

Data Poisoning Attack Defense and Evolutionary Domain Adaptation for Federated Medical Image Segmentation

(Appendix)

1 The appendix of this study provides comprehensive details
 2 that support the main framework, methodology, and experimental results presented in the paper. Below is a summary of
 3 each section: Section A provides a detailed explanation of the core algorithms of AdaShield-FL, including an overview of the framework, malicious client purification, evolutionary domain adaptation, and attack detection and aggregation. Section B presents additional quantitative and qualitative results for segmentation, disentangling encoder, independence loss, and ablation studies. Section C presents qualitative segmentation results of AdaShield-FL and other methods. Section D describes how external attackers inject adversarial perturbations into clients' local datasets to degrade the global model's performance. Section E provides the limitations of the study and potential directions for future work.

16 A Method Algorithms

17 A.1 Overview

18 Algorithm 1 provides a comprehensive overview of the AdaShield-FL framework, encompassing key components such as malicious client purification, evolutionary domain adaptation, and attack detection and aggregation.

22 A.2 Malicious Client Purification

23 Algorithm 2 illustrates the purification process for mitigating the influence of malicious clients by refining their local datasets, ensuring the patient data diversity.

26 A.3 Evolutionary Domain Adaptation

27 Algorithm 3 represents the evolutionary domain adaptation process, which enhances model robustness across diverse and heterogeneous client data distributions.

30 A.4 Attack Detection and Aggregation

31 Algorithm 4 shows the process of identifying malicious clients and aggregating model weights and covariance matrices.

Algorithm 1 AdaShield-FL Overall Algorithm

Input: Number of clients N , local datasets I_i , global model weights in r -th round θ^r , malicious status λ_i^r ;
Output: Final global model weights θ^R ;

- 1: **for** each round $r = 1$ to R **do**
- 2: **for** each client $i \in \{1, \dots, N\}$ **do**
- 3: Download $\theta_{global}^r, C_{global}^r, \lambda_i^r$
- 4: **if** $\lambda_i^r = 1$ **then**
- 5: $I'_i \leftarrow \text{Malicious Client Purification } (I_i, \theta_{global}^r)$
- 6: $(\theta_i^{r+1}, C_i^{r+1}, L_i^{r+1}) \leftarrow \text{Evolutionary Domain Adaptation } (I'_i, \theta_i^r, C_{global}^r)$
- 7: **else**
- 8: $(\theta_i^{r+1}, C_i^{r+1}, L_i^{r+1}) \leftarrow \text{Evolutionary Domain Adaptation } (I_i, \theta_i^r, C_{global}^r)$
- 9: **end if**
- 10: Upload $\theta_i^{r+1}, C_i^{r+1}, L_i^{r+1}$ to the server
- 11: **end for**
- 12: $(\theta_{global}^{r+1}, C_{global}^{r+1}, \lambda_{1,2,\dots,N}^{r+1}) \leftarrow \text{Attack Detection and Aggregation } (\theta_{1,2,\dots,N}^{r+1}, C_{1,2,\dots,N}^{r+1}, L_{1,2,\dots,N}^{r+1})$
- 13: Broadcast $\theta_{global}^{r+1}, C_{global}^{r+1}$, and $\lambda_{1,2,\dots,N}^{r+1}$ to all clients
- 14: **end for**
- 15: **return** Final global model weights θ^R

Algorithm 2 Malicious Client Purification

Input: Local data I_i , global model weights θ_{global} ;
Output: Purified data I'_i ;

- 1: Initialize DRS using θ_{global}
- 2: Input I_i into the DRS
- 3: **for** each iteration $l \in \{0, \dots, L - 1\}$ **do**
- 4: Transform l -th purification data I_i^l to k -space using FFT:
$$s(k_x, k_y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I(x, y) e^{-i2\pi(xk_x + yk_y)} dx dy$$
- 5: Extract f_R and f_S using disentangling encoder in the k -space
- 6: Transform f_R and f_S to f'_R and f'_S using IFFT
- 7: Obtain $(l + 1)$ -th purification data I_i^{l+1} using f'_R and reconstruction decoder
- 8: **end for**
- 9: $I'_i \leftarrow I_i^L$
- 10: **return** I'_i

Algorithm 3 Evolutionary Domain Adaptation

Input: Local dataset I_i , global covariance matrix C_{global}^r , global model weights θ_{global}^r ;

Output: Updated model weights θ_i^{r+1} , covariance matrix C_i^{r+1} , and training loss sequence L_i^{r+1} ;

- 1: Initialize θ_{global}^r to segmentation model
- 2: **for** each training iteration $t \in \{1, \dots, T\}$ **do**
- 3: Define total loss \mathcal{L}_{total} :
$$\mathcal{L}_{total} = (1 - \beta^t) \mathcal{L}_{dice} + \beta^t \mathcal{L}_{cm}$$
- 4: **for** each generation $g \in \{0, \dots, G - 1\}$ **do**
- 5: Sample candidate solutions using:
$$x_k^{g+1} = m^g + \sigma^g \cdot \mathcal{N}(0, \Sigma^g), \quad k = 1, \dots, 20$$
- 6: Compute total loss and select top-10 candidates using $\mathcal{L}_{total}(x_k^{g+1})$:
- 7: Update mean m^{g+1} using top-10 candidates.
- 8: **if** $\|m^{g+1} - m^g\| \leq \epsilon_{cma}$ **then**
- 9: Set $\beta^t = m^{g+1}$ and stop evolution
- 10: Update segmentation model using \mathcal{L}_{total}
- 11: **end if**
- 12: Update covariance matrix Σ^{g+1} using top-10 candidates
- 13: Sample new candidate solutions with Σ^{g+1} and m^{g+1}
- 14: **end for**
- 15: **end for**
- 16: Obtain final segmentation model weight θ_i^{r+1}
- 17: Compute covariance matrix C_i^{r+1} using disentangled segmentation features f_S
- 18: Compute training dice loss sequence L_i^{r+1}
- 19: **return** $\theta_i^{r+1}, C_i^{r+1}, L_i^{r+1}$

Algorithm 4 Attack Detection and Aggregation

Input: Training dice loss sequence L_i^r , local model weights $\theta_{1,2,\dots,N}^r$, local covariance matrices $C_{1,2,\dots,N}^r$;

Output: Updated global model weights θ^r , attack status $\lambda_{1,\dots,N}^r$;

- 1: Extract gradient latent code $\Delta L_{1,2,\dots,N}^r$
- 2: Extract curvature latent code $\Delta^2 L_{1,2,\dots,N}^r$
- 3: Concatenate differential loss patterns
- 4: Obtain the malicious status $\lambda_{1,2,\dots,N}^r$ by applying anomaly detection to the concatenated differential loss patterns.
- 5: Aggregate global model weights θ^r using $\theta_{1,2,\dots,N}^r$ and $\lambda_{1,2,\dots,N}^r$:

$$\theta^r \leftarrow \theta^{r-1} - \eta \frac{1}{N - \sum_{i=1}^N \lambda_i^r} \sum_{i=1}^N \frac{g_i^{r-1} \cdot (1 - \lambda_i^r)}{\sqrt{G^{r-1}} + \epsilon}$$

- 6: Aggregate global covariance matrix C_{global}^r using $C_{1,2,\dots,N}^r$ and $\lambda_{1,2,\dots,N}^r$:

$$C_{global}^r = \frac{1}{\sum_{i=1}^N \rho_i (1 - \lambda_i^r)} \left[\sum_{i=1}^N \rho_i C_i^r (1 - \lambda_i^r) + \sum_{i=1}^N \rho_i (\mu_i^r - \mu^r) (\mu_i^r - \mu^r)^T (1 - \lambda_i^r) \right]$$

- 7: **return** Malicious status $\lambda_{1,\dots,N}^r$, global covariance C_{global}^r , and updated weights θ^r to clients
-

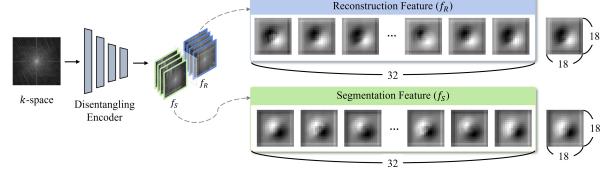


Figure 1: Visualization of disentangled reconstruction feature (f_R) and segmentation feature (f_S) extracted from our disentangling encoder in k -space.

	Dice Score			
	$L = 1$	$L = 2$	$L = 3$	$L = 4$
AdaShield-FL	78.2	80.3	82.7	81.1

Table 1: Impact of iterations (L) on dice score in DRS on the M&Ms dataset under the FGSM attack.

B Additional Experiments

B.1 Disentangling Encoder

In the DRS module, a disentangling encoder [Zhou *et al.*, 2021] is composed of four residual blocks, with each block consisting of two convolutional layers, a batch normalization layer, and a Leaky ReLU activation function. In our framework, the reconstruction and segmentation decoder mirror this structure symmetrically; however, its input vector is half the size of the encoder output.

Figure 1 illustrates the disentangled features obtained from k -space data through a disentangling encoder. It shows that the disentangled features exhibit distinct spatial distributions with the reconstruction feature (f_R) focusing on specific regions, while the segmentation feature (f_S) concentrates on distinct areas, indicating that the two features are directed toward different parts of k -space data. This directional separation indicates that the encoder effectively disentangles the features by learning task-specific representations tailored for reconstruction and segmentation. This separation highlights the effectiveness of feature disentanglement in facilitating distinct objectives in medical imaging.

B.2 the Number Iterations L

In addition, Table 1 demonstrates the impact of the number of iterations (L) on the dice score in the DRS module, evaluated on the M&Ms [Campello *et al.*, 2021] dataset under a fast gradient sign method (FGSM) [Goodfellow *et al.*, 2014] attack. AdaShield-FL achieves its optimal performance when L is set to 3, providing the best balance between reconstruction and segmentation accuracy. In Fig. 2, as L increases, the segmentation maps progressively improve, with false-negative cases decreasing significantly. Early iterations ($L=0$ and $L=1$) show notable incorrect segmentation in regions such as the left ventricle and myocardium. However, as the purification process continues, these false-negative cases are gradually corrected, leading to segmentation maps that closely align with the ground truth. Notably, when $L = 3$, the segmentation map closely resembles the ground truth, indicating that three iterations yield the most balanced and accurate results in terms of segmentation quality. This result

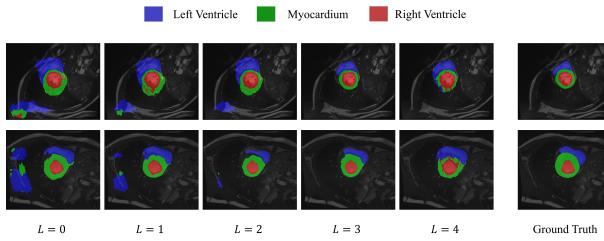


Figure 2: Qualitative examples of the segmentation map during L iterations in disentangled reconstruction and segmentation on ACDC dataset perturbed by PGD in "hypertrophic cardiomyopathy" patients.

Model	Dice Score		
	FGSM	PGD	C&W
FedAvg [McMahan <i>et al.</i> , 2017]	55.3	54.8	55.1
[Yi <i>et al.</i> , 2024]	61.7	63.1	62.5
IOS [Wu <i>et al.</i> , 2023]	62.3	61.5	60.9
[Karimireddy <i>et al.</i> , 2021]	61.1	60.8	61.6
[Li <i>et al.</i> , 2020]	60.3	60.7	60.9
FedRBN [Hong <i>et al.</i> , 2023]	61.2	62.3	61.0
AdaShield-FL	73.1	72.8	73.5

Table 2: Comparison with prior methods on the M&Ms validation set in terms of dice score. Each FL method trained with **3 clients** on the perturbed data generated by each adversarial attack.

Model	Dice Score		
	FGSM	PGD	C&W
FedAvg [McMahan <i>et al.</i> , 2017]	47.3	47.1	46.9
[Yi <i>et al.</i> , 2024]	50.3	52.8	51.7
IOS [Wu <i>et al.</i> , 2023]	52.0	51.9	51.3
[Karimireddy <i>et al.</i> , 2021]	51.3	51.7	51.1
[Li <i>et al.</i> , 2020]	50.9	51.3	50.4
FedRBN [Hong <i>et al.</i> , 2023]	50.1	51.7	50.3
AdaShield-FL	61.3	63.7	61.9

Table 3: Comparison with prior methods on the M&Ms validation set in terms of dice score. Each FL method trained with **4 clients** on the perturbed data generated by each adversarial attack.

underscores the appropriateness of using three iterations to achieve optimal segmentation performance while maintaining robustness under adversarial attacks.

B.3 Number of Malicious Clients

Tables 2 and 3 demonstrate the performance of AdaShield-FL with 3 and 4 malicious clients on the M&Ms dataset. This analysis specifically measures the impact of increasing malicious clients on model accuracy, providing insights into the effectiveness of purification and detection mechanisms. Despite the increase in the number of malicious clients, AdaShield-FL shows an ability to maintain performance stability, which is a critical indicator of the model's robustness against data poisoning attacks in federated learning (FL) environments.

Feature Variation	Adversarial Attacks		
	FGSM	PGD	C&W
Δf_R	3.182	2.674	2.945
Δf_R (w/o \mathcal{L}_{ind})	4.866	5.131	5.298
Δf_S	9.233	8.415	8.193
Δf_S (w/o \mathcal{L}_{ind})	7.162	7.064	6.883

Table 4: MAE between features extracted from clean and perturbed data, with and without the **independence loss** \mathcal{L}_{ind} on the M&Ms dataset.

Model	Dice Score		
	FGSM	PGD	C&W
FedAvg [McMahan <i>et al.</i> , 2017]	51.2	50.3	50.6
[Yi <i>et al.</i> , 2024]	60.2	60.3	60.6
IOS [Wu <i>et al.</i> , 2023]	59.3	60.1	59.8
[Karimireddy <i>et al.</i> , 2021]	59.5	59.9	60.5
[Li <i>et al.</i> , 2020]	61.3	60.7	60.0
FedRBN [Hong <i>et al.</i> , 2023]	59.7	61.2	59.8
AdaShield-FL	68.3	68.2	68.7

Table 5: Comparison with prior methods on the **CTICH dataset** in terms of dice score. Each FL method trained on the perturbed data generated by each adversarial attack.

B.4 Independence Loss

Table 4 shows the mean absolute error (MAE) between two features, f_R and f_S , extracted from clean and perturbed data, both with and without the \mathcal{L}_{ind} loss term. The results demonstrate that incorporating \mathcal{L}_{ind} reduces perturbations affecting f_R , indicating enhanced adversarial robustness. This underscores the effectiveness of \mathcal{L}_{ind} in improving model robustness by minimizing feature variation in f_R while directing perturbations toward f_S .

B.5 CT Image Segmentation for other Modality

In Table 5, we utilized the "Computed Tomography Images for Intracranial Haemorrhage Detection and Segmentation" (CTICH) dataset [Hssayeni *et al.*, 2020], which is publicly available on PhysioNet. The dataset comprises 82 computed tomography (CT) scans from patients with traumatic brain injuries, including 36 scans where intracranial hemorrhage was diagnosed. Each CT scan consists of approximately 30 slices per patient. The average age of the patients is 27.8 years with a standard deviation of 19.5 years; 46 patients are male, and 36 are female. Moreover, Table 5 demonstrates the superior performance of AdaShield-FL compared to other FL methods on the CTICH dataset, a different modality involving CT images. The results indicate that AdaShield-FL consistently achieves superior Dice scores across all adversarial attacks, including FGSM [Goodfellow *et al.*, 2014], PGD [Madry *et al.*, 2017], and C&W [Carlini and Wagner, 2017], highlighting its robustness and generalizability to diverse medical imaging modalities. This result underscores the adaptability and broad applicability of AdaShield-FL to various imaging modalities beyond the original training domain.

B.6 Ablation Study

Table 6 presents an additional ablation study of the individual contributions of each module in the AdaShield-FL frame-

Differential Loss-based Attack Detection	Evolutionary DA	Malicious Client Purification	Dice Score
✓	✓	✓	82.6
✓	✓		74.9
✓		✓	81.3
	✓		71.8
✓			75.5
			67.8

Table 6: **Additional ablation study** for AdaShield-FL on the ACDC dataset perturbed by PGD attack in terms of dice coefficient score.

Differential Loss-based Attack Detection		Precision			Recall		
First-order Derivative Extractor	Second-order Derivative Extractor	FGSM	PGD	C&W	FGSM	PGD	C&W
✓		0.88	0.85	0.89	0.84	0.82	0.85
✓	✓	0.82	0.83	0.88	0.90	0.83	0.81

Table 7: Ablation study for differential loss-based attack detection performance in terms of recall and precision for detecting malicious clients on the M&Ms datasets perturbed by each adversarial attack.

work. The results demonstrate that each module significantly enhances the overall performance of AdaShield-FL. Furthermore, the combination of these modules operates synergistically to maximize segmentation performance, effectively addressing challenges posed by data poisoning attacks and data heterogeneity.

In addition, Table 7 compares the precision and recall of the differential loss-based attack detection framework under different adversarial attack scenarios. The results are presented for two configurations of the differential loss-based attack detection module: one using the first-order derivative extractor and the other combining the first-order and second-order derivative extractors. The configuration that utilizes both first-order and second-order derivative extractors demonstrates superior performance across all metrics. In contrast, the first-order-only configuration exhibits lower performance, with precision and recall. These results indicate that incorporating second-order derivatives enhances the robustness and accuracy of AdaShield-FL’s attack detection mechanism. Moreover, these results highlight that abnormal patterns such as oscillations and slow convergence speed are more likely to be caused by malicious clients than by the diversity of patients, vendors, and diseases.

C Visual Segmentation Results

Figures 3–6 provide qualitative segmentation results of segmentation results achieved by AdaShield-FL under diverse scenarios, highlighting its robustness and accuracy across different patient conditions.

Specifically, Figure 3 depicts segmentation results for patients with “hypertrophic cardiomyopathy”. The results of AdaShield-FL exhibit exceptional segmentation, similar to the ground truth, accurately capturing the thickness of myocardial regions without predicting false-positive segmentation in unrelated areas. This precise segmentation facilitates the accurate diagnosis of hypertrophic cardiomyopathy, a disease characterized by the thickening of the myocardium. In comparison, [Li *et al.*, 2020] and [Karimireddy *et al.*, 2021]

may lead to incorrect diagnoses as a result of inaccurate segmentation of the myocardial region.

Figure 4 also presents cases of “normal” cardiac patients. AdaShield-FL provides clean segmentation maps, preserving anatomical boundaries with high precision. FedAvg and [McMahan *et al.*, 2017] often provide artifacts from noise or fail to clearly delineate boundaries, resulting in false positive regions and [Yi *et al.*, 2024] also provide incorrect right ventricle regions, potentially diagnosing a normal patient as a patient who has another disease. AdaShield-FL avoids these issues to maintain robust performance.

Moreover, Figure 5 illustrates segmentation results for patients with an “abnormal right ventricle,” characterized by significant enlargement of the right ventricle. Unlike other methods, AdaShield-FL accurately captures the thickened right ventricle and the thinned myocardium regions. IOS [Wu *et al.*, 2023] and [Li *et al.*, 2020] are particularly inaccurate in predicting the right ventricle and myocardial regions, making it challenging to assess the thickness of these regions, which is a key characteristic of the disease. As a result, this method leads to inaccurate diagnoses, unlike AdaShield-FL.

Furthermore, Figure 6 illustrates segmentation maps from “dilated cardiomyopathy” patients, where the left ventricle is significantly enlarged. AdaShield-FL accurately provides the left ventricle without false-negative segmentation, unlike other methods. For instance, methods like FedRBN [Hong *et al.*, 2023] or IOS [Wu *et al.*, 2023] often struggle with false-negative segmentation of the left ventricle, incorrectly predicting part of the ventricle as background. This incorrect prediction of left ventricle regions has potentially diagnosed a patient with “dilated cardiomyopathy” as a normal patient. In contrast, AdaShield-FL ensures more reliable patient diagnoses.

Additionally, as shown in the Fig 7, we provide qualitative segmentation results for an additional modality. We note that CT images are transformed into the frequency domain using fast Fourier transform (FFT), where the transformed frequency domain corresponds to the k -space in magnetic resonance imaging (MRI). This domain representation allows AdaShield-FL to effectively separate the data, applying the same mechanisms as those used for k -space. Figure 7 illustrates segmentation results on the CTICH dataset for “traumatic brain injury” patients perturbed by FGSM. AdaShield-FL achieves accurate segmentation of intracranial hemorrhage regions, closely matching the ground truth, even under adversarial perturbations. FedAvg [McMahan *et al.*, 2017], [Yi *et al.*, 2024] and IOS [Wu *et al.*, 2023], which fail to consistently identify hemorrhage regions and often misclassify critical areas as background or noise. Moreover, FedRBN [Hong *et al.*, 2023] often makes incorrect predictions or misses parts of the hemorrhage, potentially leading to diagnostic errors. In contrast, AdaShield-FL avoids false positives and negatives by focusing on key hemorrhage regions, demonstrating robust segmentation performance and reliability in challenging adversarial conditions. This ensures greater accuracy and dependability for medical diagnoses in clinical scenarios involving traumatic brain injuries.

█ Left Ventricle █ Myocardium █ Right Ventricle

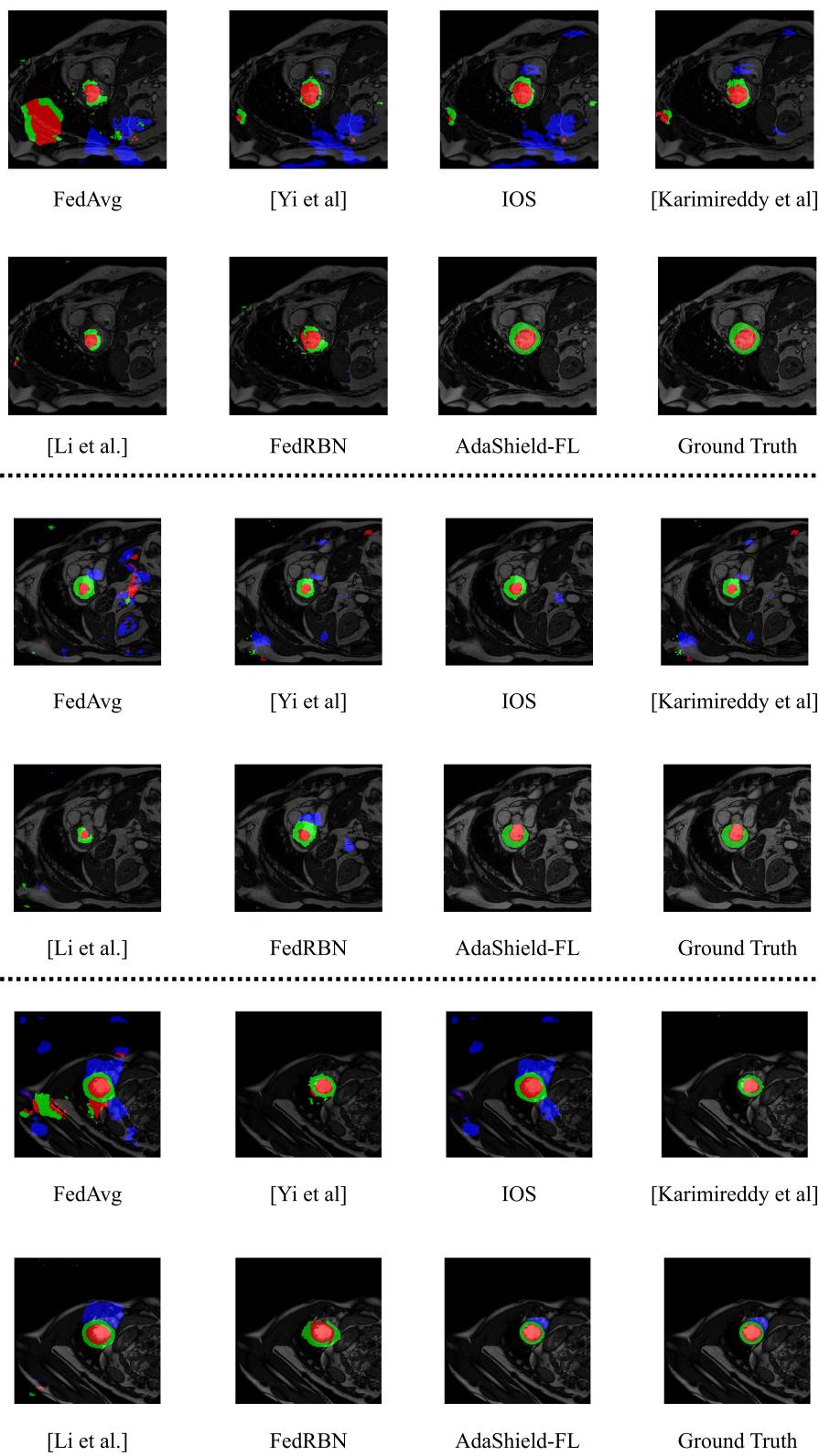


Figure 3: Qualitative segmentation results on M&Ms dataset perturbed by FGSM in “hypertrophic cardiomyopathy” patients.

█ Left Ventricle █ Myocardium █ Right Ventricle

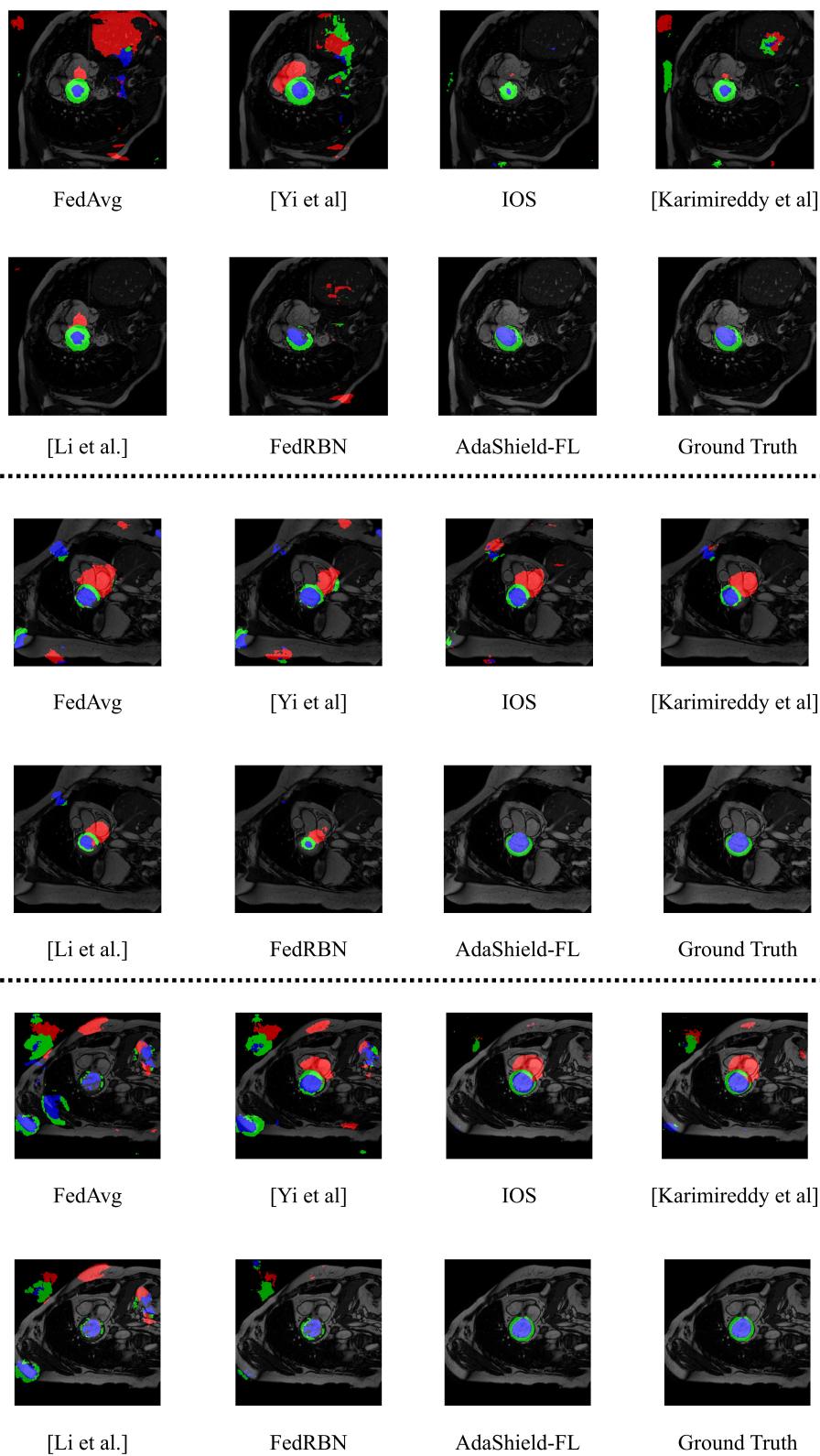
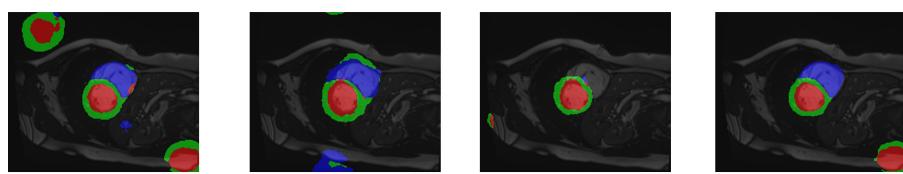


Figure 4: Qualitative segmentation results on M&Ms dataset perturbed by PGD in “normal” patients.

■ Left Ventricle ■ Myocardium ■ Right Ventricle

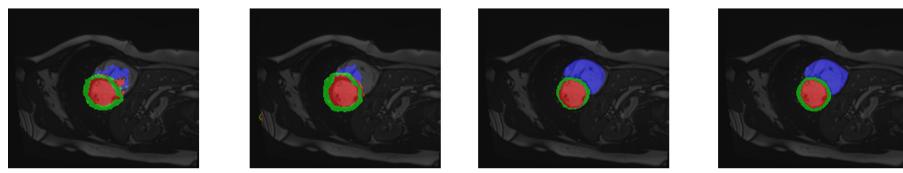


FedAvg

[Yi et al]

IOS

[Karimireddy et al]

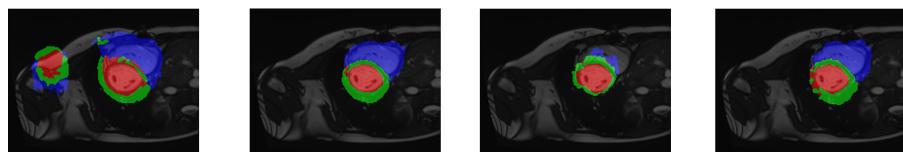


[Li et al.]

FedRBN

AdaShield-FL

Ground Truth

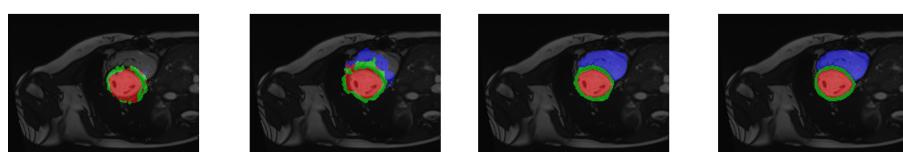


FedAvg

[Yi et al]

IOS

[Karimireddy et al]

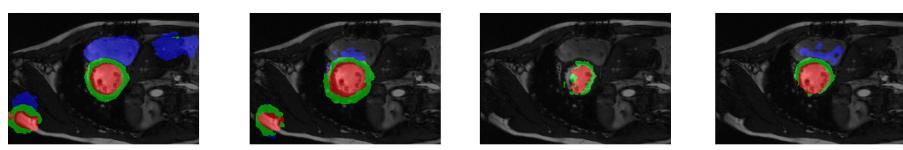


[Li et al.]

FedRBN

AdaShield-FL

Ground Truth

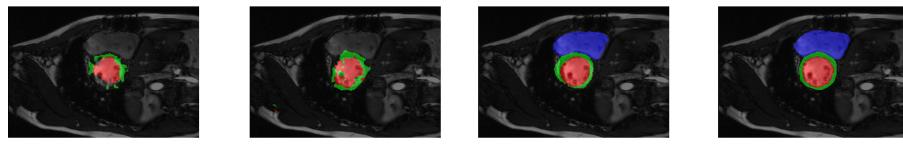


FedAvg

[Yi et al]

IOS

[Karimireddy et al]



[Li et al.]

FedRBN

AdaShield-FL

Ground Truth

Figure 5: Qualitative segmentation results on ACDC dataset perturbed by C&W in “**abnormal right ventricle**” patients.

█ Left Ventricle █ Myocardium █ Right Ventricle

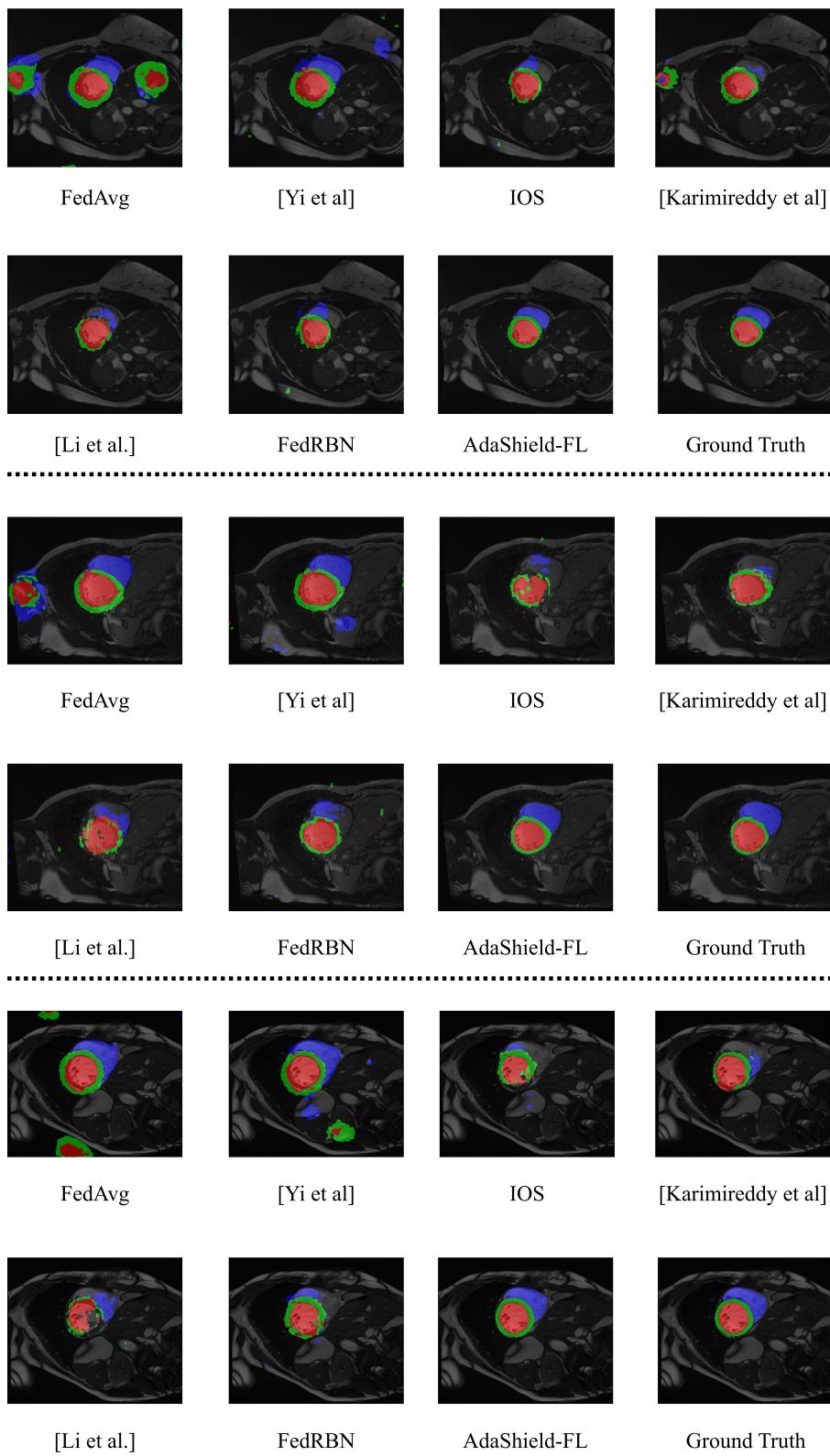


Figure 6: Qualitative segmentation results on ACDC dataset perturbed by PGD in “**dilated cardiomyopathy**” patients.

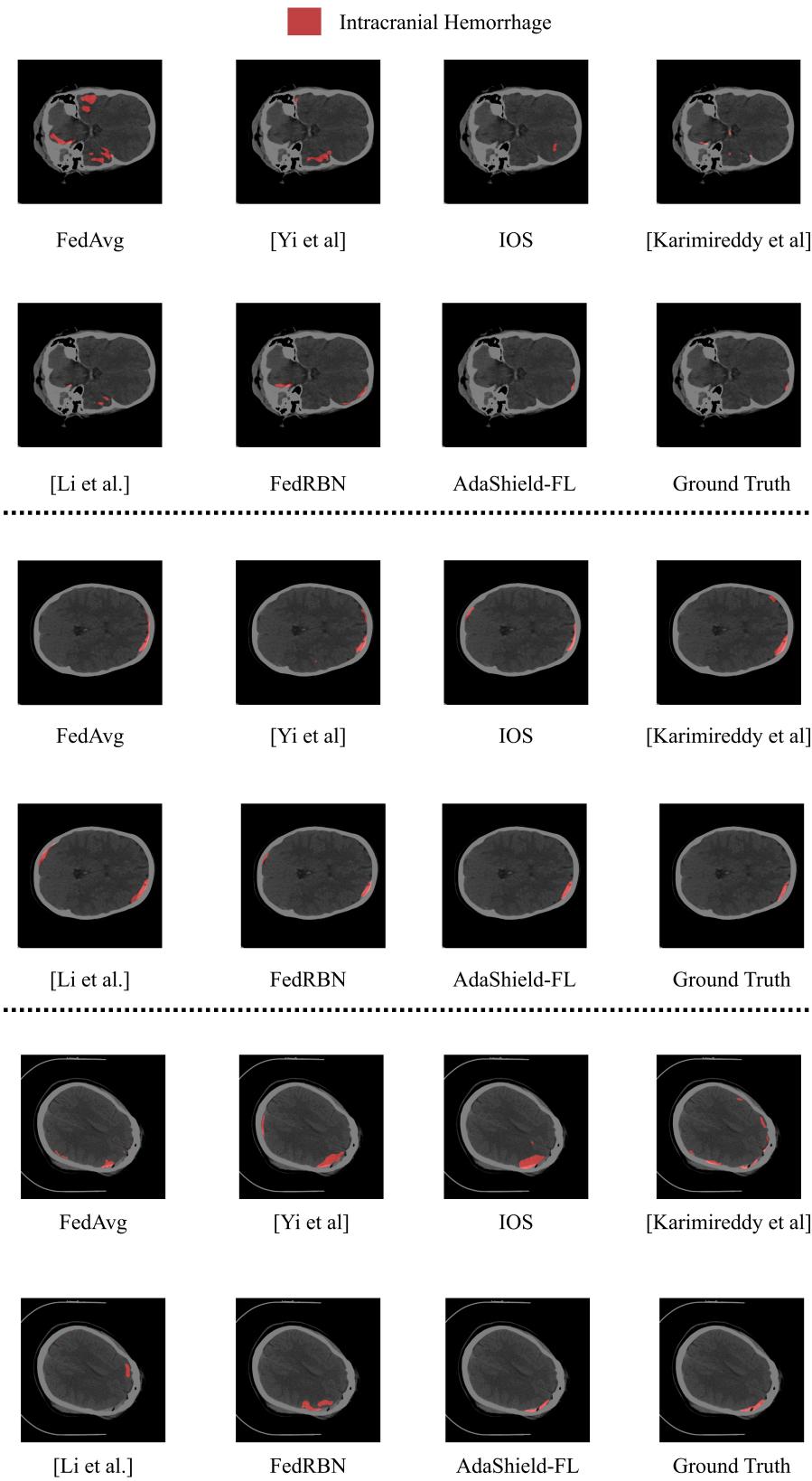


Figure 7: Qualitative segmentation results on **CTICH dataset** perturbed by FGSM in “**traumatic brain injury**” patients.

<p>214 D Adversarial Scenario</p> <p>215 D.1 Attacker’s Goal</p> <p>216 The external attacker generates adversarial examples by applying small, imperceptible perturbations to input images.</p> <p>217 These perturbations are generated using techniques such as FGSM [Goodfellow <i>et al.</i>, 2014], projected gradient descent</p> <p>218 (PGD) [Madry <i>et al.</i>, 2017], or Carlini and Wagner attack</p> <p>219 (C&W) [Carlini and Wagner, 2017], and are intended to result in incorrect segmentation maps. Moreover, the attacker</p> <p>220 exploits knowledge of the model’s parameters and gradients</p> <p>221 (assuming they have partial access to these) to optimize the</p> <p>222 perturbations for maximum impact. The result of these</p> <p>223 attacks is a reduced segmentation performance of the global</p> <p>224 model, which directly affects its clinical utility. In medical</p> <p>225 applications, such a degradation can lead to failure in identifying</p> <p>226 critical conditions, delayed treatment, or even financial</p> <p>227 fraud in insurance claims.</p> <p>228</p> <p>231 D.2 Attacker’s Capability</p> <p>232 The white-box scenario assumes that the attacker has unauthorized access to both clinical data and model parameters</p> <p>233 of certain clients. This access allows the attacker to inject</p> <p>234 data poisoning attacks. Specifically, the attacker can inject</p> <p>235 poisoned data, either by modifying input images or by manipulating output labels. This poisoned data corrupts the local</p> <p>236 training process and, when aggregated, negatively impacts</p> <p>237 the global model. The adversarial scenario assumes that the</p> <p>238 attacker can compromise more than 25% of the total client</p> <p>239 number.</p> <p>240</p> <p>242 D.3 Attacker’s Background Knowledge</p> <p>243 This work distinguishes between two types of clients. Benign</p> <p>244 clients are uncompromised clients that contribute clean</p> <p>245 data and trustworthy updates to the global model. They</p> <p>246 represent the majority of clients in the FL system and are crucial</p> <p>247 for maintaining the model’s integrity. Malicious clients have</p> <p>248 been manipulated by the external attacker and are responsible</p> <p>249 for injecting poisoned data or adversarial examples into</p> <p>250 the training process. The attacker controls these clients and</p> <p>251 uses them to execute attacks without being directly involved</p> <p>252 in the FL process. The attacker has access to the local datasets</p> <p>253 of compromised clients, which are used to generate poisoned</p> <p>254 data or adversarial examples. Since medical datasets are typically</p> <p>255 highly sensitive and diverse, this access poses a significant</p> <p>256 threat to the integrity of the global model.</p> <p>257</p> <p>E Limitation and Future Work</p> <p>258 E.1 Parameter Reduction in Reconstruction</p> <p>259 Decoder</p> <p>260 The reconstruction decoder in AdaShield-FL, while effective,</p> <p>261 involves a slightly higher number of parameters. Future</p> <p>262 work will focus on reducing the parameter count in the</p> <p>263 reconstruction decoder by employing a knowledge distillation</p> <p>264 approach. Specifically, the global reconstruction decoder</p> <p>265 will serve as the teacher model, while the student model will</p> <p>266 be designed as a lightweight local reconstruction decoder,</p> <p>267 tailored to specific client data distributions. Through this</p>	<p>method, the student model can mimic the performance of the teacher model while utilizing significantly fewer parameters, thereby optimizing efficiency and reducing computational overhead. However, this approach is not applicable to the segmentation decoder, as maintaining high segmentation performance is the primary priority. Reducing parameters in the segmentation decoder through knowledge distillation or similar methods could compromise its accuracy, making it unsuitable for this component of the framework. Thus, parameter optimization efforts will be focused exclusively on the reconstruction decoder while preserving the segmentation decoder’s robust performance.</p> <p>268</p> <p>269</p> <p>270</p> <p>271</p> <p>272</p> <p>273</p> <p>274</p> <p>275</p> <p>276</p> <p>277</p> <p>278</p> <p>279</p> <p>280</p> <p>281</p> <p>282</p> <p>283</p> <p>284</p> <p>285</p> <p>286</p> <p>287</p> <p>288</p> <p>289</p> <p>290</p> <p>291</p> <p>292</p> <p>293</p> <p>294</p> <p>295</p> <p>296</p> <p>297</p> <p>298</p> <p>299</p> <p>300</p> <p>301</p> <p>302</p> <p>303</p> <p>304</p> <p>305</p> <p>306</p> <p>307</p> <p>308</p> <p>309</p> <p>310</p> <p>311</p> <p>312</p> <p>313</p> <p>314</p> <p>315</p> <p>316</p> <p>317</p> <p>318</p> <p>319</p> <p>320</p> <p>321</p>
<p>[Campello <i>et al.</i>, 2021] Victor M Campello, Polyxeni Gkontra, Cristian Izquierdo, Carlos Martin-Isla, Alireza Sajjadi, Peter M Full, Klaus Maier-Hein, Yao Zhang, Zhiqiang He, Jun Ma, et al. Multi-centre, multi-vendor and multi-disease cardiac segmentation: the m&ms challenge. <i>IEEE Transactions on Medical Imaging</i>, 40(12):3543–3554, 2021.</p> <p>[Carlini and Wagner, 2017] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In <i>2017 ieee symposium on security and privacy (sp)</i>, pages 39–57. Ieee, 2017.</p> <p>[Goodfellow <i>et al.</i>, 2014] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. <i>arXiv preprint arXiv:1412.6572</i>, 2014.</p>	

- 322 [Hong *et al.*, 2023] Junyuan Hong, Haotao Wang,
323 Zhangyang Wang, and Jiayu Zhou. Federated ro-
324 bustness propagation: sharing adversarial robustness in
325 heterogeneous federated learning. In *Proceedings of the*
326 *AAAI Conference on Artificial Intelligence*, volume 37,
327 pages 7893–7901, 2023.
- 328 [Hssayeni *et al.*, 2020] Murtadha Hssayeni, M Croock,
329 A Salman, H Al-khafaji, Z Yahya, and B Ghoraani. Com-
330 puted tomography images for intracranial hemorrhage
331 detection and segmentation. *Intracranial hemorrhage*
332 *segmentation using a deep convolutional model. Data*,
333 5(1):14, 2020.
- 334 [Karimireddy *et al.*, 2021] Sai Praneeth Karimireddy, Lie
335 He, and Martin Jaggi. Learning from history for byzan-
336 tine robust optimization. In *International Conference on*
337 *Machine Learning*, pages 5311–5319. PMLR, 2021.
- 338 [Li *et al.*, 2020] Suyi Li, Yong Cheng, Wei Wang, Yang
339 Liu, and Tianjian Chen. Learning to detect malicious
340 clients for robust federated learning. *arXiv preprint arXiv:2002.00211*, 2020.
- 342 [Madry *et al.*, 2017] Aleksander Madry, Aleksandar
343 Makelov, Ludwig Schmidt, Dimitris Tsipras, and
344 Adrian Vladu. Towards deep learning models resistant
345 to adversarial attacks. *arXiv preprint arXiv:1706.06083*,
346 2017.
- 347 [McMahan *et al.*, 2017] Brendan McMahan, Eider Moore,
348 Daniel Ramage, Seth Hampson, and Blaise Aguera y Ar-
349 cas. Communication-efficient learning of deep networks
350 from decentralized data. In *Artificial intelligence and*
351 *statistics*, pages 1273–1282. PMLR, 2017.
- 352 [Wu *et al.*, 2023] Zhaoxian Wu, Tianyi Chen, and Qing
353 Ling. Byzantine-resilient decentralized stochastic opti-
354 mization with robust aggregation rules. *IEEE transactions*
355 *on signal processing*, 2023.
- 356 [Yi *et al.*, 2024] Yuhao Yi, Ronghui You, Hong Liu,
357 Changxin Liu, Yuan Wang, and Jiancheng Lv. Near-
358 optimal resilient aggregation rules for distributed learning
359 using 1-center and 1-mean clustering with outliers. In *Pro-*
360 *ceedings of the AAAI Conference on Artificial Intelligence*,
361 volume 38, pages 16469–16477, 2024.
- 362 [Zhou *et al.*, 2021] Hong-Yu Zhou, Jiansen Guo, Yinghao
363 Zhang, Lequan Yu, Liansheng Wang, and Yizhou Yu. nn-
364 former: Interleaved transformer for volumetric segmenta-
365 tion. *arXiv preprint arXiv:2109.03201*, 2021.