# IMPROVED HER2 TUMOR SEGMENTATION WITH SUBTYPE BALANCING USING DEEP GENERATIVE NETWORKS

*Mathias Öttl[1], Jana Mönius[2], Matthias Rübner[3], Carol I. Geppert[2], Jingna Qiu[4],*
*Frauke Wilm[1,4], Arndt Hartmann[2], Matthias W. Beckmann[3], Peter A. Fasching[3], Andreas Maier[1],*
*Ramona Erber[2], Katharina Breininger[4]*

[1] Pattern Recognition Lab, Friedrich-Alexander-Universität Erlangen-Nürnberg (FAU), Germany
[2] Institute of Pathology, University Hospital Erlangen, FAU, Germany
[3] Department of Gynecology and Obstetrics, University Hospital Erlangen, FAU, Germany
[4] Department Artificial Intelligence in Biomedical Engineering, FAU, Germany

## ABSTRACT

Tumor segmentation in histopathology images is often complicated by its composition of different histological subtypes and class imbalance. Oversampling subtypes with low prevalence features is not a satisfactory solution since it eventually leads to overfitting. We propose to create synthetic images with semantically-conditioned deep generative networks and to combine subtype-balanced synthetic images with the original dataset to achieve better segmentation performance. We show the suitability of Generative Adversarial Networks (GANs) and especially diffusion models to create realistic images based on subtype-conditioning for the use case of HER2-stained histopathology. Additionally, we show the capability of diffusion models to conditionally inpaint HER2 tumor areas with modified subtypes. Combining the original dataset with the same amount of diffusion-generated images increased the tumor Dice score from 0.833 to 0.854 and almost halved the variance between the HER2 subtype recalls. These results create the basis for more reliable automatic HER2 analysis with lower performance variance between individual HER2 subtypes.

*Index Terms*— Histopathology, HER2, Subtypes, Generative Models, Diffusion Models, Segmentation

## 1. INTRODUCTION

The extraction of distinct features to differentiate individual subtypes from one composite class can be problematic for machine learning algorithms, leading to a weakened performance for the main task and an inconsistent performance among subtypes [1, 2]. In histopathology, tumors can be subtyped by origin or reason for growth [3]. In this work, we focused on different histological subtypes among Human Epidermal growth factor Receptor 2 (HER2)-stained breast cancer samples, defined according to scoring systems for HER2. Each tumor cell in HER2-stained tissue can be scored as 0, 1+, 2+, or 3+, leading to the HER2 tumor class being a composition of these subtypes [4]. An aggregated HER2 score is usually assigned to each tumor sample, based on the composition of HER2-scored tumor tissue [5]. Treatment choices are based on the aggregated HER2 score; thus, correct results for the subtype composition are essential when automatic tumor segmentation is employed. In this work, we consider the first step in a HER2 segmentation pipeline, which is the segmentation of tumor tissue against background tissue. Different HER2 subtype characteristics combined with a different prevalence of these HER2 subtypes can lead to an inconsistent tumor segmentation performance between the underlying subtypes.

During algorithm development, this inconsistent performance across individual subtypes can be approached with oversampling of underrepresented subtypes, which can, however, quickly lead to overfitting due to the limited amount of training data available [6]. To avoid overfitting, the training set can be extended by synthetic images created by a generative model, specifically Generative Adversarial Networks (GANs) [7], which have been successfully applied for a wide range of medical applications. Chen *et al*. reviewed 105 publications in this domain, and for microscopic pathology, a performance increase was reported for all but one work [8].

Inspired by these results, we propose to use generative models to generate subtype-balanced synthetic tumor image datasets. Similar to the work of Fajardo *et al*. [9], we employ GANs for image generation, but extend it by semantic conditioning, i.e. a generative model is tasked to create an output that matches a two-dimensional label mask. Recently, diffusion models have shown great potential in image synthesis [10] and have previously outperformed GANs [11]. Therefore, we introduce diffusion models for semantic image synthesis in histopathology aiming to tackle subtype imbalances within our dataset. Furthermore, we experiment with partial image synthesis, where we use semantically-conditioned diffusion models to inpaint tumor regions, while the background remains unchanged.

We investigate how different amounts of synthetic images from these three generative methods (GAN-generated, diffusion model-generated, diffusion model-inpainted) affect the subsequent tumor segmentation performance and how well the individual subtypes are segmented.
The main contributions of this paper include the following:

- Illustration of the suitability of GANs and especially diffusion models to create realistic HER2 images using semantic subtype-conditioning. Demonstration of the suitability of diffusion models for semantic subtype-conditioned inpainting in HER2 histopathology images.
- Analysis of how different amounts of additional synthetic images influence the segmentation performance, specifically the overall tumor segmentation performance and the performance on individual subtypes.

## 2. DATASET AND METHODS

### 2.1. Dataset

The data used in this work originated from 40 breast tissue sections from 40 different patients. The tissue was histochemically stained for HER2 and digitized as Whole Slide Images (WSIs) with the 3DHistech PANNORAMIC 1000 scanner, using a $20\times$ objective lens. Despite methods to reduce the effort [12], annotation of WSI is still not reasonable; thus, ten regions-of-interest of size $1.5\,\text{mm}\times1.5\,\text{mm}$ were selected from each WSI. Cell groups of the same subtype were annotated as one tumor tissue instance using polygon contours on the EXACT platform [13]. Five subtypes of tissue were considered: the four HER2 subtypes according to the asco/cap guidelines [14] and, as a fifth subtype, Lobular Carcinoma In Situ (LCIS)/Ductal Carcinoma In Situ (DCIS) as one composite class. Figure 1 shows examples of the annotated tissue types. LCIS and DCIS describe non-invasive tumor tissues, which could be assigned to one of the four HER2 classes. This assignment was not available; therefore, these classes were handled as a composition of the four HER2 subtypes. Annotations were performed by a medical student and reviewed by a board-certified pathologist. A 24-8-8 train-validation-test split was carried out on WSI level, with an equal amount of HER2-scored sections in each set. In Figure 2, the tumor tissue composition of the different dataset splits is shown.
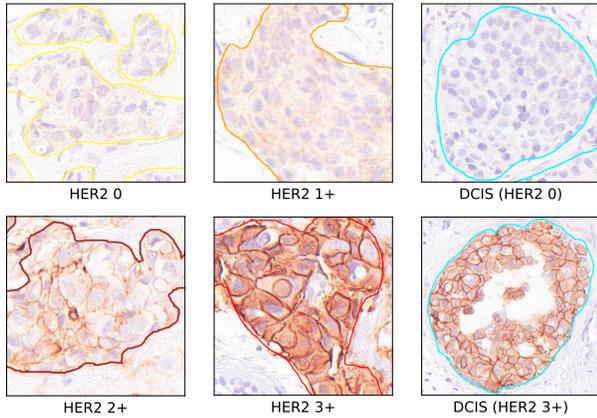


**Fig. 1**. The HER2 tumor subtypes present in the HER2 annotations, including two DCIS examples with tissue corresponding to different HER2 subtypes.
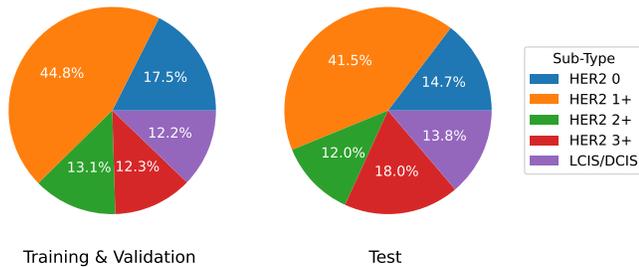


**Fig. 2**. Distribution of the different HER2 subtypes across the annotated tumor regions for the combined training and validation set, as well as for the test set.

### 2.2. Synthetic Image Generation

Figure 3 illustrates how synthetic images were created using semantic conditioning. As a first step, we modified the label masks, so that the resulting synthetic dataset was subtype-balanced. For this, we sampled an image patch and the corresponding label mask from our dataset and assigned a new, randomly selected subtype label to each tumor tissue instance. Due to a sufficiently high number of tumor tissue instances, the random assignment resulted in a subtype-balanced synthetic dataset. The modified label masks were then utilized for three methods of synthetic image generation as follows:

**GAN image generation.** We used the GAN-based architecture proposed by Park *et al*. in [15] using spatially-adaptive normalization, which enables more realistic outputs for conditioning masks where only a single class is present. Sampling of such label masks is a common finding in our dataset; therefore, this technique ensures higher-quality synthetic images. To allow various outputs for the same conditioning, we used the latent space representation of a variational autoencoder as additional input, as proposed by Park *et al*. in the Appendix of their work. For image generation, a random vector was used as additional input.

**Diffusion model image generation.** We utilized latent diffusion models, as proposed by Rombach, Blattmann, *et al*. [16]. An autoencoder, which consists of a compressing encoder and decompressing decoder, is utilized to create a lower-dimensional latent representation of the data. The image generation process takes place in the latent space, which reduces computational cost due to the compression while achieving state-of-the-art performance. We considered this architecture advantageous as we expect it to better scale for generating large(r) image patches and large datasets, although we used relatively small patches of size $512 \times 512$ in this work. To create synthetic images, the latent diffusion model was utilized with the same modified label masks as the GAN.

**Diffusion model inpainting.** We used the latent diffusion model to inpaint the tumor in existing images, conditioned by the modified HER2 subtypes. Unlike the above two methods, diffusion model inpainting only modifies the tumor tissue instances while keeping the background unchanged. Inpainting with diffusion models takes the context of an image into account, allowing the model to capture background characteristics when inpainting the new tumor tissue. Depending on how the background characteristics affect the data, this effect could be positive or negative for the subsequent task.

The generative models, as well as the subsequent segmentation task, were trained with the data described in Subsection 2.1. We followed the GAN implementation provided by Park *et al*. [1], which we trained with a batch size of 16 and the Adam optimizer (learning rate 1e-5). For latent diffusion, we adapted the implementation provided by Rombach, Blattmann, and colleagues [2]. The autoencoder for latent space computation was based on the provided vq-f4 configuration, which compressed the input by a factor of eight. The diffusion model was trained with a batch size of 16 and the Adam optimizer (learning rate 1e-6).

### 2.3. Segmentation Network

The segmentation aimed to analyze the effects of various quantities of synthetic training images from different generative methods. We used the U-Net architecture, proposed by Ronneberger *et al.* [17], for segmentation, which has proven itself a suitable architecture for segmentation in histopathology. For our experiments, we used a U-Net

---

[1]https://github.com/NVlabs/SPADE
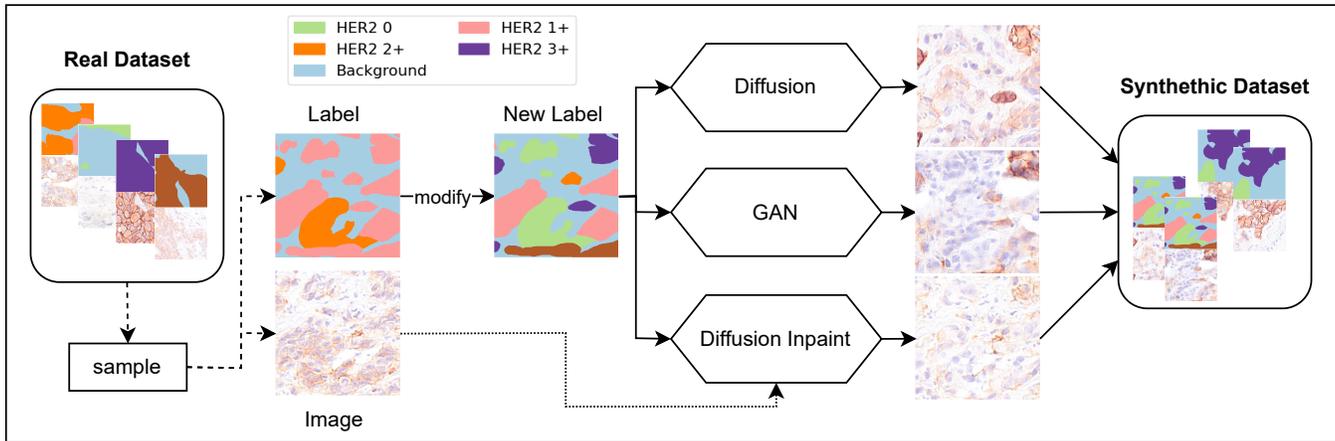[2]https://github.com/CompVis/latent-diffusion

**Fig. 3**. Illustration of how the synthetic datasets are created. An image and a corresponding label mask are sampled from the real dataset, and the subtype label of each tumor tissue instance is randomly modified. With the new label masks, synthetic images are created using a GAN, a diffusion model or diffusion inpainting. The generated images, together with their new label masks, are added to the synthetic dataset.

architecture with a ResNet-34 [18] backbone that was pre-trained on ImageNet [19]. The segmentation task consisted of two output classes (background, tumor) and was optimized with a combined Dice [20] and cross-entropy loss. Hyperparameters were optimized on the tumor subtype-sampled dataset and retained for all experiments. The Adam optimizer (learning rate 1e-6) was utilized, with a batch size of 16.

## 3. EXPERIMENTS AND RESULTS

We evaluate the proposed methods both qualitatively in terms of a visual assessment of the generated images and quantitatively with regards to the resulting segmentation performance. For quantitative evaluation of the tumor segmentation, two metrics were considered. The tumor Dice score evaluates the segmentation result independent of subtype. The variance between the HER2 tumor subtype recalls evaluates how largely the segmentation performance varies across the individual HER2 subtypes. We will call this metric *subtype variance* and lower values are favorable.

Different combinations of training data for the segmentation network were evaluated. The two baselines sampled all tumor tissue once independently of the subtypes (tumor sampled) and once uniformly across all tumor subtypes (subtype sampled). Different amounts of synthetic images extended the subtype-sampled dataset to measure the impact of synthetic data. In our experiments, we added 50%, 100%, 200%, and 400% of the original dataset size as additional synthetic images. All experiments were repeated five times to create reliable results.

### 3.1. Qualitative Results - Image Generation

Figure 4 shows an example of synthetic images created with semantic conditioning. All three synthetic images are visually similar to real HER2 histopathology images. The generated tumor structures show the staining characteristics of the subtype they were conditioned on. Signs of repeating patterns are visible in the GAN-generated images, and the background was often created as plain white without background structures. Diffusion-generated images show a high variation within all tissue types, and even rare background artifacts, like those seen in the example image, were gener-
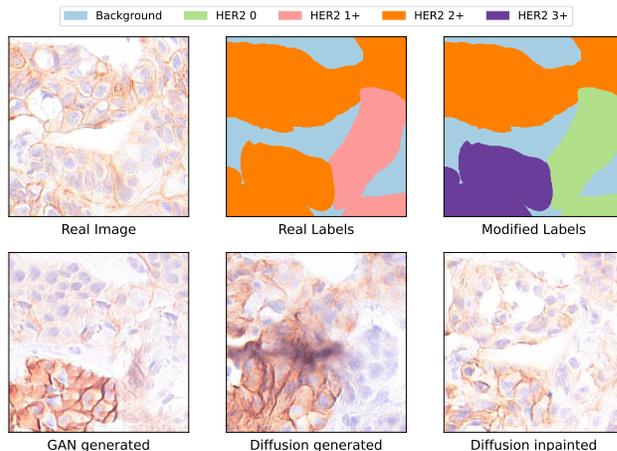


**Fig. 4**. Visual comparison of images created by generative networks.

ated. Diffusion-inpainted images also show a high variation within the created HER2 subtypes, but the network favored the creation of tissue with staining closer to the original image.

### 3.2. Quantitative Results - Tumor segmentation

The performance metrics for the tumor segmentation are visualized in Figure 5. Subtype sampling is superior to tumor sampling in both metrics. Adding synthetic data to the subtype-sampled data increased the tumor Dice score for all synthetic image methods, with diffusion-generated images performing best, GAN-generated images second, and diffusion-inpainted images last. For the *subtype variance*, diffusion-inpainted images performed inferior to subtype sampling, while GAN-generated images and diffusion-generated images were superior. Diffusion-generated images achieved a better *subtype variance* than GAN-generated images in all but one case. The best metrics were achieved when adding 100% synthetic images, where diffusion-generated images increased the Dice score from 0.833 to 0.854 and decreased the *subtype variance* by 47.8%.
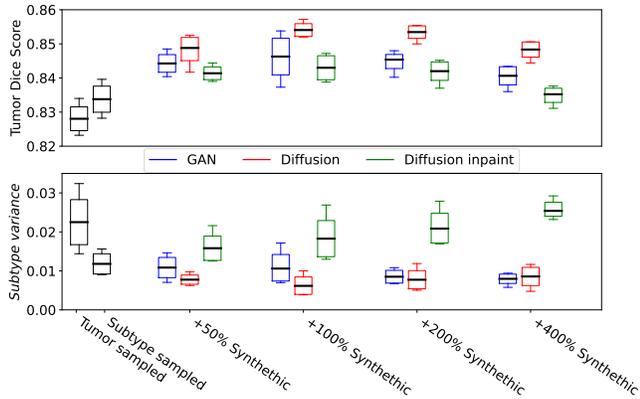
**Fig. 5**. Boxplots of the tumor Dice score and the *subtype variance* for different configurations. Mean and standard deviation are visualized with the boxplot, while the whiskers mark the minimum and maximum values.



**Fig. 6**. Row normalized confusion matrix with the tumor subtypes. Left are the averaged values from the subtype-sampled runs, while right are the results from the experiment where 100% diffusion images were added.

For the subtype-sampled dataset and the best-performing experiment, the averaged confusion matrices per subtype are shown in Figure 6. The most significant performance increase could be observed for the HER2 0 subtype, where the recall increased from 0.64 to 0.73. For the rest of the subtypes, minor improvements were achieved.

## 4. DISCUSSION

Although the qualitative results of the synthetically generated images look promising, some aspects remain open for discussion. GAN-generated images show signs of repeating patterns, a common finding among GANs, which was not fully avoided with spatially-adaptive normalization [15]. Additionally, the GAN network favored creating plain white background, which we suspect to be the network collapsing to the "easiest" solution for the background present in the training data. The diffusion-generated images were visually more compelling and even showed rare background structures, e.g. staining artifacts, which can help the segmentation network to become more robust toward these structures. Diffusion model inpainting created tumor tissue with staining intensities more similar to the original image. We suspect that the remaining background information influenced the created staining levels in an unfavourable manner.

All presented methods for synthetic image generation improved the composite tumor segmentation, although no additional annotated images were available for the generative networks. We assume that generative networks might be able to interpolate between tissue features; for example, they could combine subtype invariant features, like tumor or cell shape, with subtype-specific features, like staining intensity. This could lead to synthetic tumor tissue, which has a combination of features that is not present in the original dataset. We suspect this to be the main reason for the lower *subtype variance* when GANs and diffusion models are utilized.

The main performance benefit for diffusion models accrued for the HER2 0 subtype, which had a significantly lower recall than the other subtypes. We suspect that the lower recall was caused by the absence of the brownish staining for this class, which appears to be an easy indicator of tumor tissue. This indicates that the subtyping into four HER2 classes might not be optimal and a subtyping into non-stained (HER2 0) and stained (HER2 1+, HER2 2+,
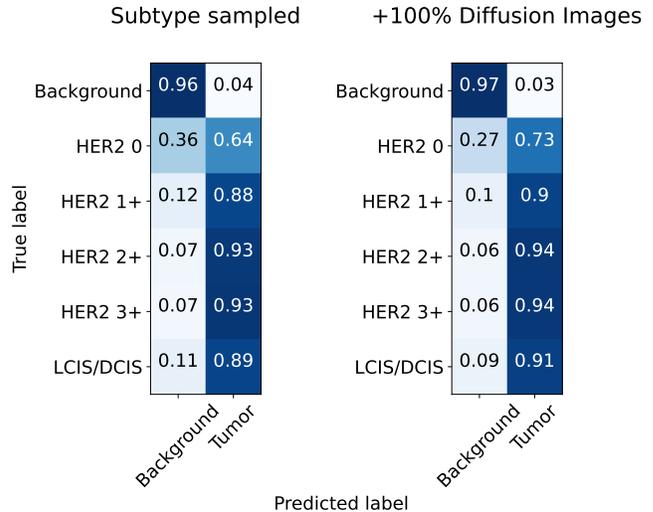
HER2 3+) would be an interesting alternative. Diffusion models appeared to be able to correctly generate non-stained HER2 tumor tissue for the learned representations, thus improving the performance for this subtype. Although we only consider tumor segmentation in this work, this more uniform performance between subtypes could lead to more reliable automatic HER2 scoring in the future, since the scoring is based on the proportion of the HER2 subtype tissue [5].

These results are promising, but some limitations have to be noted. The same persons annotated training and test data, which could introduce a bias that affected evaluation metrics. Training of GANs is notoriously unstable and hard to monitor; therefore, it is possible that the model used in this work was not perfectly adapted to the data and could have produced better results.

## 5. CONCLUSION

We proposed to subtype balance HER2 data with generative models. We showed the suitability of generative models, especially diffusion models, to generate semantically-conditioned synthetic images with a realistic appearance. Combining an equal amount of real images with diffusion model-generated images increased the Dice score of the tumor segmentation by 2.43% and reduced the variance between the tumor subtype recalls by 47.8%. This method is superior to oversampling subtypes and does not require additional annotated data.

The current approach requires the annotation of individual HER2 to improve the tumor segmentation, and other kinds of so far non-annotated subtypes might exist. Future work could explore methods to alleviate this by, for example, unsupervised tumor area clustering and balancing these clusters.

Another future work could explore the use of fully synthetic training data. Besides the subsequent algorithm performance, one interesting aspect could be whether the synthetic images can be traced back to real patients. Such work could lay the foundation for using fully synthetic datasets and thereby reduce data privacy concerns.

## 6. COMPLIANCE WITH ETHICAL STANDARDS

All procedures performed in studies involving human participants were in accordance with the ethical standards of the institutional and/or national research committee and with the 1964 Helsinki declaration and its later amendments or comparable ethical standards.

## 7. ACKNOWLEDGEMENT

## 8. REFERENCES

[1] Seyma Yucer, Furkan Tektas, Noura Al Moubayed, and Toby P. Breckon, "Measuring hidden bias within face recognition via racial phenotypes," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, January 2022, pp. 995–1004.

[2] Lisa M Koch, Christian M Schürch, Arthur Gretton, and Philipp Berens, "Hidden in plain sight: Subgroup shifts escape ood detection," in *Medical Imaging with Deep Learning*, 2021.

[3] Gautam K Malhotra, Xiangshan Zhao, Hamid Band, and Vimla Band, "Histological, molecular and functional subtypes of breast cancers," *Cancer biology & therapy*, vol. 10, no. 10, pp. 955–960, 2010.

[4] Sibylle Loibl and Luca Gianni, "Her2-positive breast cancer," *The Lancet*, vol. 389, no. 10087, pp. 2415–2429, 2017.

[5] Antonio C Wolff, M Elizabeth Hale Hammond, Kimberly H Allison, Brittany E Harvey, Pamela B Mangu, John MS Bartlett, Michael Bilous, Ian O Ellis, Patrick Fitzgibbons, Wedad Hanna, et al., "Human epidermal growth factor receptor 2 testing in breast cancer: American society of clinical oncology/college of american pathologists clinical practice guideline focused update," *Archives of pathology & laboratory medicine*, vol. 142, no. 11, pp. 1364–1382, 2018.

[6] Hui Han, Wen-Yuan Wang, and Bing-Huan Mao, "Borderline-smote: a new over-sampling method in imbalanced data sets learning," in *International conference on intelligent computing*. Springer, 2005, pp. 878–887.

[7] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio, "Generative adversarial networks," *Communications of the ACM*, vol. 63, no. 11, pp. 139–144, 2020.

[8] Yizhou Chen, Xu-Hua Yang, Zihan Wei, Ali Asghar Heidari, Nenggan Zheng, Zhicheng Li, Huiling Chen, Haigen Hu, Qianwei Zhou, and Qiu Guan, "Generative adversarial networks in medical image augmentation: a review," *Computers in Biology and Medicine*, p. 105382, 2022.

[9] Val Andrei Fajardo, David Findlay, Charu Jaiswal, Xinshang Yin, Roshanak Houmanfar, Honglei Xie, Jiaxi Liang, Xichen She, and DB Emerson, "On oversampling imbalanced data with deep conditional generative models," *Expert Systems with Applications*, vol. 169, pp. 114463, 2021.

[10] Jonathan Ho, Ajay Jain, and Pieter Abbeel, "Denoising diffusion probabilistic models," *Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, 2020.

[11] Prafulla Dhariwal and Alexander Nichol, "Diffusion models beat gans on image synthesis," *Advances in Neural Information Processing Systems*, vol. 34, pp. 8780–8794, 2021.

[12] Mathias Öttl, Jana Mönius, Christian Marzahl, Matthias Rübner, Carol I Geppert, Arndt Hartmann, Matthias W Beckmann, Peter Fasching, Andreas Maier, Ramona Erber, et al., "Superpixel pre-segmentation of her2 slides for efficient annotation," in *Bildverarbeitung für die Medizin 2022*, pp. 254–259. Springer, 2022.

[13] Christian Marzahl, Marc Aubreville, Christof A Bertram, Jennifer Maier, Christian Bergler, Christine Kröger, Jörn Voigt, Katharina Breininger, Robert Klopfleisch, and Andreas Maier, "Exact: a collaboration toolset for algorithm-aided annotation of images with annotation version control," *Scientific reports*, vol. 11, no. 1, pp. 1–11, 2021.

[14] Antonio C Wolff, ME Hammond, Jared N Schwartz, Karen L Hagerty, D Craig Allred, Richard J Cote, Mitchell Dowsett, Patrick L Fitzgibbons, Wedad M Hanna, Amy Langer, et al., "College of american pathologists guideline recommendations for human epidermal growth factor receptor 2 testing in breast cancer," *J clin oncol*, vol. 25, no. 1, pp. 118–145, 2007.

[15] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu, "Semantic image synthesis with spatially-adaptive normalization," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 2337–2346.

[16] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer, "High-resolution image synthesis with latent diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10684–10695.

[17] 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*. Springer, 2015, pp. 234–241.

[18] 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*, 2016, pp. 770–778.

[19] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.

[20] Carole H Sudre, Wenqi Li, Tom Vercauteren, Sebastien Ourselin, and M Jorge Cardoso, "Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations," in *Deep learning in medical image analysis and multimodal learning for clinical decision support*, pp. 240–248. Springer, 2017.