

A comprehensive survey on deep active learning in medical image analysis

Haoran Wang^{a,b}, Qiuye Jin^c, Shiman Li^{a,b}, Siyu Liu^{a,b}, Manning Wang^{a,b,*}, Zhijian Song^{a,b,*}

^aDigital Medical Research Center, School of Basic Medical Sciences, Fudan University, Shanghai 200032, China

^bShanghai Key Laboratory of Medical Image Computing and Computer Assisted Intervention, Shanghai 200032, China

^cComputational Bioscience Research Center (CBRC), King Abdullah University of Science and Technology (KAUST), Thuwal 23955, Saudi Arabia

Abstract

Deep learning has achieved widespread success in medical image analysis, leading to an increasing demand for large-scale expert-annotated medical image datasets. Yet, the high cost of annotating medical images severely hampers the development of deep learning in this field. To reduce annotation costs, active learning aims to select the most informative samples for annotation and train high-performance models with as few labeled samples as possible. In this survey, we review the core methods of active learning, including the evaluation of informativeness and sampling strategy. For the first time, we provide a detailed summary of the integration of active learning with other label-efficient techniques, such as semi-supervised, self-supervised learning, and so on. We also summarize active learning works that are specifically tailored to medical image analysis. Additionally, we conduct a thorough comparative analysis of the performance of different AL methods in medical image analysis with experiments. In the end, we offer our perspectives on the future trends and challenges of active learning and its applications in medical image analysis.

Keywords: Active Learning, Medical Image Analysis, Survey, Deep Learning

1. Introduction

Medical imaging visualizes anatomical structures and pathological processes. It also offers crucial information in lesion detection, diagnosis, treatment planning, and surgical intervention. In recent years, the rise of artificial intelligence (AI) has led to significant success in medical image analysis. The AI-powered systems for medical image analysis have approached the performance of human experts in certain clinical tasks. Notable examples include skin cancer classification (Esteva et al., 2017), lung cancer screening with CT (Ardila et al., 2019), polyp detection during colonoscopy (Wang et al., 2018), and prostate cancer detection in whole-slide images (Tolkach et al., 2020). Therefore, these AI-powered systems can be integrated into existing clinical workflows, which helps to improve diagnostic accuracy for clinical experts (Sim et al., 2020) and support less-experienced clinicians (Tschandl et al., 2020).

Deep learning (DL) models serve as the core of these AI-powered systems for learning complex patterns from raw images and generalizing them to more unseen cases. Leveraging their robust feature extraction and generalization capabilities, DL models have also achieved remarkable success in the field of medical image analysis (Zhou et al., 2021a). The success of DL often relies on large-scale human-annotated datasets. For example, the ImageNet dataset (Deng et al., 2009) contains tens of millions of labeled images, and it's widely used in developing DL models for computer vision (CV). The size of medical image datasets keeps expanding, but

it is still relatively smaller than that of natural image datasets. For example, the brain tumor segmentation dataset BraTS consists of multi-sequence 3D MRI scans. The BraTS dataset expanded from 65 patients in 2013 (Menze et al., 2014) to over 1,200 in 2021 (Baid et al., 2021). The latter is equivalent to more than 700,000 annotated 2D images¹. However, the high annotation cost limits the construction of large-scale medical image datasets, mainly reflected in the following two aspects:

1. Fine-grained annotation of medical images is labor-intensive and time-consuming. In clinical practice, automatic segmentation helps clinicians outline different anatomical structures and lesions more accurately. However, training such a segmentation model requires pixel-wise annotation, which is extremely tedious (Rajpurkar et al., 2022). Another case is in digital pathology. Pathologists usually require detailed examinations and interpretations of pathological tissue slices under high-magnification microscopes. Due to the complex tissue structures, pathologists must continuously adjust the microscope's magnification. As a result, it usually takes 15 to 30 minutes to examine a single slide (Qu et al., 2022). Making accurate annotations is even more challenging for pathologists. In conclusion, the annotation process in medical image analysis demands a considerable investment of time and labor.

2. The high bar for medical image annotation leads to high costs. In CV, tasks like object detection and segmentation also require fine-grained annotations. However, the

¹BraTS 2021 dataset has scans with annotations available for 1,251 patients which all have the same spatial shape. Each case features four MRI sequences, and each sequence contains 155 2D axial slices. Thus, it contains a total of $1,251 \times 4 \times 155 = 775,620$ slices.

*Corresponding authors: Manning Wang and Zhijian Song (Emails: hrwang20@fudan.edu.cn, mnwang@fudan.edu.cn, zjsong@fudan.edu.cn)

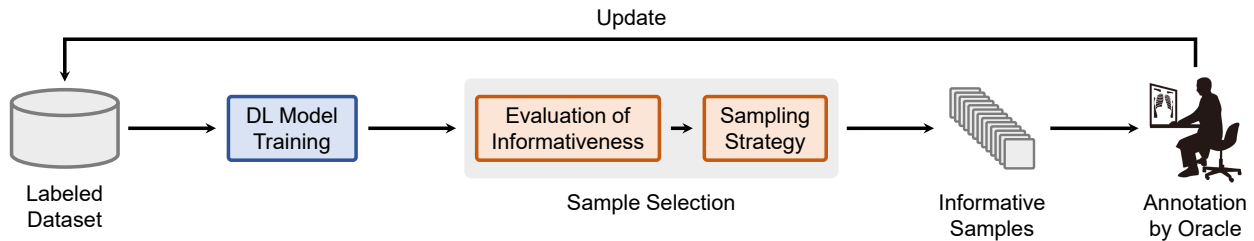


Figure 1: Illustration of the process of active learning.

widespread use of crowdsourcing platforms has significantly reduced the cost of obtaining high-quality annotations in these tasks (Kovashka et al., 2016). However, crowdsourcing platforms have certain limitations in medical image annotation. Firstly, annotating medical images demands both medical knowledge and clinical expertise. Complex cases even require discussions among multiple senior experts. Secondly, even in relatively simple tasks, crowdsourcing workers tend to provide annotations of poorer quality than professional annotators in medical image analysis. For example, results in Rädtsch et al. (2023) supported the conclusion above in annotating the segmentation mask of surgical instruments. Crowdsourcing platforms could also raise privacy concerns (Rajpurkar et al., 2022). Nevertheless, we will face new challenges when the annotators change from crowdsourcing workers to clinical experts. First of all, recruiting doctors for annotation is very expensive. For instance, a radiologist usually takes about 60 minutes to manually segment brain tumors per patient in its multi-sequence MRI volumes (Menze et al., 2014). And the median hourly rate of a radiologist is \$219 in the US². Besides, to minimize individual bias for certain scenarios, it is common to have a doctor annotate the same case multiple times or have multiple doctors annotate it. Multiple annotation rounds and annotators introduce intra- and inter-annotator variability and handling such variabilities leads to additional annotation costs (Karimi et al., 2020). In summary, high-quality annotations often require the involvement of experienced doctors, which inherently increases the annotation cost of medical images.

The high annotation cost is one of the major bottlenecks of DL in medical image analysis. Active learning (AL) is considered one of the most effective solutions for reducing annotation costs. The main idea of AL is to select the most informative samples for annotation and then train a model with these samples in a supervised way. In the general practice of AL, annotating a part of the dataset could reach comparable performance of annotating all samples. As a result, AL saves the annotation costs by querying as few informative samples for annotation as possible. The process of AL is illustrated in Fig.1, which we will detail in §2. Specifically, we refer to the AL works focusing on training a deep model as deep active learning.

Reviewing AL works in medical image analysis is essential for reducing annotation costs. Budd et al. (2021) investigated

the role of humans in developing and deploying DL in medical image analysis, where AL is considered a vital part of this process. In Tajbakhsh et al. (2020), AL was one of the solutions for training high-performance medical image segmentation models with imperfect annotation. As one of the methods in label-efficient deep learning for medical image analysis, Jin et al. (2023a) summarized AL methods from model and data uncertainty. There are also several surveys on AL in machine learning or CV. Settles (2009) provided a general introduction and comprehensive review of AL works in the machine learning era. After the advent of DL, Ren et al. (2021) reviewed the development of deep active learning and its applications in CV and natural language processing. Liu et al. (2022) summarized the model-driven and data-driven sample selectors in deep active learning. Zhan et al. (2022) reimplemented high-impact works in deep active learning with fair comparisons. Takezoe et al. (2023) reviewed recent developments of deep active learning in CV and its industrial applications.

However, the surveys mentioned above have certain limitations. Firstly, new ideas and methods are constantly emerging with the rapid development of deep active learning. Thus, a more comprehensive survey of AL is needed to cover the latest advancements. Secondly, a recent trend is combining AL with other label-efficient techniques, which is also highlighted as a future direction by related surveys (Takezoe et al., 2023; Budd et al., 2021). However, existing surveys still lack summaries and discussions on this topic. Thirdly, limited surveys have evaluated the performance of different AL methods on the medical imaging dataset, indicating a near absence of such efforts. Finally, the high annotation cost emphasizes the increased significance of AL in medical image analysis, yet related reviews still lack comprehensiveness in this regard.

This survey comprehensively reviews AL for medical image analysis, including core methods, integration with other label-efficient techniques, and AL works tailored to medical image analysis. We first searched relevant papers on Google Scholar and arXiv using the keyword “Active Learning” and expanded the search scope through citations. The included papers in this survey mainly belong to the field of medical image analysis. It should be noted that some important works of AL in the general CV field are also included since the development of AL in medical image analysis is influenced by the advance of AL in CV. Ignoring these works would flaw the logic and taxonomy of this survey. To balance the AL works of different fields, we first present the seminal works in each

²<https://www.salary.com/tools/salary-calculator/radiologist-hourly>

subsection, which may include works in the general CV field, and then provide a detailed review of the AL papers related to medical image analysis within this category. Additionally, most works in this survey are published in top-tier journals (e.g., TPAMI, TMI, MedIA, TBME, JBHI, etc.) and conferences (e.g., CVPR, ICCV, ECCV, ICML, ICLR, NeurIPS, MICCAI, ISBI, MIDL, etc.). As a result, this survey involves nearly 164 relevant AL works with 234 references. The contributions of this paper are summarized as follows:

- Through an exhaustive literature search, we provide a comprehensive survey and a novel taxonomy for AL works, especially those focusing on medical image analysis.
- While previous surveys mainly focus on evaluating informativeness, we further summarize different sampling strategies in deep active learning, such as diversity and class-balance strategies, aiming to provide references for future method improvement.
- In line with current trends, this survey is the first to provide a detailed review of the integration of AL with other label-efficient techniques, including semi-supervised learning, self-supervised learning, domain adaptation, region-based active learning, and generative models.
- For the sake of promoting research and contributing to the community, this survey evaluated the performance of several popular AL methods on multiple medical imaging datasets. The codes are also made public for better reproducibility.

The rest of this survey is organized as follows: §2 introduces problem settings and mathematical formulation of AL, §3 discusses the core methods of AL, including evaluation of informativeness (§3.1 & §3.2) and sampling strategies (§3.3), §4 reviews the integration of AL with other label-efficient techniques, §5 summarizes AL works tailored to medical image analysis. The experimental settings, results, and analysis are in §6. We discuss existing challenges and future directions of AL in §7 and conclude the whole paper in §8. The overall framework of this survey is shown in Fig. 1.

Due to the rapid development of AL, many related works are not covered in this survey. We refer readers to our constantly updated website³ for the latest progress of AL and its application in medical image analysis.

2. Problem Settings and Formulations of Active Learning

AL generally involves three problem settings: membership query synthesis, stream-based selective sampling, and pool-based active learning (Settles, 2009). In the case of membership query synthesis, we can continuously query any

samples in the input space for annotation, including synthetic samples produced by generative models (Angluin, 1988, 2004). We also refer to this setting as generative active learning in this survey. Membership query synthesis is typically suitable for low-dimensional input spaces. However, when expanded to high-dimensional spaces (e.g., images), the queried samples produced by generative models could be unidentifiable for human labelers. The recent advances of deep generative models have shown great promise in synthesizing realistic medical images, and we further discuss its combination with AL in §4.4. Stream-based selective sampling assumes that samples arrive one by one in a continuous stream, and we need to decide whether or not to request annotation for incoming samples (Cohn et al., 1994). This setting is suitable for scenarios with limited memory, such as edge computing, but it neglects sample correlations.

Most AL works follow pool-based active learning, which draw samples from a large pool of unlabeled data and requests oracle (e.g., doctors) for annotations. Moreover, if multiple samples are selected for labeling at once, we can further call this setting “batch-mode”. Deep active learning is in batch-mode by default since retraining the model every time a sample is labeled is impractical. Also, one labeled sample does not necessarily result in significant performance improvement. Therefore, unless otherwise specified, all works in this survey follow the setting of batch-mode pool-based active learning.

The flowchart of active learning is illustrated in Fig. 1. Assuming a total of T annotation rounds, active learning primarily consists of the following steps:

(1) Sample Selection: In the t -th round of annotation, $1 \leq t \leq T$, an informativeness function I is used to evaluate the informativeness of each sample in the unlabeled pool D_t^u . Then, a batch of samples is selected with a certain sampling strategy S . In medical image analysis, active learning selects a batch of images most of the time (i.e., image-wise selection). In this survey, the sampling unit of AL selection is an image (could be 2D or 3D) unless specifically stated. However, with the development of AL, region-wise (§4.3, §4.5.2) or slice-wise annotations (§5.2.1) are adopted in AL. Please refer to these sections for more details.

Specifically, the queried dataset of t -th round D_t^q is constructed as follow:

$$D_t^q = S_{D_t^q \subset D_t^u} \left(I_{x \in D_t^u}(x, f_{\theta_{t-1}}), b \right) \quad (1)$$

where x represents sample in the dataset, D_t^u and D_t^q are unlabeled and queried dataset in round t , respectively. $f_{\theta_{t-1}}$ and θ_{t-1} represent the deep model and its parameters from the previous round, respectively. The annotation budget b is the number of queried samples for each round, far less than the total count of unlabeled samples, i.e., $b = |D_t^q| \ll |D_t^u|$.

(2) Annotation by Oracle: After sample selection, the queried set D_t^q is sent to oracle (e.g., doctors) for annotation, and newly labeled samples are added into the labeled dataset D_t^l . The update of D_t^l is as follow:

³<https://github.com/LightersWang/Awesome-Active-Learning-for-Medical-Image-Analysis>

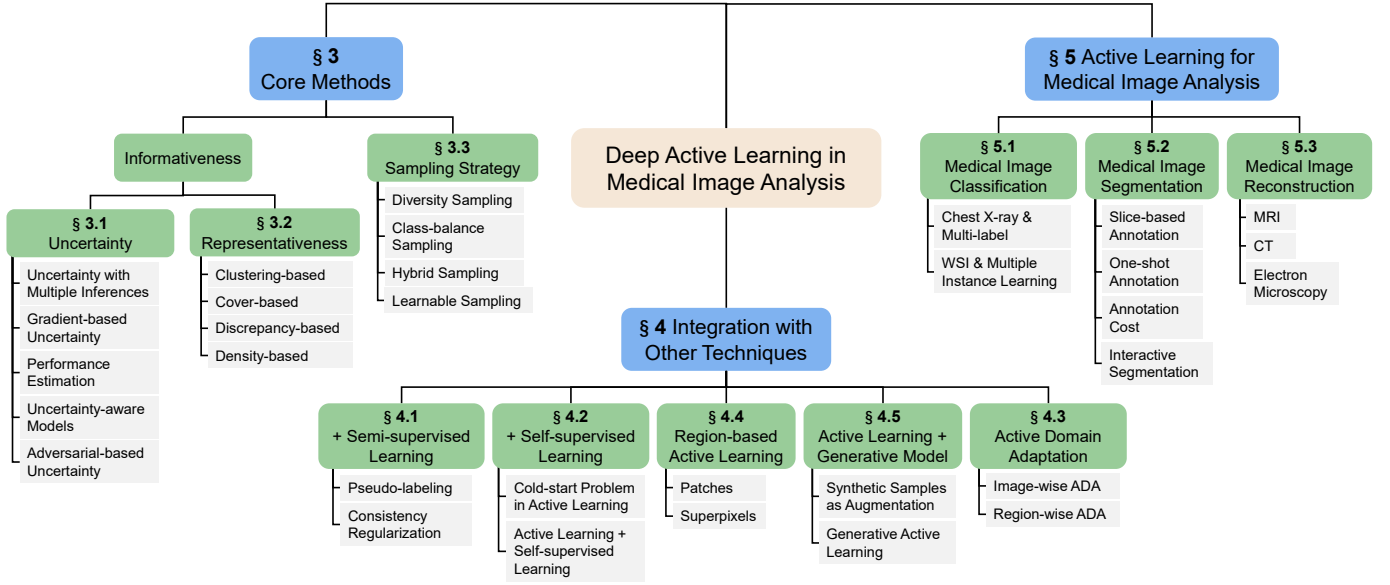


Figure 2: Overall framework of this survey.

$$D_t^l = D_{t-1}^l \cup \{(x, y) | x \in D_t^q\} \quad (2)$$

where y represents the label of x , and D_t^l and D_{t-1}^l denote the labeled sets for round t and the previous round, respectively. Besides, the queried samples should be removed from the unlabeled set D_t^u :

$$D_t^u = D_{t-1}^u \setminus \{x | x \in D_t^q\} \quad (3)$$

It is worth noting that some current works combine active learning with interactive segmentation. In interactive segmentation, the model assists experts in annotation, thereby reducing the difficulty of the annotation process. For more details, please refer to section §5.2.4.

(3) DL Model Training: After oracle annotation, we train the deep model using the labeled set of this round D_t^l in a fully supervised manner. The deep model f_{θ_t} is trained on D_t^l to obtain the optimal parameters θ_t for round t . The mathematical formulation is as follows:

$$\theta_t = \arg \min_{\theta} \mathbb{E}_{(x,y) \in D_t^l} [\mathcal{L}(f_{\theta}(x), y)] = \arg \min_{\theta} \mathbb{E}_{(x,y) \in D_t^l} [\mathcal{L}(x, y; \theta)] \quad (4)$$

where $\mathcal{L}(f_{\theta}(x), y)$ represents the loss function and it could be rewritten as $\mathcal{L}(x, y; \theta)$ for simplicity.

(4) Repeat steps 1 to 3 until the annotation budget limit or the expected performance is reached. Recently, some works adopted the one-shot fashion in active learning, which performed sample selection without multiple rounds. Please refer to §5.2.2.

It is worth noting that the model needs proper initialization to start the AL process. If the initial model f_{θ_0} is randomly initialized, it could only produce meaningless informativeness. To address this issue, most AL works randomly choose a set of

samples as initially labeled dataset D_0^l and train f_{θ_0} upon D_0^l . For more details on better initialization of AL using pre-trained models, please refer to §4.2.

3. Core Methods of Active Learning

In this survey, we consider the evaluation of informativeness and sampling strategy as the core methods of AL. Informativeness represents the value of annotating each sample. Higher informativeness often indicates a higher priority to request these samples for labeling. Typical metrics of informativeness include uncertainty and representativeness. Based on the informativeness scores, a certain sampling strategy is used to select a small number of unlabeled samples for annotation. Most AL works simply ranked these samples by their informativeness metrics and selected the highest ones according to the annotation budget (i.e., top-k selection). However, current informativeness scores are more or less flawed and they may cause issues like redundancy or class imbalance among the queried samples. Therefore, we need more advanced sampling strategies to mitigate these issues arising from imperfect informativeness metrics.

In this section, we reviewed two major informativeness metrics, including uncertainty (§3.1) and representativeness (§3.2), and sampling strategy (§3.3). As a unique contribution of this survey, we, for the first time, explicitly define sampling strategies as core methods of AL and review how to design a better sampling strategy in AL. Additionally, we provide a summarization of all the cited AL works in this survey. Methods and basic metrics of uncertainty or representativeness and sampling strategies are detailed in Table 2.

3.1. Evaluation of Informativeness: Uncertainty

Despite great progress has been made in medical image analysis, safety and interpretability are still unsolved problems

Table 1: Formulations of uncertainty metrics based on prediction probability in active learning. In the equations column, x stands for sample, f is the deep model, while C is the number of classes. In the direction column, \uparrow means higher values indicating higher uncertainty, while \downarrow means lower values indicating higher uncertainty.

Names	Equations	Direction
Prediction Probability	$\mathbf{p} = \text{Softmax}(f_\theta(x)) \in \mathbb{R}^C$, $\mathbf{p} = [p_1, p_2, \dots, p_C]$	-
Least Confidence (Lewis and Catlett, 1994)	$\max_i p_i$	\downarrow
Entropy (Joshi et al., 2009)	$-\sum_{i=1}^C p_i \log p_i$	\uparrow
Margin (Roth and Small, 2006)	$\max_i p_i - \max_{j, j \neq k} p_j, k = \arg \max_i p_i$	\downarrow
Mean Variance (Gal et al., 2017)	$-\frac{1}{C} \sum_{i=1}^C (p_i - \bar{p})^2, \bar{p} = -\frac{1}{C} \sum_{i=1}^C p_i$	\uparrow

for deploying DL models in real-world clinical practice. Due to the high variability of medical images and the limited training data, the predictions of DL models are not reliable and trusted. Correctly assessing and quantifying uncertainty in medical image analysis would allow models to alert ambiguities, artifacts, and unseen patterns in the data (Ghesu et al., 2021; Linmans et al., 2023). Such properties of uncertainty are helpful in AL since novel patterns in the unlabeled samples can be identified by the uncertainty. Therefore, uncertainty is frequently used in active learning as an informative metric. In the AL query, samples with higher uncertainty are considered hard and more likely to be misclassified by the current model. Annotating and training on these samples helps the model learn new patterns and improve performance.

Tracing back to the cause of uncertain predictions, the uncertainty can be mainly separated into two types: aleatoric uncertainty (AU) and epistemic uncertainty (EU) (Kendall and Gal, 2017). AU (i.e., data uncertainty) captures the noisy observations in data, such as motion artifacts of MRI or metal artifacts of CT in medical image analysis. AU cannot be reduced by acquiring more data. A high EU (i.e., model uncertainty) indicates that the samples contain knowledge that has not yet been mastered by the model. Therefore, the EU can be reduced by involving more data. However, most of the AL works did not consider the separation of AU and EU. So, the terminology of uncertainty in AL mainly refers to prediction uncertainty, which is the composition of AU and EU. This is because explicitly separating AU and EU is usually very difficult and may not bring much benefit in the practice of AL (Kahl et al., 2024). In this survey, unless explicitly stated otherwise, all uncertainties refer to predictive uncertainties, meaning that no separation between AU and EU has been performed.

The most straightforward uncertainty metrics in deep AL are based on prediction probabilities with a single forward pass. These metrics have been widely used in AL since the machine learning era, and their formulations are detailed in Table 1. However, directly transferring them to deep AL would be challenging due to the notorious issue of over-confidence in

deep neural networks (Mehrtash et al., 2020; Guo et al., 2017). Over-confidence refers to the model having excessively high confidence in its predictions, even though they might not be accurate. It could result in high confidence (e.g., 0.99) of the wrong class for misclassified samples. For uncertain samples, it leads to extreme confidence (e.g., 0.99 or 0.01) instead of normal one (e.g., 0.6 or 0.4) as it should. As a result, over-confidence distorts uncertainty estimation since it affects the predicted probabilities for all classes.

This section divides the uncertainty-based AL into multiple inference, gradient-based uncertainty, performance estimation, uncertainty-aware models, and adversarial-based uncertainty. The taxonomy of uncertainty-based AL is shown in Fig. 3.

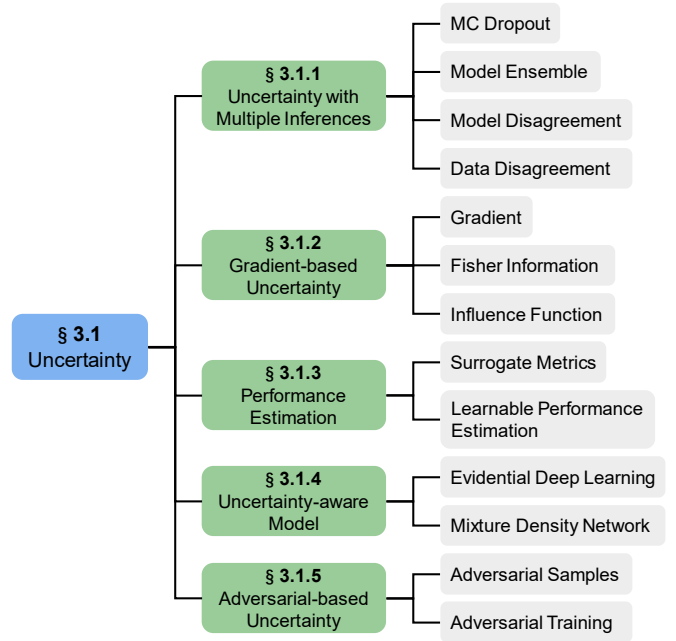


Figure 3: The taxonomy of uncertainty-based active learning.

3.1.1. Uncertainty with Multiple Inferences

To mitigate over-confidence, a common strategy for uncertainty-based AL is to run the model multiple times under perturbations. The main idea is to reduce the bias introduced by network architectures or training data. These biases often contribute to the over-confidence issue. Two approaches are often used to utilize multiple inference results for AL. The first is to calculate the classic uncertainty metrics with the average probability of multiple inferences. Averaging the prediction probabilities of multiple inferences helps to reduce individual bias that causes over-confidence. The other approach takes the disagreement between different prediction results as uncertainty quantification. Samples with higher disagreement indicate higher uncertainty and are suitable for annotation in AL.

In this section, we will introduce four types of methods for AL with multiple inferences: Monte Carlo dropout (MC dropout), model ensemble, model disagreement, and data disagreement. The first two used the average probability of results

of multiple inferences to calculate the uncertainty metrics, like entropy and margin. The last two are based on disagreement. For the source of perturbation, the first three perturb the model parameters while the last perturb the input data.

MC dropout randomly discards certain neurons in the deep model during each inference (Gal and Ghahramani, 2016). With MC dropout enabled, the model runs multiple times to get different predictions. Gal et al. (2017) was the pioneering work of deep AL. They were the first to use MC dropout in computing uncertainty metrics like entropy, standard deviation, and Bayesian active learning disagreement (BALD) (Houlsby et al., 2011). Results showed that MC Dropout could significantly improve the performance of uncertainty-based deep AL. Besides, they were also among the first to apply deep AL in medical image analysis. In the skin lesion analysis dataset ISIC 2016, they found that BALD consistently outperformed the random baseline. In brain cell type classification, Yuan et al. (2020b) calculated entropy using the average probability of multiple MC dropout runs. Gu et al. (2018) adopted the variance of multiple MC dropout runs as the uncertainty metric in the classification of confocal endomicroscopy and gastrointestinal endoscopy.

Model ensemble trains multiple models to get numerous predictions during inference. Beluch et al. (2018) conducted a detailed comparison of models ensemble and MC dropout in uncertainty-based AL. Results in the standard datasets demonstrated that the model ensemble performs better. For AL in diagnosing diabetic retinopathy, the proposed method achieved significant improvement compared to the random baseline. However, model ensemble requires significant training overhead in DL. To reduce the computational costs, snapshot ensemble (Huang et al., 2017) obtained multiple models in a single run with cyclic learning rate decay. An early attempt in Beluch et al. (2018) showed that snapshot ensemble leads to worse performance than model ensemble. Jung et al. (2023) improved the snapshot ensemble by maintaining the same optimization trajectory in different AL rounds, along with parameter regularization. Results showed that the improved snapshot ensemble outperforms the model ensemble. Nath et al. (2021) employed stein variational gradient descent to train an ensemble of models, aiming to ensure diversity. Their proposed method showed advantages to other competitors in segmenting the pancreas and tumor on CT and hippocampus on MRI.

Model disagreement: We can utilize the disagreement between the outputs of different models, which can be also referred to as Query-by-Committee (QBC) (Seung et al., 1992). This type of method was widely used in AL for medical image analysis. Suggestive annotation (SA) is the pioneering work of AL for medical image analysis (Yang et al., 2017). They trained multiple segmentation networks with bootstrapping. Variance among these models is used as the disagreement metric. SA demonstrated superior performance in segmenting glands on pathological images and lymph nodes on ultrasound images. In abdominal multi-organ segmentation, Qu et al. (2023a) trained three different segmentation models and adopted the variance between their

predictions. In carotid intima-media segmentation for ultrasound images, Tang et al. (2023) selected samples with the highest Kullback-Leibler (KL) divergence between the predictions of teacher and student models for annotation. In polyp segmentation of capsule colonoscopy, Bai et al. (2022) trained multiple decoders using class activation maps (CAMs) (Zhou et al., 2016) generated by a classification network. They further proposed model disagreement and CAM disagreement for sample selection. Model disagreement included entropy of prediction probabilities and Dice between outputs of different decoders, while CAM disagreement measured the Dice between CAMs and outputs of all decoders. This method selected samples with high model disagreement and CAM disagreement for annotation. However, samples with low model disagreement but high CAM disagreement were treated as pseudo-labels for semi-supervised training. In rib fracture detection, Huang et al. (2020) adopted Hausdorff distance to measure the disagreements between different CAMs. Besides, Mackowiak et al. (2018) adopted vote entropy between different MC dropout inferences as the disagreement metric.

Data disagreement: Since training multiple models can be computationally expensive, measuring the disagreements between different perturbations of input data is also helpful in AL. Kullback-Leibler (KL) divergence is a commonly used metric for quantifying disagreement. In COVID diagnosis, Wu et al. (2021) computed KL divergence between different versions of augmentations as the disagreement measure to select informative CT scans for annotation. Siddiqui et al. (2020) calculated the KL divergence between predictions of different viewpoints in 3D scenes to select informative regions for AL. Additionally, recent works have adopted alternative metrics to calculate disagreement. Lyu et al. (2023) proposed input-end committee, which randomly augmented the input data to get multiple predictions. They further measured the classification and localization disagreements between different predictions with cross-entropy and variance, respectively. Parvaneh et al. (2022) interpolated the unlabeled samples and labeled prototypes in the feature space. If the interpolated sample's prediction disagrees with the corresponding prototype's label, it indicates that the unlabeled samples introduce new features. Thus, these unlabeled samples should be sent for annotation. Results showed advancements across various datasets and settings.

3.1.2. Gradient-based Uncertainty

Gradient-based optimization is the cornerstone of DL-based medical image analysis. The gradient of each sample reflects its contribution to the change of model parameters. A larger gradient length indicates a tremendous change of parameters by the sample, thus implying high uncertainty. Furthermore, gradients are independent of predictive probabilities, which makes them less susceptible to over-confidence. Three metrics that are frequently used as gradient-based uncertainty: gradients, Fisher information, and influence functions. It should be noted that the gradient computation in this section did not use the ground truth labels which are unavailable for unlabeled samples. Instead, the

corresponding methods either used supervised loss with pseudo-labels (e.g. cross-entropy loss with pseudo-labels) or unsupervised loss (e.g. entropy loss), thereby making the gradient computation independent of the true labels.

Gradient: A larger gradient norm (i.e., gradient length) denotes a greater influence on model parameters, indicating higher uncertainty in AL. As an early attempt, Otálora et al. (2017) adopted the classic expected gradient length (Settles et al., 2007) to select valuable samples for annotation in exudate classification of eye fundus images. As a popular and pioneering work in the DL era, Ash et al. (2020) proposed batch active learning by diverse gradient embeddings (BADGE). They calculated the gradients only for the parameters of the network’s final layer, with the most confident classes as pseudo labels in gradient computation. Then, k-Means++ is performed on gradient embeddings for sample selection. Results showed competitive performances of BADGE across diverse datasets, network architectures, and hyperparameter settings. Gradient has been widely used in active learning of medical image analysis. Akilu and Yeung (2022) extended the BADGE framework into semantic segmentation of laparoscopic surgical images. Wang et al. (2022b) proved mathematically that a larger gradient norm corresponds to a lower upper bound of test loss. Thus, they employed expected empirical loss and entropy loss for gradient computation, which both obviate the necessity for labels. The former is the weighted sum of the losses of each class and the class probabilities, which are as follows:

$$\mathcal{L}_{exp}(x) = \sum_{i=1}^C [p_i \cdot \mathcal{L}(x, y_i; \theta)] \quad (5)$$

where y_i is the label of class i . The entropy loss is solely based on the probabilities of all classes, which are as follows:

$$\mathcal{L}_{ent}(x) = - \sum_{i=1}^C p_i \log p_i \quad (6)$$

The proposed method outperformed other comparative methods in cryo-electron tomography (cryo-ET) subtomogram classification. Besides, Dai et al. (2020) proposed a new gradients-based active learning method in MRI brain tumor segmentation. They first trained a variational autoencoder (VAE) (Kingma and Welling, 2013) to learn the data manifold. Then, they trained a segmentation model and calculated gradients of Dice loss using available labeled data. The sample selection was guided by the gradient projected onto the data manifold. Their extended work (Dai et al., 2022) further demonstrated superior performance in MRI whole brain segmentation.

Fisher information is effective in AL of machine learning models (Chaudhuri et al., 2015; Sourati et al., 2017). Fisher information (FI) reflects the overall uncertainty of model parameters according to data distribution. FI is defined as the expectation of the squared gradients with respect to the model parameter, the formulation is as follows:

$$\mathcal{I}_{Fisher}(x; \theta) = \mathbb{E}_y \left[\nabla_{\theta}^2 \mathcal{L}(x, y; \theta) \right] \quad (7)$$

where \mathcal{I} is the notation of Fisher information. The trace of the inverse of FI often serves as the objective for AL:

$$\arg \min_{D^q \subset D^u} \text{Tr} \left[\left(\sum_{x \in D^q} \mathcal{I}_{Fisher}(x; \theta) \right)^{-1} \left(\sum_{x \in D^u} \mathcal{I}_{Fisher}(x; \theta) \right) \right] \quad (8)$$

By solving Eq. 8, the selected samples could help the model converge faster toward optimal parameters. However, the computation cost of FI-based methods grows quadratically with the increase of model parameters, which is unacceptable for deep active learning. Sourati et al. (2018) and their extended work (Sourati et al., 2019) were the first to incorporate FI into deep active learning of medical image analysis. They used the average gradients of each layer to calculate the FI matrix, thus reducing the computation cost. Due to the absence of ground truth labels, they adopt the expected empirical loss (i.e., Eq. 5) for gradient computation. This method outperformed competitors in brain extraction across different age groups and pathological conditions. Additionally, Ash et al. (2021) only computed the FI matrix for the network’s last layer. The gradient computation of this work is the same as that of BADGE (Ash et al., 2020).

Influence function: Liu et al. (2021) employed influence function (Koh and Liang, 2017) to select samples that bring the most positive impact on model performance. The influence function of an unlabeled sample is defined as follows (Weisberg and Cook, 1982):

$$\mathcal{I}_{Influence}(x; D^l) = - \left(\sum_{(x,y) \in D^l} \nabla_{\theta} \mathcal{L}(x, y; \theta) \right) H_{\theta}^{-1} \nabla_{\theta} \mathcal{L}(x, y; \theta) \quad (9)$$

where H_{θ}^{-1} is the Hessian matrix of the labeled set, $H_{\theta} = \sum_{(x,y) \in D^l} \nabla_{\theta}^2 \mathcal{L}(x, y; \theta)$. In Eq. 9, the gradients of the first (i.e., the sum of gradients of labeled samples) and the second term (i.e., the Hessian matrix) can be derived with the ground truth label. For the third term, they replaced the true gradients with the gradients of the expected empirical loss due to the unavailability of the ground truth label.

3.1.3. Performance Estimation

In this section, uncertainty metrics are estimations of the current task’s performance. There are two types of such metrics: test loss or task-specific evaluation metrics. These metrics reflect the level of prediction error. For instance, a low Dice score in tumor segmentation for a patient suggests the model failed to produce accurate segmentation. Request annotations for these samples would be beneficial for improving the model’s performance. However, due to the unavailability of ground truth labels, we can only estimate these metrics instead of calculating them precisely. There are primarily two methods for estimating performance: surrogate metrics and learnable performance estimation.

Surrogate metrics are widely used in active learning of medical image analysis. For example, these metrics could be upper or lower bounds for loss or task-specific evaluation metrics. In breast cancer segmentation on immunohistochemistry images, Shen et al. (2020) calculated the intersection over union (IoU)

of all predictions by MC dropout. They found a strong linear correlation between this IoU and the real Dice coefficient. In skin lesion and X-ray hand bone segmentation, [Zhao et al. \(2021\)](#) calculated the average Dice coefficient between the predictions of the intermediate layers and final layer through deep supervision. They found a linear correlation between this average Dice and the real Dice coefficient. Besides, [Huang et al. \(2021\)](#) found that within limited training iterations, the loss of a sample is bounded by the norm of the difference between the initial and final network outputs. Inspired by this, they proposed cyclic output discrepancy (COD) as the difference in model output between two consecutive annotation rounds. Results indicated that a higher COD is associated with higher loss. Therefore, they opted for samples with high COD. They also demonstrate a linear correlation with the evaluation metrics with post-hoc validation.

Learnable performance estimation: We can train auxiliary neural network modules to predict the performance metrics. As one of the most representative works in this line of research, learning loss for active learning (LLAL) ([Yoo and Kweon, 2019](#)) trained an additional module to predict the loss value of a sample without its label. Since loss indicates the quality of network predictions, the predicted loss is a natural uncertainty metric for sample selection. Results showed that predicted and actual losses are strongly correlated. The proposed method also outperformed several AL baselines. In lung nodule detection with CT scans, [Liu et al. \(2020\)](#) built upon LLAL to predict the loss of each sample and bounding box. In COVID diagnosis, [Wu et al. \(2021\)](#) adopted both the predicted loss and the disagreements between different predictions for sample selection. [Wu et al. \(2022c\)](#) further combined the loss prediction and sample diversity in the federated active learning of COVID diagnosis and colonoscopy polyp analysis. Since AL focuses only on uncertainty ranking of the unlabeled samples, [Kim et al. \(2021\)](#) relaxed the loss regression to loss ranking prediction. Thus, they replaced the loss regressor in LLAL with the ranker in RankCGAN ([Saquil et al., 2018](#)). Results showed that loss ranking prediction outperforms the actual loss regression in LLAL. [Zhou et al. \(2021b\)](#) and their subsequent work ([Zhou et al., 2022](#)) introduced a quality assessment module to provide a predicted average IoU score for each slice. They interactively selected slices with the lowest scores in each volume for annotation.

3.1.4. Uncertainty-aware Model

In the above sections, uncertainty is derived based on the commonly used deterministic model in DL. However, some models can inherently capture uncertainty, such as VAE or probabilistic U-Net for medical image analysis ([Kohl et al., 2018](#)). In this way, they no longer output a point estimate but instead a distribution of possible predictions, thus mitigating over-confidence. We refer to them as uncertainty-aware models in this survey. They only require a single forward pass of the deep model, thus significantly reducing computational and time costs during inference. Evidential deep learning (EDL) and mixture density networks (MDN) are often used for uncertainty-aware models in AL.

Evidential deep learning replaces the Softmax distribution with a Dirichlet distribution ([Sensoy et al., 2018](#)). The network’s output is interpreted as the parameters of a Dirichlet distribution, so the predictions followed the Dirichlet distribution. The Dirichlet distribution will be sharp if the model is confident about the predictions. Otherwise, it will be flat. Another advantage brought by EDL is that AU and EU are easy to obtain with a Dirichlet distribution. In chest X-ray classification, [Balaram et al. \(2022\)](#) modified the EDL-based AL to accommodate the multi-label setting. Specifically, they transformed the Dirichlet distribution in EDL into multiple Beta distributions, each corresponding to one class label. They then calculated the entropy of the Beta distributions as the AU for annotation. [Chen et al. \(2023a\)](#) proposed a federated AL method based on EDL for medical image analysis. Following the federated AL setting, they kept a global model across all clients and local models for each client. AUs of the global and local models and the EUs of the global models were used for sample selection. [Park et al. \(2023\)](#) introduced a model evidence head to scale the parameters of the Dirichlet distribution adaptively in object detection, which enhanced training stability. They first calculated the EU for each detection box. Then, the sample-level uncertainty was obtained through hierarchical uncertainty aggregation. Besides, [Xie et al. \(2022c\)](#) introduced EDL into active domain adaptation. Samples with high distribution and data uncertainties are selected for annotation, which are both based on EDL.

Mixture density networks: [Choi et al. \(2021a\)](#) transformed the classification and localization heads in object detection networks to the architecture of MDN ([Bishop, 1994](#)). Besides the coordinates and class predictions of each bounding box, the MDN heads produced the variance of classification and localization. They used the variances as uncertainty metrics for sample selection. Results showed that this method is competitive with MC dropout and model ensemble while significantly reducing the inference time and model size.

3.1.5. Adversarial-based Uncertainty

Uncertainty in AL can also be estimated adversarially, including adversarial samples and adversarial training. **Adversarial samples** help measure the sample’s distance to the decision boundary implicitly, while a higher distance indicates higher uncertainty. By attacking the deep models, adding carefully designed perturbations to original samples results in adversarial samples ([Goodfellow et al., 2014b](#)). The differences between adversarial and original samples are nearly indiscernible to the human eye. However, deep models would produce extremely confident but wrong predictions for adversarial samples. The reason is that adversarial attacks push the original samples to the other side of the decision boundary with minimal cost, resulting in visually negligible changes but significantly different predictions. From this perspective, the strength of adversarial attacks reflects the sample’s distance to the decision boundary ([Heo et al., 2019](#)). A small perturbation indicates that the sample is closer to the decision boundary and, thus, is considered more uncertain. [Ducoffe and Precioso \(2018\)](#) adopted the DeepFool algorithm ([Moosavi-Dezfooli](#)

et al., 2016) for adversarial attacks. Samples with small adversarial perturbations are requested for labeling. Rangwani et al. (2021) attacked the deep model by maximizing the KL divergence between predictions of adversarial and original samples while the strength of perturbation is limited.

Adversarial training alternates between training feature extractors and classifiers with conflicting objectives, aiming to expose uncertain samples by increasing classifier disagreements. Yuan et al. (2021) and their extended work (Wan et al., 2023) implemented this with two classifiers on labeled and unlabeled datasets, first tuning classifiers while fixing the feature extractor to reveal more uncertain samples, then adjusting the feature extractor against fixed classifiers to minimize the gap between labeled and unlabeled samples. After several rounds, samples with the greatest disagreements are annotated.

3.2. Evaluation of Informativeness: Representativeness

While uncertainty-based methods play a crucial role in deep AL, they still face certain challenges: **1. Outlier selection:** The goal of using uncertainty in AL is to improve performance by querying hard samples of the current model. However, these methods could also select outliers that harm the model training (Karamcheti et al., 2021). This happens mainly because uncertainty-based methods often ignore the intrinsic characteristics of the sample itself. **2. Distribution misalignment:** In the feature space, uncertain samples are often located near the decision boundary (Settles, 2009). Therefore, the distribution of samples selected by uncertainty-based methods is usually different from the overall data distribution. This discrepancy introduces dataset bias and leads to a performance drop. This challenge could be alleviated if the relationship between different samples is carefully considered during the AL query. In summary, uncertainty-based AL lacks exploration of the visual information carried in each sample and the relationship between different samples. These challenges above call for a new informativeness metric in AL.

Representativeness is adopted in AL to overcome challenges brought by uncertainty. Representativeness-based AL aims to select a subset of samples that can represent the entire dataset. Specifically, representative samples should be visually distinctive in properties like imaging style or visual content. In medical image analysis, images are often high-dimensional and thus computation-intensive for the DL model. Besides, important information like lesions or tissues is not always directly visible or easily distinguishable. A good feature representation greatly reduces the image dimensionality and also extracts anatomical, histological, pathological, or even functional information in medical images. Therefore, the query process of representativeness-based AL is often conducted in the feature space. Besides, representative samples should also be widely distributed across the data distribution rather than concentrated in a specific region. In other words, these samples should be diverse. This is to minimize redundancy in the query result and try to keep the original data distribution as much as possible. Therefore, proper metrics of sample-wise or distribution-wise

distance and sample density are needed to assess the landscape of a dataset. Besides, the uniqueness of medical images may require different distance metrics than that of natural images. This section introduces four types of representativeness-based AL: clustering-based, cover-based, discrepancy-based, and density-based representativeness AL. The taxonomy of these methods is shown in Fig. 4.

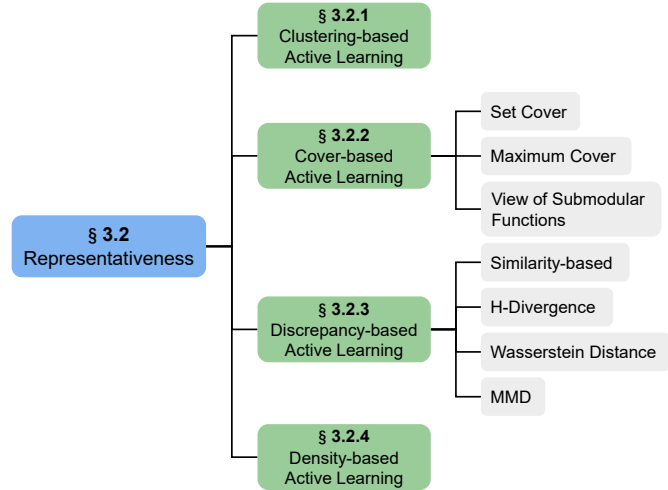


Figure 4: The taxonomy of representativeness-based active learning.

3.2.1. Clustering-based Active Learning

With the advancement of feature extraction in medical image analysis, images with similar appearances tend to group together in the feature space (Zheng et al., 2019). Therefore, a straightforward approach involves clustering on the data embeddings to select representative samples. This type of method grouped the data into several clusters and then selected the centroid samples of each cluster. It leveraged the inherent structure within the data for insightful groupings and was also very easy to implement. K-Means was the most popular choice in clustering-based active learning. Pourahmadi et al. (2021) performed k-Means on the off-the-shelf self-supervised features, then selected cluster centers for annotation. Based on the self-supervised feature, Jin et al. (2022b) adopted the k-Means++ for clustering and the silhouette coefficient to determine the optimal number of clusters. Their proposed method achieved commendable performance in lung segmentation of chest X-rays and lesion segmentation of dermoscopic images. In nuclei segmentation, Lou et al. (2023) performed coarse-level and fine-level clustering using K-Means, aiming to select informative patches from pathological images. In connectomics, Lin et al. (2020) proposed two-stream clustering for active selection. They first predicted the semantic mask for each unlabeled sample and simplified the AL task to judge the correctness of each predicted ROI. Besides, they trained two feature extractors of VAE with segmentation masks and unlabeled images, respectively. For two-stream clustering, they first applied mask-level clustering with mask features to group ROIs with similar appearances.

Within each mask cluster, image-level clustering is further performed. This method achieved excellent performance in synapse detection and mitochondria segmentation. Results also showed that two-stream clustering outperforms clustering with concatenated mask and image features by preventing the image feature from dominating the results.

3.2.2. Cover-based Active Learning

We can formulate representativeness-based AL as a covering problem. A classic example of the covering problem is facility location, such as covering all the city’s streets with billboards (Farahani and Hekmatfar, 2009). Likewise, cover-based AL uses a few samples to cover the entire dataset, which is analogous to using several spheres to cover all the samples in the feature space, with the center of each sphere being the selected sample. Ideally, these samples should be representative and contain information on other samples. These methods usually involve two settings: set cover and maximum cover. Both settings are NP-hard, meaning they cannot be optimally solved in polynomial time. However, near-optimal solutions could be achieved in linear time using greedy algorithms, which iteratively select samples that cover most of the other samples for annotation (Feige, 1998). These two variants are slightly different in the problem setting. The set cover is constrained by complete coverage, which means that it cannot omit any sample in the dataset. To achieve this goal, the radius of its covering spheres may be very large, and when there are very few samples selected, outliers might be chosen as the centers of the spheres (Yehuda et al., 2022). The goal of max cover is to cover as much of the entire dataset as possible, breaking the constraint of set cover, and thus avoiding the issue of outlier selection.

Set cover: Core-Set (Sener and Savarese, 2018) followed the setting of k-Center location (Hochbaum and Shmoys, 1985), which is also a variant of the set cover problem. They employed farthest-first traversal to solve the k-Center problem for selecting representative samples. The L2 distance of deep features is used to measure the similarity between different samples. Agarwal et al. (2020) introduced contextual diversity for AL, a metric that fused uncertainty and diversity of samples spatially and semantically. They replaced the L2 distance with contextual diversity and adopted the same farthest-first traversal for sample selection as Sener and Savarese (2018). Caramalau et al. (2021) adopted graph convolutional networks (GCN) to model the relationships between labeled and unlabeled samples. GCNs improved the feature representation of unlabeled samples with the labeled dataset. Enhanced feature representation was further used for Core-Set sampling.

Maximum cover: As a pioneering work of AL in medical image analysis, SA (Yang et al., 2017) stands out as one of the initial endeavors to introduce the concept of representativeness into AL. SA first selected highly uncertain samples and then chose representative samples for annotation. The formulation of the representativeness part in SA followed the setting of maximum cover. The representativeness metric was based on the cosine similarity of deep features. Specifically, sample x is

represented by the most similar sample from queried dataset D_t^q

$$r(D_t^q, x) = \max_{x' \in D_t^q} \text{sim}(x', x) \quad (10)$$

where r is the representativeness of sample x with respect to D_t^q and $\text{sim}(\cdot, \cdot)$ represents cosine similarity. Besides, representativeness R between D_t^q and the unlabeled set D_t^u is as follow:

$$R(D_t^q, D_t^u) = \sum_{x \in D_t^u} r(D_t^q, x) \quad (11)$$

where a larger $R(D_t^q, D_t^u)$ indicates that D_t^q better represents D_t^u . It should be noted that SA is a generalization of the maximum cover problem since the cosine similarity ranges from 0 to 1. But they still employed a greedy algorithm to find sample x that maximizing $R(D_t^q \cup x, D_t^u) - R(D_t^q, D_t^u)$. Many subsequent AL works built their framework of cover-based AL on SA, especially in the field of medical image analysis. Xu et al. (2018) quantized the segmentation networks in SA and found that it improved the accuracy of gland segmentation while significantly reducing memory usage. Zheng et al. (2019) proposed representative annotation (RA), which omits the uncertainty query in SA. RA trained a VAE for feature extraction and partitioned the feature space using hierarchical clustering. They selected representative samples in each cluster using a similar strategy to SA. RA achieved superior performance in gland segmentation on histological images, fungus segmentation on electron microscopy images, and whole heart segmentation on MRI. In breast cancer segmentation on immunohistochemistry images, Shen et al. (2020) changed the similarity measure in SA from $\text{sim}(\cdot, \cdot)$ to $1 - \text{sim}(\cdot, \cdot)$, which enhanced the diversity of the selected samples. Additionally, some works follow different formulations than those of SA in maximum cover. In keypoint detection of medical images, Quan et al. (2022) proposed a representative method to select template images for few-shot learning. First, they trained a feature extractor using self-supervised learning and applied the scale-invariant feature transform descriptor for initial keypoint detection. Next, they calculated the average cosine similarity between template images and the entire dataset. Finally, they picked the template combination with the highest similarity for annotation. Yehuda et al. (2022) found that Core-Set (Sener and Savarese, 2018), which followed the setting of the set cover, tends to select outliers, especially when the annotation budget is low. To address this issue, they proposed ProbCover, which changed the setting from set cover to maximum cover. With the help of self-supervised deep features and a graph-based greedy algorithm, ProbCover effectively avoided outlier selection in cover-based AL.

View of submodular functions: Both set cover and maximum cover can be formulated from the perspective of submodular set functions (Fujishige, 2005). These functions show diminishing returns. Specifically, given two sets A and B ,

$A \subset B$, for every element z that not in B , a submodular set function g has that $g(A \cup z) - g(A) \geq g(B \cup z) - g(B)$. This property makes submodular set functions suitable for AL. Suppose the informativeness function I is submodular. It means that each newly queried sample brings less informativeness gain than the previous one, which indicates that highly informative samples should be queried first. Besides, if we can formulate optimization problems in terms of monotonic and submodular functions, we can use a greedy algorithm to get near-optimal solutions in linear time. For AL, if I is submodular and monotonic, it means that we could greedily select the samples that maximize I . In cover-based AL, methods like SA and RA followed the setting of submodular functions, but the authors didn't present their methods from this perspective. Introducing submodular functions would extend the formulation of AL and ensure the selected samples are both representative and diverse. Typical steps for this type of method involve calculating sample similarities, constructing a submodular optimization problem, and solving it using a greedy algorithm (Wei et al., 2015). Kothawade et al. (2021) introduced an AL framework based on submodular information measures, effectively addressing issues such as scarcity of rare class, redundancy, and out-of-distribution data. In object detection, Kothawade et al. (2022a) focused on samples of minority classes. They first constructed a reference dataset containing samples of certain classes of interest. Then, unlabeled samples similar to the reference set for annotation through submodular mutual information (SMI). SMI is used to measure the similarity between two sets. Suppose two sets A , B and a submodular function g , the SMI is defined as $\mathcal{I}_{SMI} = g(A) + g(B) - g(A \cup B)$. Please refer to Kothawade et al. (2022b) for more detailed definitions of SMI.

3.2.3. Discrepancy-based Active Learning

In discrepancy-based AL, unlabeled samples farthest from the labeled set are considered the most representative. The main idea is that if we queried such samples for multiple rounds, the discrepancy between the distributions of labeled and unlabeled sets would be significantly reduced. Therefore, a small set of samples could well represent the entire dataset. The key to these methods is measuring the discrepancy (i.e., distance) between two high-dimensional distributions. In this section, we present four discrepancies between probability distributions: similarity-based discrepancy, H-divergence, Wasserstein distance, and maximum mean discrepancy (MMD).

Similarity-based discrepancy: As a practical and easy-to-implement metric, we can approximate the distance between distributions based on sample similarity. In gland and MRI infant brain segmentation, Li and Yin (2020) adopted the average cosine similarity as the distance between two datasets. They selected samples far from the labeled set and close to the unlabeled set. Caramalau et al. (2021) proposed UncertainGCN, which employed GCN to model the relationship between labeled and unlabeled samples. They selected the unlabeled samples with the lowest similarity to the labeled set. In object detection, Wu et al. (2022a) constructed prototypes

with sample features and prediction entropy. They selected unlabeled samples that were far from the labeled prototype.

H-divergence estimates the distance of distribution with the help of the discriminator from generative adversarial networks (GAN) (Goodfellow et al., 2014a). More specifically, the discriminator tries to distinguish between labeled and unlabeled samples, and there is a close relationship between H-divergence and the discriminator's output (Gissin and Shalev-Shwartz, 2019). Variational adversarial active learning (VAAL) (Sinha et al., 2019) combined VAE with a discriminator for discrepancy-based AL. In VAAL, the VAE mapped samples to a latent space while the discriminator distinguished whether samples were labeled. These two are mutually influenced by adversarial training. VAE tried to fool the discriminator into judging all samples as labeled while the discriminator attempted to differentiate between labeled and unlabeled samples correctly. After multiple rounds of adversarial training, VAAL selected samples that the discriminator deemed most likely to be unlabeled for annotation. VAAL inspired many subsequent works. Khanal et al. (2023) adopted multimodal information to improve VAAL. For multimodal medical images, they modified the VAE to reconstruct images of both modalities using the latent code of only one modality. The proposed method was evaluated on brain tumor segmentation, classification, and chest X-ray classification. Gissin and Shalev-Shwartz (2019) trained the discriminator without adversarial training. Zhang et al. (2020) replaced the discriminator's binary label with sample uncertainty. They also combined features of VAE with features from the supervised model. Wang et al. (2020c) adopted a neural network module for sample selection. To train such a module, they added another discriminator on top of VAAL, which aimed to differentiate between the real and VAE-reconstructed features for unlabeled samples. After adversarial training of both discriminators, the module selected uncertain and representative samples. Kim et al. (2021) combined learning loss for active learning with VAAL, feeding both loss ranking predictions and VAE features into the discriminator.

Wasserstein distance is widely used for computing distribution distances. Shui et al. (2020) indicated that H-divergence compromises the diversity of sample selection, while Wasserstein distance ensures the queried samples are representative and diverse. They further proposed Wasserstein adversarial active learning (WAAL), which built upon VAAL and adopted an additional module for sample selection. They trained this module by minimizing the Wasserstein distance between labeled and unlabeled sets. WAAL selected samples that are highly uncertain and most likely to be unlabeled for annotation. Mahmood et al. (2022) formulated AL as an optimal transport problem. They aimed at minimizing the Wasserstein distance between the labeled and unlabeled sets with self-supervised features. They further adopted mixed-integer programming that guarantees global convergence for diverse sample selection. Moreover, Xie et al. (2023b) considered the candidates as continuously optimizable variables based on self-supervised features. They randomly initialized the candidate samples at first. Then, they

maximized the similarity between candidates and their nearest neighbors while minimizing the similarity between candidates and labeled samples. Finally, they selected the nearest neighbors of the final candidates for annotation. They proved the objective is equivalent to minimizing the Wasserstein distance between the labeled and unlabeled samples.

Maximum mean discrepancy measures the distance of two distributions as the distances between their mean features with kernel trick (Gretton et al., 2012). In active domain adaptation (will be detailed in §4.5), Hwang et al. (2022) adopted MMD to measure the distance between the source and target domain. Then, MMD was used to select representative and diverse samples in the target domain. It should be noted that the Wasserstein distance belongs to the family of integral probability metrics (IPM), while MMD simultaneously falls into the range of IPM and previously mentioned H-divergence. Please refer to Zhao et al. (2022) for a more detailed taxonomy of the discrepancy between probability distributions.

3.2.4. Density-based Active Learning

Density-based active learning tends to select samples from the most densely populated area of the data distribution. It employs density estimation to characterize the data distribution in a high-dimensional feature space. The likelihood is the estimated density of the data distribution, and a more densely populated area indicates a higher likelihood. In this case, representative samples are samples with high likelihood. However, such methods can easily cause redundancy in sample selection. As a result, techniques like clustering are frequently used to improve diversity in sample selection. Density-based AL directly estimates the data distribution, which prevents the need to solve complex optimization problems. In shoulder MRI musculoskeletal segmentation, Ozdemir et al. (2021) adopted infoVAE (Zhao et al., 2017) to estimate the density of each sample in the labeled dataset and unlabeled pool. Specifically, MMD replaced the KL divergence as the regularization term in the training of infoVAE. The posterior probability by the encoder was used as the density metric. Samples with higher density regarding the unlabeled pool and lower density regarding the labeled dataset were selected for annotation. TypiClust (Hacohen et al., 2022) projected samples to a high-dimensional feature space via a self-supervised encoder. The density of a sample was defined as the reciprocal of the L2 distances to its k-nearest neighbors. Additionally, TypiClust performed clustering beforehand to ensure the diversity of selected samples. Wang et al. (2022c) proposed two variants of density-based AL. The first variant fixed the feature representation. The process was similar to TypiClust, but they maximized the distances between selected samples to ensure diversity. The other variant was in an end-to-end fashion. Feature representation and sample selection were trained simultaneously. This variant used a learnable k-Means clustering to jointly optimize cluster assignment and feature representation with a local smoothness constraint.

It is worth noting that cover-based and density-based AL differ in both concept and methodology. In concept, samples in cover-based AL tend to cover the entire dataset. However, they

do not have to lie in the densest area of the data distribution. For example, Yehuda et al. (2022) showed that Core-Set (Sener and Savarese, 2018), a popular cover-based method, tends to select outliers in the low-budget regime. In this case, cover-based AL is opposite to density-based AL, which also indicates that density-based AL may be a better choice in the low-budget regime. From the perspective of methodology, cover-based AL needs to solve an NP-hard problem in linear time with a greedy algorithm. Although this algorithm results in acceptable solutions, it's almost impossible to know how the AL performance would be if the optimal solutions could be achieved. For density-based AL, the NP-hard problem is replaced with density estimation, which is more computation-efficient.

3.3. Sampling Strategy

With a well-developed informativeness metric, most deep AL works simply adopted top-k to select samples with the highest informativeness for annotation. However, existing informativeness metrics face several issues, such as redundancy and class imbalance. These issues are exacerbated due to the unique characteristics of medical images. Despite the high variability, medical images of the same region-of-interest (ROI) could be classified into several groups, and images within each group share a high similarity (Zheng et al., 2019). Also, class imbalance is notorious in medical image analysis since the healthy objects often outnumber the diseased ones. Instead of proposing a better informativeness metric, we can improve the sampling strategy upon the top-k selection to effectively resolve these issues above. Besides, specific sampling strategies can also be used for combining multiple informativeness metrics. Furthermore, with the recent development of deep AL, more and more studies directly employ neural networks for sample selection. In this context, we no longer evaluate informativeness but directly choose informative samples with neural networks. Regrettably, despite the importance of sampling strategies in AL, prior works or surveys have seldom discussed their specific attributes. As one of the contributions of this survey, we systematically summarize different sampling strategies in AL, including diversity sampling, class-balanced sampling, hybrid sampling, and learnable sampling. The taxonomy of different sampling strategies in AL is shown in Fig. 5.

3.3.1. Diversity Sampling

Diversity strategies aim to reduce sampling redundancy in active learning, meaning certain selected samples are highly similar to each other. The lack of diversity leads to the waste of the annotation budget. Besides, redundancy in the training set causes the deep models to overfit to limited training samples, thus leading to a performance drop. Therefore, many AL methods employ diversity sampling to mitigate the redundancy in selected samples. In this section, we discuss four strategies of diversity sampling, including clustering, farthest-first traversal, determinantal point process (DPP), and specific strategies tailored to certain informativeness metrics.

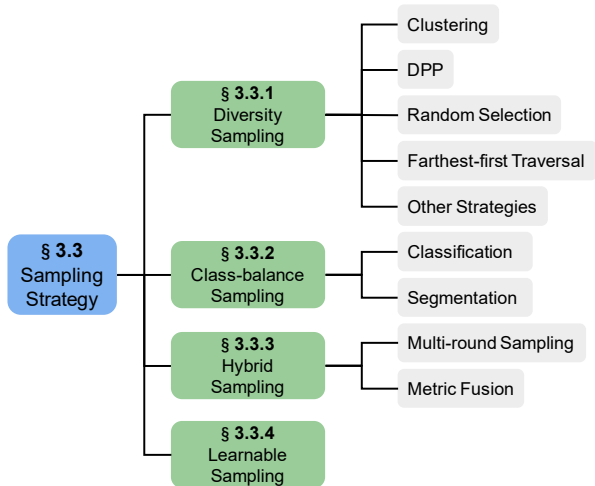


Figure 5: The taxonomy of different sampling strategies in active learning.

Clustering is one of the most commonly used strategies of diversity sampling. This strategy improves the coverage of the entire feature space, thereby easily boosting diversity. Ash et al. (2020) employed k-Means++ clustering on gradient embeddings to select diverse uncertain samples. Besides, Citovsky et al. (2021) boosted margin-based uncertainty sampling with hierarchical clustering. They selected samples with the smallest margins within each cluster. When the number of queries exceeded the number of clusters, samples from smaller clusters were prioritized. This method can scale to a huge annotation budget (e.g., one million). Zheng et al. (2019) incorporated hierarchical clustering with cover-based AL. In their experiments, clustering showed consistent performance improvement in multiple medical imaging datasets, which demonstrated that clustering does improve the sampling diversity.

It is important to highlight that clustering in this section is different from that of in §3.2.1. To ensure that the selected samples are sufficiently representative, clustering-based AL generally chooses samples that are closest to the cluster centers. However, when clustering is used to enhance diversity, we can not only select samples closest to the cluster centers, but also select samples with the highest uncertainty, or even randomly select within each cluster. Therefore, clustering can serve as a plug-and-play technique to conveniently enhance the sampling diversity in AL.

Determinantal point process is a stochastic probability model for selecting subsets from a larger set. DPP reduces the probability of sampling similar elements to ensure diversity in the results. Bıyık et al. (2019) employed two DPPs for sample selection: Uncertainty DPP is based on uncertainty scores, while Exploration DPP aims to find samples near decision boundaries. Then, sampling results from both DPPs were sent for expert annotation. However, DPP is more computationally intensive compared to clustering. Ash et al. (2020) compared the performance and time cost of using k-Means++ and k-DPP. Results showed that their performance is similar, but the time

cost for k-Means++ is significantly lower than that for k-DPP. Besides, Mi et al. (2020) adopted DPP in AL for medical image reconstruction, please refer to §5.3 for details.

Random Selection could also be used for better diversity. In prostate segmentation of MRI, Gaillochet et al. (2023a) randomly partitioned the entire dataset into different batches, which were referred to as ‘stochastic batches’. Batches with the highest uncertainty scores were selected for annotation. Experimental results showed that the stochastic batches consistently improved the performance of various uncertainty-based AL methods under an extremely low budget. Their extended works (Gaillochet et al., 2023b) further illustrated the effectiveness of stochastic batches on the anterior and posterior hippocampus segmentation.

Farthest-first Traversal is also a widely used strategy for diverse queries, which was first adopted by Sener and Savarese (2018). This strategy requires the distance between sampling points to be as large as possible, which leads to a more uniform distribution of selected samples in the feature space. Li et al. (2023) adopted a farthest-first traversal strategy with cosine distance for a diverse initial labeled dataset. Experiments on breast ultrasound, liver CT, and chest X-ray segmentation showed significant effectiveness of the farthest-first traversal. Besides, Agarwal et al. (2020) and Caramalau et al. (2021) improved the diversity with farthest-first traversal with their proposed contextual diversity and GNN-augmented features, respectively.

Other strategies: In uncertainty-based AL, BatchBALD (Kirsch et al., 2019) extended BALD-based uncertainty AL to batch mode. Results showed that BatchBALD improved the sampling diversity compared to (Gal et al., 2017). FI-based methods formulated AL as a semi-definite programming (SDP) problem to improve sampling diversity and various methods were employed for solving SDP. Sourati et al. (2019) used a commercial solver to solve SDP, while Ash et al. (2021) proposed a greedy algorithm to adapt to high-dimensional feature space. In skin lesion analysis, Shi et al. (2019) introduced image hashing for diversity sampling. In their proposed method, the first principal component of each image was used for feature representation. Then they mapped similar images into the same buckets using local sensitivity hashing. Samples were uniformly selected from each bucket for human annotation.

3.3.2. Class-balance Sampling

Class imbalance is a common issue for DL in medical image analysis, where a small set of classes have many samples while the others only contain a few samples (Zhang et al., 2023b). For example, long tail distribution of classes existed in almost all tasks of medical image classification, such as skin lesion classification and whole-slide image classification. Training on imbalanced datasets can lead to the overfitting of the majority classes and underfitting of the minority classes. Apart from dealing with class imbalance during training, AL mitigates class imbalance by avoiding over-annotation of the majority classes and enhancing the annotation of the minority classes during dataset construction.

Classification: In a class-imbalanced COVID-19 dataset, [Chong et al. \(2021\)](#) evaluated multiple informativeness scores and sampling strategies. Results showed that diversity sampling is more favored for class-imbalance. [Jin et al. \(2022c\)](#) assumed that samples closer to the tail of the distribution are more likely to belong to the minority classes. Thus, the tail probability is equivalent to the likelihood of minority classes. Specifically, they trained a VAE for feature extraction and adopted copula to estimate the tail probabilities upon VAE features. Finally, informative samples were selected with clustering and unequal probability sampling. The proposed method was validated on the ISIC 2020 dataset, which has a long-tailed distribution. [Kothawade et al. \(2022c\)](#) used submodular mutual information to focus more on samples of minority classes. They achieved excellent results on medical classification datasets in five different modalities, including X-rays, pathology, and dermoscopy. In blood cell detection under microscopy, [Sadafi et al. \(2019\)](#) requested expert annotation of a sample whenever its classification probability of the minority class exceeded 0.2. Besides, [Choi et al. \(2021b\)](#) directly estimated the probability of a classifier making a mistake for a given sample and decomposed it into three terms using Bayesian rules. First, they trained a VAE to estimate the likelihood of the data given a predicted class. Then, an additional classifier was trained upon VAE features to estimate class prior probabilities and the probability of mislabeling a specific class. By considering all three probabilities, they successfully mitigated class imbalance in AL. The proposed method achieved good performance on stepwise class-imbalanced CIFAR-10 and CIFAR-100 datasets. For uncertainty-based methods, [Bengar et al. \(2022\)](#) introduced an optimization framework to maintain class balance. They compensated the query of minority classes with the most confident samples of that class, leading to a more balanced class distribution in the queried dataset.

Segmentation: Due to certain AL methods selecting regions instead of the entire image for annotation, there is a need to ensure that the selected regions contain rare or small objects (e.g., optic chiasma or optic nerve in head and neck multi-organ segmentation). [Cai et al. \(2021\)](#) and [Wu et al. \(2022b\)](#) both proposed class-balanced sampling strategies for such scenarios, as detailed in §4.3.

3.3.3. Hybrid Sampling

In AL, more and more works use multiple informativeness metrics simultaneously. However, how to effectively integrate multiple metrics remains a critical issue. This issue is addressed by the hybrid sampling discussed in this section. Two approaches to hybrid sampling are often used, including multi-round sampling and metric fusion.

Multi-round sampling first selects a subset of samples based on one particular informativeness metric and continues sample selection within this subset based on another informativeness metric. Multi-round sampling is widely used in AL for medical image analysis for its convenience ([Shen et al., 2020](#); [Li and Yin, 2020](#); [Wang and Yin, 2021](#)). For example, SA ([Yang et al., 2017](#)) performed representativeness sampling

based on uncertainty to reduce redundancy in the sampled set. Besides, [Wu et al. \(2022b\)](#) employed an adaptive strategy that sets dynamic weights to adjust the budget of representativeness and uncertainty sampling. The weight of representativeness sampling is larger initially, while the situation is reversed in the latter phase. This is because representativeness methods can quickly spot typical data, while uncertainty methods continuously improve the model by querying samples with erroneous predictions.

Metrics fusion is another widely used approach of hybrid sampling. It directly combines different informativeness metrics. For example, one could directly sum up all metrics and select the samples with the highest values for annotation. Metrics fusion is also widely used in AL of the medical domain ([Li et al., 2024](#); [Zhou et al., 2021c](#); [Wu et al., 2021](#)). Besides, ranked batch-mode ([Cardoso et al., 2017](#)) can adaptively fuse multiple metrics in AL.

3.3.4. Learnable Sampling

Previously mentioned AL methods typically follow a “two-step” paradigm, which first involves the evaluation of informativeness and then selects samples based on specific heuristics (i.e., sampling strategy). However, learnable sampling skips the informativeness evaluation and directly uses neural networks for sample selection. In this context, the neural network is known as a “neural selector”.

One of the most common methods of learnable sampling is to formulate sample selection as a reinforcement learning (RL) problem, where the learner and the dataset are considered the environment, and the neural selector serves as the agent. The agent interacts with the environment by selecting a limited number of samples for annotation, and the environment returns a reward to train the neural selector. In medical image classification, [Wang et al. \(2020b\)](#) employed an actor-critic framework where the critic network is used to evaluate the quality of the samples selected by the neural selector. This method has performed excellently in lung CT disease classification and diabetic retinopathy classification of fundus images. Besides, [Haußmann et al. \(2019\)](#) adopted a probabilistic policy network as the neural selector. The rewards returned by the environment encouraged the neural selector to choose diverse and representative samples. The neural selector is trained using the REINFORCE algorithm ([Williams, 1992](#)). [Agarwal et al. \(2020\)](#) utilized contextual diversity as RL rewards and trains a bidirectional long short-term memory network as the neural selector.

For more works on learnable sampling in AL, such as formulating AL as few-shot learning or training neural selectors by meta-learning, please refer to the survey of [Liu et al. \(2022\)](#).

4. Integration of Active Learning and Other Label-Efficient Techniques

As discussed in §1, the high annotation cost has severely dragged down the development of DL in medical image

analysis. Despite the wide use of AL in medical image analysis, various methods have been proposed to reduce the large amount of labeled data required for training deep models, such as semi-supervised and self-supervised learning, etc. These methods, including active learning, are collectively called label-efficient deep learning (Jin et al., 2023a). Label-efficient learning is a broad concept that includes all related technologies designed to improve annotation efficiency. In the real-world practice of AL in medical image analysis, there is still room for higher label efficiency by integrating AL with other label-efficient techniques. For the example of AL in medical image segmentation, since many samples were left unlabeled in the cycle of AL, we could further include them to achieve better performance by integrating AL with semi-supervised learning. The rapid development of self-supervised learning in medical image analysis introduced many powerful pre-trained models (Taleb et al., 2020). These models are also valuable in AL of medical image analysis for their superior ability of feature extraction. For another circumstance, since the ROI in medical imaging is usually small, we could select and annotate the informative regions that contain the ROI in AL instead of annotating the whole image. As a result, integrating active learning with other label-efficient techniques holds significant potential to increase annotation efficiency. However, existing surveys have not yet systematically organized and categorized this line of research. Hence, as one of the main contributions of this survey, we comprehensively reviewed the integration of AL with other label-efficient techniques, including semi-supervised learning, self-supervised learning, domain adaptation, region-based annotation, and generative models. Additionally, how each surveyed work integrated with other label-efficient techniques is summarized in Table 2.

4.1. Semi-supervised Learning: Utilizing Unlabeled Data

Semi-supervised learning (Chen et al., 2022a; Han et al., 2024) aims to boost performance by utilizing unlabeled data upon supervised training. The need for tedious human annotation can be further reduced by integrating AL and semi-supervised learning in medical image analysis. The reason is that AL and semi-supervised learning complement each other. Specifically, a large pool of unlabeled images should be collected from the hospital information systems to train a DL model for some clinical applications. With the help of AL, the DL model is trained on an optimal labeled dataset constructed with a certain AL method, which reduces the annotation workload for doctors. However, massive unlabeled samples sit idle during the model training in the AL cycle. By combining the AL with semi-supervised learning, the model can be trained on both labeled and unlabeled samples. (Jiménez et al., 2023) This section will introduce the integration of AL and semi-supervised learning from the perspectives of pseudo-labeling and consistency regularization.

4.1.1. Pseudo-Labeling

Pseudo-labeling (Lee et al., 2013) is one of the most straightforward methods in semi-supervised learning. It uses

the model’s predictions of unlabeled data as pseudo-labels and combines them with labeled data for supervised training. Although it’s possible to assign pseudo-labels to all unlabeled samples for training, it could introduce noise. To mitigate this, Wang et al. (2017) proposed cost-effective active learning (CEAL), integrating pseudo-labeling with uncertainty-based AL. Specifically, CEAL sent the most uncertain samples for expert annotation and assigned pseudo-labels to the most confident samples. Many subsequent works have built upon the ideas of CEAL. Gorriz et al. (2017) adopted the CEAL framework in the melanoma segmentation and used the MC dropout for uncertainty estimation. In medical image segmentation, Zhao et al. (2021) refined the pseudo-labels with dense conditional random fields. Additionally, Li et al. (2022) proposed a new approach for selecting samples for oracle annotation and pseudo-labeling in the Gleason grading of prostate cancer with histopathology images. They employed curriculum learning to categorize all samples into hard and easy. Hard samples were all sent for oracle annotation. For the easy samples, they evaluated the presence of label noise based on the training loss. Easy samples with low training loss were used for pseudo-labels to assist training, whereas easy samples with high loss were considered noisy and excluded from training.

4.1.2. Consistency Regularization

Consistency regularization aims to enforce similar outputs under perturbations of input data or model parameters. Maximizing consistency serves as an unsupervised loss for unlabeled samples, which helps improve the robustness, reduces overfitting, and improves model performance. Many works integrated existing consistency-based semi-supervised methods into the training process of AL. In chest X-ray classification, Balarum et al. (2022) incorporated several semi-supervised methods with AL to further reduce annotation costs, including MeanTeacher (Tarvainen and Valpola, 2017), VAT (Miyato et al., 2018) and NoTeacher (Unnikrishnan et al., 2021). Huang et al. (2021) combined their proposed COD with MeanTeacher (Tarvainen and Valpola, 2017), demonstrating superior performance. Wang et al. (2022c) combined density-based AL with different existing semi-supervised methods. Results showed that the proposed method outperforms other active learning methods and excels in semi-supervised learning.

Consistency could be also used for sample selection. Gao et al. (2020) introduced a semi-supervised active learning framework. Consistency here was used for both semi-supervised training and evaluating informativeness. In this framework, samples are fed into the model multiple times with random augmentations. The consistency loss of unlabeled samples was implemented by minimizing the variance between multiple outputs. They further selected less consistent samples for annotation. Results showed that combining AL with semi-supervised learning significantly improves performance.

Besides, Zhang et al. (2022a) combined AL with both pseudo-labeling and consistency regularization. The unlabelled images first underwent both strong and weak data augmentations. When the confidence level of the weakly augmented images exceeded a certain threshold, they used

Table 2: Methodology summarization of surveyed active learning works.

Year	Vvenues	Uncertainty		Representativeness		Sampling Strategy	SemiSL	SelfSL	ADA	Region	Generative
		Method	Basic Metrics	Method	Basic Metrics						
Zhu and Bento (2017)	2017	arXiv	Single Model	Distance to Decision Boundary	-	Top-k					✓
Zhou et al. (2017)	2017	CVPR	Single Model Multiple Inferences - Data Disagreement	Entropy KL Divergence	-	Hybrid - Fusion					
Gal et al. (2017)	2017	ICML	Multiple Inferences - MC Dropout	Entropy, BALD, Least Confidence, Variance	-	Top-k					
Yang et al. (2017)	2017	MICCAI	Multiple Inferences - Model Disagreement	Variance	Cover-based	Hybrid - Multi-round					
Wang et al. (2017)	2017	TCSVT	Single Model	Least Confidence, Margin, Entropy	-	Top-k				Pseudo-label	
Ducicco and Precioso (2018)	2018	arXiv	Adversarial Samples	Distance to Decision Boundary	-	Top-k					
Mackowiak et al. (2018)	2018	BMVC	Multiple Inferences - Model Disagreement	Vote Entropy	-	Top-k					Patch
Xu et al. (2018)	2018	CVPR	Multiple Inferences - Model Ensemble	Variance	Cover-based	Hybrid - Multi-round					
Beluch et al. (2018)	2018	CVPR	Multiple Inferences - Model Ensemble	Entropy, BALD, Least Confidence, Variance	-	Top-k					
Sourati et al. (2018)	2018	DLMIA	Gradient-based Uncertainty	Fisher Information	-	Diversity - Solve Programming Problem					✓
Sener and Savaresi (2018)	2018	ICLR	-	-	Cover-based	L2 Distance				Diversity - Farthest-first Traversal	
Kuo et al. (2018)	2018	MICCAI	Multiple Inferences - Model Disagreement	JS Divergence	-	Diversity - Solve Programming Problem					
Mahapatra et al. (2018)	2018	MICCAI	Multiple Inferences - MC Dropout	Variance	-	Top-k					✓
Hauffmann et al. (2019)	2019	ICAI	-	-	-	Learnable - Reinforcement Learning					
Zheng et al. (2019)	2019	AAAI	-	-	Cover-based	Diversity - Clustering					
Gissin and Shalev-Shwartz (2019)	2019	arXiv	-	-	Discrepancy-based	H-Divergence					
Yoo and Kweon (2019)	2019	CVPR	Performance Estimation - Learnable	Loss	-	Top-k					
Sinha et al. (2019)	2019	ICCV	-	-	Discrepancy-based	H-Divergence					
Tran et al. (2019)	2019	ICML	Multiple Inferences - MC Dropout	BALD	-	Top-k				Pseudo-label	✓
Qi et al. (2019)	2019	JBHI	Single Model	Entropy	-	Top-k					
Sadafi et al. (2019)	2019	MICCAI	Multiple Inferences - MC Dropout	Average IoU, Class Frequency	-	Class-balance Hybrid - Fusion					
Kirsch et al. (2019)	2019	NeurIPS	Multiple Inferences - MC Dropout	BALD	-	Top-k					
Sourati et al. (2019)	2019	TMI	Gradient-based Uncertainty	Fisher Information	-	Diversity - Solve Programming Problem					
Kasaria et al. (2019)	2019	WACV	Single Model	Entropy	-	Top-k					
Zheng et al. (2020)	2020	AAAI	-	-	Cover-based	Cosine Similarity				Pseudo-label	
Shui et al. (2020)	2020	AISTATS	Single Model	Entropy, Least Confidence	Discrepancy-based	Wasserstein Distance					✓
Siddiqui et al. (2020)	2020	CVPR	Multiple Inferences - MC Dropout Multiple Inferences - Data Disagreement	Entropy KL Divergence	-	Hybrid - Fusion					Superpixel
Zhang et al. (2020)	2020	CVPR	Single Model	Variance	Discrepancy-based	H-Divergence				Diversity - Farthest-first Traversal	
Gao et al. (2020)	2020	ECCV	Multiple Inferences - Data Disagreement	Variance	-	Top-k				Consistency	
Wang et al. (2020c)	2020	ECCV	-	-	Discrepancy-based	H-Divergence				Learnable	
Agarwal et al. (2020)	2020	ECCV	-	-	Cover-based	Contextual Diversity				Diversity - Farthest-first Traversal Learnable - Reinforcement Learning	
Lin et al. (2020)	2020	ECCV	-	-	Clustering-based	L2 Distance				Diversity - Clustering	
Ash et al. (2020)	2020	ICLR	Gradient-based Uncertainty	Gradient	-	Diversity - Clustering					Patch
Casanova et al. (2020)	2020	ICLR	-	-	-	Learnable - Reinforcement Learning					
Dai et al. (2020)	2020	MICCAI	Gradient-based Uncertainty	Gradient	-	Latent Space Optimization & Nearest Neighbour Search					Slice

Table 2: Methodology summarization of surveyed active learning works.

Year	Venues	Uncertainty		Representativeness		Sampling Strategy	SemiSL	SelfSL	ADA	Region	Generative
		Method	Basic Metrics	Method	Basic Metrics						
2020	MICCAI	Multiple Inferences - MC Dropout Performance Estimation - Surrogate	Entropy IoU of all result	Cover-based	Cosine Similarity	Hybrid - Multi-round					
2020	MICCAI	Performance Estimation - Learnable	Loss	-	-	Top-k					
2020	MICCAI	Multiple Inferences - Model Ensemble	Margin	Discrepancy-based	Cosine Similarity	Hybrid - Multi-round					
2020	MICCAI	-	-	-	-	Learnable - Reinforcement Learning					
2020	TMI	Multiple Inferences - MC Dropout	Variance	Cover-based	Cosine Similarity	Hybrid - Multi-round					Slice, Pixel
2020	TMI	Multiple Inferences - Model Disagreement	Hausdorff Distance	-	-	Top-k					✓
2020	WACV	Single Model	Entropy	Discrepancy-based	H-Divergence	Hybrid - Fusion					✓
2021	CVPR	Probability of Misclassification	-	-	-	Class-balance					
2021	CVPR	Adversarial Training	Disagreement of Classifiers, margin	Discrepancy-based	H-Divergence	Hybrid - Fusion					✓
2021	CVPR	Performance Estimation - Learnable	Rank of Loss	Discrepancy-based	H-Divergence	Top-k					
2021	CVPR	Adversarial Training	Disagreement of Classifiers	-	-	Top-k					
2021	CVPR	Single Model	BSB	-	-	Class-balance					Superpixel
2021	CVPR	Single Model (w/ GNN)	Margin	Cover-based	L2 Distance of GCN-augmented Features	Top-k					
2021	ICCV	Single Model	Entropy	-	-	Diversity - Farthest-first Traversal					✓
2021	ICCV	-	-	Discrepancy-based	L2 Distance	Diversity - Clustering					
2021	ICCV	Performance Estimation - Surrogate	Temporal Output Discrepancy	-	-	Diversity - Clustering					Consistency
2021	ICCV	-	-	Discrepancy-based	Semantic and distinctive scores	Hybrid - Fusion					✓
2021	ICCV	Multiple Inferences - Model Disagreement	Inequality	-	-	Diversity - Clustering					Pseudo-label
2021	ICCV	Adversarial Samples	KL Divergence	Cover-based - Submodular	KL Divergence Bhattacharya Coefficient	Hybrid - Fusion					✓
2021	ICCV	Uncertainty-aware Models - MDN	Variance	-	-	Top-k					
2021	ICCV	Gradient-based Uncertainty	Influence	-	-	Top-k					
2021	JBHI	Performance Estimation - Surrogate	Dice	-	-	Top-k					Pseudo-label
2021	Media	Single Model	Entropy	-	-	Hybrid - Fusion					
2021	Media	Multiple Inferences - Data Disagreement	KL Divergence	-	-	Hybrid - Fusion					
2021	Media	Performance Estimation - Learnable	Loss	-	-	Hybrid - Fusion					
2021	MICCAI	Multiple Inferences - Data Disagreement	KL Divergence	-	-	Hybrid - Fusion					
2021	MICCAI	Performance Estimation - Learnable	Dice	-	-	Top-k					
2021	MICCAI	Single Model	Distance to Mean Probability	-	-	Top-k					Consistency
2021	MICCAI	Multiple Inferences - Model Ensemble	Variance	Discrepancy-based	Cosine Similarity	Diversity - Clustering Hybrid - Multi-round					Consistency
2021	MIDL	Single Model	Entropy	Cover-based	L2 Distance	Diversity - Clustering					Pseudo-label
2021	NeurIPS	Gradient-based Uncertainty	Fisher Information	-	-	Diversity - Solve Programming Problem					✓
2021	NeurIPS	-	-	Cover-based - Submodular	Gradient	Top-k					
2021	NeurIPS	Single Model	Margin	-	-	Diversity - Clustering					
2021	TMI	Multiple Inferences - Model Ensemble	Entropy	Discrepancy-based	Mutual Information	Hybrid - Fusion					
2021	TMI	-	-	Saliency Maps	Kurtosis Multivariate Radnomics Features Deep Saliency Features	Top-k					
2021	TPAMI	Single Model (in Feature Space)	Entropy	-	-	Top-k					✓

Table 2: Methodology summarization of surveyed active learning works.

Year	Venues	Uncertainty		Basic Metrics		Method	Representativeness		Sampling Strategy	SemiSL	SelfSL	ADA	Region	Generative
		Method	Method	Basic Metrics	Basic Metrics									
2022	Kothawade et al. (2022b)	-	-	Margin	-	Cover-based - Submodular	Gradient	-	Top-k	-	-	-	-	-
2022	Xie et al. (2022b)	Single Model	Single Model	Margin	-	Density-based	Energy	-	Hybrid - Multi-round	-	-	✓	-	-
2022	Wang et al. (2022b)	Gradient-based Uncertainty	Gradient	Gradient	-	-	-	-	Top-k	-	-	-	-	-
2022	Xie et al. (2022d)	Single Model	Margin, Gradient	Margin, Gradient	-	-	-	-	Top-k	-	-	✓	-	-
2022	Zhang et al. (2022a)	Single Model	Entropy	Entropy	-	Density-based	Mean Cosine Similarity of KNN	-	Hybrid - Equal Split	-	-	-	-	-
2022	Zhang et al. (2022b)	Single Model	Entropy	Entropy	-	-	-	-	Top-k	-	✓	-	-	-
2022	Parvaneh et al. (2022)	Multiple Inferences - Data Disagreement	Inequality	Inequality	-	-	-	-	Diversity - Clustering	-	-	-	-	-
2022	Xie et al. (2022a)	Single Model	Entropy	Entropy	-	-	-	-	Top-k	-	✓	-	-	-
2022	Qian et al. (2022)	-	-	-	-	Cover-based	Cosine Similarity	-	Top-k	-	-	-	-	-
2022	Wu et al. (2022a)	Single Model	Entropy	Entropy	-	Discrepancy-based	Cosine Similarity	-	Class-balance	-	-	-	-	-
2022	Wang et al. (2022c)	-	-	-	-	Density-based	KNN Density	-	Diversity - Clustering w/ Regularization	-	✓	-	-	-
2022	Kothawade et al. (2022a)	-	-	-	-	Cover-based - Submodular	Cosine Similarity	-	Top-k	-	-	-	-	-
2022	Chen et al. (2022b)	Gradient-based Uncertainty	Gradient	Gradient	-	-	-	-	Top-k	-	-	-	-	✓
2022	Hwang et al. (2022)	Single Model	Margin	Margin	-	Discrepancy-based	MMD	-	Hybrid - Multi-round	-	-	✓	-	-
2022	Yi et al. (2022)	Single Model	Least Confidence	Least Confidence	-	Self-supervised Learning	Loss of Pretext Task	-	Hybrid - Multi-round	-	✓	-	-	-
2022	Wu et al. (2022b)	Single Model	Entropy	Entropy	-	Density-based	GMM	-	Hybrid - Multi-round	-	✓	-	-	-
2022	Mahmood et al. (2022)	-	-	-	-	Discrepancy-based	Wasserstein Distance	-	Diversity - Solve Programming Problem	-	✓	-	-	-
2022	Hacohen et al. (2022)	-	-	-	-	Density-based	Inverse Average Distance to KNN samples	-	Diversity - Clustering	-	✓	-	-	-
2022	Jin et al. (2022a)	-	-	-	-	Clustering-based	Cosine Similarity	-	Diversity - Clustering	-	✓	-	-	-
2022	Jin et al. (2022c)	-	-	-	-	Clustering-based	L2 Distance	-	Class-balance	-	✓	-	-	-
2022	Jin et al. (2022b)	-	-	-	-	Clustering-based	L2 Distance	-	Diversity - Farthest-first Traversal	-	✓	-	-	-
2022	Dai et al. (2022)	Gradient-based Uncertainty	Gradient	Gradient	-	-	-	-	Latent Space Optimization & Nearest Neighbour Search	-	-	-	-	Slice
2022	Zhou et al. (2022)	Performance Estimation - Learnable	Dice	Dice	-	-	-	-	Top-k	-	-	-	-	-
2022	Azemi et al. (2022)	Performance Estimation - Surrogate	Dice	Dice	-	-	-	-	Top-k	-	-	-	-	-
2022	Nath et al. (2022)	Multiple Inferences - MC Dropout	Entropy	Entropy	-	-	-	-	Top-k	-	-	-	-	Pseudo-label
2022	Balaram et al. (2022)	Uncertainty-aware Model - EDL	Entropy	Entropy	-	-	-	-	Top-k	-	-	-	-	Pseudo-label & Consistency
2022	Wu et al. (2022d)	-	-	-	-	Cover-based	Cosine Similarity	-	Diversity - Clustering	-	-	-	-	Slice
2022	Bai et al. (2022)	Multiple Inferences - Model Disagreement	Entropy-weighted Dice Distance	Entropy-weighted Dice Distance	-	-	-	-	Hybrid - Fusion	-	-	-	-	Pseudo-label
2022	Kothawade et al. (2022c)	-	-	-	-	Cover-based - Submodular	Gradient	-	Top-k	-	-	-	-	-
2022	Yehuda et al. (2022)	-	-	-	-	Cover-based	L2 Distance	-	Graph-based Algorithm	-	-	-	-	-
2022	Mahapatra et al. (2022)	-	-	-	-	Saliency Maps	Graph-based Methods	-	Top-k	-	-	-	-	-
2022	Li et al. (2022)	Curriculum Learning & Noisy Sample Detection	Entropy	Entropy	-	-	-	-	Top-k	-	-	-	-	Pseudo-label
2022	Bengar et al. (2022)	Single Model	Entropy	Entropy	-	-	-	-	Class-balance	-	-	-	-	-
2023	Xie et al. (2023b)	-	-	-	-	Discrepancy-based	Wasserstein Distance	-	Latent Space Optimization & Nearest Neighbour Search	-	✓	-	-	-
2023	Lyu et al. (2023)	Multiple Inferences - Data Disagreement	Cross Entropy, Variance	Cross Entropy, Variance	-	-	-	-	Hybrid - Fusion	-	-	-	-	Pseudo-label
														Box

Table 2: Methodology summarization of surveyed active learning works.

Year	Venues	Uncertainty		Representativeness			SemiSL	SelfSL	ADA	Region	Generative
		Method	Method	Basic Metrics	Method	Basic Metrics					
Jung et al. (2023)	ICLR	Multiple Inferences - Model Ensemble	Entropy, Variance Ratio, BALD, Margin	-	-	-	-	-	-	Top-k	-
Xie et al. (2022c)	ICLR	Uncertainty-aware Model - EDL	Mutual Information & Entropy Expectation of Dirichlet Distribution	-	-	-	-	-	-	Hybrid - Multi-round	✓
Kim et al. (2023)	ICCV	Single Model	BvSB	-	-	-	-	-	-	Class-balance	Superpixel
Park et al. (2023)	ICLR	Uncertainty-aware Model - EDL	Mutual Information	-	-	-	-	-	-	Top-k	-
Sadafai et al. (2023)	ISBI	Multiple Inferences - MC Dropout Multiple Inferences - Model Disagreement	Variance Inequality	-	-	-	-	-	-	Hybrid - Fusion	-
Bai et al. (2023)	MICCAI	Multiple Inferences - Model Disagreement Gradient-based Uncertainty	KL divergence, gradient	Cover-based	L2 Distance	Diversity - Clustering	-	-	-	Diversity - Clustering	✓
Tang et al. (2023)	MICCAI	Multiple Inferences - Model Disagreement	KL divergence	-	-	-	-	-	-	Top-k	-
Qiu et al. (2023)	MICCAI	Single Model	Distance to 0.5	-	-	-	-	-	-	Top-k	Patch
Chen et al. (2023b)	MIDL	-	-	Loss of Self-supervised Pretext Tasks	-	-	-	-	-	Top-k	✓
Qu et al. (2023a)	NeurIPS	Multiple Inferences - Model Disagreement	Variance, Entropy, Overlap	-	-	-	-	-	-	Top-K	-
Lou et al. (2023)	TMI	-	-	Clustering-based	Consistency	Diversity - Clustering	Pseudo-label	-	-	Diversity - Clustering	✓
Du et al. (2023)	TPAMI	-	-	Discrepancy-based	Semantic and distinctive scores	Hybrid - Fusion	-	-	-	Hybrid - Fusion	✓
Wan et al. (2023)	TPAMI	Adversarial Training	Disagreement of Classifiers	-	-	-	-	-	-	Top-k	-

these samples for semi-supervised training. Specifically, predictions of the weakly augmented images were assigned as pseudo-labels, and the outputs of the strongly augmented images were forced to be consistent with the pseudo-labels. However, when the confidence level was lower than the threshold, they used these samples for AL. A balanced uncertainty selector and an adversarial instability selector were used to select samples for oracle annotation. They validated the effectiveness of their proposed method in grading metastatic epidural spinal cord compression with MRI images.

4.2. Self-supervised Learning: Utilizing Pre-trained Model

Integrating semi-supervised learning with AL has achieved successful applications. However, its effectiveness is constrained by the dataset size. This limitation is particularly evident for medical imaging datasets which are relatively small. In clinical practice, plenty of raw medical images are stored in hospital information systems without human annotation. Self-supervised learning (Krishnan et al., 2022) could be a vital tool for mining information hidden in those raw images. Its idea is to train the model with the supervision of the data itself, thus allowing pre-training on a large unlabeled dataset. Many studies have shown self-supervised pre-trained models could achieve impressive performance by finetuning on a few randomly selected labeled samples in medical image analysis (Azizi et al., 2021; Tang et al., 2022b). A natural expectation is to integrate active learning strategies with self-supervised learning, aiming for higher annotation efficiency over mere random sampling (Lüth et al., 2024). Besides, these models could also act as a powerful feature extractor, which provides good initialization for AL. In this section, we will first introduce how self-supervised models solve the cold-start problem in AL and then explore different ways of integrating AL with self-supervised learning.

4.2.1. Cold-start Problem in Active Learning

Current AL methods usually require an initial labeled dataset to train the model for start and ensure reliable informativeness evaluation. However, when the initial labeled set is small or even absent, the performance of these AL methods drops dramatically, sometimes even worse than random sampling (Chen et al., 2023b; Hacoheh et al., 2022; Yehuda et al., 2022). Studies also showed that simply integrating self-supervised learning with AL baselines leads to inferior performance than random sampling (Bengar et al., 2021; Xie et al., 2023a). This phenomenon is known as the cold-start problem which commonly exists in AL of various domains, including medical image analysis (Liu et al., 2023a). Tackling the cold-start problem is vital for improving the efficacy of AL, especially in the medical domains where annotation costs are extremely high. A key solution to the cold-start problem in AL is selecting the optimal set of initial labeled samples, which requires different strategies than the existing AL methods.

Early attempts focused on utilizing the fully supervised pre-trained models to address the cold-start problem in AL.

Zhou et al. (2017) and their subsequent work (Zhou et al., 2021c) used ImageNet pre-trained models to select samples for annotation from completely unlabeled datasets in medical image analysis. They combined entropy and disagreement as informativeness metrics, where the disagreement was the KL divergence of prediction probabilities between different patches of the same sample. They also introduced randomness to balance exploration and exploitation. Experiments on two colonoscopy datasets and a CT pulmonary embolism detection dataset showed superior performance than other competitors.

Self-supervised pre-trained models offer a good initialization for effectively tackling the cold-start problem in AL. ALPS (Yuan et al., 2020a) was the first to introduce the cold-start problem in AL and employed self-supervised pre-trained models to address this issue. Based on a contrastive learning feature extractor, CALR (Jin et al., 2022a) employed BIRCH clustering and chose the samples with maximum information density within each cluster for labeling. Compared to k-Means, BIRCH clustering is less sensitive to outliers and can further identify noisy samples. TypiClust (Hacoheh et al., 2022) theoretically proved that querying typical samples is more beneficial for a low annotation budget. Therefore, based on self-supervised features, TypiClust selected samples from high-density areas of each k-Means cluster. Beyond that, Yehuda et al. (2022) employed a graph-based greedy algorithm to select the optimal initial samples based on self-supervised features. In CT segmentation, Nath et al. (2022) proposed ProxyRank which designed new pretext tasks for self-supervised pre-training. The model was trained to learn the threshold segmentation by an abdominal soft-tissue window. Results indicated that the proposed method significantly outperforms random sampling in selecting initial samples. To benchmark the effectiveness of different cold-start AL methods in 3D medical image segmentation, Liu et al. (2023a) reproduced ALPS, CALR, TypiClust, and ProxyRank on five MSD datasets (Antonelli et al., 2022). Results showed that TypiClust stands out from the four competitors. However, no method consistently outperformed random selection on all five datasets, which calls for further exploration of cold-start AL in medical image analysis.

4.2.2. Combination of Active Learning and Self-supervised Learning

Features: The simplest way is leveraging the high-quality features of the self-supervised pre-trained models. Many studies are based on a powerful self-supervised feature extractor (Pourahmadi et al., 2021; Jin et al., 2022a; Hacoheh et al., 2022; Yehuda et al., 2022).

Pretext tasks in self-supervised learning are designed to derive supervision directly from the data itself. Solving these pretext tasks on large-scale unlabeled data, the model acquires useful feature representations that reflect data characteristics. Different pretext tasks correspond to different pre-training paradigms, the typical ones including rotation prediction (Gidaris et al., 2018), contrastive learning (He et al., 2020), and masked image modeling (He et al., 2022), etc. Related works generally employed the loss of pretext task for AL. In

Chen et al. (2023b), the loss of contrastive learning was used to tackle cold-start problem for AL in medical image analysis. They assumed that samples with higher losses are more representative of the data distribution. Specifically, they pre-trained on the target dataset with momentum contrastive learning (He et al., 2020), and then used k-Means clustering to partition the unlabeled data into multiple clusters, selecting the samples with the highest contrastive loss within each cluster for annotation. They then selected samples with the highest contrastive loss in each cluster for annotation. The proposed method addressed the class imbalance caused by the bias of traditional AL methods, and the failure to detect anomalies when the number of initially labeled datasets was limited. This method showed superior performance in PathMNIST, OrganMNIST, and BloodMNIST (Yang et al., 2023). Yi et al. (2022) found a strong correlation between the loss of pretext tasks and the loss of downstream tasks. Thus, they initially focused on annotating samples with higher loss of pretext tasks and later shifted to those with lower loss. Results showed that rotation prediction performed the best among different pretext tasks.

Others: Furthermore, recent works leverage self-supervised learning in other ways for AL. Zhang et al. (2022b) introduced one-bit annotation into AL for classification tasks. In this setting, oracles only returned whether the prediction was right or wrong rather than its specific class label. Contrastive learning was adopted to pull the correct predictions closer and push away wrong predictions from their predicted classes. Results indicated that the proposed method outperforms other AL methods regarding bit information. Du et al. (2021) integrated contrastive learning into AL to address the class distribution mismatch, where unlabeled data includes samples out of the class distribution of the labeled dataset. In this work, contrastive learning was adopted to filter samples of mismatched classes and highlight sample informativeness by carefully setting negative samples. Their extended work Du et al. (2022) provided more theoretical analysis and experimental results and also integrated existing label information into the proposed framework.

4.3. Region-based Active Learning: Smaller Labeling Unit

Most AL works require the oracle to label the full image in medical image analysis. However, labeling a full image can introduce redundancy in fine-grained tasks like segmentation or detection, resulting in an inefficient use of the annotation budget. For the example of abdomen multi-organ segmentation, large organs that are easy to segment (e.g., liver or spleen) do not need exhaustive annotation. Instead, those budgets would be better spent on small organs that are hard to segment, like the esophagus and adrenal glands. To address this issue, images could be divided into non-overlap regions for higher annotation efficiency, and experts can opt to annotate specific regions within an image, which is termed “region-based active learning”. This section introduces region-based active learning from the perspectives of patches and superpixels, which means that AL methods mentioned in this section selected either patches or superpixels within an image for annotation.

4.3.1. Patches

Patches are most commonly used in region-based active learning, generally represented as square boxes. Mackowiak et al. (2018) combined uncertainty and annotation cost to select the informative patches for annotation. In retinal blood vessels segmentation of fundus images, Xu et al. (2021) selected patches with the highest uncertainty for annotation. Furthermore, they utilized latent-space mixup to encourage linearization between labeled and unlabeled samples, thus leveraging unlabeled data to improve performance. Casanova et al. (2020) employed deep reinforcement learning to automatically select informative patches for annotation. In grey matter and white matter segmentation of pathology images, Lai et al. (2021) first split the whole slide image into multiple patches. With the confidence (i.e., maximum predictive probability) of each patch, a mean filter of size 5x5 is used to aggregate the confidence of the neighboring patches. As a result, one aggregated metric corresponded to a region of 5x5 patches, and regions with the highest uncertainty were selected for annotation. Besides, Qiu et al. (2023) adopted an adaptive region selection with non-square patches for whole-slide images. Instead of sampling square patches, they dynamically determined the size of each non-square patch by carefully locating an informative area on each slide. The proposed method demonstrated improvement in annotation efficiency and robustness to AL hyperparameters compared to the square patch baseline.

4.3.2. Superpixels

Superpixels are also widely used in region-based active learning. Superpixel-based AL initially pre-segments the images with superpixel generation algorithms based on color and texture (Achanta et al., 2012; Van den Bergh et al., 2012), and then calculates the informativeness of each superpixel. The informativeness metric of each superpixel is the average of its constituent pixels. Siddiqui et al. (2020) adopted uncertainty and disagreement between different viewpoints to select informative superpixels for annotation. In OCT segmentation, Kadir et al. (2023) proposed edge-based entropy and divergence to select highly uncertain superpixels for annotation. Experiments on three datasets were conducted to illustrate the effectiveness of their method. For superpixel selection, Cai et al. (2021) proposed dominant labeling which is the majority class label of all pixels in the superpixel. They assigned the dominant labeling to every pixel within a superpixel, thus eliminating the need for detailed delineation. They further introduced a class-balanced sampling strategy to better select superpixels containing minority classes. Results showed that dominant labeling with superpixels significantly outperforms precise labeling with patches under the same number of labeling clicks. As a follow-up work, Kim et al. (2023) proposed to adaptively merge and split spatially adjacent, similar, and complex superpixels, respectively. This approach yielded better performance than Cai et al. (2021) with dominant labeling. Li et al. (2023) utilized the superpixels to estimate the regional consistency which is the difference between the prediction and the dominant class of

each superpixel. Combining other metrics like entropy and diversity, they selected the most uncertain foreground and background superpixel to reduce the annotation cost.

4.4. Generative Model: Data Augmentation and Generative Active Learning

In recent years, the advancement of deep generative models enabled high-quality generation and flexible conditional generation. For example, a trained model could generate the corresponding lung X-ray scan when conditioned on a lung mask. By integrating generative models, we can further improve the annotation efficiency of AL. In this section, we discuss how AL can be combined with generative models from two aspects: data augmentation and generative active learning.

4.4.1. Synthetic Samples as Data Augmentation

The simplest approach considers the synthetic sample produced by generative models as advanced data augmentation. These methods utilize label-conditioned generative models. As a result, it's guaranteed that all synthetic samples are correctly labeled since specifying the labels is a prerequisite for data generation. This method enables us to acquire more labeled samples without any additional annotations. [Tran et al. \(2019\)](#) argued that most synthetic samples produced by generative models are not highly informative. Therefore, they first adopted the BALD uncertainty to select samples for annotation, then trained a VAE-ACGAN on these labeled data to generate more informative synthetic samples. [Mahapatra et al. \(2018\)](#) used conditional GANs to generate chest X-rays with varying diseases to augment the labeled dataset. Then, MC Dropout was used to select and annotate highly uncertain samples. With the help of AL and synthetic samples, they achieved performance near fully supervised using only 35% of the data. Training conditional generative models requires a large amount of labeled data, while the labeled dataset in AL is often relatively small. To address this issue, [Lou et al. \(2023\)](#) proposed a conditional SinGAN ([Shaham et al., 2019](#)) that only requires one pair of images and masks for training. The SinGAN improved the annotation efficiency for nuclei segmentation. [Chen et al. \(2022b\)](#) integrated implicit semantic data augmentation (ISDA) ([Wang et al., 2021](#)) into AL. They initially used ISDA to augment unlabeled samples, then selected samples with large diversity between different data augmentations for annotation. The model is trained on both the original data and its augmentations. [Mahapatra et al. \(2024\)](#) trained a VAE for synthesizing informative and non-redundant samples. These samples are generated by first sampling in the latent space of VAE and feeding them to the VAE decoder. Besides, scores of label preservation and redundancy avoidance were adopted to pick the most informative synthetic samples. The proposed method was tested in chest X-ray classification and multiple toy datasets from MedMNIST ([Yang et al., 2023](#)).

4.4.2. Generative Active Learning

Generative active learning selects synthetic samples produced by generative models for oracle annotation, thus

without requiring a large unlabeled sample pool. The advantage of this approach lies in its ability to continuously search the data manifold through generative models. It's worth noting that works in this section follow the setting of membership query synthesis, while works in the last section follow the setting of pool-based active learning. This distinction arises because generative models in the last section were solely utilized to augment existing labeled datasets. [Zhu and Bento \(2017\)](#) attempted to generate uncertain samples with GAN for expert annotation. Unfortunately, the quality of the generated samples was low and included many samples with indistinguishable classes. Since experts find it difficult to annotate low-quality synthetic samples, alternative methods are needed to annotate these samples. [Chen et al. \(2021\)](#) first trained a bidirectional GAN to learn the data manifold. They then selected uncertain areas in the feature space and generated images within these regions using bidirectional GAN. Finally, they used physics-based simulation to provide labels for the generated samples. In calcification level prediction in aortic stenosis of CT, they improved annotation efficiency by up to 10 times compared to random generation.

4.5. Active Domain Adaptation: Tackling Distribution Shift

Domain Adaptation (DA) ([Guan and Liu, 2021](#)) has wide applications in medical image analysis. It aims to transfer knowledge from the source to the target domain, thus minimizing annotation costs. Currently, the most common setting of DA is unsupervised domain adaptation (UDA), in which the source domain is labeled while the target domain is unlabeled. For the example of abdominal multi-organ segmentation, we can train a domain-adaptive segmentation model with labeled MR images alongside unlabeled CT images to achieve good performance on the CT domain ([Liu et al., 2023b](#)). However, the performance of UDA still lags behind fully supervised learning in the target domain. Selecting and annotating informative samples would be beneficial to bridge this gap. This setting is known as active domain adaptation (ADA). For better queries in ADA, one should consider both uncertainty and representativeness regarding the target domain. The latter is commonly referred to as domainness or targetness in ADA. This section reviews the image-wise and region-wise ADA.

4.5.1. Image-wise Active Domain Adaptation

In this section, ADA methods performed an image-level selection, which involves most of the ADA works. [Su et al. \(2020\)](#) was the first to introduce the concept of ADA and combined domain adversarial learning with AL. Through a domain discriminator and task model, they performed importance sampling to select target domain samples that are uncertain and highly different from the source domain. [Fu et al. \(2021\)](#) combined query-by-committee, uncertainty, and domainness in ADA. They adopted a domain discriminator to select samples with high domainness and employed Gaussian kernels to filter out anomalous and source-similar samples of the target domain. Random sampling was also used to improve diversity. [Prabhu et al. \(2021\)](#) performed k-Means clustering

on target domain samples and selected cluster centers for annotation. The cluster centers were weighted by uncertainty, thus ensuring that selected samples were uncertain and diverse. In segmentation tasks, [Ning et al. \(2021\)](#) introduced the idea of anchors in ADA. They concatenated features of different classes from the source domain images. Cluster centers of these concatenations were referred to as anchors. They then computed the distance between each target sample and its nearest anchor. Target samples with the highest distance were requested for annotation. [Xie et al. \(2022b\)](#) introduced the concept of energy ([LeCun et al., 2006](#)) into ADA. The energy is inversely proportional to the likelihood of the data distribution. In this work, the model trained on the source domain was used to calculate the energy of target domain samples. Samples with high energy were selected for annotation, which suggested they are representative of the target domain and substantially different from the source data. [Huang et al. \(2023\)](#) selected samples with high uncertainty and prediction inconsistency to their nearest prototypes. In the context of medical image analysis, [Chen et al. \(2023a\)](#) tackled domain shifts in the setting of federated active learning. They proposed an EDL-based framework with a global model across all clients and local models for each client. In this work, the EU is related to domain shifts between the global model and local data. Therefore, the AUs of the global and local models were calibrated by the EU, thus improving performance. Results on multiple medical imaging datasets showed its effectiveness in reducing annotation costs. In nasopharyngeal carcinoma tumor segmentation, [Wang et al. \(2023\)](#) proposed a source-domain and target-domain dual-reference strategy to select informative samples for annotation. Specifically, the features of the source samples were clustered and the cluster centers were reference samples. Target samples with the highest and lowest similarity to the references are selected for annotation, which were treated as domain-invariant and domain-specific samples, respectively.

4.5.2. Region-wise Active Domain Adaptation

To better utilize the annotation budget, some ADA works also selected patches or superpixels within an image for annotation. [Shin et al. \(2021\)](#) proposed LabOR, which first used a UDA pre-trained model to generate pseudo-labels for target samples, which was used to train two segmentation heads. They maximized the disagreements between the two heads and annotated regions that exhibited the most disagreement. LabOR achieved performance close to full supervision with only 2.2% of target domain annotations. In [Xie et al. \(2022a\)](#), uncertainty and regional impurity were used to select and annotate the most informative patches. Regional impurity measured the number of unique predicted classes within the neighborhood of a pixel, which presents the edge information. They used extremely small patches (e.g., size of 3x3) for annotation and achieved performance close to full supervision with only 5% of the annotation cost. [Wu et al. \(2022b\)](#) proposed a density-based method to select the most representative superpixels in the target domain for annotation. They employed Gaussian mixture models (GMM) as density

estimators for superpixels in both the source and target domains, aiming to select those with high density in the target domain and low density in the source domain.

5. Active Learning for Medical Image Analysis

Due to the potential of significantly reducing annotation costs, AL is receiving increasing attention in medical image analysis. The unique traits of medical imaging require us to design specialized AL methods. Building on the foundation of the previous two sections, this section will focus on introducing AL works tailored to medical image analysis across different tasks, including classification, segmentation, and reconstruction.

Additionally, in Table 3, we list all the AL works related to medical image analysis in this survey, providing the name of the used dataset, its modality, ROIs, and corresponding clinical and technical tasks.

5.1. Active Learning for Medical Image Classification

Common clinical tasks like disease diagnosis, cancer staging, and prognostic prediction can be formulated as medical image classification. Most AL works in medical imaging classification directly employ general methods, such as using class-balancing sampling in §3.3.2 to mitigate the long-tail effect of medical imaging datasets. However, specialized design of AL algorithms is required for certain modalities of medical image classification. For example, the classification of chest X-rays often involves the idea of multi-label. Besides, classifying pathological whole-slide images typically needs to be formulated as a multiple-instance learning problem. This section will introduce AL works specifically targeted at classification problems in chest X-rays and pathological whole-slide images.

5.1.1. Chest X-ray and Multi-label Classification

Chest X-ray examinations are crucial for screening and diagnosing lung, cardiovascular, skeletal, and other thoracic diseases. Computer-aided diagnosis in this domain has been extensively researched, including AL works aimed at reducing annotation costs for physicians. [Mahapatra et al. \(2021\)](#) introduced saliency maps to select informative samples for annotation. To aggregate the per-pixel saliency maps into a single scalar, they explored three different approaches, including computing the kurtosis of the saliency map, utilizing multivariate radiomic features, and combining deep features of autoencoders and clustering. Results demonstrated that the aggregation using deep features performs the best. [Nguyen et al. \(2021\)](#) introduced a gist-set to select samples near the decision boundary. Besides, uncertain samples with high entropy were sent for annotation, while the confident samples were assigned as pseudo-labels. Additionally, they adopted momentum updates to enhance the stability of the sample predictions. To handle the annotation noise, [Bernhardt et al. \(2022\)](#) proposed a framework called ‘active label cleaning’. This framework ranked samples based on estimated label

Table 3: Surveyed Works of Active Learning related to Medical Image Analysis. “-” stands for such information is not available for not provided by the authors or not in the case.

Year	Venues	Modality	ROIs	Dataset	Sampling Unit	Size (Train+Val/Test) or (Train+Val/Test)	Initial Pool Size	Budget per Round	Clinical Task	Technical Task
Goniz et al. (2017)	2017	arXiv	Skin	ISIC 2017	image	2,000 (1,600/400)	600	35	Skin Cancer Diagnosis	Classification
Zhou et al. (2017)	2017	CVPR	Colon Colon Lung	in-house in-house in-house	image image PE candidate	4,000 (2,000/2,000) 28,250 (16,300/11,950) 6,255 (3,840/2,415)	-	-	Image Quality Assessment Polyp Detection Pulmonary Embolism Detection	Classification Classification Classification
Gai et al. (2017)	2017	ICML	Skin	ISIC 2016	image	400 (200/200)	100	20	Skin Cancer Diagnosis	Classification
Yang et al. (2017)	2017	MICCAI	Colon Lymph Node	GIaS in-house	patch image	165 images (80/5/80) 74 (37/37)	8	8	Gland Segmentation Lymph Node Segmentation	Segmentation Segmentation
Odion et al. (2017)	2017	MICCAI	Eye	e-ophtha	patch	20,148 (17,520/666/1,972)	160	32	Exudate Classification	Classification
Beluch et al. (2018)	2018	CVPR	Eye	EyePacs	image	88,702 (67,961/13,000/17,741)	1,000	5,000	Diabetic Retinopathy Detection	Classification
Xu et al. (2018)	2018	CVPR	Colon	GIaS	patch	165 images (80/5/80)	8	8	Gland Segmentation	Segmentation
Sourati et al. (2018)	2018	DLMIA	Brain	dHCP in-house	patch	66 patients 25 patients	3 patients	50	Brain Extraction	Segmentation
Sourati et al. (2019)	2019	TMI	Head	in-house	region, image, volume	1,247 volumes (934/313)	1/32	double	Intracranial Hemorrhage Detection	Segmentation
Mahapatra et al. (2018)	2018	MICCAI	Chest	SCR & Chestx-ray8	image	647 (247/400)	10%	5%	Lung Segmentation Thoracic Disease Diagnosis	Segmentation Classification
Zheng et al. (2019)	2019	AAAI	Colon Fungus Heart	GIaS in-house HWSMR 2016	patch slice	1,530 784 20 volumes (10/10)	0	30% & 50% every 2, 10, 20, 40, 80 slices	Gland Segmentation Fungus Segmentation Whole-heart Segmentation	Segmentation Segmentation Segmentation
Qi et al. (2019)	2019	JBHI	Breast	BreakHis	image	See paper for details	-	-	Breast Cancer Diagnosis	Classification
Sudafi et al. (2019)	2019	MICCAI	Blood	in-house	image	208 (188/20)	30	-	Red Blood Cell Detection	Object Detection
Shi et al. (2019)	2019	MICCAI	Skin	ISIC 2017	image	4,332 (3,582/1,50/600)	10%	10%	Skin Lesion Diagnosis	Classification
Gu et al. (2018)	2019	TBME	Breast Colon	Gu et al. (2017) Ye et al. (2016)	image	1,366 (1,042/324) 7,847 (3,947/3,900)	-	-	Endomicroscopy Mosaic Classification Gastrointestinal Image Classification	Classification Classification
Zheng et al. (2020)	2020	AAAI	Heart Mouse	HWSMR 2016 Lee et al. (2015)	slice slice	20 volumes (10/10) 4 volumes (3/1)	0	every 5, 10, 20, 40, 80 slices every 4, 16, 64 slices	Whole-heart Segmentation Neuron Boundary Segmentation	Segmentation Segmentation
Lin et al. (2020)	2020	ECCV	Mouse Synapses Mouse Synapses & Mitochondria	EM-R50	image	48.7K (28.7K/20K) 15K (10K/5K)	-	1,280	Synapse Detection Mitochondria Segmentation	Segmentation Segmentation
Dai et al. (2020)	2020	MICCAI	Brain	BrATS 2019	volume slice	335 patients (260/75)	10 500	10 500	Brain Tumor Segmentation	Segmentation
Li and Yin (2020)	2020	MICCAI	Colon Brain	GIaS ISeg	2D patch 3D patch	27,200 16,380	10%	20% 50%	Gland Segmentation Infant Brain Segmentation	Segmentation Segmentation
Liu et al. (2020)	2020	MICCAI	Lung	DeepLusion	volume	1,281 (1,000/281)	10%	10%	Pulmonary Nodule Detection	Object Detection
Shen et al. (2020)	2020	MICCAI	Breast	in-house	image	3,441 (2,767/306+368)	10%	10%	Breast Cancer Region Segmentation	Segmentation
Wang et al. (2020b)	2020	MICCAI	Lung Eye	Towhee EyePacs	image image	7K (3.5K/3.5K) 4,460 (2,230/2,230)	5% 10%	about 7% about 11%	Lung Disease Detection Diabetic Retinopathy Detection	Classification Classification
Huang et al. (2020)	2020	TMI	Hip & Thigh	TCIA & in-house	slice	20 volumes	5%	5%	Muscle Segmentation	Segmentation
Shen et al. (2021)	2021	ISBI	Chest	in-house	region	10,966 images (71:12)	0%	5% & 10%	Rib Fracture Recognition	Object Detection
Shen et al. (2021)	2021	ISBI	Brain	BrATS2018	slice	285 volumes (233:52)	5 volumes	100 & 200	Brain Tumor Segmentation	Segmentation
Zhao et al. (2021)	2021	JBHI	Skin Hand	ISIC 2017 RSNA Bone Age Dataset	image image	2,000 (1,600/400) 209 (139/20/50)	600 10	100 10	Skin Lesion Segmentation Finger Bone Segmentation	Segmentation Segmentation
Ozdemir et al. (2021)	2021	KBS	Lower Extremities	in-house	slice volume	36 volumes (25/2/9)	64 1	32 1	Musculoskeletal Segmentation	Segmentation
Wu et al. (2021)	2021	Media	Lung	CC-CUII	volume	962 (7:1:2)	-	10	COVID-19 Diagnosis	Classification
Zhou et al. (2021c)	2021	Media	Colon Colon Lung	in-house in-house in-house	image image PE candidate	4,000 (2,000/2,000) 28,250 (16,300/11,950) 6,255 (3,840/2,415)	-	-	Image Quality Assessment Polyp Detection Pulmonary Embolism Detection	Classification Classification Classification
Wang and Yin (2021)	2021	MICCAI	Bacterial Cells Human Bone Marrow Human Adipocyte Cells Various Tissues & Spectres	VGG Cell MBM ADI DCC	image	150 (50/100) 29 (15/14) 150 (50/100) 176 (100/76)	-	10%	Cell Counting	Keypoint Localization
Xu et al. (2021)	2021	MICCAI	Eye Eye	DRIVE ROSE-1	region	40 images (20/20) 117 images	-	10% patches in selected images	Retina Vessel Segmentation	Segmentation

Table 3: Methodology summarization of surveyed active learning works. “-” stands for such information is not available for not provided by the authors or not in the case.

Year	Venues	Modality	ROIs	Dataset	Sampling Unit	Size (Train/Val/Test) or (Train+Val/Test)	Initial Pool Size	Budget per Round	Clinical Task	Technical Task
Zhou et al. (2021b)	MICCAI	CT	Lung Colon Kidney	MSD KITS 19 (Rahman et al., 2021)	scribble, bounding box, extreme point	96 volumes (64/32) 190 volumes (126/64) 210 volumes (168/42)	-	-	Tumor Segmentation Kidney & Tumor Segmentation	Segmentation
Chong et al. (2021)	MICCAIW	CT	Chest	in-house	image	15,153 (8:1:1)	10%	10%	COVID-19 Diagnosis	Classification
Nguyen et al. (2021)	MIDL	X-ray	Chest	in-house	image	135,309 (131,030/4,279)	6,550	≤6,550	Diagnosis of Ainspae Opacity & Lung Lesion Diagnosis of Pneumonia Detection of Pleural Effusion	Classification Classification Classification
Mahapatra et al. (2021)	TMI	X-ray Histopathology	Chest Colon	ChestX-ray8 GAS	image patch	112,120 (7:1:2) 165 images	10%	10%	Thoracic Disease Diagnosis Gland Segmentation	Classification Segmentation
Nath et al. (2021)	TMI	CT MRI	Pancreas Hippocampus	MSD	volume	281 (221/50/30) 263 (163/50/50)	20 10	5 1	Pancreas & Tumor Segmentation Hippocampus Segmentation	Segmentation Segmentation
Chen et al. (2021)	TPAMI	CT	Heart	in-house	volume	168 (126/42)	-	-	Calcification Level Prediction in Aortic Stenosis	Classification
Kothawade et al. (2022b)	AAAI	X-ray	Chest	PneumoniaMNIST	image	5,856 (4,708/524/624)	-	-	Pneumonia & Normal Classification	Classification
Wang et al. (2022b)	AAAI	cryo-ET (simulated)	-	SHREC'19	image	25,000 (24,000/1,000)	10%	5%	Subtomogram Classification	Classification
Qian et al. (2022)	CVPR	X-ray	Head Hand	Kaggle Payer et al. (2019)	image	400 (150/250) 909 (609/300)	-	-	Cephalometric Landmark Detection Hand Landmark Detection	Keypoint Localization Keypoint Localization
Zhang et al. (2022a)	CVPR	MRI	Spine	in-house	image	7,295 (7:1:2)	10%	5%	Diagnosis of Metastatic Epidural Spinal Cord Compression	Classification
Jin et al. (2022c)	Knowledge-based Systems	Dermoscopy	Skin	ISIC 2020	image	33,126 (8:1:1)	0	10%, 30% & 50%	Skin Lesion Classification	Classification
Jin et al. (2022b)	Knowledge-based Systems	Dermoscopy X-ray	Skin Chest	ISIC 2018 Jaeger et al. (2013)	image image	2,594 (3:1:1) 707 (3:1:1)	0	40%-60% with step of 2.5% 30%-50% with step of 2.5%	Skin Lesion Segmentation Lung Segmentation	Segmentation Segmentation
Aizeni et al. (2022)	Media	MRI Histology	Brain Brain	SATA in-house	boundary pixel	35 volumes 15 stacks	-	-	Brain Structure Segmentation Brain Structure Segmentation	Segmentation Segmentation
Dai et al. (2022)	Media	MRI	Brain	Brats2019 MALC	image image	335 patients (260/75) 30 patients (20/10)	0.5% 6%	1% 6%	Brain Tumor Segmentation Brain Structure Segmentation	Segmentation
Zhou et al. (2022)	Media	CT Colonoscopy	Lung Colon Kidney Colon	MSD KITS 19 CVC-ClinicDB	scribble, bounding box, extreme point	96 volumes (64/32) 190 volumes (126/64) 210 volumes (168/42) 29 sequences (23/5/3)	-	-	Tumor Segmentation Kidney & Tumor Segmentation Polyp Segmentation	Segmentation
Nath et al. (2022)	MICCAI	CT	Liver Hepatic Vessels	MSD	volume	131 (105/26) 303 (242/61)	0	5% 2%	Liver & Tumor Segmentation Hepatic Vessels & Tumor Segmentation	Segmentation
Bai et al. (2022)	MICCAI	Wireless Capsule Endoscopy	Colon	CAD-CAP	image	1,812 (4:1)	0	10%	Polyp Segmentation	Segmentation
Balaran et al. (2022)	MICCAI	X-ray	Chest	Chestx-ray8	image	112,120 (7:1:2)	2% 5%	0.5% 1%	Thoracic Disease Diagnosis	Classification
Wu et al. (2022a)	MICCAI	CT MRI CT CT	Liver	LITS CHAOS Shiwei07 MSD	slice	130 (min) CT-20 (train) MRI: 120 (0/60/0) 20 (test) 131 (test)	-	-	Liver & Tumor Segmentation	Segmentation
Kothawade et al. (2022c)	MICCAIW	X-ray Histopathology Microscopy Dermoscopy Fundus	Chest Colon Peripheral Blood Skin Eye	PneumoniaMNIST PathMNIST BloodMNIST ISIC 2018 APTOS-2019	image	5,856 (4,708/524/624) 107,180 (89,996/10,004/7,180) 17,092 (11,959/1,712/3,421) 2,594	100 50 238 1947 673	10 500 20 20 40	Pneumonia & Normal Classification Survival Prediction Cell Type Classification Skin Lesion Diagnosis Diabetic Retinopathy Detection	Classification
Alkhila and Yeung (2022)	MLHC	Laparoscopy	Gallbladder	CholecSeg8k	image	8,080 (4,640/1,600/1,640)	10	10%	Segmentation of Laparoscopic Surgical Images	Segmentation
Bernhardt et al. (2022)	Nature Communications	X-Ray	Chest	Noisy CXR	image	26.6K	-	-	Thoracic Disease Diagnosis	Classification
Li et al. (2022)	TMI	Histopathology	Prostate	PANDA	image	11,000	10%	10%	Gleason Grading of Prostate Cancer	Classification
Wu et al. (2022c)	TMI	CT Colonoscopy	Lung Colon	CC-CCH HyperKvasir	slice image	750 108,676 1000 2,717	20%	5%	Lung Lesions Segmentation COVID-19 and Pneumonia Diagnosis Polyp Segmentation Colonoscopic Lesion Classification	Segmentation Classification Segmentation Classification
Mahapatra et al. (2022)	TMI	X-ray	Chest	ChestXpert	image	65,240 patients (7:1:2)	10%	10%	Thoracic Disease Diagnosis	Classification

Table 3: Methodology summarization of surveyed active learning works. “-” stands for such information is not available for not provided by the authors or not in the case.

Year	Venues	Modality	ROIs	Dataset	Sampling Unit	Size (Train/Val/Test) or (Train/Val/Test)	Initial Pool Size	Budget per Round	Clinical Task	Technical Task
Khanal et al. (2023)	2023 arXiv	MRI X-Ray	Brain Chest	BraTS 2018 COVID-QU-Ex	slice image	5,846 (3,675/1,009/1,164) 5,826 (3,728/932/1,166)	200 100	100 100	Brain Tumor Analysis COVID-19 Diagnosis	Classification & Segmentation Classification
Chen et al. (2023a)	2023 arXiv	Dermoscopy Colonoscopy MRI Fundus	Skin Histopathology Colon Prostate Eye	Fed-BIC Fed-Camelyon Fed-PolyP Fed-Prostate Fed-Fundus	image image	21,989 (17,591/4,398) 455,954 (364,761/91,193) 2,187 (1,751/436) 1,867 (1,541/326) 1,060 (849/211)	500 50 20	500 50 20	Skin Lesion Diagnosis Detection of Cancer Metastases Lymph Nodes Polyp Segmentation Prostate Segmentation Retina Vessel Segmentation	Classification Segmentation
Wang et al. (2023)	2023 arXiv	MRI	Nose	in-house	slice	1,057 patients (7:1:2)	-	20%	Nasopharyngeal Carcinoma Tumor Segmentation	Segmentation
Jin et al. (2023b)	2023 EAAI	Dermoscopy X-ray	Skin Chest	ISIC 2018 Jaeger et al. (2013)	image	2,594 patients (3:1:1) 704 patients (3:1:1)	0 0	10%-30% with step of 2.5% 30%-50% with step of 2.5%	Skin Lesion Segmentation Lung Segmentation	Segmentation
Jiménez et al. (2023)	2023 ICCVW	Histopathology	Gland	GlaS	patch	165 images (85/16/64)	43%	5%	Gland Segmentation	Segmentation
Sudati et al. (2023)	2023 ISBI	Histopathology	Breast	CAMELYON17	WSI	500	0	2	Detection of Cancer Metastases Lymph Nodes	Classification
Gaillhofer et al. (2023b)	2023 Media	MRI MRI	Prostate Hippocampus	PROMISE 2012 MSD	slice	1,377 (1,020/109/248) 9,270 (7,163/501/1,571)	10 10	10 10	Prostate Segmentation Anterior and Posterior Hippocampus Segmentation	Segmentation
Li et al. (2023)	2023 Media	Ultrasound CT Ultrasound X-Ray	Breast Liver Breast Chest	in-house in-house BUSI CXRSat	in-house supersized volume or patch	3,200 images (2,600/300/300) 8,797 images (8,000/400/397) 647 images (400/123/124) 704 images (500/102/102)	260 200 40 10	130 200 40 10	Breast Tumor Segmentation Liver Segmentation Breast Tumor Segmentation Lung Segmentation	Segmentation
Bai et al. (2023)	2023 MICCAI	CT	Liver	in-house	volumes	941 (752/189)	0	20	Liver Tumor Segmentation	Segmentation
Liu et al. (2023a)	2023 MICCAI	MRI CT MRI CT CT	Heart Liver Hippocampus Pancreas Spleen	MSD	volume or patch	20 (1/64) 260 (208/52) 131 (105/26) 281 (225/56) 260 (208/52)	0	3 5	Left Atrial Segmentation Liver & Tumor Segmentation Anterior and Posterior Hippocampus Segmentation Pancreas & Tumor Segmentation Spleen Segmentation	Segmentation
Kadir et al. (2023)	2023 MICCAI	OCT	Eye	Duke AROI UMN	supersized	100 scans (6:2:2) 1136 scans (6:2:2) 725 scans (6:2:2)	2%	10%	OCT Layer Segmentation	Segmentation
Tang et al. (2023)	2023 MICCAI	Ultrasound	Carotid	CUBS in-house	image	3,220 (2,016/1,204) 350 (test)	159	200	Carotid Intima-Media Segmentation	Segmentation
Qu et al. (2023b)	2023 MICCAI	Histopathology	Colon	NCT-CRC-HE-100K	image	100,000	0	5%	Colorectal Cancer Diagnosis	Classification
Qu et al. (2023)	2023 MICCAI	Histopathology	Breast	CAMELYON16	patch	398 WSIs (270/128)	-	patch size of 4096, 8192, 12288 patch per WSI of 1, 3, 5	Detection of Cancer Metastases Lymph Nodes	Classification
Chen et al. (2023b)	2023 MIDL	Histopathology CT Microscopy	Colon Abdomen Peripheral Blood	PathMNIST OrganMNIST BloodMNIST	image	107,180 (89,996/10,047/1,180) 58,830 (44,561/6,491/17,778) 17,092 (11,959/1,712/3,421)	20	10	Survival Prediction Classification of Body Organs Cell Type Classification	Classification
Gaillhofer et al. (2023a)	2023 MIDL	MRI	Prostate	PROMISE 2012	slice	1,377 (1,020/109/248)	10	10	Prostate Segmentation	Segmentation
Qu et al. (2023a)	2023 NeurIPS	CT	Abdomen	AbdomenAtlas-8K	volume	8,448	-	-	Multi-organ Segmentation	Segmentation
Lüthi et al. (2024)	2023 NeurIPS	Dermoscopy	Skin	25,331 (15,200/3,799/6,332)	image	TCGA-KUMAR TNBC MoNuSeg	45,225, 900	45, 225, 900	Skin Lesion Diagnosis	Classification
Lou et al. (2023)	2023 TMI	Histopathology	Seven Organs Breast Seven Organs	MoNuSeg	patch	30 images (12/4/14) 50 images (30/7/13) 44 image (30/14)	-	5% 5% 7%	Nuclei Segmentation	Segmentation
Hu et al. (2023)	2023 TMI	Histopathology	Colon Lung & Colon	NCT-CRC-HE-100K LC25000	image	100,000 (4:1) 25,000 (4:1)	1%, 2% 2%	1%, 2% 2%	Colorectal Cancer Diagnosis Pulmonary & Colorectal Cancer Diagnosis	Classification
Li et al. (2024)	2024 IJCAIS	MRI CT	Lower Extremities	in-house	volume or slice	119 volumes (90/9/20) 30 volumes (25/1/4)	1 volume or 190 slices 1 volume or 540 slices	1 volume or 190 slices 1 volume or 540 slices	Musculoskeletal Segmentation	Segmentation
Mahapatra et al. (2024)	2024 Media	X-Ray Ultrasound Dermoscopy Fundus Eye Microscopy	Chest Breast Skin Eye Kidney	CheXpert ChestXray14 BreastMNIST BermanMNIST RetinaMNIST TissueMNIST	image	224,114 (223,414/200/500) 11,219 (test) 796 10,015 (7,007/1,003/2,005) 1,600 (1,080/1,200/400) 236,386 (165,466/23,640/47,280)	10	10	Detection of Pleural Effusion Thoracic Disease Diagnosis Breast Cancer Diagnosis Skin Lesion Diagnosis Diabetic Retinopathy Severity Grading Kidney Cortex Disease Classification	Classification

correctness and labeling difficulty. Experiments on the chest X-ray dataset showed that the proposed method improves performance by efficiently reducing label noise with fewer expert annotations compared to random selection.

However, multiple diseases and abnormalities often coexist simultaneously in diagnosing chest X-rays. Therefore, multi-label classification has been introduced, allowing each sample to be categorized into multiple classes (Baltruschat et al., 2019). Consequently, AL algorithms for chest X-ray classification must adapt to the multi-label setting. Built upon saliency maps, Mahapatra et al. (2022) further introduced GNN to model the inter-relationships between different labels. In this work, each class was treated as a node in a graph, with the relationships between classes represented as edges. They employed various techniques to aggregate information between different classes. As a follow-up work, Mahapatra et al. (2024) further introduced graph multiset transformers (Baek et al., 2020) for more powerful inter-label relationships than GNN.

5.1.2. Pathological Whole-slide Images and Multiple Instance Learning

Compared to modalities like X-ray, CT, and MRI, pathological whole-slide images (WSIs) provide microscopic details at the cellular level, making them critically important for tasks such as cancer staging and prognostic prediction. However, WSIs are very large, with maximum resolutions reaching $100,000 \times 100,000$ pixels. To handle these large images for deep learning, WSIs are usually divided into many small patches. Fully supervised methods require patch-level or even cell-level annotations, resulting in high annotation costs. AL can effectively improve annotation efficiency. For instance, in classifying breast pathological images, Qi et al. (2019) used entropy as the uncertainty metric. Uncertain patches were sent for annotation, whereas those with low entropy were given pseudo-labels to assist training. In AL of patch-level histological tissue classification, Hu et al. (2023) proposed category-wise curriculum querying to dynamically adjust the weight of uncertainty sampling of each class. They further proposed negative pre-training with wrong predictions to better distinguish the visually similar classes. To obtain fine-grained cellular annotation from WSI, van der Wal et al. (2021) proposed a human-augmenting AI-based labeling system with the help of AL. An active learner was used to select the next best patch for annotation and a classifier was trained for suggesting annotation. Specifically, Core-Set (Sener and Savarese, 2018) was used as the active learner. Experiments with pathologists demonstrate its ability to reduce workload by around 90% and slightly improve data quality across various cellular labeling tasks.

Nevertheless, pathologists might only provide WSI-level annotations in real-world clinical scenarios. Consequently, a prevailing direction in research is to formulate WSI classification as the weakly-supervised multi-instance learning (MIL) (Qu et al., 2022). In this framework, the entire WSI is viewed as a bag, and patches within each WSI are treated as instances within that bag. A well-trained MIL learner can

automatically identify relevant patches based on WSI-level labels, thus significantly reducing annotation costs. For example, a trained MIL classifier can automatically spot related patches by annotating whether or not cancer metastasis is present in a WSI. Nonetheless, task-relevant patches are often outnumbered by irrelevant ones, making MIL convergence more challenging. In MIL-based pathological WSI classification, AL filters out irrelevant patches and selects informative patches for annotation. Based on attention-based MIL, Sadafi et al. (2023) adopted MC Dropout to estimate both attention and classification uncertainties of each patch, then sent the most uncertain patches in each WSI for expert annotation. Qu et al. (2023b) found that in addition to patches related to the target (e.g., tumors, lymph nodes, and normal cells), WSIs contain many irrelevant patches (e.g., fat, stroma, and debris). Therefore, they adopted the open-set AL (Ning et al., 2022), in which the unlabeled pool contained both target and non-target class samples. They combined feature distributions with prediction uncertainty to select informative and relevant patches of the target class for annotation.

5.2. Active Learning for Medical Image Segmentation

Segmentation is one of the most common tasks in medical image analysis, capable of precisely locating anatomical structures or pathological lesions. However, training a segmentation model requires pixel-level annotation, which is time-consuming and labor-intensive for doctors. Therefore, active learning has been widely used in medical image segmentation and has become an important method to reduce annotation costs. Based on the unique traits of medical imaging, this section will focus on specialized designs in AL for medical image segmentation, including slice-based annotation, one-shot annotation, and annotation cost.

5.2.1. Slice-based Annotation

In 3D modalities like CT and MRI, adjacent 2D slices often exhibit significant semantic redundancy. Consequently, annotating only the key slices of each sample can reduce annotation costs. AL works mentioned in this section select 2D slices within a 3D volume for annotation. Representativeness-based methods have been widely applied in this line of work. For instance, Zheng et al. (2020) utilized autoencoders to learn the semantic features of each slice, then selected and annotated key slices from axial, sagittal, and coronal planes with a strategy similar to RA (Zheng et al., 2019). Specifically, they initially trained three 2D segmentation networks and one 3D segmentation network, where the inputs for the 2D networks are slices from different planes. These segmentation networks were used to generate four sets of pseudo-labels and subsequently to train the final 3D segmentation network. Results showed that this slice-based strategy outperforms uniform sampling. Building upon this method, Peng et al. (2022) adopted a similar strategy in 3D knee cartilage and bone segmentation. Besides, Wu et al. (2022d) incorporated a self-attention module into the autoencoder to enhance slice-level feature learning.

Uncertainty methods have also been introduced for selecting key slices. Zhou et al. (2021b) introduced a quality assessment module to select slices with the highest predicted average IoU score. In muscle segmentation of CT images, Hiasa et al. (2020) selected key slices and key regions. This work adopted clustering to select key slices and further selected regions with high uncertainty within each key slice for annotation.

In recent years, hybrid strategies combining both uncertainty and representativeness were proposed for slice-based annotation. In shoulder MRI musculoskeletal segmentation, Ozdemir et al. (2021) adopted the variance of multiple MC dropout runs as the uncertainty metric. The posterior probability estimated by infoVAE (Zhao et al., 2017) was used as the representativeness metric. Li et al. (2024) proposed a hybrid strategy to select informative slices in musculoskeletal segmentation of lower extremities, where uncertainty was estimated with a Bayesian U-net while the representativeness was based on cosine similarity. They further adopted mutual information to minimize the sample redundancy following Nath et al. (2021). The proposed method achieved impressive performances on both the MRI and CT datasets.

5.2.2. One-shot Annotation

Currently, most AL works require multiple rounds of annotation. However, this setting could be impractical in medical image segmentation. Multi-round annotation requires physicians to be readily available for each round of labeling, which is unrealistic in practice. If physicians cannot complete the annotations on time, the AL process must be suspended. In contrast, one-shot annotation eliminates the need for multiple interactions with physicians. It also allows for selecting valuable samples in a single round, thus reducing time costs. Both one-shot annotation and cold-start AL aim to select the most optimal initial annotations. However, the former allows for a higher annotation budget and strictly limits the number of interactions with experts to just one. Most relevant works combine self-supervised features and specific sampling strategies to achieve one-shot annotation. For example, RA (Zheng et al., 2019) is one of the earliest works in one-shot AL for medical image segmentation. They applied the VAE feature and a representativeness strategy to select informative samples for annotation in one shot. RA performed excellently in gland segmentation of pathological images, whole-heart MRI images, and fungal of electron microscopic images. Wu et al. (2022d) proposed a representativeness-based framework for selecting key slices for annotation in one shot. They adopted self-learning to learn the semantic representation of each slice and used it to propagate the expert annotation to different slices. Jin et al. (2022b) combined features of contrastive learning with farthest-first sampling to achieve one-shot annotation. The proposed method demonstrated effectiveness on the ISIC 2018 and lung segmentation datasets. Additionally, Jin et al. (2023b) utilized auto-encoding transformations for self-supervised feature learning. They selected and annotated samples with high density based on reachable distance.

5.2.3. Annotation Cost

Current AL works often assume equal annotation costs for each sample. Yet, this is not the case in medical image segmentation, where the time to annotate different samples can differ greatly. AL techniques can better support physicians by considering annotation costs (e.g., annotation time). In detecting intracranial hemorrhage of CT scans, Kuo et al. (2018) combined predictive disagreement with annotation time to select samples for annotation. Specifically, they adopted the Jensen-Shannon divergence to measure the disagreement between the outputs of multiple models. Annotation time for each sample was estimated by the length of the segmentation boundary and the number of connected components. In this work, AL was framed as a 0-1 knapsack problem, and dynamic programming is used to solve this problem for selecting informative samples. In brain tumor segmentation, Shen et al. (2021) derived the annotation cost of a slice based on the distance between the queried slices and the already-labeled slices. Specifically, lower distance represented lower annotation cost. The rationale is that the annotation cost of labeling a similar slice would be cheaper than that of the unfamiliar slices. In brain structure segmentation, Atzeni et al. (2022) further considered the spatial relationships between multiple regions of interest to more accurately estimate the annotation cost. Moreover, the average Dice coefficient of previous rounds was used to predict the average Dice for current segmentation results. They selected and annotated regions that can maximize the average Dice.

5.2.4. Interactive Segmentation

Despite the success of automatic segmentation in medical imaging, there is still a potential for errors in clinical applications due to domain shifts or unseen ROIs. Interactive segmentation (Budd et al., 2021; Luo et al., 2021) could produce real-time adjustment of the current segmentations based on the user inputs of clicks, bounding boxes, or scribbles. As a result, interactive segmentation could rapidly tune the model towards current clinical applications with the guidance of doctors. For the sake of flexibility, current interactive segmentation methods accept annotations for any position. However, such a paradigm would be more efficient when the model itself could suggest where to annotate, which is exactly what active learning is good at. Therefore, combining AL and interactive segmentation would further reduce the annotation cost. In this section, all the mentioned papers worked interactively with different labeling units. Before the DL era, Su et al. (2015) had already integrated AL in the interactive cell segmentation. They selected the most informative superpixels for interactive annotation with expected prediction error. In MRI fetal brain segmentation, Wang et al. (2020a) proposed an uncertainty-guided framework for interactive refinement. They developed a novel network architecture to produce multiple segmentation results simultaneously, and the variance between different predictions served as the uncertainty metric. Slices with the highest uncertainty were fetched for interactive refinement by human experts. In interactive segmentation of 3D medical images, (Zhou et al., 2022) proposed a quality

predictor, which produced a predicted IoU score with the current segmentation for each slice. With the interactive segmentation network, the quality predictor suggested slices with lower scores for expert annotation which could be in the form of scribble, bounding box, or extreme clicking. In [Li et al. \(2023\)](#), the most informative foreground and background superpixels were selected for interactive annotation.

5.3. Active Learning for Medical Image Reconstruction

AL can also be applied in medical image reconstruction. AL methods can help minimize the observations needed for modalities that require a long imaging time. This accelerates the imaging process and shortens the waiting period for patients. In this section, we'll explore the application of AL in the reconstruction of MRI, CT, and electron microscopy. Please refer to [Table 4](#) for more detail.

Deep learning has been applied to accelerate MRI acquisition and reconstruction. A common practice is to reduce k-space sampling through a fixed mask and use a deep model to reconstruct the undersampled MRI ([Qin et al., 2018](#)). To further improve the imaging speed, learnable sampling in AL can be applied to select the next measurement locations in k-space. For example, [Zhang et al. \(2019\)](#) adopted adversarial learning to train an evaluator for selecting the next row in k-space. [Pineda et al. \(2020\)](#) utilized reinforcement learning to train a dual deep Q-network for active sampling in k-space. [Bakker et al. \(2020\)](#) adopted policy gradient in reinforcement learning to train a policy network for adaptive sampling in k-space. The reward for the policy network was based on the improvement in structural similarity before and after the acquisition. Additionally, [Bakker et al. \(2022\)](#) explored how to jointly optimize the reconstruction and acquisition networks.

In addition to MRI imaging, AL has been employed in CT reconstruction as illustrated by [Wang et al. \(2022a\)](#). They adaptively chose the scanning angles tailored to individual patients, leading to a reduction in both radiation exposure and scanning duration. In electron microscopy, [Mi et al. \(2020\)](#) initially enhanced low-resolution images to high-resolution and then predicted the location of region-of-interest and reconstruction error. A weighted DPP based on reconstruction error was applied to select pixels that needed to be rescanned. Results showed that weighted DPP maintained both low reconstruction error and spatial diversity.

6. Performance Evaluation of Active Learning in Medical Image Analysis

In the field of medical image analysis, there is now an increasing amount of AL works. Despite its rapid development, AL for medical image analysis still faces several issues that limit its application in real-world clinical tasks. On one hand, there is a lack of comprehensive evaluation of AL methods on the medical imaging datasets. Most AL works conducted experiments on standard datasets, such as CIFAR-10, CIFAR-100, or MedMNIST. However, real-world medical imaging datasets often contain less available data and

higher complexity for analysis ([Varoquaux and Cheplygina, 2022](#)). Some AL works focused on a specific domain of medical imaging and achieved excellent performance, but their potential to generalize to a wider aspect of applications remains questionable. Beyond that, different AL methods exhibit inconsistent performance and may not necessarily outperform random sampling. In AL of classification tasks, [Munjal et al. \(2022\)](#) highlighted the absence of a consistently outperforming AL method and the fact that random sampling performs relatively well. In [§4.2](#), we have also mentioned that the AL methods are inferior to random sampling when the annotation budget is low. As a result, we are uncertain about which AL method works according to our requirements and whether this method could outperform the most straightforward baseline, random sampling.

To clarify the aforementioned issues, we have conducted a comprehensive evaluation of different AL methods on multiple medical imaging datasets. The adopted three datasets are widely used by the entire medical imaging community. They also correspond to different modalities, organs, and tasks (e.g., classification and segmentation). We choose the most representative and popular AL methods for their evaluations of the medical imaging datasets. Besides, we provided the details of dataset splits, network architecture, and training hyperparameters for better reproducibility. Codes are also available on our accompanying website ⁴.

6.1. Experimental Setups

6.1.1. Datasets

In this survey, we chose three medical imaging datasets for the performance evaluation of the AL methods, including two classification datasets and one segmentation dataset. The descriptions and dataset split are presented as follows, where a summarized table of the dataset split is in [Table 5](#).

NCT-CRC-HE-100K ([Kather et al., 2019](#)): This dataset contains 100,000 patches from 86 hematoxylin & eosin (H&E) stained histological slides of human colorectal cancer and normal tissue. All the patches are 224×224 at 0.5 microns per pixel. The patches are grouped into nine classes of different tissues, including adipose (ADI), background (BACK), debris (DEB), lymphocytes (LYM), mucus (MUC), smooth muscle (MUS), normal colon mucosa (NORM), cancer-associated stroma (STR), and colorectal adenocarcinoma epithelium (TUM). For the dataset split, we divided the dataset into training and validation sets with a ratio of 9:1 and utilized an additional dataset CRC-VAL-HE-7K which is provided by the same authors as the testing set. CRC-VAL-HE-7K shares the same acquisition protocol and tissue classes with NCT-CRC-HE-100K but contains 7,180 patches from 50 patients other than the patients of NCT-CRC-HE-100K.

ISIC 2020 ([Rotemberg et al., 2021](#)): ISIC 2020 is composed of 33,126 dermoscopic images from over 2,000 patients. Each image is labeled as benign or malignant by

⁴<https://github.com/LightersWang/Awesome-Active-Learning-for-Medical-Image-Analysis/tree/main/code>

Table 4: Summarization of surveyed works of active learning in medical image reconstruction.

	Year	Venues	Modality	ROIs	Dataset	Clinical Task
Jin et al. (2019)	2019	arXiv	MRI	Heart Knee	Cardiac Atlas Project fastMRI	MRI Reconstruction MRI Reconstruction
Zhang et al. (2019)	2019	CVPR	MRI	Knee	fastMRI	MRI Reconstruction
Mi et al. (2020)	2020	MICCAI	Electron Microscopy	Mouse Cortex Human Cerebrum	SNEMI3D in-house	Accelerated Acquisition of Electron Microscopy
Pineda et al. (2020)	2020	MICCAI	MRI	Knee	fastMRI	MRI Reconstruction
Bakker et al. (2020)	2020	NeurIPS	MRI	Knee Brain	fastMRI fastMRI	MRI Reconstruction MRI Reconstruction
Wang et al. (2022a)	2022	arXiv	CT	Lung Spine	AAPM VerSe	CT Reconstruction CT Reconstruction

Table 5: Training, validation, and testing splits of each dataset. Unless specified otherwise, the figures presented in this table represent the number of images of each split.

	NCT-CRC-HE-100K	ISIC 2020	ACDC
Training	90,000	20,869	656 (slices)
Validation	10,000	2,319	10 (volumes)
Testing	7,180 (CRC-VAL-HE-7K)	9,938	20 (volumes)

either doctors, long-term follow-up, or histopathology. ISIC 2020 contains 32,542 images of benign lesions and only 584 images of malignant lesions. We split this dataset of training, validation, and testing sets with a ratio of 6:1:3.

ACDC (Bernard et al., 2018): This dataset contains short-axis cardiac cine-MR images from 100 patients. In this survey, we only adopted the end-diastolic frame of each patient for evaluating different AL methods, which resulted in a total of 100 scans. Each scan corresponds to a human-annotated segmentation mask of the left ventricle (LV), myocardium (MYO), and right ventricle (RV). We followed the split from Luo et al. (2022), which contains 70, 10, and 20 scans in the training, validation, and testing sets, respectively. Due to the large spacing along the z-axis, 2D segmentation is more appropriate compared to 3D segmentation. Therefore, we trained the segmentation model with 2D slices and evaluated it with 3D volumes following Bai et al. (2017). Therefore, the training set is composed of 656 slices.

6.1.2. Evaluation Metrics

We employed different evaluation metrics for the task of each dataset. For the multi-class classification of NCT-CRC-HE-100K, we adopted the accuracy (ACC) to evaluate the classification performance. Due to the heavy class-imbalance of the binary classification task of ISIC 2020, we adopted the area under the receiver operating characteristic curve (AUC) for evaluation. For the segmentation task of ACDC, two well-known metrics of Dice similarity coefficient (DSC) and average surface distance (ASD) are used. DCS ranges from 0% (non-overlap) to 100% (perfect segmentation), and the lower ASD indicates a better alignment between the segmentation prediction and ground truth. To evaluate the

overall performance, we presented the mean DSC and ASD of LV, MYO, and RV. Evaluation metrics on the testing set were reported as the final results.

6.1.3. Active Learning Settings

In this study, we performed $T = 5$ rounds of annotation. To investigate how the number of queried samples in each round (i.e., annotation budget) affects the performance of different AL methods, we set different levels of annotation budgets b for each dataset. Following Lüth et al. (2024), high ($b = 1000$) and low ($b = 50$) budgets were used for classification tasks of NCT-CRC-HE-100K and ISIC 2020. Following Gaillochet et al. (2023b), we adopted a budget of 10 slices ($b = 10$) for ACDC segmentation. Considering the sizes of the involved datasets, the low budgets for classification ($b = 50$) and segmentation ($b = 10$) provide a chance to look into the performance of different AL methods in the low data regime. Before the active learning process started, we randomly selected the initial labeled pool to train an initial model. The size of the initial pool is equal to the annotation budget. The model training and sample selection are seeded. We ran each active learning method of a specific budget and dataset five times with different random seeds and reported the mean and standard deviation as the result.

6.1.4. Comparison Methods

For fairness and reproducibility, we conducted the evaluation using the following methods: **Random**: the baseline of active learning, which randomly draws the unlabeled samples. **Confidence, Entropy, and Margin Lewis and Catlett (1994)**; **Joshi et al. (2009)**; **Roth and Small (2006)**: these methods are all classic uncertainty-based AL methods, which calculated the confidence, entropy, and margin with the prediction probability as the uncertainty scores. Lower confidence, higher entropy, and lower margin indicate higher uncertainty. **DBAL Gal et al. (2017)**: This method integrated entropy and MC dropout for better uncertainty estimation. During sample selection, the model runs multiple times with all dropout layers activated, and the average probability of all MC dropout runs is used for calculating entropy. **BALD Gal et al. (2017)**: This method calculated BALD as the uncertainty score, which aims to maximize the mutual information between predictions and model parameters. MC dropout was also used in this method. **Core-Set Sener and Savarese (2018)**:

Table 6: Accuracy of different active learning methods for multi-class pathological tissue classification. We reported the test performance of the initially labeled dataset and other active learning rounds with mean and standard deviation. The best and second-best results are bolded in red and blue, respectively.

(a) Low budget ($b = 50$)							(b) High budget ($b = 1000$)						
ACC (%)	50	100	150	200	250	300	ACC (%)	1000	2000	3000	4000	5000	6000
Random		64.39 ± 0.34	66.80 ± 3.66	71.77 ± 3.35	73.55 ± 2.38	75.87 ± 1.21	Random		87.16 ± 1.43	87.68 ± 1.15	89.00 ± 0.91	89.04 ± 2.41	89.09 ± 1.08
Confidence		63.61 ± 2.49	66.16 ± 3.12	66.29 ± 3.84	70.65 ± 2.85	71.69 ± 0.97	Confidence		87.93 ± 0.76	87.46 ± 0.96	88.45 ± 0.64	90.01 ± 1.22	87.93 ± 3.04
Entropy		62.65 ± 0.79	65.66 ± 1.31	68.78 ± 2.34	69.12 ± 2.99	68.86 ± 3.84	Entropy		89.39 ± 0.23	86.54 ± 2.33	87.70 ± 2.25	88.04 ± 2.89	89.06 ± 1.57
Margin		70.65 ± 2.20	71.33 ± 1.54	75.93 ± 3.00	74.64 ± 4.21	76.61 ± 2.90	Margin		88.60 ± 0.62	88.05 ± 1.83	88.13 ± 1.74	89.90 ± 0.63	86.89 ± 2.99
BALD	60.46	62.13 ± 1.09	64.17 ± 4.20	69.34 ± 4.53	66.94 ± 2.50	68.49 ± 2.75	BALD	83.38	87.01 ± 0.99	85.85 ± 1.96	90.11 ± 0.53	89.19 ± 2.46	88.24 ± 2.21
DBAL		61.73 ± 2.43	63.76 ± 5.33	67.8 ± 1.78	67.90 ± 3.06	71.48 ± 3.38	DBAL		89.55 ± 0.75	86.56 ± 1.86	86.92 ± 3.62	88.76 ± 1.97	88.44 ± 0.53
BADGE		68.52 ± 6.57	72.27 ± 3.45	74.69 ± 0.85	75.96 ± 2.65	75.74 ± 2.22	BADGE		87.31 ± 0.51	87.56 ± 0.90	88.66 ± 2.33	88.39 ± 0.97	87.17 ± 1.30
Core-Set-L2		62.73 ± 4.33	66.26 ± 2.44	64.72 ± 3.94	66.25 ± 2.98	64.66 ± 1.24	Core-Set-L2		86.82 ± 0.37	88.99 ± 1.74	89.92 ± 0.98	89.34 ± 1.87	88.64 ± 3.10
Core-Set-Cosine		63.07 ± 1.84	63.12 ± 4.20	69.32 ± 1.37	68.25 ± 2.30	68.21 ± 3.48	Core-Set-Cosine		87.40 ± 0.49	88.09 ± 2.35	89.51 ± 1.20	89.80 ± 3.01	87.19 ± 2.78
Fully supervised				93.36			Fully supervised				93.36		

This method performed cover-based sampling using the feature embedding of each sample. For the balance of computation time and performance, we used the k-Center-Greedy for sample selection. To investigate how the distance metric affected the AL performance, we proposed a variant of Core-Set named “Core-Set-Cosine” which replaced the original L2 distance with the cosine distance. The original Core-Set was referred to as “Core-Set-L2” to avoid confusion. **BADGE** Ash et al. (2020): This method applied gradient as uncertainty estimation and utilized k-Means++ to improve diversity. Specifically, the gradient of cross-entropy loss was used in the classification task while the segmentation task used the gradient of the sum of the Dice loss and cross-entropy loss.

It should be noted that uncertainty-based methods in segmentation are slightly different from those of the classification. Specifically, we first produced the pixel-wise scores and then utilized the averaged scores for sample selection in the segmentation task.

6.1.5. Implementation Details

Classification: For all the classification task, we used ResNet-18 (He et al., 2016) as the backbone and the loss function is cross-entropy. We trained the model using stochastic gradient descent with momentum for 100 epochs with a batch size of 128. The learning rate and momentum were set as 0.01 and 0.9, respectively. Also, the cosine learning rate decay was adopted for smoother convergence. Data augmentations of the input image are different in the two classification datasets. In NCT-CRC-HE-100K, we only used the random horizontal flip. For ISIC 2020, we followed the data augmentation from Zhuang et al. (2018) which includes random crop, flip, rotation, affine transforms, and color jittering.

Segmentation: We used a 5-level U-Net (Ronneberger et al., 2015) for segmentation. Each level of the encoder or decoder contains two blocks. Each block consists of a 2D convolution, dropout layer with a probability of 0.1, batch normalization, and leaky ReLU activation. The segmentation loss is the combination of cross-entropy loss and Dice loss. We trained the model using the Adam optimizer (Kingma and Ba, 2014) for 4,000 iterations with a batch size of 32. The learning rate is 0.001 while decaying along the training iterations with the polynomial schedule. Data augmentation includes random flip, rotation of 90 degrees, and rotation of arbitrary degrees.

Experiments were conducted on NVIDIA GeForce RTX

3090 and 4090 GPUs and the CUDA version is 11.3. Codes are implemented using Python (version 3.8.10) and the PyTorch framework (version 1.11.0).

6.2. Experimental Results and Performance Analysis

6.2.1. Active Learning Results for Pathological Tissue Classification

We first evaluated the active learning performance on the pathological tissue classification task. The results of the testing accuracy are shown in Table 6. Margin performed well in the low-budget scenario. The reason for that is this method exploits the information of the wrong predictions of a similar class, which is in line with the finding in Hu et al. (2023). BADGE performs well in the low-budget scenarios largely due to the k-Means++ clustering. However, its performance drops in the high-budget scenario, which may indicate that the gradient embeddings are less suitable in AL when there is a distribution shift between the training and testing sets. The results of this section call for a more in-depth investigation of the generalizability to distribution shift of the AL methods.

6.2.2. Active Learning Results for Skin Lesion Classification

We have also conducted thorough evaluations on the ISIC 2020 dataset, which corresponds to a binary classification problem with severe class-imbalance. The AUC of the test split is shown in Table 7. It should be noted that Confidence and Margin are equivalent in the binary setting, so we only report the results of the former. In the low-budget scenario, Core-Set and its variant achieved better performances compared to the uncertainty-based methods. It indicates that representativeness-based methods or methods with improved diversity are more favored than uncertainty-based ones when the budget is low and the task is extremely difficult. Among all the uncertainty-based methods, BADGE stands out in certain rounds for its clustering operations that enhance diversity. For the high budget, the performance of the Core-Set variants is still competitive. However, the performance of the uncertainty-based methods improved. Results here demonstrated how the annotation budget affects the performance of the uncertainty-based and representativeness-based methods, in which the former fits a higher budget while the latter fits a lower budget.

Table 7: AUC of different active learning methods for the binary skin lesion classification. We reported the test performance on ISIC 2020 of the initially labeled dataset and other active learning rounds with mean and standard deviation. The best and second-best results are bolded in red and blue, respectively.

(a) Low budget ($b = 50$)							(b) High budget ($b = 1000$)						
AUC	50	100	150	200	250	300	AUC	1000	2000	3000	4000	5000	6000
Random		0.5609 ± 0.0325	0.5408 ± 0.0228	0.5807 ± 0.0391	0.5679 ± 0.0632	0.6143 ± 0.0610	Random		0.6898 ± 0.0216	0.7258 ± 0.0249	0.7556 ± 0.0236	0.7551 ± 0.0157	0.7718 ± 0.0124
Confidence		0.5432 ± 0.0233	0.5713 ± 0.0472	0.5622 ± 0.0302	0.5584 ± 0.0171	0.5649 ± 0.0794	Confidence		0.7058 ± 0.0120	0.7525 ± 0.0069	0.7760 ± 0.0171	0.7957 ± 0.0184	0.7949 ± 0.0113
Entropy		0.5375 ± 0.0182	0.5432 ± 0.0209	0.5919 ± 0.0413	0.5537 ± 0.0383	0.5543 ± 0.0220	Entropy		0.7063 ± 0.0144	0.7478 ± 0.0090	0.7692 ± 0.0117	0.7838 ± 0.0062	0.7875 ± 0.0101
BALD		0.5865 ± 0.0235	0.5828 ± 0.0499	0.5818 ± 0.0800	0.5977 ± 0.0273	0.5682 ± 0.0313	BALD		0.6851 ± 0.0156	0.7489 ± 0.0144	0.7683 ± 0.0097	0.7947 ± 0.0098	0.7837 ± 0.0073
DBAL	0.574	0.5322 ± 0.0300	0.5344 ± 0.0226	0.5601 ± 0.0341	0.5879 ± 0.0274	0.5442 ± 0.0350	DBAL	0.644	0.7120 ± 0.0195	0.7613 ± 0.0053	0.7616 ± 0.0261	0.7896 ± 0.0120	0.7883 ± 0.0127
BADGE		0.5601 ± 0.0373	0.5648 ± 0.0333	0.6219 ± 0.0757	0.6246 ± 0.0894	0.6172 ± 0.0890	BADGE		0.7074 ± 0.0264	0.7456 ± 0.0103	0.7734 ± 0.0150	0.7983 ± 0.0100	0.7917 ± 0.0064
Core-Set-L2		0.5834 ± 0.0344	0.6761 ± 0.0382	0.5956 ± 0.0610	0.5934 ± 0.0464	0.6313 ± 0.0867	Core-Set-L2		0.7304 ± 0.3033	0.7511 ± 0.0156	0.7499 ± 0.0106	0.7588 ± 0.0056	0.7575 ± 0.0067
Core-Set-Cosine		0.5459 ± 0.0367	0.6361 ± 0.0621	0.6032 ± 0.0402	0.6127 ± 0.0586	0.6317 ± 0.0797	Core-Set-Cosine		0.7288 ± 0.0236	0.7593 ± 0.0177	0.7807 ± 0.0182	0.7925 ± 0.0135	0.7917 ± 0.0109
Fully supervised				0.8424			Fully supervised				0.8424		

Table 8: Mean DSC and ASD of different active learning methods for cardiac MRI segmentation. We reported the test metrics on ACDC of the initially labeled dataset and other active learning rounds with mean and standard deviation. The best and second-best results are bolded in red and blue, respectively.

(a) Mean DSC							(b) Mean ASD						
DSC	10	20	30	40	50	60	ASD (mm)	10	20	30	40	50	60
Random		0.7872 ± 0.0552	0.8285 ± 0.0379	0.8657 ± 0.0063	0.8656 ± 0.0089	0.8721 ± 0.0047	Random		4.71 ± 1.13	2.82 ± 0.37	3.13 ± 0.48	2.05 ± 0.70	2.21 ± 0.73
Confidence		0.6805 ± 0.0205	0.7499 ± 0.0336	0.8347 ± 0.0203	0.8696 ± 0.0057	0.8717 ± 0.0065	Confidence		6.63 ± 2.05	3.60 ± 1.39	2.86 ± 0.39	2.21 ± 0.98	2.22 ± 0.62
Entropy		0.6888 ± 0.0088	0.7712 ± 0.0393	0.8538 ± 0.0117	0.8639 ± 0.0062	0.8714 ± 0.0040	Entropy	11.36	7.80 ± 1.70	3.27 ± 1.22	3.00 ± 1.08	2.24 ± 1.23	2.67 ± 0.84
Margin	0.4966	0.6841 ± 0.0124	0.7360 ± 0.0338	0.8214 ± 0.0535	0.8640 ± 0.0096	0.8734 ± 0.0053	Margin		6.64 ± 3.17	4.65 ± 1.94	2.66 ± 0.76	3.08 ± 1.02	3.22 ± 1.41
BADGE		0.8192 ± 0.0144	0.8486 ± 0.0151	0.8632 ± 0.0060	0.8731 ± 0.0032	0.8698 ± 0.0080	BADGE		4.22 ± 1.10	4.14 ± 0.32	3.10 ± 0.70	2.62 ± 0.81	2.56 ± 0.97
Core-Set-L2		0.8112 ± 0.0210	0.8491 ± 0.0144	0.8570 ± 0.0074	0.8572 ± 0.0057	0.8660 ± 0.0073	Core-Set-L2		3.22 ± 0.69	2.84 ± 0.66	2.71 ± 0.76	3.40 ± 1.31	2.86 ± 0.76
Core-Set-Cosine		0.7344 ± 0.0234	0.8248 ± 0.0109	0.8458 ± 0.0225	0.8549 ± 0.0204	0.8678 ± 0.0089	Core-Set-Cosine		4.75 ± 1.57	2.74 ± 0.52	3.50 ± 1.36	2.33 ± 0.88	1.85 ± 0.31
Fully supervised				0.9045			Fully supervised				2.09		

6.2.3. Active Learning Results for MRI Cardiac Segmentation

For segmentation, we evaluated different AL methods on the ACDC dataset. We reported the mean DSC and ASD of the segmentation results in Table 8. BADGE achieved the best or second-best performance in mean DSC for multiple rounds. Core-Set performed well on both the mean DSC and ASD in the early rounds of the lower budget scenario. Both two methods improve sampling diversity to some extent. However, for the later rounds, the performance of the uncertainty-based methods and random sampling improves in mean DSC and ASD. This result aligns with the findings in the previous section.

6.2.4. Effectiveness on Different Distances between Images

Distance measurement plays an important role in AL which may significantly impact the performance of AL algorithms. In this section, we evaluated the performance of the two most popular distances in AL, which are L2 and cosine distances. These two distances are based on the feature embedding. Assume x stands for the sample itself and its corresponding feature embedding is $\mathbf{z} = [z_1, z_2, \dots, z_d]$, d is the feature dimension. Based on the feature embedding, the L2 distance between two images x^a and x^b is as follow:

$$L2(x^a, x^b) = L2(\mathbf{z}^a, \mathbf{z}^b) = \sqrt{\sum_{i=1}^d (z_i^a - z_i^b)^2} \quad (12)$$

while the cosine distance is:

$$\text{Cosine}(x^a, x^b) = 1 - \frac{\mathbf{z}^a \cdot \mathbf{z}^b}{\|\mathbf{z}^a\| \cdot \|\mathbf{z}^b\|} = 1 - \frac{\sum_{i=1}^d z_i^a z_i^b}{\sqrt{\sum_{i=1}^d z_i^a{}^2} \sqrt{\sum_{i=1}^d z_i^b{}^2}} \quad (13)$$

To conduct experiments, we replaced the L2 distance in Core-Set with the cosine distance.

Performance comparisons between Core-Set-L2 and Core-Set-Cosine are illustrated in Fig.6. Results on the

NCT-CRC-HE-100K dataset showed no significant difference between the L2 and cosine distances across all budget levels. In ISIC 2020, the L2 distance tends to be better in the early rounds, indicating its ability to rapidly start the model, while Core-Set-Cosine significantly outperformed the Core-Set-L2 when the budget was high. On the ACDC dataset, Core-Set-L2 also outperformed Core-Set-Cosine in the early rounds but their performance after selecting more samples is similar. These results indicate that distance metrics play an important role in the performance of AL methods used in medical image analysis, and they should be carefully chosen according to the target tasks and the budgets. Generally speaking, L2 distance is more suitable for the low-budget scenario while the cosine distance might be a better choice when the budget is high.

7. Challenges and Future Perspectives

Currently, annotation scarcity is a significant bottleneck hindering the development of medical image analysis. AL improves annotation efficiency by selectively querying the most informative samples for annotation. This survey reviews the recent developments in deep active learning, focusing on the evaluation of informativeness, sampling strategies, integration with other label-efficient techniques, and the application of AL in medical image analysis. In this section, we will discuss the existing challenges faced by AL in medical image analysis and its future perspectives.

7.1. Towards Active Learning with Better Uncertainty

In AL, uncertainty plays a pivotal role. However, it would be beneficial if the uncertainty more directly highlighted the model's mistakes. We can enhance the model's performance by querying samples with inaccurate predictions.

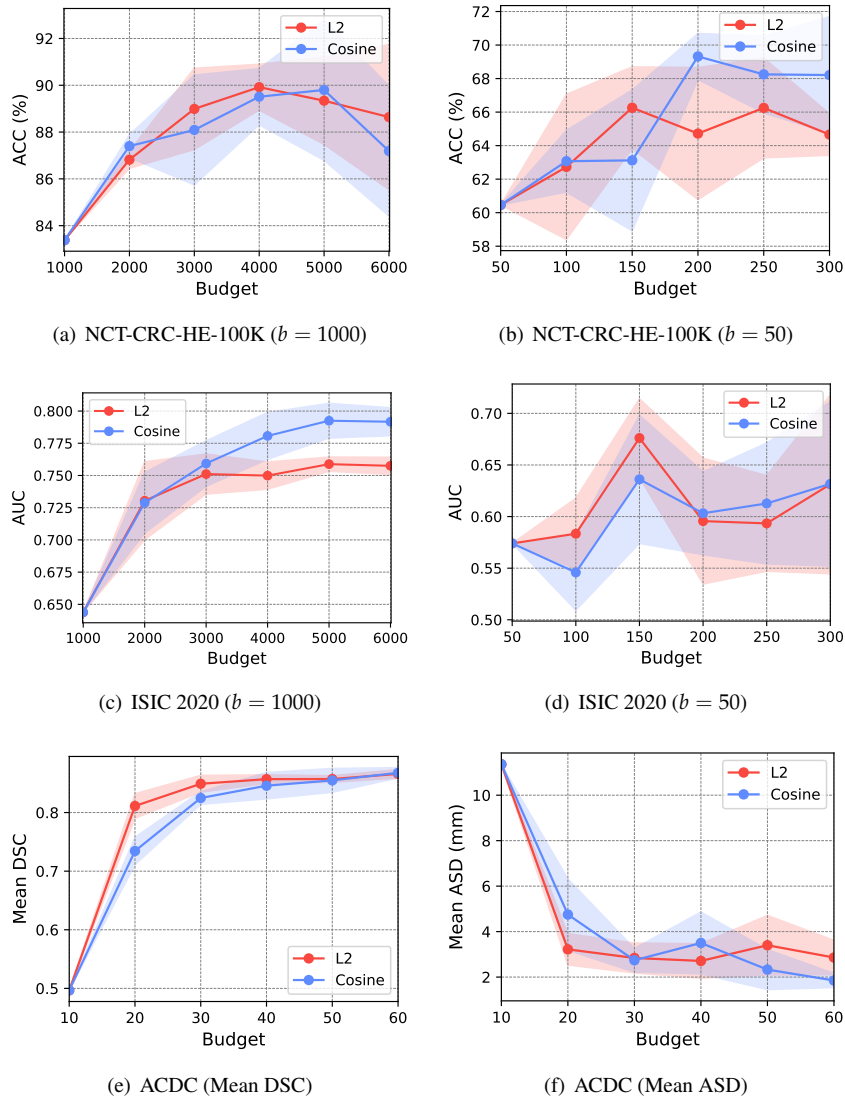


Figure 6: Performance comparisons between the L2 and cosine distances on all three datasets. We used Core-Set as the base AL method.

Recently, many works have adopted learnable performance estimation for quality control of deep model outputs. For instance, the recently proposed segment anything model (SAM) (Kirillov et al., 2023) provides IoU estimates for each mask to evaluate its quality. In medical image analysis, automated quality control is critical to ensure the reliability and safety of the deep model outputs (Kohlberger et al., 2012). For example, Wang et al. (2020d) employed deep generative models for learnable quality control in cardiac MRI segmentation, where the predicted Dice scores showed a strong linear relationship with the real ones. Additionally, Billot et al. (2023) used an additional neural network to predict the Dice coefficient of brain tissue segmentation results. Overall, learnable performance estimation can accurately predict the quality of model outputs. Hence, delving deeper into their potential for uncertainty-based AL is crucial to effectively tackle the issue of over-confidence.

Moreover, improving the probability calibration of model

prediction is a promising way to mitigate the over-confidence issue. Calibration (Guo et al., 2017; Mehtash et al., 2020) reflects the consistency between model prediction probabilities and the ground truth. A well-calibrated model should display a strong correlation between confidence and accuracy. For instance, if a perfect-calibrated polyp classifier gives an average confidence score of 0.9 on a dataset, it means that 90% of those samples should indeed have polyps. In reality, deep models generally suffer from the issue of over-confidence, which essentially means that they are not well-calibrated. Currently, only a few uncertainty-based AL works have considered probability calibration. For instance, Beluch et al. (2018) found that the model ensemble has better calibration than MC Dropout. Xie et al. (2022c) mitigated miscalibration by considering all possible prediction outcomes in the Dirichlet distribution. However, these methods are limited to proposing a better uncertainty metric and validating the calibration quality post-hoc. Existing calibration methods (Guo et al., 2017; Ding et al.,

2021) directly adjusted the distribution of prediction probabilities. However, these methods require an additional labeled dataset, thus limiting their practical applicability. Therefore, integrating probability calibration into uncertainty-based AL represents a valuable research direction worth exploring.

Among all the mentioned methods in §3.1, adversarial-based uncertainty currently has limited applications in AL of medical image analysis. Since the adversarial samples tend to be close to the classification boundary, they can be regarded as uncertain samples, and selecting them for training can potentially improve the trained model’s robustness. Exploring such ideas in medical image analysis, especially in the federated learning scenario, could be an interesting topic for future work.

7.2. Towards Active Learning with Better Representativeness

Representativeness-based AL effectively utilizes feature representations and data distributions for sample selection. Cover-based and discrepancy-based AL methods implicitly capture the data distribution, whereas density-based AL explicitly estimates it. However, the latter requires supplementary strategies to ensure diversity. For discrepancy-based AL, we can opt for a better metric of the distance between two probability distributions (Zhao et al., 2022). Besides, discrepancy-based AL has limited applications in medical image analysis currently. Finding a proper metric for medical images considering their special characteristics could be a promising direction for the future development of AL in the medical imaging domain.

As the core of density-based AL, density estimation in high-dimensional spaces has always been challenging. Popular density estimation methods, such as kernel density estimation and GMM, can encounter challenges when applied in high-dimensional spaces. In future research, we can consider introducing density estimators tailored to high-dimensional spaces. Advanced tools like normalizing flow (Papamakarios et al., 2021) could be an appropriate choice in density estimation in high-dimensional spaces.

7.3. Towards Active Learning with Weak Annotation

In §4.3, we discuss region-based active learning, which only requires region-level annotation of a sample. However, annotating all pixels within the region is still needed. Several existing works have incorporated weak annotations with AL to simplify the task for annotators. In object detection tasks, Vo et al. (2022) trained deep models with image-level annotation. They selected samples with box-in-box prediction results and annotated them with bounding boxes. Moreover, Lyu et al. (2023) adopted disagreement to choose which objects are worth annotating. Rather than annotating all objects within the image, they only required box-level annotations for a subset of objects. In AL of instance segmentation, Tang et al. (2022a) only required annotations for each object’s class label and bounding box, without the annotation of fine-grained segmentation masks. In future research, AL based on weak annotations is a direction worthy of in-depth exploration.

7.4. Towards Active Learning with Better Generative Models

In §4.4, we summarize the applications of generative models in AL. However, existing works have mainly focused on using GANs as sample generators. Recently, diffusion models (Kazerouni et al., 2023) have advanced in achieving state-of-the-art generative quality. Furthermore, text-to-image diffusion models, represented by Stable Diffusion (Rombach et al., 2022), have revolutionized the image generation domain. Their high-quality, text-guided generation results enable a more flexible image generation. With the use of ControlNet Zhang et al. (2023a), the diffusion models could learn to follow a more detailed condition like a sketch or segmentation mask. Exploring the potential of diffusion models in deep AL is a promising avenue for future research.

7.5. Towards Active Learning with Foundation Models

With the rise of visual foundational models, such as contrastive language-image pre-training (CLIP) (Radford et al., 2021) and SAM (Kirillov et al., 2023), and large language models (LLMs) like GPT-4 (OpenAI, 2023), deep learning in medical image analysis and computer vision is undergoing a paradigm shift. These foundational models (Bommasani et al., 2021) offer new opportunities for the development of AL.

AL is closely related to the training paradigms in deep learning of computer vision and medical image analysis. From the initial approach of train-from-scratch to the “pre-train-finetune” strategy using supervised or self-supervised pre-trained models, these paradigms usually require fine-tuning the entire network. Foundation models contain a wealth of knowledge. When combined with recently emerging parameter-efficient fine tuning (PEFT) or prompt tuning techniques (Hu et al., 2021; Jia et al., 2022), we can tune only a minimal subset of model weights (e.g., 5%) for rapid transfer to downstream tasks. As the number of fine-tuned parameters decreases, AL has the potential to further reduce the number of required annotated samples. Bai et al. (2023) integrated prompt tuning with AL in liver tumor segmentation. A segmentation model trained on publicly available datasets was transferred to in-house datasets via a novel prompt updater. With a mixed AL strategy of uncertainty and diversity, the proposed method reached the comparable performance of fully supervised tuning using around 5% of samples and 6% of tunable parameters. Therefore, it is essential to investigate the applicability of existing AL under PEFT or prompt tuning and explore the most suitable AL strategies for PEFT.

In natural language processing, LLMs have already taken a dominant role. Since most researchers cannot tune the LLMs, they rely on in-context learning, which provides LLMs with limited examples to transfer to downstream tasks. We believe that visual in-context learning will play a vital role in future research. Therefore, selecting the most suitable prompts for visual in-context learning will become an important research direction of AL.

8. Conclusion

Active learning is important to deep learning in medical image analysis since it effectively reduces the annotation costs incurred by human experts. This survey comprehensively reviews the core methods in deep active learning, its integration with different label-efficient techniques, and active learning works tailored to medical image analysis. We further discuss its current challenges and future perspectives. In summary, we believe that deep active learning and its application in medical image analysis hold important academic value and clinical potential, with ample room for further development.

Acknowledgments

This study was supported by the National Natural Science Foundation of China (Grant 82372097 and 82072021) and the Science and Technology Innovation Plan of Shanghai Science and Technology Commission (Grant 23S41900400).

References

- Achanta, R., Shaji, A., Smith, K., Lucchi, A., Fua, P., Süsstrunk, S., 2012. Slic superpixels compared to state-of-the-art superpixel methods. *IEEE transactions on pattern analysis and machine intelligence* 34, 2274–2282.
- Agarwal, S., Arora, H., Anand, S., Arora, C., 2020. Contextual diversity for active learning, in: *Computer Vision – ECCV 2020*. Springer International Publishing, Cham. volume 12361, pp. 137–153. doi:10.1007/978-3-030-58517-4_9.
- Aklilu, J., Yeung, S., 2022. Alges: active learning with gradient embeddings for semantic segmentation of laparoscopic surgical images, in: *Machine Learning for Healthcare Conference*, PMLR. pp. 892–911.
- Angluin, D., 1988. Queries and concept learning. *Machine Learning* 2, 319–342. doi:10.1023/A:1022821128753.
- Angluin, D., 2004. Queries revisited. *Theoretical Computer Science* 313, 175–194. doi:10.1016/j.tcs.2003.11.004.
- Antonelli, M., Reinke, A., Bakas, S., Farahani, K., Kopp-Schneider, A., Landman, B.A., Litjens, G., Menze, B., Ronneberger, O., Summers, R.M., et al., 2022. The medical segmentation decathlon. *Nature communications* 13, 4128.
- Ardila, D., Kiraly, A.P., Bharadwaj, S., Choi, B., Reicher, J.J., Peng, L., Tse, D., Etemadi, M., Ye, W., Corrado, G., et al., 2019. End-to-end lung cancer screening with three-dimensional deep learning on low-dose chest computed tomography. *Nature medicine* 25, 954–961.
- Ash, J., Goel, S., Krishnamurthy, A., Kakade, S., 2021. Gone fishing: Neural active learning with fisher embeddings, in: *Advances in Neural Information Processing Systems*, Curran Associates, Inc.. pp. 8927–8939.
- Ash, J.T., Zhang, C., Krishnamurthy, A., Langford, J., Agarwal, A., 2020. Deep batch active learning by diverse, uncertain gradient lower bounds, in: *International Conference on Learning Representations*.
- Atzeni, A., Peter, L., Robinson, E., Blackburn, E., Althonayan, J., Alexander, D.C., Iglesias, J.E., 2022. Deep active learning for suggestive segmentation of biomedical image stacks via optimisation of dice scores and traced boundary length. *Medical Image Analysis* 81, 102549. doi:10.1016/j.media.2022.102549.
- Azizi, S., Mustafa, B., Ryan, F., Beaver, Z., Freyberg, J., Deaton, J., Loh, A., Karthikesalingam, A., Kornblith, S., Chen, T., et al., 2021. Big self-supervised models advance medical image classification, in: *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 3478–3488.
- Baek, J., Kang, M., Hwang, S.J., 2020. Accurate learning of graph representations with graph multiset pooling, in: *International Conference on Learning Representations*.
- Bai, F., Xing, X., Shen, Y., Ma, H., Meng, M.Q.H., 2022. Discrepancy-based active learning for weakly supervised bleeding segmentation in wireless capsule endoscopy images, in: Wang, L., Dou, Q., Fletcher, P.T., Speidel, S., Li, S. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2022*, Springer Nature Switzerland, Cham. pp. 24–34. doi:10.1007/978-3-031-16452-1_3.
- Bai, F., Yan, K., Bai, X., Mao, X., Yin, X., Zhou, J., Shi, Y., Lu, L., Meng, M.Q.H., 2023. Slpt: Selective labeling meets prompt tuning on label-limited lesion segmentation, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer. pp. 14–24.
- Bai, W., Oktay, O., Sinclair, M., Suzuki, H., Rajchl, M., Tarroni, G., Glocker, B., King, A., Matthews, P.M., Rueckert, D., 2017. Semi-supervised learning for network-based cardiac mr image segmentation, in: *Medical Image Computing and Computer-Assisted Intervention- MICCAI 2017: 20th International Conference, Quebec City, QC, Canada, September 11-13, 2017, Proceedings, Part II 20*, Springer. pp. 253–260.
- Baid, U., Ghodasara, S., Mohan, S., Bilello, M., Calabrese, E., Colak, E., Farahani, K., Kalpathy-Cramer, J., Kitamura, F.C., Pati, S., et al., 2021. The rsna-asnr-miccai brats 2021 benchmark on brain tumor segmentation and radiogenomic classification. *arXiv preprint arXiv:2107.02314*.
- Bakker, T., Muckley, M., Romero-Soriano, A., Drozdal, M., Pineda, L., 2022. On learning adaptive acquisition policies for undersampled multi-coil mri reconstruction, in: *Proceedings of The 5th International Conference on Medical Imaging with Deep Learning*, PMLR. pp. 63–85.
- Bakker, T., van Hoof, H., Welling, M., 2020. Experimental design for mri by greedy policy search, in: *Advances in Neural Information Processing Systems*, Curran Associates, Inc.. pp. 18954–18966.
- Balam, S., Nguyen, C.M., Kassim, A., Krishnaswamy, P., 2022. Consistency-based semi-supervised evidential active learning for diagnostic radiograph classification, in: Wang, L., Dou, Q., Fletcher, P.T., Speidel, S., Li, S. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2022*, Springer Nature Switzerland, Cham. pp. 675–685. doi:10.1007/978-3-031-16431-6_64.
- Baltruschat, I.M., Nickisch, H., Grass, M., Knopp, T., Saalbach, A., 2019. Comparison of deep learning approaches for multi-label chest x-ray classification. *Scientific reports* 9, 6381.
- Beluch, W.H., Genewein, T., Nürnberger, A., Köhler, J.M., 2018. The power of ensembles for active learning in image classification, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 9368–9377.
- Bengar, J.Z., van de Weijer, J., Fuentes, L.L., Raducanu, B., 2022. Class-balanced active learning for image classification, in: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 1536–1545.
- Bengar, J.Z., van de Weijer, J., Twardowski, B., Raducanu, B., 2021. Reducing label effort: Self-supervised meets active learning, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1631–1639.
- Van den Bergh, M., Boix, X., Roig, G., De Capitani, B., Van Gool, L., 2012. Seeds: Superpixels extracted via energy-driven sampling, in: *Computer Vision–ECCV 2012: 12th European Conference on Computer Vision, Florence, Italy, October 7-13, 2012, Proceedings, Part VII 12*, Springer. pp. 13–26.
- Bernard, O., Lalonde, A., Zotti, C., Cervenansky, F., Yang, X., Heng, P.A., Cetin, I., Lekadir, K., Camara, O., Ballester, M.A.G., et al., 2018. Deep learning techniques for automatic mri cardiac multi-structures segmentation and diagnosis: is the problem solved? *IEEE transactions on medical imaging* 37, 2514–2525.
- Bernhardt, M., Castro, D.C., Tanno, R., Schwaighofer, A., Tezcan, K.C., Monteiro, M., Bannur, S., Lungren, M.P., Nori, A., Glocker, B., et al., 2022. Active label cleaning for improved dataset quality under resource constraints. *Nature communications* 13, 1161.
- Billot, B., Magdamo, C., Cheng, Y., Arnold, S.E., Das, S., Iglesias, J.E., 2023. Robust machine learning segmentation for large-scale analysis of heterogeneous clinical brain mri datasets. *Proceedings of the National Academy of Sciences* 120, e2216399120.
- Bishop, C.M., 1994. *Mixture density networks*.
- Bryk, E., Wang, K., Anari, N., Sadigh, D., 2019. Batch active learning using determinantal point processes. *arXiv preprint arXiv:1906.07975*.
- Bommasani, R., Hudson, D.A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M.S., Bohg, J., Bosselut, A., Brunskill, E., et al., 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.
- Budd, S., Robinson, E.C., Kainz, B., 2021. A survey on active learning and human-in-the-loop deep learning for medical image analysis. *Medical Image*

- Analysis 71, 102062.
- Cai, L., Xu, X., Liew, J.H., Foo, C.S., 2021. Revisiting superpixels for active learning in semantic segmentation with realistic annotation costs, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10988–10997.
- Caramalau, R., Bhattarai, B., Kim, T.K., 2021. Sequential graph convolutional network for active learning, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 9583–9592.
- Cardoso, T.N.C., Silva, R.M., Canuto, S., Moro, M.M., Gonçalves, M.A., 2017. Ranked batch-mode active learning. *Information Sciences* 379, 313–337. doi:10.1016/j.ins.2016.10.037.
- Casanova, A., Pinheiro, P.O., Rostamzadeh, N., Pal, C.J., 2020. Reinforced active learning for image segmentation, in: International Conference on Learning Representations.
- Chaudhuri, K., Kakade, S.M., Netrapalli, P., Sanghavi, S., 2015. Convergence rates of active learning for maximum likelihood estimation, in: Advances in Neural Information Processing Systems, Curran Associates, Inc.
- Chen, J., Ma, B., Cui, H., Xia, Y., Cheng, K.T., 2023a. Think twice before selection: Federated evidential active learning for medical image analysis with domain shifts. arXiv preprint arXiv:2312.02567.
- Chen, J., Xie, Y., Wang, K., Zhang, C., Vannan, M.A., Wang, B., Qian, Z., 2021. Active image synthesis for efficient labeling. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 43, 3770–3781. doi:10.1109/TPAMI.2020.2993221.
- Chen, L., Bai, Y., Huang, S., Lu, Y., Wen, B., Yuille, A.L., Zhou, Z., 2023b. Making your first choice: To address cold start problem in vision active learning, in: Medical Imaging with Deep Learning.
- Chen, Y., Mancini, M., Zhu, X., Akata, Z., 2022a. Semi-supervised and unsupervised deep visual learning: A survey. *IEEE transactions on pattern analysis and machine intelligence*.
- Chen, Z., Zhang, J., Wang, P., Chen, J., Li, J., 2022b. When active learning meets implicit semantic data augmentation, in: Avidan, S., Brostow, G., Cissé, M., Farinella, G.M., Hassner, T. (Eds.), *Computer Vision – ECCV 2022*. Springer Nature Switzerland, Cham. volume 13685, pp. 56–72. doi:10.1007/978-3-031-19806-9_4.
- Choi, J., Elezi, I., Lee, H.J., Farabet, C., Alvarez, J.M., 2021a. Active learning for deep object detection via probabilistic modeling, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 10264–10273.
- Choi, J., Yi, K.M., Kim, J., Choo, J., Kim, B., Chang, J., Gwon, Y., Chang, H.J., 2021b. Vab-al: Incorporating class imbalance and difficulty with variational bayes for active learning, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 6749–6758.
- Chong, Q.Z., Knottenbelt, W.J., Bhatia, K.K., 2021. Evaluation of active learning techniques on medical image classification with unbalanced data distributions, in: First Workshop, DGM4MICCAI 2021, and First Workshop, DALI 2021, Held in Conjunction with MICCAI 2021, Strasbourg, France, October 1, 2021, Proceedings 1, Springer. pp. 235–242.
- Citovsky, G., DeSalvo, G., Gentile, C., Karydas, L., Rajagopalan, A., Rostamzadeh, A., Kumar, S., 2021. Batch active learning at scale, in: Advances in Neural Information Processing Systems, Curran Associates, Inc.. pp. 11933–11944.
- Cohn, D., Atlas, L., Ladner, R., 1994. Improving generalization with active learning. *Machine Learning* 15, 201–221. doi:10.1007/BF00993277.
- Dai, C., Wang, S., Mo, Y., Angelini, E., Guo, Y., Bai, W., 2022. Suggestive annotation of brain mr images with gradient-guided sampling. *Medical Image Analysis* 77, 102373. doi:10.1016/j.media.2022.102373.
- Dai, C., Wang, S., Mo, Y., Zhou, K., Angelini, E., Guo, Y., Bai, W., 2020. Suggestive annotation of brain tumour images with gradient-guided sampling, in: Martel, A.L., Abolmaesumi, P., Stoyanov, D., Mateus, D., Zuluaga, M.A., Zhou, S.K., Racoceanu, D., Joskowicz, L. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2020*. Springer International Publishing, Cham. volume 12264, pp. 156–165. doi:10.1007/978-3-030-59719-1_16.
- Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L., 2009. Imagenet: A large-scale hierarchical image database, in: 2009 IEEE conference on computer vision and pattern recognition, Ieee. pp. 248–255.
- Ding, Z., Han, X., Liu, P., Niethammer, M., 2021. Local temperature scaling for probability calibration, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 6889–6899.
- Du, P., Chen, H., Zhao, S., Chai, S., Chen, H., Li, C., 2022. Contrastive active learning under class distribution mismatch. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 1–13doi:10.1109/TPAMI.2022.3188807.
- Du, P., Zhao, S., Chen, H., Chai, S., Chen, H., Li, C., 2021. Contrastive coding for active learning under class distribution mismatch, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 8927–8936.
- Ducoffe, M., Precioso, F., 2018. Adversarial active learning for deep networks: a margin based approach. arXiv preprint arXiv:1802.09841.
- Esteva, A., Kuprel, B., Novoa, R.A., Ko, J., Swetter, S.M., Blau, H.M., Thrun, S., 2017. Dermatologist-level classification of skin cancer with deep neural networks. *nature* 542, 115–118.
- Farahani, R.Z., Hekmatfar, M., 2009. Facility location: concepts, models, algorithms and case studies. Springer Science & Business Media.
- Feige, U., 1998. A threshold of $\ln n$ for approximating set cover. *Journal of the ACM (JACM)* 45, 634–652.
- Fu, B., Cao, Z., Wang, J., Long, M., 2021. Transferable query selection for active domain adaptation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 7272–7281.
- Fujishige, S., 2005. Submodular functions and optimization. Elsevier.
- Gaillochet, M., Desrosiers, C., Lombaert, H., 2023a. Active learning for medical image segmentation with stochastic batches, in: Medical Imaging with Deep Learning, short paper track.
- Gaillochet, M., Desrosiers, C., Lombaert, H., 2023b. Active learning for medical image segmentation with stochastic batches. *Medical Image Analysis* 90, 102958. doi:https://doi.org/10.1016/j.media.2023.102958.
- Gal, Y., Ghahramani, Z., 2016. Dropout as a bayesian approximation: Representing model uncertainty in deep learning, in: international conference on machine learning, PMLR. pp. 1050–1059.
- Gal, Y., Islam, R., Ghahramani, Z., 2017. Deep bayesian active learning with image data, in: Proceedings of the 34th International Conference on Machine Learning, PMLR. pp. 1183–1192.
- Gao, M., Zhang, Z., Yu, G., Arik, S.Ö., Davis, L.S., Pfister, T., 2020. Consistency-based semi-supervised active learning: Towards minimizing labeling cost, in: Vedaldi, A., Bischof, H., Brox, T., Frahm, J.M. (Eds.), *Computer Vision – ECCV 2020*. Springer International Publishing, Cham. volume 12355, pp. 510–526. doi:10.1007/978-3-030-58607-2_30.
- Ghesu, F.C., Georgescu, B., Mansoor, A., Yoo, Y., Gibson, E., Vishwanath, R., Balachandran, A., Balter, J.M., Cao, Y., Singh, R., et al., 2021. Quantifying and leveraging predictive uncertainty for medical image assessment. *Medical Image Analysis* 68, 101855.
- Gidaris, S., Singh, P., Komodakis, N., 2018. Unsupervised representation learning by predicting image rotations, in: International Conference on Learning Representations.
- Gissin, D., Shalev-Shwartz, S., 2019. Discriminative active learning. arXiv preprint arXiv:1907.06347.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y., 2014a. Generative adversarial nets. *Advances in neural information processing systems* 27.
- Goodfellow, I.J., Shlens, J., Szegedy, C., 2014b. Explaining and harnessing adversarial examples. arXiv preprint arXiv:1412.6572.
- Gorriz, M., Carlier, A., Faure, E., Giró-i Nieto, X., 2017. Cost-effective active learning for melanoma segmentation. arXiv preprint arXiv:1711.09168.
- Gretton, A., Borgwardt, K.M., Rasch, M.J., Schölkopf, B., Smola, A., 2012. A kernel two-sample test. *The Journal of Machine Learning Research* 13, 723–773.
- Gu, Y., Shen, M., Yang, J., Yang, G.Z., 2018. Reliable label-efficient learning for biomedical image recognition. *IEEE Transactions on Biomedical Engineering* 66, 2423–2432.
- Gu, Y., Vyas, K., Yang, J., Yang, G.Z., 2017. Unsupervised feature learning for endomicroscopy image retrieval, in: Medical Image Computing and Computer Assisted Intervention- MICCAI 2017: 20th International Conference, Quebec City, QC, Canada, September 11-13, 2017, Proceedings, Part III 20, Springer. pp. 64–71.
- Guan, H., Liu, M., 2021. Domain adaptation for medical image analysis: a survey. *IEEE Transactions on Biomedical Engineering* 69, 1173–1185.
- Guo, C., Pleiss, G., Sun, Y., Weinberger, K.Q., 2017. On calibration of modern neural networks, in: International conference on machine learning, PMLR. pp. 1321–1330.
- Hacohen, G., Dekel, A., Weinshall, D., 2022. Active learning on a budget: Opposite strategies suit high and low budgets, in: International Conference

- on Machine Learning, PMLR. pp. 8175–8195.
- Han, K., Sheng, V.S., Song, Y., Liu, Y., Qiu, C., Ma, S., Liu, Z., 2024. Deep semi-supervised learning for medical image segmentation: A review. *Expert Systems with Applications*, 123052.
- Haußmann, M., Hamprecht, F., Kandemir, M., 2019. Deep active learning with adaptive acquisition, in: *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, pp. 2470–2476.
- He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R., 2022. Masked auto-encoders are scalable vision learners, in: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16000–16009.
- He, K., Fan, H., Wu, Y., Xie, S., Girshick, R., 2020. Momentum contrast for unsupervised visual representation learning, in: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9729–9738.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778.
- Heo, B., Lee, M., Yun, S., Choi, J.Y., 2019. Knowledge distillation with adversarial samples supporting decision boundary, in: *Proceedings of the AAAI conference on artificial intelligence*, pp. 3771–3778.
- Hiasa, Y., Otake, Y., Takao, M., Ogawa, T., Sugano, N., Sato, Y., 2020. Automated muscle segmentation from clinical ct using bayesian u-net for personalized musculoskeletal modeling. *IEEE Transactions on Medical Imaging* 39, 1030–1040. doi:10.1109/TMI.2019.2940555.
- Hochbaum, D.S., Shmoys, D.B., 1985. A best possible heuristic for the k-center problem. *Mathematics of operations research* 10, 180–184.
- Houlsby, N., Huszar, F., Ghahramani, Z., Lengyel, M., 2011. Bayesian active learning for classification and preference learning. arXiv preprint arXiv:1112.5745.
- Hu, E.J., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W., et al., 2021. Lora: Low-rank adaptation of large language models, in: *International Conference on Learning Representations*.
- Hu, W., Cheng, L., Huang, G., Yuan, X., Zhong, G., Pun, C.M., Zhou, J., Cai, M., 2023. Learning from incorrectness: Active learning with negative pre-training and curriculum querying for histological tissue classification. *IEEE Transactions on Medical Imaging*.
- Huang, D., Li, J., Chen, W., Huang, J., Chai, Z., Li, G., 2023. Divide and adapt: Active domain adaptation via customized learning, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7651–7660.
- Huang, G., Li, Y., Pleiss, G., Liu, Z., Hopcroft, J.E., Weinberger, K.Q., 2017. Snapshot ensembles: Train 1, get m for free. arXiv preprint arXiv:1704.00109.
- Huang, S., Wang, T., Xiong, H., Huan, J., Dou, D., 2021. Semi-supervised active learning with temporal output discrepancy, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3447–3456.
- Huang, Y.J., Liu, W., Wang, X., Fang, Q., Wang, R., Wang, Y., Chen, H., Chen, H., Meng, D., Wang, L., 2020. Rectifying supporting regions with mixed and active supervision for rib fracture recognition. *IEEE Transactions on Medical Imaging* 39, 3843–3854. doi:10.1109/TMI.2020.3006138.
- Hwang, S., Lee, S., Kim, S., Ok, J., Kwak, S., 2022. Combating label distribution shift for active domain adaptation, in: *European Conference on Computer Vision*, Springer. pp. 549–566.
- Jaeger, S., Karargyris, A., Candemir, S., Folio, L., Siegelman, J., Callaghan, F., Xue, Z., Palaniappan, K., Singh, R.K., Antani, S., et al., 2013. Automatic tuberculosis screening using chest radiographs. *IEEE transactions on medical imaging* 33, 233–245.
- Jia, M., Tang, L., Chen, B.C., Cardie, C., Belongie, S., Hariharan, B., Lim, S.N., 2022. Visual prompt tuning, in: *European Conference on Computer Vision*, Springer. pp. 709–727.
- Jiménez, L.G., Dierckx, L., Amodei, M., Khosroshahi, H.R., Chidambaram, N., Ho, A.T.P., Franzin, A., 2023. Computational evaluation of the combination of semi-supervised and active learning for histopathology image segmentation with missing annotations, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pp. 2552–2563.
- Jin, C., Guo, Z., Lin, Y., Luo, L., Chen, H., 2023a. Label-efficient deep learning in medical image analysis: Challenges and future directions. arXiv preprint arXiv:2303.12484.
- Jin, K.H., Unser, M., Yi, K.M., 2019. Self-supervised deep active accelerated mri. arXiv preprint arXiv:1901.04547.
- Jin, Q., Li, S., Du, X., Yuan, M., Wang, M., Song, Z., 2023b. Density-based one-shot active learning for image segmentation. *Engineering Applications of Artificial Intelligence* 126, 106805. doi:10.1016/j.engappai.2023.106805.
- Jin, Q., Yuan, M., Li, S., Wang, H., Wang, M., Song, Z., 2022a. Cold-start active learning for image classification. *Information Sciences* 616, 16–36. doi:10.1016/j.ins.2022.10.066.
- Jin, Q., Yuan, M., Qiao, Q., Song, Z., 2022b. One-shot active learning for image segmentation via contrastive learning and diversity-based sampling. *Knowledge-Based Systems* 241, 108278. doi:10.1016/j.knosys.2022.108278.
- Jin, Q., Yuan, M., Wang, H., Wang, M., Song, Z., 2022c. Deep active learning models for imbalanced image classification. *Knowledge-Based Systems* 257, 109817. doi:10.1016/j.knosys.2022.109817.
- Joshi, A.J., Porikli, F., Papanikolopoulos, N., 2009. Multi-class active learning for image classification, in: *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2372–2379. doi:10.1109/CVPR.2009.5206627.
- Jung, S., Kim, S., Lee, J., 2023. A simple yet powerful deep active learning with snapshots ensembles, in: *International Conference on Learning Representations*.
- Kadir, M.A., Alam, H.M.T., Sonntag, D., 2023. Edgeal: An edge estimation based active learning approach for oct segmentation, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer. pp. 79–89.
- Kahl, K.C., Lüth, C.T., Zenk, M., Maier-Hein, K., Jaeger, P.F., 2024. Values: A framework for systematic validation of uncertainty estimation in semantic segmentation. arXiv preprint arXiv:2401.08501.
- Karamcheti, S., Krishna, R., Fei-Fei, L., Manning, C., 2021. Mind your outliers! investigating the negative impact of outliers on active learning for visual question answering, in: *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, Association for Computational Linguistics, Online. pp. 7265–7281. doi:10.18653/v1/2021.acl-long.564.
- Karimi, D., Dou, H., Warfield, S.K., Gholipour, A., 2020. Deep learning with noisy labels: Exploring techniques and remedies in medical image analysis. *Medical image analysis* 65, 101759.
- Kasarla, T., Nagendar, G., Hegde, G.M., Balasubramanian, V., Jawahar, C., 2019. Region-based active learning for efficient labeling in semantic segmentation, in: *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pp. 1109–1117. doi:10.1109/WACV.2019.00123.
- Kather, J.N., Krisam, J., Charoentong, P., Luedde, T., Herpel, E., Weis, C.A., Gaiser, T., Marx, A., Valous, N.A., Ferber, D., et al., 2019. Predicting survival from colorectal cancer histology slides using deep learning: A retrospective multicenter study. *PLoS medicine* 16, e1002730.
- Kazerouni, A., Aghdam, E.K., Heidari, M., Azad, R., Fayyaz, M., Hachililoglu, I., Merhof, D., 2023. Diffusion models in medical imaging: A comprehensive survey. *Medical Image Analysis*, 102846.
- Kendall, A., Gal, Y., 2017. What uncertainties do we need in bayesian deep learning for computer vision? *Advances in neural information processing systems* 30.
- Khanal, B., Bhattarai, B., Khanal, B., Stoyanov, D., Linte, C.A., 2023. M-vaal: Multimodal variational adversarial active learning for downstream medical image analysis tasks. arXiv preprint arXiv:2306.12376.
- Kim, H., Oh, M., Hwang, S., Kwak, S., Ok, J., 2023. Adaptive superpixel for active learning in semantic segmentation, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 943–953.
- Kim, K., Park, D., Kim, K.I., Chun, S.Y., 2021. Task-aware variational adversarial active learning, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8166–8175.
- Kingma, D.P., Ba, J., 2014. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980.
- Kingma, D.P., Welling, M., 2013. Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114.
- Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.Y., et al., 2023. Segment anything. arXiv preprint arXiv:2304.02643.
- Kirsch, A., van Amersfoort, J., Gal, Y., 2019. Batchbald: Efficient and diverse batch acquisition for deep bayesian active learning, in: *Advances in Neural Information Processing Systems*, Curran Associates, Inc.
- Koh, P.W., Liang, P., 2017. Understanding black-box predictions via influence functions, in: *International conference on machine learning*, PMLR.

- pp. 1885–1894.
- Kohl, S., Romera-Paredes, B., Meyer, C., De Fauw, J., Ledsam, J.R., Maier-Hein, K., Eslami, S., Jimenez Rezende, D., Ronneberger, O., 2018. A probabilistic u-net for segmentation of ambiguous images. *Advances in neural information processing systems* 31.
- Kohlberger, T., Singh, V., Alvino, C., Bahlmann, C., Grady, L., 2012. Evaluating segmentation error without ground truth, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer. pp. 528–536.
- Kothawade, S., Beck, N., Killamsetty, K., Iyer, R., 2021. Similar: Submodular information measures based active learning in realistic scenarios, in: *Advances in Neural Information Processing Systems*, Curran Associates, Inc.. pp. 18685–18697.
- Kothawade, S., Ghosh, S., Shekhar, S., Xiang, Y., Iyer, R., 2022a. Talisman: Targeted active learning for object detection with rare classes and slices using submodular mutual information, in: *Computer Vision – ECCV 2022*, Springer, Cham. pp. 1–16. doi:10.1007/978-3-031-19839-7_1.
- Kothawade, S., Kaushal, V., Ramakrishnan, G., Bilmes, J., Iyer, R., 2022b. Prism: A rich class of parameterized submodular information measures for guided data subset selection, in: *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 10238–10246.
- Kothawade, S., Savarkar, A., Iyer, V., Ramakrishnan, G., Iyer, R., 2022c. Clinical: Targeted active learning for imbalanced medical image classification, in: *Workshop on Medical Image Learning with Limited and Noisy Data*, Springer. pp. 119–129.
- Kovashka, A., Russakovsky, O., Fei-Fei, L., Grauman, K., et al., 2016. Crowdsourcing in computer vision. *Foundations and Trends® in computer graphics and Vision* 10, 177–243.
- Krishnan, R., Rajpurkar, P., Topol, E.J., 2022. Self-supervised learning in medicine and healthcare. *Nature Biomedical Engineering* 6, 1346–1352.
- Kuo, W., Häne, C., Yuh, E., Mukherjee, P., Malik, J., 2018. Cost-sensitive active learning for intracranial hemorrhage detection, in: Frangi, A.F., Schnabel, J.A., Davatzikos, C., Alberola-López, C., Fichtinger, G. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2018*. Springer International Publishing, Cham. volume 11072, pp. 715–723. doi:10.1007/978-3-030-00931-1_82.
- Lai, Z., Wang, C., Oliveira, L.C., Dugger, B.N., Cheung, S.C., Chuah, C.N., 2021. Joint semi-supervised and active learning for segmentation of gigapixel pathology images with cost-effective labeling, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pp. 591–600.
- LeCun, Y., Chopra, S., Hadsell, R., Ranzato, M., Huang, F., 2006. A tutorial on energy-based learning. *Predicting structured data* 1.
- Lee, D.H., et al., 2013. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks, in: *Workshop on challenges in representation learning*, ICML, Atlanta. p. 896.
- Lee, K., Zlateski, A., Ashwin, V., Seung, H.S., 2015. Recursive training of 2d-3d convolutional networks for neuronal boundary prediction. *Advances in Neural Information Processing Systems* 28.
- Lewis, D.D., Catlett, J., 1994. Heterogeneous uncertainty sampling for supervised learning, in: Cohen, W.W., Hirsh, H. (Eds.), *Machine Learning Proceedings 1994*. Morgan Kaufmann, San Francisco (CA), pp. 148–156. doi:10.1016/B978-1-55860-335-6.50026-x.
- Li, G., Otake, Y., Soufi, M., Taniguchi, M., Yagi, M., Ichihashi, N., Uemura, K., Takao, M., Sugano, N., Sato, Y., 2024. Hybrid representation-enhanced sampling for bayesian active learning in musculoskeletal segmentation of lower extremities. *International Journal of Computer Assisted Radiology and Surgery* , 1–10.
- Li, H., Yin, Z., 2020. Attention, suggestion and annotation: A deep active learning framework for biomedical image segmentation, in: Martel, A.L., Abolmaesumi, P., Stoyanov, D., Mateus, D., Zuluaga, M.A., Zhou, S.K., Racoceanu, D., Joskowicz, L. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2020*, Springer International Publishing, Cham. pp. 3–13. doi:10.1007/978-3-030-59710-8_1.
- Li, W., Li, J., Wang, Z., Polson, J., Sisk, A.E., Sajed, D.P., Speier, W., Arnold, C.W., 2022. Pathal: An active learning framework for histopathology image analysis. *IEEE Transactions on Medical Imaging* 41, 1176–1187. doi:10.1109/TMI.2021.3135002.
- Li, X., Xia, M., Jiao, J., Zhou, S., Chang, C., Wang, Y., Guo, Y., 2023. Hal-ia: A hybrid active learning framework using interactive annotation for medical image segmentation. *Medical Image Analysis* , 102862.
- Lin, Z., Wei, D., Jang, W.D., Zhou, S., Chen, X., Wang, X., Schalek, R., Berger, D., Matejek, B., Kamensky, L., Peleg, A., Haehn, D., Jones, T., Parag, T., Lichtman, J., Pfister, H., 2020. Two stream active query suggestion for active learning in connectomics, in: Vedaldi, A., Bischof, H., Brox, T., Frahm, J.M. (Eds.), *Computer Vision – ECCV 2020*, Springer International Publishing, Cham. pp. 103–120. doi:10.1007/978-3-030-58523-5_7.
- Linmans, J., Elfving, S., van der Laak, J., Litjens, G., 2023. Predictive uncertainty estimation for out-of-distribution detection in digital pathology. *Medical Image Analysis* 83, 102655.
- Liu, H., Li, H., Yao, X., Fan, Y., Hu, D., Dawant, B.M., Nath, V., Xu, Z., Oguz, I., 2023a. Colossal: A benchmark for cold-start active learning for 3d medical image segmentation, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer. pp. 25–34.
- Liu, J., Cao, L., Tian, Y., 2020. Deep active learning for effective pulmonary nodule detection, in: Martel, A.L., Abolmaesumi, P., Stoyanov, D., Mateus, D., Zuluaga, M.A., Zhou, S.K., Racoceanu, D., Joskowicz, L. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2020*, Springer International Publishing, Cham. pp. 609–618. doi:10.1007/978-3-030-59725-2_59.
- Liu, P., Wang, L., Ranjan, R., He, G., Zhao, L., 2022. A survey on active deep learning: from model driven to data driven. *ACM Computing Surveys (CSUR)* 54, 1–34.
- Liu, S., Yin, S., Qu, L., Wang, M., Song, Z., 2023b. A structure-aware framework of unsupervised cross-modality domain adaptation via frequency and spatial knowledge distillation. *IEEE Transactions on Medical Imaging* .
- Liu, Z., Ding, H., Zhong, H., Li, W., Dai, J., He, C., 2021. Influence selection for active learning, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9274–9283.
- Lou, W., Li, H., Li, G., Han, X., Wan, X., 2023. Which pixel to annotate: A label-efficient nuclei segmentation framework. *IEEE Transactions on Medical Imaging* 42, 947–958. doi:10.1109/TMI.2022.3221666.
- Luo, X., Hu, M., Song, T., Wang, G., Zhang, S., 2022. Semi-supervised medical image segmentation via cross teaching between cnn and transformer, in: *International Conference on Medical Imaging with Deep Learning*, PMLR. pp. 820–833.
- Luo, X., Wang, G., Song, T., Zhang, J., Aertsen, M., Deprest, J., Ourselin, S., Vercauteren, T., Zhang, S., 2021. Mideepseg: Minimally interactive segmentation of unseen objects from medical images using deep learning. *Medical image analysis* 72, 102102.
- Lüth, C., Bungert, T., Klein, L., Jaeger, P., 2024. Navigating the pitfalls of active learning evaluation: A systematic framework for meaningful performance assessment. *Advances in Neural Information Processing Systems* 36.
- Lyu, M., Zhou, J., Chen, H., Huang, Y., Yu, D., Li, Y., Guo, Y., Guo, Y., Xiang, L., Ding, G., 2023. Box-level active detection, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23766–23775.
- Mackowiak, R., Lenz, P., Ghorri, O., Diego, F., Lange, O., Rother, C., 2018. Cereals - cost-effective region-based active learning for semantic segmentation, in: *29th British Machine Vision Conference*.
- Mahapatra, D., Bozorgtabar, B., Ge, Z., Reyes, M., 2024. Gandalf: Graph-based transformer and data augmentation active learning framework with interpretable features for multi-label chest xray classification. *Medical image analysis* 93, 103075.
- Mahapatra, D., Bozorgtabar, B., Thiran, J.P., Reyes, M., 2018. Efficient active learning for image classification and segmentation using a sample selection and conditional generative adversarial network, in: Frangi, A.F., Schnabel, J.A., Davatzikos, C., Alberola-López, C., Fichtinger, G. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2018*. Springer International Publishing, Cham. volume 11071, pp. 580–588. doi:10.1007/978-3-030-00934-2_65.
- Mahapatra, D., Poellinger, A., Reyes, M., 2022. Graph node based interpretability guided sample selection for active learning. *IEEE Transactions on Medical Imaging* , 1–1doi:10.1109/TMI.2022.3215017.
- Mahapatra, D., Poellinger, A., Shao, L., Reyes, M., 2021. Interpretability-driven sample selection using self supervised learning for disease classification and segmentation. *IEEE Transactions on Medical Imaging* 40, 2548–2562. doi:10.1109/TMI.2021.3061724.
- Mahmood, R., Fidler, S., Law, M.T., 2022. Low-budget active learning via wasserstein distance: An integer programming approach, in: *International Conference on Learning Representations*.

- Mehrtash, A., Wells, W.M., Tempany, C.M., Abolmaesumi, P., Kapur, T., 2020. Confidence calibration and predictive uncertainty estimation for deep medical image segmentation. *IEEE transactions on medical imaging* 39, 3868–3878.
- Menze, B.H., Jakab, A., Bauer, S., Kalpathy-Cramer, J., Farahani, K., Kirby, J., Burren, Y., Porz, N., Slotboom, J., Wiest, R., et al., 2014. The multimodal brain tumor image segmentation benchmark (brats). *IEEE transactions on medical imaging* 34, 1993–2024.
- Mi, L., Wang, H., Meirovitch, Y., Schalek, R., Turaga, S.C., Lichtman, J.W., Samuel, A.D.T., Shavit, N., 2020. Learning guided electron microscopy with active acquisition, in: Martel, A.L., Abolmaesumi, P., Stoyanov, D., Mateus, D., Zuluaga, M.A., Zhou, S.K., Racoceanu, D., Joskowicz, L. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2020*, Springer International Publishing, Cham. pp. 77–87. doi:10.1007/978-3-030-59722-1_8.
- Miyato, T., Maeda, S.I., Koyama, M., Ishii, S., 2018. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and machine intelligence* 41, 1979–1993.
- Moosavi-Dezfooli, S.M., Fawzi, A., Frossard, P., 2016. Deepfool: a simple and accurate method to fool deep neural networks, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2574–2582.
- Munjal, P., Hayat, N., Hayat, M., Sourati, J., Khan, S., 2022. Towards robust and reproducible active learning using neural networks, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 223–232.
- Nath, V., Yang, D., Landman, B.A., Xu, D., Roth, H.R., 2021. Diminishing uncertainty within the training pool: Active learning for medical image segmentation. *IEEE Transactions on Medical Imaging* 40, 2534–2547. doi:10.1109/TMI.2020.3048055.
- Nath, V., Yang, D., Roth, H.R., Xu, D., 2022. Warm start active learning with proxy labels and selection via semi-supervised fine-tuning, in: Wang, L., Dou, Q., Fletcher, P.T., Speidel, S., Li, S. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2022*, Springer Nature Switzerland, Cham. pp. 297–308. doi:10.1007/978-3-031-16452-1_29.
- Nguyen, C., Huynh, M.T., Tran, M.Q., Nguyen, N.H., Jain, M., Ngo, V.D., Vo, T.D., Bui, T., Truong, S.Q.H., 2021. Goal: Gist-set online active learning for efficient chest x-ray image annotation, in: *Medical Imaging with Deep Learning*.
- Ning, K.P., Zhao, X., Li, Y., Huang, S.J., 2022. Active learning for open-set annotation, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 41–49.
- Ning, M., Lu, D., Wei, D., Bian, C., Yuan, C., Yu, S., Ma, K., Zheng, Y., 2021. Multi-anchor active domain adaptation for semantic segmentation, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9112–9122.
- OpenAI, 2023. Gpt-4 technical report. [arXiv:2303.08774](https://arxiv.org/abs/2303.08774).
- Otáloro, S., Perdomo, O., González, F., Müller, H., 2017. Training deep convolutional neural networks with active learning for exudate classification in eye fundus images, in: *6th Joint International Workshops, CVII-STENT 2017 and Second International Workshop, LABELS 2017, Held in Conjunction with MICCAI 2017, Québec City, QC, Canada, September 10–14, 2017, Proceedings 2*, Springer. pp. 146–154.
- Ozdemir, F., Peng, Z., Fuernstahl, P., Tanner, C., Goksel, O., 2021. Active learning for segmentation based on bayesian sample queries. *Knowledge-Based Systems* 214, 106531.
- Papamakarios, G., Nalisnick, E., Rezende, D.J., Mohamed, S., Lakshminarayanan, B., 2021. Normalizing flows for probabilistic modeling and inference. *The Journal of Machine Learning Research* 22, 2617–2680.
- Park, Y., Kim, S., Choi, W., Han, D.J., Moon, J., 2023. Active learning for object detection with evidential deep learning and hierarchical uncertainty aggregation, in: *International Conference on Learning Representations*.
- Parvaneh, A., Abbasnejad, E., Teney, D., Haffari, G.R., van den Hengel, A., Shi, J.Q., 2022. Active learning by feature mixing, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12237–12246.
- Payer, C., Štern, D., Bischof, H., Urschler, M., 2019. Integrating spatial configuration into heatmap regression based cnns for landmark localization. *Medical image analysis* 54, 207–219.
- Peng, Y., Zheng, H., Liang, P., Zhang, L., Zaman, F., Wu, X., Sonka, M., Chen, D.Z., 2022. Kcb-net: A 3d knee cartilage and bone segmentation network via sparse annotation. *Medical image analysis* 82, 102574.
- Pineda, L., Basu, S., Romero, A., Calandra, R., Drozdal, M., 2020. Active mr k-space sampling with reinforcement learning, in: Martel, A.L., Abolmaesumi, P., Stoyanov, D., Mateus, D., Zuluaga, M.A., Zhou, S.K., Racoceanu, D., Joskowicz, L. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2020*, Springer International Publishing, Cham. pp. 23–33. doi:10.1007/978-3-030-59713-9_3.
- Pourahmadi, K., Nooralinejad, P., Pirsiavash, H., 2021. A simple baseline for low-budget active learning. *arXiv preprint arXiv:2110.12033*.
- Prabhu, V., Chandrasekaran, A., Saenko, K., Hoffman, J., 2021. Active domain adaptation via clustering uncertainty-weighted embeddings, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 8505–8514.
- Qi, Q., Li, Y., Wang, J., Zheng, H., Huang, Y., Ding, X., Rohde, G.K., 2019. Label-efficient breast cancer histopathological image classification. *IEEE Journal of Biomedical and Health Informatics* 23, 2108–2116. doi:10.1109/JBHI.2018.2885134.
- Qin, C., Schlemper, J., Caballero, J., Price, A.N., Hajnal, J.V., Rueckert, D., 2018. Convolutional recurrent neural networks for dynamic mr image reconstruction. *IEEE transactions on medical imaging* 38, 280–290.
- Qiu, J., Wilm, F., Öttl, M., Schlereth, M., Liu, C., Heimann, T., Aubreville, M., Breininger, K., 2023. Adaptive region selection for active learning in whole slide image semantic segmentation, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer. pp. 90–100.
- Qu, C., Zhang, T., Qiao, H., Liu, J., Tang, Y., Yuille, A., Zhou, Z., 2023a. Annotating 8,000 abdominal ct volumes for multi-organ segmentation in three weeks. *arXiv preprint arXiv:2305.09666*.
- Qu, L., Liu, S., Liu, X., Wang, M., Song, Z., 2022. Towards label-efficient automatic diagnosis and analysis: a comprehensive survey of advanced deep learning-based weakly-supervised, semi-supervised and self-supervised techniques in histopathological image analysis. *Physics in Medicine & Biology*.
- Qu, L., Ma, Y., Yang, Z., Wang, M., Song, Z., 2023b. Openal: An efficient deep active learning framework for open-set pathology image classification, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer. pp. 3–13.
- Quan, Q., Yao, Q., Li, J., Zhou, S.K., 2022. Which images to label for few-shot medical landmark detection?, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 20606–20616.
- Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al., 2021. Learning transferable visual models from natural language supervision, in: *International conference on machine learning*, PMLR. pp. 8748–8763.
- Rädsch, T., Reinke, A., Weru, V., Tizabi, M.D., Schreck, N., Kavur, A.E., Pekdemir, B., Roß, T., Kopp-Schneider, A., Maier-Hein, L., 2023. Labelling instructions matter in biomedical image analysis. *Nature Machine Intelligence* 5, 273–283.
- Rahman, T., Khandakar, A., Qiblawey, Y., Tahir, A., Kiranyaz, S., Kashem, S.B.A., Islam, M.T., Al Maadeed, S., Zughair, S.M., Khan, M.S., et al., 2021. Exploring the effect of image enhancement techniques on covid-19 detection using chest x-ray images. *Computers in biology and medicine* 132, 104319.
- Rajpurkar, P., Chen, E., Banerjee, O., Topol, E.J., 2022. Ai in health and medicine. *Nature medicine* 28, 31–38.
- Rangwani, H., Jain, A., Aithal, S.K., Babu, R.V., 2021. S3vaada: Submodular subset selection for virtual adversarial active domain adaptation, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 7516–7525.
- Ren, P., Xiao, Y., Chang, X., Huang, P.Y., Li, Z., Gupta, B.B., Chen, X., Wang, X., 2021. A survey of deep active learning. *ACM computing surveys (CSUR)* 54, 1–40.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B., 2022. High-resolution image synthesis with latent diffusion models, in: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695.
- Ronneberger, O., Fischer, P., Brox, T., 2015. U-net: Convolutional networks for biomedical image segmentation, in: *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*,

Springer. pp. 234–241.

Rotemberg, V., Kurtansky, N., Betz-Stablein, B., Caffery, L., Chousakos, E., Codella, N., Combalia, M., Dusza, S., Guitera, P., Gutman, D., et al., 2021. A patient-centric dataset of images and metadata for identifying melanomas using clinical context. *Scientific data* 8, 34.

Roth, D., Small, K., 2006. Margin-based active learning for structured output spaces, in: Hutchison, D., Kanade, T., Kittler, J., Kleinberg, J.M., Matern, F., Mitchell, J.C., Naor, M., Nierstrasz, O., Pandu Rangan, C., Steffen, B., Sudan, M., Terzopoulos, D., Tygar, D., Vardi, M.Y., Weikum, G., Fürnkranz, J., Scheffer, T., Spiliopoulou, M. (Eds.), *Machine Learning: ECML 2006*. Springer Berlin Heidelberg, Berlin, Heidelberg. volume 4212, pp. 413–424. doi:10.1007/11871842_40.

Sadafi, A., Koehler, N., Makhro, A., Bogdanova, A., Navab, N., Marr, C., Peng, T., 2019. Multiclass deep active learning for detecting red blood cell subtypes in brightfield microscopy, in: Shen, D., Liu, T., Peters, T.M., Staib, L.H., Essert, C., Zhou, S., Yap, P.T., Khan, A. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2019*, Springer International Publishing, Cham. pp. 685–693. doi:10.1007/978-3-030-32239-7_76.

Sadafi, A., Navab, N., Marr, C., 2023. Active learning enhances classification of histopathology whole slide images with attention-based multiple instance learning, in: 2023 IEEE 20th International Symposium on Biomedical Imaging (ISBI), pp. 1–5. doi:10.1109/ISBI53787.2023.10230685.

Saquil, Y., Kim, K.I., Hall, P., 2018. Ranking cgans: Subjective control over semantic image attributes. arXiv preprint arXiv:1804.04082 .

Sener, O., Savarese, S., 2018. Active learning for convolutional neural networks: A core-set approach, in: International Conference on Learning Representations.

Sensoy, M., Kaplan, L., Kandemir, M., 2018. Evidential deep learning to quantify classification uncertainty. *Advances in neural information processing systems* 31.

Settles, B., 2009. Active learning literature survey .

Settles, B., Craven, M., Ray, S., 2007. Multiple-instance active learning. *Advances in neural information processing systems* 20.

Seung, H.S., Opper, M., Sompolinsky, H., 1992. Query by committee, in: Proceedings of the Fifth Annual Workshop on Computational Learning Theory, Association for Computing Machinery, New York, NY, USA. pp. 287–294. doi:10.1145/130385.130417.

Shaham, T.R., Dekel, T., Michaeli, T., 2019. Singan: Learning a generative model from a single natural image, in: Proceedings of the IEEE/CVF international conference on computer vision, pp. 4570–4580.

Shen, H., Tian, K., Dong, P., Zhang, J., Yan, K., Che, S., Yao, J., Luo, P., Han, X., 2020. Deep active learning for breast cancer segmentation on immunohistochemistry images, in: Martel, A.L., Abolmaesumi, P., Stoyanov, D., Mateus, D., Zuluaga, M.A., Zhou, S.K., Racoceanu, D., Joskowicz, L. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2020*, Springer International Publishing, Cham. pp. 509–518. doi:10.1007/978-3-030-59722-1_49.

Shen, M., Zhang, J.Y., Chen, L., Yan, W., Jani, N., Sutton, B., Koyejo, O., 2021. Labeling cost sensitive batch active learning for brain tumor segmentation, in: 2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI), IEEE. pp. 1269–1273.

Shi, X., Dou, Q., Xue, C., Qin, J., Chen, H., Heng, P.A., 2019. An active learning approach for reducing annotation cost in skin lesion analysis, in: *Machine Learning in Medical Imaging: 10th International Workshop, MLMI 2019, Held in Conjunction with MICCAI 2019, Shenzhen, China, October 13, 2019, Proceedings 10*, Springer. pp. 628–636.

Shin, I., Kim, D.J., Cho, J.W., Woo, S., Park, K., Kweon, I.S., 2021. Labor: Labeling only if required for domain adaptive semantic segmentation, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 8588–8598.

Shui, C., Zhou, F., Gagné, C., Wang, B., 2020. Deep active learning: Unified and principled method for query and training, in: Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics, PMLR. pp. 1308–1318.

Siddiqui, Y., Valentin, J., Niessner, M., 2020. Viewal: Active learning with viewpoint entropy for semantic segmentation, in: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, Seattle, WA, USA. pp. 9430–9440. doi:10.1109/CVPR42600.2020.00945.

Sim, Y., Chung, M.J., Kotter, E., Yune, S., Kim, M., Do, S., Han, K., Kim, H., Yang, S., Lee, D.J., et al., 2020. Deep convolutional neural network-based software improves radiologist detection of malignant lung nodules on chest radiographs. *Radiology* 294, 199–209.

Sinha, S., Ebrahimi, S., Darrell, T., 2019. Variational adversarial active learning, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 5972–5981.

Sourati, J., Akcakaya, M., Leen, T.K., Erdogmus, D., Dy, J.G., 2017. Asymptotic analysis of objectives based on fisher information in active learning. *The Journal of Machine Learning Research* 18, 1123–1163.

Sourati, J., Gholipour, A., Dy, J.G., Kurugol, S., Warfield, S.K., 2018. Active deep learning with fisher information for patch-wise semantic segmentation, in: Stoyanov, D., Taylor, Z., Carneiro, G., Syeda-Mahmood, T., Martel, A., Maier-Hein, L., Tavares, J.M.R., Bradley, A., Papa, J.P., Belagiannis, V., Nascimento, J.C., Lu, Z., Conjeti, S., Moradi, M., Greenspan, H., Madabhushi, A. (Eds.), *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support*. Springer International Publishing, Cham. volume 11045, pp. 83–91. doi:10.1007/978-3-030-00889-5_10.

Sourati, J., Gholipour, A., Dy, J.G., Tomas-Fernandez, X., Kurugol, S., Warfield, S.K., 2019. Intelligent labeling based on fisher information for medical image segmentation using deep learning. *IEEE Transactions on Medical Imaging* 38, 2642–2653. doi:10.1109/TMI.2019.2907805.

Su, H., Yin, Z., Huh, S., Kanade, T., Zhu, J., 2015. Interactive cell segmentation based on active and semi-supervised learning. *IEEE transactions on medical imaging* 35, 762–777.

Su, J.C., Tsai, Y.H., Sohn, K., Liu, B., Maji, S., Chandraker, M., 2020. Active adversarial domain adaptation, in: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pp. 739–748.

Tajbakhsh, N., Jeyaseelan, L., Li, Q., Chiang, J.N., Wu, Z., Ding, X., 2020. Embracing imperfect datasets: A review of deep learning solutions for medical image segmentation. *Medical Image Analysis* 63, 101693.

Takezoe, R., Liu, X., Mao, S., Chen, M.T., Feng, Z., Zhang, S., Wang, X., et al., 2023. Deep active learning for computer vision: Past and future. *APSIPA Transactions on Signal and Information Processing* 12.

Taleb, A., Loetzsch, W., Danz, N., Severin, J., Gaertner, T., Bergner, B., Lippert, C., 2020. 3d self-supervised methods for medical imaging. *Advances in neural information processing systems* 33, 18158–18172.

Tang, C., Xie, L., Zhang, G., Zhang, X., Tian, Q., Hu, X., 2022a. Active pointly-supervised instance segmentation, in: European Conference on Computer Vision, Springer. pp. 606–623.

Tang, Y., Hu, Y., Li, J., Lin, H., Xu, X., Huang, K., Lin, H., 2023. Pld-al: Pseudo-label divergence-based active learning in carotid intima-media segmentation for ultrasound images, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer. pp. 57–67.

Tang, Y., Yang, D., Li, W., Roth, H.R., Landman, B., Xu, D., Nath, V., Hatamizadeh, A., 2022b. Self-supervised pre-training of swin transformers for 3d medical image analysis, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 20730–20740.

Tarvainen, A., Valpola, H., 2017. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Advances in neural information processing systems* 30.

Tolkach, Y., Dohmgörge, T., Toma, M., Kristiansen, G., 2020. High-accuracy prostate cancer pathology using deep learning. *Nature Machine Intelligence* 2, 411–418.

Tran, T., Do, T.T., Reid, I., Carneiro, G., 2019. Bayesian generative active deep learning, in: Proceedings of the 36th International Conference on Machine Learning, PMLR. pp. 6295–6304.

Tschandl, P., Rinner, C., Apalla, Z., Argenziano, G., Codella, N., Halpern, A., Janda, M., Lallas, A., Longo, C., Malvehy, J., et al., 2020. Human–computer collaboration for skin cancer recognition. *Nature Medicine* 26, 1229–1234.

Unnikrishnan, B., Nguyen, C., Balaram, S., Li, C., Foo, C.S., Krishnaswamy, P., 2021. Semi-supervised classification of radiology images with noteacher: A teacher that is not mean. *Medical Image Analysis* 73, 102148.

Varoquaux, G., Cheplygina, V., 2022. Machine learning for medical imaging: methodological failures and recommendations for the future. *NPJ digital medicine* 5, 48.

Vo, H.V., Siméoni, O., Gidaris, S., Bursuc, A., Pérez, P., Ponce, J., 2022. Active learning strategies for weakly-supervised object detection, in: European Conference on Computer Vision, Springer. pp. 211–230.

van der Wal, D., Jhun, I., Lakloul, I., Nirschl, J., Richer, L., Rojansky, R., Theparee, T., Wheeler, J., Sander, J., Feng, F., et al., 2021. Biological data

- annotation via a human-augmenting ai-based labeling system. *NPJ digital medicine* 4, 145.
- Wan, F., Ye, Q., Yuan, T., Xu, S., Liu, J., Ji, X., Huang, Q., 2023. Multiple instance differentiation learning for active object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 1–15doi:10.1109/TPAMI.2023.3277738.
- Wang, C., Shang, K., Zhang, H., Zhao, S., Liang, D., Zhou, S.K., 2022a. Active ct reconstruction with a learned sampling policy. *arXiv preprint arXiv:2211.01670*.
- Wang, G., Aertsen, M., Deprest, J., Ourselin, S., Vercauteren, T., Zhang, S., 2020a. Uncertainty-guided efficient interactive refinement of fetal brain segmentation from stacks of mri slices, in: *Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part IV 23*, Springer. pp. 279–288.
- Wang, H., Chen, J., Zhang, S., He, Y., Xu, J., Wu, M., He, J., Liao, W., Luo, X., 2023. Dual-reference source-free active domain adaptation for nasopharyngeal carcinoma tumor segmentation across multiple hospitals. *arXiv preprint arXiv:2309.13401*.
- Wang, J., Yan, Y., Zhang, Y., Cao, G., Yang, M., Ng, M.K., 2020b. Deep reinforcement active learning for medical image classification, in: *Martel, A.L., Abolmaesumi, P., Stoyanov, D., Mateus, D., Zuluaga, M.A., Zhou, S.K., Racoceanu, D., Joskowicz, L. (Eds.), Medical Image Computing and Computer Assisted Intervention – MICCAI 2020, Springer International Publishing, Cham*. pp. 33–42. doi:10.1007/978-3-030-59710-8_4.
- Wang, K., Zhang, D., Li, Y., Zhang, R., Lin, L., 2017. Cost-effective active learning for deep image classification. *IEEE Transactions on Circuits and Systems for Video Technology* 27, 2591–2600. doi:10.1109/TCSVT.2016.2589879.
- Wang, P., Xiao, X., Glissen Brown, J.R., Berzin, T.M., Tu, M., Xiong, F., Hu, X., Liu, P., Song, Y., Zhang, D., et al., 2018. Development and validation of a deep-learning algorithm for the detection of polyps during colonoscopy. *Nature biomedical engineering* 2, 741–748.
- Wang, S., Li, Y., Ma, K., Ma, R., Guan, H., Zheng, Y., 2020c. Dual adversarial network for deep active learning, in: *Vedaldi, A., Bischof, H., Brox, T., Frahm, J.M. (Eds.), Computer Vision – ECCV 2020, Springer International Publishing, Cham*. pp. 680–696. doi:10.1007/978-3-030-58586-0_40.
- Wang, S., Tarroni, G., Qin, C., Mo, Y., Dai, C., Chen, C., Glocker, B., Guo, Y., Rueckert, D., Bai, W., 2020d. Deep generative model-based quality control for cardiac mri segmentation, in: *Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part IV 23*, Springer. pp. 88–97.
- Wang, T., Li, X., Yang, P., Hu, G., Zeng, X., Huang, S., Xu, C.Z., Xu, M., 2022b. Boosting active learning via improving test performance. *Proceedings of the AAAI Conference on Artificial Intelligence* 36, 8566–8574. doi:10.1609/aaai.v36i8.20834.
- Wang, X., Lian, L., Yu, S.X., 2022c. Unsupervised selective labeling for more effective semi-supervised learning, in: *European Conference on Computer Vision, Springer*. pp. 427–445.
- Wang, Y., Huang, G., Song, S., Pan, X., Xia, Y., Wu, C., 2021. Regularizing deep networks with semantic data augmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44, 3733–3748.
- Wang, Z., Yin, Z., 2021. Annotation-efficient cell counting, in: *de Bruijne, M., Cattin, P.C., Cotin, S., Padoy, N., Speidel, S., Zheng, Y., Essert, C. (Eds.), Medical Image Computing and Computer Assisted Intervention – MICCAI 2021, Springer International Publishing, Cham*. pp. 405–414. doi:10.1007/978-3-030-87237-3_39.
- Wei, K., Iyer, R., Bilmes, J., 2015. Submodularity in data subset selection and active learning, in: *Proceedings of the 32nd International Conference on Machine Learning, PMLR*. pp. 1954–1963.
- Weisberg, S., Cook, R.D., 1982. Residuals and influence in regression.
- Williams, R.J., 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning* 8, 229–256.
- Wu, J., Chen, J., Huang, D., 2022a. Entropy-based active learning for object detection with progressive diversity constraint, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9397–9406.
- Wu, T.H., Liou, Y.S., Yuan, S.J., Lee, H.Y., Chen, T.I., Huang, K.C., Hsu, W.H., 2022b. D2ada: Dynamic density-aware active domain adaptation for semantic segmentation, in: *European Conference on Computer Vision, Springer*. pp. 449–467.
- Wu, X., Chen, C., Zhong, M., Wang, J., Shi, J., 2021. Covid-al: The diagnosis of covid-19 with deep active learning. *Medical Image Analysis* 68, 101913. doi:10.1016/j.media.2020.101913.
- Wu, X., Pei, J., Chen, C., Zhu, Y., Wang, J., Qian, Q., Zhang, J., Sun, Q., Guo, Y., 2022c. Federated active learning for multicenter collaborative disease diagnosis. *IEEE Transactions on Medical Imaging*.
- Wu, Y., Zheng, B., Chen, J., Chen, D.Z., Wu, J., 2022d. Self-learning and one-shot learning based single-slice annotation for 3d medical image segmentation, in: *Wang, L., Dou, Q., Fletcher, P.T., Speidel, S., Li, S. (Eds.), Medical Image Computing and Computer Assisted Intervention – MICCAI 2022, Springer Nature Switzerland, Cham*. pp. 244–254. doi:10.1007/978-3-031-16452-1_24.
- Xie, B., Yuan, L., Li, S., Liu, C.H., Cheng, X., 2022a. Towards fewer annotations: Active learning via region impurity and prediction uncertainty for domain adaptive semantic segmentation, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8068–8078.
- Xie, B., Yuan, L., Li, S., Liu, C.H., Cheng, X., Wang, G., 2022b. Active learning for domain adaptation: An energy-based approach, in: *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 8708–8716. doi:10.1609/aaai.v36i8.20850.
- Xie, M., Li, S., Zhang, R., Liu, C.H., 2022c. Dirichlet-based uncertainty calibration for active domain adaptation, in: *The Eleventh International Conference on Learning Representations*.
- Xie, M., Li, Y., Wang, Y., Luo, Z., Gan, Z., Sun, Z., Chi, M., Wang, C., Wang, P., 2022d. Learning distinctive margin toward active domain adaptation, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7993–8002.
- Xie, Y., Ding, M., Tomizuka, M., Zhan, W., 2023a. Towards free data selection with general-purpose models, in: *Thirty-seventh Conference on Neural Information Processing Systems*. URL: <https://openreview.net/forum?id=KBXcDaaZE7>.
- Xie, Y., Lu, H., Yan, J., Yang, X., Tomizuka, M., Zhan, W., 2023b. Active fine-tuning: Exploiting annotation budget in the pretraining-finetuning paradigm, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23715–23724.
- Xu, X., Lu, Q., Yang, L., Hu, S., Chen, D., Hu, Y., Shi, Y., 2018. Quantization of fully convolutional networks for accurate biomedical image segmentation, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8300–8308.
- Xu, Y., Xu, X., Jin, L., Gao, S., Goh, R.S.M., Ting, D.S.W., Liu, Y., 2021. Partially-supervised learning for vessel segmentation in ocular images, in: *de Bruijne, M., Cattin, P.C., Cotin, S., Padoy, N., Speidel, S., Zheng, Y., Essert, C. (Eds.), Medical Image Computing and Computer Assisted Intervention – MICCAI 2021, Springer International Publishing, Cham*. pp. 271–281. doi:10.1007/978-3-030-87193-2_26.
- Yang, J., Shi, R., Wei, D., Liu, Z., Zhao, L., Ke, B., Pfister, H., Ni, B., 2023. Medmnist v2-a large-scale lightweight benchmark for 2d and 3d biomedical image classification. *Scientific Data* 10, 41.
- Yang, L., Zhang, Y., Chen, J., Zhang, S., Chen, D.Z., 2017. Suggestive annotation: A deep active learning framework for biomedical image segmentation, in: *Descoteaux, M., Maier-Hein, L., Franz, A., Jannin, P., Collins, D.L., Duchesne, S. (Eds.), Medical Image Computing and Computer Assisted Intervention - MICCAI 2017, Springer International Publishing, Cham*. pp. 399–407. doi:10.1007/978-3-319-66179-7_46.
- Ye, M., Giannarou, S., Meining, A., Yang, G.Z., 2016. Online tracking and retargeting with applications to optical biopsy in gastrointestinal endoscopic examinations. *Medical image analysis* 30, 144–157.
- Yehuda, O., Dekel, A., Hacoen, G., Weinshall, D., 2022. Active learning through a covering lens, in: *Advances in Neural Information Processing Systems*.
- Yi, J.S.K., Seo, M., Park, J., Choi, D.G., 2022. Pt4al: Using self-supervised pretext tasks for active learning, in: *Avidan, S., Brostow, G., Cissé, M., Farinella, G.M., Hassner, T. (Eds.), Computer Vision – ECCV 2022, Springer Nature Switzerland, Cham*. pp. 596–612. doi:10.1007/978-3-031-19809-0_34.
- Yoo, D., Kweon, I.S., 2019. Learning loss for active learning, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 93–102.
- Yuan, M., Lin, H.T., Boyd-Graber, J., 2020a. Cold-start active learning through self-supervised language modeling, in: *Proceedings of the 2020 Conference*

- on Empirical Methods in Natural Language Processing (EMNLP), Association for Computational Linguistics, Online. pp. 7935–7948. doi:[10.18653/v1/2020.emnlp-main.637](https://doi.org/10.18653/v1/2020.emnlp-main.637).
- Yuan, P., Mobiny, A., Jahaniipour, J., Li, X., Cicalese, P.A., Roysam, B., Patel, V.M., Dragan, M., Van Nguyen, H., 2020b. Few is enough: task-augmented active meta-learning for brain cell classification, in: Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part I 23, Springer. pp. 367–377.
- Yuan, T., Wan, F., Fu, M., Liu, J., Xu, S., Ji, X., Ye, Q., 2021. Multiple instance active learning for object detection, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5330–5339.
- Zhan, X., Wang, Q., Huang, K.h., Xiong, H., Dou, D., Chan, A.B., 2022. A comparative survey of deep active learning. arXiv preprint arXiv:2203.13450.
- Zhang, B., Li, L., Yang, S., Wang, S., Zha, Z.J., Huang, Q., 2020. State-relabeling adversarial active learning, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8756–8765.
- Zhang, L., Rao, A., Agrawala, M., 2023a. Adding conditional control to text-to-image diffusion models, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 3836–3847.
- Zhang, W., Zhu, L., Hallinan, J., Zhang, S., Makmur, A., Cai, Q., Ooi, B.C., 2022a. Boostmis: Boosting medical image semi-supervised learning with adaptive pseudo labeling and informative active annotation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 20666–20676.
- Zhang, Y., Kang, B., Hooi, B., Yan, S., Feng, J., 2023b. Deep long-tailed learning: A survey. IEEE Transactions on Pattern Analysis and Machine Intelligence.
- Zhang, Y., Zhang, X., Xie, L., Li, J., Qiu, R.C., Hu, H., Tian, Q., 2022b. One-bit active query with contrastive pairs, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 9697–9705.
- Zhang, Z., Romero, A., Muckley, M.J., Vincent, P., Yang, L., Drozdal, M., 2019. Reducing uncertainty in undersampled mri reconstruction with active acquisition, in: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, Long Beach, CA, USA. pp. 2049–2053. doi:[10.1109/CVPR.2019.00215](https://doi.org/10.1109/CVPR.2019.00215).
- Zhao, S., Sinha, A., He, Y., Perreault, A., Song, J., Ermon, S., 2022. Comparing distributions by measuring differences that affect decision making, in: International Conference on Learning Representations.
- Zhao, S., Song, J., Ermon, S., 2017. Infovae: Information maximizing variational autoencoders. arXiv preprint arXiv:1706.02262.
- Zhao, Z., Zeng, Z., Xu, K., Chen, C., Guan, C., 2021. Dsal: Deeply supervised active learning from strong and weak labelers for biomedical image segmentation. IEEE Journal of Biomedical and Health Informatics 25, 3744–3751. doi:[10.1109/JBHI.2021.3052320](https://doi.org/10.1109/JBHI.2021.3052320).
- Zheng, H., Yang, L., Chen, J., Han, J., Zhang, Y., Liang, P., Zhao, Z., Wang, C., Chen, D.Z., 2019. Biomedical image segmentation via representative annotation. Proceedings of the AAAI Conference on Artificial Intelligence 33, 5901–5908. doi:[10.1609/aaai.v33i01.33015901](https://doi.org/10.1609/aaai.v33i01.33015901).
- Zheng, H., Zhang, Y., Yang, L., Wang, C., Chen, D.Z., 2020. An annotation sparsification strategy for 3d medical image segmentation via representative selection and self-training. Proceedings of the AAAI Conference on Artificial Intelligence 34, 6925–6932. doi:[10.1609/aaai.v34i04.6175](https://doi.org/10.1609/aaai.v34i04.6175).
- Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., Torralba, A., 2016. Learning deep features for discriminative localization, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2921–2929.
- Zhou, S.K., Greenspan, H., Davatzikos, C., Duncan, J.S., Van Ginneken, B., Madabhushi, A., Prince, J.L., Rueckert, D., Summers, R.M., 2021a. A review of deep learning in medical imaging: Imaging traits, technology trends, case studies with progress highlights, and future promises. Proceedings of the IEEE 109, 820–838.
- Zhou, T., Li, L., Bredell, G., Li, J., Konukoglu, E., 2021b. Quality-aware memory network for interactive volumetric image segmentation, in: de Bruijne, M., Cattin, P.C., Cotin, S., Padoy, N., Speidel, S., Zheng, Y., Essert, C. (Eds.), Medical Image Computing and Computer Assisted Intervention – MICCAI 2021, Springer International Publishing, Cham. pp. 560–570. doi:[10.1007/978-3-030-87196-3_52](https://doi.org/10.1007/978-3-030-87196-3_52).
- Zhou, T., Li, L., Bredell, G., Li, J., Konukoglu, E., 2022. Volumetric memory network for interactive medical image segmentation. Medical Image Analysis, 102599doi:[10.1016/j.media.2022.102599](https://doi.org/10.1016/j.media.2022.102599).
- Zhou, Z., Shin, J., Zhang, L., Gurudu, S., Gotway, M., Liang, J., 2017. Fine-tuning convolutional neural networks for biomedical image analysis: Actively and incrementally, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 7340–7351.
- Zhou, Z., Shin, J.Y., Gurudu, S.R., Gotway, M.B., Liang, J., 2021c. Active, continual fine tuning of convolutional neural networks for reducing annotation efforts. Medical Image Analysis 71, 101997. doi:[10.1016/j.media.2021.101997](https://doi.org/10.1016/j.media.2021.101997).
- Zhu, J.J., Bento, J., 2017. Generative adversarial active learning. arXiv preprint arXiv:1702.07956.
- Zhuang, J., Li, W., Manivannan, S., Wang, R., Zhang, J.J.G., Pan, J., Jiang, G., Yin, Z., 2018. Skin lesion analysis towards melanoma detection using deep neural network ensemble. ISIC Challenge 2018 2.