

# Anatomy of Domain Shift Impact on U-Net Layers in MRI Segmentation

Ivan Zakazov <sup>\*1,2</sup>, Boris Shirokikh <sup>\*2</sup>, Alexey Chernyavskiy<sup>1</sup>, and Mikhail Belyaev<sup>2</sup>

<sup>1</sup> Philips Research, Moscow, Russia

<sup>2</sup> Skolkovo Institute of Science and Technology, Moscow, Russia  
[boris.shirokikh@skoltech.ru](mailto:boris.shirokikh@skoltech.ru)

**Abstract.** Domain Adaptation (DA) methods are widely used in medical image segmentation tasks to tackle the problem of differently distributed train (source) and test (target) data. We consider the supervised DA task with a limited number of annotated samples from the target domain. It corresponds to one of the most relevant clinical setups: building a sufficiently accurate model on the minimum possible amount of annotated data. Existing methods mostly fine-tune specific layers of the pretrained Convolutional Neural Network (CNN). However, there is no consensus on which layers are better to fine-tune, e.g. the first layers for images with low-level domain shift or the deeper layers for images with high-level domain shift. To this end, we propose SpotTUnet – a CNN architecture that automatically chooses the layers which should be optimally fine-tuned. More specifically, on the target domain, our method additionally learns the policy that indicates whether a specific layer should be fine-tuned or reused from the pretrained network. We show that our method performs at the same level as the best of the non-flexible fine-tuning methods even under the extreme scarcity of annotated data. Secondly, we show that SpotTUnet policy provides a layer-wise visualization of the domain shift impact on the network, which could be further used to develop robust domain generalization methods. In order to extensively evaluate SpotTUnet performance, we use a publicly available dataset of brain MR images (CC359), characterized by explicit domain shift. We release a reproducible experimental pipeline<sup>3</sup>.

**Keywords:** Domain Adaptation, Deep Learning, MRI, Segmentation

## 1 Introduction

Whenever a model, trained on one distribution, is given some data, belonging to another distribution, a detrimental effect called *domain shift* might pop up and decrease the inference quality [18]. This problem is especially acute in the field of

---

\* equal contribution

<sup>3</sup> <https://github.com/neuro-ml/domain-shift-anatomy>

medical imaging since any data collection instance (e.g., MRI apparatus) might appear to sample data, belonging to the domain of its own due to peculiarities of a model or the scanning protocol [5]. Besides, preserving quality on the new data (i.e., *domain adaptation*) is of utter importance because of the industry standards. A realistic set-up is that of *supervised domain adaptation (sDA)*: some data from the new (*target*) domain is labeled and should be utilized for fine-tuning the model, pre-trained on the *source* domain.

The central question of sDA research is *how* should the net be fine-tuned. A great number of works adopt the transfer learning approach of fine-tuning the last layers only, which is underpinned by the notion of feature complexity increasing with depth [19]. It is assumed, that low-level features should be shared across domains, while high-level ones are more prone to domain shift and should be therefore fine-tuned. However, in a number of Domain Adaptation papers the presence of low-level domain shift is demonstrated [3, 12, 21].

SpotTune [6] is a *Transfer Learning (TL)* approach, allowing for Adaptive Fine-tuning, providing, therefore, a domain shift stamp of each task by learning the corresponding fine-tuning policy. We employ SpotTune approach for getting better insight into the *anatomy of domain shift* in the problem of Domain Adaptation.

Our contribution is threefold:

- To the best of our knowledge, we are the first to propose SpotTUNet: SpotTune adapted for supervised DA in medical image segmentation
- We introduce interpretable regularization, with which SpotTune performs on par with the best of the alternative fine-tuning methods across the entire data availability spectrum
- We study the optimal fine-tuning strategies for different data availability scenarios and provide intuition for the obtained results

## 2 Related work

In this paper, we focus on the supervised DA setup in medical image segmentation. Our approach is motivated by SpotTune [6], which learns a policy to choose between fine-tuning and reusing pretrained network layers and at the same time fine-tunes chosen layers on the Target data (see details in Sec. 3). Previous approaches in both supervised and unsupervised medical image DA mostly rely on the explicit choice of layers to fine-tune or to split the adversarial head from. We detail the most relevant approaches below.

The authors of [19] have extensively evaluated features transferability in the image classification task. Their study suggests that features from the first layers of the network could be transferred between tasks, while the last layers should be re-trained. Contrary, convolutional filter reconstruction [1] is designed to tackle the low-level domain shift under unsupervised DA setup showing that the first layers are susceptible to domain shift more than the later ones.

DA methods for medical image segmentation also tackle domain shift problem differently: the earlier approaches follow the motivation of [19], thus fine-tune the

later layers of the network. According to the approach of [10] fine-tuning of only the last CNN layer is performed, which yields improvement over transferring without adaptation. However, no comparison is provided with other supervised DA methods or various transferring strategies. Similarly, in [15] the last CNN layer is fine-tuned, but the authors focus more on the training cases selection procedure rather than on the fine-tuning method development. Several works [4, 16] provide a comparison between the outcomes of various numbers of the later layers being fine-tuned. Notably, in [4] fine-tuning of the whole network is added to comparison, with better results demonstrated for a smaller number of the later layers fine-tuned. In the unsupervised DA setup, [8] achieves better results adapting features from all the layers except for the first block, but the layers choice strategy remains unlearnable.

In contrast, later approaches follow the motivation of [1], arguing that medical images (e.g. MRI) mostly contain low-level domain shift, which is due to varying intensity profiles but similar high-level structures, thus the first layers should be targeted. In [12] fine-tuning of the first layers is compared to fine-tuning of the last layers and of the whole network. The conclusion is that fine-tuning of the first layers is superior to fine-tuning of the last ones and is even preferable to fine-tuning of the whole network in the case of annotated Target data scarcity. In [9], an adaptive image normalization module is developed for the unsupervised test-time DA, which is conceptually close to fine-tuning of the first layers. Approaches of [3, 21] are also motivated by the same hypothesis and adapt the first layers in the unsupervised DA setup.

To this end, we compare SpotTUNet with the best unlearnable layer choice strategies within the supervised DA setup and show it to be a reliable tool for domain shift analysis. While authors of [13] demonstrate SpotTune to perform worse than histogram matching preprocessing in the medical image classification task, we argue that histogram matching is a task-specific method and show its extremely poor segmentation quality in our task. Many approaches competitive to SpotTune have been developed recently, but their focus is more narrow: obtaining the best score on the Target domain rather than analyzing domain shift properties. Therefore, we further study only the SpotTune-based approach.

### 3 Method

The majority of supervised DA methods in medical image segmentation are based on the explicit choice of layers to fine-tune. However, as indicated in Sec.2, it is not always clear, whether the first or the last layers are to be targeted. Moreover, state-of-the-art architectures consist of skip-connections and residual paths [2], while residual networks behave like ensembles of *shallow* networks [17]. Therefore, it is also unclear which layers are actually the first across the most meaningful shallow sub-parts of the residual network.

We introduce an extension of SpotTune [6] on the supervised DA for medical image segmentation called **SpotTUNet**. SpotTUNet consists of two copies of the main (segmentation) network and a policy network (see Fig. 1). The main

network is pretrained on the Source domain and then duplicated: the first copy has frozen weights (Fig. 1, blue blocks), while the second copy is fine-tuned on the Target domain (Fig. 1, orange blocks). The policy network predicts  $N$  pairs of logits for each of  $N$  segmentation network blocks (residual blocks or separate convolutions). For each pair of logits, we apply softmax and interpret the result as probabilities in a 2-class classification task: class 0 corresponds to the choice of a frozen block, while class 1 means choosing to fine-tune the unfrozen copy. Then, for the  $l$ -th level of the network (frozen block is denoted  $F_l$  and fine-tuned block  $\tilde{F}_l$ ) we define its output as  $x_l = I_l(x)F_l(x_{l-1}) + (1 - I_l(x))\tilde{F}_l(x_{l-1})$ , where  $I_l(x)$  is the indicator of choosing the frozen block (i.e. class 0 probability  $> 0.5$ ). Here, we use Gumbel-Softmax to propagate the gradients through the binary indicator  $I_l(x)$  exactly reproducing the methodology of SpotTune [6]. Thus, we simultaneously train the policy network and fine-tune the duplicated layers.

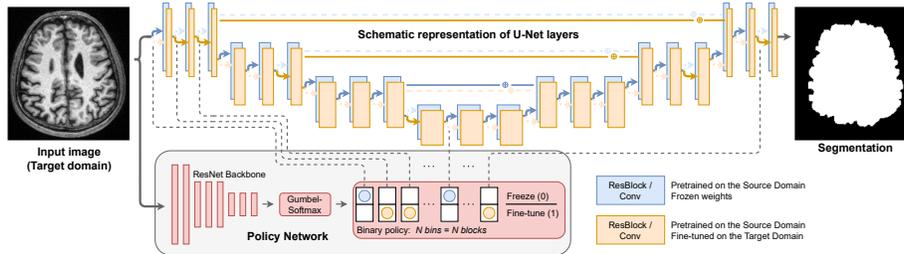


Fig. 1: SpotTUNet architecture for the supervised DA in medical image segmentation. We use U-Net architecture from [12] as the segmentation backbone, which is pretrained on the Source domain. The pretrained segmentation network is frozen (blue blocks) and has a copy (orange blocks) that is fine-tuned on the Target domain. The policy network is simultaneously trained on the Target domain to output binary decisions for each pair of blocks from the segmentation networks: use the frozen block (blue) *vs* use the fine-tuned block (orange).

Authors of SpotTune also propose a compact global policy (Global-k), which, via additional losses, constrains all the images to fine-tune the same  $k$  blocks and aims at reducing the memory and computational costs while possibly decreasing the overall quality. We propose a simplified regularization approach that moreover yields higher quality in the cases of annotated data scarcity. Effectively, we apply  $\mathbb{L}_1$  regularization term in the Global-0 SpotTune case, simultaneously minimizing the total number of fine-tuned blocks and achieving a more deterministic policy (exactly 0 or 1, which is due to  $\mathbb{L}_1$  properties). The resulting loss is

$$\mathcal{L} = \mathcal{L}_{segm} + \lambda \sum_{l=1}^N (1 - I_l(x)), \quad (1)$$

where  $\lambda$  is the balance parameter of the regularization term. Our motivation is different from the original: we assume that fewer blocks should be optimally fine-tuned in case of limited annotated data so that to avoid possible overfitting. Our regularization has only 1 parameter compared to 3 parameters of the original SpotTune regularization, and this parameter ( $\lambda$ ) could be optimized during preliminary validation. The intuition behind  $\lambda$  is simple: the less annotated data is available the larger  $\lambda$  value should be.

## 4 Experiments

### 4.1 Technical details

**Data.** We report our results on a publicly available dataset CC359 [14]. CC359 consists of 359 MR images of head with the task being skull stripping. The dataset is split into 6 equal domains which are shown [12] to contain domain shift resulting in a severe score deterioration. The data preprocessing steps are: interpolation to  $1 \times 1 \times 1$  mm voxel spacing and scaling intensities into 0 to 1 interval. All splits and setups are detailed in Sec. 4.2.

**Metric.** To evaluate different approaches, we use surface Dice Score [11] at the tolerance of 1 mm. While preserving the consistency with the methodology of [12], we also find surface Dice Score to be a more representative metric for the brain segmentation task than the standard Dice Score.

**Architecture and training.** The experimental evaluation provided in [12] shows that neither architecture nor training procedure variations, e.g. augmentation, affect the relative performance of conceptually different approaches. Therefore, in all our experiments we consistently use 2D U-Net architecture implementation from [12]. We also use a ResNet architecture [7] as the policy network backbone. In all the experiments we minimize Binary Cross-Entropy loss ( $\mathcal{L}_{segm}$ ) via stochastic gradient descent. Baseline and oracle (see Sec. 4.2) are trained for 100 epochs (100 iterations per epoch) with the learning rate of  $10^{-2}$  reduced to  $10^{-3}$  at the 80-th epoch. All fine-tuning methods are trained for 60 epochs with the learning rate of  $10^{-3}$  reduced to  $10^{-4}$  at the 45-th epoch. We ensure all the models reach the loss plateau. All models are trained with batch size 16. The training takes about 4 hours on a 16GB nVidia Tesla V100 GPU [20].

### 4.2 Experimental setup

*Baseline and oracle* The models trained on a single domain form the *baseline* of our study. The transfer of such a model (without fine-tuning) on the other 5 unseen domains results in quality deterioration, which is also shown in [12]. We also obtain scores within each domain (the *oracle*) via 3-fold cross-validation, thereby setting the upper bound for various DA methods.

*SpotTUNet validation* Six domains yield 30 Source-Target pairs, thus, 30 supervised DA experiments. In order to avoid overfitting, we separate one Source domain (Siemens, 1.5T) and, correspondingly, 5 Source-Target pairs for SpotTUNet

validation and use the other 25 pairs for testing various DA approaches. On the 5 validation pairs, we firstly adjust the temperature parameter of Gumbel-Softmax ( $\tau$ ) via grid-search over  $\tau \in \{.01, .1, .5, 1, 2, 5\}$ . The number of annotated slices from the Target domain, in this case, is 270 (one 3D image). Secondly, we search for the optimal  $\lambda$  for each amount of annotated Target data considered via grid-search over  $\lambda \in \{0, 1, 3, 5, 7, 10, 12, 15, 20 (\times 10^{-3})\}$ . In both validation and testing experiments, we study the same setups of the annotated Target data scarcity: 8, 12, 24, 45, 90, 270, and 800 slices available for fine-tuning. The optimal value of  $\lambda$  is fixed for each data scarcity setup and used for SpotTUNet when testing on the remaining 25 pairs.

*Supervised DA methods* On the rest of the 25 testing pairs, we compare 4 methods: *fine-tuning of the first network layers*, *fine-tuning of the whole network* from [12], *histogram matching* from [13], and SpotTUNet. We load a *baseline* model pretrained on the corresponding Source domain and then fine-tune it via one of the methods or preprocess the Target data in case of histogram matching. We compare methods by averaging surface Dice Scores over the Target images, separated for test (omitted when fine-tuning).

### 4.3 Results and discussion

We firstly find and fix hyperparameters for SpotTUNet through validation: the temperature of Gumbel-Softmax is set to  $\tau = 0.1$  and the optimal regularization  $\lambda$  is set to the optimal value for each data scarcity setup (see Fig. 2). We also note that Gumbel-Softmax training stability is extremely sensitive to  $\tau$  choice, thus suggest to validate different values of  $\tau$  at one of the first stages of deploying SpotTune-like architectures. Positive regularization term benefits almost all data scarcity setups: surface Dice Score of the optimal  $\lambda$  is significantly higher ( $p < 10^{-3}$ , one-sided Wilcoxon signed-rank test) than the surface Dice Score of the corresponding model without the regularization. The only exception is the case of 800 available Target slices (or three 3D images), where the optimal  $\lambda$  is close to 0, and the quality drops with the increase of  $\lambda$  (see Fig. 2). We conclude that while SpotTUNet learns the optimal policy without regularization when there is no Target data shortage, regularization improves DA performance significantly in case of Target data scarcity.

We further compare SpotTUNet with *Fine-Tuning All Layers*, *Fine-Tuning the First Layers*, and *histogram matching*. We present both distributions of the surface Dice Score (violin plots) and their average values (corresponding lines) in Fig. 3. Histogram matching achieves only 0.29 average surface Dice Score, which is even lower than the *baseline* average score of 0.55; we exclude both methods from the comparison in Fig. 3. Here, we show that SpotTune performs at the same level as the best of the other methods regardless of the Target data scarcity severity.

Finally, we track SpotTUNet policy on the test data: for each layer, we calculate the frequency of choosing the fine-tuned one instead of the pretrained and frozen one. The layer-wise visualization is given in Fig. 4. We find that blocks

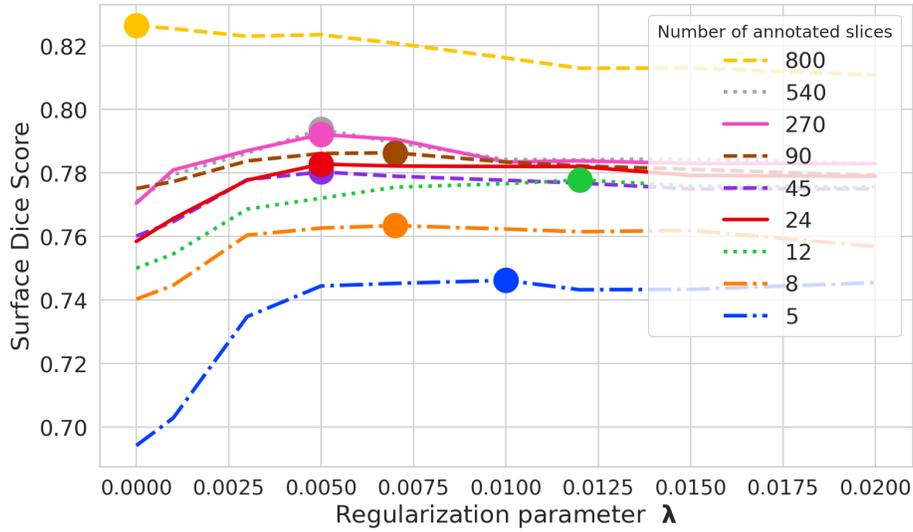


Fig. 2: The validation performance of SpoTUNet dependence on the regularization parameter  $\lambda$ . In each point, the average surface Dice Score over 5 validation experiments is calculated. Each line corresponds to some amount of available annotated data from the Target domain. The bold points indicate the optimal  $\lambda$  values in terms of surface Dice Score.

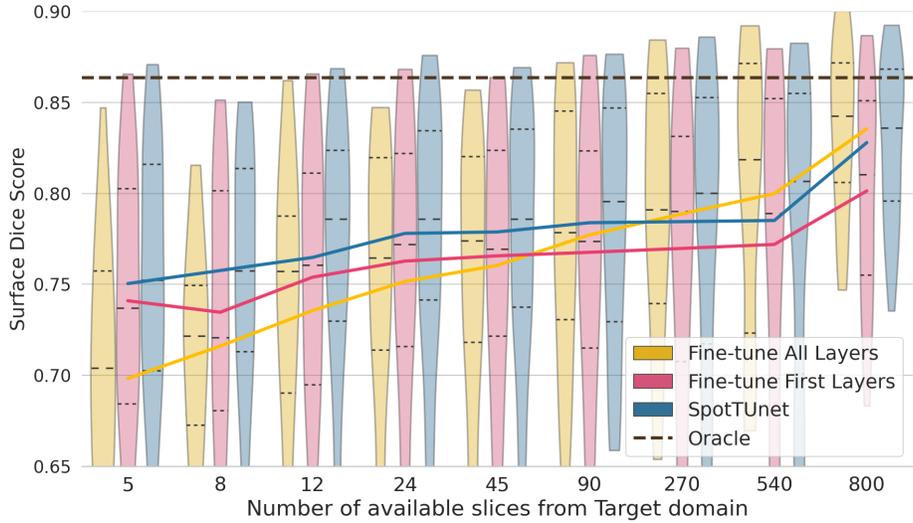


Fig. 3: The methods performance dependence on the amount of available Target data. Baseline and histogram matching yield poor quality, thus not included.

from the encoder part of U-Net are more likely to be fine-tuned, especially in the case of annotated data scarcity. However, these are not exactly the first layers (e.g., preserving the original resolution), which opposes the conclusion of [12]. We further hypothesize, that SpotTUNet policy indicates layers that should be fine-tuned for the optimal solution. Consequently, feature maps preceding these frequently fine-tuned layers might be marked with drastic domain shift. We note that it is worth evaluating if unsupervised DA approaches [8, 21] would benefit from passing these SpotTUNet indicated domain shift reach feature maps to the adversarial heads and leave this hypothesis validation for the future research.

We attribute the difference between the policies observed and those presented in the original SpotTune paper [6] (mostly final layers fine-tuned) to the fundamental difference between Transfer Learning (TL) and DA. In TL one deals with data of varying nature, thus the later layers should be addressed; in DA, the datasets contain semantically homogeneous data (e.g., brain MRI scans), thus domain shift is mostly low-level and the first layers should be targeted.

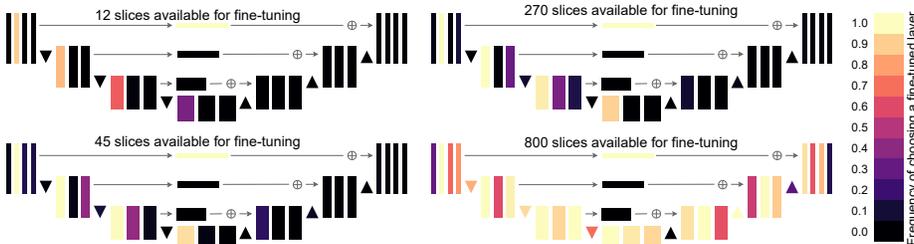


Fig. 4: SpotTUNet learnt policy visualization for the cases of 12 (upper-left), 45 (bottom-left), 270 (upper-right), and 800 (bottom-right) available Target slices. Colored blocks correspond to either residual blocks or convolutions. Triangular blocks are the convolutions that perform  $\times 2$  up- or down-sampling.

## 5 Conclusion

We propose a fine-tuning approach for supervised DA in medical image segmentation called SpotTUNet. Our experiments demonstrate SpotTUNet to preserve the quality of the alternative methods while eliminating the need for switching between various methods depending on the Target data availability. Besides, it learns automatically, which layers are to be optimally fine-tuned on the target domain, therefore providing a policy, indicative of the network layers most susceptible to domain shift. We believe that SpotTUNet generated policy might be used for developing more robust unsupervised DA methods, which is the goal of our future research.

## References

1. Aljundi, R., Tuytelaars, T.: Lightweight unsupervised domain adaptation by convolutional filter reconstruction. In: European Conference on Computer Vision. pp. 508–515. Springer (2016)
2. Bakas, S., Reyes, M., Jakab, A., Bauer, S., Rempfler, M., Crimi, A., Shinohara, R.T., Berger, C., Ha, S.M., Rozycki, M., et al.: Identifying the best machine learning algorithms for brain tumor segmentation, progression assessment, and overall survival prediction in the brats challenge. arXiv preprint arXiv:1811.02629 (2018)
3. Dou, Q., Ouyang, C., Chen, C., Chen, H., Heng, P.A.: Unsupervised cross-modality domain adaptation of convnets for biomedical image segmentations with adversarial loss. In: Proceedings of the 27th International Joint Conference on Artificial Intelligence. p. 691–697. IJCAI’18, AAAI Press (2018)
4. Ghafoorian, M., Mehrtash, A., Kapur, T., Karssemeijer, N., Marchiori, E., Pesteie, M., Guttman, C.R., de Leeuw, F.E., Tempny, C.M., Van Ginneken, B., et al.: Transfer learning for domain adaptation in mri: Application in brain lesion segmentation. In: International conference on medical image computing and computer-assisted intervention. pp. 516–524. Springer (2017)
5. Glocker, B., Robinson, R., Castro, D.C., Dou, Q., Konukoglu, E.: Machine learning with multi-site imaging data: An empirical study on the impact of scanner effects (2019)
6. Guo, Y., Shi, H., Kumar, A., Grauman, K., Rosing, T., Feris, R.: Spottune: transfer learning through adaptive fine-tuning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4805–4814 (2019)
7. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
8. Kamnitsas, K., Baumgartner, C., Ledig, C., Newcombe, V., Simpson, J., Kane, A., Menon, D., Nori, A., Criminisi, A., Rueckert, D., et al.: Unsupervised domain adaptation in brain lesion segmentation with adversarial networks. In: International conference on information processing in medical imaging. pp. 597–609. Springer (2017)
9. Karani, N., Erdil, E., Chaitanya, K., Konukoglu, E.: Test-time adaptable neural networks for robust medical image segmentation. *Medical Image Analysis* **68**, 101907 (2021)
10. Kushibar, K., Valverde, S., González-Villà, S., Bernal, J., Cabezas, M., Oliver, A., Lladó, X.: Supervised domain adaptation for automatic sub-cortical brain structure segmentation with minimal user interaction. *Scientific reports* **9**(1), 1–15 (2019)
11. Nikolov, S., Blackwell, S., Mendes, R., De Fauw, J., Meyer, C., Hughes, C., Askham, H., Romera-Paredes, B., Karthikesalingam, A., Chu, C., et al.: Deep learning to achieve clinically applicable segmentation of head and neck anatomy for radiotherapy. arXiv preprint arXiv:1809.04430 (2018)
12. Shirokikh, B., Zakazov, I., Chernyavskiy, A., Fedulova, I., Belyaev, M.: First u-net layers contain more domain specific information than the last ones. In: Domain Adaptation and Representation Transfer, and Distributed and Collaborative Learning, pp. 117–126. Springer (2020)
13. Singh, S., Matthews, T.P., Shah, M., Mombourquette, B., Tsue, T., Long, A., Almhosen, R., Pedemonte, S., Su, J.: Adaptation of a deep learning malignancy model from full-field digital mammography to digital breast tomosynthesis. In: *Medical Imaging 2020: Computer-Aided Diagnosis*. vol. 11314, p. 1131406. International Society for Optics and Photonics (2020)

14. Souza, R., Lucena, O., Garrafa, J., Gobbi, D., Saluzzi, M., Appenzeller, S., Ritter, L., Frayne, R., Lotufo, R.: An open, multi-vendor, multi-field-strength brain mr dataset and analysis of publicly available skull stripping methods agreement. *NeuroImage* **170**, 482–494 (2018)
15. Valindria, V.V., Lavdas, I., Bai, W., Kamnitsas, K., Aboagye, E.O., Rockall, A.G., Rueckert, D., Glocker, B.: Domain adaptation for mri organ segmentation using reverse classification accuracy. arXiv preprint arXiv:1806.00363 (2018)
16. Valverde, S., Salem, M., Cabezas, M., Pareto, D., Vilanova, J.C., Ramió-Torrentà, L., Rovira, À., Salvi, J., Oliver, A., Lladó, X.: One-shot domain adaptation in multiple sclerosis lesion segmentation using convolutional neural networks. *NeuroImage: Clinical* **21**, 101638 (2019)
17. Veit, A., Wilber, M.J., Belongie, S.J.: Residual networks behave like ensembles of relatively shallow networks. In: Lee, D.D., Sugiyama, M., von Luxburg, U., Guyon, I., Garnett, R. (eds.) *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*. pp. 550–558 (2016)
18. Wang, M., Deng, W.: Deep visual domain adaptation: A survey. *Neurocomputing* **312**, 135–153 (2018)
19. Yosinski, J., Clune, J., Bengio, Y., Lipson, H.: How transferable are features in deep neural networks? In: *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2*. p. 3320–3328. NIPS’14, MIT Press, Cambridge, MA, USA (2014)
20. Zacharov, I., Arslanov, R., Gunin, M., Stefonishin, D., Pavlov, S., Panarin, O., Maliutin, A., Rykovanov, S.G., Fedorov, M.: ‘zhores’ – petaflops supercomputer for data-driven modeling, machine learning and artificial intelligence installed in skolkovo institute of science and technology. *Open Engineering* **9**, 512 – 520 (2019)
21. Zhao, X., Sicilia, A., Minhas, D., O’Connor, E.E., Aizenstein, H., Klunk, W., Tudorascu, D., Hwang, S.J.: Robust white matter hyperintensity segmentation on unseen domain. 2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI) pp. 1047–1051 (2021)