

Training-Free Adversarial Robustness in Computational MRI

Mahdi Saberi^{1,2}, Chi Zhang^{1,2}, Mehmet Akçakaya^{1,2}

¹ Department of Electrical and Computer Engineering, University of Minnesota, MN, USA

² Center for Magnetic Resonance Research, University of Minnesota, MN, USA.

- In MRI, raw measurements are collected in frequency domain (known as k-space)
- Accelerated MRI techniques, acquire subsampled data:

- Most modern reconstruction solves a regularized least squares problem

$$y_{\Omega} = E_{\Omega}x + n$$

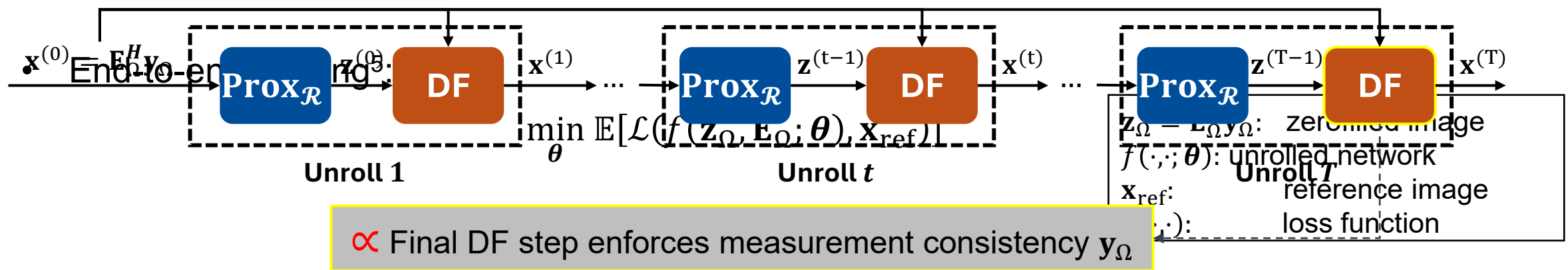
Acquired sub-sampled data \rightarrow y_{Ω} \leftarrow Measurements noise

Multi-coil encoding matrix \rightarrow E_{Ω}

$$\arg \min_x \underbrace{\|y_{\Omega} - E_{\Omega}x\|_2^2}_{\text{Data Fidelity (DF)}} + \underbrace{\mathcal{R}(x)}_{\text{Regularizer}}$$

conventionally solved iteratively, by alternating between DF and regularization¹

- Physics-driven deep learning (PD-DL) employ algorithm unrolling^{2,3}
- Unrolling variable splitting with quadratic penalty (VSQP)¹, as MoDL⁴



Adversarial Attacks in MRI

- PD-DL MRI reconstructions are vulnerable to adversarial attacks^{1,2,3}
- Small perturbations \rightarrow large reconstruction errors^{1,2,3}
- Worst-case perturbations in two domains:

- Image domain (ℓ_∞ attacks):

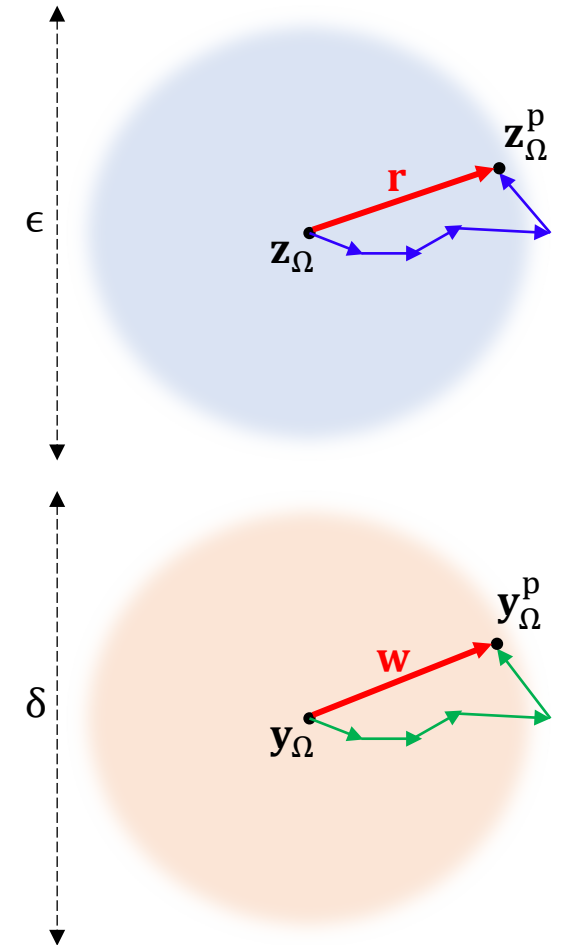
$$\arg \max_{\mathbf{r}: \|\mathbf{r}\|_\infty < \epsilon} \mathcal{L}(f(\mathbf{z}_\Omega + \mathbf{r}, \mathbf{E}_\Omega; \boldsymbol{\theta}), f(\mathbf{z}_\Omega, \mathbf{E}_\Omega; \boldsymbol{\theta}))$$

- k-space domain (ℓ_2 attacks):

$$\arg \max_{\mathbf{w}: \|\mathbf{w}\|_2 < \delta} \mathcal{L}(f(\mathbf{E}_\Omega^H(\mathbf{y}_\Omega + \mathbf{w}), \mathbf{E}_\Omega; \boldsymbol{\theta}), f(\mathbf{E}_\Omega^H \mathbf{y}_\Omega, \mathbf{E}_\Omega; \boldsymbol{\theta}))$$

- Typically solved by:

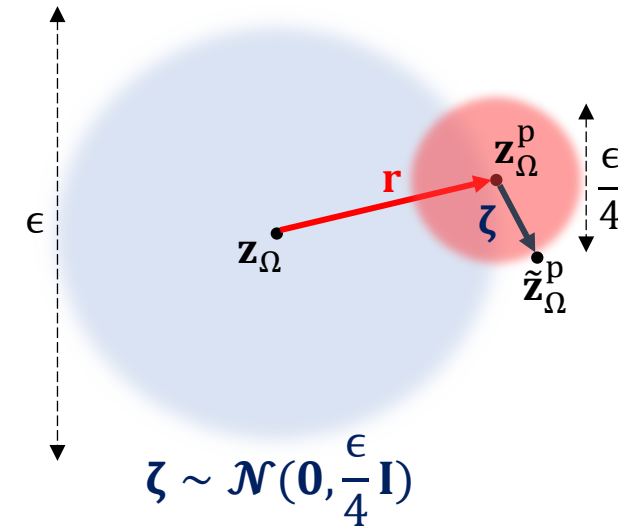
- Fast gradient sign method (FGSM)⁴
- Projected gradient descent (PGD)⁵



Motivation – Why Adversary in MRI?

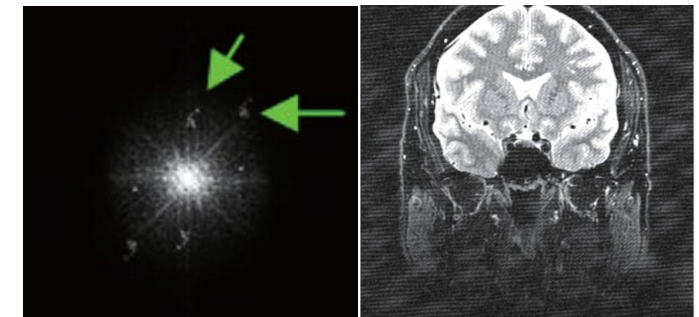
- **Worse-case perturbations – Non-zero probability,**

- Practical¹ and theoretical² evidence
- Sample from the perturbed input,
 - Another failed reconstruction¹
- Sample Gaussian noise from perturbed input,
 - Instability with non-zero probability²



- **Herringbone artifacts,**

- Other MRI measurements perturbations: Body motion³, Hardware issues⁴
- Electromagnetic spikes:
 - Gradient power fluctuation
 - Inadequate room shielding
- Impulse noise in k-space → Herringbone artifact in image domain^{5,6}



- **Training-based defenses:**

- Adversarial training (**AT**)^{1,2}
- Smooth unrolling (**SMUG**)³

- **Training-free defenses:**

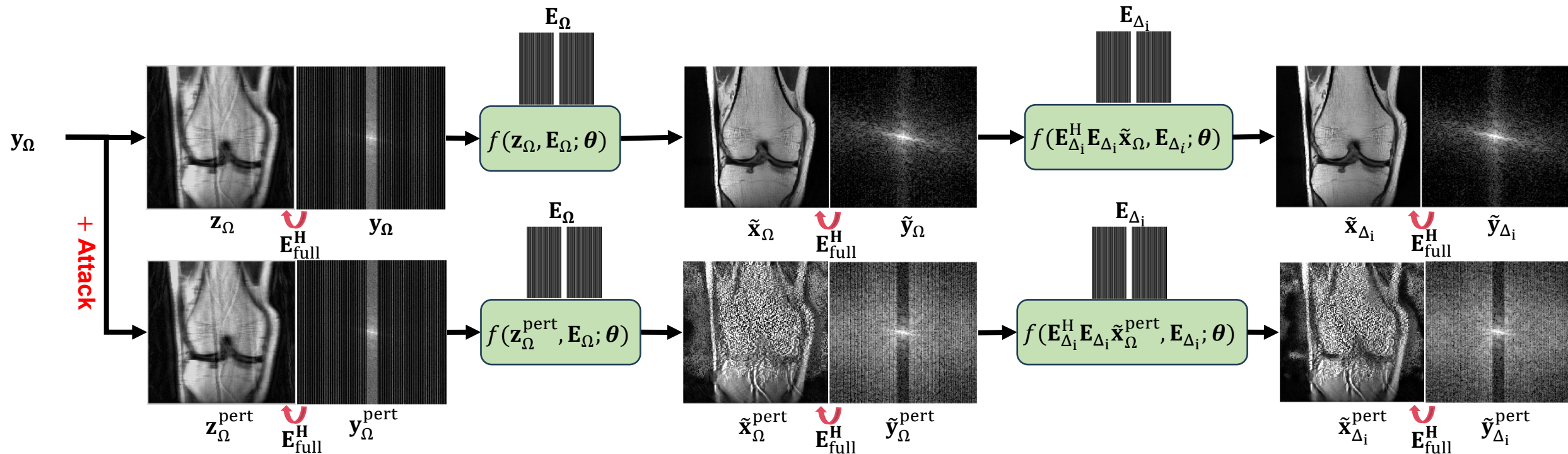
- Randomized smoothing (**RS**)^{3,4}
- **JPEG** compression⁵
- Total variation (**TV**)⁶ denoising

Attack Propagation in Simulated k-space

- A well-trained model generalizes to undersampling patterns drawn from similar distribution¹
- For Ω , Δ_i drawn with same acceleration rate, and number of ACS lines:

$$\tilde{\mathbf{x}}_{\Omega} = f(\mathbf{z}_{\Omega}, \mathbf{E}_{\Omega}; \boldsymbol{\theta})$$

$$\tilde{\mathbf{x}}_{\Delta_i} = f(\mathbf{E}_{\Delta_i}^H \mathbf{E}_{\Delta_i} \tilde{\mathbf{x}}_{\Omega}, \mathbf{E}_{\Delta_i}; \boldsymbol{\theta})$$

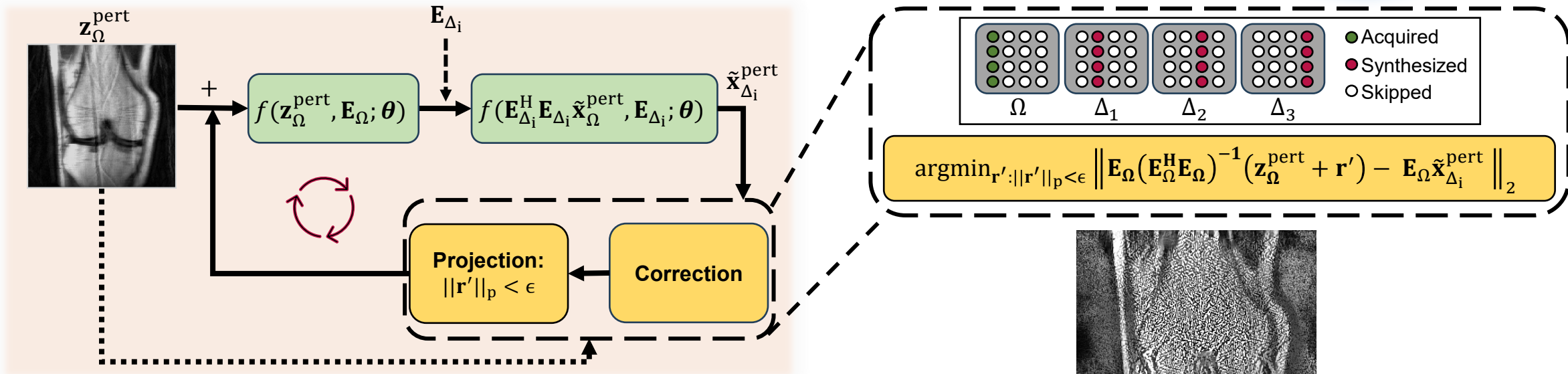
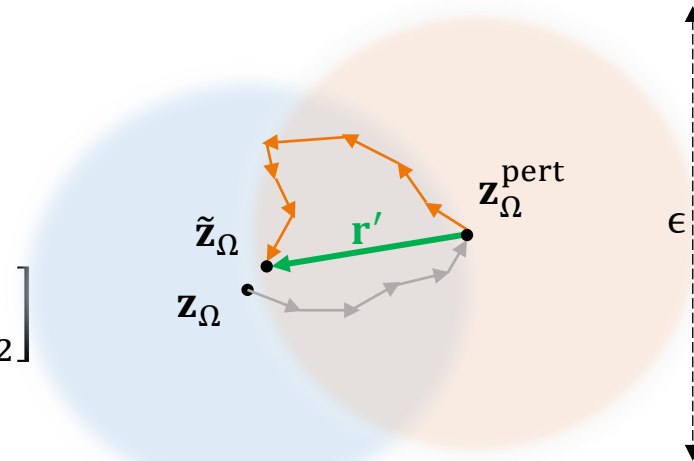


Attack Mitigation with Cyclic Consistency

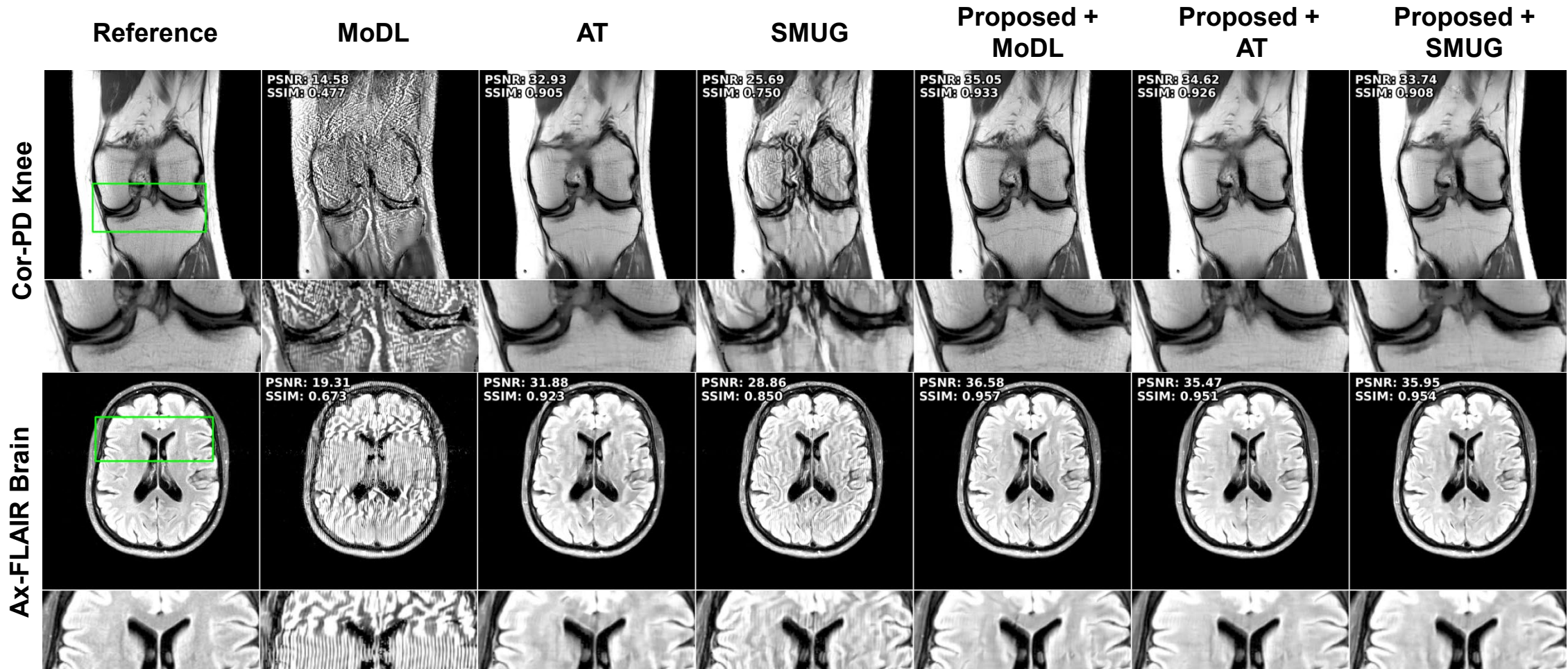
- Adversarial attack: Small perturbation within a ball around clean input
- **Goal:** reverse this process
- We devise a novel objective based on **cyclic consistency**:

$$\arg \min_{\mathbf{r}': \|\mathbf{r}'\|_p < \epsilon} \mathbb{E}_{\Delta} \left[\left\| (\mathbf{E}_{\Omega}^H)^{\dagger} (\mathbf{z}_{\Omega}^{\text{pert}} + \mathbf{r}') - \mathbf{E}_{\Omega} f(\mathbf{E}_{\Delta_i}^H (\mathbf{E}_{\Delta_i} f(\mathbf{z}_{\Omega}^{\text{pert}} + \mathbf{r}', \mathbf{E}_{\Omega}; \theta) + \tilde{\mathbf{n}}), \mathbf{E}_{\Delta_i}; \theta) \right\|_2 \right]$$

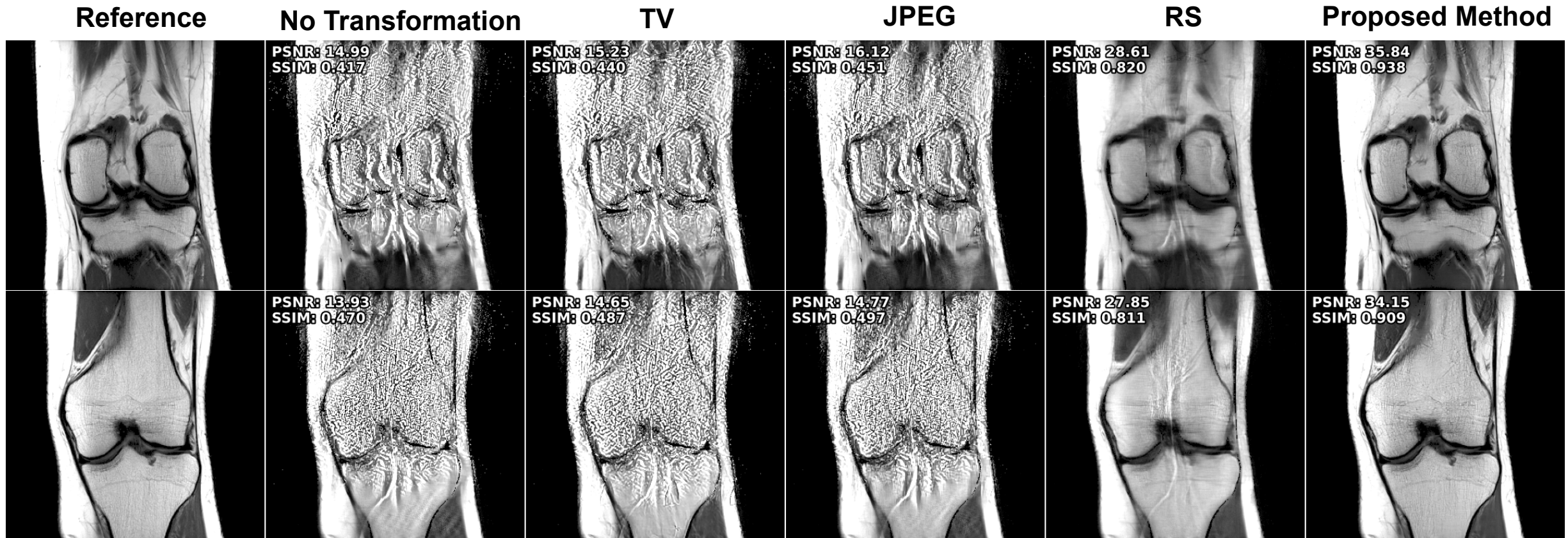
$(\mathbf{E}_{\Omega}^H)^{\dagger} (\mathbf{z}_{\Omega}^{\text{pert}} + \mathbf{r}')$: minimum ℓ_2 k-space solution of zerofilled image¹



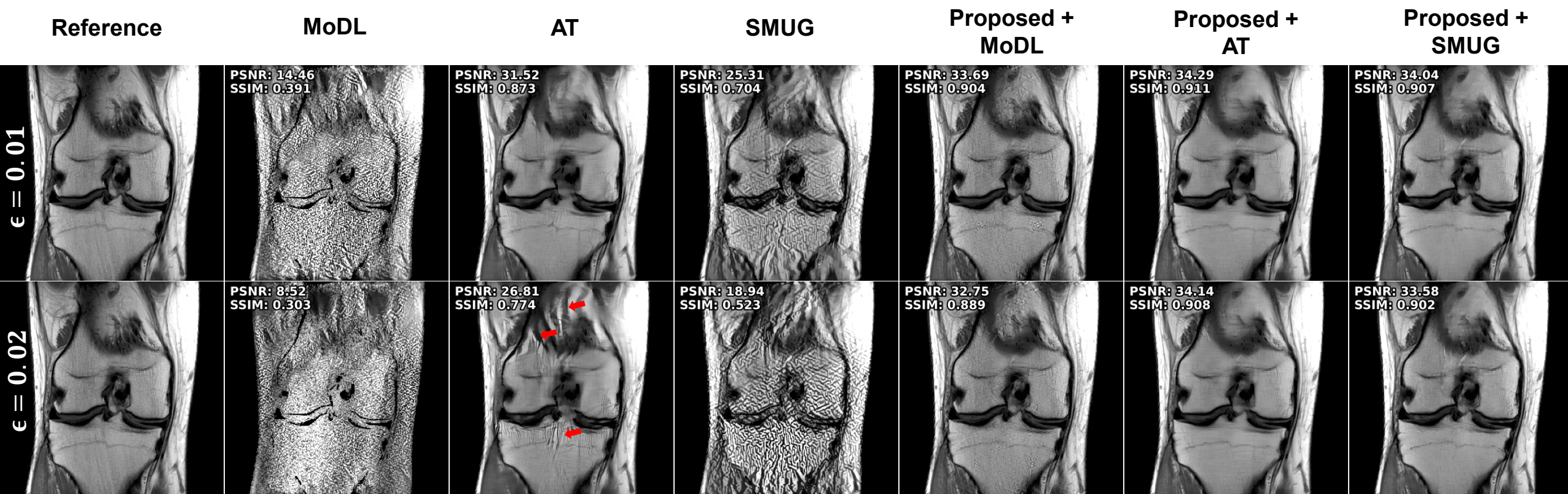
Ours vs. Training-Based Defenses








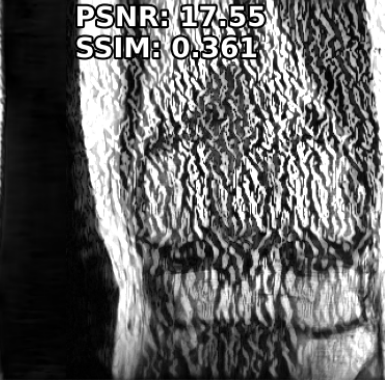




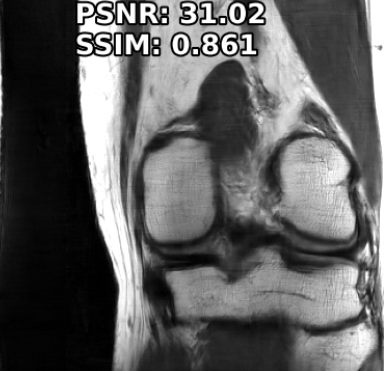

Ours vs. Training-Free Defenses



Results on Different Attack Strengths



Ours with Various Unrolled Networks

	XPDNet	RIM	E2E-VarNet	Recurrent VarNet
No Attack	 <p>PSNR: 29.66 SSIM: 0.809</p>	 <p>PSNR: 37.47 SSIM: 0.940</p>	 <p>PSNR: 32.41 SSIM: 0.890</p>	 <p>PSNR: 32.03 SSIM: 0.885</p>
PGD Attack	 <p>PSNR: 24.94 SSIM: 0.663</p>	 <p>PSNR: 17.55 SSIM: 0.361</p>	 <p>PSNR: 24.98 SSIM: 0.633</p>	 <p>PSNR: 27.41 SSIM: 0.698</p>
Proposed	 <p>PSNR: 28.30 SSIM: 0.780</p>	 <p>PSNR: 33.04 SSIM: 0.890</p>	 <p>PSNR: 31.02 SSIM: 0.861</p>	 <p>PSNR: 30.23 SSIM: 0.842</p>

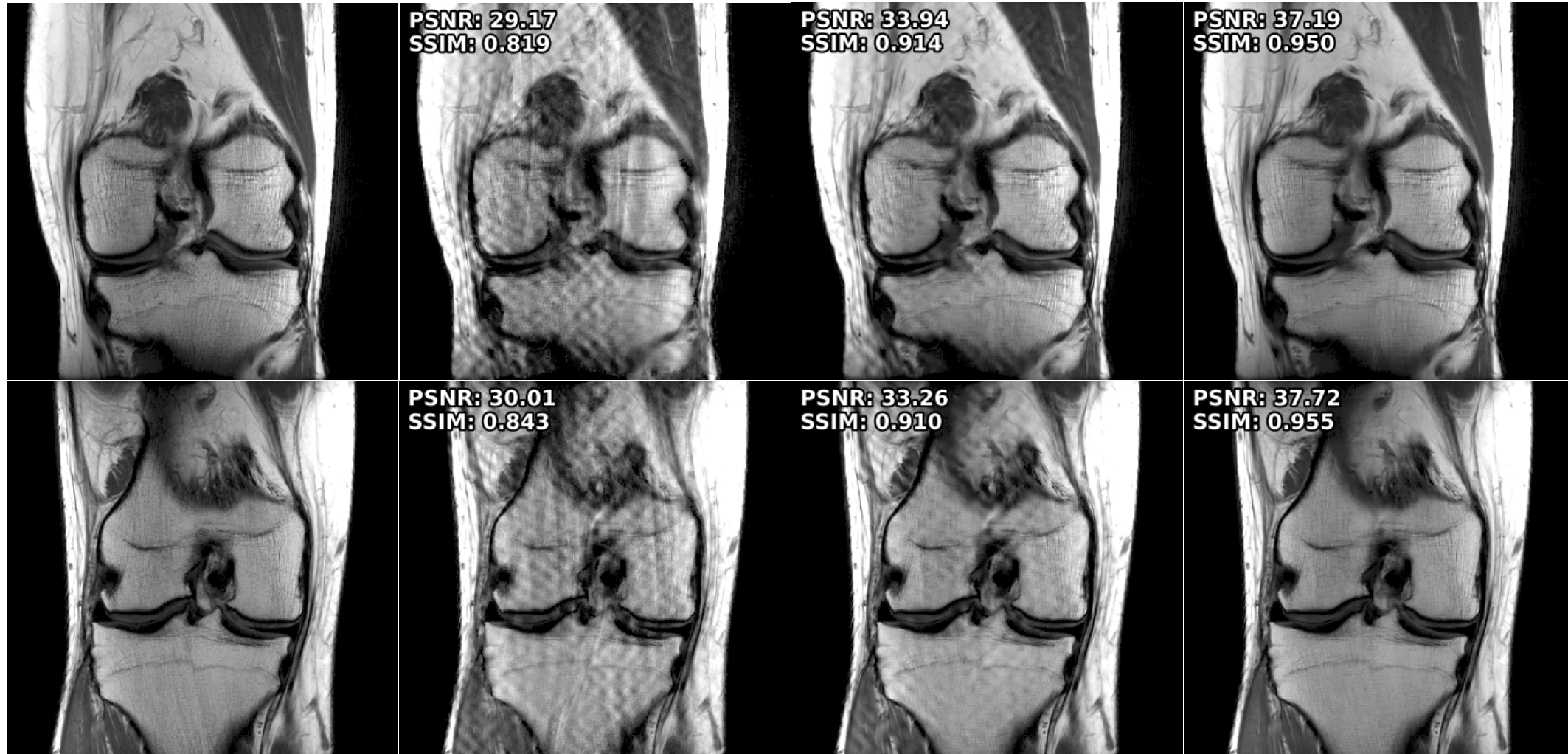
Practical Use Case— Herringbone Artifacts

Reference

CG-SENSE

MoDL

Ours + MoDL



- In this paper, we also show our method works for many cases:
 - Adaptive attacks
 - Supervised attacks
 - ℓ_2 attacks in k-space
 - Non-uniform undersampling patterns
 - Blind setup with no predefined ϵ
 - Non-optimal reconstruction conditions
 - Other inverse problems (i.e., inpainting)

- **In this work, we proposed a novel mitigation strategy for computational MRI,**
 - Leverages cyclic measurement consistency
 - Optimizes the input within a small neighborhood of the attack
 - Requires no training or modification of network
- **Results demonstrate,**
 - Our method is robust across different,
 - Datasets
 - Inverse problems
 - Unrolled networks
 - Attack strengths, Adaptive attacks
 - Our method can,
 - Perform in a blind setup
 - Be combined with existing training-based methods

Thank you!



IMAGINE Lab,
University of Minnesota



Questions?



Code & Paper



Supported by,

*NIH R01EB032830
NIH P41EB027061*