

Tackling **Structural** **Hallucination** in Image Translation with Local Diffusion

Seunghoi Kim*, Chen Jin*, Tom Dieth, Matteo Figini, Henry F. J. Tregidgo, Asher Mullokandov, Philip Teare, and Daniel C. Alexander

Introduction

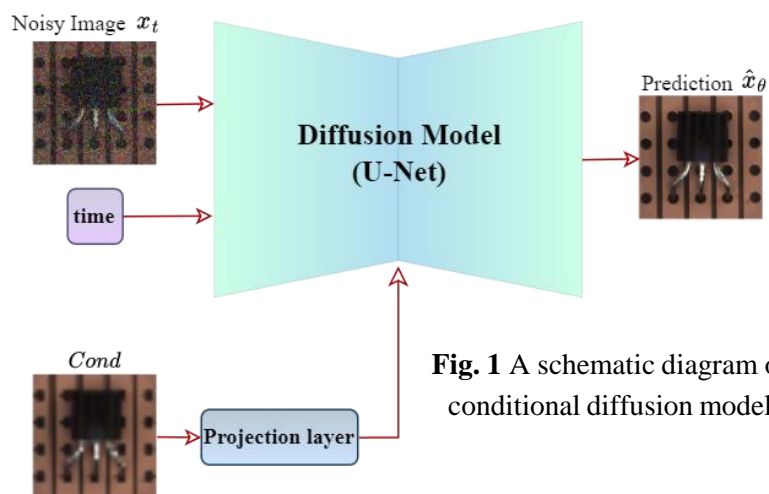


Fig. 1 A schematic diagram of conditional diffusion model

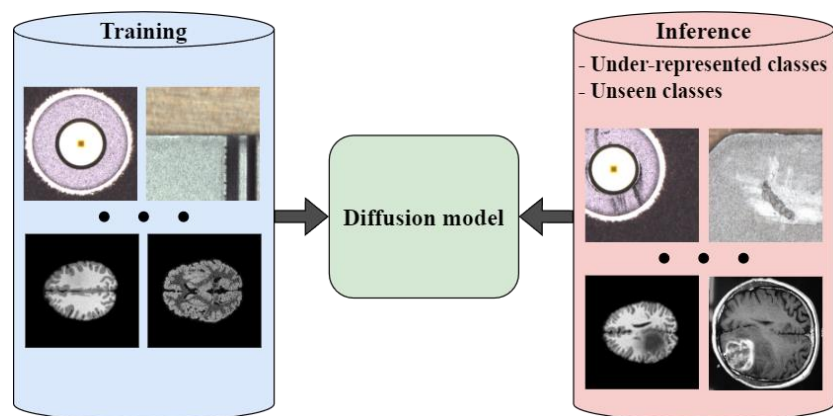


Fig. 2 Illustration of data-shift problem

- Conditional diffusion models have attained state-of-the-art performances in various image translation tasks
- Generative model f_θ learns joint distribution of train data and its condition
- Many applications with significant under-represented classes (e.g. rare diseases, defects) exist
- Performs well on in-distribution (IND) data, but what if the condition contains out-of-distribution (OOD) region?

$Error(f_\theta(c_{ood\&ind})) \gg Error(f_\theta(c_{ind}))$ or at worst hallucination!

Problem: Hallucination

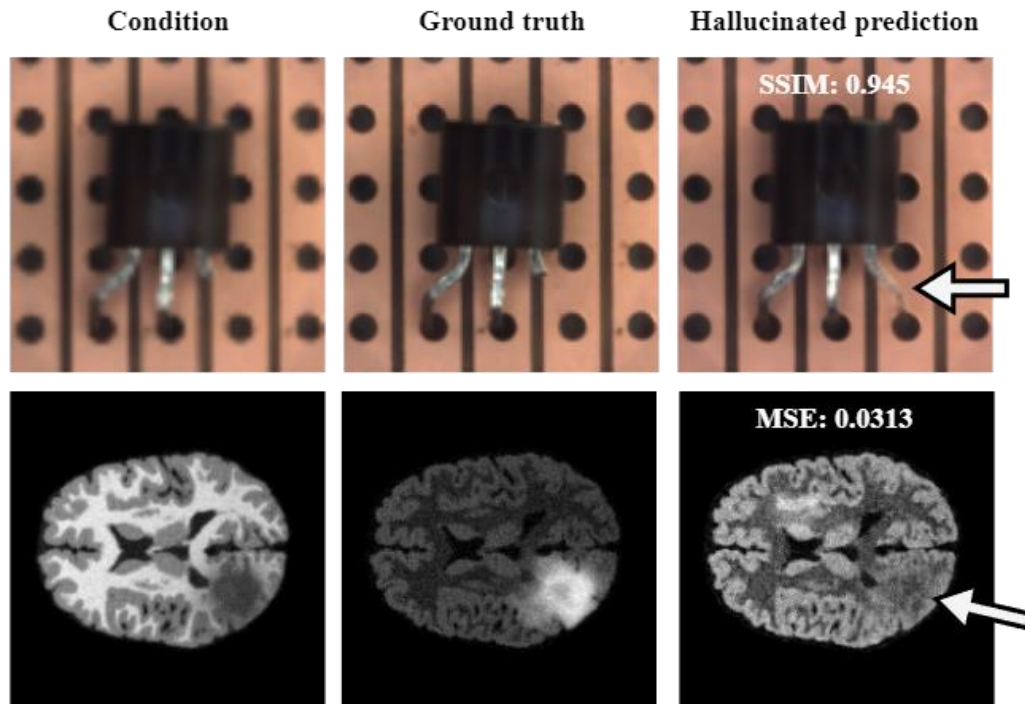


Fig. 3 Visual examples of structural hallucination

- What is structural hallucination and why should we care?
 - Realistic-looking but inaccurately reconstructed features, leading to discrepancies with the actual structure
 - Misinterpretation => patient misdiagnosis, machine failure, increase in time and cost
 - Often **insensitive** using **standard image quality metrics** (e.g. MSE, SSIM)

Solution?

- Simple way to fix => fine-tuning, but **expensive**

Hypothesis & Verification

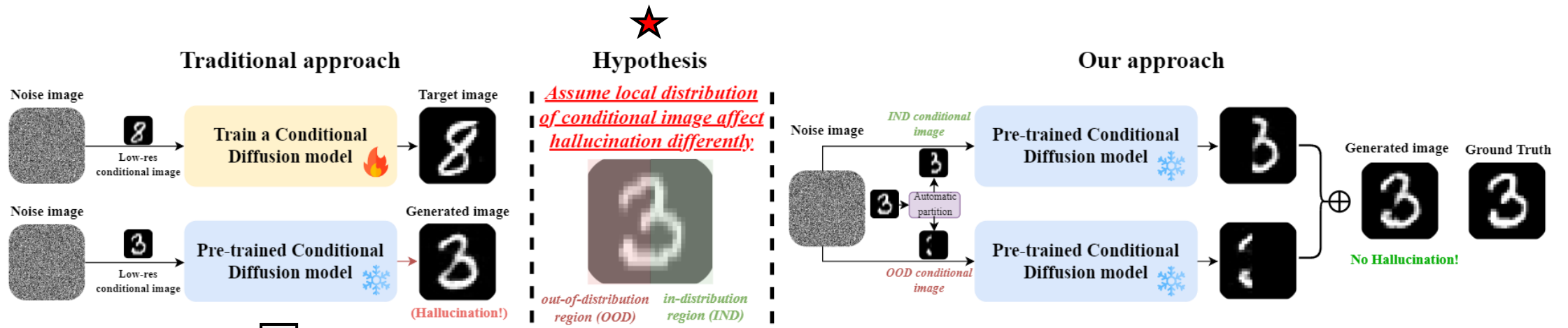


Fig. 4 Traditional conditional diffusion process vs. our OOD/IND Local Diffusion

Condition	Ground truth	Original pred.
		Cls: 8

Traditional approach leads to hallucination!

Can OOD-based Local Image Generation Help to Reduce Hallucination?

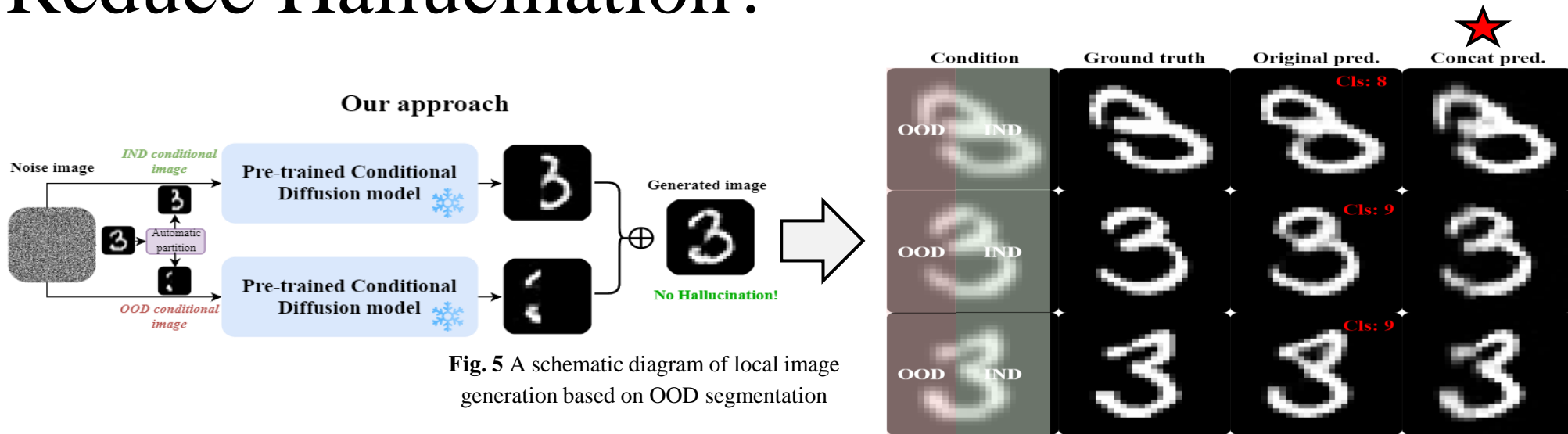
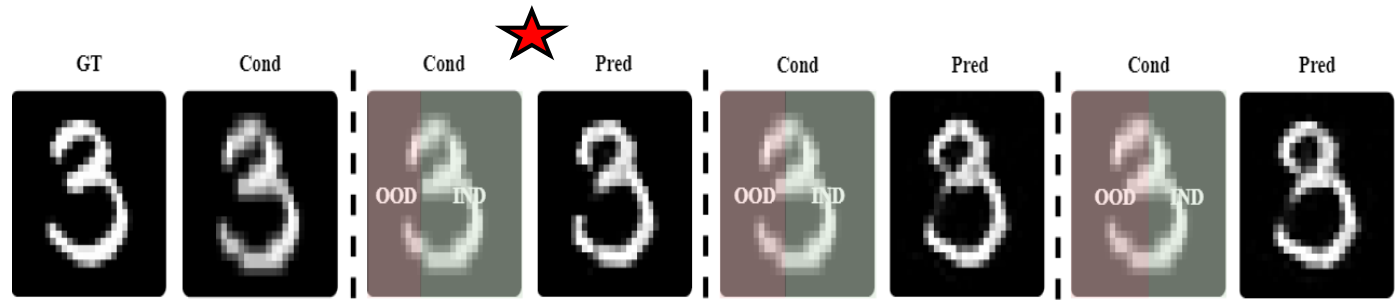


Fig. 5 A schematic diagram of local image generation based on OOD segmentation



Precise segmentation of OOD is crucial for hallucination mitigation!

Fig. 6 Visual illustration of the impact of shifting OOD boundary

Identifying Hallucination Hotspots in Diffusion Models

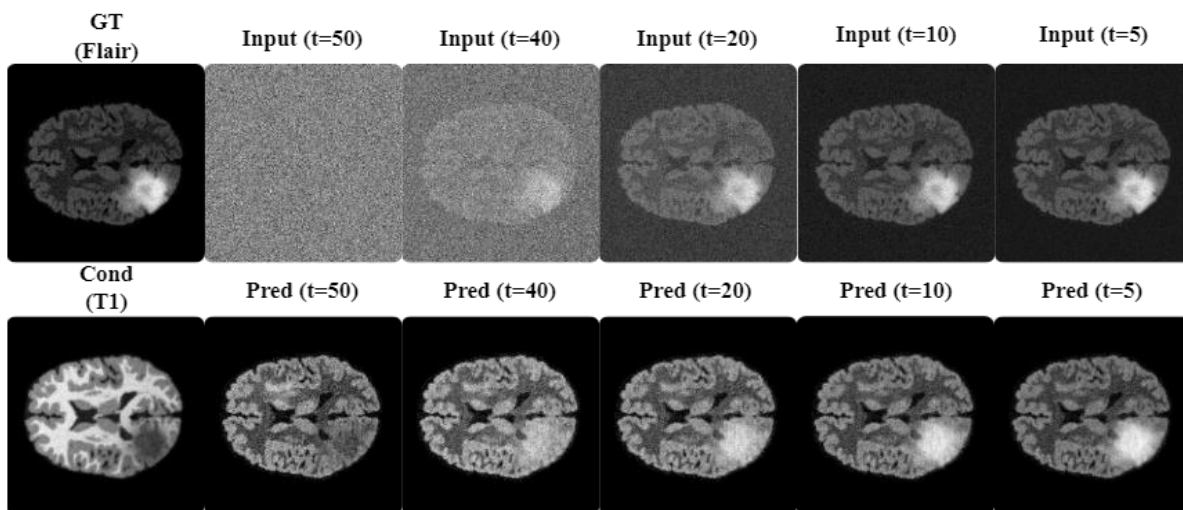


Fig. 7 Qualitative comparisons of predictions starting from different intermediate time points. We sample noisy GT (Flair) and perform a reverse process from it by conditioning the corresponding T1 image.

Less hallucination! 

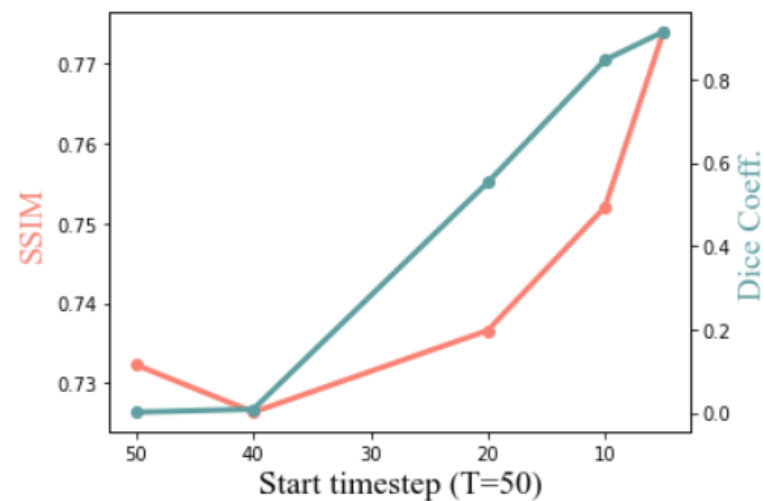
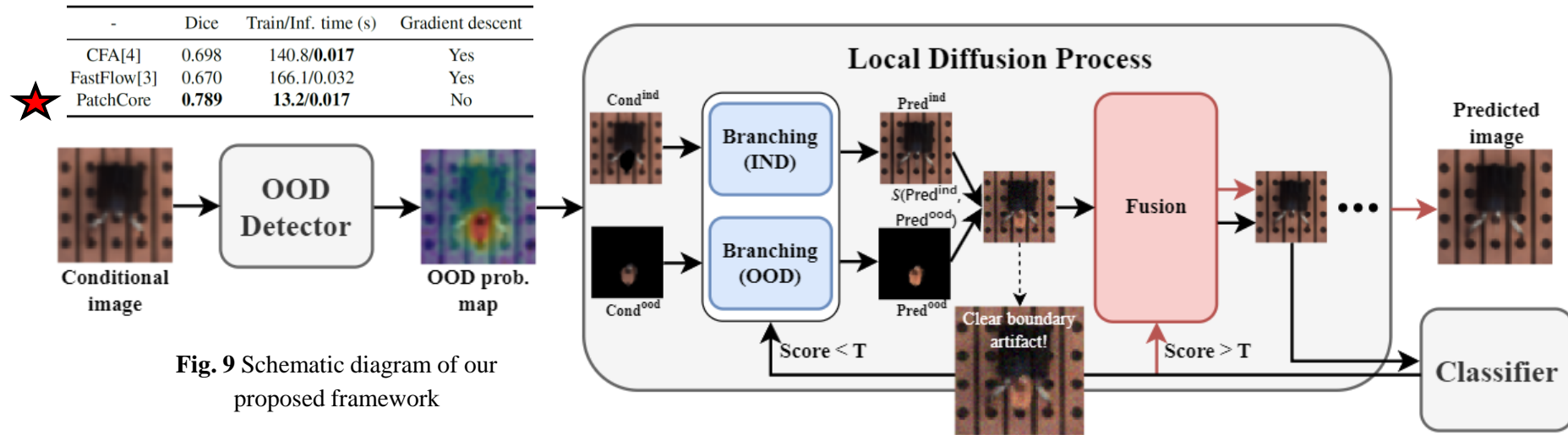


Fig. 8 Dice score of tumor segmentation and SSIM of predicted images

Methods: Local Diffusion



- **OOD estimation:** One-class classification anomaly detector (PatchCore [CVPR'22])
- **Branching:** Separate local image generation based on OOD probability map
- **Fusion:** Fuse the OOD/IND predictions for more cohesive image generation
- **Classifier:** Checks if the prediction at intermediate time step contains hallucination

Main Results (Quantitative)

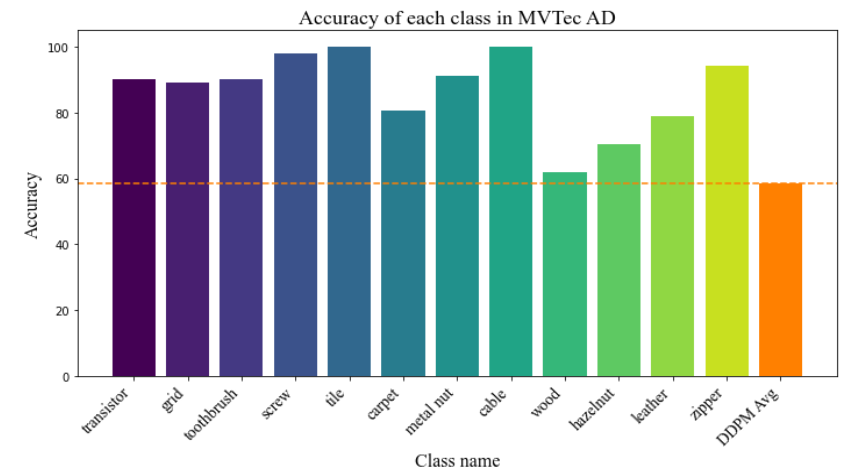
	MNIST		BraTS		MVTec AD	
	<i>PSNR</i> (\uparrow)	<i>SSIM</i> (\uparrow)	<i>PSNR</i> (\uparrow)	<i>SSIM</i> (\uparrow)	<i>PSNR</i> (\uparrow)	<i>SSIM</i> (\uparrow)
DDPM [11]	20.8 \pm 1.90	0.897 \pm 0.04	19.5 \pm 1.75	0.709 \pm 0.04	26.8 \pm 3.10	0.839 \pm 0.10
DDIM [32] (0.1T)	20.8 \pm 1.88	0.895 \pm 0.04	19.8 \pm 1.53	0.715 \pm 0.05	27.3 \pm 3.00	0.844 \pm 0.10
DDIM [32] (0.5T)	20.6 \pm 1.86	0.899 \pm 0.04	19.7 \pm 1.61	0.703 \pm 0.04	27.1 \pm 2.98	0.840 \pm 0.10
DDPM + DSI [38]	20.8 \pm 1.90	0.898 \pm 0.04	18.7 \pm 1.74	0.695 \pm 0.04	26.6 \pm 3.40	0.815 \pm 0.10
DDIM + Ours (0.1T)	20.9\pm1.87	0.897 \pm 0.04	20.7 \pm 1.52	0.720\pm0.05	27.5\pm3.02	0.847\pm0.09
DDIM + Ours (0.5T)	20.8 \pm 1.82	0.902\pm0.03	20.9 \pm 1.61	0.711 \pm 0.04	27.4 \pm 2.97	0.844 \pm 0.10
DDPM + Ours	20.9\pm1.85	0.900 \pm 0.04	21.2\pm1.74	0.720\pm0.03	27.0 \pm 3.05	0.843 \pm 0.10
p-value	0.014	0.016	≈ 0	0.02	0.001	0.004

Tab. 1 Quantitative comparisons of overall image quality across various datasets, where an upward arrow signifies that a higher value is better. T represents the total number of time steps for sampling.

	MNIST	BraTS	MVTec AD
	<i>Accuracy</i> (%) (\uparrow)	<i>Dice Coefficient</i> (\uparrow)	<i>Accuracy</i> (%) (\uparrow)
DDPM [11]	95.7 \pm 0.61	0.194 \pm 0.10	58.4 \pm 19.9
DDIM [32] (0.1T)	95.8 \pm 0.61	0.256 \pm 0.11	53.9 \pm 19.5
DDIM [32] (0.5T)	95.9 \pm 0.60	0.246 \pm 0.13	59.1 \pm 19.4
DDPM + DSI [38]	96.0 \pm 0.50	0.100 \pm 0.06	60.1 \pm 18.6
DDIM + Ours (0.1T)	96.1 \pm 0.30	0.447 \pm 0.10	66.8 \pm 17.8
DDIM + Ours (0.5T)	97.1 \pm 0.28	0.537 \pm 0.06	83.2 \pm 14.3
DDPM + Ours	98.0\pm0.26	0.590\pm0.09	85.0\pm13.2
Avg. gain	+1.5%	+0.293	+21.2%

Tab. 2 Quantitative results on downstream tasks to measure hallucination

Downstream task performance is more important in our task to evaluate the level of hallucination!



+26%!

Main Results (Qualitative)

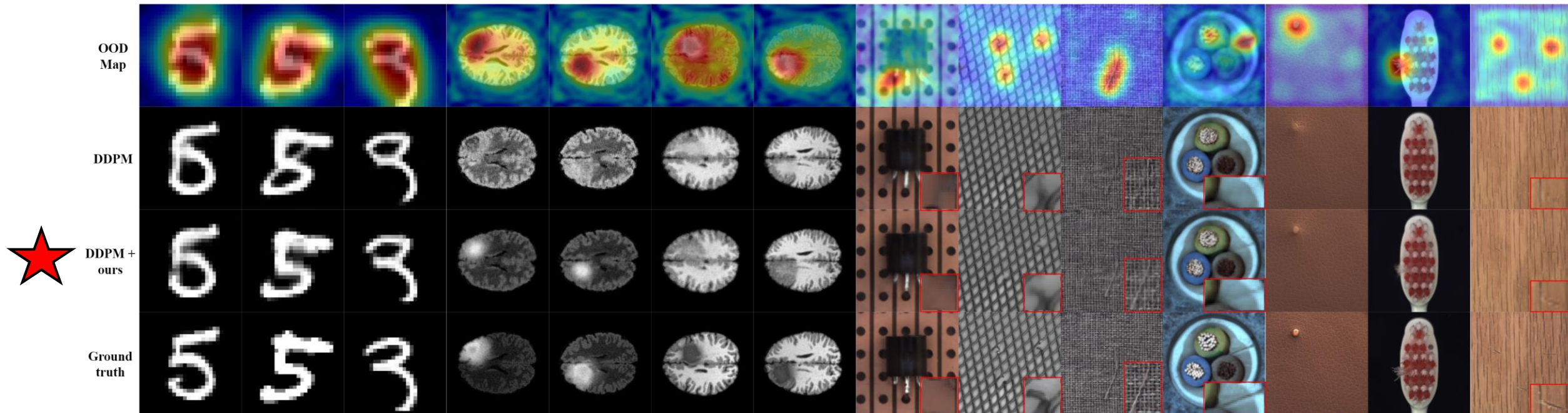


Fig. 10 Qualitative comparison on MNIST, BraTS and MVTEC (From top: predicted OOD map, DDPM, DDPM with ours and ground truth).

Further Analysis

Does Local Diffusion work on various OOD?

-	Single		Multiple	
	DDPM	Ours	DDPM	Ours
<i>PSNR</i> (\uparrow)	19.2	22.7	19.1	21.4
<i>Dice coeff.</i> (\uparrow)	0.032	0.667	0.055	0.540

(a) Model's performance on single and multiple OOD in a single image

-	Small		Large	
	DDPM	Ours	DDPM	Ours
<i>PSNR</i> (\uparrow)	20.7	23.4	18.0	21.5
<i>Dice coeff.</i> (\uparrow)	0.050	0.685	0.030	0.580

(b) Model's performance on small (< 1.5%) and large OOD (> 3%) regions

Tab. 3 Quantitative comparisons on various types of OOD

Does Local Diffusion have negative impact on IND region?

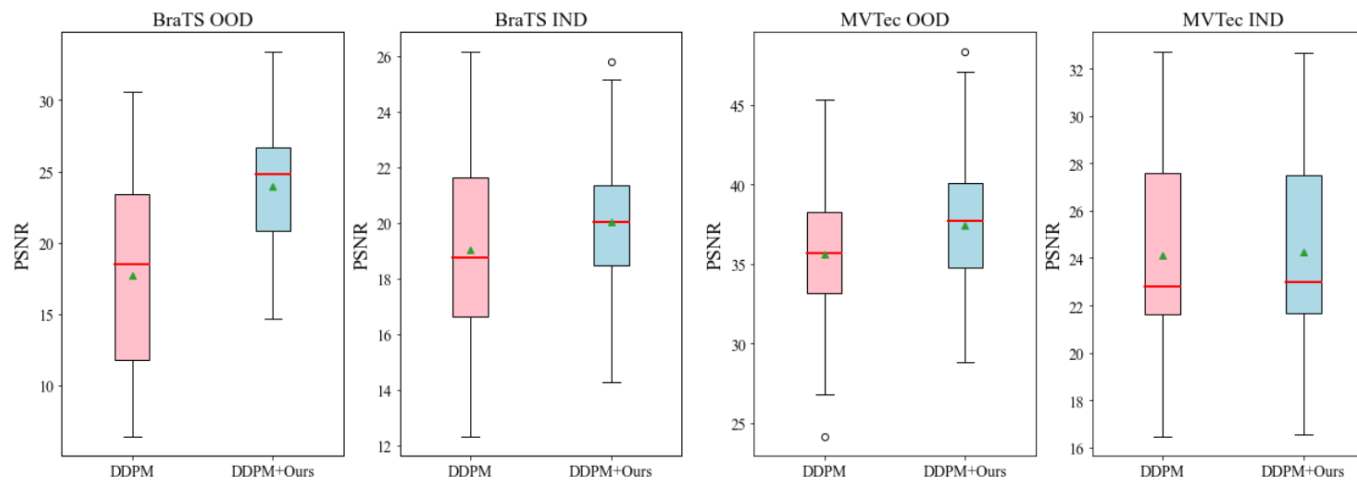


Fig. 11 Comparative analysis of performance across individual OOD/IND regions, The red lines and green dots represent the median and mean of each box, respectively

Thank you!

Arxiv



Github

