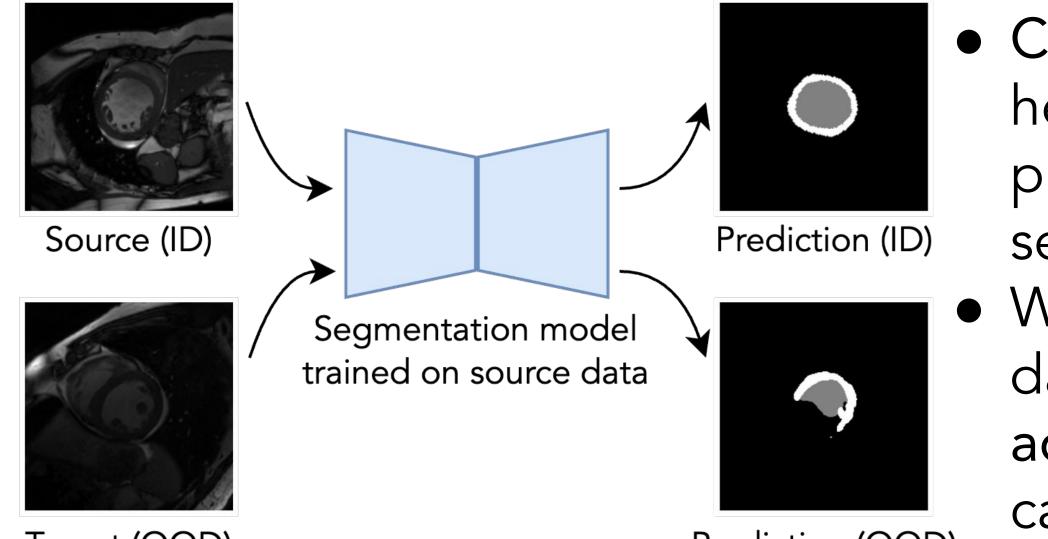# Progressive Test Time Energy Adaptation for Medical Image Segmentation

Xiaoran Zhang[1], Byung-Woo Hong[2], Hyoungseob Park[1], Daniel H. Pak[1], Anne-Marie Rickmann[1], Lawrence H. Staib[1], James S. Duncan[1*], Alex Wong[1*]

1. Yale University 2. Chung-Ang University *Joint supervision

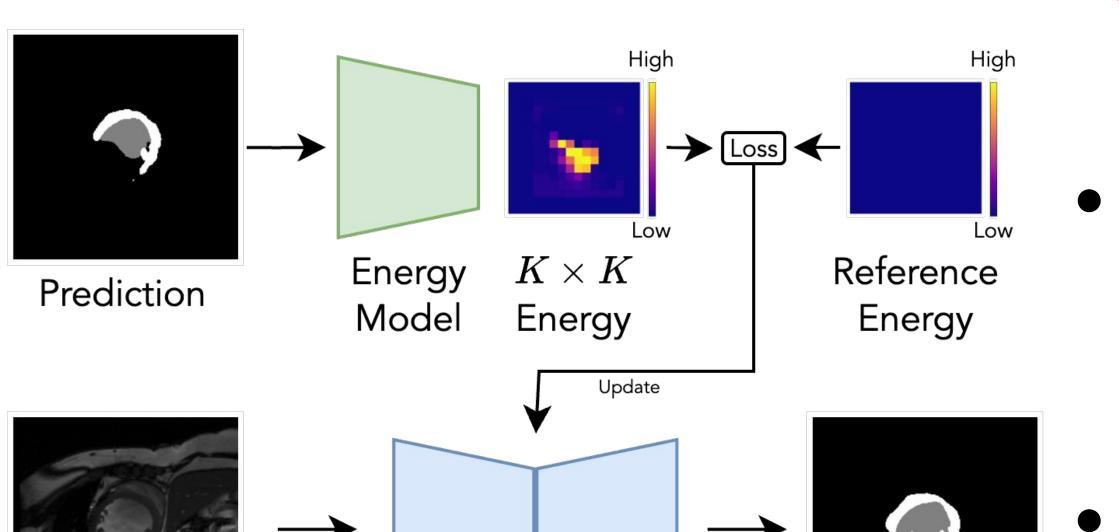**ICCV** OCT 19-23, 2025 HONOLULU HAWAII

## Background



- Covariate shifts caused by nuisances such as heteroscedastic noise and inconsistent imaging protocols limit the fidelity of medical image segmentation models.
- Without assuming access to a pre-collected target dataset, which is often impractical, **test-time adaptation (TTA)** offers a practical solution to calibrate models on-the-fly during inference.

## Problem formulation

Assuming a segmentation model is solely trained on source dataset, our goal is to adapt the model to target data without access to the entire target dataset.

## Key components

### Region-based Shape Energy Model



**Why energy? Quantifies distribution misalignment at test-time.**

- Energy-based models naturally capture distribution changes by reflecting sample likelihood, make them suited for TTA.
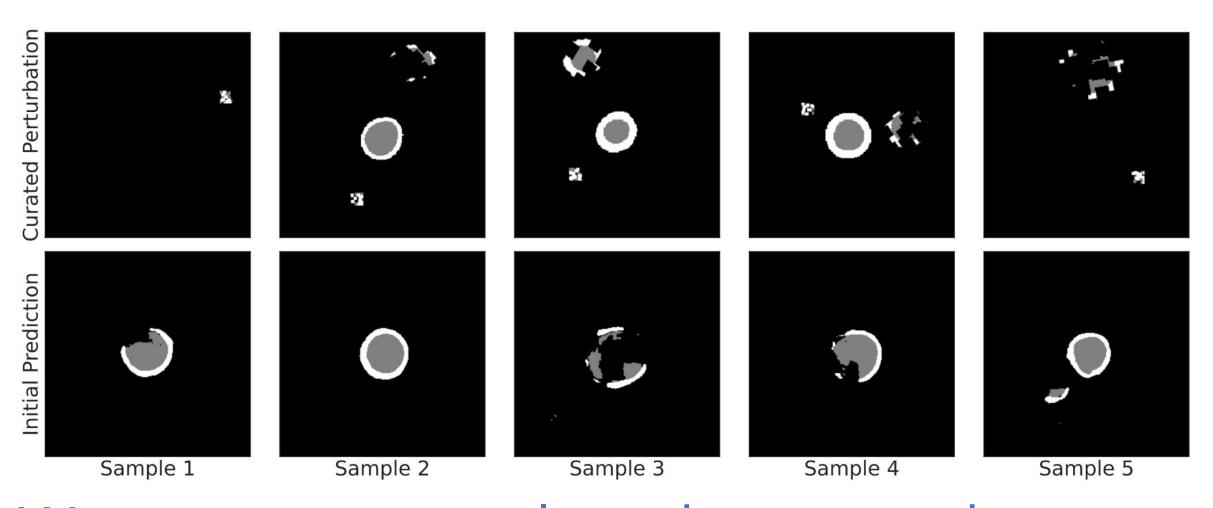
$$p_\phi(x) = \frac{\exp(-E_\phi(x))}{Z(\phi)}$$

- We propose a shape energy model trained on source data, which assigns an energy score at the region level:
  - low energy -> ID (accurate) shapes
  - high energy -> OOD (erroneous) predictions

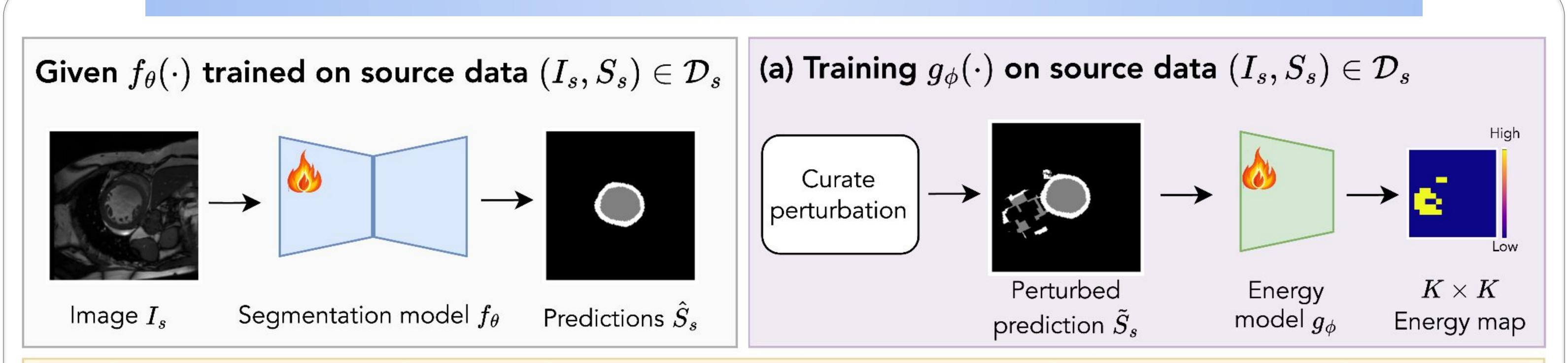### Curate negative examples for energy model training

**Impossible to have real negative examples? Curate them instead.**

- Our formulation assumes two collections of examples, one following the desired distribution (of shapes) and the other out-of-distribution (OOD).
- In addition, the input distribution to the energy model is constrained by the predictions afforded by the segmentation model.



We propose to explore data space by probing the segmentation model with inputs optimized to simulate OOD examples.

## Methods

**Given $f_\theta(\cdot)$ trained on source data $(I_s, S_s) \in \mathcal{D}_s$**



Image $I_s$ → Segmentation model $f_\theta$ → Predictions $\hat{S}_s$

**(a) Training $g_\phi(\cdot)$ on source data $(I_s, S_s) \in \mathcal{D}_s$**



Curate perturbation → Perturbed prediction $\tilde{S}_s$ → Energy model $g_\phi$ → $K \times K$ Energy map

**(b) Progressive test time adaptation on target data $I_t \in \mathcal{D}_t$**



Image $I_t$ — Iterations — Segmentation model $f_\theta$ — Adapted prediction $\hat{S}_t$ — Energy model $g_\phi$ — $K \times K$ Energy map — Reference energy $0_{K \times K}$

### Perturbation curation

We generate negative (implausible) examples by applying FGSM adversarial noise and spatial affine transformations to the input images.

- Apply FGSM adversarial noise: $\epsilon = \delta \operatorname{sign}\left(\nabla_{I_s}\mathcal{L}(f_\theta(I_s), S_s)\right)$
- Apply random affine transforms
- Generate perturbed segmentation: $\tilde{S}_s = f_\theta(I_s + \epsilon)$

### Shape energy model training

A region-based model learns patchwise energy values, assigning high energy to implausible regions and low energy to anatomically valid ones.

Loss: $\mathcal{L}_\phi = \frac{1}{N_p}\sum_{i=1}^{N_p}\left(-y_s^i \log \sigma(-g_\phi(s_s^i)) - (1-y_s^i)\log(1-\sigma(-g_\phi(s_s^i)))\right)$

### Label curation

For each perturbed segmentation, we compare it with ground truth and assign categorical energy labels to each region, where regions dissimilar to the ground truth are labeled as high-energy. $y_s = 1 - \mathbf{1}(d(\tilde{s}_s, s_s) < \tau)$

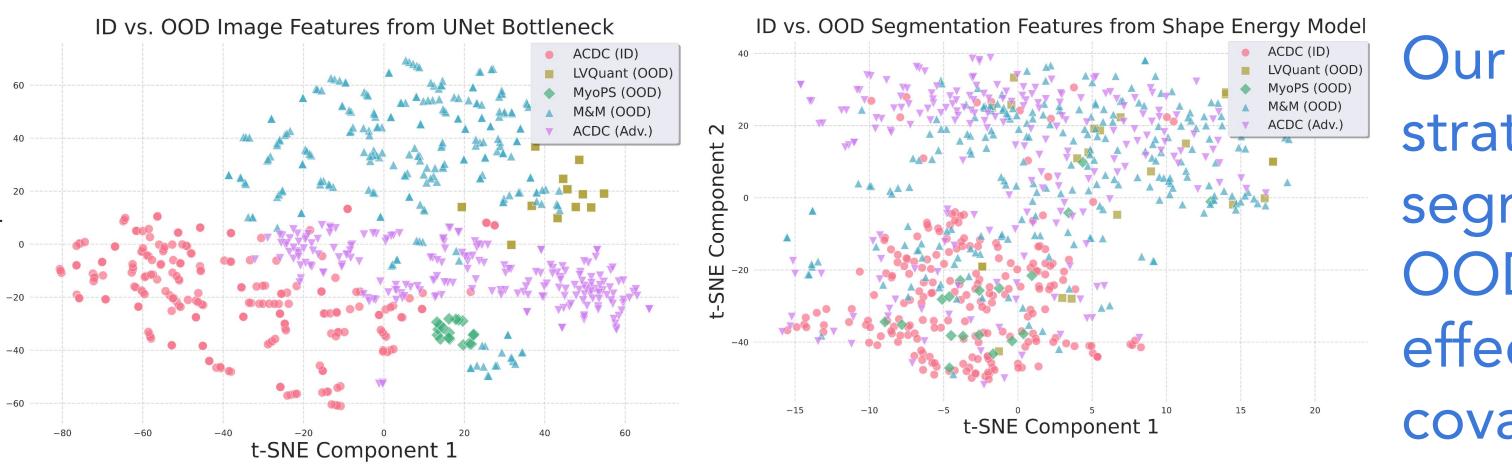### Progressive test-time adaptation

At inference, the segmentation model is iteratively updated to minimize the predicted energy, aligning outputs with plausible anatomical shapes.

Update rule: $\theta^* = \arg\min_\theta -\sum_{i=1}^{B_t}\log(1-\sigma(-g_\phi(\hat{s}_t^i)))$

Datasets: (1) Cardiac (2D MRI): ACDC, LVQuant, MyoPS, M&M (2) Spinal cord (2D MRI): GMSC (sites 1-4) (3) Lung (2D X-ray): CHN, MCU, JSRT. Metrics: (1) Dice coefficient score (DSC, %) (2) average surface distance (ASD, %)

## Results
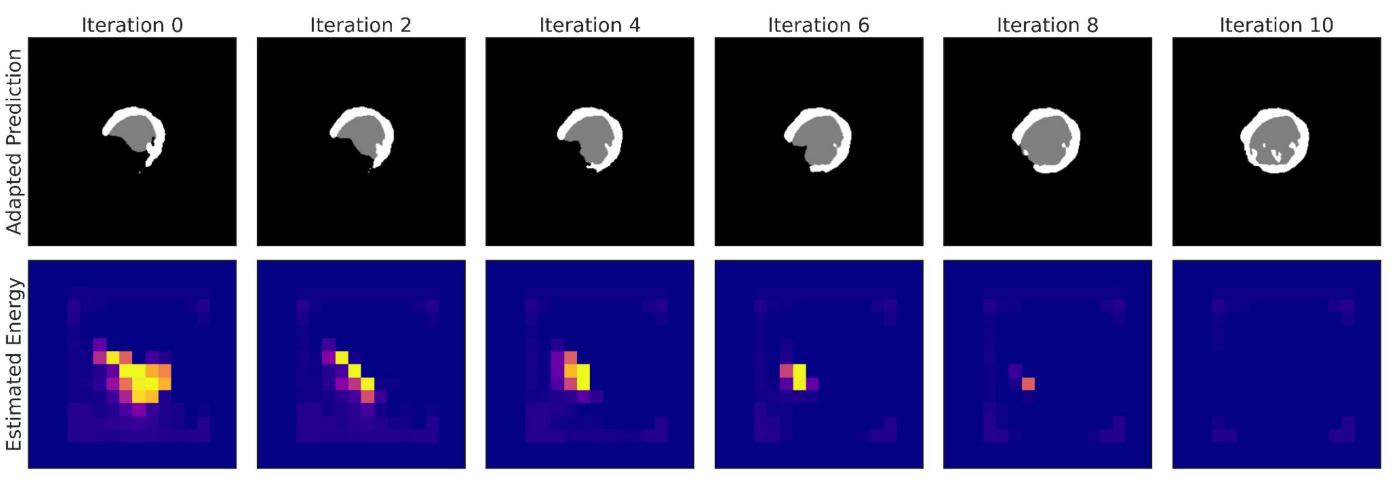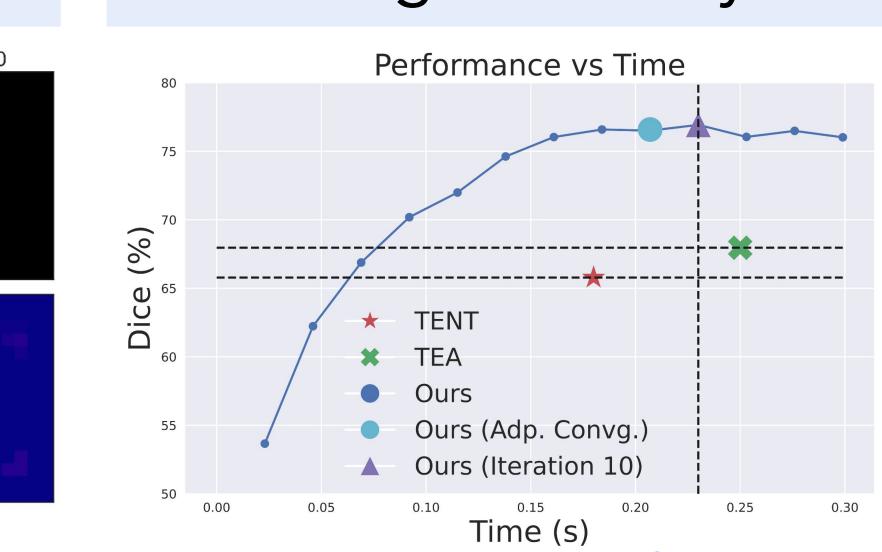
### T-SNE analysis of curated perturbation



Our adversarial perturbation strategy produces images and segmentations that align with OOD cases, validating its effectiveness in modeling real covariate shifts.

### Progressive update visualization
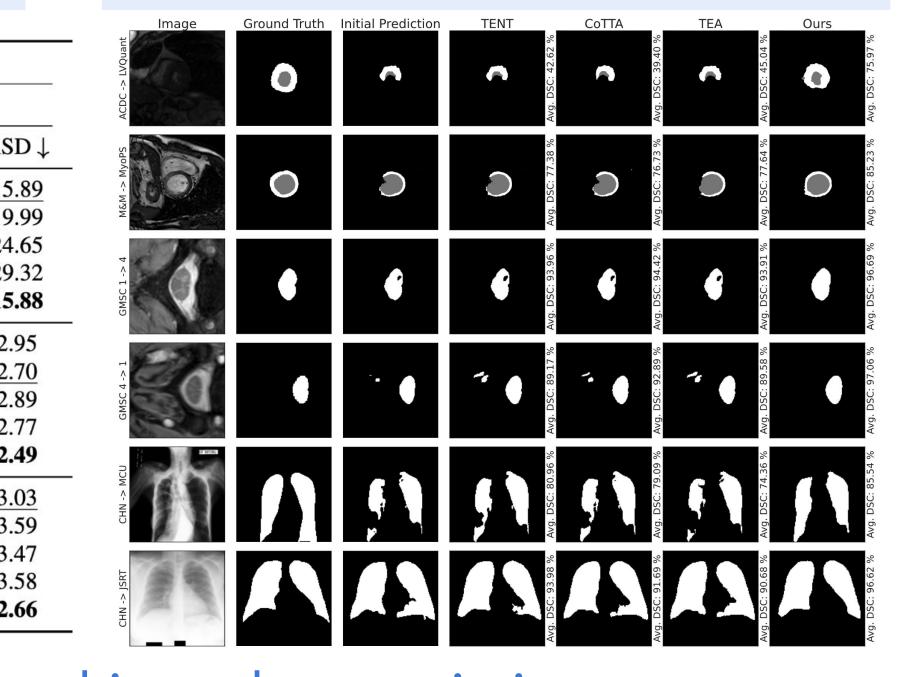


### Convergence analysis



Our method progressively refines segmentation quality over iterations (left), while achieving better convergence under the same time budget (right).

### Quantitative evaluation



Our proposed approach can be plugged-and-played into three existing architectures and we consistently outperform baselines in eight datasets.

### Adaptation visualization



### High-energy corresponds to test-time segmentation errors

| Method | ACDC↦LVQuant | ACDC↦MyoPS | ACDC↦M&M |
|---|---|---|---|
| UNet | 93.64 | 96.55 | 94.53 |
| MedNeXt | 92.17 | 95.83 | 93.66 |
| SwinUNETR | 92.01 | 96.42 | 93.97 |

Our shape energy model achieves over 92% accuracy across different segmentation models, confirming its effectiveness in identifying errors at test-time.

### Hyperparameter sensitivity

| | | 1↦2 | 1↦3 | 1↦4 | 4↦1 | 4↦2 | 4↦3 | Avg. |
|---|---|---|---|---|---|---|---|---|
| Patch Size | 4 × 4 | 69.1 | 73.2 | 93.4 | 89.4 | 42.5 | 85.3 | 75.5 |
| | 9 × 9† | **73.6** | 77.7 | 95.3 | **95.1** | 56.2 | 87.2 | **80.9** |
| | 18 × 18 | 73.0 | **77.9** | 94.5 | 94.7 | **57.2** | **87.7** | 80.8 |
| | 36 × 36 | 69.4 | 75.5 | 93.1 | 88.4 | 45.9 | 87.4 | 76.6 |
| Threshold | $\tau = 25$ | 70.5 | 73.1 | 91.3 | 94.2 | 54.1 | 86.3 | 78.3 |
| | $\tau = 50$† | **73.6** | 77.7 | 95.3 | **95.1** | **56.2** | 87.2 | **80.9** |
| | $\tau = 75$ | 73.0 | **79.1** | 95.1 | 94.7 | 52.6 | **87.3** | 80.3 |
| | $\tau = 100$ | 71.6 | 75.7 | 95.1 | 94.6 | 52.5 | 86.7 | 79.4 |
| Pert. Mag. | $\delta = 0.1$ | 70.8 | 76.2 | 94.2 | 95.1 | 54.5 | 87.0 | 79.6 |
| | $\delta = 0.05$† | **73.6** | 77.7 | **95.3** | **95.1** | **56.2** | **87.2** | **80.9** |
| | $\delta = 0.01$ | 73.3 | 73.1 | 93.6 | 94.8 | 53.5 | 86.4 | 79.1 |

† is proposed.