



Toward Fair and Accurate Cross-Domain Medical Image Segmentation: A VLM-Driven Active Domain Adaptation Paradigm

Supplementary Material

001 This is the supplementary material for *Toward Fair and*
002 *Accurate Cross-Domain Medical Image Segmentation: A*
003 *VLM-Driven Active Domain Adaptation Paradigm*. We
004 present the following materials:

- 005 • Sec. 1 The more details of evaluation metrics we used.
- 006 • Sec. 2 The more experiments (about different attributes)
- 007 and related ablation studies.
- 008 • Sec. 3 A brief introduction to our comparative methods.
- 009 • Sec. 4 More visualizations of VLM-attribute learning.

010 1. Details of Evaluation Metrics

011 For quantifying fairness in medical image segmentation,
012 we adopt the same metrics used in previous studies [1, 5,
013 7]: Dice, IoU, and the important composite metrics, ES-
014 Dice and ES-IoU, to assess fairness alongside performance.
015 More detailed information will be added as follows.

016 For Equity-Scaled metrics, we first need to compute a
017 performance discrepancy Δ for every sensitive attribute.
018 This discrepancy is characterized by the cumulative differ-
019 ence between each demographic subgroup’s metric and the
020 overall performance. It is articulated in the following for-
021 mulation:

$$\Delta = \sum_{A \in \text{attrs}} |M(\{\hat{y}, y\}) - M(\{\hat{y}, y, a\} | a = A)|, \quad (1)$$

022 where *attrs* represent demographic groups such as
023 {Female, Male}, {Asian, Black, White}, or {Hispanic,
024 Non-Hispanic, Unknown}, *M* denotes a specific metric
025 (e.g., Dice or IoU), and \hat{y} is the ground truth. A positive Δ
026 value implies that smaller values correspond to reduced per-
027 formance disparities among demographic groups relative to
028 the overall performance, indicating improved fairness. The
029 Equity-Scaled metric can be formulated as follows:
030

$$\text{ESM} = \frac{M(\{\hat{y}, y\})}{1 + \Delta}. \quad (2)$$

031 Through the above steps, we can calculate the ES-Dice and
032 ES-IoU metrics.
033

034 2. More Experiments

035 **More Experiments on Other Attribute.** To further vali-
036 date the generalizability of our method, we conducted ex-
037 periments across another sensitive attribute (Gender) and
038 performed a comprehensive comparison with state-of-the-

art Domain Adaptation (DA) and Active Domain Adapta-
039 tion (ADA) approaches. All ADA methods were evalu-
040 ated under the same labeling quota to ensure fair compar-
041 isons. As shown in Table 1, our method significantly out-
042 performs all DA and ADA methods across key metrics. For
043 instance, in rim segmentation tasks, our approach achieves
044 an ES-Dice score of 0.785 and ES-IoU of 0.661, surpass-
045 ing previous methods by a notable margin. Statistical anal-
046 ysis (P -value < 0.05) further confirms the significant su-
047 periority of our framework in improving performance and
048 fairness. These results highlight our method’s ability to
049 address cross-domain challenges while ensuring equitable
050 outcomes for diverse demographic subgroups.
051

Detailed Ablation Studies. Our ablation study, conducted
052 after VLM learning, encompasses four main configura-
053 tions: 1) Fair-Base: This configuration adheres to the Fair
054 Quota principle by randomly selecting samples with an
055 equal quota across different subgroups of the same sen-
056 sitive attribute, resulting in the final selection list *S*. 2)
057 Fair-Attr: Building on the Fair Quota framework, this setup
058 exclusively employs the Attribute representative selection
059 method to choose samples from each subgroup, culminat-
060 ing in the selection list *S*. 3) Fair-Poly: Maintaining the same
061 quota, this configuration focuses on selecting samples with
062 Polysemy representatives within each subgroup to form the
063 final list *S*. 4) Our Fair-AP: Our approach enhances the
064 balanced allocation strategy of Fair-Base by prioritizing the
065 selection of samples that integrate both Attribute and Pol-
066 ysemy considerations, thereby augmenting representation
067 and diversity, and resulting in the final selection list *S*. The
068 ablation study results for each sensitive attribute (gender,
069 race, and ethnicity) are displayed in Table 2, Table 3, and
070 Table 4, respectively.
071

Gender Attribute: In Table 2, the standalone application of
072 Fair-Attr or Fair-Poly demonstrates potential for enhancing
073 certain metrics, while the observed improvements exhibit
074 variability, suggesting limited consistency in performance
075 across all evaluation criteria. However, it is evident that
076 our Fair-AP comprehensively considers both attribute rep-
077 resentative and polysemy representative, ensuring superior
078 performance and fairness.
079

Race Attribute: The findings from the ablation study on
080 racial attributes indicate that if the distribution of various
081 races is imbalanced, both Fair-Attr and Fair-Poly show im-
082 provements in some performance measures (like ES-IoU
083 and IoU for Cup), as illustrated in Table 3. Our Fair-AP
084

Table 1. Cup and rim segmentation performance on the FairDomain-Segmentation benchmark using different DA and ADA methods with **Gender** as the demographic attribute. The * denotes p -value < 0.05 in all paired t -test, indicating statistically significant differences.

		Method	Venue	Overall ES-Dice \uparrow	Overall Dice \uparrow	Overall ES-IoU \uparrow	Overall IoU \uparrow	Male Dice \uparrow	Female Dice \uparrow	Male IoU \uparrow	Female IoU \uparrow
Cup	Baseline (Source)	-	-	0.885	0.888	0.806	0.808	0.886	0.889	0.807	0.810
	Baseline (Target)	-	-	0.688	0.700	0.535	0.555	0.693	0.711	0.557	0.574
Rim	Baseline (Source)	-	-	0.854	0.861	0.753	0.762	0.864	0.856	0.767	0.755
	Baseline (Target)	-	-	0.485	0.495	0.336	0.342	0.486	0.507	0.334	0.353
DA	Cup	PixMatch [6]	CVPR'21	0.768	0.775	0.650	0.660	0.772	0.769	0.645	0.660
		DAFormer [4]	CVPR'22	0.781	0.785	0.676	0.680	0.783	0.789	0.678	0.684
		DAFormer-FIA [7]	ECCV'24	0.802	0.810	0.692	0.700	0.806	0.816	0.695	0.706
	Rim	PixMatch [6]	CVPR'21	0.660	0.673	0.519	0.523	0.669	0.688	0.519	0.528
		DAFormer [4]	CVPR'22	0.344	0.345	0.212	0.213	0.344	0.347	0.212	0.214
		DAFormer-FIA [7]	ECCV'21	0.528	0.531	0.367	0.369	0.533	0.528	0.372	0.366
ADA	Cup	Random	-	0.828	0.834	0.729	0.734	0.838	0.831	0.738	0.731
		Entropy [10]	CVPR'19	0.819	0.823	0.717	0.721	0.826	0.821	0.724	0.719
		MHPL [2]	CVPR'23	0.829	0.834	0.728	0.733	0.838	0.832	0.738	0.730
		DML-Core [9]	NeurIPS'24	0.831	0.839	0.733	0.740	0.844	0.835	0.746	0.736
		Detective [11]	CVPR'24	0.831	0.837	0.733	0.740	0.841	0.834	0.745	0.736
		STDR [3]	IEEE TMI'24	0.831	0.837	0.731	0.738	0.842	0.834	0.743	0.734
	Our FairAP	-	0.839*	0.843*	0.742*	0.746*	0.845	0.841*	0.749*	0.744*	
	Rim	Random	-	0.778	0.780	0.652	0.654	0.779	0.782	0.652	0.655
		Entropy [10]	CVPR'19	0.768	0.772	0.639	0.643	0.769	0.774	0.640	0.645
		MHPL [2]	CVPR'23	0.779	0.782	0.654	0.656	0.780	0.784	0.654	0.658
DML-Core [9]		NeurIPS'24	0.782	0.783	0.658	0.658	0.782	0.784	0.658	0.658	
Detective [11]	CVPR'24	0.781	0.784	0.657	0.659	0.782	0.786	0.658	0.661		
STDR [3]	IEEE TMI'24	0.777	0.781	0.652	0.656	0.778	0.783	0.653	0.658		
Our FairAP	-	0.785*	0.787*	0.661*	0.663*	0.786*	0.789*	0.661*	0.664*		

Table 2. Ablation study of Optic cup and rim segmentation performance on the FairDomain-Segmentation dataset with **Gender** as the demographic attribute.

		Method	Overall ES-Dice \uparrow	Overall Dice \uparrow	Overall ES-IoU \uparrow	Overall IoU \uparrow	Male Dice \uparrow	Female Dice \uparrow	Male IoU \uparrow	Female IoU \uparrow	
Cup	Fair-Base		0.832	0.836	0.734	0.737	0.838	0.834	0.739	0.735	
	Fair-Attr		0.831	0.836	0.732	0.738	0.841	0.833	0.743	0.735	
	Fair-Poly		0.834	0.841	0.735	0.743	0.846	0.837	0.749	0.739	
	Ours Fair-AP		0.839	0.843	0.742	0.746	0.845	0.841	0.749	0.744	
ADA	Rim	Fair-Base		0.777	0.783	0.652	0.658	0.778	0.787	0.653	0.662
		Fair-Attr		0.774	0.780	0.649	0.655	0.776	0.783	0.650	0.659
		Fair-Poly		0.781	0.784	0.657	0.659	0.782	0.785	0.657	0.660
		Ours Fair-AP		0.785	0.787	0.661	0.663	0.786	0.789	0.661	0.664

085 not only performs well in segmentation outcomes but also
 086 effectively ensures algorithmic fairness.
 087 **Ethnicity Attribute:** In Table 4, a similar trend is ob-
 088 served: our Fair-AP model preserves fairness robustness
 089 while achieving competitive performance gains in segmen-
 090 tation tasks, even under imbalanced distribution scenarios.

091 3. Comparison Methods

092 We conducted a series comparisons of Fair-AP with leading
 093 DA techniques, including Pixmatch [6], DAFormer [4], and
 094 DAFormer-FIA [7], as well as ADA approaches like En-
 095 tropy [10], MHPL [2], DML-Core [9], Detective [11], and
 096 STDR [3]. The details of these methods are summarized as

follows:

Pixmatch [6]: Pixmatch introduces a novel unsupervised
 domain adaptation framework that enhances model perfor-
 mance on target domains by ensuring consistency in predic-
 tions under small input perturbations.

DAFormer [4]: DAFormer introduces a Transformer-based
 encoder-decoder architecture with three key training strate-
 gies—Rare Class Sampling, Thing-Class ImageNet Feature
 Distance, and learning rate warmup—to stabilize training
 and mitigate overfitting in unsupervised domain adaptation
 (UDA) for semantic segmentation.

DAFormer-FIA [7]: DAFormer-FIA introduces a novel
 Fair Identity Attention (FIA) module and the first fairness-

Table 3. Ablation study of Optic cup and rim segmentation performance on the FairDomain-Segmentation dataset with **Race** as the demographic attribute.

Method		Overall ES-Dice↑	Overall Dice↑	Overall ES-IoU↑	Overall IoU↑	Asian Dice↑	Black Dice↑	White Dice↑	Asian IoU↑	Black IoU↑	White IoU↑
Cup	Fair-Base	0.816	0.836	0.714	0.736	0.815	0.837	0.838	0.712	0.742	0.738
	Fair-Attr	0.810	0.836	0.719	0.739	0.818	0.825	0.840	0.721	0.733	0.742
	Fair-Poly	0.819	0.840	0.721	0.742	0.818	0.839	0.842	0.718	0.744	0.744
	Ours Fair-AP	0.828	0.843	0.731	0.747	0.827	0.843	0.845	0.730	0.751	0.748
ADA Rim	Fair-Base	0.701	0.786	0.581	0.661	0.739	0.729	0.803	0.608	0.596	0.680
	Fair-Attr	0.690	0.778	0.572	0.651	0.738	0.709	0.796	0.610	0.574	0.672
	Fair-Poly	0.699	0.789	0.579	0.664	0.740	0.726	0.807	0.611	0.592	0.685
	Ours Fair-AP	0.711	0.791	0.592	0.667	0.750	0.735	0.807	0.624	0.602	0.686

Table 4. Ablation study of Optic cup and rim segmentation performance on the FairDomain-Segmentation dataset with **Ethnicity** as the demographic attribute.

Method		Overall ES-Dice↑	Overall Dice↑	Overall ES-IoU↑	Overall IoU↑	Hispanic Dice↑	Non-Hispanic Dice↑	Hispanic IoU↑	Non-Hispanic IoU↑
Cup	Fair-Base	0.822	0.838	0.719	0.738	0.854	0.836	0.762	0.735
	Fair-Attr	0.822	0.833	0.717	0.732	0.845	0.832	0.751	0.730
	Fair-Poly	0.822	0.835	0.717	0.735	0.849	0.833	0.758	0.733
	Ours Fair-AP	0.836	0.840	0.730	0.743	0.844	0.839	0.759	0.741
ADA Rim	Fair-Base	0.778	0.786	0.656	0.661	0.794	0.784	0.667	0.659
	Fair-Attr	0.781	0.784	0.657	0.659	0.786	0.782	0.658	0.658
	Fair-Poly	0.777	0.786	0.655	0.663	0.797	0.784	0.672	0.661
	Ours Fair-AP	0.784	0.790	0.663	0.668	0.796	0.788	0.673	0.666

110 focused paired imaging dataset to enhance algorithmic fair-
 111 ness under domain shifts in medical AI, significantly im-
 112 proving performance across demographics in both domain
 113 adaptation (DA) and domain generalization (DG) tasks.

114 **Entropy [10]:**The AdvEnt method [10] is employed to cal-
 115 culate the prediction map entropy for each sample within
 116 the target domain, and those samples with the highest en-
 117 tropy are selected for manual annotation.

118 **MHPL [2]:** MHPL introduces a novel active source-free
 119 domain adaptation approach by identifying and exploiting
 120 "minimum happy points" through tailored selection strate-
 121 gies and a neighbor focal loss, achieving significant perfor-
 122 mance gains with minimal labeling effort.

123 **DML-Core [9]:**DML-Core introduces a slice-based active
 124 learning method integrating deep metric learning with Core-
 125 set, significantly reducing annotation costs for 3D medical
 126 segmentation while achieving high performance under low
 127 annotation budgets.

128 **Detective [11]:**Detective introduces a dynamic domain
 129 adaptation model with evidential uncertainty valuation and
 130 contextual diversity enhancement, effectively selecting in-
 131 formative target samples by jointly evaluating domain shifts
 132 and prediction confidence.

133 **STDR [3]:** STDR improves gross tumor volume seg-
 134 mentation for nasopharyngeal carcinoma by employing a
 135 dual reference strategy to select and annotate representative
 136 target-domain samples, enabling effective source-free do-
 137 main adaptation while ensuring data privacy and minimiz-

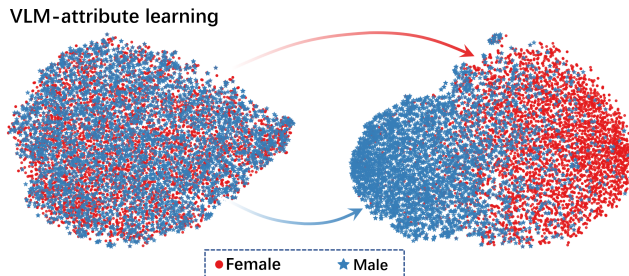


Figure 1. Visualization of gender feature distribution before and after feature alignment via VLM-attribute learning.

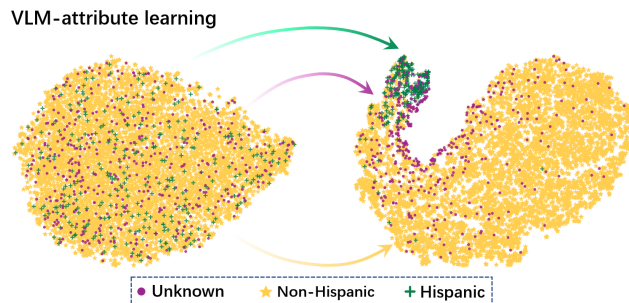


Figure 2. Visualization of ethnic feature distribution before and after feature alignment via VLM-attribute learning.

ing the annotation effort.

4. More Visualizations

The t-SNE visualizations [8] of sample distributions be-
 fore and after feature alignment via VLM-attribute learning

142 for gender and ethnic attributes are depicted in Fig. 1 and
 143 Fig. 2. Our VLM-attribute learning results in a more con-
 144 centrated distribution of features across different subgroups,
 145 thus improving the distinguishability of features associated
 146 with various attributes. In Fig. 1, we note that male and
 147 female samples were effectively distinguished after VLM-
 148 attribute learning. We also notice a balanced representation
 149 of male and female samples, which aids in the subsequent
 150 active selection strategy. Conversely, in Fig. 2, the ethnic
 151 proportions are uneven, and the distribution of Unknown
 152 samples is irregular, causing significant degradation in the
 153 ablation study results when either representativeness or di-
 154 versity is considered in isolation. Our proposed Fair-AP
 155 method achieves success in both segmentation performance
 156 and fairness, even when subgroup proportions are imbal-
 157 anced and sample distributions are irregular.

158 **References**

159 [1] Fairseg: A large-scale medical image segmentation dataset
 160 for fairness learning using segment anything model with fair
 161 error-bound scaling. In *ICLR*. 1
 162 [2] Mhpl: Minimum happy points learning for active source free
 163 domain adaptation. In *CVPR*, 2023. 2, 3
 164 [3] Dual-reference source-free active domain adaptation for na-
 165 sopharyngeal carcinoma tumor segmentation across multiple
 166 hospitals. *IEEE Transactions on Medical Imaging*, 2024. 2,
 167 3
 168 [4] Lukas Hoyer, Dengxin Dai, and Luc Van Gool. Daformer:
 169 Improving network architectures and training strategies for
 170 domain-adaptive semantic segmentation. In *Proceedings of*
 171 *the IEEE/CVF conference on computer vision and pattern*
 172 *recognition*, pages 9924–9935, 2022. 2
 173 [5] Yan Luo, Yu Tian, Min Shi, Louis R Pasquale, Lucy Q Shen,
 174 Nazlee Zebardast, Tobias Elze, and Mengyu Wang. Harvard
 175 glaucoma fairness: a retinal nerve disease dataset for fairness
 176 learning and fair identity normalization. *IEEE Transactions*
 177 *on Medical Imaging*, 2024. 1
 178 [6] Luke Melas-Kyriazi and Arjun K Manrai. Pixmatch: Unsu-
 179 pervised domain adaptation via pixelwise consistency train-
 180 ing. In *Proceedings of the IEEE/CVF Conference on Com-*
 181 *puter Vision and Pattern Recognition*, pages 12435–12445,
 182 2021. 2
 183 [7] Yu Tian, Congcong Wen, Min Shi, Muhammad Muneeb
 184 Afzal, Hao Huang, Muhammad Osama Khan, Yan Luo, Yi
 185 Fang, and Mengyu Wang. Fairdomain: Achieving fairness
 186 in cross-domain medical image segmentation and classifica-
 187 tion. In *European Conference on Computer Vision*, pages
 188 251–271. Springer, 2024. 1, 2
 189 [8] Laurens Van der Maaten and Geoffrey Hinton. Visualizing
 190 data using t-sne. *Journal of machine learning research*, 9
 191 (11), 2008. 3
 192 [9] Arvind Vepa, Zukang Yang, Andrew Choi, Jungseock Joo,
 193 Fabien Scalzo, and Yizhou Sun. Integrating deep metric
 194 learning with coreset for active learning in 3d segmenta-
 195 tion. *Advances in Neural Information Processing Systems*,
 196 37:71643–71671, 2024. 2, 3

[10] Tuan-Hung Vu, Himalaya Jain, Maxime Bucher, Matthieu
 Cord, and Patrick Pérez. Advent: Adversarial entropy min-
 imization for domain adaptation in semantic segmentation.
 In *Proceedings of the IEEE/CVF conference on computer vi-*
sion and pattern recognition, pages 2517–2526, 2019. 2, 3
 [11] Wenqiao Zhang, Zheqi Lv, Hao Zhou, Jia-Wei Liu, Juncheng
 Li, Mengze Li, Yunfei Li, Dongping Zhang, Yueting Zhuang,
 and Siliang Tang. Revisiting the domain shift and sample
 uncertainty in multi-source active domain transfer. In *Pro-*
ceedings of the IEEE/CVF Conference on Computer Vision
and Pattern Recognition, pages 16751–16761, 2024. 2, 3