

Interpreting Optical Diffractive Neural Networks (ODNN) via Gradient-based Physical Activation Mapping (GPAM)

Hong-Yuech Huang¹, Cheng-Jun Huang¹, Ting-Siang Ou Yang¹, Jui-Chien Lin¹, Ming-Hua Lin¹ and Li-Bang Wang¹

¹ Department of Physics, National Tsing Hua University, Hsinchu, Taiwan 30013, Republic of China

Optical diffractive neural networks (ODNNs) compute through light propagation and diffraction, but their internal physical decision mechanism is difficult to interpret. We propose Gradient-based Physical Activation Mapping (GPAM), which estimates the contribution of each diffractive neuron by back-propagating gradients through physical wave propagation. We further establish a quantitative evaluation framework including localization, energy concentration, weakly-supervised segmentation, and faithfulness analysis using mask ratio. Experiments on MNIST and Fashion-MNIST demonstrate that GPAM reveals physically meaningful regions and that masking high-relevance regions significantly reduces confidence and accuracy. The proposed framework provides a tool for interpreting optical neural networks and analyzing optical system behavior.

Keywords: Optical diffractive neural networks, Interpretability, Optical computing, Saliency map, Optical neural networks

ODNN Setup & Datasets

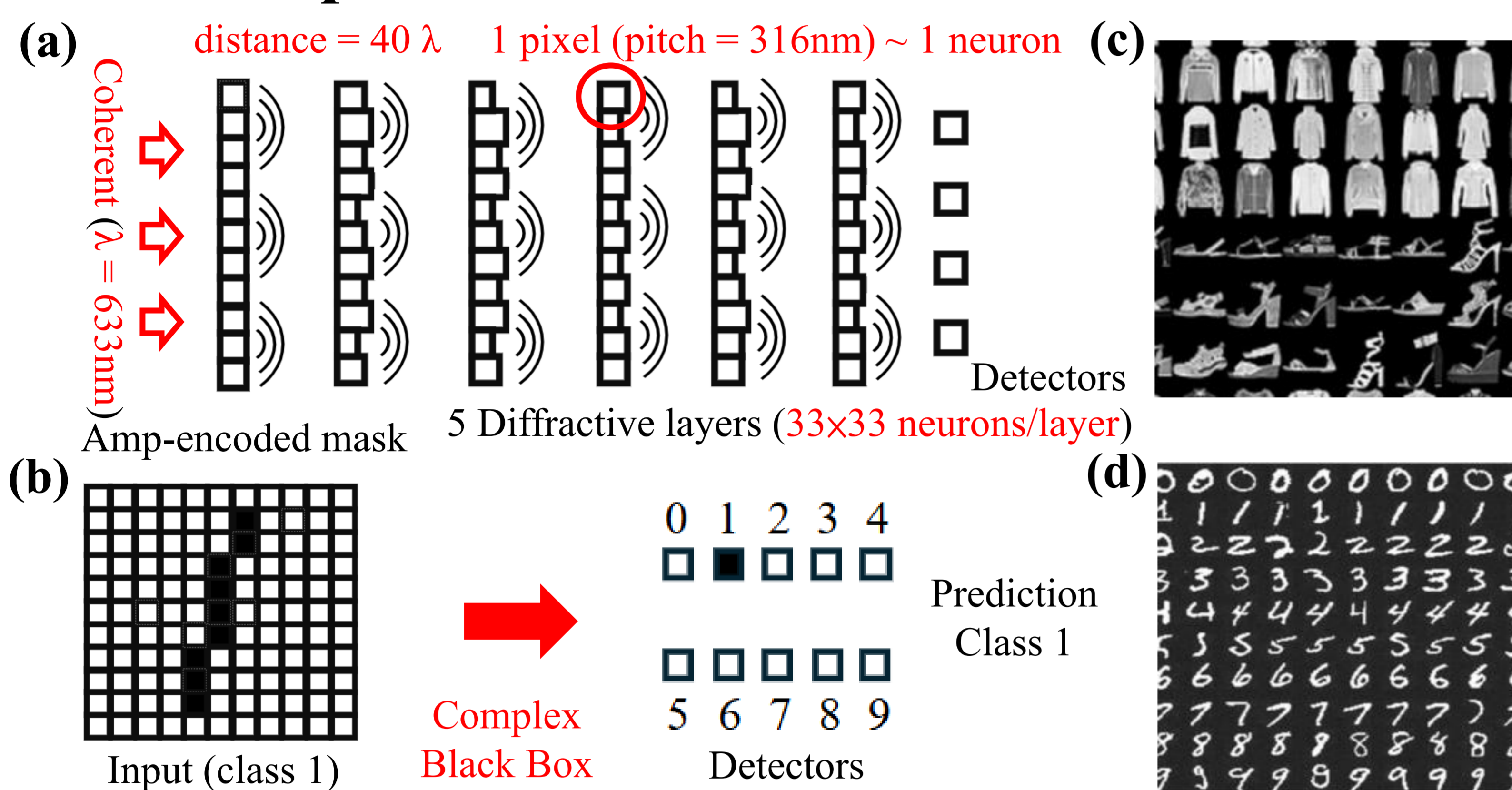


Fig. 1. ODNN system and datasets. (a) Schematic of the ODNN architecture. (b) Working principle of the system, where the ODNN acts as a complex black-box mapping from input to detector outputs. (c) Fashion-MNIST dataset. (d) MNIST dataset. (**Train/Test = 6w/1w**)

GPAM Method

Optical Forward Model

$$U_{l+1} = P_l D_l U_l$$

$$U_L = P_{L-1} D_{L-1} \dots P_1 D_1 U_0$$

$$I(x, y) = |U_L(x, y)|^2$$

$$y_k = \sum_{(x,y) \in D_k} I(x, y)$$

Backward Optical Field

$$G_L = M_k \cdot U_L^*$$

$$G_1 = P_1^\dagger G_{L+1}$$

Notation: $U_l(x, y)$ optical field at layer l , $U_L(x, y)$ detector field, $G_l(x, y)$ backward field, $R_l(x, y)$ relevance map, P_l propagation operator, P_l^\dagger adjoint propagation, D_l phase layer, M_k detector mask, $I(x, y) = |U_L(x, y)|^2$ intensity, y_k detector output.

Meaning: The relevance map is the spatial overlap between the forward optical field and the back-propagated detector field, indicating where optical information contributes to the final classification.

Pixel Perturbation Concept

$$U_1(x, y) \rightarrow U_1(x, y) + \delta U$$

$$\delta y_k = \frac{\partial y_k}{\partial U_1(x, y)} \delta U$$

$$\frac{\partial y_k}{\partial U_1(x, y)}$$

GPAM Relevance Definition

$$R_1(x, y) = \left| U_1(x, y) \cdot \frac{\partial y_k}{\partial U_1(x, y)} \right|$$

$$\frac{\partial y_k}{\partial U_1} = G_1$$

$$R_1(x, y) = |U_1(x, y) \cdot G_1(x, y)|$$

Results

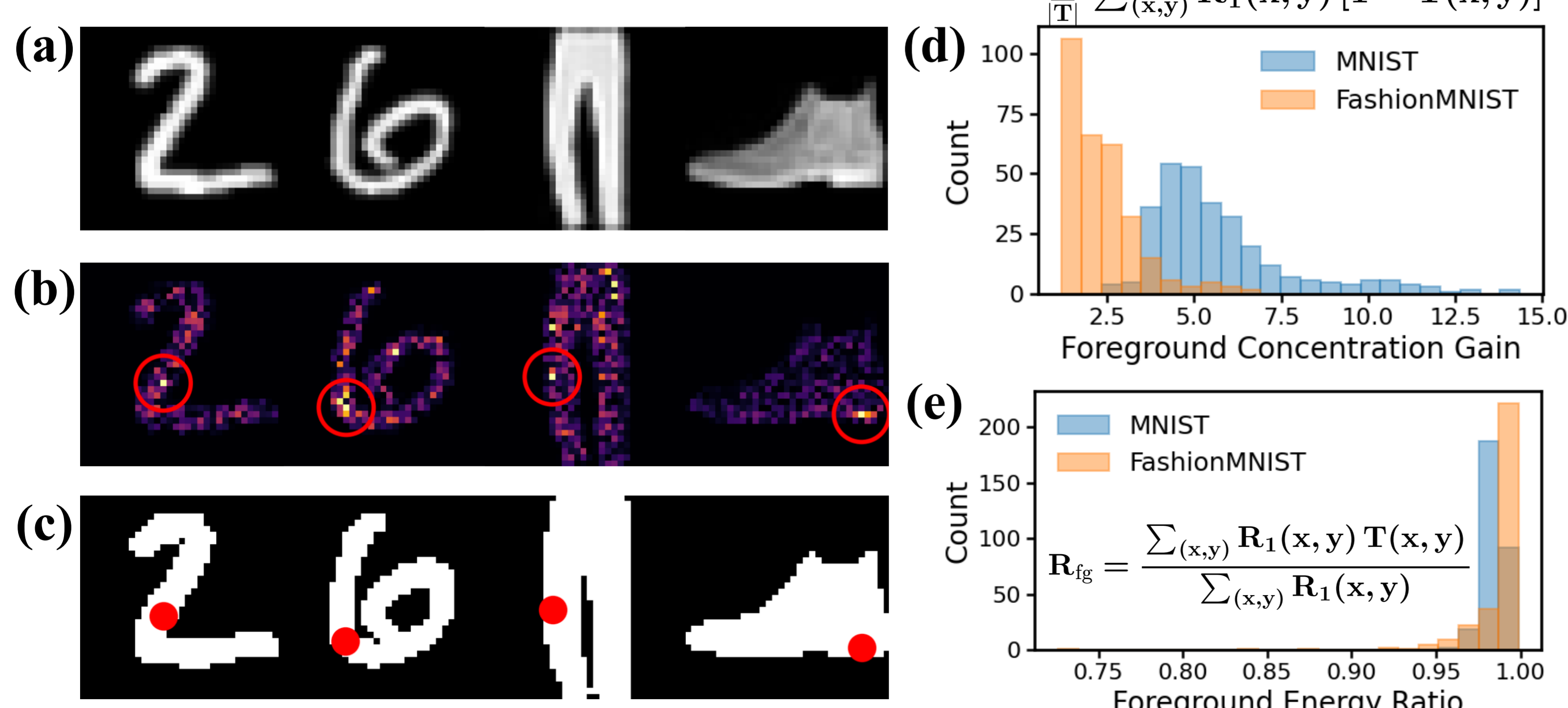


Fig. 2. GPAM localization and foreground energy analysis. Accuracy: 94.37% (MNIST), 92.75% (FashionMNIST). (a) Input images. (b) GPAM relevance maps R_1 with peak locations. (c) Ground truth masks T ; pointing game success = 100% for correct predictions. (d) Foreground concentration gain. (e) Foreground energy ratio.

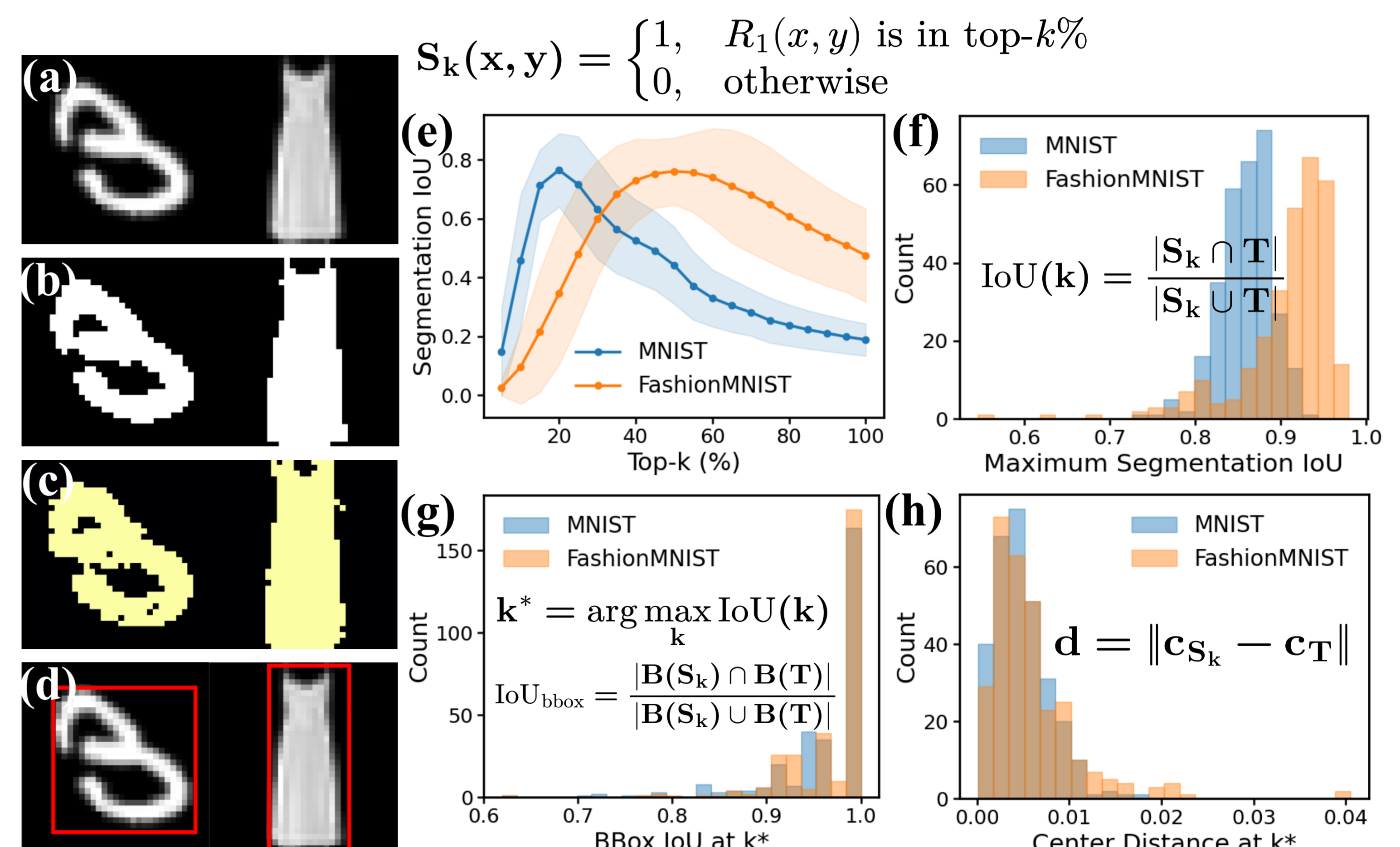


Fig. 3. GPAM weakly-supervised localization evaluation. (a) Input. (b) Ground truth. (c) Segmentation from GPAM top- k regions. (d) Bounding box derived from (c). (e) IoU vs top- k . (f) Maximum IoU. (g) BBox IoU at k^* . (h) Center distance at k^* .

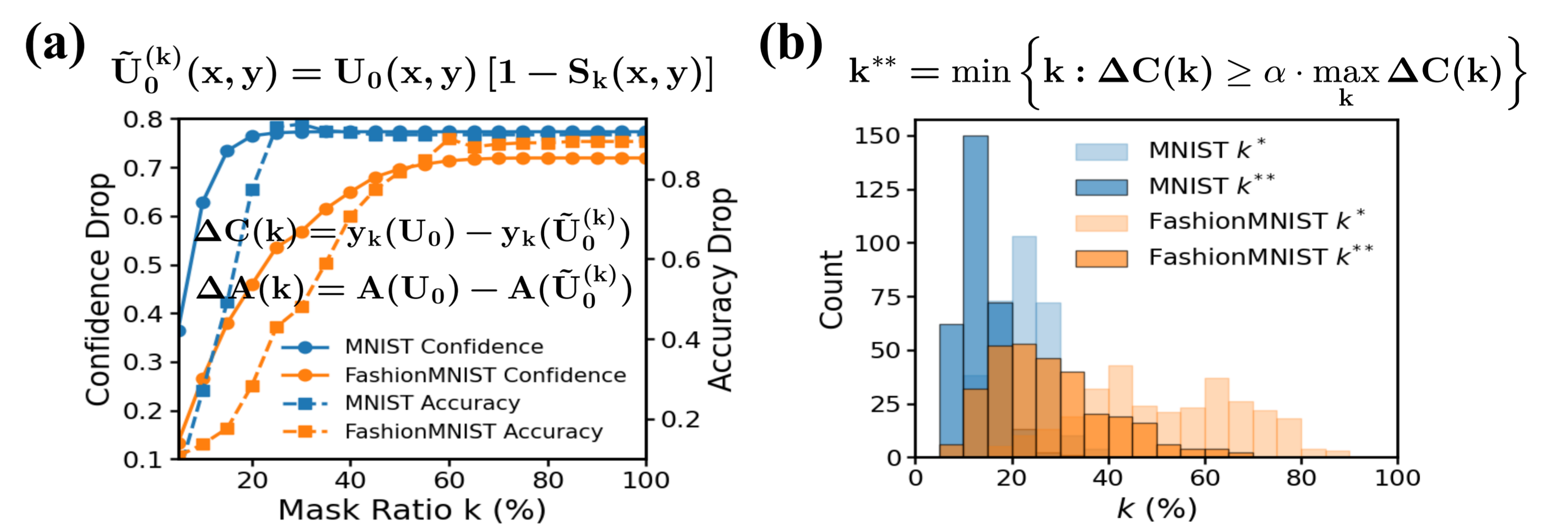


Fig. 4. Faithfulness and critical mask ratio analysis. (a) Confidence drop and accuracy drop as a function of mask ratio k when masking the top- k % GPAM relevance pixels. (b) Distribution of the critical mask ratios k^* (segmentation optimum) and k^{**} (faithfulness threshold) for MNIST and Fashion-MNIST datasets.

Discussion

Datasets images contain sparse foreground strokes on a dark background, which may bias localization methods toward foreground regions. As a control experiment, future work will introduce background illumination bias while maintaining classification accuracy, to verify that GPAM localization reflects decision-relevant regions rather than input sparsity.

GPAM is not only a visualization tool, but also a design and analysis tool for optical systems. **Potential Applications** were shown below.

- Physical AI interpretability
- Optical neural network debugging
- Optical computing optimization
- Optical system/inverse design
- Robust optical neural networks
- Optical microscopy & sensing

Conclusion

We proposed GPAM to interpret optical diffractive neural networks and established a quantitative framework to evaluate localization, segmentation, and faithfulness. Results show that GPAM reveals physically meaningful regions and that masking high-relevance regions significantly reduces prediction confidence and accuracy. This work moves optical neural networks from black-box models toward physically interpretable computing systems.

GPAM enables explainable optical neural networks !

What the system looks at when making a prediction?