# FREQUENCY DOMAIN ADVERSARIAL TRAINING FOR ROBUST VOLUMETRIC MEDICAL SEGMENTATION

Asif Hanif, Muzammal Naseer, Salman Khan, Mubarak Shah, Fahad Shahbaz Khan

Mohamed Bin Zayed University of Artificial Intelligence (MBZUAI), UAE

University of Central Florida (UCF), USA

Linköping University, Sweden

# Introduction

- Semantic segmentation of organs, anatomical structures, or anomalies in medical images (e.g. CT or MRI scans) remains one of the fundamental tasks in medical image analysis.

- Deep learning and volumetric medical image segmentation

- Adversarial vulnerability of volumetric medical image segmentation

- Adversarial robustness of medical image models

- Understanding the vulnerability and robustness of medical image models

# Contributions

■ Volumetric Adversarial Frequency Attack (VAFA)

■ Volumetric Adversarial Frequency Training (VAFT)

# Voxel-domain Attacks  vs. Frequency-domain Attacks

- Voxel-domain Attacks: Directly perturb input space

- Frequency-domain Attacks: Preturb the frequency-domain representation of an image e.g. DCT

# Method

- Volumetric Adversarial Frequency domain Attack (**VAFA**)

- Volumetric Adversarial Frequency domain Training (**VAFT**)

# Volumetric Adversarial Frequency Domain Attack (VAFA)

# Volumetric Adversarial Frequency Domain Attack (VAFA)

$$\mathbf{x} \mapsto \mathcal{D}(\mathbf{x}) \mapsto \underbrace{\varphi(\mathcal{D}(\mathbf{x}), \mathbf{q})}_{\substack{\text{quantization,} \\ \text{rounding and} \\ \text{de-quantization}}} \mapsto \mathcal{D}_I(\varphi(\cdot)) \mapsto \mathbf{x}'$$

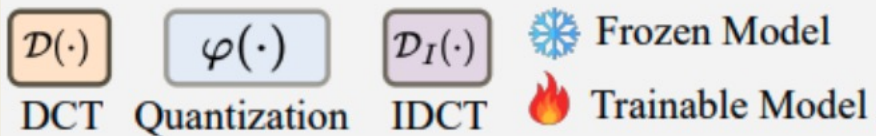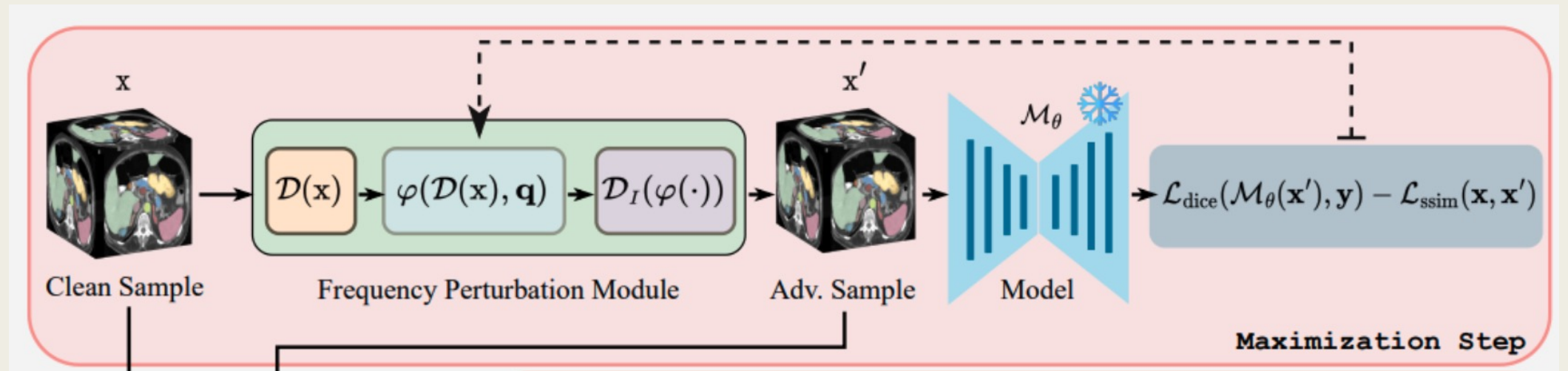# Volumetric Adversarial Frequency Domain Attack (VAFA)

$$\varphi(\mathcal{D}(\mathbf{x}), \mathbf{q}) := \lfloor \frac{\mathcal{D}(\mathbf{x})}{\mathbf{q}} \rfloor \odot \mathbf{q}$$

$$\mathbf{q} \in \mathbb{Z}^{h \times w \times d}$$

# Volumetric Adversarial Frequency Domain Attack (VAFA)

$$\underset{\mathbf{q}}{\text{maximize}} \quad \mathcal{L}_{\text{dice}}(\mathcal{M}_\theta(X'), Y) - \mathcal{L}_{\text{ssim}}(X, X')$$

$$\text{s.t.} \quad \|\mathbf{q}\|_\infty \leq q_{\text{max}},$$

# Volumetric Adversarial Frequency Domain Attack (VAFA)

# Volumetric Adversarial Frequency Domain Attack (VAFA)

---

**Algorithm 1** Volumetric Adversarial Frequency Attack (**VAFA**)

---

1: Number of Steps: $T$, Quantization Threshold: $q_{max}$
2: **Input:** $X \in \mathbb{R}^{H \times W \times D}, Y \in \{0, 1\}^{\mathrm{NumClass} \times H \times W \times D}$    **Output:** $X' \in \mathbb{R}^{H \times W \times D}$
3: **function VAFA**(X,Y)
4:      $\mathbf{q}_i \leftarrow \mathbf{1}$      $\forall\, i \in \{1, 2, \ldots, n\}$      ▷ Initialize all *quantization tables* with ones.
5:      **for** $t \leftarrow 1$ to $T$ **do**
6:          $\{\mathbf{x}_i\}_{i=1}^n \leftarrow \mathrm{Split}(X)$      ▷ Split X into 3D patches of size $(h \times w \times d)$
7:          $\mathbf{x}_i' \leftarrow \mathcal{D}_I\big(\, \varphi(\mathcal{D}(\mathbf{x}_i), \mathbf{q}_i)\,\big)$    $\forall\, i \in \{1, 2, \ldots, n\}$      ▷ Frequency Perturbation
8:          $X' \leftarrow \mathrm{Merge}(\{\mathbf{x}_i'\}_{i=1}^n)$      ▷ Merge all adversarial patches to form $X'$
9:          $\mathcal{L}(X, X', Y) = \mathcal{L}_{\mathrm{dice}}(\mathcal{M}_\theta(X'), Y) - \mathcal{L}_{\mathrm{ssim}}(X, X')$
10:        $\mathbf{q}_i \leftarrow \mathbf{q}_i + \mathrm{sign}(\nabla_{\mathbf{q}_i}\mathcal{L})$      $\forall\, i \in \{1, 2, \ldots, n\}$
11:        $\mathbf{q}_i \leftarrow \mathrm{clip}(\mathbf{q}_i,\ \min=1,\ \max=q_{max})$      $\forall\, i \in \{1, 2, \ldots, n\}$
12:      **end for**
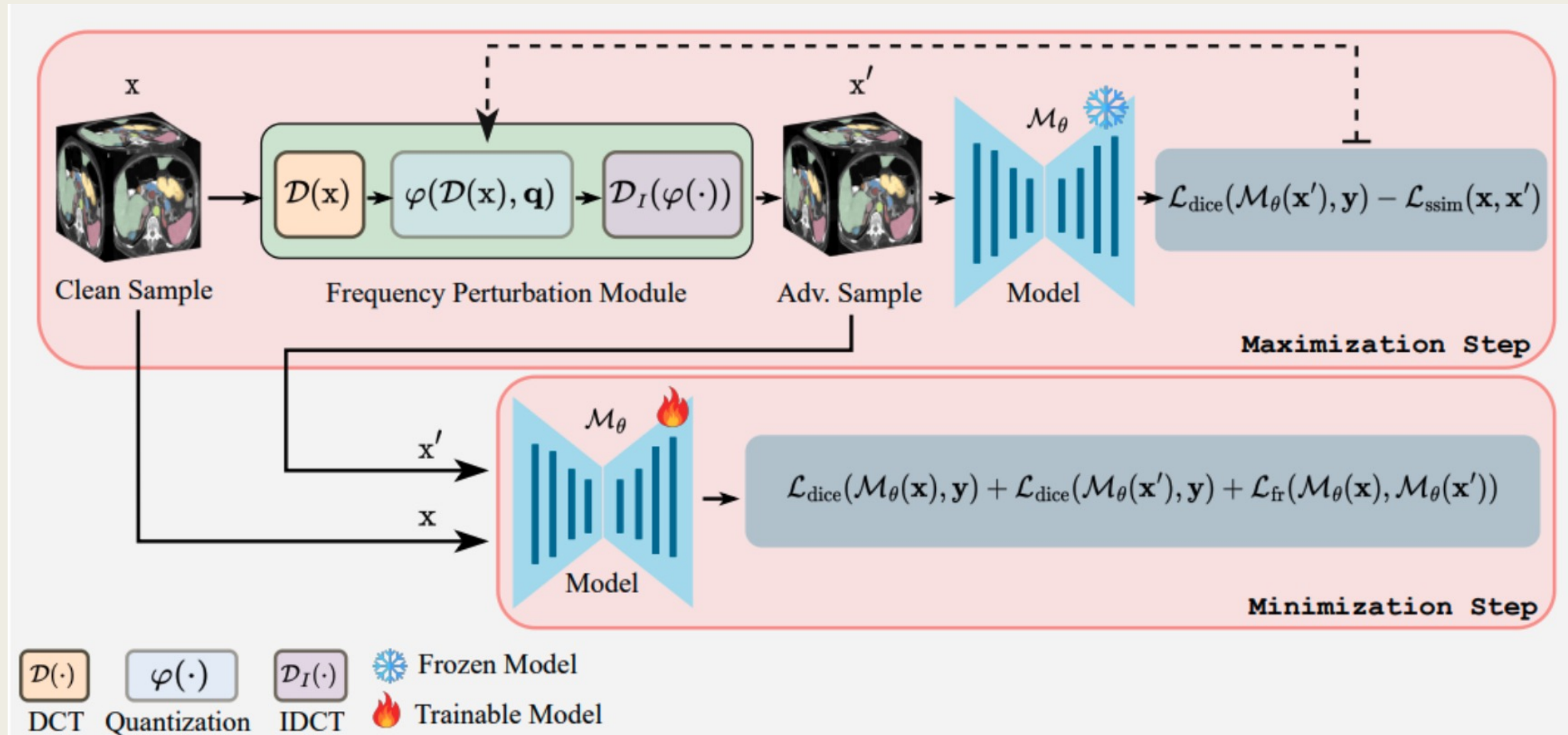13: **end function**
14: **Return** $X'$

# Volumetric Adversarial Frequency Domain Training (VAFT)

# Volumetric Adversarial Frequency Domain Training (VAFT)

$$\underset{\theta}{\text{minimize}} \ \mathcal{L}_{\text{dice}}(\mathcal{M}_\theta(X), Y) + \mathcal{L}_{\text{dice}}(\mathcal{M}_\theta(X'), Y) + \mathcal{L}_{\text{fr}}(\mathcal{M}_\theta(X), \mathcal{M}_\theta(X')),$$

$$\mathcal{L}_{\text{fr}}(\mathcal{M}_\theta(X), \mathcal{M}_\theta(X')) = \|\mathcal{D}(\mathcal{M}_\theta(X)) - \mathcal{D}(\mathcal{M}_\theta(X'))\|_1$$

# Volumetric Adversarial Frequency Domain Training (VAFT)

# Volumetric Adversarial Frequency Domain Training (VAFT)

---

**Algorithm 2** Volumetric Adversarial Frequency Training (**VAFT**)

---

1: Train Dataset: $\mathcal{X} = \{(X_i, Y_i)\}_{i=1}^{N}$, $X_i \in \mathbb{R}^{H \times W \times D}$, $Y_i \in \{0, 1\}^{\text{NumClass} \times H \times W \times D}$

2: NumSamples=$N$, BatchSize=$B$, Target Model: $\mathcal{M}_\theta$, AT Robust Model: $\mathcal{M}_{\bullet}$

3: **for** $i \leftarrow 1$ to NumEpochs **do**

4:     **for** $j \leftarrow 1$ to $\lfloor N/B \rfloor$ **do**

5:         Sample a mini-batch $\mathcal{B} \subseteq \mathcal{X}$ of size $B$

6:         $X' \leftarrow \mathbf{VAFA}(X, Y) \quad \forall (X, Y) \in \mathcal{B}$     ▷ Adv. Freq. Attack on clean images.

7:         $\mathcal{L} = \mathcal{L}_{\text{dice}}(\mathcal{M}_\theta(X), Y) + \mathcal{L}_{\text{dice}}(\mathcal{M}_\theta(X'), Y) + \mathcal{L}_{\text{fr}}(\mathcal{M}_\theta(X), \mathcal{M}_\theta(X'))$

8:         Backward pass and update $\mathcal{M}_\theta$

9:     **end for**

10: **end for**

11: $\mathcal{M}_{\bullet} \leftarrow \mathcal{M}_\theta$         ▷ AT robust model after training completion.
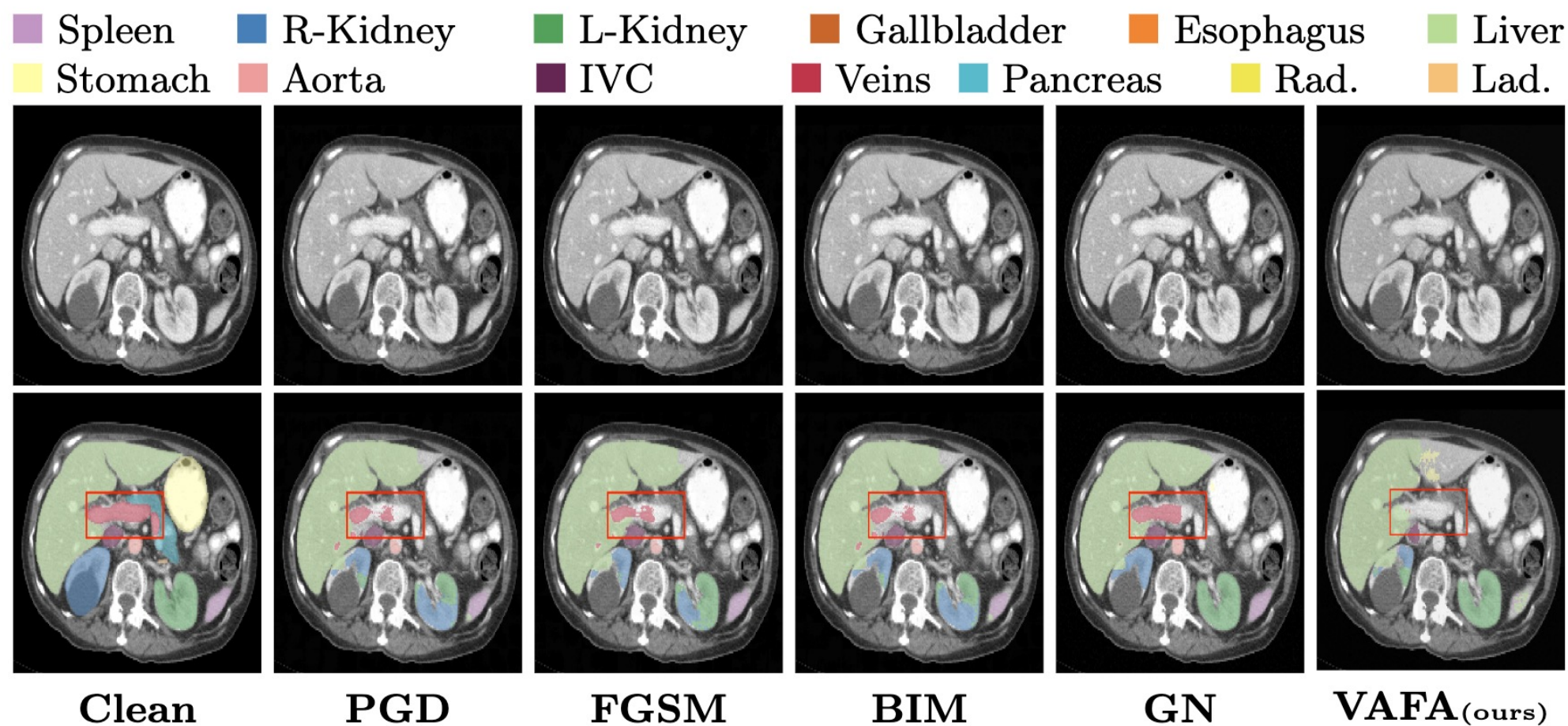
12: **Return** $\mathcal{M}_{\bullet}$

---

# Results

- **Segmentation Models:** UNETR and UNETR++

- **Datasets:** Synapse and ACDC

- **Baseline Attacks:** PGD, FGSM, BIM, GN

- **Evaluation Metrics:** Dice Score, HD95 Distance, LPIPS

- **Programming Framework:** Pytorch

# Volumetric Adversarial Frequency Domain Attack (VAFA)

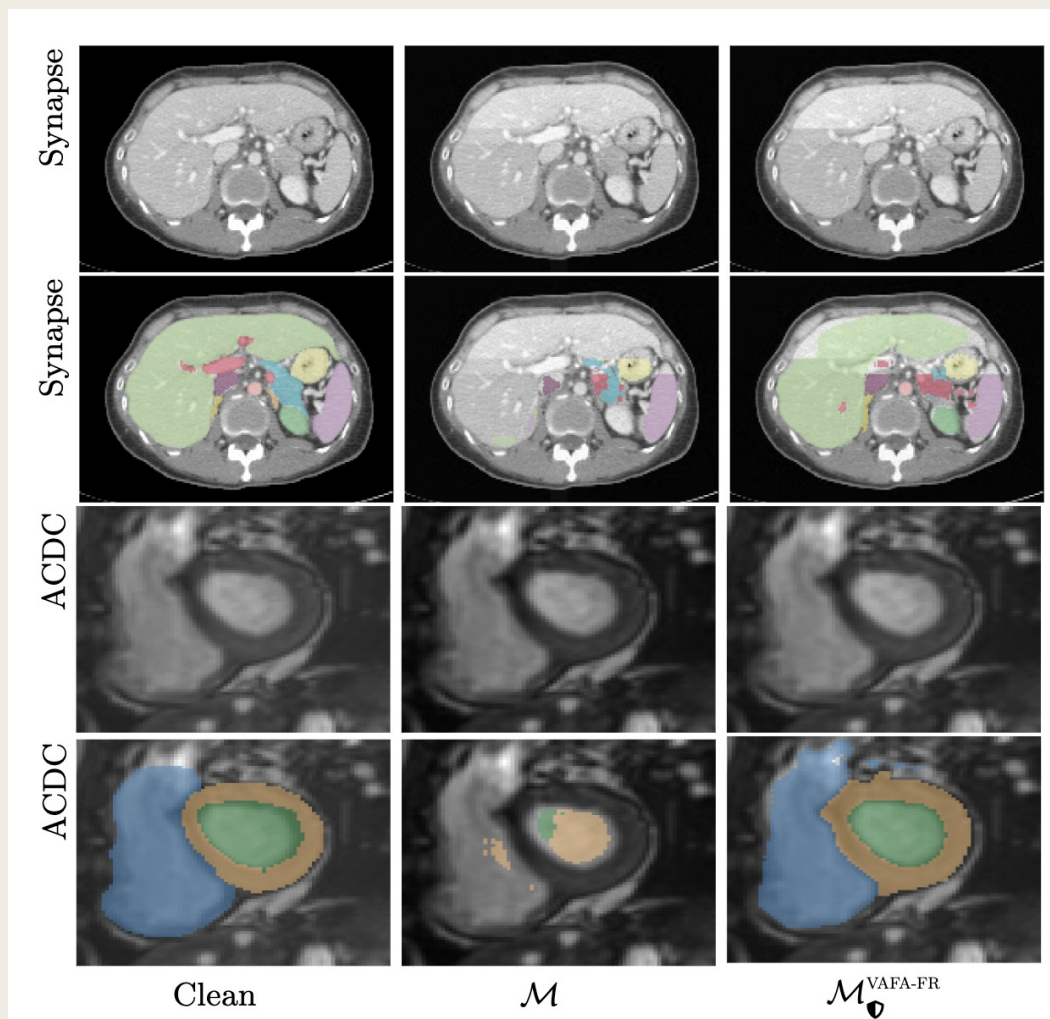| Models → Attacks ↓ | UNETR | | | UNETR++ | | |
|---|---|---|---|---|---|---|
| | DSC↓ | HD95↑ | LPIPS↑ | DSC↓ | HD95↑ | LPIPS↑ |
| Clean Images | 74.3 | 14.0 | - | 84.7 | 12.7 | - |
| PGD $(\epsilon = 4/8)$ | 62.7/50.8 | 40.4/64.5 | **98.9**/95.3 | 77.5/67.1 | 48.1/78.3 | 95.7/85.1 |
| FGSM $(\epsilon = 4/8)$ | 62.8/53.9 | 34.8/48.7 | 98.8/94.7 | 73.1/67.1 | 37.3/43.2 | 94.7/82.2 |
| BIM $(\epsilon = 4/8)$ | 62.8/50.7 | 39.9/**65.8** | 98.8/95.3 | 77.3/66.8 | 46.6/78.1 | **95.8**/85.3 |
| GN $(\sigma = 4/8)$ | 74.2/73.9 | 17.0/15.4 | 97.7/91.1 | 84.7/84.3 | 12.3/13.4 | 93.3/78.2 |
| VAFA $(q_{max} = 20/30)$ | **32.2/29.8** | **57.6**/59.9 | 97.5/**96.9** | **45.3/39.3** | **73.9/85.2** | 94.2/**94.7** |

# Volumetric Adversarial Frequency Domain Attack (VAFA)

# Volumetric Adversarial Frequency Domain Training (VAFT)

| Attacks → Models ↓ | UNETR | | | | | UNETR++ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Clean | PGD | FGSM | BIM | VAFA | Clean | PGD | FGSM | BIM | VAFA |
| Synapse $\mathcal{M}^{PGD}$ | 73.47 | 65.53 | 65.68 | 65.51 | 42.47 | 75.43 | 67.81 | 67.82 | 67.80 | 38.22 |
| $\mathcal{M}^{FGSM}$ | 72.44 | 64.80 | 66.31 | 64.76 | 39.02 | 81.06 | 73.84 | 74.76 | 73.77 | 37.48 |
| $\mathcal{M}^{BIM}$ | 75.12 | 67.78 | 68.32 | 67.78 | 45.97 | 74.80 | 67.58 | 67.46 | 67.57 | 35.72 |
| $\mathcal{M}^{GN}$ | 73.17 | 61.40 | 61.77 | 61.29 | 30.00 | 80.05 | 76.23 | 70.96 | 74.51 | 41.44 |
| $\mathcal{M}^{VAFA}$ | 74.67 | 64.83 | 65.49 | 64.73 | 66.31 | 81.88 | 69.09 | 65.40 | 68.90 | 76.47 |
| $\mathcal{M}^{VAFA-FR}$ | **75.66** | 65.90 | 66.79 | 65.83 | 66.33 | **82.65** | 70.61 | 67.00 | 70.41 | 78.19 |
| ACDC $\mathcal{M}^{VAFA}$ | 81.95 | 60.77 | 68.16 | 60.75 | 69.76 | 89.00 | 76.28 | 80.41 | 76.56 | 88.45 |
| $\mathcal{M}^{VAFA-FR}$ | **83.44** | 60.63 | 69.33 | 60.61 | 73.05 | **91.36** | 85.42 | 87.42 | 83.90 | 91.23 |

# Volumetric Adversarial Frequency Domain Training (VAFT)

# Github Repository

https://github.com/asif-hanif/vafa

# Thank you !

# VAFA – Ablation Study

- Impact of quantization threshold

- Impact of steps

- Impact of patch size

# VAFA – Ablation Study
# Impact of Quantization Threshold: $q_{max}$

| $q_{max}$ | DSC↓ | HD95↑ | LPIPS↑ |
|---|---|---|---|
| - | 74.31 | 14.03 | - |
| 10 | 65.95 | 26.25 | **99.10** |
| 20 | 56.24 | 35.92 | 98.70 |
| 30 | 50.96 | 44.09 | 98.33 |
| 40 | 49.58 | 43.66 | 97.90 |
| 60 | 48.83 | 44.55 | 96.60 |
| 80 | **48.76** | **45.30** | 94.50 |

# VAFA – Ablation Study
# Impact of Steps

| Steps | DSC↓ | HD95↑ | LPIPS↑ |
|:-----:|:-----:|:-----:|:------:|
| - | 74.31 | 14.03 | - |
| 10 | 61.33 | 33.20 | **98.85** |
| 20 | 56.24 | 35.92 | 98.70 |
| 30 | 54.37 | 38.00 | 98.64 |
| 40 | 53.31 | 37.76 | 98.59 |
| 50 | 52.97 | **39.23** | 98.54 |
| 60 | **52.25** | 39.19 | 98.52 |

# VAFA – Ablation Study
# Impact of Patch-Size

| Patch Size | DSC↓ | HD95↑ | LPIPS↑ |
|---|---|---|---|
| - | 74.31 | 14.03 | - |
| ( 4 × 4 × 4) | 63.48 | 32.63 | **98.90** |
| ( 8 × 8 × 8) | 56.24 | 35.92 | 98.70 |
| (16 × 16 × 16) | 41.30 | 45.98 | 98.14 |
| (32 × 32 × 32) | 32.40 | 56.64 | 97.49 |
| (48 × 48 × 48) | 28.19 | **66.08** | 97.16 |
| (96 × 96 × 96) | **28.08** | 59.09 | 96.47 |