# Self-Rule to Multi-Adapt: Generalized Multi-source Feature Learning Using Unsupervised Domain Adaptation for Colorectal Cancer Tissue Detection

**Christian Abbet**[*]
Signal Processing Lab 5 (LTS5)
EPFL
Switzerland
christian.abbet@epfl.ch

**Linda Studer**[*]
Documents, Image and Video Analysis (DIVA)
University of Fribourg
Switzerland
linda.studer@unifr.ch

**Andreas Fischer**
Documents, Image and Video Analysis (DIVA)
University of Fribourg
Switzerland

**Heather Dawson**
Institute of Pathology
University of Bern
Switzerland

**Inti Zlobec**
Institute of Pathology
University of Bern
Switzerland

**Behzad Bozorgtabar**
Signal Processing Lab 5 (LTS5)
EPFL
Switzerland

**Jean-Philippe Thian**
Signal Processing Lab 5 (LTS5)
EPFL
Switzerland

January 20, 2022

## Abstract

Supervised learning is constrained by the availability of labeled data, which are especially expensive to acquire in the field of digital pathology. Making use of open-source data for pre-training or using domain adaptation can be a way to overcome this issue. However, pre-trained networks often fail to generalize to new test domains that are not distributed identically due to tissue stainings, types, and textures variations. Additionally, current domain adaptation methods mainly rely on fully-labeled source datasets. In this work, we propose Self-Rule to Multi-Adapt (SRMA), which takes advantage of self-supervised learning to perform domain adaptation, and removes the necessity of fully-labeled source datasets. SRMA can effectively transfer the discriminative knowledge obtained from a few labeled source domain's data to a new target domain without requiring additional tissue annotations. Our method harnesses both domains' structures by capturing visual similarity with intra-domain and cross-domain self-supervision. Moreover, we present a generalized formulation of our approach that allows the framework to learn from multiple source domains. We show that our proposed method outperforms baselines for domain adaptation of colorectal tissue type classification in single and multi-source settings, and further validate our approach on an in-house clinical cohort. The code and trained models are available open-source: `https://github.com/christianabbet/SRA`.

***Keywords*** Computational pathology · self-supervised learning · unsupervised domain adaptation · colorectal cancer

---

[*]Co-first author. Christian Abbet and Linda Studer contributed equally.

# 1 Introduction

Colorectal cancer (CRC) is one of the most common cancers worldwide, and its understanding through computational pathology techniques can significantly improve the chances of effective treatment [Geessink et al., 2019, Smit and Mesker, 2020] by refining disease prognosis and assisting pathologists in their daily routine. The data used in computational pathology most often consists of Hematoxylin and Eosin (H&E) stained whole slide images (WSIs) [Hegde et al., 2019, Lu et al., 2021] and tissue microarrays (TMAs) [Arvaniti et al., 2018, Nguyen et al., 2021]

Although fully supervised deep learning models have been widely used for a variety of tasks, including tissue classification [Kather et al., 2019] and semantic segmentation [Qaiser et al., 2019, Chan et al., 2019], in practice, it is time-consuming and expensive to obtain fully-labeled data as it involves expert pathologists. This hinders the applicability of supervised machine learning models to real-world scenarios. Weakly supervised learning is a less demanding approach that does not depend on large labeled cohorts. Examples of this approach applied to digital pathology include WSIs classification [Tellez et al., 2018, Silva-Rodríguez et al., 2021] and Multiple-Instance Learning (MIL) based on diagnostic reports [Campanella et al., 2019]. However, these methods still need an adequate training set to initialize the learning process, limiting the gain that can be achieved from using unlabeled samples.

Self-supervised learning was proposed to address limitations linked to labeled data availability. It involves a training scheme where "*the data creates its own supervision*"[Pieter et al., 2020] to learn rich features from structured unlabeled data. Applications of this approach in computational pathology include multiple tasks such as survival analysis [Abbet et al., 2020], WSIs classification [Li et al., 2021], and image retrieval [Gildenblat and Klaiman, 2019].

Over the years, various large data banks have been made available online containing samples from a variety of organs [Weinstein et al., 2013, Litjens et al., 2018, Veta et al., 2019], such as the colon and rectum [Kather et al., 2016, Shanah et al., 2016a,b, Kather et al., 2019]. This opens up possibilities for transfer learning and domain adaptation. Yet, using these data banks to develop computational pathology-based models for real-world scenarios remains challenging because of the domain gap, as these images were created under different imaging scenarios. A tissue sample's visual appearance can be heavily affected by the staining procedure [Otálora et al., 2019], the type of scanner used [Cheng et al., 2019], or other artifacts such as folded tissues [Komura and Ishikawa, 2018].

To tackle this issue, color normalization techniques [Macenko et al., 2009, Zanjani et al., 2018, Anand et al., 2019] have been widely adopted. Nevertheless, these techniques solely rely on image color information, while the morphological structure of the tissue is not taken into account [Tam et al., 2016, Zarella et al., 2017]. This could lead to unpredictable results in the presence of substantial staining variations and dark staining due to densely clustered tumor cells.

Another field of research that aims to improve the classification of heterogeneous WSIs is unsupervised domain adaptation (UDA). These methods work by learning from a rich source domain together with the label-free target domain to have a well-performing model on the target domain at inference time. UDA allows models to include a large variety of constraints to match relevant morphological features across the source and target domains.

DANN [Ganin and Lempitsky, 2015], for example, uses gradient reversal layers to learn domain-invariant features. Self-Path [Koohbanani et al., 2021] combines the DANN approach with self-supervised auxiliary tasks. The selected tasks reflect the structure of the tissue and are assumed to improve the stability of the framework when working with histopathological images. Such auxiliary tasks include hematoxylin channel prediction, Jigsaw puzzle-solving, and magnification prediction. Another example is CycleGAN [Zhu et al., 2017], which takes advantage of adversarial learning to map images between the source and target domain cyclically. However, adversarial approaches can fall short because they do not consider task-specific decision boundaries and only try to distinguish the features as either coming from the source or target domain [Saito et al., 2018a].

A further issue is that most UDA methods consider fully-labeled source datasets [Dou et al., 2019] for domain adaptation. However, digital pathology mainly relies on unlabeled or partly-labeled data as the acquisition of fully labeled cohorts is often unfeasible. In addition, recent approaches tend to treat domain adaptation as a closed-set scenario [Carlucci et al., 2019], which assumes that all target samples belong to classes present in the source domain, even though this is often not the case in a real-world scenario.

To overcome this, OSDA [Saito et al., 2018b] proposes an adversarial open-set domain adaptation approach, where the feature generator has the option to reject mistrusted or unknown target samples as an additional class. In another recent work, SSDA [Xu et al., 2019] uses self-supervised domain adaptation methods that combine auxiliary tasks, adversarial loss, and batch normalization calibration across the source and target domains.

Another domain adaptation framework DANCE [Saito et al., 2020] proposes a universal domain adaptation method to address arbitrary category shifts based on neighborhood clustering on the unlabeled target domain in a self-supervised way. Then, entropy-based optimization is utilized for feature alignment of known categories and unknown ones are

rejected, based on their entropy. The recently proposed method SENTRY [Prabhu et al., 2021] uses unsupervised domain adaptation based on selective entropy optimization, in which the target domain samples are selected based on their predictive consistency under a set of randomly augmented views. Then, SENTRY selectively optimizes the model's entropy on these samples based on their consistency to induce the domain alignment. Finally, some approaches take advantage of multiple source datasets to learn features that are discriminant under varying modalities. In Matsuura and Harada [2020], domain-agnostic features are generated by combining a domain discriminator as well as a hard clustering approach.

In this work, we propose a label-efficient framework called Self-Rule to Multi-Adapt (SRMA) for tissue type recognition in histological images and attempt to overcome the issues mentioned above by combining self-supervised learning approaches with UDA. We present an entropy-based approach that progressively learns domain invariant features, thus making our model more robust to class definition inconsistencies as well as the presence of unseen tissue classes when performing domain adaptation. SRMA is able to accurately identify tissue types in H&E stained images, which is an important step for many downstream tasks. Our proposed method achieves this by using few labeled open-source datasets and unlabeled data which are abundant in digital pathology, thus reducing the annotation workload for pathologists. We show that our method outperforms previous domain adaptation approaches in a few-label setting and highlight the potential use for clinical application in the diagnostics of CRC.

This study is an extension of the work we presented at the Medical Imaging with Deep Learning (MIDL) 2021 conference [Abbet et al., 2021]. Here, we provide a more in-depth explanation and analysis of our proposed entropy-based easy-to-hard (E2H) learning strategy. Additionally, we reformulate the entropy-based cross-domain matching used by the E2H learning strategy which improves the prediction robustness when dealing with complex tissue structures. Moreover, we also provide the generalization of the previously proposed Self-Rule to Adapt (SRA) framework to multi-source domain adaptation by including an additional public dataset and performing further experiments to assess the model's performance. Thus, we name this improved framework Self-Rule to Multi-Adapt (SRMA).

## 2 Methods

In our unsupervised domain adaptation scenario, we have access to a small set of labeled data sampled from a source domain distribution and a set of unlabeled data from a target distribution. The goal is to learn a hypothesis function (for example, a classifier) on the source domain that provides a good generalization in the target domain.

To this end, we propose a novel self-supervised cross-domain adaptation setting called SRMA, which is described in more detail below. We first introduce the architecture in a single-source setting and then present the generalization to the multi-source setting in Section 2.4. Figure 1 gives an overview of the proposed framework, and Algorithm 1 presents the pseudo-code of our SRMA method in the single-source setting.

To train our framework, we rely on a set of images $\mathcal{D} = \mathcal{D}_s \cup \mathcal{D}_t$ that is the aggregation of a set of source images $\mathcal{D}_s$ and a set of target images $\mathcal{D}_t$. The model takes as input an RGB image $\mathbf{x} \in \mathbb{R}^{H \times W \times 3}$ sampled from $\mathcal{D}$ where $H$ and $W$ denote the height and width of the image, respectively. When sampling from $\mathcal{D}$, there is an equal probability to draw a sample from either the source or the target domain. After sampling, two sets of random transformations are applied to the image $\mathbf{x}$ using an image transformer $f_T : \mathbb{R}^{H \times W \times 3} \to \mathbb{R}^{H \times W \times 3}$. This generates a pair of augmented views $\tilde{\mathbf{x}}, \tilde{\mathbf{x}}^+ \in \mathbb{R}^{H \times W \times 3}$ that are assumed to share similar content as they are both different augmentations of the same sampled input image. Each image of the pair $\tilde{\mathbf{x}}, \tilde{\mathbf{x}}^+$ is then fed to its respective encoder $f_\Phi : \mathbb{R}^{H \times W \times 3} \to \mathbb{R}^d$ and $f_\Psi : \mathbb{R}^{H \times W \times 3} \to \mathbb{R}^d$ to compute the query $\mathbf{z} \in \mathbb{R}^d$ and key $\mathbf{z}^+ \in \mathbb{R}^d$ embeddings of the input image. Here, $d$ represents the dimension of the embedding space. For notational simplicity, when sampling an image $\mathbf{x}$, we directly assume its embedding as $\mathbf{z}, \mathbf{z}^+ \in \mathcal{D}$.

Each network's branch consists of a residual encoder followed by two linear layers based on the state-of-the-art (SOTA) architecture proposed in Chen et al. [2020a] (MoCoV2). We use the key embeddings $\mathbf{z}^+$ to maintain a queue of negative samples $\mathcal{Q} = \{\mathbf{q}_l \in \mathbb{R}^d\}_{l=1}^{|\mathcal{Q}|}$ in a first-in, first-out fashion. When updating the queue with a new negative sample, not only the sampled image's embedding is stored but also its domain of origin (source or target). It allows the architecture to know at anytime the domain of origin of each queue sample.

The queue provides a large number of examples which alleviates the need for a large batch size [Chen et al., 2020b] or the use of a memory bank [Kim et al., 2020]. In addition, it enables the model to scale more easily as $\mathcal{D}$ grows as the size of the queue does not depend on it. Moreover, $f_\Psi$ is updated using a momentum approach, combining its weights with those of $f_\Phi$. This approach ensures that $f_\Psi$ generates a slowly shifting and, therefore, coherent embedding.

Motivated by Ge et al. [2020], Kim et al. [2020], Abbet et al. [2021], we extend the domain adaptation learning procedure to our model definition and task. Hence, we split the loss terms into two distinct tasks, namely the in-domain
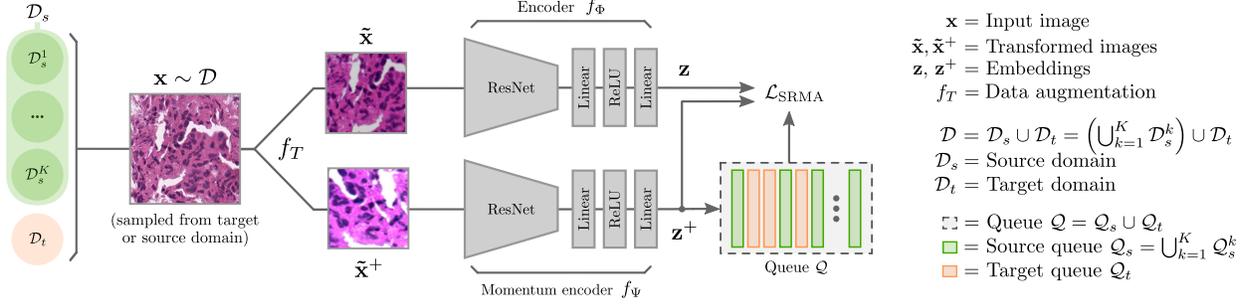
Figure 1: Schematic overview of the Self-Rule to Multi-Adapt (SRMA) framework for a given input image $\mathbf{x}$ sampled from $\mathcal{D} = \mathcal{D}_s \cup \mathcal{D}_t = \bigcup_{k=1}^K \mathcal{D}_s^k \cup \mathcal{D}_t$. Each encoder receives a different augmented version of the input image, generated by $f_T$. The loss $\mathcal{L}_{\text{SRMA}} = \mathcal{L}_{\text{IND}} + \mathcal{L}_{\text{CRD}}$ is the composition of the in-domain loss $\mathcal{L}_{\text{IND}}$ and cross-domain loss $\mathcal{L}_{\text{CRD}}$, which aims at reducing the domain gap between the source and target domains. The queue $\mathcal{Q}$ keeps track of previous samples' embeddings and their set of origin (source or target).

$\mathcal{L}_{\text{IND}}$ and cross-domain $\mathcal{L}_{\text{CRD}}$ representation learning. The objective loss $\mathcal{L}_{\text{SRMA}}$ is the summation of both terms, which are described in more detail below.

$$\mathcal{L}_{\text{SRMA}} = \mathcal{L}_{\text{IND}} + \mathcal{L}_{\text{CRD}}, \tag{1}$$

## 2.1 In-domain Loss

The first objective $\mathcal{L}_{\text{IND}}$ aims at learning the distribution of each the source and the target domain features individually. We want to keep the two domains independent as their alignment is optimized separately by the cross-domain loss term. For each embedding vector $\mathbf{z}$, there is a paired embedding vector $\mathbf{z}^+$ that is generated from the same sampled tissue image and therefore is, by definition, similar. As a result, their similarity can be jointly optimized using a contrastive learning approach [Oord et al., 2018]. Here, we strongly benefit from data augmentation to create discriminant features that match both $\mathbf{z}$ and $\mathbf{z}^+$, making them more robust to outliers. By selecting data augmentations suited to histology [Tellez et al., 2019, Faryna et al., 2021], we can ensure that the learned features are consistent with naturally occurring data variations in histology, and therefore guide the model towards histopathologically meaningful representations. This approach differs from Kim et al. [2020], where a memory bank is used instead of the combination of a queue and data augmentation to keep track of past samples. Therefore, the in-domain loss, as expressed in Equations 2-4 constrains the representation of the embedding space for each domain separately.

$$p_{\text{IND}}(\mathbf{z}, \mathbf{z}^+, \mathcal{Q}) = \frac{\exp(\mathbf{z}^\top \mathbf{z}^+ / \tau)}{\exp(\mathbf{z}^\top \mathbf{z}^+ / \tau) + \sum\limits_{\mathbf{q}_l \in \mathcal{Q}} \exp(\mathbf{z}^\top \mathbf{q}_l / \tau)}. \tag{2}$$

$$l_{\text{IND}}(\mathcal{D}, \mathcal{Q}) = \sum_{\mathbf{z}, \mathbf{z}^+ \in \mathcal{D}} \log \left[ p_{\text{IND}}(\mathbf{z}, \mathbf{z}^+, \mathcal{Q}) \right]. \tag{3}$$

$$\mathcal{L}_{\text{IND}} = \frac{-1}{|\mathcal{D}_s| + |\mathcal{D}_t|} \left[ l_{\text{IND}}(\mathcal{D}_s, \mathcal{Q}_s) + l_{\text{IND}}(\mathcal{D}_t, \mathcal{Q}_t) \right]. \tag{4}$$

We denote $\mathcal{Q}_s, \mathcal{Q}_t \subset \mathcal{Q}$ as the sets of indexed samples of the queue that were previously drawn from the corresponding domain $\mathcal{D}_s, \mathcal{D}_t \subset \mathcal{D}$, and $\tau \in \mathbb{R}$ as the temperature. The temperature is typically small ($\tau \ll 1$), thus sharpening the signal and helping the model to make confident predictions. For all images of each dataset $\mathcal{D}_s, \mathcal{D}_t$, we want to minimize the distance between $\mathbf{z}$ and $\mathbf{z}^+$ while maximizing the distance to the previously generated negative samples from the corresponding sets $\mathcal{Q}_s, \mathcal{Q}_t$. The samples in the queue are considered reliable negative candidates as they are generated by $f_\Psi$ whose weights are slowly optimized due to its momentum update procedure.

## 2.2 Cross-domain Loss

We can see the cross-domain matching task as the generation of features that are discriminative across both sets. In other words, two samples that are visually similar but are drawn from the source $\mathcal{D}_s$ and target $\mathcal{D}_t$ domain, respectively,

---

**Algorithm 1:** Pseudocode for the single-source Self-Rule to Multi-Adapt (SRMA) framework

---

Initialize queue $\mathcal{Q}$ by sampling from normal distribution $\mathcal{N}(0, 1)$;
Normalize queue entries $\{q_i\} \in Q$;
**for** e = 0 **to** $N_{\text{epochs}} - 1$ **do**
    Create $\mathcal{D}$ by uniformly sampling from $\mathcal{D}_s$ and $\mathcal{D}_t$;
    Update easy-to-hard coefficient $r$ using Equation 10;
    **for** batch $\{\mathbf{x}_i\}_{i=1}^{B}$ in $\mathcal{D}$ **do**
        Get augmented samples $\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_i^+$ using $f_T$;
        Perform forward pass $\mathbf{z}_i = f_\phi(\tilde{\mathbf{x}}_i)$, $\mathbf{z}_i^+ = f_\psi(\tilde{\mathbf{x}}_i^+)$ ;
        Normalize vectors $\mathbf{z}_i, \mathbf{z}_i^+$ ;
        Compute in-domain loss $\mathcal{L}_{\text{IND}}$ using Equation 4;
        Calculate cross-entropy $\bar{H}$ using Equation 7 ;
        Compute easy-to-hard $\mathcal{R}_s, \mathcal{R}_t$ sets using Equation 11 ;
        Evaluate cross-domain loss $\mathcal{L}_{\text{CRD}}$ by replacing $\mathcal{D}_s, \mathcal{D}_t$ with $\mathcal{R}_s, \mathcal{R}_t$ in Equation 9, respectively;
        Compute $\mathcal{L}_{\text{SRA}} = \mathcal{L}_{\text{IND}} + \mathcal{L}_{\text{CRD}}$;
        Update $f_\Phi$ weights with backpropagation ;
        Update $f_\Psi$ weights with momentum ;
        Update queue $\mathcal{Q}$ with $\mathbf{z}_i^+$;
    **end**
**end**

---

should have a similar embedding. On the other hand, when comparing these samples to the remaining candidates of the opposite domain, their resulting embeddings should be far apart. Practically, performing cross-domain matching using the number of available candidates within a batch might deteriorate the quality of the domain matching process due to the limited amount of negative samples. Therefore, we use the queue to find negative samples for domain matching. Hence, we compute the similarity and cross-entropy of each query pair $\mathbf{z}, \mathbf{z}^+$ drawn from one set (for example $\mathcal{D}_s$) to the stored queue samples from the other set (for example $\mathcal{Q}_t$):

$$p_{\text{CRD}}(\mathbf{z}, \mathbf{q}, \mathcal{Q}) = \frac{\exp(\mathbf{z}^\top \mathbf{q}/\tau)}{\sum\limits_{\mathbf{q}_l \in \mathcal{Q}} \exp(\mathbf{z}^\top \mathbf{q}_l/\tau)}, \tag{5}$$

$$H(\mathbf{z}, \mathbf{z}^+, \mathcal{Q}) = -\sum_{\mathbf{q} \in \mathcal{Q}} p_{\text{CRD}}(\mathbf{z}, \mathbf{q}, \mathcal{Q}) \log \left[ p_{\text{CRD}}(\mathbf{z}^+, \mathbf{q}, \mathcal{Q}) \right], \tag{6}$$

$$\bar{H}(\mathbf{z}, \mathbf{z}^+, \mathcal{Q}) = \frac{1}{2} \left[ H(\mathbf{z}, \mathbf{z}^+, \mathcal{Q}) + H(\mathbf{z}^+, \mathbf{z}, \mathcal{Q}) \right]. \tag{7}$$

A low cross-entropy $H$ means that the selected query pair $\mathbf{z}, \mathbf{z}^+$ from one domain matches with a limited number of samples from another domain. Moreover, we update our initial definition of $H$ [Abbet et al., 2021], where solely $\mathbf{z}$ is used. By taking the average cross-entropy $\bar{H}$, the model is now also penalized when the predictions from $\mathbf{z}, \mathbf{z}^+$ of the same image are different. This improves the consistency of the domain matching[Assran et al., 2021]. As a result, the loss $\mathcal{L}_{\text{CRD}}$ aims to minimize the averaged cross-entropy of the similarity distributions, assisting the model in making confident predictions:

$$l_{\text{CRD}}(\mathcal{D}, \mathcal{Q}) = \sum_{\mathbf{z}, \mathbf{z}^+ \in \mathcal{D}} \bar{H}(\mathbf{z}, \mathbf{z}^+, \mathcal{Q}), \tag{8}$$

$$\mathcal{L}_{\text{CRD}} = \frac{1}{|\mathcal{D}_s| + |\mathcal{D}_t|} \left[ l_{\text{CRD}}(\mathcal{D}_s, \mathcal{Q}_t) + l_{\text{CRD}}(\mathcal{D}_t, \mathcal{Q}_s) \right]. \tag{9}$$

## 2.3 Easy-to-hard Learning

There are two main pitfalls that can hamper the performance of the cross-domain entropy minimization.
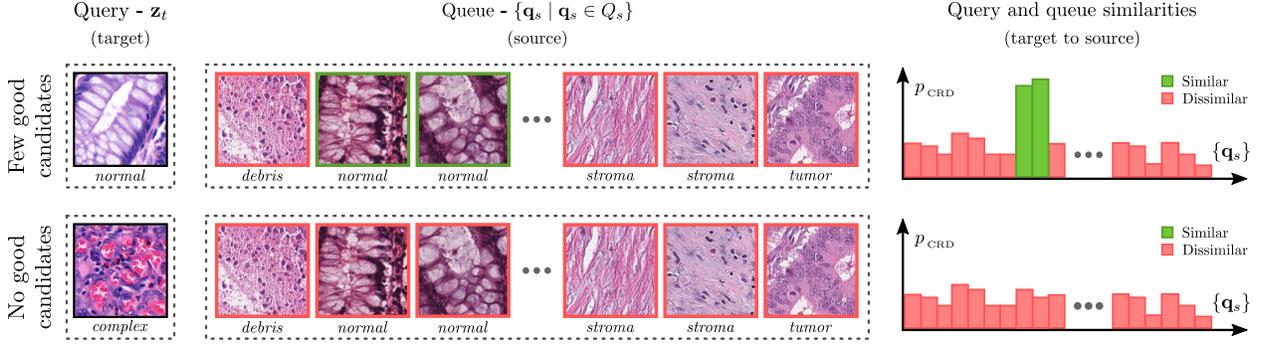
Figure 2: Toy example of the cross-domain matching of different target queries to a fixed source queue. The first column shows two example target query images with computed embedding $\mathbf{z}_t$. The second column depicts the source queue images maintained by the model and their corresponding embeddings $\{\mathbf{q}_s\}$ In the third column, the distribution of the computed similarities $p_{\mathrm{CRD}}$ between the queries and each queue sample are plotted. Similar and dissimilar patterns with respect to the query are displayed in green and red. The top row highlights the case where the model is able to find at least a subset of elements of the queue that match the query (low entropy), as opposed to the bottom row where none of the queue samples match the presented query (high entropy). The class labels in this figure have been added for ease of reading and are not available during training.

Firstly, at the start of the learning process, the similarity measure between samples and the queue is unclear as the model weights are initialized randomly, which does not guarantee proper feature descriptions. As a result, the optimization of their relative entropy and the loss term $\mathcal{L}_{\mathrm{CRD}}$ is ambiguous in the first epochs.

Secondly, being able to find matching samples for all input queries across datasets is a strong assumption. In clinical applications, we often rely on open-source datasets with a limited number of classes to annotate complex tissue databases. More specifically, challenging tissue types such as complex stroma subtypes are often not present in public datasets while being frequent in the WSIs encountered in daily diagnostics. This example is illustrated in Figure 2. The top row shows the case where for a given target query $\mathbf{z_t}$ there are samples with a similar pattern in the source queue, i.e., the distribution of similarities $p_{\mathrm{CRD}}$ has low entropy. The second row highlights the opposite scenario where no queue elements match the query, generating a quasi-uniform distribution of similarities and, therefore, a high entropy. In other words, optimizing Equation 7 for all samples will result in a performance drop as the loss will try to find cross-domain candidates even if there are none to be found.

To tackle both of these issues, we introduce an easy-to-hard (E2H) learning scheme. The model starts with easy to match samples (low cross-entropy) samples and progressively includes harder (high cross-entropy) samples as the training progresses. We assume that the model becomes more robust after each iteration and is more likely to properly process harder examples in later stages. Formally, we substitute the domains $\mathcal{D}_s, \mathcal{D}_t$ in Equation 9 with the corresponding set of candidates $\mathcal{R}_s, \mathcal{R}_t$ defined as:

$$r = \left\lfloor \frac{e}{N_{\mathrm{epochs}} \cdot s_w} \right\rfloor \cdot s_h, \tag{10}$$

$$\begin{aligned}
\mathcal{R}_s &= \{\mathbf{z}_s, \mathbf{z}_s^+ \in \mathcal{D}_s \mid \bar{H}(\mathbf{z}_s, \mathbf{z}_s^+, \mathcal{Q}_t) \text{ is reverse top-}r\}, \\
\mathcal{R}_t &= \{\mathbf{z}_t, \mathbf{z}_t^+ \in \mathcal{D}_t \mid \bar{H}(\mathbf{z}_t, \mathbf{z}_t^+, \mathcal{Q}_s) \text{ is reverse top-}r\},
\end{aligned} \tag{11}$$

where the ratio $r$ is gradually increased during training using a step function. We denote $s_w, s_h$ as the width and height of the step, respectively, $N_{\mathrm{epochs}}$ as the total number of epochs, and $e$ the current epoch. The term reverse top-r indicates the ranking of cross-entropy terms in reverse order (low to high values). For example, $r = 0.2$ will capture the top 20% of the samples with the lowest cross-entropy. This definition ensures that as long as $r = 0$ (i.e. $e < N_{\mathrm{epochs}} \cdot s_w$) we only use the in-domain loss $\mathcal{L}_{\mathrm{IND}}$ for backpropagation, and the cross-domain loss term $\mathcal{L}_{\mathrm{CRD}}$ is not considered. This allows us to first only learn feature representations based on the in-domain feature distribution. Moreover, with the tuning of the parameter $s_h$ we can control the range of $r$ and thus ensure that its value never reaches $r = 1$ to avoid systematic cross-domain matching where no candidates are available.

6

(a) In-domain optimization scenarios ($\mathcal{L}_{\text{IND}}$)  (b) Cross-domain matching scenarios ($\mathcal{L}_{\text{CRD}}$)
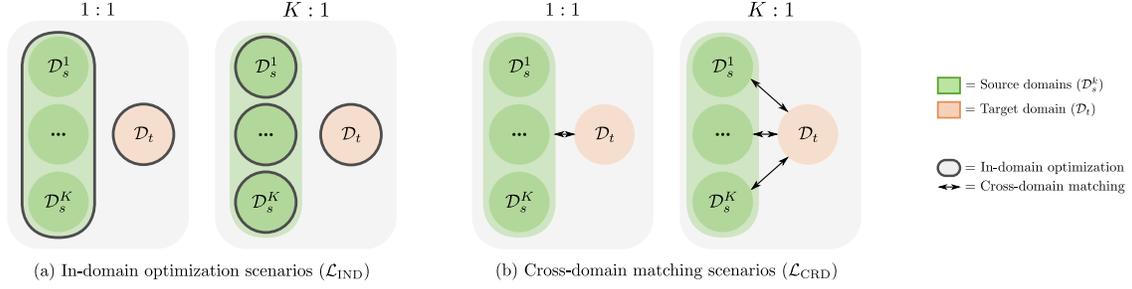
Figure 3: Proposed multi-source scenarios for the in-domain $\mathcal{L}_{\text{IND}}$ (a) and cross-domain $\mathcal{L}_{\text{CRD}}$ (b) optimization. With the one-to-one settings ($1 : 1$), we treat all source sets $\mathcal{D}_s^k$ as a single set $\mathcal{D}_s$. In the K-to-one ($K : 1$) setting, each source domain is considered as an independent set. Note that there are no restrictions regarding the combination of the loss terms. For example, the source set can be considered as a single set for the in-domain optimization while being considered as multiple sets for the cross-domain matching.

## 2.4   Generalization to Multiple Source Scenario

Our proposed SRMA framework can be generalized to consider multiple datasets in the source domain. This is especially useful if the available source datasets overlap in terms of class definitions, which increases the diversity of the visual appearance in the source data. More formally, we rely on $K$ source datasets denoted $\mathcal{D}_s^k$ where $\bigcup_{k=1}^{K} \mathcal{D}_s^k = \mathcal{D}_s$, and $\mathcal{D} = \mathcal{D}_s \cup \mathcal{D}_t$. The same is valid for the source queues $\mathcal{Q}_s^k$ where $\bigcup_{k=1}^{K} \mathcal{Q}_s^k = \mathcal{Q}_s$, and $\mathcal{Q} = \mathcal{Q}_s \cup \mathcal{Q}_t$. For both the in-domain and cross-domain loss we present two multi-source scenarios as depicted in Figure 3.

One option is to consider the whole source domain as a single domain $\mathcal{D}_s = \bigcup_{k=1}^{K} \mathcal{D}_s^k$ for the in-domain loss:

$$\mathcal{L}_{\text{IND}}^{1:1} = \frac{-1}{|\mathcal{D}_s| + |\mathcal{D}_t|} \left[ l_{\text{IND}}(\bigcup_{k=1}^{K} \mathcal{D}_s^k, \bigcup_{k=1}^{K} \mathcal{Q}_s^k) + l_{\text{IND}}(\mathcal{D}_t, \mathcal{Q}_t) \right]. \tag{12}$$

Here, we make no distinction between the source sets and consider a one-to-one features representation importance ($1 : 1$) between the source and target domain. This definition is equivalent to the single source in-domain adaptation.

Alternatively, we can consider each source and the target domain as independent sets as in Equation 13. With this K-to-one ($K : 1$) scenario, we have $K + 1$ separate in-domain optimizations:
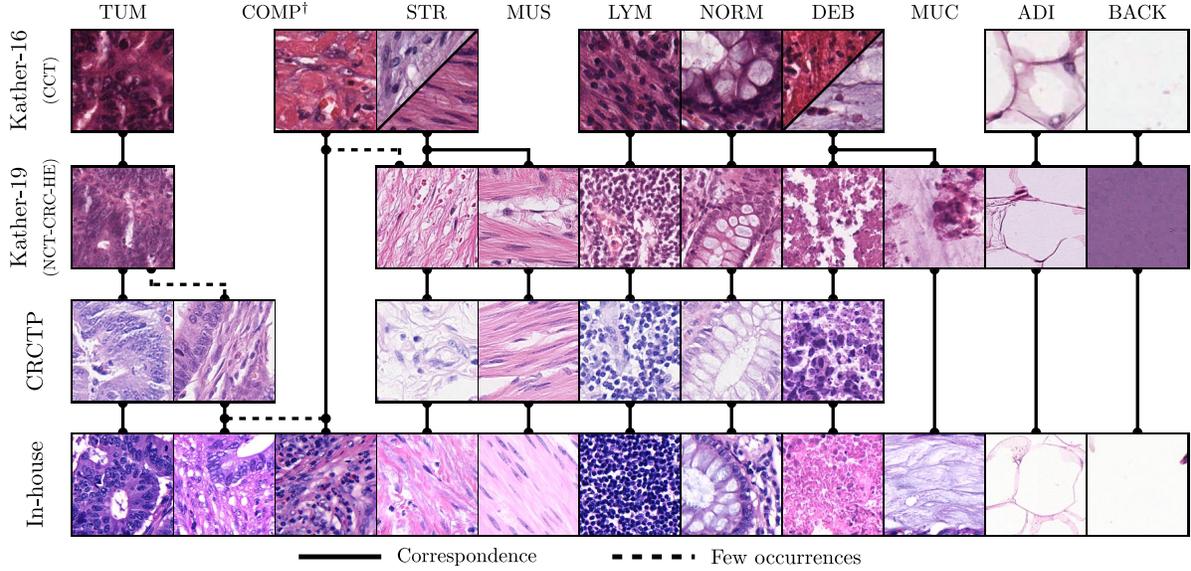
$$\mathcal{L}_{\text{IND}}^{K:1} = \frac{-1}{|\mathcal{D}_s| + |\mathcal{D}_t|} \left[ \sum_{k=1}^{K} l_{\text{IND}}(\mathcal{D}_s^k, \mathcal{Q}_s) + l_{\text{IND}}(\mathcal{D}_t, \mathcal{Q}_t) \right]. \tag{13}$$

The same logic applies to the cross-domain matching. We can either consider a one-to-one correspondence between the unified source domain and the target domain as in Equation 14, or match each of the individual source domains to the target as in Equation 15.

$$\mathcal{L}_{\text{CRD}}^{1:1} = \frac{-1}{|\mathcal{D}_s| + |\mathcal{D}_t|} \left[ l_{\text{CRD}}(\bigcup_{k=1}^{K} \mathcal{D}_s^k, \mathcal{Q}_t) + l_{\text{CRD}}(\mathcal{D}_t, \bigcup_{k=1}^{K} \mathcal{Q}_s^k) \right]. \tag{14}$$

$$\mathcal{L}_{\text{CRD}}^{K:1} = \frac{-\frac{1}{K}}{|\mathcal{D}_s| + |\mathcal{D}_t|} \sum_{k=0}^{K-1} \left[ l_{\text{CRD}}(\mathcal{D}_s^k, \mathcal{Q}_t) + l_{\text{CRD}}(\mathcal{D}_t, \mathcal{Q}_s^k) \right]. \tag{15}$$

The formulation of the E2H learning procedure has to be updated to comply with multi-source domain definition. For the one-to-one setting, sets $\mathcal{R}_s$, $\mathcal{R}_t$ remain unchanged as we make no distinction between the different source sets. However, for the K-to-one setting, the model seeks to match the target domain to the source domain without taking into

Figure 4: Example images of the different tissue types present in the used datasets and their association. The labeled datasets Kather-16 (K16), Kather-19 (K19), and Colorectal Cancer Tissue Phenotyping (CRC-TP) are publicly available. Examples from the in-house dataset are manually picked for comparison but are not labeled. We use the following abbreviations: TUM: tumor epithelium, STR: simple stroma, COMP: complex stroma, LYM: lymphocytes, NORM: normal mucosal glands, DEB: debris/necrosis, MUS: muscle, MUC: mucus, ADI: adipose tissue, BACK: background. The solid and dashed lines indicate classes correspondences and reported overlaps (also see Section 3.5).

consideration that there are multiple available source domains. We replace the domains $\mathcal{D}_s^k, \mathcal{D}_t$ in Equation 15 with the corresponding set of candidates $\mathcal{R}_s^k, \mathcal{R}_t$ defined as:

$$\mathcal{R}_s^k = \{\mathbf{z}_s, \mathbf{z}_s^+ \in \mathcal{D}_s^k \mid \bar{H}(\mathbf{z}_s, \mathbf{z}_s^+, \mathcal{Q}_t) \text{ is reverse top-}r\},$$
$$\mathcal{R}_t = \{\mathbf{z}_t, \mathbf{z}_t^+ \in \mathcal{D}_t \mid \bar{H}(\mathbf{z}_t, \mathbf{z}_t^+, \mathcal{Q}_s^k) \text{ is reverse top-}r\}.$$

(16)

The overall loss $\mathcal{L}_{\mathrm{SRMA}}$ for the multi-source setting is the combination of the in-domain loss ($\mathcal{L}_{\mathrm{IND}}^{1:1}$ or $\mathcal{L}_{\mathrm{IND}}^{K:1}$) and the cross-domain loss ($\mathcal{L}_{\mathrm{CRD}}^{1:1}$ or $\mathcal{L}_{\mathrm{CRD}}^{K:1}$).

# 3   Datasets

In this study, we use three publicly available datasets, Kather-16 (K16), Kather-19 (K19) and Colorectal Cancer Tissue Phenotyping (CRC-TP), that contain patches extracted from H&E-stained WSIs of different tissue types found in the human gastrointestinal tract. We also use an in-house CRC cohort, which does not have patch-level labels, and evaluate our method on three regions of interest (ROIs). More details on the datasets can be found below.

Figure 4 shows the occurrence and relationship of different tissue types across all four datasets. The displayed crops of the in-house WSI datasets are cherry-picked for comparison purposes.

## 3.1   Kather-16 Dataset

The K16 dataset [Kather et al., 2016] contains $5,000$ patches ($150 \times 150$ pixels, $74\mu m \times 74\mu m$) from multiple H&E WSIs. All images are digitized using a scanner magnification of 20x ($0.495\mu m$ per pixel). There are eight classes of tissue phenotypes, namely tumor epithelium, simple stroma (homogeneous composition, and smooth muscle), complex stroma (stroma containing single tumor cells and/or few immune cells), immune cells, debris (including necrosis, erythrocytes, and mucus), normal mucosal glands, adipose tissue, and background (no tissue). The dataset is balanced with 625 patches per class.

### 3.2 Kather-19 Dataset

The K19 dataset [Kather et al., 2019] consists of patches depicting nine different tissue types: cancerous tissue, stroma, normal colon mucosa, adipose tissue, lymphocytes, mucus, smooth muscle, debris, and background. Each class is roughly equally represented in the dataset. In total, there are $100,000$ patches ($224 \times 224$ pixels, $112\mu m \times 112\mu m$) in the training set. All images are digitized using a scanner at a magnification of 20x ($0.5\mu m$ per pixel).

### 3.3 Colorectal Cancer Tissue Phenotyping Dataset

The CRC-TP [Javed et al., 2020] dataset contains a total of $196,000$ patches depicting seven different tissue phenotypes (tumor, inflammatory, stroma, complex stroma, necrotic, benign, and smooth muscle). The different phenotypes are roughly equally represented in the dataset. For tumor, complex stroma, stroma, and smooth muscle, there are 35,000 patches per class, for benign and inflammatory, there are 21,000, and for debris, there are 14,000. The patches ($150 \times 150$ pixels) are extracted at 20x resolution from 20 H&E WSIs, each one coming from a different patient. For each class, only a subset of the WSIs is used to extract the patches. The annotations are made by two expert pathologists. Out of the two dataset splits available, we use the training set of the patient-level split.

### 3.4 In-house Dataset

Our cohort is composed of 665 H&E-stained WSIs from our local CRC patient cohort at the Institute of Pathology, University of Bern, Switzerland. The slides originate from 378 unique patients diagnosed with adenocarcinoma and are scanned at a resolution of $0.248\mu m$ per pixel (40x). None of the selected slides originated from patients that underwent preoperative treatment.

From each WSI we uniformly sample 300 ($448 \times 448$ pixels, $111\mu m \times 111\mu m$) regions from the foreground masks to reduce the computational complexity of the proposed approach. This creates a dataset with a total of $199,500$ unique and unlabeled patches. We assume that these randomly selected samples are a good estimation of the tissue complexity and heterogeneity of our cohort.

We also select three ROIs of size $5 \times 5mm$ ($\simeq 20,000 \times 20,000$ pixel), which are annotated by an expert pathologist according to the definitions used in the K19 dataset, and use them for evaluation. The regions are selected such that, overall, they represent all tissue types, as well as challenging cases such as late cancer stage (ROI 1), mucinous carcinoma (ROI 2), and torn tissue (ROI 3).

### 3.5 Discrepancies in Class Definitions Between Datasets

The class definitions are not homogeneous across the datasets and they also do not contain the same number of tissue classes. Following a discussion with expert pathologists, we group stroma/muscle and debris/mucus as stroma and debris, respectively, to create a corresponding adaptation between K19 and K16.

Moreover, the complex stroma class definition between K16 and CRC-TP is not identical. The CRC-TP complex stroma class contains tiles from the tumor border region and is more consistent with the tumor class in the K16 and K19 dataset. In K16, the complex stroma is not limited to the tumor border surroundings and is defined as the desmoplastic reaction area, which is usually composed of a mixture of debris, lymphocytes, single tumor cells, and tumor cell clusters.

As a result, the complex stroma class is kept for training but excluded from the evaluation process when performing adaptation on K16 and CRC-TP. With this problem definition, we fall into an open-set scenario where the class distribution of the two domains does not rigorously match, as opposed to a closed set adaptation scheme.

## 4 Results and Discussion

In this section, we present and discuss the experimental results. The general experimental setup is described in Section 4.1. We validate our proposed self-supervised domain adaptation approach using publicly available datasets and compare it to current SOTA methods for UDA in Section 4.2. Additionally, we assess the performance in a clinically relevant use case by validating our model on WSI sections from our in-house cohort in Section 4.3. We perform an ablation study in Section 4.4 for the single-source setting as well as additional experiments on the importance of the E2H learning procedure in Section 4.5. These experiments are further extended to a multi-source application in Section 4.6-4.7 on both publicly available datasets and WSI sections. To help future research, the implementation and trained models are available open-source[2].

---

[2]Code available on `https://github.com/christianabbet/SRA`.

## 4.1 General Experimental Setup

In this section, we present the general setup that is used in all experiments. First, the architecture is trained in an unsupervised fashion, and in a second step, a linear classifier is trained on top as described by Chen et al. [2020b].

For the unsupervised learning step, the architecture of the feature extractors, $f_\Phi$ and $f_\Psi$, are composed of a ResNet18 [He et al., 2016] followed by two fully connected layers (projection head) using rectified linear activation units (ReLUs). The output dimension of the multi-layer projection head is $d = 128$. We update the weights of $f_\Phi$ as $\theta_\Phi$ using standard backpropagation and $f_\Psi$ as $\theta_\Psi$ with momentum $m = 0.999$, as described in He et al. [2020].

The model is trained from scratch for $N_{\text{epochs}} = 200$ epochs until convergence using the stochastic gradient descent (SGD) optimizer (momentum $= 0.9$, weight decay $= 10^{-4}$), a learning rate of $\lambda = 0.03$, and a batch size of $B = 128$. The size of the queue is fixed to $|\mathcal{Q}| = 2^{16} = 65,536$ samples. For the similarity learning we set $\tau = 0.2$. We apply random cropping, grayscale transformation, horizontal/vertical flipping, rotation, grid distortion, ISO noise, Gaussian noise, and color jittering as data augmentations $f_T$. At each epoch, we sample $50,000$ examples with replacement from both the source and target dataset to create $\mathcal{D}$ with a total of $N = 100,000$ samples. The ratio $r$ is updated between each epoch, while the sets $\mathcal{R}_s$, $\mathcal{R}_t$ for cross-domain matching are computed batch-wise.

During the second phase, the momentum encoder branch is discarded as it is not used for inference. The classification performance is evaluated using a linear classifier, which is placed on top of the frozen feature extractor. The linear classifier directly matches the output of the embedding $d$ to the number of classes. It is trained for $N_{\text{epochs}} = 100$ epochs until convergence using the SGD optimizer (momentum $= 0.9$, weight decay $= 0$), a batch size of $B = 128$, and a learning rate of $\lambda = 1$. We use only few randomly selected source labels to train this classification layer in order to simulate the clinical application, where we usually rely on a large quantity of unlabeled data and only have access to few labeled samples. More precisely, we use $n = 1,000$ samples (i.e., $1\%$) to train the linear classifier with K19 and $n = 500$ samples (i.e., $10\%$) when training with K16. While training the linear classifier, we multi-run 10 times to obtain statistically relevant results. The set of selected source labels varies between these runs, as they are randomly sampled for each run. If not specified otherwise, we use $s_w = 0.25$ and $s_h = 0.15$ for E2H learning.

For a fair comparison, we also use a ResNet18 backbone for all the presented baselines.

## 4.2 Cross-Domain Patch Classification

In this task, we use the larger dataset K19 as the source dataset and adapt it to K16. We motivate the selection of K19 as the source set by the fact that it is closer to the clinical scenario where we mainly rely on a large quantity of unlabeled data and only a few labeled ones, by using only $1\%$ of the labels in K19. We evaluate the performance of the model with the patch classification task on the K16 dataset. The mucin and muscle in K19 are grouped with debris and stroma, respectively, to allow comparison with the K16 class definitions. We use $70\%$ of K16 to train the unsupervised domain adaptation. The remaining $30\%$ are used to test the performance of the linear classifier trained on top of the self-supervised model.

The results of our proposed SRMA method are presented in Table 1, in comparison with baselines and SOTA algorithms for domain adaption. As the lower bound, we consider both MoCoV2 where the source and the target domain are merged into a single set and direct transfer learning (source only), where the model is trained in a supervised fashion on the source data only. We use the same logic for the upper bound by training on the target domain data (fully supervised approach). The performances on complex stroma are not reported as the class is not present in K19. Figure 5 shows the t-SNE projection and alignment of the domain adaptation for the transfer learning (source only), the top-performing baselines (OSDA, SSDA with jigsaw solving), and our method (SRMA). Complementary results can be found in A and B.

MoCoV2 fails to generalize knowledge between source and target domain. Since the model is not constrained, it learns two distinct embeddings for each domain. The experiment highlights the limitations of contrastive learning without domain adaptation.

Stain normalization slightly decreases the performances, compared to the source only baseline, as it introduces color artifacts that are very challenging for the network classifier. This mainly comes from the distribution of target samples, namely K16, that are composed of dark stained patches which are difficult to normalize properly.

CycleGAN suffers from performance degradation for the lymphocytes predictions. Like color normalization, it tends to create saturated images. In addition, the model alters the shape of the lymphocytes nuclei, thus fooling the classifier toward either debris or tumor classification.

Table 1: Results of the domain adaptation from K19 (source) to K16 (target). 1% of the source domain labels are used and the target domain labels are unknown. Complex stroma is excluded as the class is not present in K19. The mucin and muscle class in K19 are grouped with debris and stroma, respectively, as they overlap in K16. The top results for the domain adaptation methods are highlighted in bold. We report the F1 score for each class as well as the overall weighted F1 score averaged over 10 runs.

| Methods | TUM | COMP | STR | LYM | DEB | NORM | ADI | BACK | ALL |
|---|---|---|---|---|---|---|---|---|---|
| MoCoV2 [Chen et al., 2020a][†] | 36.8[**] | - | 45.4[**] | 27.1[**] | 30.8[**] | 45.2[**] | 43.1[**] | 43.6[**] | 38.9[**] |
| Source only[‡] | 74.0[**] | - | 77.4[**] | 75.3[**] | 50.5[**] | 66.9[**] | 87.0[**] | 93.1[**] | 75.1[**] |
| DANN [Ganin and Lempitsky, 2015] | 65.8[**] | - | 60.8[**] | 42.3[**] | 47.8[**] | 61.9[**] | 64.1[**] | 62.3[**] | 57.8[**] |
| Stain norm. [Macenko et al., 2009] | 77.8[**] | - | 75.9[**] | 68.2[**] | 42.1[**] | 75.1[**] | 77.4[**] | 87.6[**] | 72.2[**] |
| CylceGAN [Zhu et al., 2017] | 70.7[**] | - | 71.6[**] | 62.3[**] | 47.6[**] | 75.5[**] | 89.0[**] | 88.2[**] | 72.4[**] |
| SelfPath [Koohbanani et al., 2021] | 71.5[**] | - | 68.8[**] | 68.1[**] | 57.6[**] | 77.6[**] | 82.3[**] | 85.5[**] | 73.1[**] |
| OSDA [Saito et al., 2018b] | 82.0[**] | - | 78.2[*] | 83.6[*] | 63.8[**] | 80.3[**] | 90.8[**] | 93.2[*] | 81.7[**] |
| SSDA - Rot [Xu et al., 2019] | 85.1[**] | - | 78.5[**] | 81.3[**] | **68.2** | 88.7[**] | 93.9[**] | 96.5[**] | 84.7[**] |
| SSDA - Jigsaw Xu et al. [2019] | 90.0[**] | - | **81.2** | 79.5[**] | 64.4[**] | 88.3[**] | 94.2[**] | 95.7[*] | 84.9[**] |
| SENTRY [Prabhu et al., 2021] | 88.7[**] | - | 74.4[**] | **86.0** | 65.5[+] | 91.5[**] | 94.1[**] | **97.9**[+] | 85.7[**] |
| SRA [Abbet et al., 2021] | 93.4[**] | - | 72.9[**] | 82.7[*] | 67.9[+] | 96.5[*] | **97.0**[+] | 97.2[+] | 86.9[*] |
| SRMA (ours) | **97.3** | - | 79.3[+] | 80.2[**] | 62.2[**] | **98.7** | **97.6** | **98.1** | **87.7** |
| Target only[§] | 94.6[**] | - | 83.6[**] | 92.6[**] | 88.7[**] | 95.4[**] | 97.8[+] | 98.5[+] | 93.0[**] |

[†] Source and target dataset are merged and trained using contrastive learning.
[‡] Direct transfer learning: trained on the source domain only, no adaptation (lower bound).
[§] Fully supervised: trained knowing the labels of the target domain (upper bound).
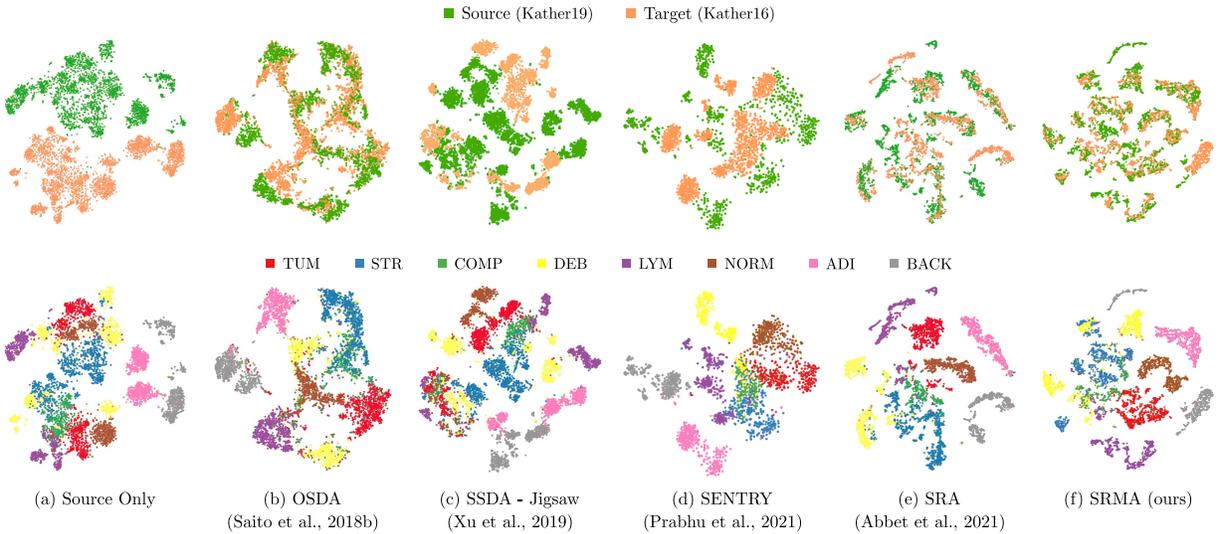[+] $p \geq 0.05$; [*] $p < 0.05$; [**] $p < 0.001$; unpaired t-test with respect to the top result.



Figure 5: The t-SNE projection of the source (K19) and target (K16) domain embeddings. The top row shows the alignment between the source and target domain, while the bottom row highlights the representations of the different classes. We compare our approach (f) to other UDA methods (b-e), and the fully supervised, transfer learning baseline (source only) (a).

In our setup, we observe that the use of the gradient reversal layer leads to an unstable loss optimization for both Self-Path and DANN, which explains the large performance drops when training. Heavier data augmentations partially solve this issue.

OSDA benefits from the open-set definition of the approach and achieves very good performance for lymphocytes detections.
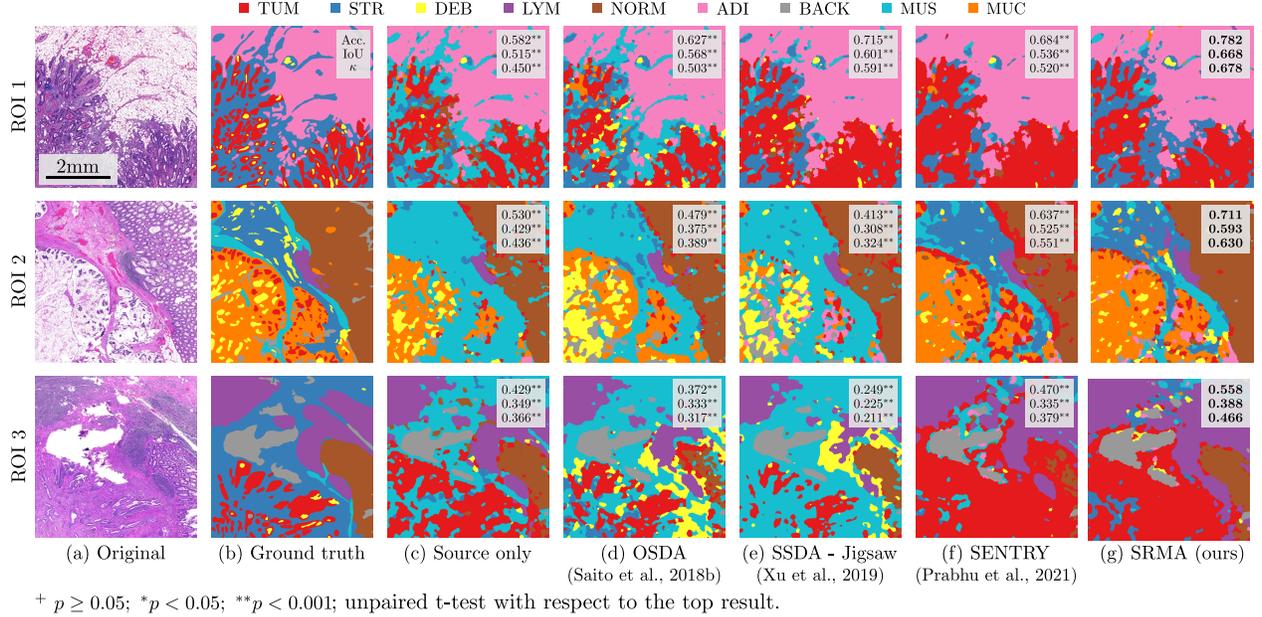
**Figure 6:** Quantitative results of the domain adaptation from K19 to our unlabeled in-house dataset based on three selected regions of interest (ROIs). (a-b) show the original ROIs from the WSIs and their ground truth, respectively. We compare the performance of our Self-Rule to Multi-Adapt (SRMA) algorithm (g) to the lower bound and the top-performing SOTA methods (c-f). We report the pixel-wise accuracy, the weighted intersection over union, and the pixel-wise Cohen's kappa ($\kappa$) score averaged over 10 runs.

SSDA achieves similar results when using either rotation or jigsaw puzzle-solving as an auxiliary task. Due to the rotational invariance structure of the tissue and selected large magnification for tilling, rotation and jigsaw puzzle-solving are not optimal auxiliary tasks for digital pathology.

Out of the presented baselines, SENTRY achieves top competitive results on almost all classes. The main limitation appears to be the distinction between tumor and normal mucosa.

Our proposed SRMA method shows an excellent alignment between the same class clusters of the source and target distributions and outperforms SOTA approaches in terms of weighted F1 score. Notably, our approach is even able to match the upper bound model for normal and tumor tissue identification. The embedding of complex stroma, which only exists in the target domain, is represented as a single cluster with no matching candidates, which highlights the model's ability to reject unmatchable samples from domain alignment.

Furthermore, the cluster representation is more compact compared to other presented methods, where for example, normal mucosa tends to be aligned with complex stroma and tumor. Our approach suffers a drop in performance for stroma detection, which can be explained by the presence of lymphocytes in numerous stroma tissue examples, causing a higher rate of misclassification. Moreover, the presence of loose tissue that has a similar structure as stroma in the debris class is challenging. The overlap is also observed in the embedding projection.

### 4.3   Use Case: Cross-Domain Segmentation of WSIs

To further validate our approach in a real case scenario, we perform domain adaptation using our proposed model from K19 to our in-house dataset and validate it on WSIs regions of interest (ROIs).

The results are presented in Figure 6, alongside the original H&E ROIs, their corresponding ground truth annotations, direct transfer learning (source only), as well as comparative results of the top-scoring SOTA approaches. We use a tile-based approach to predict classes on each ROIs and use conditional random fields as in Chan et al. [2019] to smooth the prediction map. The number of available labeled tissue regions is limited to the presented ROIs.

For all models, stroma and muscle are poorly differentiated as both have similar visual features without contextual information. This phenomenon is even more apparent in the source only setting, where muscle tissue is almost systematically interpreted as stroma. Moreover, due to the lack of domain adaptation, the boundary between tumor and normal tissues is not well defined, leading to incorrect predictions of these classes.

(a) t-SNE visualization of target patches distribution

(b) t-SNE labeled source samples projection

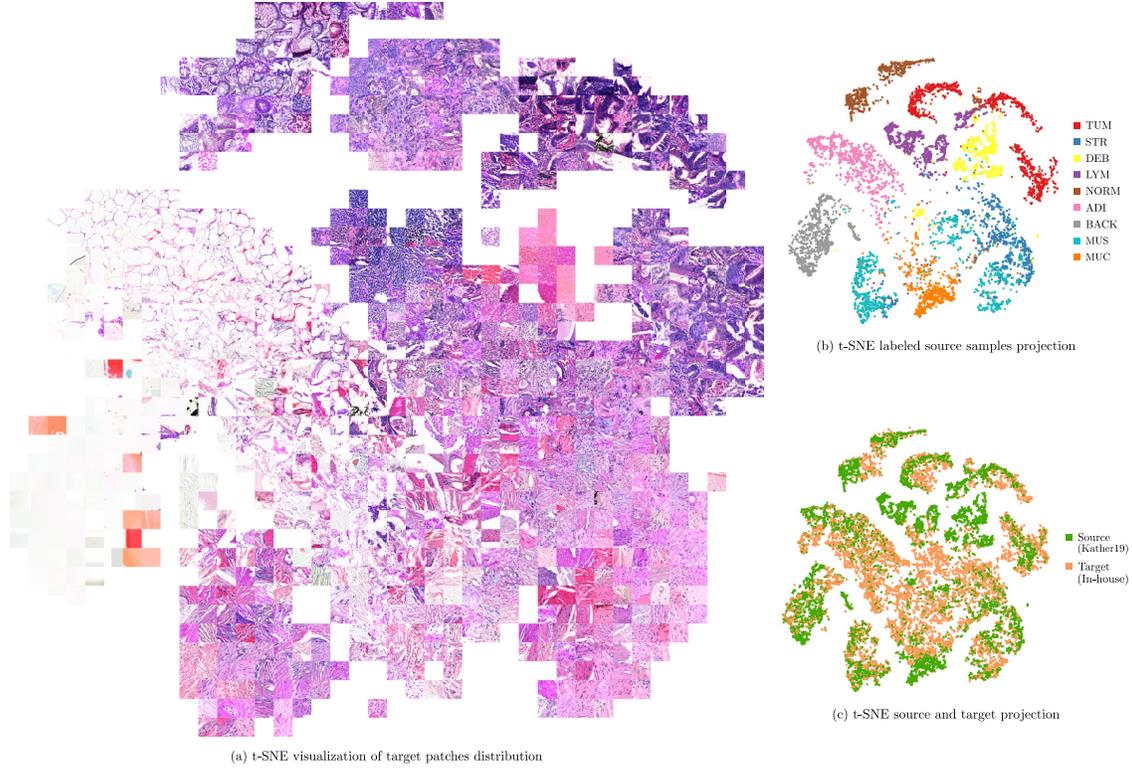(c) t-SNE source and target projection

Figure 7: The t-SNE visualization of the SRMA model trained on K19 and our in-house data. All sub-figures depict the same embedding. (a) Patch-based visualization of the embedding. (b) Distribution of the labeled source samples. (c) The relative alignment of the source and target domain samples.

OSDA, on the other hand, fails to adapt and generalize to new tumor examples while trying to reject mistrusted samples. This phenomenon is most visible in ROI 3, where the model interprets the surroundings of the cancerous region as a mixture of debris, stroma, and muscle.

SSDA tends to predict lymphocyte aggregates as debris. This can be explained by the model's sensitivity to staining variations as well as both classes' similarly dotted structure. Moreover, the model struggles to properly embed the representations of mucin. The scarcity of mucinous examples in the target domain makes the representation of this class difficult.

As in the patch classification task, SENTRY is as the top performing baseline. However, the model is still limited by its capacity to distinguish between tumor and normal mucosa due to the few label setting. Also, the detection of the stroma area appears less detailed compared to other approaches such as OSDA or SRMA.

Our approach outperforms the other SOTA domain adaptation methods in terms of pixel-wise accuracy, weighed intersection over union (IoU) and pixel-wise Cohen's kappa score $\kappa$. Regions with mixed tissue types (e.g., lymphocytes + stroma or stroma + isolated tumor cells) are challenging cases because the samples available in the public cohorts mainly contain homogeneous tissue textures and few examples of class mixtures. Subsequently, domain adaptation methods naturally struggle to align features resulting in a biased classification. We observe that thinner or torn stroma regions, where the background behind is well visible, are often misclassified as adipose tissue by SRMA, which is most likely due to their similar appearance. However, our SRMA model is able to correctly distinguish between normal mucosa and tumor, which are tissue regions with very relevant information downstream tasks such as survival analysis.

Figure 7 presents a qualitative visualization of the model's embedding space. The figure shows the actual visual distribution of the target patches, the source domain label arrangement, and the overlap of the source and target domain. The patch visualization also shows a smooth transition between class representations where for example, neighboring samples of the debris cluster include a mixture of tissue and debris. The embedding reveals a large area in the center of the visualization that does not match with the source domain. The area mostly includes loose connective tissue and stroma, which are both under-represented in the training examples. Also, mucin is improperly matched to the loose

stroma, which explains the misclassification of stromal tissue in the ROI 2. The scarcity of mucinous examples in our in-house cohort makes it difficult for the model to find good candidates.

Table 2: Ablation study for the proposed Self-Rule to Multi-Adapt (SRMA) approach. We denote $\mathcal{L}_{\mathrm{IND}}$ as the in-domain loss, $\mathcal{L}_{\mathrm{CRD}}$ as the cross-domain loss, and E2H as easy-to-hard. We train the domain adaptation from Kather-19 (source) to Kather-16 (target). Only $1\%$ of the source domain labels are used, and no labels for the target domain. We report the F1 and weighted F1 score for the individual classes and the overall mean performance (all) (average over 10 runs).

| Methods | $\mathcal{L}_{\mathrm{IND}}$ | $\mathcal{L}_{\mathrm{CRD}}$ | E2H | TUM | STR | LYM | DEB | NORM | ADI | BACK | ALL |
|---|---|---|---|---|---|---|---|---|---|---|---|
| MoCoV2[†] | - | - | - | 36.8** | 45.4** | 27.1** | 30.8** | 45.2** | 43.1** | 43.6** | 38.9** |
| SRA[‡] | ✓ | - | - | 88.1** | 72.8** | 78.0* | 71.8* | 89.9** | 93.4* | 86.0* | 82.9** |
| SRA[‡] | - | ✓ | - | 14.1** | 9.1** | 0.2** | 10.1** | 4.9** | 0.0** | 61.5** | 14.4** |
| SRA[‡] | ✓ | ✓ | - | 63.0** | 69.9** | 85.1 | 57.7** | 98.2+ | 97.9 | 90.0** | 80.3** |
| SRA[‡] | ✓ | ✓ | ✓ | 93.4** | 72.9** | 82.7* | 67.9 | 96.5** | 97.0** | 97.2* | 86.9* |
| SRMA | - | ✓ | - | 35.3** | 3.6** | 0.0** | 2.1 | 15.6** | 64.0** | 16.5** | 19.8** |
| SRMA | ✓ | ✓ | - | 93.3** | 77.4+ | 80.5** | 66.2+ | 91.4** | 97.8+ | 98.3 | 86.5* |
| SRMA | ✓ | ✓ | ✓ | 97.3 | 79.3 | 80.2** | 62.2** | 98.7 | 97.6+ | 98.1+ | 87.7 |

[†] Chen et al. [2020a]. Source and target dataset are merged and trained using contrastive learning.
[‡] Abbet et al. [2021].
+ $p \geq 0.05$; * $p < 0.05$; ** $p < 0.001$; unpaired t-test with respect to top result.

Table 3: Ablation study for the proposed Self-Rule to Multi-Adapt (SRMA) approach. We denote $\mathcal{L}_{\mathrm{IND}}$ as the in-domain loss, $\mathcal{L}_{\mathrm{CRD}}$ as the cross-domain loss, and E2H as easy-to-hard. We train the domain adaptation from Kather-19 (source) to our in-house dataset (target). Only $1\%$ of the source domain labels are used, and no labels for the target domain. We report the pixel-wise accuracy, the weighted intersection over union, and the pixel-wise Cohen's kappa ($\kappa$) score for three manually annotated regions of interest (ROIs) (average over 10 runs).

| Methods | $\mathcal{L}_{\mathrm{IND}}$ | $\mathcal{L}_{\mathrm{CRD}}$ | E2H | ROI 1 | | | ROI 2 | | | ROI 3 | | | ROI 1-3 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Acc. | IoU | $\kappa$ | Acc. | IoU | $\kappa$ | Acc. | IoU | $\kappa$ | Acc. | IoU | $\kappa$ |
| MoCoV2 [†] | - | - | - | 0.556** | 0.470** | 0.417** | 0.298** | 0.198** | 0.220** | 0.321** | 0.255** | 0.240** | 0.399** | 0.301** | 0.319** |
| SRA [‡] | ✓ | - | - | 0.754* | 0.655+ | 0.646* | 0.679* | 0.551* | 0.594* | 0.498** | 0.357** | 0.415** | 0.644** | 0.497** | 0.590** |
| SRA [‡] | - | ✓ | - | 0.108** | 0.022** | 0.000** | 0.060** | 0.004** | 0.000** | 0.061** | 0.006** | 0.000** | 0.076** | 0.008** | 0.000** |
| SRA [‡] | ✓ | ✓ | - | 0.766* | 0.660+ | 0.658* | 0.701* | 0.582* | 0.619* | 0.526** | 0.368* | 0.438** | 0.664** | 0.526* | 0.615** |
| SRA [‡] | ✓ | ✓ | ✓ | 0.752* | 0.638* | 0.639* | 0.689** | 0.574** | 0.607** | 0.541* | 0.373* | 0.448** | 0.661** | 0.521** | 0.611** |
| SRMA | - | ✓ | - | 0.593** | 0.471** | 0.429** | 0.096** | 0.019** | 0.029** | 0.261** | 0.118** | 0.080** | 0.322** | 0.166** | 0.196** |
| SRMA | ✓ | ✓ | - | 0.724** | 0.634** | 0.608** | 0.706+ | 0.591+ | 0.630+ | 0.518** | 0.319** | 0.415** | 0.650** | 0.484** | 0.599** |
| SRMA | ✓ | ✓ | ✓ | 0.782 | 0.668 | 0.678 | 0.711 | 0.593 | 0.630 | 0.558 | 0.388 | 0.466 | 0.684 | 0.535 | 0.635 |

[†] Chen et al. [2020a]. Source and target dataset are merged and trained using contrastive learning.
[‡] Abbet et al. [2021].
+ $p \geq 0.05$; * $p < 0.05$; ** $p < 0.001$; unpaired t-test with respect to top result.

## 4.4 Ablation Study of the Proposed Loss Function

In this section, we present the ablation study of our SRMA approach. We denote $\mathcal{L}_{\mathrm{IND}}$ as the in-domain loss, $\mathcal{L}_{\mathrm{CRD}}$ as the cross-domain loss, and E2H as the easy-to-hard learning scheme. We evaluate the performance of our model on two tasks. The first one is the domain alignment between K19 (source) and K16 (target), which follows the experimental setting described in Section 4.2. The results are presented in Table 2. The second task is the domain adaptation of K19 (source) to ROIs from our in-house dataset (target), as presented in Section 4.3. Table 3 shows the results of these experiments. The following section jointly discusses the results of both tasks.

We use MoCoV2 [Chen et al., 2020a] as a baseline, where both domains are merged to a dataset $\mathcal{D}$, and train following a contrastive learning approach. We also compare our proposed approach SRMA to our previous work SRA [Abbet et al., 2021]. For the single-source domain adaptation, the difference between SRA and the proposed extension SRMA lies in the reformulation of the cross-entropy matching. As a result, only the entropy-related terms, namely $\mathcal{L}_{\mathrm{CRD}}$ and E2H, are affected. Thus, training SRA and SRMA using only the in-domain loss $\mathcal{L}_{\mathrm{IND}}$ is the same set-up.

The baseline fails to learn discriminant features that match both sets leading to poor performances in both cross-domain adaptation tasks. This shows that, if not constrained, the model is not able to generalize the knowledge and ends up learning two distinct feature spaces, one for the source and one for the target domain.

| Images | Metrics | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Acc. | IoU | $\kappa$ | Acc. | IoU | $\kappa$ | Acc. | IoU | $\kappa$ |
| | $s_w = 0.125 , s_h = 0.075$ | | | $s_w = 0.125 , s_h = 0.1$ | | | $s_w = 0.125 , s_h = 0.125$ | | |
| ROI 1 | **0.777**$^+$ | 0.658$^*$ | **0.670**$^+$ | 0.758$^*$ | 0.642$^{**}$ | 0.646$^{**}$ | 0.752$^{**}$ | 0.652$^*$ | 0.643$^{**}$ |
| ROI 2 | 0.686$^{**}$ | 0.567$^{**}$ | 0.602$^{**}$ | 0.653$^{**}$ | 0.527$^{**}$ | 0.561$^{**}$ | 0.697$^*$ | 0.565$^*$ | 0.613$^*$ |
| ROI 3 | 0.544$^*$ | 0.375$^*$ | 0.452$^*$ | 0.542$^*$ | **0.388**$^{**}$ | **0.458**$^+$ | 0.546$^*$ | 0.369$^*$ | 0.454$^*$ |
| ALL | 0.669$^*$ | 0.518$^{**}$ | 0.618$^{**}$ | 0.651$^{**}$ | 0.495$^{**}$ | 0.599$^{**}$ | 0.665$^{**}$ | 0.509$^{**}$ | 0.615$^{**}$ |
| | $s_w = 0.25 , s_h = 0.15$ | | | $s_w = 0.25 , s_h = 0.2$ | | | $s_w = 0.25 , s_h = 0.25$ | | |
| ROI 1 | **0.782**$^+$ | **0.668**$^+$ | **0.678**$^+$ | 0.764$^*$ | 0.642$^{**}$ | 0.654$^*$ | 0.756$^*$ | 0.633$^{**}$ | 0.642$^{**}$ |
| ROI 2 | **0.711**$^+$ | **0.593** | **0.630**$^+$ | **0.709**$^+$ | 0.581$^*$ | **0.626**$^+$ | 0.703$^*$ | 0.573$^*$ | **0.620**$^+$ |
| ROI 3 | 0.558 | 0.388 | 0.466 | 0.552$^+$ | **0.379**$^*$ | 0.464$^+$ | 0.542$^*$ | **0.384**$^*$ | 0.459$^+$ |
| ALL | **0.684** | **0.535**$^+$ | **0.635** | 0.675$^*$ | 0.521$^{**}$ | 0.626$^{**}$ | 0.667$^*$ | 0.511$^{**}$ | 0.617$^{**}$ |
| | $s_w = 0.5 , s_h = 0.45$ | | | $s_w = 0.5 , s_h = 0.6$ | | | $s_w = 0.5 , s_h = 0.75$ | | |
| ROI 1 | **0.786** | **0.680** | **0.684** | 0.758$^{**}$ | 0.641$^{**}$ | 0.646$^{**}$ | 0.745$^{**}$ | 0.626$^{**}$ | 0.629$^{**}$ |
| ROI 2 | **0.714** | **0.589**$^+$ | **0.631** | 0.697$^*$ | 0.563$^{**}$ | 0.610$^*$ | 0.697$^*$ | 0.571$^*$ | 0.614$^*$ |
| ROI 3 | 0.534$^{**}$ | **0.380**$^+$ | 0.447$^*$ | 0.524$^{**}$ | 0.370$^*$ | 0.439$^{**}$ | 0.520$^{**}$ | 0.364$^*$ | 0.438$^{**}$ |
| ALL | **0.678**$^+$ | **0.539** | **0.629**$^+$ | 0.659$^{**}$ | 0.510$^{**}$ | 0.609$^{**}$ | 0.654$^{**}$ | 0.496$^{**}$ | 0.603$^{**}$ |

$^+$ $p \geq 0.05$; $^*$ $p < 0.05$; $^{**}$ $p < 0.001$; unpaired t-test with respect to top result.



Classification performance for different $s_w$, $s_h$ values on the ROIs.       Profile of the easy-to-hard ratio $r$.

Figure 8: Importance of $s_w$ and $s_h$ parameter tuning for the easy-to-hard learning scheme. (left) Performance of the model on the three regions of interest (ROIs) for each parameter pair. (right) Corresponding profiles of the step function $r$ (Equation 10) as a function of the current epoch. The variable $r$ represents the fraction of the "trusted" samples included for cross-domain matching, based the cross-entropy.

Training with using only $\mathcal{L}_{\text{IND}}$ achieves relatively good performances but fails to generalize knowledge to classes where textures and staining strongly vary. In the patch classification task for example, this is apparent for the background and tumor class. For the second evaluation task, we can observe the same trend in the ROI 3 where tumor and normal stroma are mixed.

Using only $\mathcal{L}_{\text{CRD}}$ does not help and creates an unstable model. As we do not impose domain representation, the model converges toward incorrect solutions where random sets of samples are matched between the source and target datasets. Moreover, it can create degenerated solutions where examples from the source and target domain are perfectly matched even though they do not present any visual similarity. The reformulation of the entropy, however, slightly improves the cross domain matching.

Even the combination of the in-domain and cross-domain loss is not sufficient to improve the capability of the model. When performing a class-wise analysis, we observe that the performance on tumor and debris detection drastically dropped without the entropy reformulation. Both classes are forced to match samples from other classes, thus worsening the representation of the embedding.

The introduction of the E2H procedure improves the overall classification as well as most of the per-class performance for the first task. In the second task, it improves the performance across all metrics in all three ROIs. The importance of the E2H learning is evaluated and discussed in more detail in the next section.

Overall, the updated definition of the entropy improves the model's performance for both the cross-domain patch classification and WSIs segmentation task. It helps to ensure that both model branches output a similar distribution, thus providing better cross-domain candidates. The improvement is most visible for the tumor and stroma predictions.

## 4.5 Evaluation of the E2H Learning Scheme

In this section, we discuss the usefulness and robustness of the E2H learning. The learning procedure is based on $r$, and the two contributing variables $s_w$ and $s_h$:

$$r = \left\lfloor \frac{e}{N_{\text{epochs}} \cdot s_w} \right\rfloor \cdot s_h, \qquad \text{(10 revisited)}$$

In Figure 8, we show the impact of different combinations of these parameters on the single cross-domain segmentation task (see Section 4.3). We report the pixel-wise accuracy, the weighted intersection over union, and the pixel-wise
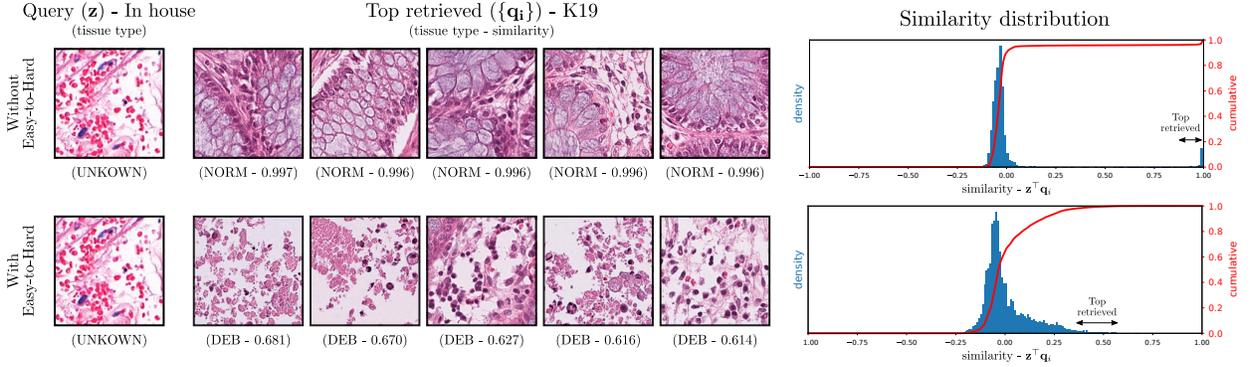
Figure 9: Importance of the easy-to-hard (E2H) learning scheme for the cross-domain image retrieval. The first column shows the input query image $\mathbf{z}$ from our in-house cohort (target domain), the second column presents the retrieved samples from K19 that have the highest similarity in the source queue $\{\mathbf{q}_i\}$, and the third column shows the density distribution (blue) of similarities across the source queue as well as its cumulative profile (red). We list the retrieved examples with their assigned classes. The query class is unknown. The top and bottom rows highlight the result of training without and with E2H learning, respectively. Without E2H, the model tries to optimize $\mathcal{L}_{\text{CRD}}$ at any cost, which creates out-of-distribution samples (seen at the very right). With E2H the model predicts samples with lower confidence, but that are still visually similar.

Cohen's kappa ($\kappa$) score for the presented ROIs. For each parameter pair, we also display the profile of the variable $r$ as a function of the of the epoch $e$.

Firstly, we observe that the model is more robust when $s_h$ is low. The variable is an indicator of the ratio of samples used for cross-domain matching. In other words, the architecture benefits from a small $s_h$ that allows it to focus on examples with high similarity/confidence while avoiding complex samples without properly matching candidates. Secondly, the selection of $s_w$ is also crucial to the stability of the prediction. This quantity measures the number of epochs to wait before considering more complex examples in the cross-domain matching optimization. For small $s_w$ values, the model has no time to learn the feature representation properly before encountering more difficult samples. This is especially true for the first few epochs after initialization, where the architecture is not yet able to optimally embed features. Furthermore, using large $s_w$ weakens the model capability to progressively learn from more complex samples.

Figure 9 shows an example patch from the training phase and highlights the usefulness of the E2H scheme. When dealing with a heterogeneous target data cohort, some tissue types might not have relevant candidates in the other set (open-set scenario). The presented example shows an example composed of a vein and blood cells. Such a tissue structure is absent from the source cohort, and thus does not have matching sample in the target domain.

Without the E2H learning, the model is forced to find matching candidates for the query $\mathbf{z}$, here normal mucosa (NORM), to minimize the cross-entropy term $\bar{H}$. When plotting the similarity distribution, the matched samples form an out-of-distribution cluster with a high similarity to the query ($\mathbf{z}^\top \mathbf{q}_i \simeq 1$). This phenomenon is even more visible with the cumulative function (red) that tends to the step function.

When training with the E2H scheme, we observe a continuous transition in the distribution of samples similarities. Here, the top retrieved samples share the same granular structure as the query. Still, we have to be careful as they do not represent the same type of tissue. The retrieved samples are examples of necrosis, whereas the query shows red blood cells. The fact that the architecture is less confident (i.e., the similarity is lower for the top retrieved samples) is a good indicator of its robustness and ability to process complex queries.

As a result, the introduction of the E2H process prevents the model from learning degenerated solutions. We also observe this with other open-set tissue classes such as complex stroma and loose connective tissue, which are absent in the source domain.

## 4.6 Multi-Source Patch Classification

We explore the benefit of using multiple source domains with different distributions to perform domain adaptation for the patch classification task. To do so, we select K19 and K16 as the source sets and CRC-TP as the target set. To learn the feature representations, the model is trained in an unsupervised fashion using both source domains as well as the unlabeled target domain. For the evaluation, we train a linear classifier on top of the frozen features with few randomly

Table 4: Performance of the Self-Rule to Multi-Adapt (SRMA) framework on the CRC-TP dataset in a multi-source domain setting. We show the results for different combinations of K16 and K19 used for the self-supervised pre-training as well as training the classification header. For the source domains K19 and K16, only $1\%$ and $10\%$ of the labeled data are used, respectively. We also compare the performance of the $1:1$ with the $K:1$ setting for the loss definitions (see Equations 12-15). We report the F1 score for the individual classes and weighted F1 score for the overall mean performance (all) (averaged over 10 runs). Some classes have been merged due to overlapping definitions.

| | Pretraining | | Classification | | Multi-source | | | | | | | |
| Methods | K19 | K16 | K19 | K16 | $\mathcal{L}_{\text{IND}}$ | $\mathcal{L}_{\text{CRD}}$ | TUM | STR$^\dagger$ | LYM | NORM | DEB$^\dagger$ | ALL |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *Single source:* | | | | | | | | | | | | |
| SRA [Abbet et al., 2021] | - | ✓ | - | ✓ | - | - | **82.2** | **69.3** | 62.5$^{**}$ | **69.8** | 47.4$^{*}$ | **69.4** |
| SRMA | - | ✓ | - | ✓ | - | - | 82.0$^{+}$ | 63.5$^{**}$ | **66.3** | 51.9$^{**}$ | **50.3** | 65.2$^{**}$ |
| SRA [Abbet et al., 2021] | ✓ | - | ✓ | - | - | - | 91.0$^{*}$ | 84.9$^{**}$ | 62.0$^{*}$ | **71.7** | 58.5$^{+}$ | 79.2$^{**}$ |
| SRMA | ✓ | - | ✓ | - | - | - | **91.7** | **86.7** | **65.4** | 68.6$^{**}$ | **58.9** | **80.2** |
| *Multi source:* | | | | | | | | | | | | |
| DeepAll [Dou et al., 2019] | ✓ | ✓ | - | ✓ | - | - | 52.4$^{**}$ | 64.1$^{**}$ | 36.5$^{**}$ | 14.2$^{**}$ | 13.8$^{**}$ | 47.1$^{**}$ |
| SRA [Abbet et al., 2021] | ✓ | ✓ | - | ✓ | $1:1$ | $1:1$ | 70.9$^{**}$ | 68.5$^{**}$ | 45.6$^{**}$ | 72.2$^{**}$ | 19.1$^{**}$ | 62.2$^{**}$ |
| SRMA | ✓ | ✓ | - | ✓ | $1:1$ | $1:1$ | 76.6$^{**}$ | 69.3$^{**}$ | 48.7$^{**}$ | 74.5$^{**}$ | 18.2$^{**}$ | 64.4$^{**}$ |
| SRMA | ✓ | ✓ | - | ✓ | $K:1$ | $1:1$ | **89.4$^{+}$** | **74.9$^{+}$** | **66.8** | **75.6** | **43.7** | **74.4** |
| SRMA | ✓ | ✓ | - | ✓ | $1:1$ | $K:1$ | 75.9$^{**}$ | 73.3$^{*}$ | 45.9$^{**}$ | 73.0$^{**}$ | 22.6$^{**}$ | 65.8$^{**}$ |
| SRMA | ✓ | ✓ | - | ✓ | $K:1$ | $K:1$ | **89.8** | 75.2 | 64.5$^{**}$ | 74.1$^{**}$ | 25.7$^{**}$ | 72.5$^{**}$ |
| DeepAll [Dou et al., 2019] | ✓ | ✓ | ✓ | - | - | - | 72.4$^{**}$ | 88.6$^{**}$ | 43.6$^{**}$ | 53.2$^{**}$ | 71.8$^{**}$ | 73.2$^{**}$ |
| SRA [Abbet et al., 2021] | ✓ | ✓ | ✓ | - | $1:1$ | $1:1$ | 86.2$^{**}$ | 87.6$^{**}$ | 66.7$^{**}$ | 71.0$^{**}$ | **80.5** | 81.8$^{**}$ |
| SRMA | ✓ | ✓ | ✓ | - | $1:1$ | $1:1$ | **92.5** | 88.4$^{**}$ | 68.7$^{**}$ | 68.3$^{**}$ | 74.2$^{*}$ | 82.9$^{*}$ |
| SRMA | ✓ | ✓ | ✓ | - | $K:1$ | $1:1$ | 91.5$^{*}$ | 87.6$^{**}$ | 70.7 | **75.0** | 65.7$^{**}$ | 82.7$^{*}$ |
| SRMA | ✓ | ✓ | ✓ | - | $1:1$ | $K:1$ | 90.1$^{**}$ | **90.1** | 69.6$^{+}$ | 72.9$^{**}$ | 71.6$^{**}$ | **83.6** |
| SRMA | ✓ | ✓ | ✓ | - | $K:1$ | $K:1$ | 91.6$^{+}$ | 87.4$^{**}$ | 68.7$^{**}$ | 73.9$^{**}$ | 53.3$^{**}$ | 81.2$^{**}$ |
| DeepAll [Dou et al., 2019] | ✓ | ✓ | ✓ | ✓ | - | - | 81.4$^{**}$ | **85.7$^{+}$** | 50.9$^{**}$ | 50.1$^{**}$ | 51.5$^{**}$ | 72.6$^{**}$ |
| SRA [Abbet et al., 2021] | ✓ | ✓ | ✓ | ✓ | $1:1$ | $1:1$ | 85.8$^{**}$ | 85.9 | 72.9$^{*}$ | 72.1$^{**}$ | **59.2** | **80.1** |
| SRMA | ✓ | ✓ | ✓ | ✓ | $1:1$ | $1:1$ | **92.9** | 82.4$^{**}$ | 72.1$^{*}$ | 70.8$^{**}$ | 53.7$^{**}$ | 79.3$^{*}$ |
| SRMA | ✓ | ✓ | ✓ | ✓ | $K:1$ | $1:1$ | 92.8$^{+}$ | 81.7$^{**}$ | **73.5** | **74.6** | 49.8$^{**}$ | 79.3$^{*}$ |
| SRMA | ✓ | ✓ | ✓ | ✓ | $1:1$ | $K:1$ | 89.6$^{**}$ | 84.7$^{*}$ | 72.5 | 74.4$^{+}$ | 52.1$^{**}$ | 80.0$^{+}$ |
| SRMA | ✓ | ✓ | ✓ | ✓ | $K:1$ | $K:1$ | 92.5$^{*}$ | 80.6$^{**}$ | 70.5$^{**}$ | 73.9$^{**}$ | 39.4$^{**}$ | 77.4$^{**}$ |

$^\dagger$ The STR and MUS classes are merged as STR class; DEB and MUC classes as DEB.
$^{+}$ $p \geq 0.05$; $^{*}$ $p < 0.05$; $^{**}$ $p < 0.001$; unpaired t-test with respect to top result.

selected labeled samples from the source domains (1000 samples from K19 (1%), and 500 samples from K16 (10%)). By using only little labeled data, we aim to reduce the annotation workload for pathologists while still achieving good classifications performances. The set of labeled data differs between each run, as they are randomly sampled for each individual run.

The three datasets K19, K16, and CRC-TP do not have one-to-one classes correspondence. Thus, for the evaluation of the target set, we only consider the classes present in all datasets, namely, tumor (TUM), stroma (STR), lymphocytes (LYM), normal mucosa (NORM), and debris / necrotic tissue (DEB). Still, during the unsupervised pre-training we consider all classes, including those who do not have matching candidates across the sets, such as background (BACK) and adipose (ADI). This setup creates an open-set scenario for the cross-domain matching and allows the model to learn more robust features representations.

For comparison purposes, we use the same hyper-parameters as in the single source domain patch classification setting with $s_w = 0.25$, $s_h = 0.15$. The probability of drawing a sample $\mathbf{x}$ from the source or the target domain is the same and is given by $p(\mathbf{x} \in \mathcal{D}_s) = Kp(\mathbf{x} \in \mathcal{D}_s^k) = p(\mathbf{x} \in \mathcal{D}_t)$, where K is the number of source domains.

The results are presented in Table 4. We compare the performance of different experimental setups in regards to the used datasets and multi-source scenario for our SRMA. We show three scenarios where we use either K16, K19, or the combination of the two (K16 and K19) to train the classification layer. To evaluate the impact of the multi-source scenario, where we investigate all possibilities for the in-domain ($\mathcal{L}_{\text{IND}}^{1:1}$, $\mathcal{L}_{\text{IND}}^{K:1}$) and cross-domain ($\mathcal{L}_{\text{CRD}}^{1:1}$, $\mathcal{L}_{\text{CRD}}^{K:1}$) loss definitions, as introduced in Equations 12-15. As baselines, we consider the single source setting of the presented SRMA model, our previous SRA work, as well as the DeepAll approach that uses aggregation of all the source tissue data into a single training set [Dou et al., 2019].

The SRMA and SRA single source baselines both show a better performance for K19 compared to K16. This is most likely due to the fact that the K16 subset for training the classification header is relatively small, with only $5,000$ different examples. Also, SRMA outperforms our previous SRA work for all classes except one, which is an indicator of the robustness of the entropy reformulation.

For the multi-source adaptation, we show three scenarios where we use either K16, K19, or the combination of the two (K16 and K19) to train the classification layer. When using solely K16, we can observe that the debris classification tends to have lower performances across all models. Debris examples in K16 appear highly saturated, which makes the generalization of the class a challenging task. Only the proposed SRMA approach is able to achieve better performances compared to the single source baselines. Using K19 for the classification of target patches gives overall the best performance. Interestingly, using both K19 and K16 leads to a drop in performance. This is most likely due to potential discrepancies between the class definitions, which makes it more difficult for the model to generalize the class representations across the different modalities.

When comparing the in-domain and cross-domain multi-source scenarios, we find that using $\mathcal{L}_{\text{IND}}^{1:1}$ and $\mathcal{L}_{\text{CRD}}^{K:1}$ achieves the best results across the various settings. This suggests that it is better to optimize the source domain as a single set for the in-domain representation. However, when performing cross-domain matching, considering domain to domain correspondence between each source set and the target domain yields better performances. It ensures that the model looks for relevant candidates in all individual source sets as tissue samples might have a distinct appearance in different source domains.

We also note that $\mathcal{L}_{\text{IND}}^{K:1}$ is only relevant when only using K16 to train the classification header. This is due to the fact that the cross-domain matching fails to retrieve debris samples correctly from the K16 domain, which tend to be misclassified as lymphocytes because of their similar granular appearance and as well as their hematoxylin-positive aspect. Overall the combination of both $\mathcal{L}_{\text{IND}}^{K:1}$ and $\mathcal{L}_{\text{CRD}}^{K:1}$ degrades the performance slightly.

Complementary results on the importance of the dataset ratios when sampling data for the unsupervised pre-training phase are available in C.

### 4.7   Use Case: Multi-source Segmentation of WSI

In this section, we present the results for the multi-source domain adaptation for patch-based segmentation of WSI ROIs. More specifically, we are interested in the detection of desmoplastic reactions (complex stroma), which is a prognostic factor in CRC [Ueno et al., 2021]. We use both K19 and CRC-TP as the source datasets to add complex stroma examples to the source domain. Our in-house dataset is used as the target domain.

To assess the quality of the prediction, we evaluate the models on the same ROIs as in the single-source setting. However, the previously provided annotations do not include complex stroma. We overcome this by defining a margin around the tumor tissue in the existing annotations, which is considered as the interaction area. Stroma in this region is therefore re-annotated as complex stroma. The margin is fixed to $500\mu m$ such that it includes the close tumor neighborhood and matches the definition of complex stroma in the literature [Berben et al., 2020, Nearchou et al., 2021].

As a baseline, we use DeepAll, which aggregates all the source tissue data into a single training set [Dou et al., 2019]. The model is trained in an unsupervised fashion using a standard contrastive loss to optimize the data representation of the features [Chen et al., 2020a]. In this case, no domain adaption is performed across the sets.

The results are presented in Table 5 and Figure 10. In Table 5, we compare the performance of the models with and without complex stroma detection across all three ROIs. We compare the single as well as the multi-source SRMA approaches to the baselines, DeepAll and our previously published SRA method. We report the F1-score for complex stroma, the overall weighted F1-score, the pixel-wise accuracy, the Dice score, the weighted intersection over union (IoU), and pixel-wise Cohen's kappa ($\kappa$).

Without considering the complex stroma class, the numerical results show that all the multi-source settings achieve similar performances. Including an additional dataset, namely CRC-TP, does not improve nor seriously deteriorate the classification performances on the ROIs. Furthermore, merging the source domains for in-domain optimization ($\mathcal{L}_{\text{IND}}^{1:1}$) seems to be the best setup. For the cross-domain matching, both $\mathcal{L}_{\text{CRD}}^{1:1}$ and $\mathcal{L}_{\text{CRD}}^{K:1}$ achieve similar scores.

However, the benefit of using the multi-source approach can be observed when including complex stroma detection. Here, the models which use CRC-TP as source set achieve better results. The detection of complex stroma improves by up to $20 - 25\%$. By contrast, the cross-domain matching on each subsets $\mathcal{L}_{\text{CRD}}^{K:1}$ penalizes the complex stroma detection. This can be explained by the fact that only CRC-TP contains examples of complex stroma. Therefore, imposing complex stroma retrieval in K19 is unfeasible. Another challenge is the relatively significant overlap between the complex stroma and the tumor class. The model tends to classify the tumor border area as complex stroma.

Table 5: Analysis of the performance of the Self-Rule to Multi-Adapt (SRMA) approach in regards to complex stroma detection. Multiple possible scenarios are evaluated in regard to the data included for pre-training, as well as the multi-source setting ($1 : 1$ versus $K : 1$, see Equations 12-15), as indicated in the table. Only $1\%$ of the labels are used for the classification stage. We report the F1-score for complex stroma, the overall weighed F1-score, the pixel-wise accuracy, the dice score, the weighted intersection over union (IoU), and the pixel-wise Cohen's kappa ($\kappa$) (averaged over 10 runs).

| | Pretraining | | Multi-source | | | | | | | |
| Model | K19 | CRCTP | $\mathcal{L}_{\text{IND}}$ | $\mathcal{L}_{\text{CRD}}$ | F1-CSTR$^\dagger$ | F1-ALL | Acc. | Dice | IoU | $\kappa$ |
|---|---|---|---|---|---|---|---|---|---|---|
| *ROI 1-3 (w/o CSTR)* | | | | | | | | | | |
| DeepAll [Dou et al., 2019] | ✓ | ✓ | - | - | - | $0.622^{**}$ | $0.615^{**}$ | $0.583^{**}$ | $0.483^{**}$ | $0.552^{**}$ |
| SRA [Abbet et al., 2021] | ✓ | - | - | - | - | $0.648^{**}$ | $0.661^{**}$ | $0.632^{**}$ | $0.521^{**}$ | $0.611^{**}$ |
| SRMA | ✓ | - | - | - | - | $\mathbf{0.667}^{+}$ | $\mathbf{0.684}^{+}$ | $0.647^{**}$ | $\mathbf{0.536}^{+}$ | $\mathbf{0.636}^{+}$ |
| SRMA | ✓ | ✓ | $1 : 1$ | $1 : 1$ | - | $\mathbf{0.673}$ | $\mathbf{0.685}$ | $\mathbf{0.669}$ | $\mathbf{0.541}$ | $\mathbf{0.636}$ |
| SRMA | ✓ | ✓ | $K : 1$ | $1 : 1$ | - | $0.644^{**}$ | $0.665^{**}$ | $0.637^{**}$ | $0.516^{**}$ | $0.615^{**}$ |
| SRMA | ✓ | ✓ | $1 : 1$ | $K : 1$ | - | $\mathbf{0.662}^{+}$ | $\mathbf{0.678}^{+}$ | $0.652^{**}$ | $0.528^{*}$ | $\mathbf{0.629}^{+}$ |
| SRMA | ✓ | ✓ | $K : 1$ | $K : 1$ | - | $0.638^{**}$ | $0.660^{**}$ | $0.632^{**}$ | $0.509^{**}$ | $0.609^{**}$ |
| *ROI 1-3 (w/ CSTR)* | | | | | | | | | | |
| DeepAll [Dou et al., 2019] | ✓ | ✓ | - | - | $0.001^{**}$ | $0.505^{**}$ | $0.539^{**}$ | $0.496^{**}$ | $0.399^{**}$ | $0.479^{**}$ |
| SRA [Abbet et al., 2021] | ✓ | - | - | - | $0.214^{**}$ | $0.600^{**}$ | $0.624^{**}$ | $0.582^{**}$ | $0.490^{**}$ | $0.577^{**}$ |
| SRMA | ✓ | - | - | - | $0.263^{**}$ | $0.614^{**}$ | $0.641^{**}$ | $0.595^{**}$ | $0.498^{**}$ | $0.594^{**}$ |
| SRMA | ✓ | ✓ | $1 : 1$ | $1 : 1$ | $\mathbf{0.479}^{+}$ | $\mathbf{0.647}^{+}$ | $0.659^{*}$ | $\mathbf{0.631}$ | $\mathbf{0.524}^{+}$ | $0.613^{**}$ |
| SRMA | ✓ | ✓ | $K : 1$ | $1 : 1$ | $\mathbf{0.492}$ | $\mathbf{0.650}$ | $\mathbf{0.669}$ | $0.618^{**}$ | $\mathbf{0.524}$ | $\mathbf{0.624}$ |
| SRMA | ✓ | ✓ | $1 : 1$ | $K : 1$ | $\mathbf{0.464}^{+}$ | $\mathbf{0.640}^{+}$ | $0.651^{**}$ | $0.619^{**}$ | $0.513^{*}$ | $0.604^{**}$ |
| SRMA | ✓ | ✓ | $K : 1$ | $K : 1$ | $0.366^{**}$ | $0.623^{**}$ | $0.646^{**}$ | $0.597^{**}$ | $0.500^{**}$ | $0.599^{**}$ |

$^\dagger$ Performances are only available with extended annotations (w/CSTR).
$^{+}$ $p \geq 0.05$; $^{*}$ $p < 0.05$; $^{**}$ $p < 0.001$; unpaired t-test with respect to top

In Figure 10, we display the visual results of the complex stroma detection on ROI 1 and 3, where desmoplastic reactions, and thus complex stroma, are present. We show, from left to right, the reference images, the original ground truth labels, the extended ground truth labels with complex stroma, the DeepAll baseline, our previous SRA work, and as well the results of the presented SRMA model ($\mathcal{L}_{\text{IND}}^{1:1}$ and $\mathcal{L}_{\text{CRD}}^{K:1}$ setting).

SRMA outperforms the baselines in terms of pixel-wise accuracy, Jaccard index (IoU), and Cohen's kappa score $\kappa$. Notably, the detection of the tumor is much more detailed compared to the single-source approach in both ROIs. Parts of the tissue previously considered as tumor can now be properly matched, thanks to the introduction of the complex stroma class.

Another interesting result in ROI 3 is that all the stromal areas are now considered as either complex stroma, tumor, or lymphocytes by all models. This highlights how challenging the classification of complex stroma is without access to the higher-level context. Pathologists also find this difficult, as they rely not only on the tissue morphology for this assessment but also on the spatial relations, i.e., the proximity to the tumor area. Here, according to our extended ground truth, the complex stroma only surrounds the tumor region. However, the tissue tear disconnected some of the tumor surrounding regions, which suggests that the complex stroma area, in reality, spans even further. This correlates with the prediction of both models, which identify the whole region as complex stroma.

Lastly, using the multi-source setting allows the introduction of a new class such as complex stroma to the detection task. In the presented setting, the source domains do not need one-to-one class correspondences for the model to learn meaningful cross-domain features. Here, CRC-TP does not include mucin, background, and adipose while K19 does not contain complex stroma. This is an interesting outcome, as it shows that new data that might even be acquired under different circumstances can be added with additional tissue classes without interfering with or altering the performance of the existing classes.

A visualization of the multi-source domain embedding space as well as the patch-based segmentation of a full WSI image are available in D-E.
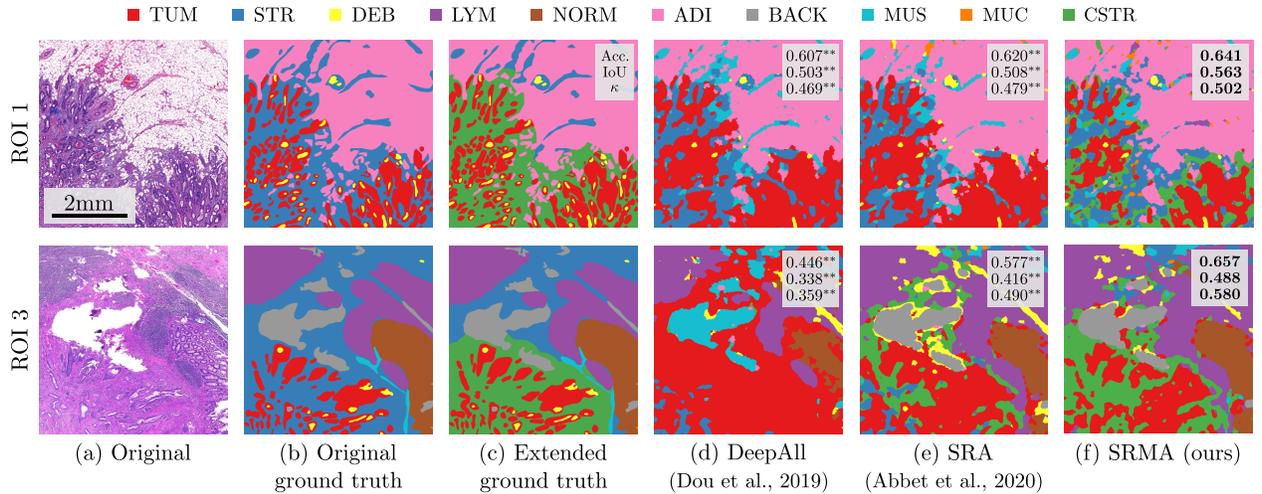
Figure 10: Results of the multi-source domain adaptation from K19 and CRC-TP to our in-house dataset. (a-c) show the original regions of interest (ROIs) from the WSIs, their original ground truth (without CSTR), and the extended ground truth (with CSTR), respectively. We compare the performance of our SRMA framework (f) to our previous work SRA (e) and to the DeepAll baseline (d). For the multi-source optimization, we use the $1:1$ and $K:1$ approach for the in-domain and cross-domain, respectively. We report the pixel-wise accuracy, the weighted intersection over union, and the pixel-wise Cohen's kappa ($\kappa$) score averaged over 10 runs.

## 5   Conclusion and Future Work

In this work, we explore the usefulness of self-supervised learning and UDA for the identification of histological tissue types. Motivated by the difficulty of obtaining expert annotations, we explore different UDA models using a variety of label-scarce colorectal cancer histopathology datasets.

As our main contribution, we present a new label transferring approach from partially labeled, public datasets (source domain) to unlabeled target domains. This is more practical than most previous UDA approaches which are often tailored to fully annotated source domain data or tied to additional network branches dedicated to auxiliary tasks. Instead, we perform progressive cross-entropy minimization based on the similarity distribution among the unlabeled target and source domain samples, yielding discriminative and domain-agnostic features for domain adaptation.

In reality, not all tissue types are equally present in a WSI, and some are quite rare. Thus, the extracted patches are imbalanced in regards to class labels (categories), which imposes significant challenges for the trained models to generalize well. For example, mucin is frequently present in mucinous carcinoma but is scarcely found in adenocarcinomas. Throughout various label transfer tasks, we show that our proposed Self-Rule to Multi-Adapt (SRMA) method can discover the relevant semantic information even in the presence of few labeled source samples, and yields a better generalization on different target domain datasets. Moreover, we show that our model definition can be generalized to a multi-source setting. As a result, the proposed model is able to learn rich data representation using multiple source domains.

Another example is the complex stroma class, which can be further divided into three subcategories (immature, intermediate, or mature), whose occurrences are highly variable and which are linked to patients' prognostic factor [Okuyama et al., 2020]. Possible future work could take this class imbalance across WSIs into account and aim to improve the quality and variety of the provided positive and negative examples.

In addition, publicly available datasets are so far mostly composed of curated and thus homogeneous patches in terms of tissue types. This data, however, do not capture the heterogeneity and complexity of patches extracted from images in the diagnostic routine. This can lead to erroneous detections, e.g., background and stroma interaction being interpreted as adipose tissue. Thus, finding a self-supervised learning approach that can also properly embed mixed patches is a possible future extension of this work.

Furthermore, the SRMA framework is also highly modular and can thus be used for similar problems in other image analysis research fields. The selected backbone can be replaced, and the used data augmentations adapted to better fit with the task and data at hand.

Lastly, the patch-based segmentation using our method can also be applied in a clinical context. Many clinically relevant downstream tasks depend on accurate tissue segmentation, such as tumor-stroma ratio calculation, disease-free survival prediction, or adjuvant treatment decision-making.

## Acknowledgments

## References

Oscar GF Geessink, Alexi Baidoshvili, Joost M Klaase, Babak Ehteshami Bejnordi, Geert JS Litjens, Gabi W van Pelt, Wilma E Mesker, Iris D Nagtegaal, Francesco Ciompi, and Jeroen AWM van der Laak. Computer aided quantification of intratumoral stroma yields an independent prognosticator in rectal cancer. *Cellular Oncology*, 42(3):331–341, 2019.

Marloes A Smit and Wilma E Mesker. The role of artificial intelligence to quantify the tumour-stroma ratio for survival in colorectal cancer. *EBioMedicine*, 61, 2020.

Narayan Hegde, Jason D Hipp, Yun Liu, Michael Emmert-Buck, Emily Reif, Daniel Smilkov, Michael Terry, Carrie J Cai, Mahul B Amin, Craig H Mermel, et al. Similar image search for histopathology: Smily. *NPJ digital medicine*, 2 (1):1–9, 2019.

Ming Y Lu, Drew FK Williamson, Tiffany Y Chen, Richard J Chen, Matteo Barbieri, and Faisal Mahmood. Data-efficient and weakly supervised computational pathology on whole-slide images. *Nature Biomedical Engineering*, 5 (6):555–570, 2021.

Eirini Arvaniti, Kim S Fricker, Michael Moret, Niels Rupp, Thomas Hermanns, Christian Fankhauser, Norbert Wey, Peter J Wild, Jan H Rueschoff, and Manfred Claassen. Automated gleason grading of prostate cancer tissue microarrays via deep learning. *Scientific reports*, 8(1):1–11, 2018.

Huu-Giao Nguyen, Annika Blank, Heather E Dawson, Alessandro Lugli, and Inti Zlobec. Classification of colorectal tissue images from high throughput tissue microarrays by ensemble deep learning methods. *Scientific Reports*, 11(1): 1–11, 2021.

Jakob Nikolas Kather, Johannes Krisam, Pornpimol Charoentong, Tom Luedde, Esther Herpel, Cleo-Aron Weis, Timo Gaiser, Alexander Marx, Nektarios A Valous, Dyke Ferber, et al. Predicting survival from colorectal cancer histology slides using deep learning: A retrospective multicenter study. *PLoS medicine*, 16(1), 2019.

Talha Qaiser, Yee-Wah Tsang, Daiki Taniyama, Naoya Sakamoto, Kazuaki Nakane, David Epstein, and Nasir Rajpoot. Fast and accurate tumor segmentation of histology images using persistent homology and deep convolutional features. *Medical image analysis*, 55:1–14, 2019.

Lyndon Chan, Mahdi S Hosseini, Corwyn Rowsell, Konstantinos N Plataniotis, and Savvas Damaskinos. Histosegnet: Semantic segmentation of histological tissue type in whole slide images. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 10662–10671, 2019.

David Tellez, Jeroen van der Laak, and Francesco Ciompi. Gigapixel whole-slide image classification using unsupervised image compression and contrastive training. *Medical Imaging with Deep Learning*, 2018.

Julio Silva-Rodríguez, Adrián Colomer, and Valery Naranjo. Weglenet: A weakly-supervised convolutional neural network for the semantic segmentation of gleason grades in prostate histology images. *Computerized Medical Imaging and Graphics*, 88:101846, 2021.

Gabriele Campanella, Matthew G Hanna, Luke Geneslaw, Allen Miraflor, Vitor Werneck Krauss Silva, Klaus J Busam, Edi Brogi, Victor E Reuter, David S Klimstra, and Thomas J Fuchs. Clinical-grade computational pathology using weakly supervised deep learning on whole slide images. *Nature medicine*, 25(8):1301–1309, 2019.

Chen Pieter, Abbeel anf Xi, Ho Jonathan, Srinivas Aravind, Li Alex, and Yan Wilson. Cs 294-158. deep unsupervised learning, February 2020.

Christian Abbet, Inti Zlobec, Behzad Bozorgtabar, and Jean-Philippe Thiran. Divide-and-rule: Self-supervised learning for survival analysis in colorectal cancer. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 480–489. Springer, 2020.

Bin Li, Yin Li, and Kevin W Eliceiri. Dual-stream multiple instance learning network for whole slide image classification with self-supervised contrastive learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14318–14328, 2021.

Jacob Gildenblat and Eldad Klaiman. Self-supervised similarity learning for digital pathology. *2nd COMPAY Workshop at MICCAI 2019*, 2019.

John N Weinstein, Eric A Collisson, Gordon B Mills, Kenna R Mills Shaw, Brad A Ozenberger, Kyle Ellrott, Ilya Shmulevich, Chris Sander, Joshua M Stuart, Cancer Genome Atlas Research Network, et al. The cancer genome atlas pan-cancer analysis project. *Nature genetics*, 45(10):1113, 2013.

Geert Litjens, Peter Bandi, Babak Ehteshami Bejnordi, Oscar Geessink, Maschenka Balkenhol, Peter Bult, Altuna Halilovic, Meyke Hermsen, Rob van de Loo, Rob Vogels, et al. 1399 h&e-stained sentinel lymph node sections of breast cancer patients: the camelyon dataset. *GigaScience*, 7(6):giy065, 2018.

Mitko Veta, Yujing J Heng, Nikolas Stathonikos, Babak Ehteshami Bejnordi, Francisco Beca, Thomas Wollmann, Karl Rohr, Manan A Shah, Dayong Wang, Mikael Rousson, et al. Predicting breast tumor proliferation from whole-slide images: the tupac16 challenge. *Medical image analysis*, 54:111–121, 2019.

Jakob Nikolas Kather, Cleo-Aron Weis, Francesco Bianconi, Susanne M Melchers, Lothar R Schad, Timo Gaiser, Alexander Marx, and Frank Gerrit Zöllner. Multi-class texture analysis in colorectal cancer histology. *Scientific reports*, 6:27988, 2016.

Kirk Shanah, Sadow Cheryl A., Smith J. Keith, Levine Seth, Roche Charles, Bonaccio Ermalinda, and Filippini Joe. Radiology data from the cancer genome atlas colon adenocarcinoma [tcga-coad] collection, 2016a.

Kirk Shanah, Sadow Cheryl A., and Levine Seth. Radiology data from the cancer genome atlas rectum adenocarcinoma [tcga-read] collection, 2016b.

Sebastian Otálora, Manfredo Atzori, Vincent Andrearczyk, Amjad Khan, and Henning Müller. Staining invariant features for improving generalization of deep convolutional neural networks in computational pathology. *Frontiers in bioengineering and biotechnology*, 7:198, 2019.

Wei-Chung Cheng, Firdous Saleheen, and Aldo Badano. Assessing color performance of whole-slide imaging scanners for digital pathology. *Color Research & Application*, 44(3):322–334, 2019.

Daisuke Komura and Shumpei Ishikawa. Machine learning methods for histopathological image analysis. *Computational and structural biotechnology journal*, 16:34–42, 2018.

Marc Macenko, Marc Niethammer, James S Marron, David Borland, John T Woosley, Xiaojun Guan, Charles Schmitt, and Nancy E Thomas. A method for normalizing histology slides for quantitative analysis. In *2009 IEEE International Symposium on Biomedical Imaging: From Nano to Macro*, pages 1107–1110. IEEE, 2009.

Farhad Ghazvinian Zanjani, Svitlana Zinger, et al. Deep convolutional gaussian mixture model for stain-color normalization of histopathological images. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 274–282. Springer, 2018.

Deepak Anand, Goutham Ramakrishnan, and Amit Sethi. Fast gpu-enabled color normalization for digital pathology. In *2019 International Conference on Systems, Signals and Image Processing (IWSSIP)*, pages 219–224. IEEE, 2019.

Allison Tam, Jocelyn Barker, and Daniel Rubin. A method for normalizing pathology images to improve feature extraction for quantitative pathology. *Medical physics*, 43(1):528–537, 2016.

Mark D Zarella, Chan Yeoh, David E Breen, and Fernando U Garcia. An alternative reference space for h&e color normalization. *PloS one*, 12(3):e0174489, 2017.

Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning*, pages 1180–1189, 2015.

Navid Alemi Koohbanani, Balagopal Unnikrishnan, Syed Ali Khurram, Pavitra Krishnaswamy, and Nasir Rajpoot. Self-path: Self-supervision for classification of pathology images with limited annotations. *IEEE Transactions on Medical Imaging*, 2021.

Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017.

Kuniaki Saito, Kohei Watanabe, Yoshitaka Ushiku, and Tatsuya Harada. Maximum classifier discrepancy for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3723–3732, 2018a.

Qi Dou, Daniel Coelho de Castro, Konstantinos Kamnitsas, and Ben Glocker. Domain generalization via model-agnostic learning of semantic features. *Advances in Neural Information Processing Systems*, 32:6450–6461, 2019.

Fabio M Carlucci, Antonio D'Innocente, Silvia Bucci, Barbara Caputo, and Tatiana Tommasi. Domain generalization by solving jigsaw puzzles. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2229–2238, 2019.

Kuniaki Saito, Shohei Yamamoto, Yoshitaka Ushiku, and Tatsuya Harada. Open set domain adaptation by backpropagation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 153–168, 2018b.

Jiaolong Xu, Liang Xiao, and Antonio M López. Self-supervised domain adaptation for computer vision tasks. *IEEE Access*, 7:156694–156706, 2019.

Kuniaki Saito, Donghyun Kim, Stan Sclaroff, and Kate Saenko. Universal domain adaptation through self supervision. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 16282–16292. Curran Associates, Inc., 2020.

Viraj Prabhu, Shivam Khare, Deeksha Kartik, and Judy Hoffman. Sentry: Selective entropy optimization via committee consistency for unsupervised domain adaptation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8558–8567, 2021.

Toshihiko Matsuura and Tatsuya Harada. Domain generalization using a mixture of multiple latent domains. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 11749–11756, 2020.

Christian Abbet, Linda Studer, Andreas Fischer, Heather Dawson, Inti Zlobec, Behzad Bozorgtabar, and Jean-Philippe Thiran. Self-rule to adapt: Learning generalized features from sparsely-labeled data using unsupervised domain adaptation for colorectal cancer tissue phenotyping. In *Medical Imaging with Deep Learning*, 2021.

Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020a.

Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020b.

Donghyun Kim, Kuniaki Saito, Tae-Hyun Oh, Bryan A Plummer, Stan Sclaroff, and Kate Saenko. Cross-domain self-supervised learning for domain adaptation with few source labels. *arXiv preprint arXiv:2003.08264*, 2020.

Yixiao Ge, Dapeng Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Self-paced contrastive learning with hybrid memory for domain adaptive object re-id. *arXiv preprint arXiv:2006.02713*, 2020.

Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.

David Tellez, Geert Litjens, Péter Bándi, Wouter Bulten, John-Melle Bokhorst, Francesco Ciompi, and Jeroen van der Laak. Quantifying the effects of data augmentation and stain color normalization in convolutional neural networks for computational pathology. *Medical image analysis*, 58:101544, 2019.

Khrystyna Faryna, Jeroen van der Laak, and Geert Litjens. Tailoring automated data augmentation to h&e-stained histopathology. In *Medical Imaging with Deep Learning*, 2021.

Mahmoud Assran, Mathilde Caron, Ishan Misra, Piotr Bojanowski, Armand Joulin, Nicolas Ballas, and Michael Rabbat. Semi-supervised learning of visual features by non-parametrically predicting view assignments with support samples. *arXiv preprint arXiv:2104.13963*, 2021.

Sajid Javed, Arif Mahmood, Muhammad Moazam Fraz, Navid Alemi Koohbanani, Ksenija Benes, Yee-Wah Tsang, Katherine Hewitt, David Epstein, David Snead, and Nasir Rajpoot. Cellular community detection for tissue phenotyping in colorectal cancer histology images. *Medical Image Analysis*, page 101696, 2020.

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.

Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020.

Hideki Ueno, Yoshiki Kajiwara, Yoich Ajioka, Tamotsu Sugai, Shigeki Sekine, Megumi Ishiguro, Atsuo Takashima, and Yukihide Kanemitsu. Histopathological atlas of desmoplastic reaction characterization in colorectal cancer. *Japanese Journal of Clinical Oncology*, 51(6):1004–1012, 2021.

Lieze Berben, Hans Wildiers, Lukas Marcelis, Asier Antoranz, Francesca Bosisio, Sigrid Hatse, and Giuseppe Floris. Computerised scoring protocol for identification and quantification of different immune cell populations in breast tumour regions by the use of qupath software. *Histopathology*, 77(1):79–91, 2020.

Ines P Nearchou, Hideki Ueno, Yoshiki Kajiwara, Kate Lillard, Satsuki Mochizuki, Kengo Takeuchi, David J Harrison, and Peter D Caie. Automated detection and classification of desmoplastic reaction at the colorectal tumour front using deep learning. *Cancers*, 13(7):1615, 2021.

Takashi Okuyama, Shinichi Sameshima, Emiko Takeshita, Takashi Mitsui, Takuji Noro, Yuko Ono, Tamaki Noie, Shinichi Ban, and Masatoshi Oya. Myxoid stroma is associated with postoperative relapse in patients with stage ii colon cancer. *BMC cancer*, 20(1):1–11, 2020.

Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *arXiv preprint arXiv:2004.11362*, 2020.

## A    Selection of Self-supervised Model

To assess which self-supervised model we should use as the backbone for the UDA, we compare the performances of several SOTA self-supervised methods (SimCLR [Chen et al., 2020b], SupContrast [Khosla et al., 2020], and MoCoV2 [Chen et al., 2020a]), as well as the performance of the standard supervised learning approach when facing different levels of data availability. The results are presented in Table 6. We report the performance of the single domain classification on K16 and K19. The supervised approach uses ImageNet pre-trained weights. The self-supervised baselines are trained from scratch. After self-supervised training, we freeze the weights, add a linear classifier on top, and train it until convergence. For SupContrast [Khosla et al., 2020] we jointly train the representation and the classification as described in the original paper.

We find that MoCoV2 [Chen et al., 2020a] outperforms the two other SOTA approaches. On K16, the model gains up to $10\%$ in terms of the F1-score with respect to the other self-supervised baselines. In addition, MoCoV2 gives competitive results with the supervised baseline that is initialized with ImageNet weights. It shows that MoCoV2 is able to efficiently learn from unlabeled data and create a generalized feature space. This mainly comes from the combination of the momentum encoder and the access to a large number of negative samples. Hence, we choose to adapt MoCoV2 for our proposed UDA method.

## B    Patch Classification - t-SNE Projection

In this section, we present the complementary results to the ones in Section 4.2 for patch classification. The embeddings of all baselines and our proposed approach are displayed in Figure 11 using t-SNE visualization. We show the alignment between the source (K19) and target (K16) embedding domain, as well as classes-wise.

With the source only approach, we can observe the lack of domain alignment between the feature spaces. Here, the model learns two distinct distributions for each set. On the other side, our approach shows a satisfactory alignment of domains compared to most baselines. The target complex stroma (K16) is linked to tumor, debris, lymphocytes, and stroma in the source domain (K19).

## C    Multi-Source Dataset Sampling Ratio

Table 6: Classification results of the different SOTA self-supervised approaches, as well as the supervised baseline on the Kather-19 (K19) and Kather-16 (K16) patch classification tasks. We present the results for different percentages of available training data. The top results are highlighted in bold. We report the weighted F1 score.

| | Kather-16 Labels fraction | | | Kather-19 Labels fraction | | |
|---|---|---|---|---|---|---|
| Methods | 10% | 20% | 50% | 1% | 2% | 5% |
| Supervised[‡] | 85.8[**] | 86.5[**] | 87.9[**] | 89.2[+] | 89.9[+] | 90.5[+] |
| SimCLR [Chen et al., 2020b] | 79.6[**] | 78.9[**] | 78.6[**] | 76.9[**] | 79.4[**] | 80.7[**] |
| SupContrast [Khosla et al., 2020] | 60.8[**] | 73.2[**] | 80.8[**] | 78.7[**] | 81.6[**] | 85.0[**] |
| MoCoV2 [Chen et al., 2020a] | **88.5** | **90.2** | **91.1** | **89.9** | **90.3** | **90.6** |

[‡] Model initialized with ImageNet pre-trained weights.
[+] $p \geq 0.05$; [*] $p < 0.05$; [**] $p < 0.001$; unpaired t-test with respect to the top result.

When performing multi-source domain adaptation, we assume that the distribution of all the source and target samples are the same. More formally, we have $p(\mathbf{x} \in \mathcal{D}_s) = Kp(\mathbf{x} \in \mathcal{D}_s^k) = p(\mathbf{x} \in \mathcal{D}_t)$. This section, we analyze the importance of balancing the source and target domains during the pre-training stage. We use K19 and K16 as source datasets and CRC-TP the target dataset. For K19 and K16, only $1\%$ and $10\%$ of the source labels are used, respectively.
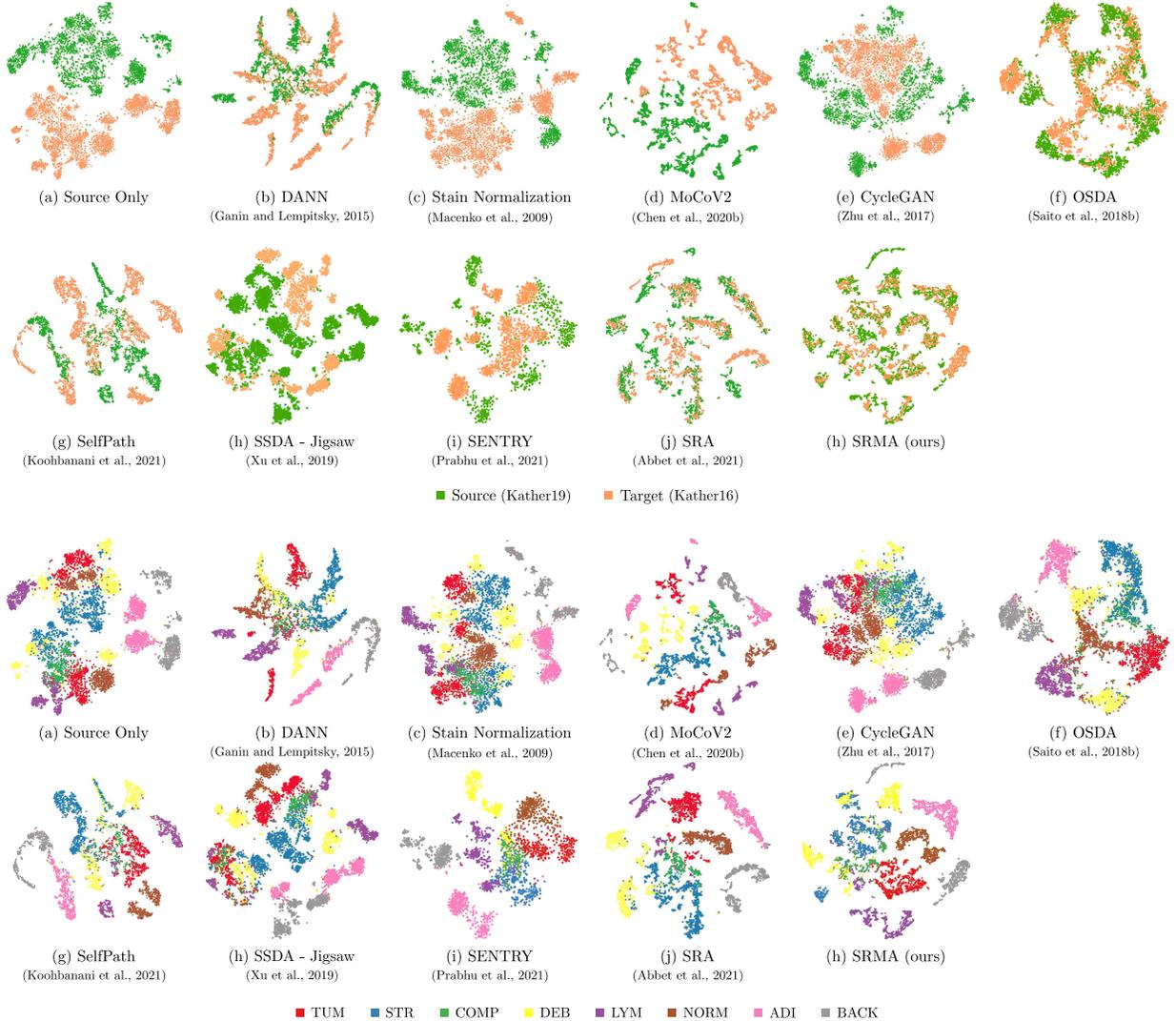
Figure 11: t-SNE projection of the source (Kather-19) and target (Kather-16) domain embeddings. We show the alignment of the embedding space as well as the individual classes for all presented models between the source and target domain. The classes of Kather-19 are merged and relabeled according to the definitions in Kather-16. The standard supervised approach is depicted in (a). We compare our approach (i) to other domain adaptation methods (b-j). Our approach (h) qualitatively show the best alignment between the source and target domain.

The results of the classification performance on the CRC-TP dataset are presented in Table 7. We indicate the multi-source scenario ($1 : 1$ or $K : 1$), the sampling probability for each of the dataset, and the batch size.

The cross-domain matching using the $K : 1$ scenario shows the highest variance and its performances can change up to 2.6%. Overall, we can observe that balanced probability between all sets, namely $\frac{1}{3}$ each, gives similar results across all multi-source scenarios. In addition, when lowering the sampling probability of K16 we can see a drop in performances. This suggests that it is important to have a balanced sampling strategy even if one of the source sets (e.i., K16 with 5,000 examples) is much smaller.

# D  Multi-source - t-SNE Projection

Figure 12 shows the visualization of the embedding for the proposed multi-source domain adaptation in Sections 4.7. It highlights the alignment of the feature space between the two source sets (K19, CRC-TP) and our in-house dataset.

Table 7: Study of the multi-source domain performance of the Self-Rule to Multi-Adapt (SRMA) approach with different sampling ratios. We use K19 and K16 as source datasets and CRC-TP as the target dataset. For K19 and K16, only $1\%$ and $10\%$ of the source labels are used, respectively. For the proposed SRMA model we compare the introduced multi-source approaches defined in Equations 12-15, where $1:1$ and $K:1$ refers to the one-to-one and $K$-to-one setting, respectively. The probability of sampling an example from each set within a batch is indicated. We report the F1 score for the individual classes and weighted F1 score as the overall mean performance (all) averaged over 10 runs.

| Model | Multi-source | | Sampling probability | | | Batch size | TUM | STR$^{\ddagger}$ | LYM | NORM | DEB | ALL |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\mathcal{L}_{\text{IND}}$ | $\mathcal{L}_{\text{CRD}}$ | K19 | K16 | CRCTP | | | | | | | |
| DeepAll [Dou et al., 2019] | - | - | - | - | - | 128 | 72.4$^{**}$ | 88.6$^{**}$ | 43.6$^{**}$ | 53.2$^{**}$ | 71.8$^{**}$ | 73.2$^{**}$ |
| SRA[Abbet et al., 2021] | 1 : 1 | 1 : 1 | 0.25 | 0.25 | 0.50 | 128 | 86.2$^{**}$ | 87.6$^{**}$ | 66.7$^{**}$ | 71.0$^{**}$ | **80.5** | 81.8$^{**}$ |
| SRMA | 1 : 1 | 1 : 1 | 0.25 | 0.25 | 0.50 | 128 | **92.5** | 88.4$^{**}$ | 68.7$^{**}$ | 68.3$^{**}$ | 74.2$^{*}$ | 82.9$^{*}$ |
| SRMA | $K$ : 1 | 1 : 1 | 0.25 | 0.25 | 0.50 | 128 | 91.5$^{*}$ | 87.6$^{**}$ | **70.7** | **75.0** | 65.7$^{**}$ | 82.7$^{*}$ |
| SRMA | 1 : 1 | $K$ : 1 | 0.25 | 0.25 | 0.50 | 128 | 90.1$^{**}$ | **90.1** | 69.6$^{+}$ | 72.9$^{**}$ | 71.6$^{**}$ | **83.6** |
| SRMA | $K$ : 1 | $K$ : 1 | 0.25 | 0.25 | 0.50 | 128 | **91.6** | 87.4$^{**}$ | 68.7$^{**}$ | 73.9$^{**}$ | 53.3$^{**}$ | 81.2$^{**}$ |
| SRMA | 1 : 1 | 1 : 1 | 0.33 | 0.33 | 0.33 | 128 | **92.9**$^{+}$ | 87.8$^{**}$ | 68.3$^{**}$ | 65.3$^{**}$ | 72.0$^{*}$ | 82.0$^{**}$ |
| SRMA | $K$ : 1 | 1 : 1 | 0.33 | 0.33 | 0.33 | 128 | **93.1** | 87.3$^{**}$ | 70.5$^{**}$ | **78.3** | 66.9$^{**}$ | 83.4$^{*}$ |
| SRMA | 1 : 1 | $K$ : 1 | 0.33 | 0.33 | 0.33 | 128 | 92.5$^{*}$ | **89.7** | **71.6** | 73.0$^{**}$ | 66.9$^{**}$ | **83.8** |
| SRMA | $K$ : 1 | $K$ : 1 | 0.33 | 0.33 | 0.33 | 128 | 92.2$^{*}$ | 88.6$^{**}$ | 66.1$^{**}$ | 74.3$^{**}$ | **74.5** | 83.4$^{*}$ |
| SRMA | 1 : 1 | 1 : 1 | 0.40 | 0.20 | 0.40 | 128 | 90.5$^{**}$ | 88.3$^{**}$ | 63.8$^{**}$ | 71.8$^{**}$ | 66.1$^{**}$ | 81.5$^{**}$ |
| SRMA | $K$ : 1 | 1 : 1 | 0.40 | 0.20 | 0.40 | 128 | 90.8$^{**}$ | **89.8** | 62.0$^{**}$ | **74.7** | 64.1$^{**}$ | 82.2$^{**}$ |
| SRMA | 1 : 1 | $K$ : 1 | 0.40 | 0.20 | 0.40 | 128 | 92.0$^{*}$ | 88.6$^{**}$ | **69.5** | 73.7$^{**}$ | 64.8$^{**}$ | 82.8$^{**}$ |
| SRMA | $K$ : 1 | $K$ : 1 | 0.40 | 0.20 | 0.40 | 128 | **92.7** | 89.3$^{**}$ | 65.8$^{**}$ | 74.7$^{+}$ | **75.2** | **83.8** |

$^{\ddagger}$ The STR and MUS classes are merged as STR class; DEB and MUC classes as DEB.
$^{+}$ $p \geq 0.05$; $^{*}$ $p < 0.05$; $^{**}$ $p < 0.001$; unpaired t-test with respect to top

We observe that for each source domain, the categories are well clustered. Moreover, we notice that the classes shared by both domains (e.i., namely tumor, stroma, debris, lymphocytes, normal mucosa, and muscle) fully overlap. In addition, the tissues that are domain-specific (e.i., adipose, background, mucin, and complex stroma) form individual groups. Subsequently, it indicates that our approach was able to properly correlate similar tissue definitions across the source domains while maintaining domain-specific tissue representation.

Looking at the source and target projection, we discern a batch of tissue (center-top) that does not align with the source domain. When associated with the patches visualization, we can recognize tiles that include loose stroma, collagen, or blood vessels representation. Rightfully, none of the mentioned classes were present in the source domain, thus proving the usefulness of the easy-to-hard approach.

# E    Patch-based Segmentation of WSI from the TCGA cohort

In this section we highlight the performance of our framework on a publicly available WSI (UUDI: 2d961af6-9f08-4db7-92b2-52b2380cd022) from the TCGA colon cohort [Shanah et al., 2016a,b]. We apply our trained SRMA framework, as described in section 4.3, where K19 is used as the source domain and our in-house domain as the target one. We show the original image, as well as the classification output and the tumor class probability map of our proposed SRMA method.

The model is able to accurately classify tissue across the whole slide. Moreover, the pipeline gives a rather detailed output which is a remarkable performance for a patch-based approach that was not specifically designed for segmentation purposes. Moreover, the model is agnostic to artifacts such as the permanent marker spots (green marks on the bottom left). The tumor prediction map gives an overview of the tumor class probability across the WSI. This class is of particular interest, as tumor detection is an important step for many downstream tasks, e.g., detection of the invasive front or the tumor stroma ratio.

(b) t-SNE Kather19 labeled source sample projection

(c) t-SNE CRCTP labeled source sample projection

(a) t-SNE visualization of target patches distribution
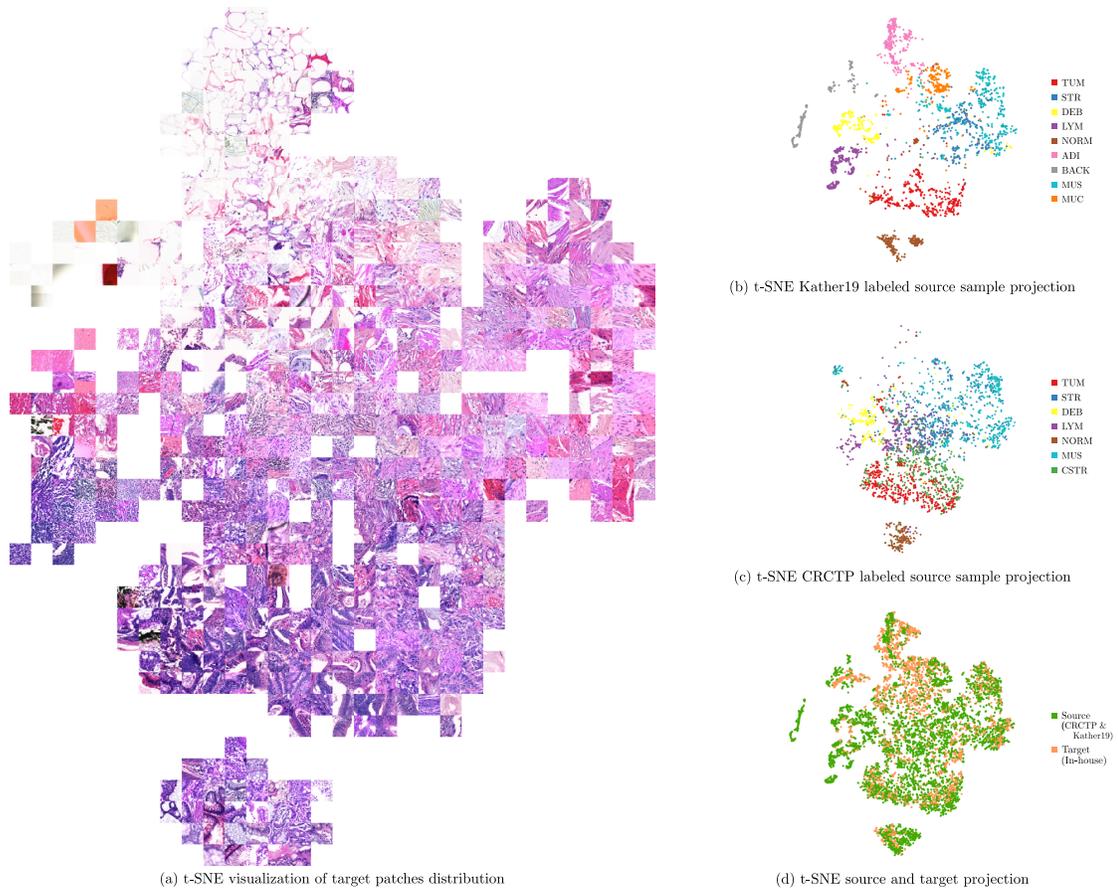
(d) t-SNE source and target projection

Figure 12: t-SNE visualization of the SRMA model trained on CRC-TP, K19 and the in-house dataset. All sub-figures depict the same embedding. (a) Patch-based visualization of the embedding. (b-c) Distribution of the labeled source samples. (d) Relative alignment of the source and target domain samples.
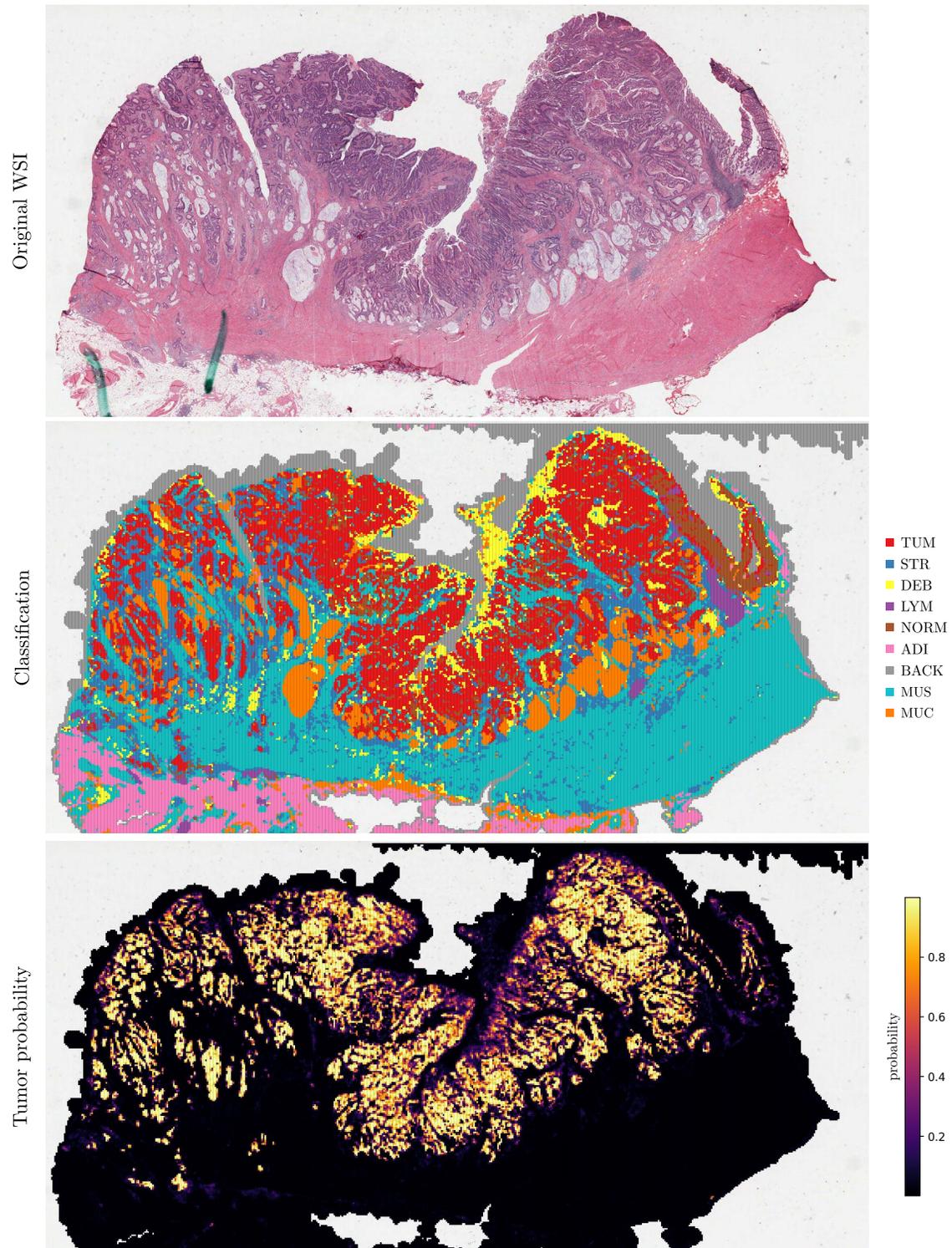
Figure 13: Segmentation results on a sample WSI from the TCGA cohort achieved by our SRMA model trained using K19 as the source dataset and our in-house set as the target dataset. From top to bottom, we show the original image, the classification output, and tumor class probability map.