

# TAAL: Test-time Augmentation for Active Learning in Medical Image Segmentation

Mélanie Gaillochet\*, Christian Desrosiers, and Hervé Lombaert

ETS Montreal, Canada

**Abstract.** Deep learning methods typically depend on the availability of labeled data, which is expensive and time-consuming to obtain. Active learning addresses such effort by prioritizing which samples are best to annotate in order to maximize the performance of the task model. While frameworks for active learning have been widely explored in the context of classification of natural images, they have been only sparsely used in medical image segmentation. The challenge resides in obtaining an uncertainty measure that reveals the best candidate data for annotation. This paper proposes Test-time Augmentation for Active Learning (TAAL), a novel semi-supervised active learning approach for segmentation that exploits the uncertainty information offered by data transformations. Our method applies cross-augmentation consistency during training and inference to both improve model learning in a semi-supervised fashion and identify the most relevant unlabeled samples to annotate next. In addition, our consistency loss uses a modified version of the JSD to further improve model performance. By relying on data transformations rather than on external modules or simple heuristics typically used in uncertainty-based strategies, TAAL emerges as a simple, yet powerful task-agnostic semi-supervised active learning approach applicable to the medical domain. Our results on a publicly-available dataset of cardiac images show that TAAL outperforms existing baseline methods in both fully-supervised and semi-supervised settings. Our implementation is publicly available on <https://github.com/melinphd/TAAL>.

## 1 Introduction

The performance of deep learning-based models improves as the number of labeled training samples increases. Yet, the burden of annotation limits the amount of data that can be labeled. One solution to that problem is offered by active learning (AL) [1]. Based on the hypothesis that all data samples have a different impact on training, active learning aims to find the best set of candidate samples to annotate in order to maximize the performance of the task model. In such context, medical image segmentation emerges as a remarkably relevant task for active learning. Indeed, medical images typically require prior expert knowledge for their analysis and annotation, an expensive and time-consuming task. Initial attempts have explored active learning in medical imaging [2], but

---

\* Corresponding author: M. Gaillochet. **Email:** [melanie.gaillochet.1@ens.etsmtl.ca](mailto:melanie.gaillochet.1@ens.etsmtl.ca)

their methodology either relied on simple uncertainty heuristics [3,4] or required heavy computations during sampling [5,6] or training [7].

**Deep active learning** Active learning has been extensively explored for the classification [8,9,10,11,12,13] or segmentation [14,15,16] of natural images. Recent deep active learning approaches based on entropy [12] or ensembles [9] adapted traditional uncertainty-based AL strategies to deep learning models. Similarly, DBAL [10] combined measures such as entropy or mutual information with Monte-Carlo dropout to suggest which samples to annotate next. Core-set selection [11] aimed to find the best batch sampling strategy for CNNs in classification, but did not scale well to high-dimensional data.

The use of auxiliary modules [13,17,18] has been similarly explored to improve AL sampling strategies. The loss prediction module of [13] measured model uncertainty with intermediate representations. Likewise, a VAE was used in VAAL [17] to learn the latent representation of the unlabeled dataset and distinguish between labeled and unlabeled samples. While these state-of-the-art methods have improved previous approaches, their dependence on auxiliary modules reduces their flexibility and increase the burden of hyperparameter tuning.

**Semi-supervised AL** Semi-supervised learning (SSL) exploits the representations of unlabeled data to improve the performance of the task model. Since semi-supervised learning and active learning are closely connected, recent works in AL have attempted to combine both domains [12,17,18,19]. For instance, CEAL [12] used pseudo-labeling of unlabeled samples to enhance the labeled set during training. VAAL [17] and TA-VAAL [18] employed a VAE to learn a latent representation of labeled and unlabeled data. The Mean Teacher framework of [19] combined a supervised loss on labeled data with an unsupervised loss on unlabeled data based on Temporal Output Discrepancy (TOD), evaluating the distance between the model’s output at different gradient steps. The model used a variant of TOD at sampling time to identify the most uncertain samples to annotate. However, these semi-supervised AL methods solely focused on classification tasks or the segmentation of natural images in very large quantities, which is a different context than medical imaging. Another recent work comparable to ours combined AL and SSL via consistency regularization [20]. The consistency loss adopted during training employed MixMatch [21] and sample selection measured inconsistency across input perturbations. However, as opposed to our work, [20] kept the consistency loss used during training and the AL inconsistency metric used for sample selection independent of each other, and the latter was quantified through variance. Furthermore, the method was only validated on classification tasks.

**Test-time augmentation** Data augmentation is a well-known regularization technique to improve generalization in low-data regimes. These augmentation techniques are particularly essential in medical imaging where datasets tend to be smaller than those of natural images. Yet most recent attempts in active learning do not exploit data augmentation during training [8,6], or only use random horizontal flipping [17,18]. Recent learning methods [22,23] have also investigated the use of augmentation at test-time in order evaluate prediction

uncertainty. Randomly augmented test images yield different model outputs. Combining these outputs can improve the overall predictions as well as generate uncertainty maps for these predictions. Uncertainty estimated through test-time augmentation was shown to be more reliable than model uncertainty measures such as test-time dropout or entropy of the output [23].

Motivated by the limitations of current active learning methods for medical image segmentation and the unused potential of active augmentation, this paper proposes a novel semi-supervised active learning strategy called Test-time Augmentation for Active Learning (TAAL).

**Our contribution:** Our method leverages the uncertainty information provided by data augmentation during both training and test-time sample selection phases. More specifically, TAAL employs a cross-augmentation consistency loss both to train the model in a semi-supervised fashion *as well as* to identify the most uncertain samples to annotate at the next cycle. TAAL comprises three key features:

1. a semi-supervised framework based on cross-augmentation consistency that exploits unlabeled samples during training and sampling;
2. a flexible task-agnostic sample selection strategy based on test-time augmentation;
3. a novel uncertainty measure based on a modified Jensen-Shannon divergence (JSD), which accounts for both cross-augmentation consistency and prediction entropy, and leads to improved performance.

## 2 Method

**Cross-augmentation consistency training** We consider a semi-supervised setting where we train a multi-class segmentation model  $f_\theta(\cdot)$  parameterized by  $\theta$  with  $N$  labeled samples and  $M$  unlabeled samples. We denote the labeled set as  $\mathcal{D}_L = \{(\mathbf{x}^{(j)}, \mathbf{y}^{(j)})\}_{j=1}^N$  and the unlabeled set as  $\mathcal{D}_U = \{\mathbf{x}_u^{(j)}\}_{j=1}^M$ , with data  $\mathbf{x}, \mathbf{x}_u \in \mathbb{R}^{H \times W}$  and segmentation mask  $\mathbf{y} \in \mathbb{R}^{C \times H \times W}$  ( $C$  is the number of classes).

The overall loss that we optimize,  $\mathcal{L} = \mathcal{L}_s + \lambda \mathcal{L}_c$ , is a combination of a supervised segmentation loss  $\mathcal{L}_s$  and an unsupervised consistency loss  $\mathcal{L}_c$  weighted by a factor  $\lambda$ . More explicitly, the objective is defined as

$$\mathcal{L} = \frac{1}{N} \sum_{j=1}^N \mathcal{L}_s(f_\theta(\mathbf{x}^{(j)}), \mathbf{y}^{(j)}) + \frac{\lambda}{M} \sum_{j=1}^M \mathcal{L}_c(f_\theta(\mathbf{x}_u^{(j)}), \Gamma), \quad (1)$$

where  $\Gamma$  are the transformations applied to  $\mathbf{x}_u^{(j)}$ . At each iteration, we apply a series of random transformations  $\{\Gamma_1, \dots, \Gamma_K\}$  to  $\mathbf{x}_u$ .  $\mathcal{L}_c$  measures the variability of segmentation predictions for different augmentations of  $\mathbf{x}_u$  measured by a function  $\mathcal{D}iv$ :

$$\mathcal{L}_c(f_\theta(\mathbf{x}_u^{(j)}), \Gamma) = \mathcal{D}iv\{\Gamma_1^{-1}[f_\theta(\Gamma_1(\mathbf{x}_u^{(j)}))], \dots, \Gamma_K^{-1}[f_\theta(\Gamma_K(\mathbf{x}_u^{(j)}))]\}. \quad (2)$$

While different measures can be used for *Div* [24], our consistency loss builds on the Jensen Shannon divergence (JSD),

$$\text{JSD}(P_1, \dots, P_K) = H\left(\frac{1}{K} \sum_{i=1}^K P_i\right) - \frac{1}{K} \sum_{i=1}^K H(P_i), \quad (3)$$

where  $H(P_i)$  is the Shannon entropy [25] for the probability distributions  $P_i$ . Minimizing the JSD reduces the entropy of the average prediction (making the predictions more similar to each other) while increasing the average of individual prediction entropies (ensuring confident predictions). In AL we typically want to select samples which have a high output entropy [12]. Selecting samples with highest JSD would thus have the opposite effect. To avoid this issue, and to control the relative importance of average prediction entropy versus entropy of individual predictions, we propose a weighted version of JSD with parameter  $\alpha$ .

$$\text{JSD}_\alpha(P_1, \dots, P_K) = \alpha H\left(\frac{1}{K} \sum_{i=1}^K P_i\right) - \frac{(1-\alpha)}{K} \sum_{i=1}^K H(P_i). \quad (4)$$

Note that using  $\alpha = 0.5$  is equivalent to using the standard JSD.

**Test-time augmentation sampling** In active learning, the goal is to select the best unlabeled samples to annotate after each training cycle to augment the next labeled training set. Hence, after each cycle, we apply our active learning strategy based on test-time augmentation to select the next samples to annotate.

For each sample  $\mathbf{x}_u \in \mathcal{D}_U$ , we apply a series of transformations  $\{\Gamma'_1, \dots, \Gamma'_{K_s}\}$ , and we compute an uncertainty score  $U_{\Gamma'}$  based on the same divergence function as the consistency loss:

$$U_{\Gamma'} = \text{JSD}_\alpha(\Gamma_1'^{-1}[f_\theta(\Gamma_1'(\mathbf{x}_u))], \dots, \Gamma_{K_s}'^{-1}[f_\theta(\Gamma_{K_s}'(\mathbf{x}_u))]). \quad (5)$$

The samples with highest uncertainty are annotated and added to the labeled training set. After sample selection, the model goes through a new training cycle.

## 3 Experiments and results

### 3.1 Implementation details

**Dataset** The publicly available ACDC dataset [26] comprises cardiac 3D cine-MRI scans from 100 patients. These are evenly distributed into 5 groups (4 pathological and 1 healthy subjects groups).

Segmentation masks identify 4 regions of interest: right-ventricle cavity, left-ventricle cavity, myocardium and background. For comparative purposes, our experiments focus on the MRI scans at the end of diastole. Preprocessing of the volumes includes resampling to a fixed  $1.0 \text{ mm} \times 1.0 \text{ mm}$  resolution in the x- and y-directions as well as a 99<sup>th</sup> percentile normalization. The 3-dimensional dataset of volumes are converted to a 2-dimensional dataset of images by extracting all

the z-axis slices for each volume. Each image is downsampled to  $128 \times 128$  pixels. Testing is performed on 181 images taken from 20 different patients, ensuring subjects are not split up across training and testing sets. The validation uses 100 randomly selected images. The same validation set is used for all experiments. In total, the available training set, both labeled and unlabeled, thus comprises 660 images.

**Implementation and training** We employ a standard 4-layer UNet [27] for our backbone segmentation model with dropout ( $p = 0.5$ ), batch normalization and a leaky ReLU activation function. For a fairer comparison in our experiments, we keep the number of training steps fixed during all cycles. We train our models for 75 epochs, each iterating over 250 batches, with  $BS = 4$ . We use the Adam optimizer [28], with  $LR = 10^{-6}$  and weight decay  $w = 10^{-4}$ . To improve convergence, we apply a gradual warmup with a cosine annealing scheduler [29,30], increasing the learning rate by a factor 200 during the first 10 epochs. During training, we apply data augmentation, using transformations similar to those utilized for the consistency loss.

In this work, we model the transformations  $\Gamma$  as a combination of  $f$ ,  $r$  and  $\epsilon$ , where  $f$  is the random variable for flipping the image along the horizontal axis,  $r$  is the number of  $90^\circ$  rotations in 2D, and  $\epsilon$  models Gaussian noise. We set  $f \sim \mathcal{U}(0, 1)$ ,  $r \sim \mathcal{U}(0, 3)$  and  $\epsilon \sim \mathcal{N}(0, 0.01)$ , and use  $K = 3$  transformations to compute the consistency loss during training.

We use the standard Dice loss as our supervised loss. In the semi-supervised case, following [31], we ramp-up the unsupervised component weight using a Gaussian ramp-up curve such that  $\lambda = \exp(-5(1-t/t_R)^2)$ , where  $t$  is the current epoch. We use a ramp-up length  $t_R$  of 10 epochs, corresponding to the learning rate gradual warmup length.

We repeat each experiment 5 times, each with a different seed determining different initialization of our model weights. For all experiments, the same initial labeled set is used for the first cycle. Experiments were run on NVIDIA PV100 GPU with CUDA 10.2 and Python 3.8.10. We implemented the methods using the PyTorch framework.

**Evaluation metrics** To evaluate the performance of the trained models, we employ the standard Dice similarity score, averaged over all non-background channels. We compute both the mean 3D Dice on test volumes and mean 2D Dice on the individual images from these volumes. We give the results as the mean Dice obtained over the repeated experiments.

### 3.2 Active learning setup

We begin each experiment with 10 labeled samples chosen uniformly at random in the training set and use a sampling budget of 1, meaning that we select one new sample to be labeled after each cycle. Following previous active learning validation settings [11], we retrain the model from scratch after each annotation

cycle. We use the same types of augmentations during training and sample selection. For test-time augmentation (TTA) sampling,  $\{I'_1, \dots, I'_{K_s}\}$  comprises all 8 combinations of flip and rotation augmentations, in order to apply similar transformations to all images, and adopts the same augmentation Gaussian noise parameters as for training. For comparative purposes, with dropout-based sampling, we also run 8 inferences with dropout to obtain different predictions. Both TTA and dropout-based sampling then evaluate uncertainty with  $U_{I'}$  computed on the different generated predictions. We set  $\alpha = 0.75$  in TAAL’s weighted JSD.

### 3.3 Comparison of active learning strategies

Our aim is to evaluate the effectiveness of our proposed semi-supervised active learning approach on a medical image segmentation task. In our active learning experiments, we compare TAAL and its unweighted version (with standard JSD) with random sampling, entropy sampling, sampling based on dropout and core-set selection. Entropy-based sampling selects the most uncertain samples based on the entropy of the output probabilities. Dropout-based sampling [10] identifies the samples with the highest JSD given multiple inferences with dropout. Finally, core-set selection [11] aims to obtain the most diverse labeled set by solving the maximum cover-set problem.

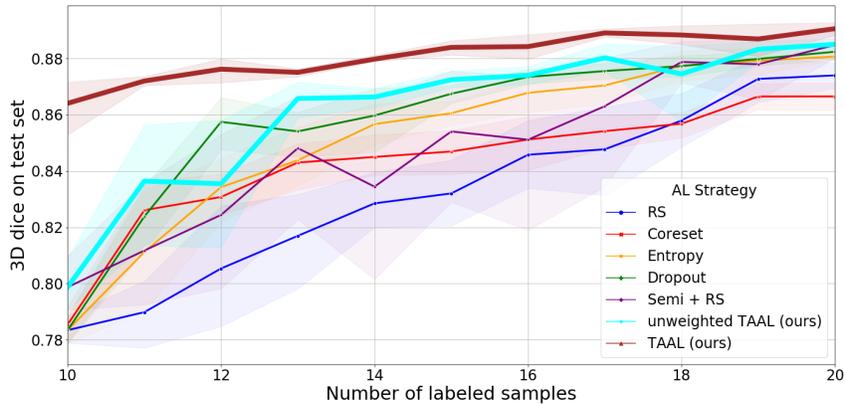


Fig. 1: Active learning results on the ACDC dataset, given as the mean 3D Dice scores on the test set and corresponding 95% confidence interval. In a fully-supervised setting: random sampling (RS), core-set selection (Coreset), uncertainty-based sampling based on entropy of output probabilities (Entropy), and uncertainty-based sampling based on JSD given multiple inferences with dropout (Dropout). In a semi-supervised setting: random sampling (Semi + RS), TAAL with standard JSD (unweighted TAAL), and TAAL with weighted JSD (TAAL). Our approach TAAL demonstrates significant improvements for low-data regimes in both fully and semi-supervised segmentation.

Figure 1 shows the segmentation performance of our proposed method with its 2 variants along with other existing active learning methods. TAAL consistently outperforms the other baselines by a large margin. We observe that our semi-supervised approach based on cross-augmentation consistency (Semi + RS) noticeably improves the fully-supervised vanilla model (RS). We notice that our unweighted version of TAAL (with standard JSD,  $\alpha=0.5$ ) already improves the performance of the semi-supervised model (Semi + RS) by selecting the most uncertain samples based on their cross-augmentation consistency loss. With higher  $\alpha=0.75$ , our proposed TAAL with weighted JSD yields the highest performance gain compared to the fully-supervised vanilla model with random sampling (RS).

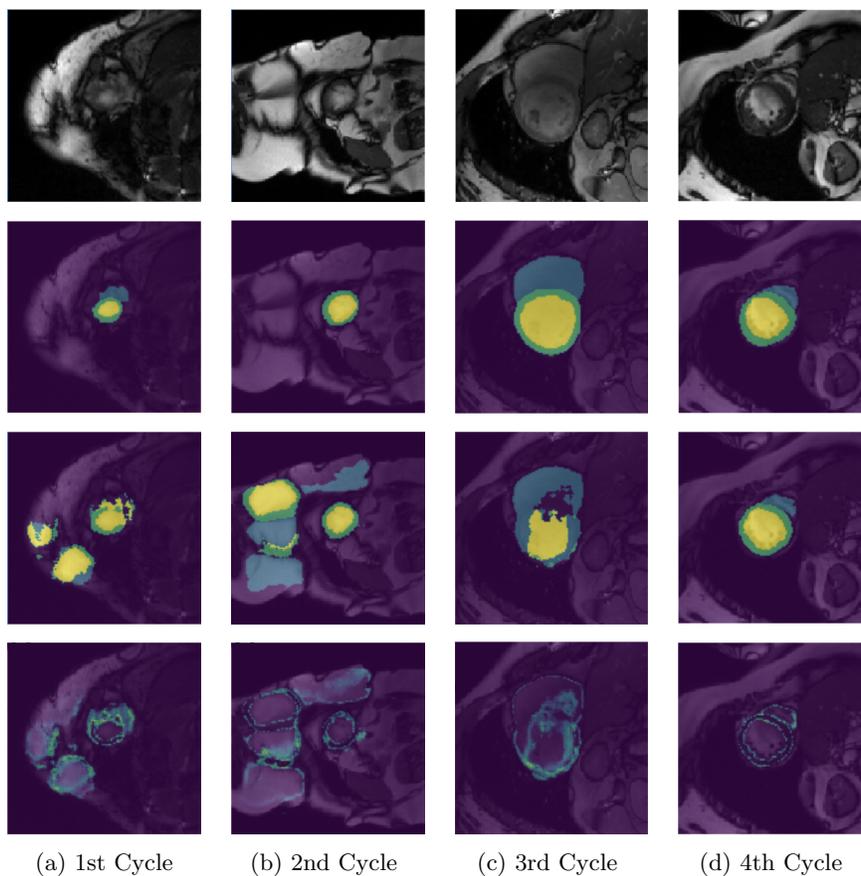


Fig. 2: Examples of images sampled by TAAL at different AL cycles. Are depicted the image sampled (row 1), the ground-truth segmentation (row 2), the segmentation prediction (row 3), and the JSD map given the different predictions from the augmented image (row 4). We observe that TAAL initially selected images with a large amount of hallucinated inaccurate predictions.

Figure 2 shows examples of images sampled by TAAL during the first 4 annotation cycles. TAAL initially selects image slices which show the apex of the heart. These samples are more difficult to learn in early stages since the areas to segment are much smaller than in the central slices of the heart and the image qualities are typically of lesser quality due to partial volume effects. Thus, we see that the choice of TAAL is first directed at samples yielding highly inaccurate predictions. The previous model has in fact even hallucinated multiple false segmentations for these samples as seen on the third row of subfigures 2a and 2b. In the next cycles, TAAL selects more central cardiac slices, which have improved predictions when compared to the ground-truth annotations. Hence, TAAL seems to first focus on correcting inaccurate predictions, before sharpening its predictions on a fine-grained level for slices with more prominent areas to segment.

Table 1: Active learning performances after doubling the number of initial labeled samples. We show the mean 2D and mean 3D Dice scores. ‘Fully’: Fully-supervised vanilla UNet. ‘Semi’: Proposed semi-supervised training with standard ( $\alpha = 0.5$ ) or weighted ( $\alpha = 0.75$ ) JSD. ‘RS’: Random sampling. ‘TTA’: Sampling with Test-time augmentation. ‘unweighted TAAL’: Our proposed method with standard JSD. ‘TAAL’: Our proposed method with weighted JSD, which finds the best candidate image to annotate.

Metric	Fully					Semi ( $\alpha = 0.5$ )		Semi ( $\alpha = 0.75$ )
	RS	Coreset	Entropy	Dropout	TTA	RS	unweighted TAAL	TAAL
2D Dice	80.69	79.95	80.99	81.32	81.67	81.51	81.90	<b>82.51</b>
3D Dice	87.40	86.65	88.07	88.24	88.48	88.48	88.50	<b>89.06</b>

Table 1 gathers the model’s segmentation performance after 10 cycles in terms of mean 2D Dice and mean 3D Dice scores over whole test volumes. In the fully-supervised setting, test-time augmentation-based sampling (TTA) outperforms random sampling, core-set selection, entropy sampling and sampling based on dropout. Similarly, unweighted TAAL and TAAL outperform random sampling in both semi-supervised and fully-supervised settings. After labeling 10 extra samples, the mean 3D Dice score attains 89.06% with TAAL while only reaching respectively 87.40% and 88.48% with random sampling in fully- and semi-supervised settings. Similar results were observed with 2D Dice on test images.

## 4 Conclusion

In this paper, we presented a simple, yet effective semi-supervised deep active learning approach for medical image segmentation. Our method, Test-time Augmentation for Active Learning (TAAL), employs a cross-augmentation consistency framework that produces both an improved training due to its unsupervised consistency loss, and a better sampling method through the uncertainty

measure it provides. TAAL also uses a modified JSD that significantly improves the model’s performance. Our results on the ACDC cardiac segmentation dataset show that, with TAAL, the trained model can reach up to 89.06% 3D Dice with 20 labeled samples when it only reaches 87.40% with random sampling. Because our approach exploits standard augmentation techniques already used in medical image segmentation tasks, TAAL emerges as a simple, yet efficient semi-supervised active learning strategy. While our method highly depends on the presence of disagreeing predictions for augmented inputs to identify the most informative samples, our observed improvements on a cardiac MRI dataset highlight promising avenues for future work, notably the investigation of more complex datasets and types of augmentations.

**Acknowledgments** – This work is supported by the Canada Research Chair on Shape Analysis in Medical Imaging, and the Research Council of Canada (NSERC). Computational resources were partially provided by Compute Canada. The authors also thank the ACDC Challenge organizers for providing the data.

## References

1. Settles, B.: Active Learning Literature Survey. Technical Report, University of Wisconsin-Madison Department of Computer Sciences (2009)
2. Budd, S., Robinson, E.C., Kainz, B.: A survey on active learning and human-in-the-loop deep learning for medical image analysis. *Medical Image Analysis* **71** (2021) 102062
3. Top, A., Hamarneh, G., Abugharbieh, R.: Active Learning for Interactive 3D Image Segmentation. In: *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*. Volume 6893. (2011) 603–610
4. Konyushkova, K., Sznitman, R., Fua, P.: Geometry in Active Learning for Binary and Multi-class Image Segmentation. *Computer Vision and Image Understanding (CVIU)* **182** (2019) 1–16
5. Sourati, J., Gholipour, A., Dy, J.G., Tomas-Fernandez, X., Kurugol, S., Warfield, S.K.: Intelligent Labeling Based on Fisher Information for Medical Image Segmentation Using Deep Learning. *IEEE Trans. Med. Imaging* **38**(11) (2019) 2642–2653
6. Nath, V., Yang, D., Landman, B.A., Xu, D., Roth, H.R.: Diminishing Uncertainty within the Training Pool: Active Learning for Medical Image Segmentation. *IEEE Trans. Med. Imaging* (2020)
7. Yang, L., Zhang, Y., Chen, J., Zhang, S., Chen, D.Z.: Suggestive Annotation: A Deep Active Learning Framework for Biomedical Image Segmentation. In: *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*. Volume 10435 of *Lecture Notes in Computer Science.*, Springer (2017) 399–407
8. Ash, J.T., Zhang, C., Krishnamurthy, A., Langford, J., Agarwal, A.: Deep Batch Active Learning by Diverse, Uncertain Gradient Lower Bounds. In: *Eighth International Conference on Learning Representations (ICLR)*. (2020)
9. Beluch, W.H., Genewein, T., Nurnberger, A., Kohler, J.M.: The Power of Ensembles for Active Learning in Image Classification. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. (2018) 9368–9377
10. Gal, Y., Islam, R., Ghahramani, Z.: Deep Bayesian Active Learning with Image Data. In: *Proceedings of the 34th International Conference on Machine Learning (ICML)*. (2017) 1183–1192

11. Sener, O., Savarese, S.: Active Learning for Convolutional Neural Networks: A Core-Set Approach. In: International Conference on Learning Representations (ICLR). (2018)
12. Wang, K., Zhang, D., Li, Y., Zhang, R., Lin, L.: Cost-Effective Active Learning for Deep Image Classification. *IEEE Trans. Circuits Syst. Video Technol.* **27**(12) (2017) 2591–2600
13. Yoo, D., Kweon, I.S.: Learning Loss for Active Learning. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). (2019) 93–102
14. Vezhnevets, A., Buhmann, J.M., Ferrari, V.: Active learning for semantic segmentation with expected change. In: 2012 IEEE conference on computer vision and pattern recognition, IEEE (2012) 3162–3169
15. Siddiqui, Y., Valentin, J., Nießner, M.: Viewal: Active learning with viewpoint entropy for semantic segmentation. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. (2020) 9433–9443
16. Casanova, A., Pinheiro, P.O., Rostamzadeh, N., Pal, C.J.: Reinforced active learning for image segmentation. In: International Conference on Learning Representations. (2019)
17. Sinha, S., Ebrahimi, S., Darrell, T.: Variational Adversarial Active Learning. In: IEEE/CVF International Conference on Computer Vision (ICCV). (2019) 5971–5980
18. Kim, K., Park, D., Kim, K.I., Chun, S.Y.: Task-Aware Variational Adversarial Active Learning. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). (2021) 8166–8175
19. Huang, S., Wang, T., Xiong, H., Huan, J., Dou, D.: Semi-Supervised Active Learning With Temporal Output Discrepancy. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). (2021) 3447–3456
20. Gao, M., Zhang, Z., Yu, G., Arik, S.O., Davis, L.S., Pfister, T.: Consistency-Based Semi-supervised Active Learning: Towards Minimizing Labeling Cost. In: European Conference on Computer Vision (ECCV). Volume 12355., Cham, Springer International Publishing (2020) 510–526 Series Title: Lecture Notes in Computer Science.
21. Berthelot, D., Carlini, N., Goodfellow, I., Papernot, N., Oliver, A., Raffel, C.A.: MixMatch: A Holistic Approach to Semi-Supervised Learning. In: Advances in Neural Information Processing Systems (NeurIPS). Volume 32., Curran Associates, Inc. (2019)
22. Ayhan, M.S., Berens, P.: Test-time Data Augmentation for Estimation of Heteroscedastic Aleatoric Uncertainty in Deep Neural Networks. (2018)
23. Wang, G., Li, W., Aertsen, M., Deprest, J., Ourselin, S., Vercauteren, T.: Aleatoric uncertainty estimation with test-time augmentation for medical image segmentation with convolutional neural networks. *Neurocomputing* **338** (2019) 34–45
24. Camarasa, R., Bos, D., Hendrikse, J., Nederkoorn, P., Kooi, E., Lugt, A.v.d., Bruijne, M.d.: Quantitative comparison of monte-carlo dropout uncertainty measures for multi-class segmentation. In: Uncertainty for Safe Utilization of Machine Learning in Medical Imaging, and Graphs in Biomedical Image Analysis. Springer (2020) 32–41
25. Shannon, C.E.: A Mathematical Theory of Communication. *The Bell System Technical Journal* **27**(3) (1948) 379–423
26. Bernard, O., Lalande, A., Zotti, C., Cervenansky, F., Yang, X., Heng, P.A., Cetin, I., Lekadir, K., Camara, O., Gonzalez Ballester, M.A., Sanroma, G., Napel, S.,

- Petersen, S., Tziritas, G., Grinias, E., Khened, M., Kollerathu, V.A., Krishnamurthi, G., Rohe, M.M., Pennec, X., Sermesant, M., Isensee, F., Jager, P., Maier-Hein, K.H., Full, P.M., Wolf, I., Engelhardt, S., Baumgartner, C.F., Koch, L.M., Wolterink, J.M., Isgum, I., Jang, Y., Hong, Y., Patravali, J., Jain, S., Humbert, O., Jodoin, P.M.: Deep Learning Techniques for Automatic MRI Cardiac Multi-Structures Segmentation and Diagnosis: Is the Problem Solved? *IEEE Transactions on Medical Imaging* **37**(11) (2018) 2514–2525
27. Ronneberger, O., Fischer, P., Brox, T.: U-Net: Convolutional Networks for Biomedical Image Segmentation. In Navab, N., Hornegger, J., Wells, W.M., Frangi, A.F., eds.: *Medical Image Computing and Computer-Assisted Intervention (MICCAI) 2015*. Lecture Notes in Computer Science, Cham, Springer International Publishing (2015) 234–241
  28. Kingma, D., Ba, J.: Adam: A method for stochastic optimization. In: *Int. Conf. on Learning Representations (ICLR)*. (2014)
  29. Loshchilov, I., Hutter, F.: SGDR: Stochastic Gradient Descent with Warm Restarts. In: *International Conference on Learning Representations (ICLR)*. (2017)
  30. Goyal, P., Dollár, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., He, K.: Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour. *arXiv:1706.02677 [cs]* (April 2018) *arXiv: 1706.02677*.
  31. Cui, W., Liu, Y., Li, Y., Guo, M., Li, Y., Li, X., Wang, T., Zeng, X., Ye, C.: Semi-supervised Brain Lesion Segmentation with an Adapted Mean Teacher Model. In Chung, A.C.S., Gee, J.C., Yushkevich, P.A., Bao, S., eds.: *Information Processing in Medical Imaging*. Lecture Notes in Computer Science, Cham, Springer International Publishing (2019) 554–565