

Domain Transfer Through Image-to-Image Translation for Uncertainty-Aware Prostate Cancer Classification

Meng Zhou^{a,b,*}, Amoon Jamzad^b, Jason Izzard^c, Alexandre Menard^c, Robert Siemens^c, Parvin Mousavi^{b,**}

^aUniversity of Toronto, Department of Computer Science, 40 St George St., Toronto, ON, Canada, M5S 2E4

^bQueen’s University, Medical Informatics Laboratory, 557 Goodwin Hall, Kingston, ON, Canada, K7L 2N8

^cKingston Health Sciences Research Centre, 76 Stuart Street, Kingston, ON, Canada, K7L 2V7

Abstract. Purpose: Prostate Cancer (PCa) is a prevalent disease among men, and multi-parametric MRIs offer a non-invasive method for its detection. While MRI-based deep learning solutions have shown promise in supporting PCa diagnosis, acquiring sufficient training data, particularly in local clinics remains challenging. One potential solution is to take advantage of publicly available datasets to pre-train deep models and fine-tune them on the local data, but multi-source MRIs can pose challenges due to cross-domain distribution differences. These limitations hinder the adoption of explainable and reliable deep-learning solutions in local clinics for PCa diagnosis. In this work, we present a novel approach for unpaired image-to-image translation of prostate multi-parametric MRIs and an uncertainty-aware training approach for classifying clinically significant PCa, to be applied in data-constrained settings such as local and small clinics. **Method:** Our approach involves a novel pipeline for translating unpaired 3.0T multi-parametric prostate MRIs to 1.5T, thereby augmenting the available training data. Additionally, we introduce an evidential deep learning approach to estimate model uncertainty and employ dataset filtering techniques during training. Furthermore, we propose a simple, yet efficient *Evidential Focal Loss*, combining focal loss with evidential uncertainty, to train our model effectively. **Results:** Our experiments demonstrate that the proposed method significantly improves the Area Under ROC Curve (AUC) by over 20% compared to the previous work (98.4% vs. 76.2%). **Conclusions:** Our proposed framework effectively translates and aligns public data with local data to increase the number of training data for deep models. Additionally, our proposed uncertainty estimation method enhances PCa detection performance. Providing prediction uncertainty to radiologists may aid in prioritizing uncertain cases, and expediting the diagnostic process effectively. Our code is available at https://github.com/med-i-lab/DT_UE_PCa.

Keywords: Deep Learning, Image Translation, Uncertainty Estimation, Prostate Cancer.

*The majority of this work was completed while M.Zhou was at the Medical Informatics Laboratory, Queen’s University.

**Parvin Mousavi, corresponding author, mousavi@queensu.ca

1 Introduction

Prostate Cancer (PCa) is a prevalent form of cancer among men [1], and the clinically significant PCa is defined by the Gleason score > 6 or the histopathology ISUP grade ≥ 2 [2, 3]. The current PCa diagnosis procedure involves a combination of the prostate-specific antigen test and the histopathology analysis of the Transrectal Ultrasound-guided biopsy (TRUS) taken from 10-12 regions on the prostate gland [4, 5]. However, the histopathology analysis on TRUS can miss up to 20% of clinically significant PCa due to the limited number of biopsy samples [1, 5].

Multi-parametric Magnetic Resonance Imaging (mp-MRI) has emerged as an effective alternative to TRUS for the early detection of PCa. mp-MRI uses a combination of anatomical and functional sequences of MRI that can

further highlight the differences between normal and abnormal (cancer) cells. The evaluation and reporting guideline of prostate mp-MRI was first introduced in the Prostate Imaging Reporting and Data System (PI-RADS) [6, 7, 8]. The guideline provides a comprehensive scoring schema for suspicious prostate lesions and mp-MRI sequences. An extensive Prostate MRI imaging study (PROMIS) [9] reported that targeted biopsy using mp-MRI has higher sensitivity and negative-predictive value (NPV) but lower specificity compared to TRUS biopsy [9, 10, 11]. The study also showed that 27% of the patients did not need to undergo biopsy, had mp-MRI been used for screening. Although PROMIS provides strong practical implications for mp-MRI in PCa diagnosis, the low specificity indicates that mp-MRI can be plausibly improved by advanced analyses.

In recent years, deep learning methods have emerged as a powerful tool for image classification tasks, and have provided promising performance in detecting and segmenting PCa on multi-parametric Prostate MRIs [12, 13, 14, 15, 16]. A more recent grand challenge, ProstateX [17], has further shown the ability of deep learning approaches in detecting clinically significant PCa on 3.0T mp-MRI data. Several groups have developed Convolutional Neural Network (CNN)-based models that achieve high performance for PCa classification [5, 17, 18, 19, 20, 21]. These methods have a great potential for clinical translations by highlighting abnormal lesions for radiologists during the PCa diagnostic process.

While deep learning has shown promising results in detecting PCa on mp-MRI, there are several substantial challenges in training and deploying deep models in clinics. Training deep models typically requires a large amount of data, which is not always available for local clinics with limited patient throughput. The alternative is to use pre-trained models or publicly available labeled data for training. However, the effectiveness of these models is significantly impacted by the differences in distribution between the public models/data and the images from local centers [5]. The primary contributing factor to this disparity is the field strength of MRI scanners.

While 3.0T scanner generally has higher image quality and spatial resolution [22] and most of the public datasets for prostate cancer are predominantly acquired with 3.0T scanners [23, 24, 25], 1.5T machines are considered standard of care and widely used in local clinical centers due to their cost-effectiveness [26]. Despite the preference for 3.0T MRI scanners, PI-RADS guidelines state that both 1.5T and 3.0T can provide adequate and reliable diagnostic examinations [27] and certain research indicates no statistically significant difference in clinical impact between 1.5T and 3.0T MRIs for diagnosis of prostate cancer [28, 29, 30]. Considering that over 80% of MRI systems in Canada operate

at a 1.5T while exhibiting a higher rate of exams per hour than their 3.0T counterparts [31], 1.5T MRIs continue to be the standard in clinical care.

Therefore, to bridge the distribution gap between the public data and clinical application, a necessary step is to translate 3.0T MRI data to 1.5T. This translation process facilitates aligning the data distributions and allows for the training of classification models using both the translated data and local data. It is worth noting that there are ongoing efforts in the literature to address related challenges in federated learning [32, 33]. However, this aspect is not the focus of this study.

Besides, classical deep models are primarily designed for predicting labels during the inference of test data, irrespective of whether the test image is within or outside the training set distribution. These models lack the capability to discern data samples belonging to unrelated distributions [34] or to quantify their confidence in predictions. These limitations make models hard to interpret, raising concerns about the reliability of such models. Hence, reusing and deploying models for local PCa detection is challenging. Addressing the above limitations and drawbacks is crucial when integrating deep learning models into real clinical routines. Thus, two main questions arise in this context:

1. For small local clinical centers, can they take advantage of the extensive high-resolution 3.0T public MRI data to enhance the classification performance on their limited low-resolution local 1.5T MRI data?
2. When deploying models in clinical centers, can we provide supplementary information regarding the confidence of the model’s predictions, beyond the final result, to enhance the reliability and explainability of the models?

In this work, we aim to answer the two important questions outlined above. To this end, we propose a novel 2-stage learning framework for clinically significant PCa classification using multi-parametric, multi-center MRI data. Our framework not only aims to enhance classification performance but also provides an estimate of predictive confidence alongside the corresponding predicted label, introducing a valuable dimension to the model’s interpretability. In the first stage, we introduce a data preprocessing pipeline that translates prostate mp-MRI data from 3.0T to 1.5T via a Generative Adversarial Network (GAN) approach to increase the number of training samples. This step addresses the challenge of limited data in local clinics with low patient throughput (refer to Section 4.1). In the second stage, we propose an uncertainty-aware PCa classification approach. Specifically, we explore various model architectures to enhance classification performance, experiment with different strategies for combining multi-parametric MRI data, and

leverage the *co-teaching* framework [35] to mitigate potential issues related to noisy labels (refer to Section 4.2.1). During the training phase, we incorporate dataset filtering using *evidential uncertainty estimation* [34] to eliminate training data samples with high prediction uncertainty, thereby enhancing the robustness of our models. Additionally, we extend the work of [34] by introducing a novel *Evidential Focal Loss* to optimize our classification models during training (see Section 4.2.2). The experimental results affirm the effectiveness of our proposed framework, demonstrating a significant improvement in classification performance compared to previous work. Our contributions not only advance the field of PCa detection but also underscore the potential of uncertainty-aware approaches in enhancing the reliability and interpretability of deep learning models in clinical settings.

Contributions: In summary, our work makes three main contributions:

1. We develop a GAN-based framework to translate unpaired prostate mp-MRIs from 3.0T to 1.5T, which we termed as domain transfer. This framework would align different data distributions and increase the number of training data for deep classification models. This is also the first attempt to translate prostate mp-MRIs in an unpaired manner.
2. We propose a novel loss function termed *Evidential Focal Loss* that can jointly compute the uncertainty for training samples and optimize the classification model. To the best of our knowledge, it is the first time that the original Focal Loss [36] is combined with the evidential uncertainty [34] for binary PCa classification.
3. By filtering the training samples based on their uncertainty value, our results outperform the state-of-the-art and improve the interpretability of model predictions. By providing confidence estimates for the predictions, radiologists can make informed decisions during the PCa diagnostic process and effectively expedite the process.

2 Related Work

2.1 Domain Adaptation

Machine learning algorithms usually perform well when training and test data share the same distribution and feature space. However, in real-world applications, the distribution of test data often shifts, leading to biased or inaccurate predictions. In addition, it is time-consuming or infeasible to acquire new training data and fully repeat training steps. Domain Adaptation (DA) is an approach that addresses this issue by mitigating the dataset bias or domain

shift problem caused by different distributions. There has been a lot of work on this topic in the past few years, which can be grouped into the following three general tasks [37]: (1) unsupervised DA tasks [38, 39, 40, 41, 42] focus on addressing the domain shift problem without requiring labeled target domain data; (2) semi-supervised DA tasks [43, 44, 45] aim to explore the partially labeled target domain data to further enhance the performance of domain adaptation algorithms; and (3) multi-source DA tasks [46, 47, 48] deal with scenarios where multiple source domains are available for adaptation. DA methods are often used to extract domain-invariant features for transferring knowledge between source and target domains. These methods incorporate various learning objectives with deep neural networks [49] for distribution matching: (1). **Discrepancy Measurement-based** methods aim to align feature distributions between two domains by fine-tuning deep models, e.g., using statistic criterion like Maximum Mean Discrepancy [50, 51, 52], and class criterion [53, 54, 55]. Some of these methods often require large labeled target domain data to diminish the domain shift problem, which is sometimes infeasible to get such medical data in the real-life scenario. (2). **Adversarial-based** methods aim to confuse domain discriminators from Generative Adversarial Networks (GANs) to enhance the invariant feature extraction [39, 56, 57]. One common scenario involves utilizing noise vectors, either with or without source images, to generate realistic target images while preserving the source features. However, training GANs are hard and sometimes results in generator degradation, e.g., mode collapse [58]. (3). **Reconstruction-based** methods, in addition to the general GANs approach from the above category, aim to reconstruct source-like images as an auxiliary task to preserve domain invariant features through an adversarial reconstruction paradigm [59, 60]. These methods usually have superior performance over the conventional GANs approach because they have an explicit reconstruction task to supervise the entire pipeline and make the training process more stable.

CycleGAN [60] is one of the state-of-the-art unsupervised adversarial reconstruction-based methods that is widely used for unpaired image-to-image translation. Its cycle consistency loss ensures the pixel-level similarity between two images through a reconstruction task, i.e., the source image s is translated to the target domain \hat{s} and then translated back \tilde{s} , where it should be identical to the original image ($s = \tilde{s}$). However, a drawback of the cycle consistency loss lies in its stringent constraint on pixel-level similarity, which will degrade the performance of GANs in some certain tasks [61]. To address this limitation, the adversarial consistency loss GAN (ACL-GAN) [61] is proposed to replace the pixel-level similarity with the distance between distributions. This modification allows ACL-GAN to retain essential features from source images while overcoming the drawbacks associated with the strict cycle consistency constraint.

Therefore, we adapt the ACL-GAN model and build our framework based on it.

In medical imaging, domain shift problems usually fall into two variations: subject-related variation (age, gender, etc.), and acquisition-related variation (MRI vendor, field strength, imaging protocol, etc.) [62]. To solve such problem, one intuitive approach is to fine-tune a model that is pre-trained on the source domain with the new data from the target domain. [63] propose to use the pre-trained VGG model on the ImageNet dataset [64] to learn robust high-level features of natural images, and then fine-tune it on the labeled MR images for the Alzheimer’s Disease (AD) classification task to achieve state-of-the-art performance. Similarly, [65] study the impact of the fine-tuning techniques on the brain lesion segmentation task, demonstrating that fine-tuning with only a small number of target domain training samples can outperform models trained from scratch. Another approach is to use domain adaptation as an intermediate step to reduce variance in image acquisition parameters from both domains and then use it for downstream tasks. Researchers have attempted to address the problem of acquisition variation in MRI data for several years. [62] propose a feature-level representation learning method to either extract acquisition-invariant features or remove acquisition-variant features from paired 1.5T and 3.0T brain MRIs. The learned features are then used for a downstream classification task. However, obtaining paired 1.5T and 3.0T MRI data in real-life scenarios is impractical. Another way to align acquisition-invariant features is to synthesize images from different types of acquisition parameters using GAN-based adversarial reconstruction methods. GANs have been applied to perform cross-modality image translation between different medical images or generate synthetic images from random noise. The objective of such translation tasks is to retain the underlying structure while changing the appearance of the image [66]. Researchers have attempted to estimate images in the target modality from the source modality, such as MRI-CT translation [67, 68, 69, 70] and X-ray to CT translation [71, 72]. Other areas that have been explored include intra-modality translation, such as T1/T2-FLAIR translation [73, 74] and pure data augmentation by generating synthetic images from random noise vectors [75, 76, 77, 78]. However, most of the works do not consider the real clinical practicality, for example, T1/T2-FLAIR MRI translation may require paired training data, which is not feasible in real clinical settings. Generating synthetic images from noise does not take advantage of the publicly available data and ignores *a-priori* information. The current limitations provide great potential for unpaired image translation for medical images, which we employ in this work.

2.2 Deep Learning for PCa Classification

The use of 3D-CNN models has gained widespread popularity for classifying PCa based on volumetric image data due to their excellent performance. [20] propose a feature fusion 3D-CNN to classify clinically significant PCa using mp-MRI data. They use ADC maps, DWI, and K^{trans} 3.0T MR data to enable the model to learn multi-modal information. Inspired by the VGG architecture [79], the model has three VGG-like feature extractors for each image modality, followed by the concatenation between outputs of each extractor and a vector represents the zonal information of the suspicious region. On the test set, the proposed model achieves the area under the receiver operating characteristic (AUC) curve of 0.80 on 140 unseen patients. [19] propose a similar VGG-like 3D-CNN architecture for the same PCa classification task. Different from [20], they only have one model for feature extraction. To obtain the multi-modal information, they stack three images from each of the ADC maps, DWI, K^{trans} into one 3-channel image as the input. The model achieves the AUC of 0.84 on the test set.

In [14], a probabilistic approach using mp-MRI data is employed for PCa classification. The authors develop an automated pipeline for the classification of clinically significant PCa using 3.0T DWI images from 427 patients. The pipeline consists of three parts: classification of each DWI slice using the pre-activated ResNet model [80], extraction and selection of first-order statistics from the CNN outputs, and final class label prediction using a random forest classifier. On the test set, the model achieves an AUC of 0.87. While the aforementioned studies may yield favorable AUCs, the reproducibility of the model might be challenging in clinics with limited patient (data) throughput. Recently, [5, 21] address the data-hungry problem by introducing a disentangled representation learning approach (SDNet) to synthesize public 3.0T MRI images into 1.5T MRI images to increase the training data size for centres with limited 1.5T data. Their approach aims to separate the anatomy- and modality-specific features present in images, subsequently merging the 1.5T modality features with the 3.0T anatomical features to generate MRI images resembling those acquired at 1.5T. Finally, a simple 3D-CNN classifier is used for the binary classification of clinically significant PCa. The model outperforms the state-of-the-art performance in PCa classification through domain alignment between different data sources.

While these methods exhibit excellent classification performance, a common limitation is the absence of a confidence score for their predictions, which hinders their interpretability in clinical practice.

2.3 Uncertainty Estimation

Recent studies in medical imaging have highlighted the detrimental impact of label noise on the performance of modern deep learning models [81]. Conventional regularization techniques such as dropout, batch normalization, weight decay, etc. fail to effectively address this challenge [82, 83]. Methods proposed to mitigate such problem can be broadly categorized into three groups [84]: (1) Robust loss functions and loss adjustments [85, 86, 87] aiming to stabilize the model performance when optimizing its parameters; (2) Sample selection [88, 89, 90] aiming to select a subset of “clean” data from a batch of samples to compute the loss; and (3) Robust architectures [35, 91] aiming to learn the same data by training multiple models with different initialization assess output stability. While these methods inherently handle the noisy label problem, they can not provide explicit uncertainty estimation in terms of confidence in their output. Moreover, the ability of deep learning models to identify irrelevant samples remains limited. For instance, when a model trained on prostate MRIs is presented with a CT scan of the prostate at the time of inference, it is unclear whether the model can provide meaningful predictions or simply indicate a lack of in-domain knowledge and perform a human-in-the-loop analysis instead. In recent years, research has been conducted on uncertainty estimation for deep learning models. [92, 93] develop the *dropout neural networks* framework to represent the prediction uncertainty of deep learning models, where the dropout layers in the model are formed by Bernoulli distributed random variables. During the test phase, predictive uncertainty is determined by enabling dropout layers and averaging the results over multiple runs, providing a valuable mechanism for uncertainty quantification in model predictions. An alternative method for modeling uncertainty in deep learning models is through the use of *evidential neural networks* [34], which formulate uncertainty by fitting a Dirichlet distribution - acting as the conjugate prior of the categorical distribution - to the class probabilities acquired from neural networks. This method considers model predictions as multinomial subjective opinions [94] or beliefs [95], which can be further modeled explicitly using subjective logic. The "evidential" approach emphasizes the ability of the model to deliver certain predictions and exhibits superiority compared to the dropout approach [92].

In clinical practice, uncertainty estimation is crucial. By integrating uncertainty information into prediction outcomes, misclassification rates can be significantly reduced. For instance, in radiograph classification task [96], the authors employ the Dempster-Shafer Theory of Evidence [95] and the principles of subjective logic [94] to develop a framework that jointly estimates per-class probabilities and provides predictive uncertainty. This approach has been

extended to abdominal ultrasound and brain MR images [97]. In the context of breast cancer classification, [98] apply the evidential neural networks approach [34] to effectively diagnose breast cancer. A similar approach is used for the same task by [99] through the evidence adjustment technique, which focuses on the difference in the risks of uncertain samples from different classes. Consequently, we build upon the work from [34] by adding uncertainty estimation to improve the robustness of the model and the interpretability of predictions.

3 Materials

3.1 Data

In this work, we use both large publicly available ProstateX data and small private local clinical data. A visualization of sample images from both datasets is presented in Figure 1.

ProstateX Grand Challenge Data (3.0T). The 3.0T data is provided by the International Society of Optics and Photonics in the “ProstateX” challenge [100]. The dataset contains T2-weighted (T2), maximum b-value diffusion, diffusion-weighted imaging (DWI) with apparent diffusion coefficient (ADC) maps, and K^{trans} images of 346 patients undergoing prostate biopsies. T2 images show the anatomical structure of the prostate, and both the ADC maps and K^{trans} could further highlight the differences between normal and abnormal (cancer) cells in the MRI scans [101, 102]. We only use 204 of the total 346 patients in this work since these are reserved as training data, and hence they are provided with the spatial location of the suspicious finding, and a binary label indicating whether or not there is cancer. The remaining 142 patients are reserved as the test set and no labels are provided, hence, we exclude those from our work.

Kingston Health Science Center Data (1.5T). The local 1.5T data is obtained from the Kingston Health Science Center (KHSC), which contains 104 patients with the corresponding biopsy-confirmed cancer and the Gleason Score. For the local data, only T2, ADC, and b-value images are available. All patients MRI have the spatial location of the suspicious finding(s), the Gleason Score, and the binary label indicating whether it is a cancer lesion or not.

Since all patients in both datasets have complete T2 and ADC data, our focus in this work is solely on these two types of images. Each MRI data in our study is associated with a single patient. Both datasets are processed similarly unless stated.

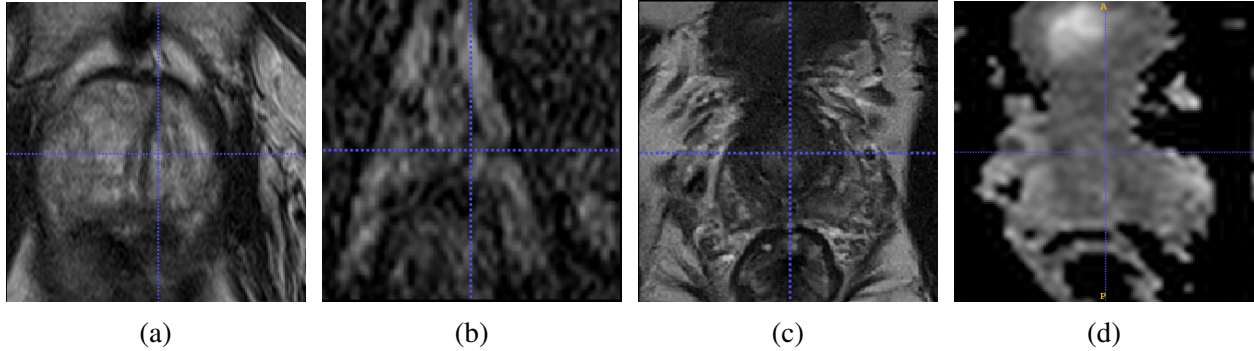


Fig 1: Visualization of sample data. **1a** and **1b** are the 1.5T T2 and ADC images from KHSC, respectively. Similarly, **1c** and **1d** are the 3.0T T2 and ADC images from the “ProstateX” Challenge, respectively.

3.2 Pre-processing

T2 and ADC sequences from both datasets are $160 \times 160 \times C$, where C is the total number of slices in the MRI. We resample all 3D data to have the same voxel spacing. To reduce aliasing artifacts, the most common voxel spacing ($0.5 \times 0.5 \times 3 \text{ mm}^3$) is used across all data, and the cosine-windowed interpolation is utilized during sampling. We normalize pixel intensities to $[-1, 1]$ for all data. For the translation purpose from 3.0T to 1.5T, we further resample all 3D data to $256 \times 256 \times C$ and split into C 2D gray-scale slices.

Augmentation: For each patient, the MRI volume undergoes rotation ranging from 0 to 100 degrees in 5-degree increments, hence expanding the data size 20-fold.

Cropped Patches: To reduce the computational cost, cropped patches of the MRI volume were employed. The process involves identifying the suspicious slice (i_s) based on the provided spatial location. Recognizing that PCA lesions can span multiple slices, two neighboring slices (i_{s-1} and i_{s+1}) are selected as well and cropped around the biopsy location to generate a patch of size $64 \times 64 \times 3$.

4 Methods

Figure 2 summarizes an overview of our proposed approach. The domain transfer framework aims to reduce the distribution-level discrepancy between two prostate MRI datasets. The framework matches the acquisition parameters of publicly available, large 3.0T prostate mp-MRI data with local, small 1.5T prostate mp-MRI data. Once all the data from 3.0T are translated to 1.5T, a subsequent classifier is trained to classify clinically significant PCA. Furthermore, during the training process, the uncertainty is calculated along with the class output. We also introduce a novel

evidential focal loss for the PCa classification task. Lastly, we utilize dataset filtering to improve robustness and accuracy by eliminating uncertain data samples from the training set.

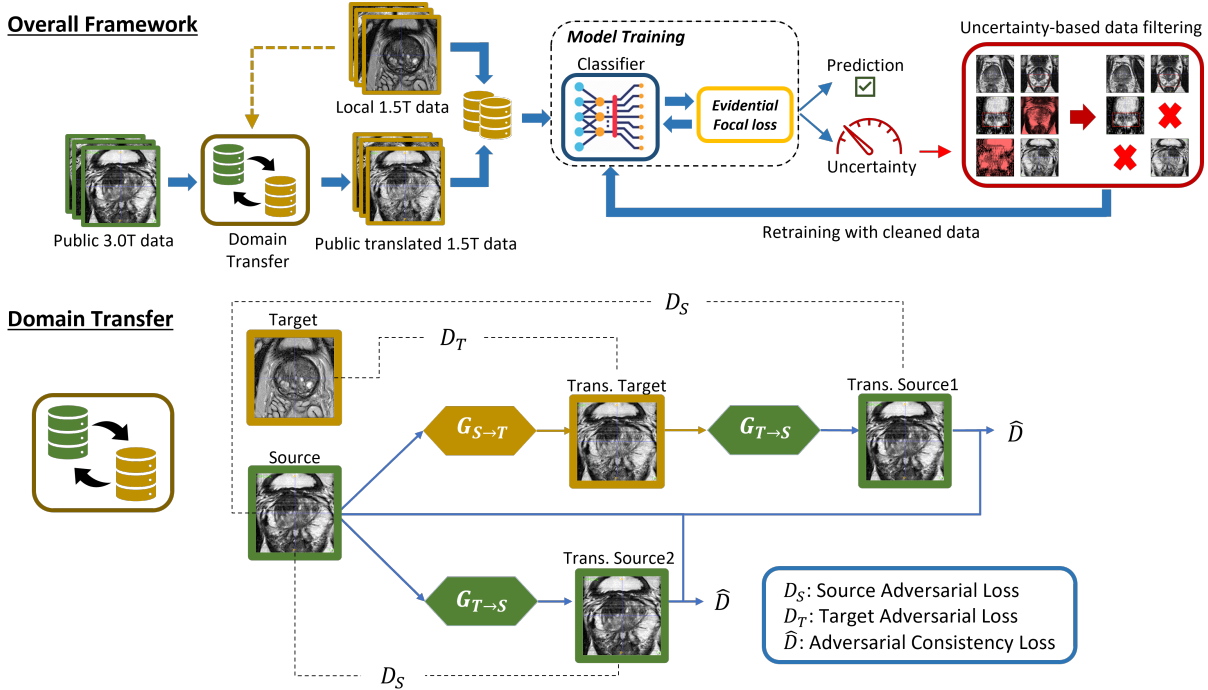


Fig 2: Detailed schematic of the proposed method. The overall framework of our proposed method contains two stages: 1), domain translation to map public 3.0T MRI with local 1.5T MRI; 2), uncertainty-aware clinically significant PCa classification. The bottom figure is the training schema for domain transfer. The upper right portion of the figure illustrates the PCa classification training process, which involves training the classifier using the Evidential Focal loss, filtering the training set based on uncertainty, and retraining the classifier on the filtered data to obtain the final classifier.

4.1 The Domain Transfer Framework

We adopt the ACL-GAN model [61] for unpaired MR image translation from 3.0T to 1.5T. We have made modifications to the architecture to adapt it to single-channel MRI slices. There are two generators in this model namely $G_{T \to S}$ and $G_{S \to T}$, $G_{T \to S}$ translates the images from the target domain (X_T) to the source domain (X_S) given the input $x \in X_T$ and a noise vector z sampled from $\mathcal{N}(0, 1)$. On the other hand, $G_{S \to T}$ performs the reverse process of $G_{T \to S}$, translating images from the source domain to the target domain. There are three discriminators, D_S , D_T , and \hat{D} in this model. The first two ensure that translated images are in their respective domains by optimizing adversarial losses. The third discriminator, \hat{D} , ensures that translated images retain anatomical features in 3.0T by distinguishing the pair (Source, Trans. Source1) and (Source, Trans. Source2), as shown in the bottom of Figure 2. The loss func-

tion of ACL-GAN contains four parts [61]. First, the adversarial \mathcal{L}_{adv} loss aims to encourage the translated image, either from source to target domain or from target to source domain, to be in the correct domain. The \mathcal{L}_{adv} is further decomposed to source-domain and target-domain adversarial loss, \mathcal{L}_{adv}^S and \mathcal{L}_{adv}^T , respectively. The \mathcal{L}_{adv}^T is given by:

$$\begin{aligned} \mathcal{L}_{adv}^T(G_{S \rightarrow T}, D_T, X_S, X_T) &= \mathbb{E}_{x_T \sim p_T} [\log D_T(x_T)] \\ &+ \mathbb{E}_{\bar{x}_T \sim p_{\{\bar{x}_T\}}} [\log(1 - D_T(\bar{x}_T))] \end{aligned} \quad (1)$$

where $\bar{x}_T = G_{S \rightarrow T}(x_S, z_1)$ and $z_1 \sim \mathcal{N}(0, 1)$. Similarly, the \mathcal{L}_{adv}^S is given by:

$$\begin{aligned} \mathcal{L}_{adv}^S(G_{T \rightarrow S}, D_S, \{\bar{x}_T\}, X_S) &= \mathbb{E}_{x_S \sim p_S} [\log D_S(x_S)] \\ &+ (\mathbb{E}_{\hat{x}_S \sim p_{\{\hat{x}_S\}}} [\log(1 - D_S(\hat{x}_S))] \\ &+ \mathbb{E}_{\tilde{x}_S \sim p_{\{\tilde{x}_S\}}} [\log(1 - D_S(\tilde{x}_S))]) / 2 \end{aligned} \quad (2)$$

where $\hat{x}_S = G_{T \rightarrow S}(\bar{x}_T, z_2)$, $\tilde{x}_S = G_{T \rightarrow S}(x_S, z_3)$, $z_2, z_3 \sim \mathcal{N}(0, 1)$. Combine these two adversarial losses we get the total adversarial loss for ACL-GAN, $\mathcal{L}_{adv} = \mathcal{L}_{adv}^T + \mathcal{L}_{adv}^S$.

A problem of \mathcal{L}_{adv} is that this loss can not encourage the translated image to the target domain \bar{x}_T is similar to the corresponding source domain image x_S , as we do not want the model to change the anatomical structure or features as we discussed previously. Hence, the adversarial consistency loss is proposed, which is given by:

$$\begin{aligned} \mathcal{L}_{acl} &= \mathbb{E}_{(x_S, \hat{x}_S) \sim p_{(x_S, \{\hat{x}_S\})}} [\log \hat{D}(x_S, \hat{x}_S)] \\ &+ \mathbb{E}_{(x_S, \tilde{x}_S) \sim p_{(x_S, \{\tilde{x}_S\})}} [\log(1 - \hat{D}(x_S, \tilde{x}_S))] \end{aligned} \quad (3)$$

where $\bar{x}_T = G_{S \rightarrow T}(x_S, z_1)$, $\hat{x}_S = G_{T \rightarrow S}(\bar{x}_T, z_2)$, $\tilde{x}_S = G_{T \rightarrow S}(x_S, z_3)$.

Next, \mathcal{L}_{idt} is the identity loss, which encourages generators to perform approximately identity mapping when images in the respective domain are provided, e.g., x_S to $G_{T \rightarrow S}$ or x_T to $G_{S \rightarrow T}$. The \mathcal{L}_{idt} is given by:

$$\mathcal{L}_{idt} = \mathbb{E}_{x_S \sim p_S} [\|x_S - x_S^{idt}\|_1] + \mathbb{E}_{x_T \sim p_T} [\|x_T - x_T^{idt}\|_1] \quad (4)$$

where $E_S^z : X_S \rightarrow Z$ and $E_T^z : X_T \rightarrow Z$ are two noise encoder networks for G_S and G_T , respectively, which map

images to noise vectors. $x_S^{idt} = G_{T \rightarrow S}(x_S, E_S^z(x_S))$ and $x_T^{idt} = G_{S \rightarrow T}(x_T, E_T^z(x_T))$.

Finally, \mathcal{L}_{mask} is used to force both generators to only modify certain regions of the source image and keep the rest of the areas unchanged. We let generators produce a two-channel image, where the first channel is one of the translated images between the source and target domain (i.e., one-channel gray-scale prostate MRIs), and the second channel is the bounded mask, whose values are between $[0, 1]$. The \mathcal{L}_{mask} is given by:

$$\begin{aligned} \mathcal{L}_{mask} = & \delta[(\max\{\sum_k x_m[k] - \delta_{max} \times P_t, 0\})^2 \\ & + (\max\{\delta_{min} \times P_t - \sum_k x_m[k], 0\})^2] \\ & + \sum_k \frac{1}{|x_m[k] - 0.5| + \epsilon} \end{aligned} \quad (5)$$

where δ , δ_{max} and δ_{min} are hyper-parameters for controlling the size of masks, $x_m[k]$ is the k-th pixel of the mask and P_t is the total number of pixels in an image. The ϵ is a very small value to avoid dividing by zero. The first term of this loss controls the size of the mask. It encourages the generator to perform sufficient modifications while preserving the background information. Here, δ_{max} and δ_{min} are the maximum and minimum proportions of the foreground in the mask, i.e., the region we want to modify. The last term of this loss encourages the mask to be binary, either 0 or 1 to segment the foreground and background of the input image.

Aggregating all loss terms together, we have:

$$\mathcal{L}_{total} = \mathcal{L}_{adv} + \lambda_{acl}\mathcal{L}_{acl} + \lambda_{idt}\mathcal{L}_{idt} + \lambda_{mask}\mathcal{L}_{mask} \quad (6)$$

where λ_{acl} , λ_{idt} , λ_{mask} are the weighting factors for \mathcal{L}_{acl} , \mathcal{L}_{idt} , \mathcal{L}_{mask} , respectively.

4.2 Uncertainty-aware PCa Classification

4.2.1 Classifier architectures

The traditional CNN approach is used for the clinically significant PCa binary classification task. Specifically, we explore three different model architectures for combinations of T2 and ADC patches as the classifier (Top row of Figure 2). The first architecture, called the multi-stream CNN (M.S. MpMRI), treats T2 and ADC patches as separate inputs,

as shown in Figure 3. The model takes 3D patches of T2 and ADC as parallel inputs, which are then processed by the same feature extractor (i.e., weights are shared) to extract deep semantic representations. The output representations of T2 and ADC are then concatenated channel-wise and fed into another convolutional layer followed by a fully connected layer to produce the class probabilities.

In the second architecture, we adopt two different ways to combine ADC and T2 patches into a single input for the network. First, we stack cropped 3D patches of T2 and ADC along the channel axis, resulting in an input data size of $64 \times 64 \times 6$. Alternatively, we consider only the located suspicious slice i_s for both T2 and ADC, and stack them along the channel axis to obtain the input data size of $64 \times 64 \times 2$. The model architecture for both combinations is similar to Figure 3, where there is only one branch and no concatenation afterward. We refer to the model with input size of $64 \times 64 \times 6$ (resp. input size $64 \times 64 \times 2$) as Vol. MpMRI (resp. MpMRI).

Lastly, we use only 3D T2 patches as input to match with the previous work [5]. The model architecture is same as the one for MpMRI, and we refer to this model as ‘‘T2-only’’.

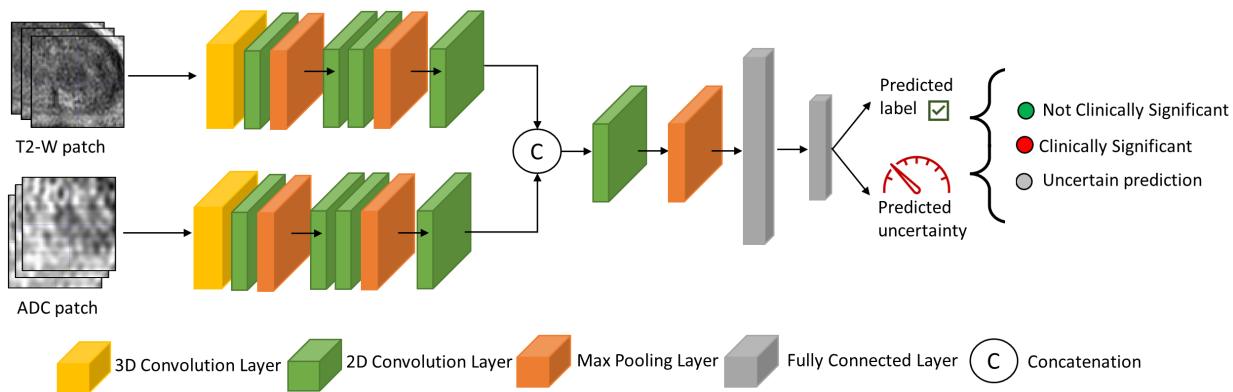


Fig 3: Detailed architecture of ‘‘M.S. MpMRI’’ model. The first sequence of CNN layers contains $1 \times$ 3D convolution layer and $4 \times$ 2D convolution layers, $2 \times$ Max Pooling layers with window size 2×2 . Both extracted feature maps of T2 and ADC are concatenated channel-wise. After that, another set of convolution-max pooling layers is utilized. Finally, the extracted 2D features are reshaped to 1D and fed into a Fully connected layer follow by a softmax layer with 2 outputs representing the probabilities of which class the input data belongs to.

We also explore with the *co-teaching* framework [35] to combat the potential noisy label issue to improve the classification performance. This framework can be directly used as our classifier in a plug-and-play manner. The key distinction between the co-teaching framework and the conventional training paradigm lies in the fact that co-teaching only influences the number of samples in the batch that contribute to the loss calculation. We follow the exact setup in

the original paper, and detailed hyperparameter settings can be found in Section 5. The backbone model we selected for this framework is the ‘‘MpMRI’’ architecture, which we depicted previously.

4.2.2 Evidential Focal Loss

Dataset filtering during the training phase could reduce the effect of the noisy label on the deep model. Followed by [96], the process of uncertainty-based filtering is shown at the top of Figure 2: Firstly, we calculate the uncertainty value for each sample in the training set. We then remove a portion of the training samples that exhibit high predictive uncertainty. Finally, we retrain the model using the remaining ‘‘clean’’ training data.

Following the previous work [34], we extend and combine the idea of subjective logic [94] with the focal loss [36] for the clinically significant PCa binary classification task. In the context of the Theory of Evidence, a belief mass is assigned to individual attributes, e.g., the possible class label of a specific data sample. The belief mass is generally calculated from the evidence collected from the observed data [95]. Let K be the number of classes, and $b_k \geq 0, k \in [1, K]$ is the belief mass for class k and $u \geq 0$ is the overall uncertainty measure. Let $e_k \geq 0$ be the evidence computed for k^{th} class, then the belief b_k and the uncertainty u are computed as follows:

$$b_k = \frac{e_k}{S} \quad \text{and} \quad u = \frac{K}{S} \quad (7)$$

where $S = \sum_{i=1}^K (e_i + 1)$. For our binary task ($K = 2$), we can further simplify Equation (7) to $b_0 = \frac{e_0}{e_0+e_1+2}$, $b_1 = \frac{e_1}{e_0+e_1+2}$, and $u = \frac{2}{e_0+e_1+2}$. The belief mass assignment, e.g., subjective opinion, corresponds to the Dirichlet distribution with parameters $\alpha_k = e_k + 1$, and $S = \sum_{i=1}^K \alpha_k$ is the Dirichlet strength. The expected probability for the k^{th} class is calculated by the mean with the associate Dirichlet distribution, e.g., $\hat{p}_k = \frac{\alpha_k}{S}, k \in [1, \dots, K]$ [34].

Now, we formally define our proposed evidential focal loss. Given the training set contains N data samples, $D := \{x_i, y_i\}_{i=1}^N$, where x_i is the i^{th} data sample and $y_i \in [0, 1]$ is the corresponding label, 0 is the negative sample and 1 is the positive sample. We further denote \mathbf{y}_i as the one-hot encoding label for sample i , e.g., $\mathbf{y}_i = [1, 0]$ for class 0 and $\mathbf{y}_i = [0, 1]$ for class 1. The focal loss [36] for binary classification is defined by $FL(p_t) = -w_t(1-p_t)^\gamma \log(p_t)$ where $p_t = p$ if $y_i = 1$ for i^{th} sample, otherwise $p_t = 1 - p$ with probability output p from the model, and w is the class weight. Let \mathbf{P}_i be a vector that contains the probability of i^{th} sample for both classes from our model output;

p_{ij} is the probability of i^{th} sample belonging to j^{th} class; K is the number of classes, and β_j is the class weight of j^{th} class. γ is the focusing parameter to reduce the loss for well-classified samples, and we fix $\gamma = 2$ in this task. We could define Evidential Focal Loss as the following:

$$\mathcal{L}_i^{cls}(\theta) = \int \sum_{j=1}^K -\beta_j (1 - p_j)^\gamma \log(p_{ij}) \frac{1}{\mathcal{B}(\alpha_i)} \prod_{j=1}^K p_{ij}^{\alpha_{ij}-1} d\mathbf{P}_i \quad (8)$$

Here, we denote $\mathcal{L}_i^{cls}(\theta)$ as the classification loss for a single sample i , α_i is the belief mass α for i^{th} sample for all classes and α_{ij} is the belief mass α for i^{th} sample and j^{th} class. $\mathcal{B}(\cdot)$ is the multinomial beta function and $\frac{1}{\mathcal{B}(\alpha_i)} \prod_{j=1}^K p_{ij}^{\alpha_{ij}-1}$ is the probability density function of Dirichlet distribution. Rewriting class probabilities in vector form, the equation (8) can be simplified to (9) by the definition of expectations:

$$\mathcal{L}_i^{cls}(\theta) = - \sum_{j=1}^K \beta_j \mathbf{E}[(1 - \mathbf{P}_i)^2 \log(p_{ij})] \quad (9)$$

Different from the original focal loss, we replace the constant term 1 with \mathbf{y}_i , to tackle the hard-to-classified samples and reduce the loss of well-classified samples *in both classes*. Recall that expected probability \hat{p}_k for the k^{th} class is α_k/S , then by the linearity of expectations and the definition of expectations of Dirichlet distribution, we have:

$$\mathcal{L}_i^{cls}(\theta) = \sum_{j=1}^K \beta_j (y_{ij} - (\alpha_j/S))^2 (\psi(S_i) - \psi(\alpha_{ij})) \quad (10)$$

where $\psi(\cdot)$ is the digamma function, y_{ij} is the j^{th} class label in the one hot encoding representation \mathbf{y}_i and β_j is the class weight vector for class j .

To ensure that highly uncertain data samples, referred to as "I do not know" decisions, do not impact the overall data fit and to minimize their associated evidence, we adopt the Kullback-Leibler (KL) loss as the regularization term to penalize the unknown predictive distributions, as done in [34]. Combining \mathcal{L}^{cls} and the KL loss yields our final *evidential focal loss* for uncertainty-aware classification:

$$\mathcal{L}^{efl}(\theta) = \sum_{i=1}^N \mathcal{L}_i^{cls}(\theta) + \lambda_t \sum_{i=1}^N KL[D(\mathbf{P}_i|\alpha_i)||D(\mathbf{P}_i|\mathbf{1})] \quad (11)$$

where $\mathbf{1}$ is an one-vector, $D(\mathbf{P}_i|\mathbf{1})$ is the uniform Dirichlet distribution. λ_t is the weighting factor of the KL divergence loss and is defined as $\lambda_t = \min(1.0, t/10) \in [0, 1]$, where t is the current number of epochs of training.

Finally, we introduce two proposed methods for filtering training samples based on the calculated uncertainty.

Patch-driven filtering: Given the uncertainty for each training patch, we simply eliminate $x\%$, $x \in [10, 20]$ of the *patches* with the highest uncertainty and retrain the model on the rest of the samples on the training set.

Patient-driven filtering: Similar to Patch-driven filtering, we first calculate the uncertainty of each training patch. To determine the uncertainty of each patient, we calculate the average uncertainty value across their corresponding patches (20 patches per patient as mentioned in Section 3.2). We then eliminate $x\%$, $x \in [10, 20]$ of the training *patients* with high uncertainty value and retrain the model on the rest on the training set.

5 Experiments & Setup

5.1 Setup

Data Split: For the domain transfer task, as mentioned in Section 3.2 and 4, the resampled T2 and ADC “images” with size 256×256 from both ProstateX and the local dataset is used for training the ACL-GAN model. Particularly, we allocate 90% of images in both datasets as training and keep 10% as the validation set to avoid overfitting the ACL-GAN model. Importantly, we ensure that the images corresponding to each patient are exclusively present in either the training or validation set, but not both. To improve the robustness and enhance the ability of the ACL-GAN to capture feature-level representations of 1.5T images, we use all data from our local hospital. However, it is important to note that this approach does not yield any additional impacts on the subsequent classification task. The model only modifies image regions that have visual differences caused by acquisition parameters of various MRI machines, but it does not alter the context of the prostate itself.

For the PCa classification task, we use cropped and augmented T2 and ADC *patches* from both datasets. As mentioned before, this includes 204 ProstateX patients translated to 1.5T, as well as 104 patients from our local hospital captured in 1.5T. Regarding the data split for the classification, we keep patches of 34 patients from our local

center as a standalone test set. From the remaining patches (70 local patients and all ProstateX patients), we allocate 80% for training and 20% for validation, assuring patches from the same patient are not included in both of these sets.

Next, we provide a brief intro to the experiments we conducted in this study.

Domain Transfer: The first experiment involved translating ProstateX MRI data from 3.0T to 1.5T using our proposed ACL-GAN model. We evaluated the effectiveness of our proposed approach by using quantitative metrics such as Fréchet Inception Distance (FID) score [103] and maximum mean discrepancy (MMD) score [104], more details could be found in Section 5.3. The translated MRI data was then employed in a downstream binary classification task for clinically significant PCa, demonstrating the superiority over the SDNet [5] baseline.

Classification: We divided our classification experiments into two folds. **(1). Conventional approach**, we used the conventional training paradigm without any filtering or uncertainty estimation. Three different model architectures were utilized for these experiments as discussed in Section 4.2.1. **(2). Uncertainty-aware approach**, we used the dataset filtering method and evidential focal loss proposed in Section 4.2.2 to train our models. Additionally, several ablation studies on data modalities, model architectures, and loss functions were conducted, with corresponding classification results detailed in Section 6. Finally, we focused on dataset filtering during deployment and examined how this technique affects the classification performance on the test data.

5.2 Experimental Details

We trained two ACL-GAN models separately for T2 and ADC images as part of our domain transfer framework. The optimizer used for both models was Stochastic Gradient Descent with Adam update rule [105], with an initial learning rate of 0.0001 and weight decay of 0.0001 to prevent overfitting. The batch size is 3 and are trained for 30,000 epochs. Moreover, when training the model for T2 images, we set the $\lambda_{mask} = 0.0025$, $\lambda_{idt} = 1$, $\lambda_{acl} = 0.2$ in Equation (6) and lower and upper mask threshold to be 0.005 and 0.1, respectively. When training the model for ADC images, the values of λ_{mask} , λ_{idt} , λ_{acl} are the same as in the T2 model with lower and upper mask thresholds set to 0.001 and 0.005, respectively. We adopt the Least-Square (LS) loss [106] for \mathcal{L}_{adv} and \mathcal{L}_{acl} in Equation (6), as done in [61].

Converting 3.0T to 1.5T: Once we have obtained two ACL-GAN models, we need to standardize the acquisition parameters of 3.0T prostate MRIs to match those of the 1.5T data in our local dataset. To achieve this, we divided the original 3.0T MRI into multiple 2D grayscale slices. For each 2D slice, we used the generator G_T and a noise vector

z randomly sampled from $\mathcal{N}(0, 1)$ to translate the 3.0T slice to 1.5T, e.g., $I_{1.5T} = G_T(I_{3.0T}, z)$ as discussed in Section 4.1. This process was repeated for all 2D slices, and the slices were stacked back together to reconstruct the 3D MRI for each patient. The voxel spacing remained unchanged before and after the translation process. This procedure was applied to both T2 and ADC data.

All classification models were trained with Stochastic Gradient Descent with Adam, and batch normalization was applied to expedite convergence. In the first category of classification experiments (Conventional approach), the traditional focal loss [36] with $\gamma = 2$ was employed. Specifically, all models except co-teaching were trained for 300 epochs with a learning rate of 0.0001, weight decay of 0.01, and a batch size of 10. For the co-teaching model, the noise rate and forget rate were set to 0.1, and the number of epochs for the linear drop rate to 10. The model was trained for 300 epochs, with a batch size set to 10, and a learning rate of 0.00001. In the second category of experiments (Uncertainty-aware approach), a learning rate of 0.0001 was used and decayed by a factor of 0.1 every 200 epochs. The weight decay was set to 0.01, the total training epochs were 300, and the batch size was 10. Further training details can be found in Appendix A.

5.3 Evaluation

To assess the quality of the translated images and validate the effectiveness of our proposed domain transfer framework for ProstateX data from 3.0T to 1.5T, we followed previous works [78, 107, 108] to compute the maximum mean discrepancy (MMD) score [104] and the Fréchet Inception Distance (FID) score [103]. We followed the implementation provided in the Project-MONAI library ¹. Both MMD and FID score measure the distribution distance between translated 1.5T images from the ProstateX dataset and real 1.5T images from our local hospital. Lower values in these metrics indicate greater fidelity to the real data distribution. The translated 1.5T images’ quality was further validated in a downstream classification task, combining them with the original 1.5T images to train a classifier distinguishing patients with and without clinically significant PCa. Traditional classification metrics, including accuracy (Acc.), sensitivity (Sen.), specificity (Spec.), and Area under ROC curve (AUC), were employed for this evaluation. Furthermore, we used “uncertainty calibration” to assess for performance of our uncertainty-based models. To compute calibration, we employed Expected Calibration Error (ECE) as done in [109, 110]. ECE measures the correspondence between the

¹<https://github.com/Project-MONAI/GenerativeModels/tree/main/generative/metrics>

confidence of predictions and the actual model accuracy. Unlike other classification metrics, smaller values of ECE indicate less miscalibration, hence indicating better model calibration capability.

Reporting the patient-level performance is more relevant to the real clinical setting. However, since patches were used as model input, aggregating individual results from these patches is necessary to calculate patient-level metrics. This involved using the classifier to predict test patches, which were then sequentially grouped into x groups, where x is the number of patients in the test dataset. For each patient, there are 20 patches hence probabilities due to the mentioned augmentation in Section 3.2. The *median* probability \tilde{p}_i was computed over 20 probabilities as the aggregated probability for each patient. Finally, a threshold of 0.5 was applied to determine whether a patient has PCa, assigning label 1 if $\tilde{p}_i > 0.5$ and 0 otherwise.

6 Results and Discussion

In this section, we analyze the translated image from our proposed domain transfer framework and report **patient-based** classification results; the performance of our methods on patches are reported in Appendix B.

6.1 Quantitative Analysis of Translated Samples

We evaluate the quality of translated T2 images using two common metrics MMD [104] and FID [103] score, as reported in Table 1. The FID score is computed over each slice in the MRI data, i.e., a $160 \times 160 \times 32$ MRI volume results in 32 individual slices and FID scores, and we take the average score across all slices as the final score. The ACL-GAN model we used outperforms the SDNet in both metrics, we attribute the observed improvements to several key factors. Unlike SDNet [5], which merges modality features from randomly selected 1.5T images with anatomical features from 3.0T, our method learns the overall data distribution in an adversarial manner. This enables it to capture the entire distribution of 1.5T images and perform the translation more effectively. Additionally, our approach ensures that the translated image retains crucial features of the original, with the generator making modifications only to specific parts of the image. These designs contribute to the superior performance of our ACL-GAN-based domain transfer framework.

Since SDNet [5] is solely trained on T2 images, a direct quantitative comparison for translated ADC images is omitted. Instead, we conducted an ablation study to highlight the improvement in classification performance of adding translated ADC images to our model, see details in Section 6.2.

	MMD ↓	FID ↓
SDNet [5]	0.0018	12.740
ACL-GAN (ours)	0.00054	11.464

Table 1: Quantitative results of translated T2 images using the baseline method SDNet and ACL-GAN in our proposed domain transfer framework. Lower values indicate better performance for all metrics. **Bold** values represent the best results.

6.2 PCa classification without filtering

Table 2 summarizes the experiment results in this section, which contains the PCa classification performance using the conventional approach, as detailed in Section 5.1. We observed that the AUC of using the co-teaching framework with MpMRI architecture as the base model achieves the best AUC and outperforms the baseline. The co-teaching model exhibits approximately a 50% increase in sensitivity while experiencing only a modest 10% decrease in specificity compared to the baseline model. This suggests that the co-teaching model demonstrates superior learning capabilities for classifying both positive and negative data samples on the test set. In the training process, we adopt a greedy approach of assuming 10% of the samples to be noisy. Consequently, both models need to designate a portion of the data in each batch as “clean” to update the parameters. This strategy allowed our model to prioritize learning from the clean data, leading to enhanced robustness.

Ablation Study: We embed the results of the ablation study in Table 2, which explores two key variations: alteration of the number of input modalities and alteration of the architecture of the model. To assess the impact of data modalities on classification performance, the T2-only model is compared with MpMRI and M.S. MpMRI models, both use T2 and ADC patches as input. The addition of the ADC modality leads to a substantial improvement in classification performance, highlighting the utility of multi-modal information in guiding the model for clinically significant PCa classification. The examination of model architecture reveals that the model with simpler inputs, MpMRI, performs better. Furthermore, the results can be further enhanced by leveraging the co-teaching framework.

6.3 PCa classification with filtration

In this section, we embark on experiments utilizing two different architectures, MpMRI and M.S. MpMRI, and incorporating training set filtering at various rates. The evidential focal loss described in Section 4.2.2 is used as the loss function to optimize the models. The co-teaching framework is excluded from this section for the following reason: while co-teaching *implicitly* handles noisy labels or samples in the training set, the training set filtering in Section

	Data	Acc.	Sen.	Spec.	AUC
SDNet[5](baseline)	T2	79.4	28.6	92.6	76.2±17.5
T2-only	T2	64.7	71.4	63.0	77.8±18.0
MpMRI	T2+ADC	79.4	85.7	77.8	84.7±15.5
Vol. MpMRI	T2+ADC	67.6	71.4	66.7	68.9±26.1
M.S. MpMRI	T2+ADC	73.5	71.4	74.1	82.5±14.3
MpMRI+co-teaching	T2+ADC	82.3	85.7	81.2	88.4±10.6

Table 2: **Patient-based results** of experiments using conventional training paradigm in Section 5.1. Standard deviations are computed from the 95% bootstrap confidence interval with $n = 3000$ samples. The best results are **bold**. All units of the numeric values are in %.

4.2.2 is an explicit alternative to dealing with them. The co-teaching framework will first update its model parameters with simpler and cleaner samples during training. However, through the filtering process, data samples with high uncertainty values are considered potentially noisy and hence are not involved in the training process. We argue that the use of co-teaching and simultaneous data filtering might be redundant. Our hypothesis for training set filtering is that by explicitly eliminating highly uncertain data samples from the set and optimizing only on the remaining "confident" samples using the evidential focal loss (Section 4.2.2), we can produce a more robust model. Therefore, to coalesce our proposed loss function with training set filtering, we do not use co-teaching and instead, we select MpMRI and M.S. MpMRI for experiments in this section. We use these two models to compute the uncertainty for all training data first, and then the filtering process can be done either patch-driven or patient-driven on the training set, as we discussed in Section 4.2.2.

In Table 3, we present the **patient-based results** on filtering 10% and 20% of training *patches* in the training set for the two selected models. Table 4 represents the same as Table 3, except we filter 10% and 20% of training *patients* in the training set. The results from both tables demonstrate that the MpMRI model performs better than the M.S.MpMRI model in each filtration rate, and the binary classification performance improves when filtering more uncertain data for both models. From the expected calibration error (ECE), we observed that the MpMRI model with 20% filtering on both patch- and patient-based has a lower ECE value compared to those for the M.S. MpMRI model, demonstrating the predicted output probabilities of the MpMRI model matches well with the actual probabilities of the ground truth. Comparing these results with those from Table 2, we can conclude that the dataset filtering method applied to the training set, together with the evidential focal loss we proposed, can effectively improve the classification performance.

	Data	F.R.	F.M.	Acc.	Sen.	Spec.	AUC	ECE ↓
MpMRI	T2+ADC	10%	patch	82.4	85.7	81.5	85.7±9.9	0.27
M.S.MpMRI	T2+ADC	10%	patch	82.4	71.4	85.2	83.6±13.5	0.15
MpMRI	T2+ADC	20%	patch	85.3	100	81.5	98.4±1.6	0.20
M.S. MpMRI	T2+ADC	20%	patch	85.3	71.4	88.9	92.6±7.4	0.22

Table 3: **Patient-based results** of experiments using evidential focal loss and *patch-based filtering*. Standard deviations are computed from the 95% bootstrap confidence interval with $n = 3000$ samples. The “F.R.” and “F.M.” represent the filtering rate and the filtering method, respectively. The best results for each filtration rate are **bold**. All units of the numeric values are in %.

	Data	F.R.	F.M.	Acc.	Sen.	Spec.	AUC	ECE ↓
MpMRI	T2+ADC	10%	patient	88.2	100	85.2	92.6±7.4	0.27
M.S. MpMRI	T2+ADC	10%	patient	85.3	100	81.5	86.2±9.0	0.24
MpMRI	T2+ADC	20%	patient	73.5	85.7	70.4	86.8±12.6	0.21
M.S. MpMRI	T2+ADC	20%	patient	73.5	71.4	74.1	84.6±14.2	0.22

Table 4: **Patient-based results** of experiments using evidential focal loss and *patient-based filtering*. Standard deviations are computed from the 95% bootstrap confidence interval with $n = 3000$ samples. The “F.R.” and “F.M.” represent the filtering rate and the filtering method, respectively. The best results for each filtration rate are in **bold**. All units of the numeric values are in %.

Moreover, an interesting observation reveals a gradual deterioration in performance with *patient-driven* filtering. This phenomenon may stem from the distinction between patch-driven and patient-driven filtering approaches. In patch-driven filtering, the exclusion of training patches with high uncertainty values is straightforward during the training process, irrespective of the patient to which these patches belong. Conversely, in the case of patient-driven filtering, we have to consider the average uncertainty of the 20 patches for each patient. In cases where the average uncertainty for a patient falls below the defined threshold, all patches for that patient are used in training, irrespective of whether specific patches exhibit exceptionally high uncertainty. Consequently, there exists a risk of erroneously retaining patches with high uncertainty if the average uncertainty for the corresponding patient is relatively low, potentially impacting the model’s performance. This nuanced difference contributes to the observed discrepancy in results, favoring the patch-driven filtration approach.

Ablation Study: As previously mentioned, Table 2 corresponds to experiments conducted without using evidential focal loss or filtering. On the other hand, Table 3 and 4 encompass experiments that incorporate both these elements. To solely examine the influence of our proposed loss, we conducted an experiment employing the evidential focal loss without any filtering (0%). These results are summarized in Table 5 in comparison with 20% *patch-based filtering* approach. As can be seen, even without any data filtering during the training, we could correctly classify all

patients with clinically significant PCa (sensitivity = 100%), which demonstrates a significant improvement compared to MpMRI using traditional focal loss in Table 2 and also the baseline result in [5]. As expected, the addition of data filtering further improves the results for all metrics. The original results based on image patches of Table 5 can be found in Appendix B.

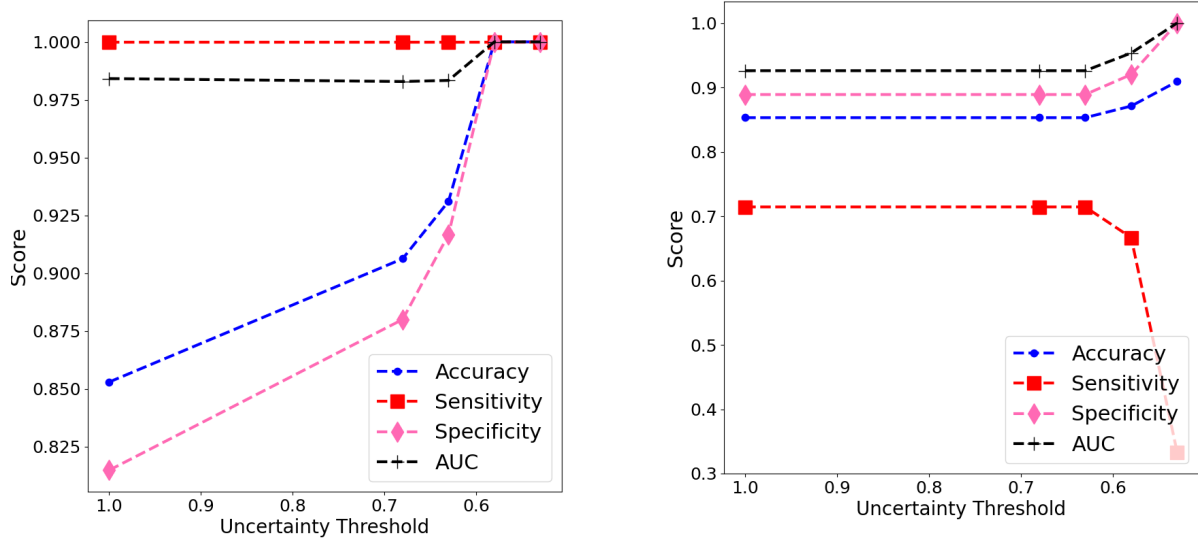
	filter 0%					filter 20%				
	Acc.	Sen.	Spec.	AUC	ECE	Acc.	Sen.	Spec.	AUC	ECE
MpMRI	82.4	100	77.8	89.4±9.1	0.29	85.3	100	81.5	98.4±1.6	0.20
M.S. MpMRI	76.5	100	70.4	82.0±12.7	0.23	85.3	71.4	88.9	92.6±7.4	0.22

Table 5: Ablation on employing proposed evidential focal loss with and without data filtering for the two selected architectures. The **Patient-based** results are reported. Standard deviations are computed from the 95% bootstrap confidence interval with $n = 3000$ samples. The best results for each filtration rate are in **bold**. All units of the numeric values are in %.

6.4 Filtering during deployment

So far, our exploration has centered on the impact of the data filtering strategy during training, aimed at enhancing model robustness and performance. It is also possible to apply filtering on the test set, i.e., when deploying the model to real clinical routines. This is equivalent to refraining from making decisions on the test samples that are identified as highly uncertain. The evaluation involves utilizing the pre-trained MpMRI and M.S. MpMRI models, each with 0% and 20% *patch-based* training filtering rates, as final models, and assessing their performance on the test set. The performance on filtering the test data based on different uncertainty thresholds using pre-trained models with 20% filtering during training is shown in Figure 4. For the performance of the other two models (0% filtering), refer to Appendix C. We started with an uncertainty threshold of 1.0 (keeping patients with $ncertainty < 1.0$); we progressively adjusted the threshold to 0.68, 0.63, 0.58, and ended at a threshold value of 0.53. These thresholds are determined based on the uncertainty value computed in the test set. We observe that the model improves its performance when filtering out highly uncertain patients from the test set, and eventually classified all patients correctly, as shown in Figure 4a.

This pragmatic approach bears relevance in real clinical settings, offering radiologists an efficient means to allocate their time, focusing on patients filtered out during the diagnostic process due to their high uncertainty values, rather than those confidently classified.



(a) MpMRI, with 20% filtering on training patches

(b) M.S. MpMRI, with 20% filtering on training patches

Fig 4: Test performance of selected models based on uncertainty threshold.

7 Conclusion

In this study, we introduced a novel approach for unpaired image-to-image translation of prostate mp-MRI and developed a robust deep-learning model for classifying clinically significant PCa using evidential focal loss. We demonstrated the effectiveness of our method on our local dataset, reinforced by a publicly available one, and showed that uncertainty-aware filtering during both training and deployment can significantly improve the PCa classification performance. The quantitative results for image translation also demonstrated that our proposed domain transfer framework using the ACL-GAN model outperforms the baseline SDNet [5]. Our approach holds promise in expediting the diagnostic process by identifying patients with high uncertainty, allowing clinicians to focus on precise diagnosis and expedite cases with high prediction certainty.

While our approach has shown promising results, there are still opportunities for improvement. One potential area for future work is to consider the spatial dependency between slices in volumetric MRI. Currently, our domain transfer framework only accepts 2D images as input and output, and we reshape the volumetric MRI into several 2D slices. However, explicitly splitting 3D images into 2D slices may eliminate the spatial dependency within each MRI data and affect the classification results. Therefore, a plausible solution could be translating the entire 3D MRI volume from 3.0T to 1.5T instead of handling individual slices.

Lastly, there is great potential for further improving the classification performance by combining more images from different MRI functional sequences, such as b-value and K^{trans} . We have already demonstrated that incorporating additional ADC images significantly enhances classification performance. We believe that if we successfully translate other images from b-value or K^{trans} acquired at 3.0T to 1.5T and incorporate them into the classification, the results could be further improved. However, the additional MRI sequences may not be available in the local 1.5T dataset. The conversion process may become feasible if we acquire those sequences from local hospitals.

Disclosures

We declare we don't have conflicts of interest.

Code, Data, and Materials Availability

ProstateX data is publicly available here ², our local KHSC data is not available due to ethical concerns. Our code is at https://github.com/med-i-lab/DT_UE_PCa, we will make the code available upon acceptance.

Acknowledgments

This work was supported by the Natural Sciences and Engineering Council of Canada, Canadian Institutes for Health Research, and Queen's University. Parvin Mousavi is supported by Canada CIFAR AI Chair and the Vector AI Institute.

References

- [1] I. Reda, A. Khalil, M. Elmogy, *et al.*, "Deep learning role in early diagnosis of prostate cancer," *Technology in cancer research & treatment* **17**, 1533034618775530 (2018).
- [2] R. P. Smith, S. B. Malkowicz, R. Whittington, *et al.*, "Identification of clinically significant prostate cancer by prostate-specific antigen screening," *Archives of internal medicine* **164**(11), 1227–1230 (2004).
- [3] M. Arif, I. G. Schoots, J. Castillo Tovar, *et al.*, "Clinically significant prostate cancer detection and segmentation in low-risk patients using a convolutional neural network on multi-parametric mri," *European radiology* **30**(12), 6582–6592 (2020).

²<https://wiki.cancerimagingarchive.net/pages/viewpage.action?pageId=23691656>

- [4] R. H. Fletcher, “Guideline: Experts recommend against prostate cancer screening with prostate-specific antigen test,” *Annals of Internal Medicine* **170**(2), JC2 (2019). PMID: 30641553.
- [5] A. Grebenisan, A. Sedghi, J. Izard, *et al.*, “Spatial decomposition for robust domain adaptation in prostate cancer detection,” in *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*, 1218–1222, IEEE (2021).
- [6] J. O. Barentsz, J. Richenberg, R. Clements, *et al.*, “Esur prostate mr guidelines 2012,” *European radiology* **22**(4), 746–757 (2012).
- [7] J. C. Weinreb, J. O. Barentsz, P. L. Choyke, *et al.*, “Pi-rads prostate imaging–reporting and data system: 2015, version 2,” *European urology* **69**(1), 16–40 (2016).
- [8] J. O. Barentsz, J. C. Weinreb, S. Verma, *et al.*, “Synopsis of the pi-rads v2 guidelines for multiparametric prostate magnetic resonance imaging and recommendations for use,” *European urology* **69**(1), 41 (2016).
- [9] A. E.-S. Bosaily, C. Parker, L. Brown, *et al.*, “Promis—prostate mr imaging study: a paired validating cohort study evaluating the role of multi-parametric mri in men with clinical suspicion of prostate cancer,” *Contemporary clinical trials* **42**, 26–40 (2015).
- [10] A. Stabile, F. Giganti, A. B. Rosenkrantz, *et al.*, “Multiparametric mri for prostate cancer diagnosis: current status and future directions,” *Nature reviews urology* **17**(1), 41–61 (2020).
- [11] H. U. Ahmed, A. E.-S. Bosaily, L. C. Brown, *et al.*, “Diagnostic accuracy of multi-parametric mri and trus biopsy in prostate cancer (promis): a paired validating confirmatory study,” *The Lancet* **389**(10071), 815–822 (2017).
- [12] A. Saha, M. Hosseinzadeh, and H. Huisman, “End-to-end prostate cancer detection in bpmri via 3d cnns: Effects of attention mechanisms, clinical priori and decoupled false positive reduction,” *Medical image analysis* **73**, 102155 (2021).
- [13] M. H. Le, J. Chen, L. Wang, *et al.*, “Automated diagnosis of prostate cancer in multi-parametric mri based on multimodal convolutional neural networks,” *Physics in Medicine & Biology* **62**(16), 6497 (2017).

- [14] S. Yoo, I. Gujrathi, M. A. Haider, *et al.*, “Prostate cancer detection using deep convolutional neural networks,” *Scientific reports* **9**(1), 1–10 (2019).
- [15] S. Iqbal, G. F. Siddiqui, A. Rehman, *et al.*, “Prostate cancer detection using deep learning and traditional techniques,” *IEEE Access* **9**, 27085–27100 (2021).
- [16] O. J. Pellicer-Valero, J. L. Marengo Jiménez, V. Gonzalez-Perez, *et al.*, “Deep learning for fully automatic detection, segmentation, and gleason grade estimation of prostate cancer in multiparametric magnetic resonance images,” *Scientific reports* **12**(1), 1–13 (2022).
- [17] S. G. Armato, H. Huisman, K. Drukker, *et al.*, “Prostatex challenges for computerized classification of prostate lesions from multiparametric magnetic resonance images,” *Journal of Medical Imaging* **5**(4), 044501 (2018).
- [18] G. Litjens, O. Debats, J. Barentsz, *et al.*, “Computer-aided detection of prostate cancer in mri,” *IEEE transactions on medical imaging* **33**(5), 1083–1092 (2014).
- [19] S. Liu, H. Zheng, Y. Feng, *et al.*, “Prostate cancer diagnosis using deep learning with 3d multiparametric mri,” in *Medical imaging 2017: computer-aided diagnosis*, **10134**, 581–584, SPIE (2017).
- [20] A. Mehrtash, A. Sedghi, M. Ghafoorian, *et al.*, “Classification of clinical significance of mri prostate findings using 3d convolutional neural networks,” in *Medical Imaging 2017: Computer-Aided Diagnosis*, **10134**, 101342A, International Society for Optics and Photonics (2017).
- [21] A. Grebenisan, A. Sedghi, A. Menard, *et al.*, “Towards democratizing ai in mr-based prostate cancer diagnosis: 3.0 to 1.5 tesla,” in *Medical Imaging 2020: Image-Guided Procedures, Robotic Interventions, and Modeling*, **11315**, 184–189, SPIE (2020).
- [22] M. E. Ladd, P. Bachert, M. Meyerspeer, *et al.*, “Pros and cons of ultra-high-field mri/mrs for human application,” *Progress in nuclear magnetic resonance spectroscopy* **109**, 1–50 (2018).
- [23] L. C. Adams, M. R. Makowski, G. Engel, *et al.*, “Dataset of prostate mri annotated for anatomical zones and cancer,” *Data in Brief* **45**, 108739 (2022).

- [24] L. C. Adams, M. R. Makowski, G. Engel, *et al.*, “Prostate158-an expert-annotated 3t mri dataset and algorithm for prostate cancer detection,” *Computers in Biology and Medicine* **148**, 105817 (2022).
- [25] T. Hulsen, “An overview of publicly available patient-centered prostate cancer datasets,” *Translational andrology and urology* **8**(Suppl 1), S64 (2019).
- [26] A. I. Mushlin, C. Mooney, R. G. Holloway, *et al.*, “The cost-effectiveness of magnetic resonance imaging for patients with equivocal neurological symptoms,” *International Journal of Technology Assessment in Health Care* **13**(1), 21–34 (1997).
- [27] B. Turkbey, A. B. Rosenkrantz, M. A. Haider, *et al.*, “Prostate imaging reporting and data system version 2.1: 2019 update of prostate imaging reporting and data system version 2,” *European Urology* **76**(3), 340–351 (2019).
- [28] Z. Ryznarova, M. Jiru, V. Vik, *et al.*, “Comparison of 1.5 t and 3t prostate mr examination using surface array coils in routine clinical practice,” *J Diagnostic Tech Biomed Anal* **7**(2) (2018).
- [29] M. Virarkar, J. Szklaruk, R. Diab, *et al.*, “Diagnostic value of 3.0 t versus 1.5 t mri in staging prostate cancer: systematic review and meta-analysis,” *Polish Journal of Radiology* **87**(1), 421–429 (2022).
- [30] A. Woernle, C. Englman, L. Dickinson, *et al.*, “Picture perfect: the status of image quality in prostate mri,” *Journal of Magnetic Resonance Imaging* (2023).
- [31] Canadian Agency for Drugs and Technologies in Health (CADTH), “Average volume of mri exams conducted per hour across canada,” *Canadian Journal of Health Technologies* **4**(1) (2024).
- [32] X. Li, Y. Gu, N. Dvornek, *et al.*, “Multi-site fmri analysis using privacy-preserving federated learning and domain adaptation: Abide results,” *Medical Image Analysis* **65**, 101765 (2020).
- [33] M. Adnan, S. Kalra, J. C. Cresswell, *et al.*, “Federated learning and differential privacy for medical image analysis,” *Scientific reports* **12**(1), 1953 (2022).
- [34] M. Sensoy, L. Kaplan, and M. Kandemir, “Evidential deep learning to quantify classification uncertainty,” *Advances in Neural Information Processing Systems* **31** (2018).

- [35] B. Han, Q. Yao, X. Yu, *et al.*, “Co-teaching: Robust training of deep neural networks with extremely noisy labels,” *Advances in neural information processing systems* **31** (2018).
- [36] T.-Y. Lin, P. Goyal, R. Girshick, *et al.*, “Focal loss for dense object detection,” in *Proceedings of the IEEE international conference on computer vision*, 2980–2988 (2017).
- [37] S. Cui, X. Jin, S. Wang, *et al.*, “Heuristic domain adaptation,” *Advances in Neural Information Processing Systems* **33**, 7571–7583 (2020).
- [38] Y. Ganin and V. Lempitsky, “Unsupervised domain adaptation by backpropagation,” in *International conference on machine learning*, 1180–1189, PMLR (2015).
- [39] Y. Ganin, E. Ustinova, H. Ajakan, *et al.*, “Domain-adversarial training of neural networks,” *The journal of machine learning research* **17**(1), 2096–2030 (2016).
- [40] M. Long, H. Zhu, J. Wang, *et al.*, “Unsupervised domain adaptation with residual transfer networks,” *Advances in neural information processing systems* **29** (2016).
- [41] K. Saito, K. Watanabe, Y. Ushiku, *et al.*, “Maximum classifier discrepancy for unsupervised domain adaptation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 3723–3732 (2018).
- [42] M. Long, J. Wang, G. Ding, *et al.*, “Transfer joint matching for unsupervised domain adaptation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 1410–1417 (2014).
- [43] T. Yao, Y. Pan, C.-W. Ngo, *et al.*, “Semi-supervised domain adaptation with subspace learning for visual recognition,” in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2142–2150 (2015).
- [44] K. Saito, D. Kim, S. Sclaroff, *et al.*, “Semi-supervised domain adaptation via minimax entropy,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 8050–8058 (2019).
- [45] B. Li, Y. Wang, S. Zhang, *et al.*, “Learning invariant representations and risks for semi-supervised domain adaptation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 1104–1113 (2021).

- [46] J. Hoffman, B. Kulis, T. Darrell, *et al.*, “Discovering latent domains for multisource domain adaptation,” in *European Conference on Computer Vision*, 702–715, Springer (2012).
- [47] R. Xu, Z. Chen, W. Zuo, *et al.*, “Deep cocktail network: Multi-source unsupervised domain adaptation with category shift,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 3964–3973 (2018).
- [48] X. Peng, Q. Bai, X. Xia, *et al.*, “Moment matching for multi-source domain adaptation,” in *Proceedings of the IEEE/CVF international conference on computer vision*, 1406–1415 (2019).
- [49] M. Wang and W. Deng, “Deep visual domain adaptation: A survey,” *Neurocomputing* **312**, 135–153 (2018).
- [50] M. Long, Y. Cao, J. Wang, *et al.*, “Learning transferable features with deep adaptation networks,” in *International conference on machine learning*, 97–105, PMLR (2015).
- [51] H. Yan, Y. Ding, P. Li, *et al.*, “Mind the class weight bias: Weighted maximum mean discrepancy for unsupervised domain adaptation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2272–2281 (2017).
- [52] A. Kumagai and T. Iwata, “Unsupervised domain adaptation by matching distributions based on the maximum mean discrepancy via unilateral transformations,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, **33**(01), 4106–4113 (2019).
- [53] E. Tzeng, J. Hoffman, T. Darrell, *et al.*, “Simultaneous deep transfer across domains and tasks,” in *Proceedings of the IEEE international conference on computer vision*, 4068–4076 (2015).
- [54] G. Hinton, O. Vinyals, J. Dean, *et al.*, “Distilling the knowledge in a neural network,” *arXiv preprint arXiv:1503.02531* **2**(7) (2015).
- [55] S. Motiian, M. Piccirilli, D. A. Adjeroh, *et al.*, “Unified deep supervised domain adaptation and generalization,” in *Proceedings of the IEEE international conference on computer vision*, 5715–5725 (2017).

- [56] K. Bousmalis, N. Silberman, D. Dohan, *et al.*, “Unsupervised pixel-level domain adaptation with generative adversarial networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 3722–3731 (2017).
- [57] W. Hong, Z. Wang, M. Yang, *et al.*, “Conditional generative adversarial network for structured domain adaptation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 1335–1344 (2018).
- [58] T. Karras, M. Aittala, J. Hellsten, *et al.*, “Training generative adversarial networks with limited data,” *Advances in Neural Information Processing Systems* **33**, 12104–12114 (2020).
- [59] J. Hoffman, E. Tzeng, T. Park, *et al.*, “Cycada: Cycle-consistent adversarial domain adaptation,” in *International conference on machine learning*, 1989–1998, PMLR (2018).
- [60] J.-Y. Zhu, T. Park, P. Isola, *et al.*, “Unpaired image-to-image translation using cycle-consistent adversarial networks,” in *Proceedings of the IEEE international conference on computer vision*, 2223–2232 (2017).
- [61] Y. Zhao, R. Wu, and H. Dong, “Unpaired image-to-image translation using adversarial consistency loss,” in *European Conference on Computer Vision*, 800–815, Springer (2020).
- [62] W. M. Kouw, M. Loog, L. W. Bartels, *et al.*, “Mr acquisition-invariant representation learning,” *arXiv preprint arXiv:1709.07944* (2017).
- [63] N. M. Khan, N. Abraham, and M. Hon, “Transfer learning with intelligent training data selection for prediction of alzheimer’s disease,” *IEEE Access* **7**, 72726–72735 (2019).
- [64] J. Deng, W. Dong, R. Socher, *et al.*, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE conference on computer vision and pattern recognition*, 248–255, IEEE (2009).
- [65] M. Ghafoorian, A. Mehrtash, T. Kapur, *et al.*, “Transfer learning for domain adaptation in mri: Application in brain lesion segmentation,” in *International conference on medical image computing and computer-assisted intervention*, 516–524, Springer (2017).
- [66] K. Armanious, C. Jiang, M. Fischer, *et al.*, “Medgan: Medical image translation using gans,” *Computerized medical imaging and graphics* **79**, 101684 (2020).

- [67] Y. Hiasa, Y. Otake, M. Takao, *et al.*, “Cross-modality image synthesis from unpaired data using cyclegan,” in *International workshop on simulation and synthesis in medical imaging*, 31–41, Springer (2018).
- [68] D. Nie, R. Trullo, J. Lian, *et al.*, “Medical image synthesis with context-aware generative adversarial networks,” in *International conference on medical image computing and computer-assisted intervention*, 417–425, Springer (2017).
- [69] R. Oulbacha and S. Kadoury, “Mri to ct synthesis of the lumbar spine from a pseudo-3d cycle gan,” in *2020 IEEE 17th international symposium on biomedical imaging (ISBI)*, 1784–1787, IEEE (2020).
- [70] K. Armanious, C. Jiang, S. Abdulatif, *et al.*, “Unsupervised medical image translation using cycle-medgan,” in *2019 27th European Signal Processing Conference (EUSIPCO)*, 1–5, IEEE (2019).
- [71] X. Ying, H. Guo, K. Ma, *et al.*, “X2ct-gan: reconstructing ct from biplanar x-rays with generative adversarial networks,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10619–10628 (2019).
- [72] R. Ge, Y. He, C. Xia, *et al.*, “X-ctrnet: 3d cervical vertebra ct reconstruction and segmentation directly from 2d x-ray images,” *Knowledge-Based Systems* **236**, 107680 (2022).
- [73] X. Hu, “Multi-texture gan: Exploring the multi-scale texture translation for brain mr images,” *arXiv preprint arXiv:2102.07225* (2021).
- [74] H. Uzunova, J. Ehrhardt, and H. Handels, “Memory-efficient gan-based domain translation of high resolution 3d medical images,” *Computerized Medical Imaging and Graphics* **86**, 101801 (2020).
- [75] A. Radford, L. Metz, and S. Chintala, “Unsupervised representation learning with deep convolutional generative adversarial networks,” *arXiv preprint arXiv:1511.06434* (2015).
- [76] M. Frid-Adar, I. Diamant, E. Klang, *et al.*, “Gan-based synthetic medical image augmentation for increased cnn performance in liver lesion classification,” *Neurocomputing* **321**, 321–331 (2018).
- [77] G. Huang and A. H. Jafari, “Enhanced balancing gan: Minority-class image generation,” *Neural Computing and Applications*, 1–10 (2021).

- [78] G. Kwon, C. Han, and D.-s. Kim, “Generation of 3d brain mri using auto-encoding generative adversarial networks,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 118–126, Springer (2019).
- [79] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556* (2014).
- [80] K. He, X. Zhang, S. Ren, *et al.*, “Identity mappings in deep residual networks,” in *European conference on computer vision*, 630–645, Springer (2016).
- [81] D. Karimi, H. Dou, S. K. Warfield, *et al.*, “Deep learning with noisy labels: Exploring techniques and remedies in medical image analysis,” *Medical Image Analysis* **65**, 101759 (2020).
- [82] D. Arpit, S. Jastrzebski, N. Ballas, *et al.*, “A closer look at memorization in deep networks,” in *International conference on machine learning*, 233–242, PMLR (2017).
- [83] C. Zhang, S. Bengio, M. Hardt, *et al.*, “Understanding deep learning (still) requires rethinking generalization,” *Communications of the ACM* **64**(3), 107–115 (2021).
- [84] H. Song, M. Kim, D. Park, *et al.*, “Learning from noisy labels with deep neural networks: A survey,” *IEEE Transactions on Neural Networks and Learning Systems* (2022).
- [85] B. Van Rooyen, A. Menon, and R. C. Williamson, “Learning with symmetric label noise: The importance of being unhinged,” *Advances in neural information processing systems* **28** (2015).
- [86] N. Charoenphakdee, J. Lee, and M. Sugiyama, “On symmetric losses for learning from corrupted labels,” in *International Conference on Machine Learning*, 961–970, PMLR (2019).
- [87] Z. Zhang and M. Sabuncu, “Generalized cross entropy loss for training deep neural networks with noisy labels,” *Advances in neural information processing systems* **31** (2018).
- [88] L. Jiang, Z. Zhou, T. Leung, *et al.*, “Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels,” in *International conference on machine learning*, 2304–2313, PMLR (2018).

- [89] E. Malach and S. Shalev-Shwartz, “Decoupling” when to update” from” how to update”,” *Advances in neural information processing systems* **30** (2017).
- [90] Z. Wang, G. Hu, and Q. Hu, “Training noise-robust deep neural networks via meta-learning,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 4524–4533 (2020).
- [91] X. Yu, B. Han, J. Yao, *et al.*, “How does disagreement help generalization against label corruption?,” in *International Conference on Machine Learning*, 7164–7173, PMLR (2019).
- [92] Y. Gal and Z. Ghahramani, “Dropout as a bayesian approximation: Representing model uncertainty in deep learning,” in *international conference on machine learning*, 1050–1059, PMLR (2016).
- [93] Y. Gal and Z. Ghahramani, “Bayesian convolutional neural networks with bernoulli approximate variational inference,” *arXiv preprint arXiv:1506.02158* (2015).
- [94] A. Jøsang, *Subjective logic*, vol. 3, Springer (2016).
- [95] A. P. Dempster, “A generalization of bayesian inference,” *Journal of the Royal Statistical Society: Series B (Methodological)* **30**(2), 205–232 (1968).
- [96] F. C. Ghesu, B. Georgescu, E. Gibson, *et al.*, “Quantifying and leveraging classification uncertainty for chest radiograph assessment,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 676–684, Springer (2019).
- [97] F. C. Ghesu, B. Georgescu, A. Mansoor, *et al.*, “Quantifying and leveraging predictive uncertainty for medical image assessment,” *Medical Image Analysis* **68**, 101855 (2021).
- [98] M. Tardy, B. Scheffer, and D. Mateus, “Uncertainty measurements for the reliable classification of mammograms,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 495–503, Springer (2019).
- [99] B. Yuan, X. Yue, Y. Lv, *et al.*, “Evidential deep neural networks for uncertain data classification,” in *International Conference on Knowledge Science, Engineering and Management*, 427–437, Springer (2020).

- [100] G. Litjens, O. Debats, J. Barentsz, *et al.*, “Prostatex challenge data,” *Cancer Imaging Arch* **10**, K9TCIA (2017).
- [101] V. Kasivisvanathan, A. S. Rannikko, M. Borghi, *et al.*, “Mri-targeted or standard biopsy for prostate-cancer diagnosis,” *New England Journal of Medicine* **378**(19), 1767–1777 (2018).
- [102] M. Kasson, M. Ortman, K. Gaitonde, *et al.*, “Imaging prostate cancer using multiparametric magnetic resonance imaging: past, present, and future,” in *Seminars in Roentgenology*, **53**(3), 200–205, Elsevier (2018).
- [103] M. Heusel, H. Ramsauer, T. Unterthiner, *et al.*, “Gans trained by a two time-scale update rule converge to a local nash equilibrium,” *Advances in neural information processing systems* **30** (2017).
- [104] A. Gretton, K. M. Borgwardt, M. J. Rasch, *et al.*, “A kernel two-sample test,” *The Journal of Machine Learning Research* **13**(1), 723–773 (2012).
- [105] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980* (2014).
- [106] X. Mao, Q. Li, H. Xie, *et al.*, “Least squares generative adversarial networks,” in *Proceedings of the IEEE international conference on computer vision*, 2794–2802 (2017).
- [107] A. Volokitin, E. Erdil, N. Karani, *et al.*, “Modelling the distribution of 3d brain mri using a 2d slice vae,” in *Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part VII 23*, 657–666, Springer (2020).
- [108] M. Zhou and F. Khalvati, “Conditional generation of 3d brain tumor regions via vqgan and temporal-agnostic masked transformer,” in *Medical Imaging with Deep Learning*, (2024).
- [109] M. Gilany, P. Wilson, A. Jamzad, *et al.*, “Towards confident detection of prostate cancer using high resolution micro-ultrasound,” in *International conference on medical image computing and computer-assisted intervention*, 411–420, Springer (2022).
- [110] A. Mehrtash, W. M. Wells, C. M. Tempany, *et al.*, “Confidence calibration and predictive uncertainty estimation for deep medical image segmentation,” *IEEE transactions on medical imaging* **39**(12), 3868–3878 (2020).

Appendix A: Experimental Details

In this section, we provide details of hyperparameter settings to train uncertainty-aware classification models described in Section 5.1.

To train the ‘‘MpMRI’’ model for patch-driven filtering, we set the learning rate to 0.0001; weight decay to 0.01; total training epochs to 300, and batch size to 10. The class weights β in Equation (10) are set to [0.25, 0.75] for filtering 10%, and [0.25, 1.25] for filtering 20% of the training data. For patient-driven filtering, all parameters are the same except the class weights β are set to [0.25, 1] for filtering 10% of the training data. Last but not least, we set the initial learning rate to 0.0001; total training epochs to 300; batch size to 10; the learning rate decayed by a factor of 0.1 for every 200 epochs, and the class weights β are set to [0.25, 1] for filtering 20% of the training data.

To train the ‘‘M.S. MpMRI’’ model for patch-driven filtering, we set the initial learning rate to 0.0001; weight decay to 0.01; total training epochs to 300. The class weights β in Equation (10) are set to [0.25, 1] for both filtering 10% and 20% of the training data. For patient-driven filtering, all parameters were the same except the class weights β are set to [0.25, 1] for filtering 10% and [0.25, 1.25] for filtering 20% of the training data.

Appendix B: Original Results

In this section, we provide all original results based on image patches (Patch-based results) of all experiments we performed in the paper. We start by presenting Table 6, which is the patch-based result for Table 2 in Section 6.2.

	Data	Acc.	Sen.	Spec.	AUC
SDNet [5](baseline)	T2	77.8	27.9	90.7	74.8±4.1
T2-only	T2	65.6	75.7	63.0	77.9±3.9
MpMRI	T2+ADC	75.1	69.3	76.7	79.9±3.8
Vol. MpMRI	T2+ADC	64.9	62.9	65.4	63.4±5.5
M.S. MpMRI	T2+ADC	70.4	74.3	69.4	80.1±3.6
MpMRI+co-teaching	T2+ADC	78.4	72.1	80.0	80.7±3.9

Table 6: **Patch-based results** of experiments using conventional training paradigm in Section 5.1. Standard deviations are computed from the 95% bootstrap confidence interval with $n = 3000$ samples. The best results are **bold**. All units of the numeric values are in %.

Then, we provide the patch-based results for patch-based filtering in Table 7, which corresponds to Table 3 in Section 6.3.

	Data	F.R.	F.M.	Acc.	Sen.	Spec.	AUC	ECE ↓
MpMRI	T2+ADC	10%	patch	77.5	80.7	76.7	82.4±3.4	0.22
M.S.MpMRI	T2+ADC	10%	patch	79.7	70.7	82.0	81.1±3.8	0.16
MpMRI	T2+ADC	20%	patch	83.8	80.0	84.8	89.7±2.6	0.16
M.S. MpMRI	T2+ADC	20%	patch	82.4	73.6	84.6	87.6±2.8	0.17

Table 7: **Patch-based results** of experiments using evidential focal loss and *patch-based filtering*. Standard deviations are computed from the 95% bootstrap confidence interval with $n = 3000$ samples. The “F.R.” and “F.M.” represent the filtering rate and the filtering method, respectively. The best results for each filtration rate are **bold**. All units of the numeric values are in %.

For patient-based filtering, we provide the patch-based results in Table 8, which corresponds to Table 4 in Section

6.3.

	Data	F.R.	F.M.	Acc.	Sen.	Spec.	AUC	ECE ↓
MpMRI	T2+ADC	10%	patient	75.3	80.7	73.9	86.1±2.9	0.21
M.S. MpMRI	T2+ADC	10%	patient	75.6	75.6	76.9	82.4±3.6	0.17
MpMRI	T2+ADC	20%	patient	72.5	74.3	72.0	81.5±3.6	0.22
M.S. MpMRI	T2+ADC	20%	patient	76.3	74.3	76.9	80.1±4.2	0.18

Table 8: **Patch-based results** of experiments using evidential focal loss and *patient-based filtering*. Standard deviations are computed from the 95% bootstrap confidence interval with $n = 3000$ samples. The “F.R.” and “F.M.” represent the filtering rate and the filtering method, respectively. The best results for each filtration rate are in **bold**. All units of the numeric values are in %.

Table 9 shows the original patch-based results based on image patches for the ablation study conducted in Section 6.3.

	filter 0%					filter 20%				
	Acc.	Sen.	Spec.	AUC	ECE	Acc.	Sen.	Spec.	AUC	ECE
MpMRI	74.1	79.3	72.8	84.8±2.9	0.21	83.8	80.0	84.8	89.7±2.6	0.16
M.S. MpMRI	73.7	86.4	70.4	80.4±3.1	0.22	82.4	73.6	84.6	87.6±2.8	0.17

Table 9: Patch-based results for disable filtration on the training set for two models.

we also provide the visualization of patch-based AUC curves for the selected experiments in Section 6.2 and 6.3, along with the 95% confidence interval against the baseline model in Figure 5.

Appendix C: Filtration while deploying

Next, we provide the performance for test set filtering using pre-trained MpMRI and M.S. MpMRI with 0% filtering rate on the training set in Figure 6.

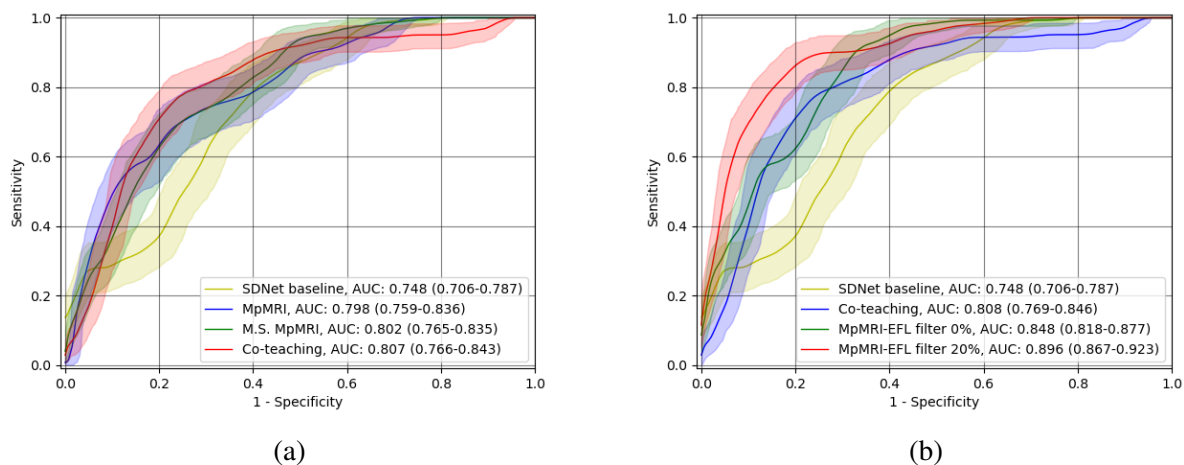


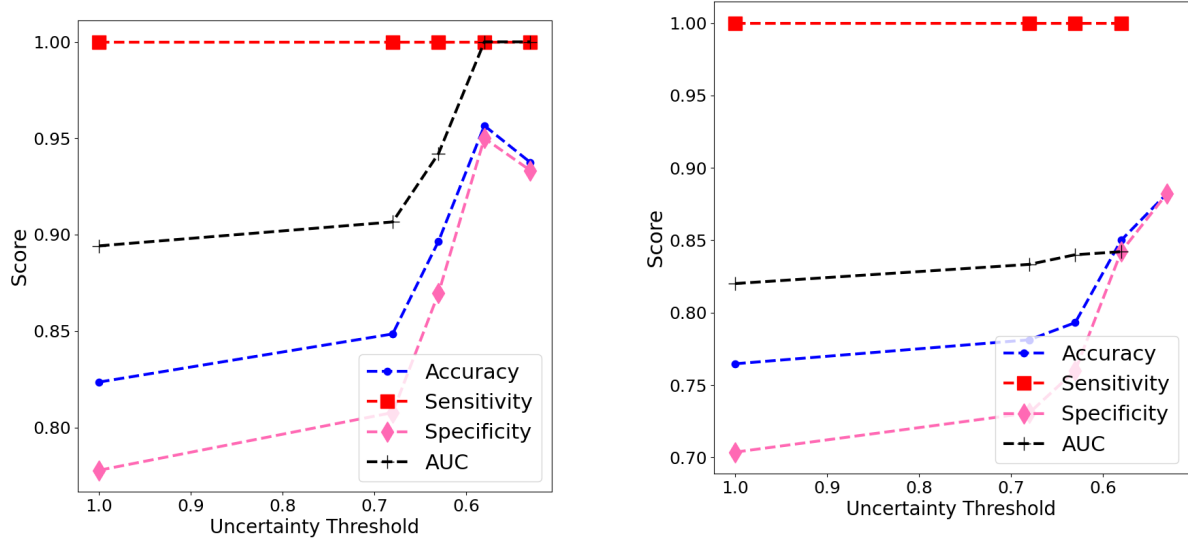
Fig 5: Both figures demonstrate the original AUC results. 5a shows the comparison of AUC curves between the baseline and the models without filtration (experiments in the first category); 5b shows the comparison of AUC curves between the baseline, the best model without filtration, and the best model with filtration on the training set. “EFL” is short for Evidential Focal Loss. The shaded areas in both figures represent the 95% confidence intervals (CI) of each model. CIs are obtained by using Bootstrap with $n = 3000$.

First Author insert biography

Biographies and photographs of the other authors are not available.

List of Figures

- 1 Visualization of sample data. 1a and 1b are the 1.5T T2 and ADC images from KHSC, respectively. Similarly, 1c and 1d are the 3.0T T2 and ADC images from the “ProstateX” Challenge, respectively.
- 2 Detailed schematic of the proposed method. The overall framework of our proposed method contains two stages: 1), domain translation to map public 3.0T MRI with local 1.5T MRI; 2), uncertainty-aware clinically significant PCa classification. The bottom figure is the training schema for domain transfer. The upper right portion of the figure illustrates the PCa classification training process, which involves training the classifier using the Evidential Focal loss, filtering the training set based on uncertainty, and retraining the classifier on the filtered data to obtain the final classifier.



(a) MpmMRI, with 0% filtering on training patches (b) M.S. MpmMRI, with 0% filtering on training patches

Fig 6: Test performance of selected models based on uncertainty threshold.

- 3 Detailed architecture of “M.S. MpmMRI” model. The first sequence of CNN layers contains $1 \times 3D$ convolution layer and $4 \times 2D$ convolution layers, $2 \times$ Max Pooling layers with window size 2×2 . Both extracted feature maps of T2 and ADC are concatenated channel-wise. After that, another set of convolution-max pooling layers is utilized. Finally, the extracted 2D features are reshaped to 1D and fed into a Fully connected layer follow by a softmax layer with 2 outputs representing the probabilities of which class the input data belongs to.
- 4 Test performance of selected models based on uncertainty threshold.
- 5 Both figures demonstrate the original AUC results. 5a shows the comparison of AUC curves between the baseline and the models without filtration (experiments in the first category); 5b shows the comparison of AUC curves between the baseline, the best model without filtration, and the best model with filtration on the training set. “EFL” is short for Evidential Focal Loss. The shaded areas in both figures represent the 95% confidence intervals (CI) of each model. CIs are obtained by using Bootstrap with $n = 3000$.
- 6 Test performance of selected models based on uncertainty threshold.

List of Tables

- 1 Quantitative results of translated T2 images using the baseline method SDNet and ACL-GAN in our proposed domain transfer framework. Lower values indicate better performance for all metrics. **Bold** values represent the best results.
- 2 **Patient-based results** of experiments using conventional training paradigm in Section 5.1. Standard deviations are computed from the 95% bootstrap confidence interval with $n = 3000$ samples. The best results are **bold**. All units of the numeric values are in %.
- 3 **Patient-based results** of experiments using evidential focal loss and *patch-based filtering*. Standard deviations are computed from the 95% bootstrap confidence interval with $n = 3000$ samples. The “F.R.” and “F.M.” represent the filtering rate and the filtering method, respectively. The best results for each filtration rate are **bold**. All units of the numeric values are in %.
- 4 **Patient-based results** of experiments using evidential focal loss and *patient-based filtering*. Standard deviations are computed from the 95% bootstrap confidence interval with $n = 3000$ samples. The “F.R.” and “F.M.” represent the filtering rate and the filtering method, respectively. The best results for each filtration rate are in **bold**. All units of the numeric values are in %.
- 5 Ablation on employing proposed evidential focal loss with and without data filtering for the two selected architectures. The **Patient-based** results are reported. Standard deviations are computed from the 95% bootstrap confidence interval with $n = 3000$ samples. The best results for each filtration rate are in **bold**. All units of the numeric values are in %.
- 6 **Patch-based results** of experiments using conventional training paradigm in Section 5.1. Standard deviations are computed from the 95% bootstrap confidence interval with $n = 3000$ samples. The best results are **bold**. All units of the numeric values are in %.
- 7 **Patch-based results** of experiments using evidential focal loss and *patch-based filtering*. Standard deviations are computed from the 95% bootstrap confidence interval with $n = 3000$ samples. The “F.R.” and “F.M.” represent the filtering rate and the filtering method, respectively. The best results for each filtration rate are **bold**. All units of the numeric values are in %.

- 8 **Patch-based results** of experiments using evidential focal loss and *patient-based filtering*. Standard deviations are computed from the 95% bootstrap confidence interval with $n = 3000$ samples. The “F.R.” and “F.M.” represent the filtering rate and the filtering method, respectively. The best results for each filtration rate are in **bold**. All units of the numeric values are in %.
- 9 Patch-based results for disable filtration on the training set for two models.