and fifth downsample blocks (variants four–six), instead of performing $F(x_t) + G(I)$. In this variant, $G(I)$ is downsampled to match the required number of channels by propagating it through a 2D-convolutional layer with a stride of two.

These variant experiments were tested by averaging nine generated instances on the Vaihingen dataset and on Cityscapes “Bus” (the performance reported for our method is therefore slightly different from those reported in Tab. 2).

The summation we introduce as a conditioning approach outperforms concatenation (variant one) on Vaihingen by

a large margin, while on Cityscapes “Bus”, the difference is small. The RRDB blocks are preferable to the FC-HardNet architecture in both datasets (variant two). Removing the encoder affects the metrics significantly (variant three), slightly more so on Vaihingen. The change in the signal’s integration position of variant four leads to a negligible difference on Vaihingen and even outperforms our full method on Cityscapes “Bus”. Variants five and six lead to a decrease in performance as the distance from the first layer increases. Fig. 7 depicts sample results for the various variants (one–four) of the Vaihingen dataset.

Parameter sensitivity For testing the stability of our proposed method, we experimented with the two hyperparam-

eters that can affect performance the most: the number of diffusion steps, and the number of RRDB blocks. To study the effect of these parameters, we varied the number of diffusion steps in the range of [25, 50, 75, 100, 150, 200], and the number of RRDB blocks in the range of [1, 3, 5, 10] for Vaihingen and [5, 10, 15, 20, 25] for Cityscapes “Bus”. We started from a baseline configuration (which was 100 diffusion steps, 3 RRDB blocks for Vaihingen, and 10 RRDB blocks for Cityscapes “Bus”) and experimented with different values around these.

The effect of the number of RRDB blocks In this part we set the number of diffusion steps to 100. As can be seen in Fig. 8, with our configuration, the optimal number of RRDB blocks is 3 for Vaihingen, and 10 for Cityscapes “Bus”. However, evidently, the number of blocks has a limited impact in the case of both Cityscapes and Vaihingen. The gap between the best and worst performance points is less than 1 mIoU for Vaihingen and less than 2 mIoU for Cityscapes “Bus”. Therefore, we conclude that this hyperparameter has a small effect on performance.

Varying the number of diffusion steps T In this part, we set the number of RRDB blocks of Vaihingen to 3 and Cityscapes “Bus” to 10. We explore the possible accuracy/runtime tradeoff with regards to the number T of diffusion steps. Results are shown in Fig. 9.

When the number of diffusion steps is increased - as we can see in Fig. 9(a) - the graph fluctuation for Vaihingen is less than 1 mIoU, and for Cityscapes “Bus” it is less than 2 mIoU.

Surprisingly, when the number of diffusion steps is reduced, even to just 25, which is a very low number compared to the literature [18, 35], the segmentation results remain stable in both datasets, with a degradation of only up to 2 mIoU for Vaihingen, and 1 mIoU for Cityscapes “Bus”. This reduction can speed up performance by a factor of four and provide a reasonable accuracy to runtime tradeoff.

The results for the generation time of one sample in seconds are presented in Fig. 9(b). As can be observed, both graphs are linear, with a different slope. The main reason for this is the difference in image size (which is 256×256 for Vaihingen and 128×128 for Cityscapes “Bus”). Another, minor reason is the difference in the number of RRDB blocks in this experiment.

6. Conclusions

A wealth of methods have been applied to image segmentation, including active contour and their deep variants, encoder-decoder architectures, and U-Nets, which - together with more recent, transformer-based methods - represent a leading approach. In this work, we propose utilizing the state-of-the-art image generation technique of diffusion models. Our diffusion model employs a U-Net architecture, which is used to incrementally improve the obtained

generation, similarly to other recent diffusion models.

In order to condition the input image, we generate another encoding path, which is similar to U-Net’s encoder-decoder use in conventional image segmentation methods. The two encoder pathways are merged by summing the activations early in the U-Net’s encoder.

Using our approach, we obtain state-of-the-art segmentation results on a diverse set of benchmarks, including street view images, aerial images, and microscopy.

7. Acknowledgments

This project has received funding from the European Research Council (ERC) under the European Unions Horizon 2020 research and innovation programme (grant ERC CoG 725974).

References

- [1] David Acuna, Huan Ling, Amlan Kar, and Sanja Fidler. Efficient interactive annotation of segmentation datasets with polygon-rnn++. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 859–868, 2018. 4, 5, 7
- [2] Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg. Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems*, 34:17981–17993, 2021. 2
- [3] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 39(12):2481–2495, 2017. 6, 7
- [4] Ping Chao, Chao-Yang Kao, Yu-Shan Ruan, Chien-Hsiang Huang, and Youn-Long Lin. Hardnet: A low memory traffic network. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3552–3561, 2019. 1, 6, 7, 9
- [5] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017. 5, 7
- [6] Nanxin Chen, Yu Zhang, Heiga Zen, Ron J Weiss, Mohammad Norouzi, and William Chan. Wavegrad: Estimating gradients for waveform generation. *arXiv preprint arXiv:2009.00713*, 2020. 2
- [7] Dominic Cheng, Renjie Liao, Sanja Fidler, and Raquel Urtasun. Darnet: Deep active ray network for building segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7431–7439, 2019. 6, 7
- [8] Jooyoung Choi, Sungwon Kim, Yonghyun Jeong, Youngjune Gwon, and Sungroh Yoon. Ilvr: Conditioning method for denoising diffusion probabilistic models. *arXiv preprint arXiv:2108.02938*, 2021. 2

- [9] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. 1, 4
- [10] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34, 2021. 1, 2
- [11] Mingyuan Fan, Shenqi Lai, Junshi Huang, Xiaoming Wei, Zhenhua Chai, Junfeng Luo, and Xiaolin Wei. Rethinking bisenet for real-time semantic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9716–9725, 2021. 5, 6, 7
- [12] Volker Fischer, Mummadi Chaithanya Kumar, Jan Hendrik Metzen, and Thomas Brox. Adversarial examples for semantic image segmentation. *arXiv preprint arXiv:1703.01101*, 2017. 1
- [13] Jun Fu, Jing Liu, Jie Jiang, Yong Li, Yongjun Bao, and Hanqing Lu. Scene segmentation with dual relation-aware attention network. *IEEE Transactions on Neural Networks and Learning Systems*, 32(6):2547–2560, 2020. 1
- [14] Shir Gur, Tal Shaharabany, and Lior Wolf. End to end trainable active contours via differentiable rendering. *arXiv preprint arXiv:1912.00367*, 2019. 5, 6, 7
- [15] Ali Hatamizadeh, Debleena Sengupta, and Demetri Terzopoulos. End-to-end trainable deep active contour models for automated image segmentation: Delineating buildings in aerial imagery. In *European Conference on Computer Vision*, pages 730–746. Springer, 2020. 6, 7
- [16] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 5
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 4
- [18] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 2, 3, 6, 11
- [19] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *Journal of Machine Learning Research*, 23(47):1–33, 2022. 1, 2, 9
- [20] Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows and multinomial diffusion: Towards non-autoregressive language models. *arXiv e-prints*, pages arXiv–2102, 2021. 2
- [21] Chin-Wei Huang, Jae Hyun Lim, and Aaron C Courville. A variational perspective on diffusion-based generative models and score matching. *Advances in Neural Information Processing Systems*, 34, 2021. 2
- [22] Zilong Huang, Xinggang Wang, Lichao Huang, Chang Huang, Yunchao Wei, and Wenyu Liu. Ccnet: Criss-cross attention for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 603–612, 2019. 1
- [23] Diederik P Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *arXiv preprint arXiv:2107.00630*, 2021. 2
- [24] Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion model for audio synthesis. *arXiv preprint arXiv:2009.09761*, 2020. 2
- [25] Neeraj Kumar, Ruchika Verma, Deepak Anand, Yanning Zhou, Omer Fahri Onder, Efstratios Tsougenis, Hao Chen, Pheng-Ann Heng, Jiahui Li, Zhiqiang Hu, et al. A multi-organ nucleus segmentation challenge. *IEEE transactions on medical imaging*, 39(5):1380–1391, 2019. 1, 6
- [26] Neeraj Kumar, Ruchika Verma, Sanuj Sharma, Surabhi Bhargava, Abhishek Vahadane, and Amit Sethi. A dataset and a technique for generalized nuclear segmentation for computational pathology. *IEEE transactions on medical imaging*, 36(7):1550–1560, 2017. 1, 6
- [27] Haoying Li, Yifan Yang, Meng Chang, Shiqi Chen, Huajun Feng, Zhihai Xu, Qi Li, and Yueting Chen. Srdiff: Single image super-resolution with diffusion probabilistic models. *Neurocomputing*, 2022. 1, 2
- [28] Justin Liang, Namdar Homayounfar, Wei-Chiu Ma, Yuwen Xiong, Rui Hu, and Raquel Urtasun. Polytransform: Deep polygon transformer for instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9131–9140, 2020. 5
- [29] Huan Ling, Jun Gao, Amlan Kar, Wenzheng Chen, and Sanja Fidler. Fast interactive object annotation with curve-gcn. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5257–5266, 2019. 5, 7
- [30] Jinglin Liu, Chengxi Li, Yi Ren, Feiyang Chen, Peng Liu, and Zhou Zhao. Diffsinger: Singing voice synthesis via shallow diffusion mechanism. *arXiv preprint arXiv:2105.02446*, 2021. 2
- [31] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015. 1, 2
- [32] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 6
- [33] Pauline Luc, Camille Couprie, Soumith Chintala, and Jakob Verbeek. Semantic segmentation using adversarial networks. *arXiv preprint arXiv:1611.08408*, 2016. 1
- [34] Diego Marcos, Devis Tuia, Benjamin Kellenberger, Lisa Zhang, Min Bai, Renjie Liao, and Raquel Urtasun. Learning deep structured active contours end-to-end. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8877–8885, 2018. 6, 7
- [35] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021. 2, 3, 4, 6, 11
- [36] Yuval Nirkin, Lior Wolf, and Tal Hassner. Hyperseg: Patch-wise hypernetwork for real-time semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4061–4070, 2021. 2

- [37] Vadim Popov, Ivan Vovk, Vladimir Gogoryan, Tasnima Sadekova, and Mikhail Kudinov. Grad-tts: A diffusion probabilistic model for text-to-speech. In *International Conference on Machine Learning*, pages 8599–8608. PMLR, 2021. [1](#), [2](#)
- [38] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. [1](#), [2](#), [3](#), [4](#), [6](#), [7](#)
- [39] Franz Rottensteiner, Gunho Sohn, Markus Gerke, and Jan D Wegner. Isprs semantic labeling contest. *ISPRS: Leopoldshöhe, Germany*, 2014. [1](#), [5](#)
- [40] Chitwan Saharia, William Chan, Huiwen Chang, Chris A Lee, Jonathan Ho, Tim Salimans, David J Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. *arXiv preprint arXiv:2111.05826*, 2021. [2](#)
- [41] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *arXiv preprint arXiv:2104.07636*, 2021. [1](#), [2](#), [9](#)
- [42] Robin San-Roman, Eliya Nachmani, and Lior Wolf. Noise estimation for generative diffusion models. *arXiv preprint arXiv:2104.02600*, 2021. [6](#)
- [43] 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. [2](#)
- [44] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7262–7272, 2021. [1](#)
- [45] Jeya Maria Jose Valanarasu, Poojan Oza, Ilker Hacihaliloglu, and Vishal M Patel. Medical transformer: Gated axial-attention for medical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 36–46. Springer, 2021. [6](#), [7](#)
- [46] Huiyu Wang, Yukun Zhu, Bradley Green, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Axial-deeplab: Stand-alone axial-attention for panoptic segmentation. In *European Conference on Computer Vision*, pages 108–126. Springer, 2020. [6](#), [7](#)
- [47] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *Proceedings of the European conference on computer vision (ECCV) workshops*, pages 0–0, 2018. [4](#)
- [48] Xiao Xiao, Shen Lian, Zhiming Luo, and Shaozi Li. Weighted res-unet for high-quality retina vessel segmentation. In *2018 9th international conference on information technology in medicine and education (ITME)*, pages 327–331. IEEE, 2018. [6](#), [7](#)
- [49] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Yuyin Zhou, Lingxi Xie, and Alan Yuille. Adversarial examples for semantic segmentation and object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 1369–1378, 2017. [1](#)
- [50] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems*, 34, 2021. [1](#), [2](#), [5](#), [6](#), [7](#)
- [51] Yuan Xue, Tao Xu, Han Zhang, L Rodney Long, and Xiaolei Huang. Segan: adversarial network with multi-scale l1 loss for medical image segmentation. *Neuroinformatics*, 16(3):383–392, 2018. [1](#)
- [52] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015. [1](#)
- [53] Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. Unet++: A nested u-net architecture for medical image segmentation. In *Deep learning in medical image analysis and multimodal learning for clinical decision support*, pages 3–11. Springer, 2018. [1](#), [6](#), [7](#)