
Self-Supervised Masked Autoencoders for Prostate Cancer Segmentation via Learned Representations

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Prostate cancer incidence is rising at an estimated rate of 6.7% annually. MRI
2 guided biopsy of suspected lesions can improve diagnosis but requires precise
3 delineation of the prostate gland and lesions. Manually delineating lesions is labor
4 intensive and remains challenging due to their morphological similarity to surround-
5 ing healthy tissue. In this study, we propose a self-supervised learning framework
6 to learn rich anatomical and pathological representations from bi-parametric MRI
7 (bpMRI) scans. Specifically, we leverage a Masked Autoencoder (MAE) architec-
8 ture trained on multimodal bpMRI to capture context-aware features of the prostate
9 region through a novel lesion-aware masking strategy. The pretrained encoder is
10 then fine-tuned for lesion segmentation using a lightweight decoder augmented
11 with skip connections from the MAE. Our fine-tuning strategy incorporates a bal-
12 anced dataset and a hybrid loss function to address class imbalance. The proposed
13 approach outperforms state-of-the-art segmentation methods, achieving over a 10%
14 improvement in Dice score on the held-out test set.

15 1 Introduction

16 Prostate cancer (PCa) is the second most commonly diagnosed cancer in men globally, with over 1.2
17 million new cases annually and more than 350,000 deaths each year. In the United States alone, the
18 American Cancer Society estimates approximately 313,780 new cases and 35,770 deaths in 2025,
19 accounting for roughly 15.4% [1], of all new cancer diagnoses in men .

20 Bi-parametric MRI (bpMRI), which includes T2-weighted (T2W) and diffusion-weighted imaging
21 (DWI), has become an essential tool in the prostate cancer (PCa) diagnostic pathway. bpMRI
22 improves lesion detection within the prostate gland, enables targeted biopsies, and is widely regarded
23 as the preferred non-invasive approach for PCa risk stratification.

24 Accurate segmentation of the prostate and suspicious lesions can guide targeted biopsies and assist
25 in automated lesion grading, thereby improving diagnostic precision. At present, these tasks rely
26 heavily on the expertise and experience of clinicians, contributing to high workload and leading to
27 variability in outcomes.

28 Lesion segmentation in prostate MRI is an inherently challenging task due to combination of
29 clinical, anatomical, and technical factors. Accurate detection of both clinically significant and
30 non-significant lesions is essential to improve diagnostic performance, yet several barriers persist.
31 Subtle intensity differences between tumors and surrounding healthy tissue make lesion boundaries
32 difficult to delineate, while anatomical complexity further complicates segmentation, as tumors
33 often exhibit irregular shapes, ill-defined borders, and multifocal patterns dispersed across the gland.
34 In addition, inter-reader variability, severe class imbalance, and variability across scanners and
35 acquisition protocols amplify these challenges. Collectively, these factors highlight the need for

36 robust, automated methods capable of learning discriminative representations from MRI data to
 37 achieve reliable lesion segmentation.

38 In view of this, the following study introduces a self-supervised (SSL) approach towards learning
 39 pathology specific representations, we introduce a Masked Autoencoder (MAE) [2], architecture
 40 on bpMRIs of the prostate. The multimodal approach allows the model to learn joint latent
 41 representations that capture both structural anatomy and pathological patterns.
 42

43 Below, we highlight the key contributions and generalizable insights of our approach:

- 44 • **Lesion-aware masking in self-supervised pretraining encourages clinical relevance**
 45 **in learned representations:** By prioritizing the retention of lesion-containing regions
 46 during masked autoencoding, the model learns to encode semantically meaningful features
 47 around pathologies—without explicit labels. This approach can generalize to other sparse
 48 abnormalities.
- 49 • **An effective way to learn from limited labelled data:** Our approach introduces an encoder-
 50 heavy architecture that leverages unlabeled data through self-supervised learning to capture
 51 rich, domain-specific representations. To adapt these representations to downstream tasks,
 52 we pair the encoder with a lightweight decoder tailored for efficient fine-tuning. This design
 53 enables rapid generalization to multiple medical imaging tasks, particularly in settings with
 54 scarce annotated data.
- 55 • **Balanced data construction addresses dual imbalance in medical tasks:** Our approach
 56 introduces a tailored data sampling strategy that enables the encoder to learn from the
 57 full, naturally distributed dataset during pretraining, while the fine-tuning stage employs a
 58 balanced subset to address the dual challenge of underrepresented positive samples (lesion-
 59 containing slices) and spatial sparsity of lesions within each image. This design ensures
 60 robust representation learning and improved segmentation performance on clinically relevant
 61 regions.

62 2 Methods

63 This section outlines our two-phase training pipeline (Figure 1), Phase 1 involves pretraining a
 64 custom MAE on 4,636 volumetric bpMRI exams, the details for the dataset can be found in Appendix
 65 A.1. The MAE is trained to reconstruct masked portions of the input using a **lesion-aware masking**
 66 **strategy**, which prioritizes the retention of lesion-containing patches to enforce richer feature learning
 67 in clinically relevant regions. The encoder learns high-level representations of the prostate anatomy
 68 and potential lesion patterns without requiring ground-truth annotations, (see Appendix A.2).

69 In Phase 2, we fine-tune the pretrained encoder for supervised segmentation. A lightweight decoder
 70 head is appended to the encoder to form a complete prediction network. The segmentation model
 71 is optimized using a hybrid loss function that combines focal loss and Dice loss to handle class
 72 imbalance and optimize overlap-based accuracy. The overall architecture is tailored to the spatial
 73 resolution and anatomical structure of the prostate region (see Appendix A.3).

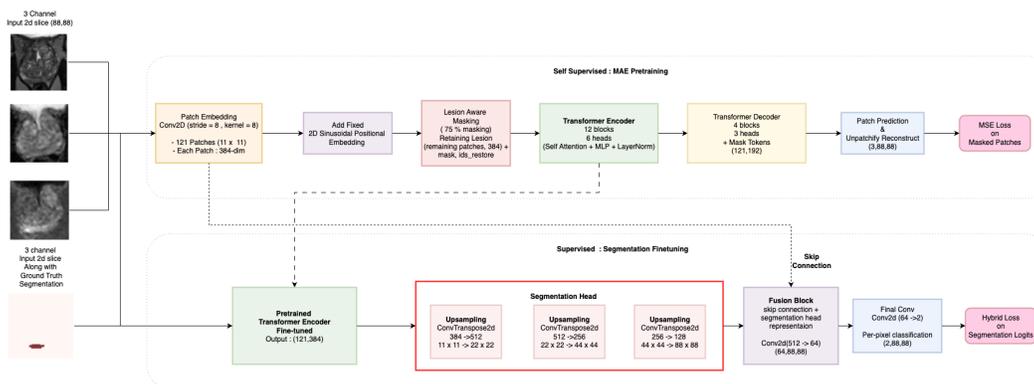


Figure 1: Multimodal Two-phase training pipeline

74 3 Results

75 To quantitatively evaluate the performance of our proposed lesion segmentation framework we reserve
76 a dedicated test set comprising 693 prostate bpMRI volumes, yielding a total of 44,352 2D slices.
77 These volumes are strictly held out during both self-supervised pretraining and supervised fine-tuning
78 to ensure an unbiased assessment of generalization capability, segmentation visualisation can be seen
79 in Appendix A.4. Table 1 reports the Dice Similarity Coefficient (Dice) and the Intersection over
80 Union (IoU), computed on a per-slice basis compared against state of art methods.

Table 1: Comparison to state-of-art approaches

Method	Dice Score	IoU
UNet [3]	0.2031	0.1429
VNet [4]	0.1891	0.1264
Attention-UNet [5]	0.2203	0.1937
TransUnet [6]	0.2601	0.2122
SwinUnet [7]	0.2891	0.2311
Our Approach	0.4213	0.3785

81 4 Discussion and Future Work

82 A persistent bottleneck in medical image analysis remains the lack of large, annotated datasets. Even
83 with increasingly powerful model architectures, the scarcity of expert-annotated labels limits the
84 performance and generalizability of supervised approaches. A common workaround is to adopt
85 transfer learning from models pretrained on generic datasets like ImageNet. However, such models
86 often fail to generalize well to medical domains due to a significant domain shift. Features learned
87 from natural images do not adequately capture the complex texture, anatomy, and modality-specific
88 patterns in medical scans.

89 Recognizing these gaps, our approach leverages self-supervised learning (SSL) to pretrain on un-
90 labelled prostate bpMRI data. This allows the model to learn prostate-specific anatomical and
91 pathological features directly from the domain of interest, without reliance on external data distribu-
92 tions.

93 Lesion segmentation in the prostate often resembles a concealed object segmentation task. Lesions
94 tend to be amorphous, poorly contrasted, and lack consistent structure, making them difficult to
95 learn through traditional pattern-based methods under supervised settings. Conventional CNN-based
96 models, which rely heavily on localized features and spatial hierarchies, struggle to generalize
97 across the morphological variability of prostate lesions. In contrast, we employ a transformer-based
98 encoder, specifically a Masked Autoencoder (MAE) which captures global context and non-local
99 dependencies, making it more suitable for learning nuanced lesion features across diverse patient
100 populations. By integrating skip connections from the MAE encoder into a lightweight CNN-based
101 decoder, our finetuned segmentation head combines rich contextual embeddings with localized
102 upsampling, enabling precise and robust lesion delineation.

103 Our proposed framework shows strong performance in segmenting both clinically significant and
104 non-significant prostate lesions, highlighting its potential to enhance diagnostic precision across the
105 spectrum of prostate cancer. Moving forward, we aim to extend this work toward anatomy-specific
106 generalizability by developing task-specific decoder heads for downstream applications, such as
107 identifying clinically significant lesions and delineating finer prostate subregions. These extensions
108 will further advance the clinical relevance of our approach and support its scalable deployment in
109 real-world settings

110 5 Potential Negative Societal Impact

111 While our proposed framework demonstrates promising performance in prostate lesion segmentation,
112 several potential negative societal impacts warrant careful consideration. Automated segmentation
113 errors could lead to missed cancerous lesions or false positives, potentially delaying treatment or

114 causing unnecessary biopsies. False negatives are particularly dangerous as they may provide false
115 reassurance to patients with actual cancer. Furthermore, clinicians may overrely on the system,
116 reducing critical human oversight in diagnosis, a concern amplified by the "black box" nature of
117 deep learning models, which makes it difficult for clinicians to understand why certain predictions
118 were made and, in cases of misdiagnosis, creates challenges in determining accountability among
119 model developers, clinicians, and institutions. The model may perform poorly on underrepresented
120 patient populations if the training data lacks diversity across different ethnicities, age groups, and
121 scanner types, and such performance disparities could exacerbate existing healthcare inequities, with
122 certain demographic groups receiving inferior diagnostic accuracy. Additionally, patients may not be
123 adequately informed about AI involvement in their diagnosis, raising concerns about informed consent
124 and transparency. To mitigate these risks, we emphasize that this system is designed to assist, not
125 replace, clinical expertise. We strongly recommend rigorous external validation across diverse patient
126 populations and clinical settings, continuous monitoring of model performance post-deployment, and
127 clear communication with patients regarding AI-assisted diagnosis. Future work should prioritize
128 model interpretability and establish robust clinical validation protocols before real-world deployment.

129 **References**

- 130 [1] Hyuna Sung, Jacques Ferlay, Rebecca L. Siegel, Mathieu Laversanne, Isabelle Soerjomataram,
131 Ahmedin Jemal, and Freddie Bray. Global cancer statistics 2020: Globocan estimates of incidence
132 and mortality worldwide for 36 cancers in 185 countries. *CA: A Cancer Journal for Clinicians*,
133 71(3):209–249, 2021.
- 134 [2] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked
135 autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on*
136 *Computer Vision and Pattern Recognition*, pages 16000–16009, 2022.
- 137 [3] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for
138 biomedical image segmentation. In *International Conference on Medical image computing and*
139 *computer-assisted intervention*, pages 234–241. Springer, 2015.
- 140 [4] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. Vnet: Fully convolutional neural
141 networks for volumetric medical image segmentation. In *2016 Fourth International Conference*
142 *on 3D Vision (3DV)*, pages 565–571. IEEE, 2016.
- 143 [5] Ozan Oktay, Jo Schlemper, Loic Le Folgoc, Matthew Lee, Mattias Heinrich, Kazunari Misawa,
144 Kensaku Mori, Steven McDonagh, Nils Y Hammerla, Bernhard Kainz, et al. Attention u-net:
145 Learning where to look for the pancreas. In *Medical Imaging with Deep Learning*, 2018.
- 146 [6] Jieneng Chen, Yongyi Lu, Qihang Yu, Xiangde Luo, Ehsan Adeli, Yan Wang, Le Lu, Alan L
147 Yuille, and Yuyin Zhou. Transunet: Transformers make strong encoders for medical image
148 segmentation. *arXiv preprint arXiv:2102.04306*, 2021.
- 149 [7] Hu Cao, Yueyue Wang, Joy Chen, Dongsheng Jiang, Xiaopeng Zhang, Qi Tian, and Manning
150 Wang. Swin-unet: Unet-like pure transformer for medical image segmentation. *arXiv preprint*
151 *arXiv:2105.05537*, 2022.
- 152 [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
153 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,
154 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image
155 recognition at scale. In *International Conference on Learning Representations*, 2021.

156 **A Appendix**

157 **A.1 Dataset & Preprocessing**

158 The dataset used in this study comprised 4,636 3T bpMRI volumetric scans, collected with institu-
159 tional review board approval and patient consent. Each scan includes axial T2W imaging, DWI, and
160 apparent diffusion coefficient (ADC) maps. High b-value images ($b = 1500$ s/mm²) were calculated
161 from acquired $b = 50$ and $b = 1000$ s/mm² data following standard protocols.

162 All scans were co-registered to the T2W space to ensure spatial alignment across modalities. The
163 volumes were then resampled to a standardized size of $(192 \times 192 \times 64)$ voxels and intensity-
164 normalized.

165 Given our focus on lesion identification within the prostate region, we trained a 3D U-Net architecture
166 from scratch to automatically segment the prostate gland. This segmentation model achieved a Dice
167 coefficient of 0.90 for prostate region identification, demonstrating robust performance for region-
168 of-interest (ROI) extraction. Using the automated prostate segmentation masks, we cropped axial
169 regions around the prostate with 8-pixel padding along both height and width dimensions to ensure
170 complete prostate coverage while accounting for potential segmentation uncertainties. Following
171 ROI extraction, the volumetric scans were resized to $(88 \times 88 \times 64)$ voxels.

172 For training, the 3D volumes were processed as 2D slices of size (88×88) along the z-direction. Each
173 slice comprised three co-registered modalities—T2W, ADC, and DWI—which were concatenated
174 along the channel dimension, resulting in a 3-channel input tensor of dimensions $(3 \times 88 \times 88)$.
175 This 2D slice-based approach enabled efficient processing while preserving the multiparametric
176 information essential for lesion characterization.

177 A.2 Representation Learning in multimodal setting

178 The MAE framework consists of two core components: a Vision Transformer (ViT)-based [8], encoder
179 and a lightweight Transformer decoder, optimized to learn anatomical and pathological representations
180 from 2D axial slices extracted from 3D volumetric scans. Each input slice is represented as a 3-
181 channel image, these three modalities are stacked channel-wise to form a tensor of shape $(3 \times 88 \times$
182 $88)$. During training, all three channels are treated equally and passed through a shared convolutional
183 patch embedding layer, which projects the full input into patch tokens without distinguishing between
184 modalities. However, because the model is trained to reconstruct all three channels jointly, it implicitly
185 learns inter-modality correlations.

186 This multi-modal joint reconstruction enables the model to:

- 187 • Leverage structural consistency in T2W images
- 188 • Learn contrast sensitivity from ADC values
- 189 • Attend to high-signal specific regions from the DWI

190 This cross-modal representation learning serves as critical step for the downstream lesion segmenta-
191 tion, to learn the subtle differences between these modalities.

192 The input image is split into 8×8 non-overlapping patches using a convolutional embedding layer. For
193 an 88×88 slice, this results in $11 \times 11 = 121$ patches. Each patch is projected into a high-dimensional
194 embedding (384 dimensions) and enriched with a 2D sinusoidal positional encoding that preserves
195 spatial structure.

196 The encoder comprises 12 Transformer blocks, each with multi-head self-attention and feedforward
197 MLPs. It operates on only the visible (unmasked) patches, enabling the model to learn from context.
198 A learnable [MASK] token is later used by the decoder to predict masked patches.

199 In conventional MAE training, a fixed proportion of input patches are randomly masked, and the
200 model is trained to reconstruct them using only the visible patches. However, this randomness can
201 result in lesion-containing patches being masked out entirely — which is suboptimal for medical
202 images, where pathologies may occupy a very small portion of the image. To improve focus on
203 clinically meaningful regions, we introduce **lesion-aware masking** designed specifically to enhance
204 the MAE’s ability to learn informative representations of prostate lesions from multi-modal MRI.

205 Our lesions aware masking uses binary lesion annotations to bias the patch sampling process, without
206 introducing label supervision into the loss. Each input slice is accompanied by a binary lesion mask
207 $M \in \{0, 1\}^{H \times W}$ where 1 denotes the lesion pixel. We downsample M to patch resolution using
208 average pooling with a kernel and stride equal to the patch size, (8×8) for our case. This yields a
209 patch-level lesion score matrix $M' \in [0, 1]^{H_p \times W_p}$ where each entry represents the fraction of lesion
210 pixels in a patch. We convert M' into sampling weights using the following transformation (1)

$$w_{ij} = 1 + \alpha \cdot M'_{ij} \quad (1)$$

211 Here α is a hyperparameter that controls the emphasis on lesion patches, we set it to α to 10 for our
 212 implementation. This ensures that lesion-containing patches are preferentially retained during MAE
 213 training, while the surrounding context is masked and must be reconstructed. As a result:

- 214 • The model is forced to reconstruct anatomical context around the lesion
- 215 • It still learns from both healthy and pathological regions
- 216 • No lesion labels are used in the reconstruction loss — only in guiding the mask.

217 Following the lesion aware masking, the decoder reconstructs the complete slice (all three channels)
 218 by predicting the pixel values of masked patches. It consists of 4 lightweight Transformer blocks
 219 followed by a linear projection layer that maps each token back to a flattened $8 \times 8 \times 3$ patch. The
 220 reconstruction loss is computed only on masked patches using mean squared error (MSE) (2). Before
 221 computing the reconstruction loss, each target patch is independently normalized to zero mean and
 222 unit variance. This patch-wise normalization ensures consistent gradient scaling across modalities
 223 and patch intensities, mitigating training instability due to high dynamic range in multi-modal inputs,
 224 like the bright DWI lesions and low-signal ADC regions. It also encourages the model to learn
 225 relative structural patterns rather than raw intensities, which improves generalizability across varying
 226 scan conditions.

$$L_{MAE} = \frac{1}{\sum m_i} \cdot \sum_i^N m_i \cdot \|\hat{x}_i - x_i\|^2 \quad (2)$$

227 where $m_i \in 0, 1$ is the binary mask for the patch i , \hat{x} and x are the reconstructed and original patches
 228 respectively.

229 We train the MAE for 250 epochs using the AdamW optimizer with a cosine learning rate scheduler,
 230 Figure 2 shows the MAE loss curve. The training data for this phase includes all axial slices from all
 231 training volumes without filtering or balancing. This ensures that the MAE is exposed to the natural
 232 distribution of normal and pathological anatomy, including slices with no lesions.

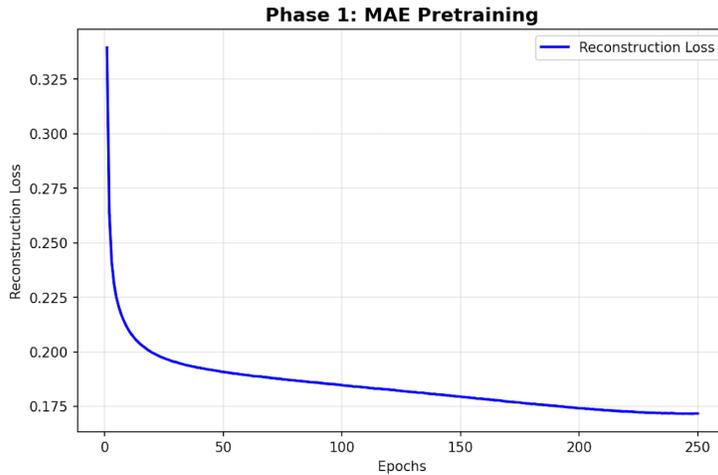


Figure 2: MAE training loss curve

233 This phase results in a pretrained encoder with structural priors to build lesion specific understanding
 234 allowing the model to learn clinical relevance in an unsupervised setting, Figure 3 & 4.

235 A.3 Finetuning for Lesion Segmentation

236 Following the SSL pretraining, we fine-tune the encoder for the prostate lesion segmentation task
 237 using the annotated slice. The pretrained encoder is frozen initially and then fine-tuned along with a

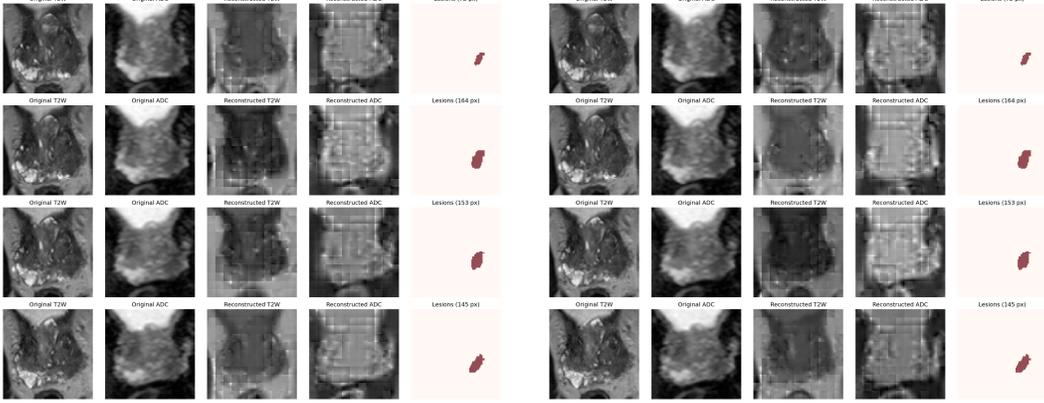


Figure 3: MAE reconstruction 100 epoch

Figure 4: MAE reconstruction 250 epoch

238 lightweight decoder. This phase leverages the rich representations learned during Phase 1 to improve
 239 lesion detection performance.

240 To recover full-resolution predictions from the patch-wise latent representations, the decoder is
 241 composed of 3 upsampling CNN blocks implemented to expand the spatial resolution from patch-
 242 level (11×11) back to full image size (88×88). We introduce a fusion block that concatenates
 243 upsampled decoder features with skip features from the encoder’s patch embedding layer, to inject
 244 high-level spatial features from the encoder layer. The decoder and skip features are fused via a 3×3
 245 convolution and passed through a final 1×1 convolutional layer to produce the two-class output mask.

246 We use differential learning rates while finetuning the encoder and decoder head, a lower learning
 247 rate, of $5.00e - 06$, for the pretrained encoder to preserve learned features while a higher of learning
 248 rate, of $5.00e - 05$, for the newly initialized decoder to adapt to the segmentation task.

249 Along with sharing close morphological features between benign and malignant structures, PCA
 250 lesion an extended challenge of the class imbalance problem at two levels, most bpMRI slices contain
 251 no lesions and even with slices containing the lesions, lesions occupy an extremely small fraction of
 252 the pixels.

253 To address the class imbalance issue, we construct a balanced lesion sampling strategy along with
 254 a hybrid loss function. While finetuning, include all slices with lesions, provided they contain at
 255 least a minimum number of lesion pixels, set to be 5 pixels for our study. We randomly sample
 256 non-lesion slices to make up the remainder of the training set, maintaining a 50%–50% ratio of
 257 lesion vs. non-lesion slices. This ensures that the model is exposed to a meaningful number of lesion
 258 samples while still learning to differentiate lesions from healthy tissue. To further combat pixel-level
 259 class imbalance, we employ a hybrid loss function (3) that combines focal loss (4) and dice loss (5)

$$L_{semgnetation} = \lambda_{focal} \cdot L_{focal} + \lambda_{dice} \cdot L_{dice}, \quad \text{where } \lambda_{focal} = \lambda_{dice} = 0.5 \quad (3)$$

260 here, L_{focal} is calculated as follows:

261 Let the cross-entropy loss be defined as:

$$CE(x, y) = -\log p_t, \quad \text{where } p_t = \begin{cases} p & \text{if } y = 1 \\ 1 - p & \text{if } y = 0 \end{cases}$$

262 Then, the focal loss is:

$$L_{focal} = \frac{1}{N} \sum_{i=1}^N \alpha \cdot (1 - p_t^{(i)})^\gamma \cdot CE(x^{(i)}, y^{(i)}) \quad (4)$$

263 where $N = B \cdot H \cdot W$ is the number of pixels in the batch, $\alpha = 0.25$, $\gamma = 2.0$ (selected hyperparam-
 264 eters) and $p_t^{(i)}$ is the model’s predicted probability

265 \mathcal{L}_{dice} for class $c \in \{0, 1\}$ (background and lesion) is computed as:

$$\mathcal{L}_{dice} = 2 - \sum_{c=0}^1 \frac{2 \cdot \sum_i p_c^{(i)} \cdot y_c^{(i)} + \epsilon}{\sum_i p_c^{(i)} + \sum_i y_c^{(i)} + \epsilon} \quad (5)$$

266 where $p_c^{(i)}$ is the predicted probability for class c at pixel i , $y_c^{(i)}$ is the one-hot encoded ground truth
267 for class c at pixel i and ϵ is a smoothing constant, set to 10^{-5} to avoid division by zero

268 The AdamW optimizer with a cosine annealing learning rate schedule is used to stabilize training.
269 Model performance is monitored using the dice score on a held-out validation set, and early stopping
270 is applied if no improvement is observed over 15 consecutive epochs. Following this two-phase
271 training strategy - combining self-supervised pretraining with supervised fine-tuning - our approach
272 achieves strong performance and outperforms several recent state-of-the-art methods for prostate
273 lesion segmentation.

274 **A.4 Results for Segmentation**

275 We report standard overlap-based segmentation metrics, including the Dice Similarity Coefficient
276 (Dice) (6) and the Intersection over Union (IoU) (7), computed on a per-slice basis.

$$\text{Dice} = \frac{2 \times \text{Area of Overlap}}{\text{Total Area}} = \frac{2 \cdot |A \cap B|}{|A| + |B|} \quad (6)$$

$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}} = \frac{|A \cap B|}{|A \cup B|} \quad (7)$$

277 where A is the predicted segmentation mask and B is the ground truth mask

278 Figure 5 and Figure 6, show the qualitative results on the test set samples showing T2W, ADC, and
279 DWI inputs with ground truth (GT) and model predictions for lesion segmentation.

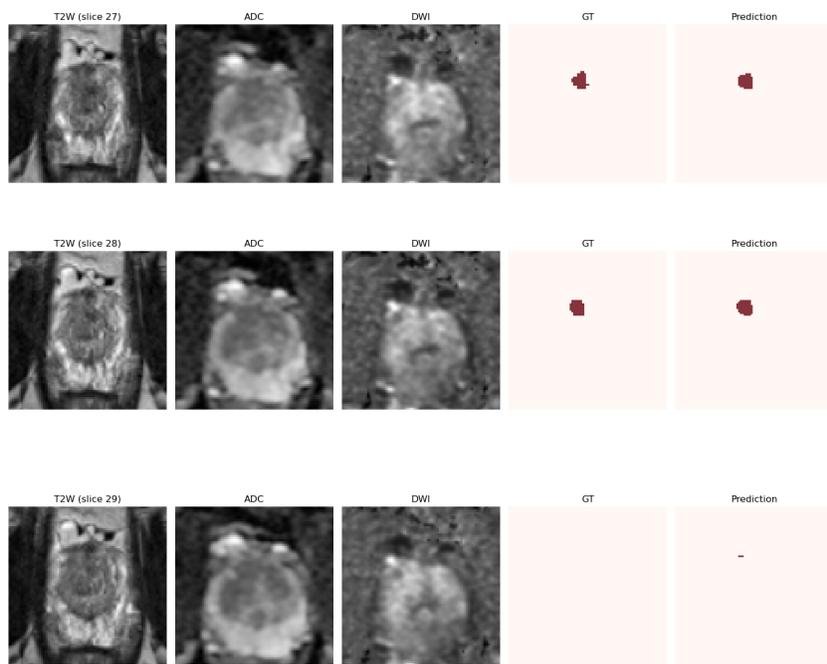


Figure 5: Qualitative Segmentation results

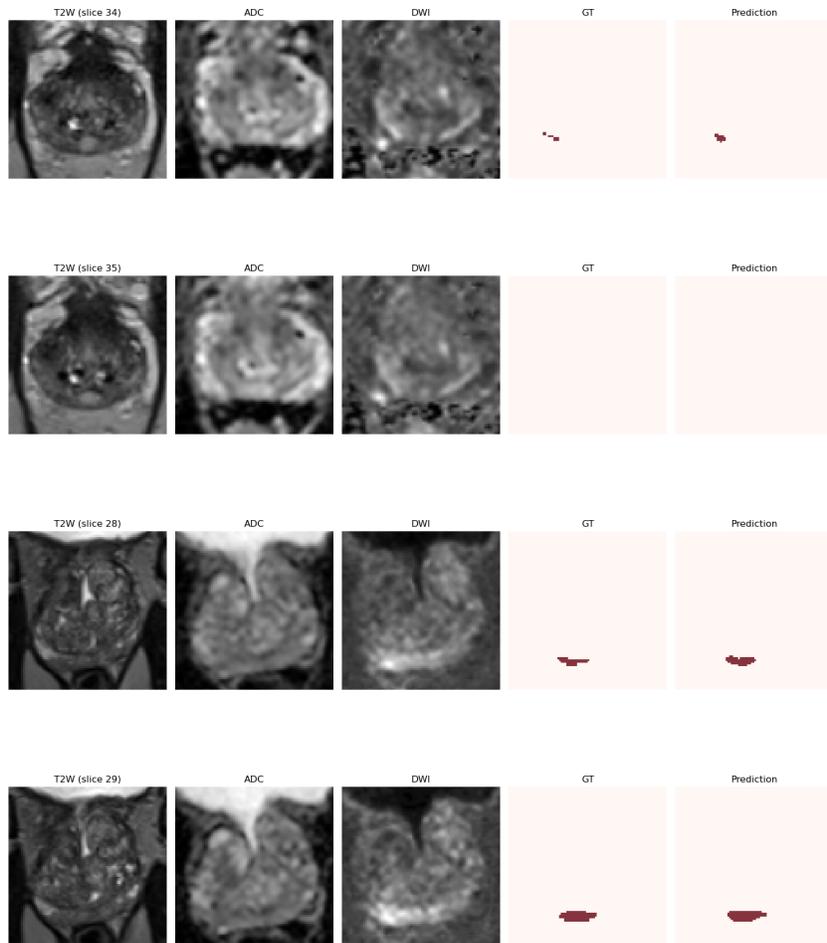


Figure 6: Qualitative Segmentation results (Part 2)