



INTRODUCTION

Deep models can be overconfident, so reliable prediction requires understanding not only how uncertain a model is, but also why it is uncertain [1]. Uncertainty in deep learning mainly comes from two sources: aleatoric uncertainty, caused by noise or ambiguity in the data, and epistemic uncertainty, caused by limited model knowledge or insufficient training coverage [2]. Distinguishing them is important for trustworthy decision-making [3]. Existing methods often estimate only one uncertainty type, or they require modifying and retraining the base model. This increases computational cost and limits compatibility with pretrained systems.

We propose **CUPID**, a lightweight plug-in module that can be inserted into a pretrained network to jointly estimate aleatoric and epistemic uncertainty without changing or retraining the base model. It also enables layer-wise analysis of how uncertainty evolves inside the network.

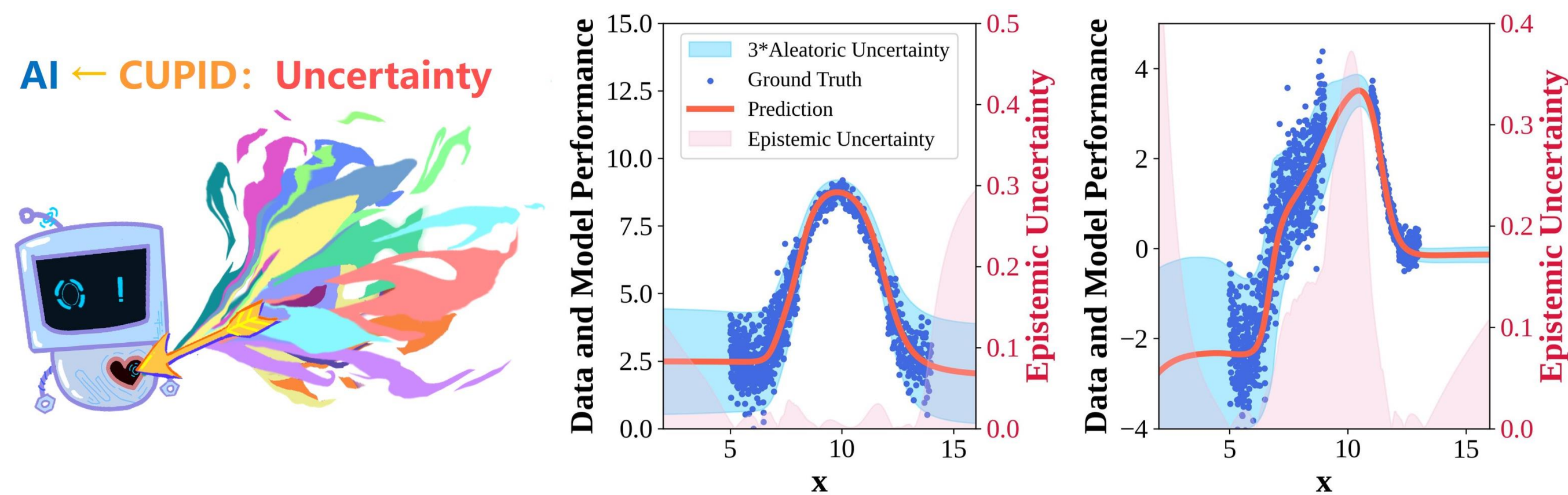


Figure 1: CUPID uncertainty estimation on a 1D regression toy problem. CUPID is inserted into an MLP-based predictive model. CUPID captures both aleatoric (blue) and epistemic (red) uncertainty.

METHODS

CUPID is a plug-in module inserted at any intermediate layer of a pretrained model to jointly estimate aleatoric uncertainty from a variance branch and epistemic uncertainty from prediction changes caused by feature perturbation.

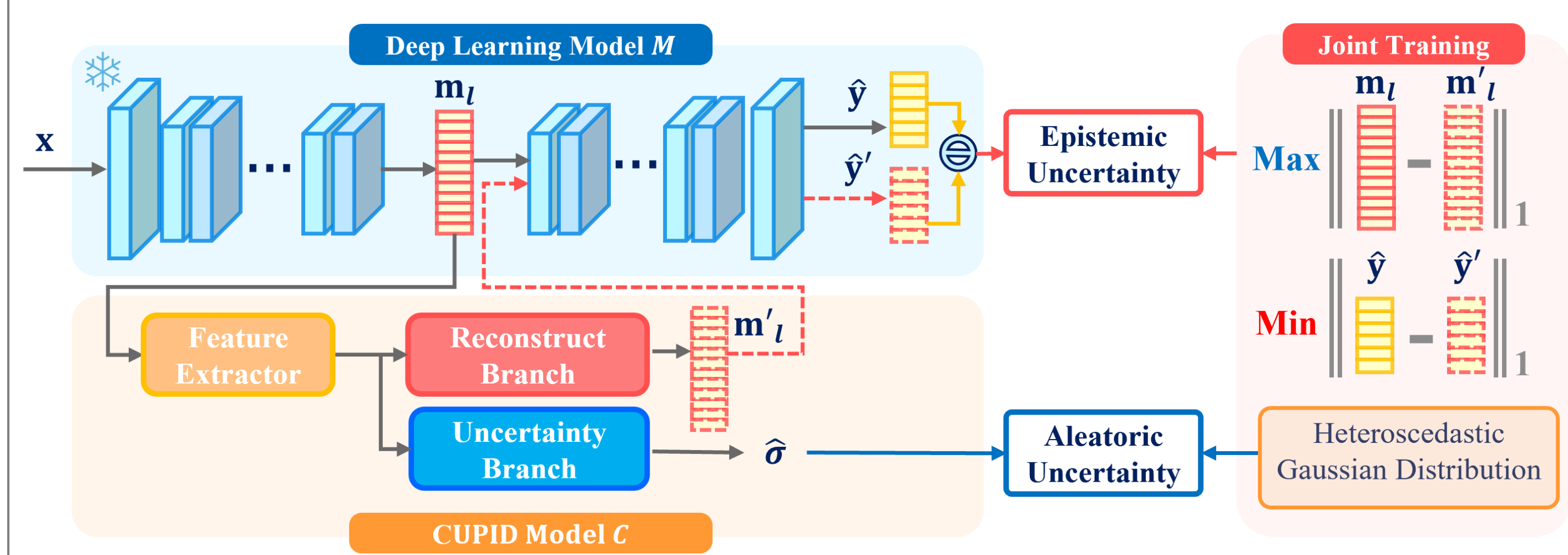


Figure 2: The CUPID pipeline. Aleatoric uncertainty is estimated using a dedicated Uncertainty Branch, while epistemic uncertainty is captured by measuring the variance between the original model output \hat{y} and the perturbed output \hat{y}' .

- **Plug-in formulation.** Let $M(x) = F_l(B_l(x))$, where B_l extracts the intermediate feature and F_l maps it to the final prediction. CUPID model takes $m_l = B_l(x)$ and produces $(m_l', \hat{\sigma})$.
- **Aleatoric uncertainty branch.** The uncertainty branch models input-dependent noise with a heteroscedastic distribution. The predicted variance is used as the aleatoric uncertainty score. For training, CUPID predicts the log-variance for numerical stability.
- **Epistemic uncertainty branch.** The reconstruction branch learns a feature perturbation that changes the internal representation while preserving the final prediction. Epistemic uncertainty is measured by the discrepancy between the original output and the perturbed output.
- **Joint objective.** Both branches are optimized jointly in a single model. This lets CUPID estimate both uncertainty types without modifying or retraining the base predictor.

$$(m_{l,n}', \hat{\sigma}_n) = C(m_{l,n}; \omega)$$

$$U_{\text{alea}}(x_n) = \hat{\sigma}_n^2$$

$$U_{\text{epis}}(x) = \|\hat{y}_n - \hat{y}'_n\|_1$$

$$L_{\text{CUPID}} = L_{\text{epis}} + \lambda_2 L_{\text{alea}}$$

REFERENCES

[1] K. Zou, Z. Chen, X. Yuan, X. Shen, M. Wang, and H. Fu, "A review of uncertainty estimation and its application in medical imaging," *Meta-Radiology*, p. 100003, 2023.

[2] J. Gawlikowski, C. R. N. Tassi, M. Ali, J. Lee, M. Humt, J. Feng, A. Kruspe, R. Triebel, P. Jung, R. Roscher, et al., "A survey of uncertainty in deep neural networks," *Artificial Intelligence Review*, vol. 56, suppl. 1, pp. 1513–1589, 2023.

[3] M. Chan, M. Molina, and C. Metzler, "Estimating epistemic and aleatoric uncertainty with a single model," in *Advances in Neural Information Processing Systems*, vol. 37, pp. 109845–109870, 2024.

RESULTS

We evaluate CUPID on three complementary settings: medical image misclassification detection, out-of-distribution (OOD) detection, and image super-resolution as a regression task. These experiments are designed to test whether CUPID can separate aleatoric and epistemic uncertainty across both classification and regression problems. We further include an ablation study on insertion location to analyze how uncertainty emerges across network depth. All experiments are repeated three times, and results are reported as mean \pm standard deviation. The best model for each metric is in bold, and the second best is underlined.

Table 1: Performance of misclassification detection (misclassified samples as positive). CUPID Aleatoric achieved the best performance on GLV2, while CUPID Epistemic performed best on HAM10000, suggesting different dominant sources of uncertainty across datasets.

Method	GLV2			HAM10000		
	AUC (\uparrow)	AURC (\downarrow)	Spearman (\uparrow)	AUC (\uparrow)	AURC (\downarrow)	Spearman (\uparrow)
CUPID Alea.	0.870 \pm 0.002	0.018 \pm 0.001	0.941 \pm 0.004	0.769 \pm 0.023	0.067 \pm 0.007	0.722 \pm 0.014
CUPID Epis.	0.769 \pm 0.015	0.034 \pm 0.002	0.701 \pm 0.051	0.855 \pm 0.006	0.047 \pm 0.001	0.907 \pm 0.001
MC Dropout	0.768 \pm 0.006	0.027 \pm 0.001	0.888 \pm 0.005	0.829 \pm 0.001	0.076 \pm 0.001	0.861 \pm 0.002
Rate-in	0.815 \pm 0.006	0.024 \pm 0.001	0.816 \pm 0.004	0.846 \pm 0.001	0.048 \pm 0.000	0.915 \pm 0.000
IGRUE	0.642 \pm 0.007	0.058 \pm 0.002	0.199 \pm 0.004	0.548 \pm 0.004	0.157 \pm 0.002	0.027 \pm 0.018
PostNet Alea.	0.671 \pm 0.006	0.182 \pm 0.004	0.641 \pm 0.011	0.793 \pm 0.007	0.142 \pm 0.003	0.764 \pm 0.006
PostNet Epis.	0.559 \pm 0.031	0.238 \pm 0.019	0.284 \pm 0.054	0.751 \pm 0.017	0.158 \pm 0.010	0.698 \pm 0.033
BNN	0.829 \pm 0.018	0.025 \pm 0.003	0.954 \pm 0.007	0.793 \pm 0.006	0.096 \pm 0.004	0.821 \pm 0.009
DEC	0.503 \pm 0.012	0.192 \pm 0.006	0.803 \pm 0.139	0.837 \pm 0.017	0.082 \pm 0.004	0.874 \pm 0.007

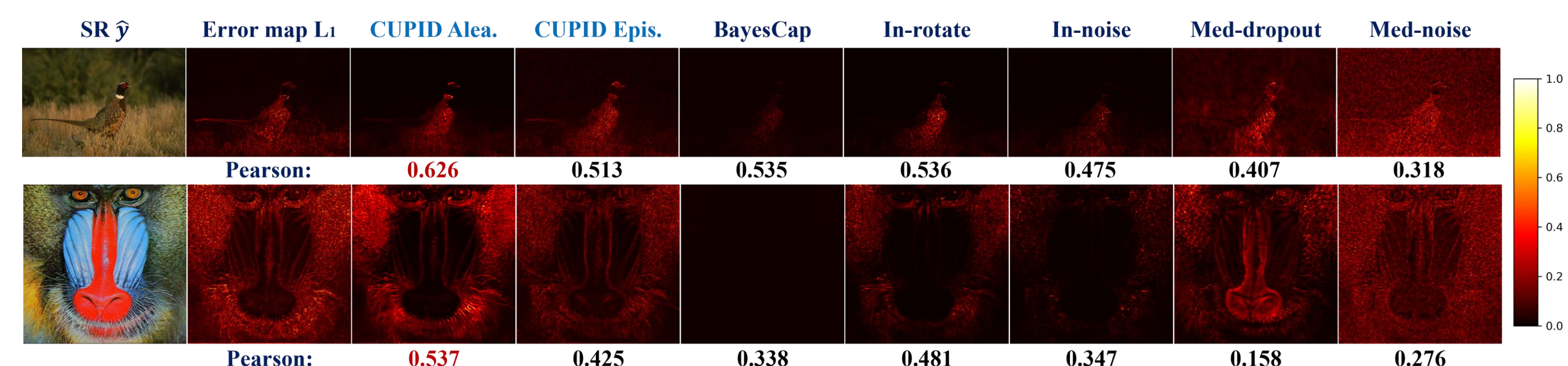
Table 2: Performance of OOD detection (OOD samples as positive). PAPILA and ACRIMA share the same research problem (glaucoma detection) with the ID dataset while CIFAR10 is a general classification dataset. CUPID Epistemic performs best on PAPILA and ACRIMA, suggesting that it is more sensitive to subtle distribution shifts within related clinical domains. In contrast, CUPID Aleatoric performs best on CIFAR-10, showing stronger response to extreme domain mismatch.

Method	PAPILA		ACRIMA		CIFAR10	
	AUC(\uparrow)	AUPR(\uparrow)	AUC(\uparrow)	AUPR(\uparrow)	AUC(\uparrow)	AUPR(\uparrow)
CUPID Alea.	0.379 \pm 0.027	0.333 \pm 0.007	0.717 \pm 0.029	0.661 \pm 0.027	0.983 \pm 0.005	0.998 \pm 0.001
CUPID Epis.	0.877 \pm 0.032	0.854 \pm 0.027	0.978 \pm 0.010	0.984 \pm 0.007	0.898 \pm 0.054	0.991 \pm 0.005
MC Dropout	0.733 \pm 0.002	0.586 \pm 0.007	0.869 \pm 0.003	0.816 \pm 0.009	0.887 \pm 0.004	0.986 \pm 0.001
Rate-in	0.328 \pm 0.005	0.329 \pm 0.008	0.363 \pm 0.003	0.390 \pm 0.003	0.620 \pm 0.001	0.927 \pm 0.002
IGRUE	0.636 \pm 0.114	0.486 \pm 0.097	0.941 \pm 0.008	0.944 \pm 0.008	0.978 \pm 0.005	0.998 \pm 0.001
PostNet Alea.	0.638 \pm 0.060	0.487 \pm 0.067	0.549 \pm 0.040	0.487 \pm 0.040	0.657 \pm 0.032	0.952 \pm 0.005
PostNet Epis.	0.577 \pm 0.097	0.425 \pm 0.088	0.685 \pm 0.154	0.654 \pm 0.151	0.773 \pm 0.082	0.976 \pm 0.011
BNN	0.707 \pm 0.040	0.612 \pm 0.050	0.708 \pm 0.073	0.699 \pm 0.042	0.643 \pm 0.108	0.959 \pm 0.013
DEC	0.515 \pm 0.024	0.457 \pm 0.024	0.680 \pm 0.003	0.685 \pm 0.012	0.660 \pm 0.015	0.963 \pm 0.003

Table 3: Performance on natural image datasets (Set5, Set14, BSDS100) and medical imaging dataset IXI (MRI scans). CUPID Aleatoric gives the best overall quantitative performance on the natural-image benchmarks, suggesting that reconstruction uncertainty is mainly driven by data-related ambiguity in this setting. On the IXI dataset, however, CUPID Epistemic becomes more informative, which is consistent with the larger domain shift from the training distribution.

Method	Set5			Set14		
	Pearson (\uparrow)	AUSE (\downarrow)	UCE (\downarrow)	Pearson (\uparrow)	AUSE (\downarrow)	UCE (\downarrow)
CUPID Alea.	0.528 \pm 0.006	0.010 \pm 0.000	0.045 \pm 0.018	0.527 \pm 0.002	0.012 \pm 0.000	0.049 \pm 0.005
CUPID Epis.	0.416 \pm 0.004	0.018 \pm 0.001	0.266 \pm 0.007	0.449 \pm 0.005	0.019 \pm 0.000	0.226 \pm 0.003
BayesCap	0.485 \pm 0.038	0.010 \pm 0.000	0.098 \pm 0.001	0.422 \pm 0.064	0.012 \pm 0.000	0.100 \pm 0.000
in-rotate	0.493 \pm 0.000	0.010 \pm 0.000	0.071 \pm 0.000	0.490 \pm 0.000	0.013 \pm 0.000	0.072 \pm 0.000
in-noise	0.370 \pm 0.006	0.019 \pm 0.000	0.051 \pm 0.035	0.354 \pm 0.001	0.022 \pm 0.000	0.826 \pm 0.006
med-dropout	0.219 \pm 0.023	0.030 \pm 0.001	0.680 \pm 0.043	0.271 \pm 0.012	0.024 \pm 0.000	0.292 \pm 0.022
med-noise	0.312 \pm 0.003	0.022 \pm 0.000	0.826 \pm 0.006	0.293 \pm 0.002	0.022 \pm 0.000	0.826 \pm 0.006
Method	BSDS100			IXI		
	Pearson (\uparrow)	AUSE (\downarrow)	UCE (\downarrow)	Pearson (\uparrow)	AUSE (\downarrow)	UCE (\downarrow)
CUPID Alea.	0.536 \pm 0.001	0.012 \pm 0.000	0.042 \pm 0.012	0.677 \pm 0.008	0.004 \pm 0.000	0.021 \pm 0.004
CUPID Epis.	0.464 \pm 0.007	0.018 \pm 0.000	0.185 \pm 0.007	0.734 \pm 0.018	0.004 \pm 0.000	0.298 \pm 0.013
BayesCap	0.427 \pm 0.034	0.011 \pm 0.000	0.100 \pm 0.000	0.447 \pm 0.034	0.004 \pm 0.000	0.100 \pm 0.000
in-rotate	0.465 \pm 0.000	0.012 \pm 0.000	0.077 \pm 0.000	0.598 \pm 0.000	0.004 \pm 0.000	0.093 \pm 0.000
in-noise	0.353 \pm 0.001	0.022 \pm 0.000	0.826 \pm 0.006	0.461 \pm 0.001	0.005 \pm 0.000	0.091 \pm 0.002
med-dropout	0.397 \pm 0.002	0.020 \pm 0.000	0.136 \pm 0.008	0.570 \pm 0.001	0.007 \pm 0.000	0.337 \pm 0.026
med-noise	0.293 \pm 0.000	0.024 \pm 0.000	0.700 \pm 0.002	0.439 \pm 0.000	0.006 \pm 0.000	0.859 \pm 0.002

Figure 3: Comparison of visual results between error and uncertainty maps. CUPID Aleatoric shows the best texture alignment and highest correlation with error maps.



The ablation results provide insight into how uncertainty develops inside the network. As shown in Figure 4, aleatoric uncertainty is estimated more effectively when CUPID is placed closer to the output, whereas epistemic uncertainty benefits more from earlier insertion points. This trend suggests that the two uncertainty types emerge differently across network depth, and that CUPID can be used not only for uncertainty estimation but also for analyzing uncertainty propagation.

Figure 4: Performance of CUPID inserted at varying locations: misclassification detection (Left) and super-resolution (Right).

