# MIND THE GAP: SCANNER-INDUCED DOMAIN SHIFTS POSE CHALLENGES FOR REPRESENTATION LEARNING IN HISTOPATHOLOGY

*Frauke Wilm*[1,2] *, Marco Fragoso*[3]*, Christof A. Bertram*[4]*, Nikolas Stathonikos*[5]*, Mathias Öttl*[1]*,*
*Jingna Qiu*[2]*, Robert Klopfleisch*[3]*, Andreas Maier*[1]*, Marc Aubreville*[6,†]*, Katharina Breininger*[2,†]

[1] Pattern Recognition Lab, Friedrich-Alexander-Universität Erlangen-Nürnberg, Germany
[2] Department AIBE, Friedrich-Alexander-Universität Erlangen-Nürnberg, Germany
[3] Institute of Veterinary Pathology, Freie Universität Berlin, Germany
[4] Institute of Pathology, University of Veterinary Medicine, Vienna, Austria
[5] Pathology Department, University Medical Centre Utrecht, The Netherlands
[6] Technische Hochschule Ingolstadt, Ingolstadt, Germany

## ABSTRACT

Computer-aided systems in histopathology are often challenged by various sources of domain shift that impact the performance of these algorithms considerably. We investigated the potential of using self-supervised pre-training to overcome scanner-induced domain shifts for the downstream task of tumor segmentation. For this, we present the *Barlow Triplets* to learn scanner-invariant representations from a multi-scanner dataset with local image correspondences. We show that self-supervised pre-training successfully aligned different scanner representations, which, interestingly only results in a limited benefit for our downstream task. We thereby provide insights into the influence of scanner characteristics for downstream applications and contribute to a better understanding of why established self-supervised methods have not yet shown the same success on histopathology data as they have for natural images.

***Index Terms***— Histopathology, Domain Shift, Representation Learning, Barlow Twins

## 1. INTRODUCTION

Machine learning-based image analysis has experienced an upswing in histopathology, to which the ease of creating digital whole slide images (WSIs) using slide scanners greatly contributed. Whilst computer-aided systems often excel within their training domain - sometimes even surpassing the performance of experienced pathologists [1] - their performance considerably decreases for out-of-distribution samples [2, 3, 4]. Representation learning has previously shown great success for domain generalization [5] by pre-training models with cross-domain image pairs and aligning their representations directly in the embedding space. These cross-domain image pairs are usually created without supervision

through augmentations that simulate possible domain shifts. Recent work has shown that self-supervised pre-training can help to extract task-relevant information and increase the performance for downstream tasks even without supervision for these tasks at the pre-training stage [6, 7, 8]. These works, however, see a high dependency of the downstream performance on the amount and type of information the constructed views share. If two views share too little task-relevant information, performance on the downstream task can degrade, whilst if they share too many tasks-irrelevant (nuisance) features, self-supervised pre-training can lead to the extraction of shortcut features [7, 9]. Based on these observations, Tian *et al*. stated the *InfoMin* objective of finding a sweet spot where two views retain enough task-relevant information without sharing irrelevant nuisances [7].

In this work, we try to find this sweet spot for creating views that help to mitigate scanner-induced domain shifts. For this, we utilized a multi-scanner dataset with local image correspondences. By digitizing the same sample with multiple slide scanning systems, we preserved task-relevant features (e.g. tissue morphology) while introducing "intrinsic" domain shifts that might exceed apparent differences such as color or resolution. For representation learning, we adapted the *Barlow Twin* [10] architecture to accommodate multi-scanner tuples and evaluated this use of self-supervised pre-training for the downstream task of tumor segmentation on canine skin cancer specimens. Our experiments thereby extend the work of Stacke *et al*. who have applied representation learning for the task of tissue phenotyping in the presence of cross-organ and cross-laboratory domain shifts. The authors, however, faced challenges when applying established methods of self-supervised learning to histopathology and argued that "the subtleties of the difference between classes in histology data makes it more challenging to find effective augmentations" [9], which we circumvent by using a multi-scanner dataset that inherently possesses these differences.
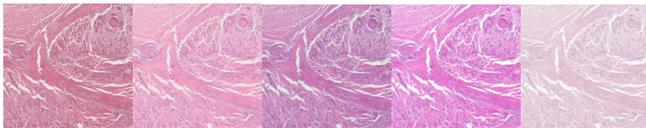
---

† Shared senior authors.

Even though we show that self-supervised pre-training helped to align the scanner embeddings, we only observed a limited benefit for our downstream task. We suspect that feature alignment enhanced some scanner characteristics that might be beneficial for tumor segmentation, but that the small inter-class variance of histopathology limits the benefit for the downstream task. We thereby contribute to a better understanding of why self-supervised learning has not yet shown the same potential for histopathology data as it has for natural images [9] and provide recommendations for using representation learning to approach the domain shift inherent in multi-scanner histopathology datasets.

## 2. MATERIALS

Our experiments were performed on the squamous cell carcinoma (SCC) subset of the publicly available CATCH dataset [11]. Use of these samples was approved by the local governmental authorities (State Office of Health and Social Affairs of Berlin, approval ID: StN 011/20). The specimens were originally digitized with the Aperio ScanScope CS2 (Leica, Germany) at a resolution of $0.25\,\mu\text{m/px}$ using a $40\times$ objective lens. To create a multi-scanner dataset with local correspondences, we digitized the samples with four additional slide scanners (exemplary patches in Figure 1) :

- NanoZoomer S210 (Hamamatsu, Japan), $0.22\,\mu\text{m/px}$
- NanoZoomer 2.0-HT (Hamamatsu, Japan), $0.23\,\mu\text{m/px}$
- Pannoramic 1000 (3DHISTECH, Hungary), $0.25\,\mu\text{m/px}$
- Aperio GT 450 (Leica, Germany), $0.26\,\mu\text{m/px}$



(a) CS2  (b) NZ 210  (c) NZ 2.0  (d) P 1000  (e) GT 450

**Fig. 1**: Exemplary patch of the multi-scanner dataset.

All Aperio ScanScope CS2 WSIs were manually annotated for tumor and three skin tissue classes (epidermis, dermis, subcutis) using the open-source software SlideRunner [12]. For automatic background detection, the images were converted to grayscale and non-annotated regions with a grayscale value above 235 were labeled as background. All remaining non-annotated regions were excluded from training and evaluation. Due to the image-to-image correspondences, all annotations could be transferred to the remaining scanners using the WSI registration algorithm by Marzahl et al. [13]. The registration success was visually validated by overlaying the transformed polygon annotations onto the WSIs. Due to severe scanning artifacts in at least one of the scans, six specimens were excluded from the dataset, resulting in a total of 220 WSIs (44 samples digitized with five scanners each). For algorithm development, a slide-level split into 30-5-9

train-validation-test was performed. The GT 450 and P 1000 scanners were selected as hold-out test scanners, to test the model's capability of generalizing to unseen domains.

## 3. METHODS

Given the multi-scanner dataset with local image correspondences, the goal was to extract scanner-invariant features for a subsequent tumor segmentation task. For this, we followed a two-step training pipeline: We first pre-trained a feature extractor in a self-supervised fashion and then switched to a fully-supervised training setup for the segmentation task.

For pre-training, we used the *Barlow Twin* architecture [10], which originally uses a Siamese structure to create representations for two augmented versions of the same input image, that are then projected into a higher dimensional feature space $\mathbb{R}^d$, where the cross-correlation matrix $\mathcal{C}$ of both embeddings is computed. The *Barlow Twin* loss $\mathcal{L}_{\mathcal{BT}}$ enforces this matrix to be similar to the identity matrix:

$$\mathcal{L}_{\mathcal{BT}} = \sum_{i=0}^{d}(1 - \mathcal{C}_{ii})^2 + \lambda \sum_{i=0}^{d}\sum_{\substack{j=0 \\ j\neq i}}^{d}\mathcal{C}_{ij}{}^2 \qquad (1)$$

By enforcing high values on the main diagonal, features become invariant to the applied distortions (or in our case: the scanner-induced domain shift), whilst low values elsewhere disentangle features and thereby reduce redundancy.
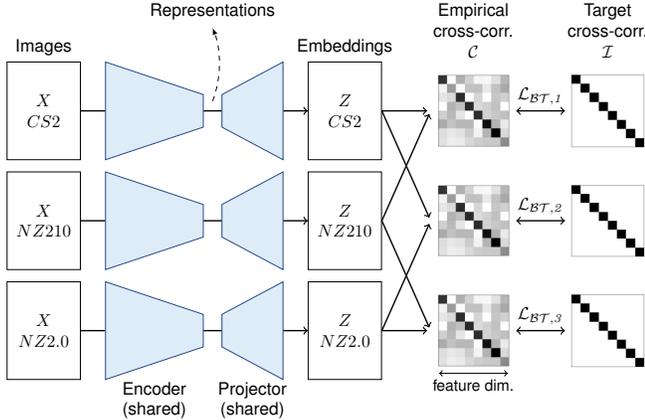
We adapted the *Barlow Twins* to generalize to input tuples composed of one patch per training scanner. Figure 2 visualizes our adapted version for the case of three training scanners - the *Barlow Triplets*. We performed a straightforward extension of the *Barlow Twin* loss function to a *Barlow Tuple* loss $\mathcal{L}_{\mathcal{BT}^\star}$, that computes the joint loss of all unique pairs of a given input tuple as:

$$\mathcal{L}_{\mathcal{BT}^\star} = \frac{1}{K(n,2)}\sum_{k=1}^{K(n,2)}\mathcal{L}_{\mathcal{BT},\,k} \qquad , \qquad (2)$$

where $K(n,2) = \frac{n!}{(n-2)!2!}$ computes the number of unique pairs that can be constructed from $n$ scanner embeddings, e.g. $K(3,2) = \frac{6}{2} = 3$ combinations for scanner triplets.

As encoder, we used a ResNet18 [14] pre-trained on ImageNet [15]. For the projector, we followed the implementation of Zbontar et al. and used three linear layers with a four-fold up-scaling in feature dimension [10], i.e. 2048 output units for a ResNet18 encoder. Similar to Zbontar et al. the first two linear layers were followed by batch normalization and rectified linear units.

After self-supervised pre-training, we used the annotation database to perform a fully-supervised segmentation training into background, tumor, and a third non-tumor class that combined all skin tissue classes. To construct a segmentation model, the encoder was connected to a decoding branch using skip connections in a U-Net-like fashion [16].

**Fig. 2**: *Barlow Triplets*. Figure adapted from [10].

**Training parameters.** The network was trained on $256 \times 256$ pixel-sized image patches. To cover more tissue context, we extracted the patches at a lower resolution of $4\,\mu\mathrm{m/px}$. For each epoch, we sampled 50 patches from each WSI using an equal weighting of tumor and non-tumor and 10% background patches. To increase the diversity of training data, this guided selection of patches was repeated for each training epoch but kept unchanged for the validation set. All patches were z-score normalized using the mean and standard deviation of all tissue-containing areas of the CS2 training WSIs.
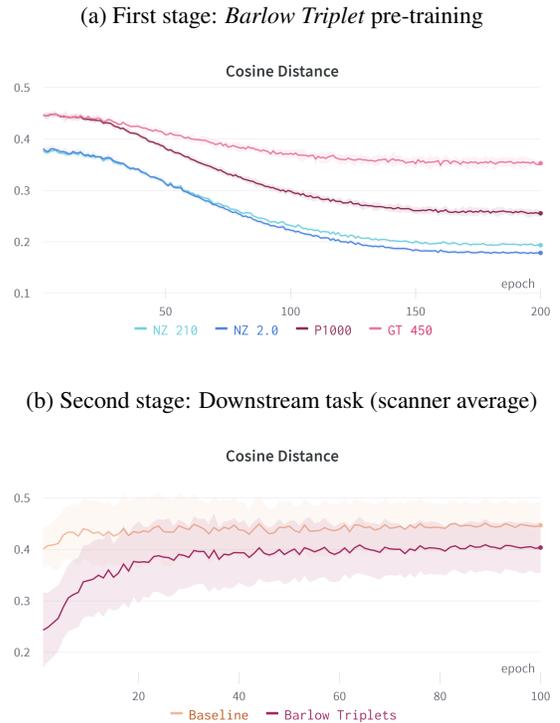
During self-supervised pre-training, we trained the encoder for 200 epochs, after which we observed model convergence. We used the Adam optimizer, a batch size of 64, and a cyclic learning rate with a maximum of $10^{-6}$. Following Zbontar *et al.* [10], we selected $\lambda = 5 \times 10^{-3}$ for $\mathcal{L}_{\mathcal{BT}\star}$. After pre-training, we switched to a fully-supervised setup and trained the extended models for another 100 epochs using a cyclic learning rate with a maximum of $10^{-4}$. For this, we used a reduced batch size of 8 (due to memory constraints) and a combination of cross-entropy and Dice loss [17]. Model selection was guided by the highest mean intersection over union (mIoU) score on the validation patches.

## 4. EXPERIMENTS AND RESULTS

To evaluate our representations, we compared the performance of our pre-trained encoder to a baseline initialized with ImageNet [15] weights for the downstream task of tumor segmentation. We evaluated two different settings: A single-domain approach, where only the CS2 patches were used for fully-supervised training of tumor segmentation (single), and a multi-domain approach, where the CS2 and both NanoZoomer scanners were used (multi). For each experiment, we repeated the fully-supervised training three times and averaged the test performance. For a fair comparison we used seeding to initialize the optimizers and data loaders to ensure that for each seed, all models obtained the same

ordering of randomly sampled training patches and the same validation patches for the fully-supervised training stage. We used the same hyperparameters to train all models.

To evaluate whether the *Barlow Triplets* successfully aligned the pair-wise scanner representations, we monitored the mean cosine distance of the CS2 embeddings to the other scanners at the encoder bottleneck, as visualized in Figure 3a. The plot shows a steady decrease especially for the two NanoZoomers incorporated in the self-supervised pre-training, but also for the two scanners that were not seen during training (P 1000 and GT 450). Figure 3b compares the mean patch embedding distance to the CS2 patches (averaged across all scanners and across single- and multi-domain experiments) for the second stage of the training pipeline. The plot shows that the pre-trained *Barlow Triplets* started with a considerably lower cosine distance than the ImageNet-initialized baseline. Even though during this second stage, the pair-wise cosine distance partially increased again, it converged to an overall lower value than the baseline.

(a) First stage: *Barlow Triplet* pre-training



(b) Second stage: Downstream task (scanner average)



**Fig. 3**: Cosine distance of CS2 scanner to seen (NZ 210, NZ 2.0) and unseen (P 1000, GT 450) scanners ($\mu \pm \sigma$ of three repetitions).
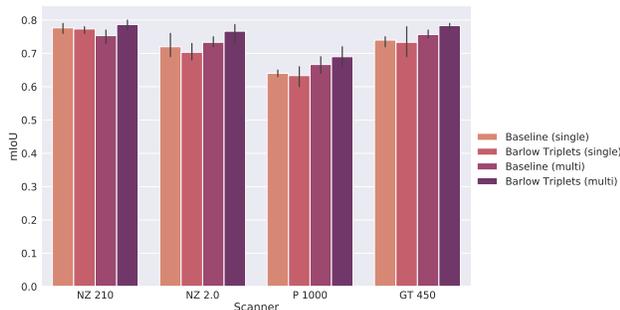
Table 1 summarizes the test set performance of the baseline and the pre-trained model in the single-domain and multi-domain setup of the fully-supervised training stage ($\mu \pm \sigma$ of three repetitions). For inference, we used a moving-window approach with a 128-pixel overlap and center-cropped the

segmentation output. We accumulated the confusion matrix across all WSIs of our test set and calculated the mIoU. The performance of the single-domain baseline highlights the domain shift within the cross-scanner dataset, as the mIoU varies between $0.76$ and $0.84$. The out-of-domain mIoU was increased when using the self-supervised pre-trained *Barlow Triplets* in this single-domain setup. Incorporating the two NanoZoomers into the fully-supervised training improved the performance across these two scanners but slightly degraded the performance on the CS2 scanner. Interestingly, the performance on the unseen P1000 scanner was improved the most. The pre-trained *Barlow Triplets* improved the baseline performance of the multi-domain setup overall by either increasing the mean performance or reducing the variance.

| mIoU | CS2 | NZ 210 | NZ 2.0 | P 1000 | GT 450 |
|---|---|---|---|---|---|
| Baseline (single) | $0.76 \pm 0.02$ | $0.79 \pm 0.02$ | $0.84 \pm 0.01$ | $0.78 \pm 0.02$ | $0.84 \pm 0.02$ |
| *Barlow Triplets* (single) | $0.75 \pm 0.01$ | $0.81 \pm 0.02$ | $0.85 \pm 0.01$ | $0.79 \pm 0.02$ | $0.85 \pm 0.01$ |
| Baseline (multi) | $0.71 \pm 0.02$ | $0.80 \pm 0.02$ | $0.88 \pm 0.01$ | $0.85 \pm 0.01$ | $0.83 \pm 0.02$ |
| *Barlow Triplets* (multi) | $0.75 \pm 0.02$ | $0.82 \pm 0.01$ | $0.88 \pm 0.00$ | $0.85 \pm 0.01$ | $0.83 \pm 0.01$ |

**Table 1**: Class-averaged (tumor, non-tumor, background) mean intersection over union ($\mu \pm \sigma$ of three repetitions) of baselines and *Barlow Triplets*.

Figure 4 visualizes the agreement of the CS2 prediction and target scanner predictions. For the single-domain setup, the baseline and the pre-trained model yielded a similar averaged concordance across all scanners. The multi-scanner fully-supervised training (Baseline (multi)) helped to align the segmentation outputs of all scanners and self-supervised pre-training increased the concordance even further.



**Fig. 4**: Concordance of segmentation outputs to the CS2 prediction measured as mean intersection over union of masks.

## 5. DISCUSSION

The large discrepancy in segmentation performance across the different scanners has highlighted the scanner-induced domain shift within our dataset. This gap could neither be completely closed by training the U-Net fully supervised on multiple scanner domains nor by initializing the model with the pre-trained encoder using the *Barlow Triplets*.

During pre-training we observed a continuous decrease of the pair-wise cosine distance at the encoder bottleneck, indicating that the *Barlow Triplets* helped to align the scanner representations. When integrating this encoder into a U-Net-like architecture, however, these bottleneck embeddings are passed through the decoder and enriched by features from the skip connection. We here see the potential risk of scanner-specific features being bypassed through the skip connections of the U-Net. Future work could therefore use representation learning strategies on multiple encoder levels.

Interestingly, all models performed worst on the CS2 scanner, even though this scanner was seen for all experiments (single- and multi-domain). A closer evaluation of the segmentation masks showed that the model considerably underestimated the tumor area on the CS2 WSIs compared to all other scanners. This indicates that some scanner characteristics were inherently more beneficial for the downstream task of tumor segmentation than others, e.g. different color representations that facilitate the separation of hematoxylin and eosin components or differences in contrast or sharpness. While the results are not fully conclusive, the performance increase in the (multi) setting could indicate that the CS2 indeed benefits from a feature alignment that enhances these characteristics. Overall, we however only observed a slight increase in mIoU compared to the respective baseline. Although the multi-scanner views should fulfill the *InfoMin* principle [7] and put aside the need for finding effective augmentations [9], the small inter-class variance in histopathology may make it difficult to fully exploit this. Since we observe a stronger effect on the segmentation concordance (Figure 4) than on the performance (Table 1), it will be interesting to understand the effect of label noise on feature alignment.

Previous work has mostly employed representation learning as self-supervised pre-training due to limited annotations for downstream tasks. One of the key strengths of the CATCH dataset is its extensive annotation database, which allowed the fully-supervised training of tumor segmentation with high performance on the test set [18]. To transfer this performance across different scanner domains, future work could explore a joint or alternating self-supervised training of the encoder and fully-supervised training for the downstream task.

## 6. CONCLUSION

In this work, we investigated self-supervised pre-training for scanner-induced domain shifts in histopathology. Our experiments show that self-supervised pre-training is generally applicable to the task of cross-scanner representation alignment but did not yield a significant performance boost for our downstream task. Our results indicate that some scanner-specific characteristics might be relevant for the downstream task of tumor segmentation, which has to be considered when employing representation learning to mitigate scanner-induced domain shifts.

# 7. COMPLIANCE WITH ETHICAL STANDARDS

All specimens were submitted by veterinary clinics or surgeries for routine diagnostic examination of neoplastic disease. As to local regulations, no ethical vote is required for these samples.

# 8. ACKNOWLEDGMENTS

# 9. REFERENCES

[1] Marc Aubreville et al., "Deep learning algorithms outperform veterinary pathologists in detecting the mitotically most active tumor region," *Scientific Reports*, vol. 10, no. 16447, pp. 1–11, 2020.

[2] Marc Aubreville et al., "Quantifying the scanner-induced domain gap in mitosis detection," in *Medical Imaging with Deep Learning (MIDL), Lübeck*, 2021.

[3] Maxime W. Lafarge, Josien P.W. Pluim, Koen A.J. Eppenhof, and Mitko Veta, "Learning domain-invariant representations of histological images," *Frontiers in Medicine*, vol. 6, pp. 162, 2019.

[4] Karin Stacke, Gabriel Eilertsen, Jonas Unger, and Claes Lundström, "Measuring domain shift for deep learning in histopathology," *IEEE Journal of Biomedical and Health Informatics*, vol. 25, no. 2, pp. 325–336, 2020.

[5] Saeid Motiian, Marco Piccirilli, Donald A Adjeroh, and Gianfranco Doretto, "Unified deep supervised domain adaptation and generalization," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 5715–5725.

[6] Yao-Hung Hubert Tsai, Yue Wu, Ruslan Salakhutdinov, and Louis-Philippe Morency, "Self-supervised learning from a multi-view perspective," in *Proceedings of the 9th International Conference on Learning Representations (ICLR)*, 2021.

[7] Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola, "What makes for good views for contrastive learning?," *Advances in Neural Information Processing Systems*, vol. 33, pp. 6827–6839, 2020.

[8] Jason D Lee, Qi Lei, Nikunj Saunshi, and Jiacheng Zhuo, "Predicting what you already know helps: Provable self-supervised learning," *Advances in Neural Information Processing Systems*, vol. 34, pp. 309–323, 2021.

[9] Karin Stacke, Jonas Unger, Claes Lundström, and Gabriel Eilertsen, "Learning representations with contrastive self-supervised learning for histopathology applications," *Machine Learning for Biomedical Imaging*, vol. 1, 2022.

[10] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny, "Barlow twins: Self-supervised learning via redundancy reduction," in *International Conference on Machine Learning*. PMLR, 2021, pp. 12310–12320.

[11] Frauke Wilm et al., "CAnine CuTaneous Cancer Histology dataset (version 1)," *The Cancer Imaging Archive*, 2022, https://doi.org/10.7937/TCIA.2M93-FX66.

[12] Marc Aubreville, Christof Bertram, Robert Klopfleisch, and Andreas Maier, "SlideRunner," in *Bildverarbeitung für die Medizin 2018*, pp. 309–314. Springer, 2018.

[13] Christian Marzahl et al., "Robust quad-tree based registration on whole slide images," in *MICCAI Workshop on Computational Pathology*. PMLR, 2021, pp. 181–190.

[14] 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.

[15] Olga Russakovsky et al., "ImageNet Large Scale Visual Recognition Challenge," *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.

[16] 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.

[17] 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.

[18] Frauke Wilm et al., "Pan-tumor CAnine cuTaneous Cancer Histology (CATCH) dataset," *Scientific Data*, vol. 9, no. 588, pp. 1–13, 2022.