## Supporting information

Stochastic Photo-responsive Memristive Neuron for an In-sensor Visual System based on Restricted Boltzmann Machine

**Jin Hong Kim[a]\*, Hyun Wook Kim[a], Min Jung Chung[a], Dong Hoon Shin[a], Yeong Rok Kim[a], Jaehyun Kim[a], Yoon Ho Jang[a], Sun Woo Cheong[a], Soo Hyung Lee[a], Janguk Han[a], Hyung Jun Park[a], Joon-Kyu Han[b]\*, and Cheol Seong Hwang[a]\***

[a.] *Department of Materials Science and Engineering and Inter-University Semiconductor Research Center, Seoul National University, Gwanak-ro 1, Gwanak-gu, Seoul 08826, Republic of Korea.*

[b.] *System Semiconductor Engineering and Department of Electronic Engineering, Sogang University, 35 Baekbeom-ro, Mapo-gu, Seoul 04107, Republic of Korea.*

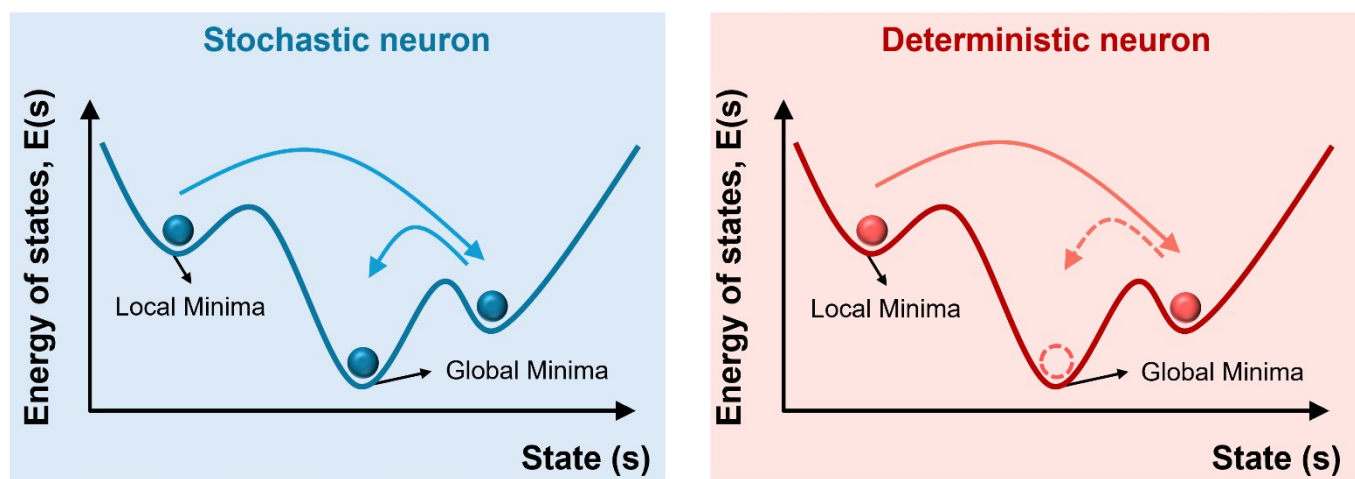*Corresponding author: joonkyuhan@sogang.ac.kr, cheolsh@snu.ac.kr*



**Fig. S1.** Comparison of local minima escapes when using stochastic neuron and deterministic neuron.

**Conventional ex-sensor visual system with a von Neumann processor**

| Light | Photodetector array | Signal transducer (ADC) | Memory | Processor |

Large data    Large data    Large data

**Ex-sensor visual system based on restricted Boltzmann machine (RBM)**

| Light | Photodetector array | Signal transducer (ADC) | RBM with a random number generator | Supportive neural network |

Large data    Large data    Reduced data

**In-sensor visual system based on RBM with a stochastic photo-responsive memristive neuron**

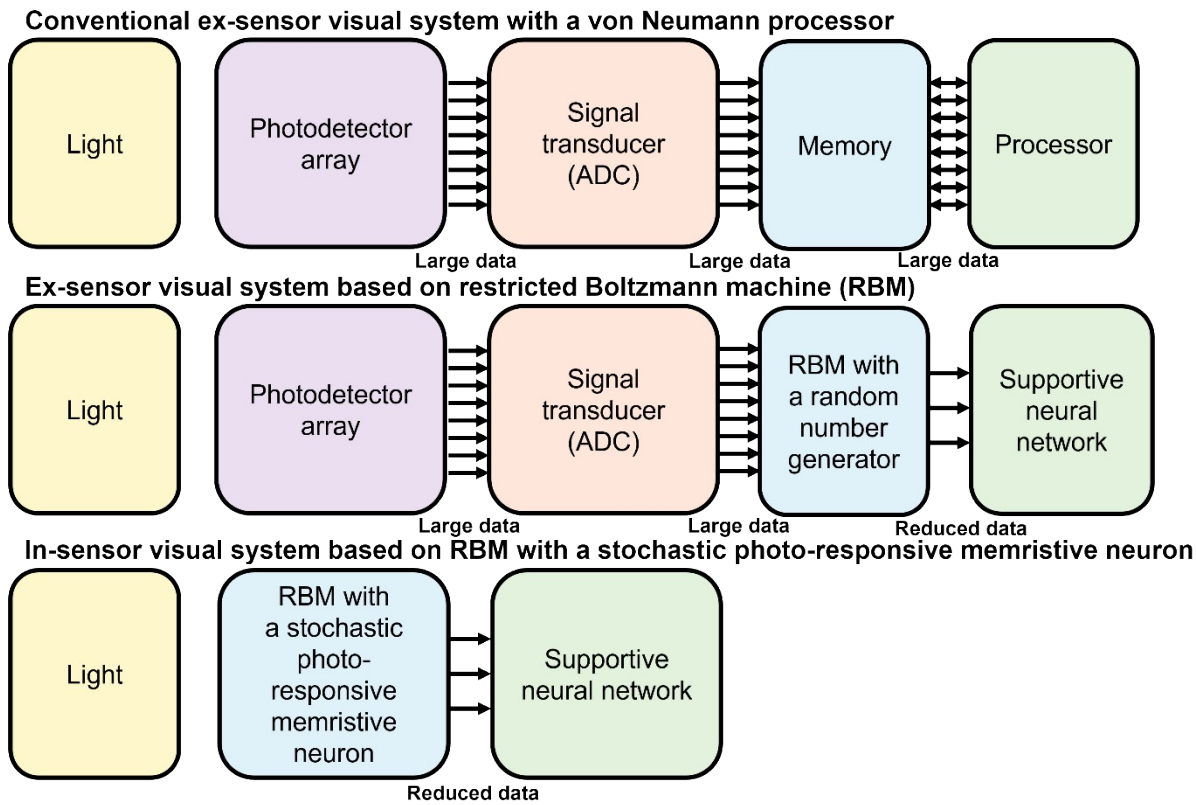| Light | RBM with a stochastic photo-responsive memristive neuron | Supportive neural network |

Reduced data

**Fig. S2.** Block diagram comparison of a conventional ex-sensor visual system with a von Neumann processor, ex-sensor visual system based on RBM, and the proposed in-sensor visual system based on RBM.

   The block diagram in **Fig. S2** compares the structure of a conventional ex-sensor visual system with conventional von Neumann architecture, an ex-sensor visual system with RBM, and a proposed in-sensor visual system with RBM. In the conventional ex-sensor visual system, a photodetector array receives light, and an analog-to-digital converter (ADC) converts the analog signal from the photodetector array into the digital signal.[1,2] This information is then stored in memory and provided to a processor for data processing. However, significant bottlenecks present as the massive data moves from the photodetector array to the memory through the ADC or the memory to the processor.[3–5]   Using the RBM eliminates bottlenecks between the memory and processor. Significantly, the RBM enables faster convergence and lower computational cost than other neural networks, making it desirable in energy-efficient visual systems.[6] However, the bottleneck still exists at the ADC in the ex-sensor visual system with RBM because there is an additional random number generator, such as ring oscillators, to stochastically update the neuron states after receiving digital signals from the ADC.[7] Therefore, an in-sensor visual system based on RBM significantly reduces the physical and computational distances between the photodetector array and the input neuron of RBM by coupling each element using a stochastic photo-responsive neuron. Consequently, prominently reduced data transmission ensures low latency and low power consumption.
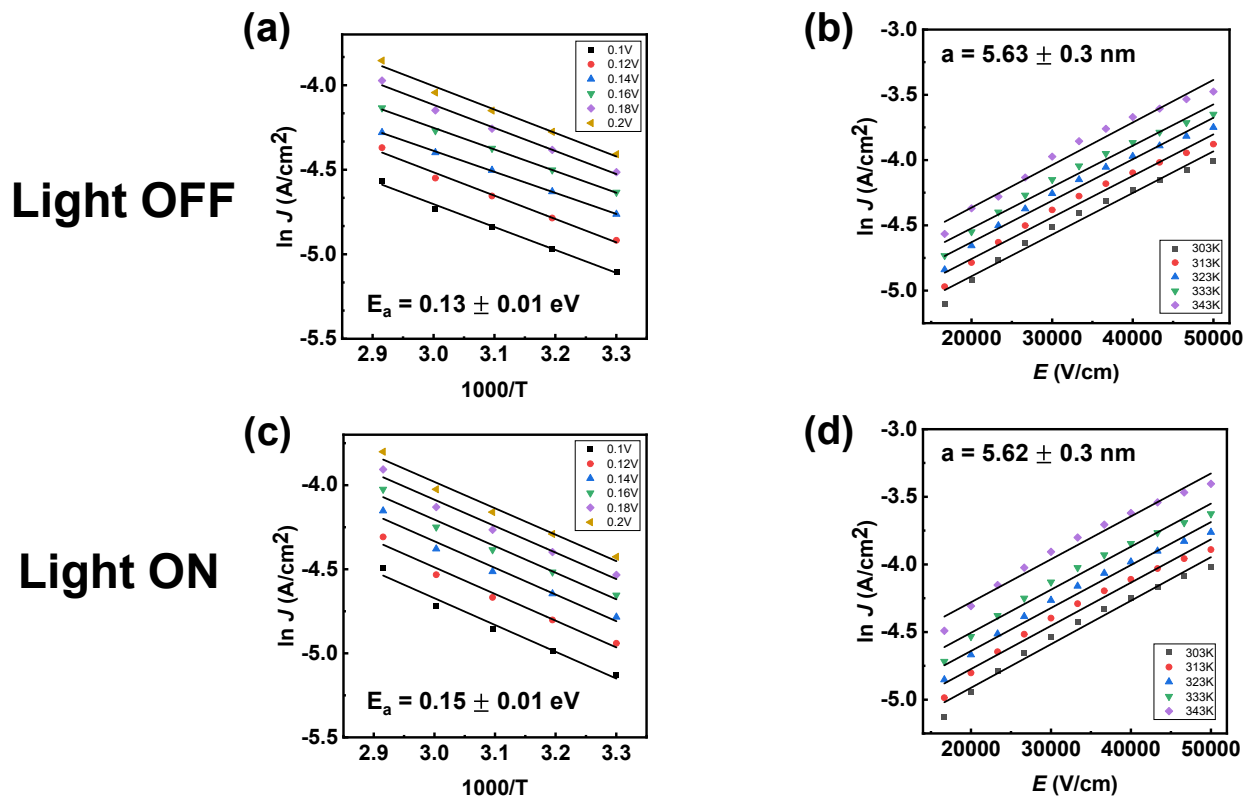
**Fig. S3.** (a) Arrhenius plot of the hopping conduction at low fields of the TIT device at light OFF state. (b) ln(*J*)-*E* characteristics of hopping conduction of the TIT device at light OFF state. (c) Arrhenius plot of the hopping conduction at low fields of the TIT device at light ON state. (d) ln(*J*)-*E* characteristics of hopping conduction of the TIT device at light ON state.
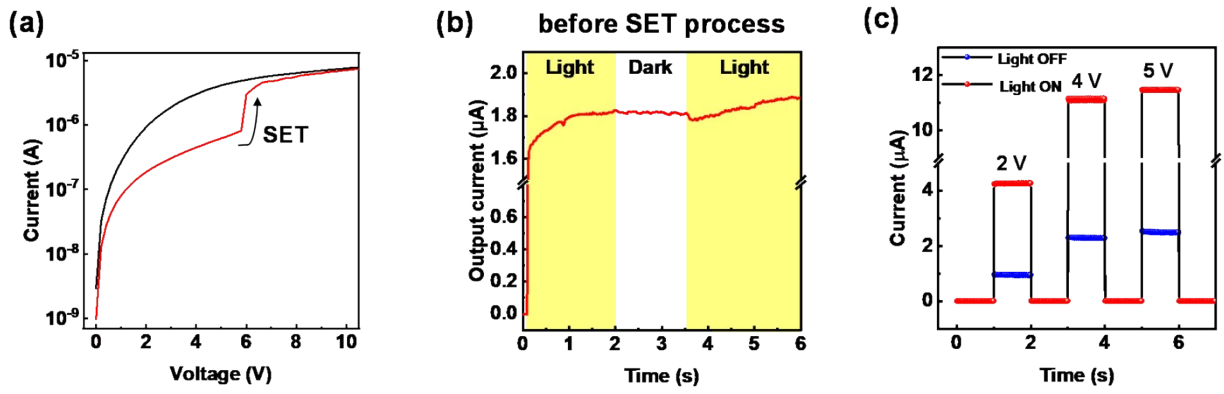
**Fig. S4.** (a) SET process of the TIT optoelectronic memristor. The black and red curves represent the *I-V* curves before and after the SET operation. (b) Photo-response of the device before the SET process. (c) Electrical pulse test at various conditions. Blue and red circles indicate the output current at light OFF and ON conditions, respectively. A 2 V, 4 V, and 5 V pulse was applied in sequential order.
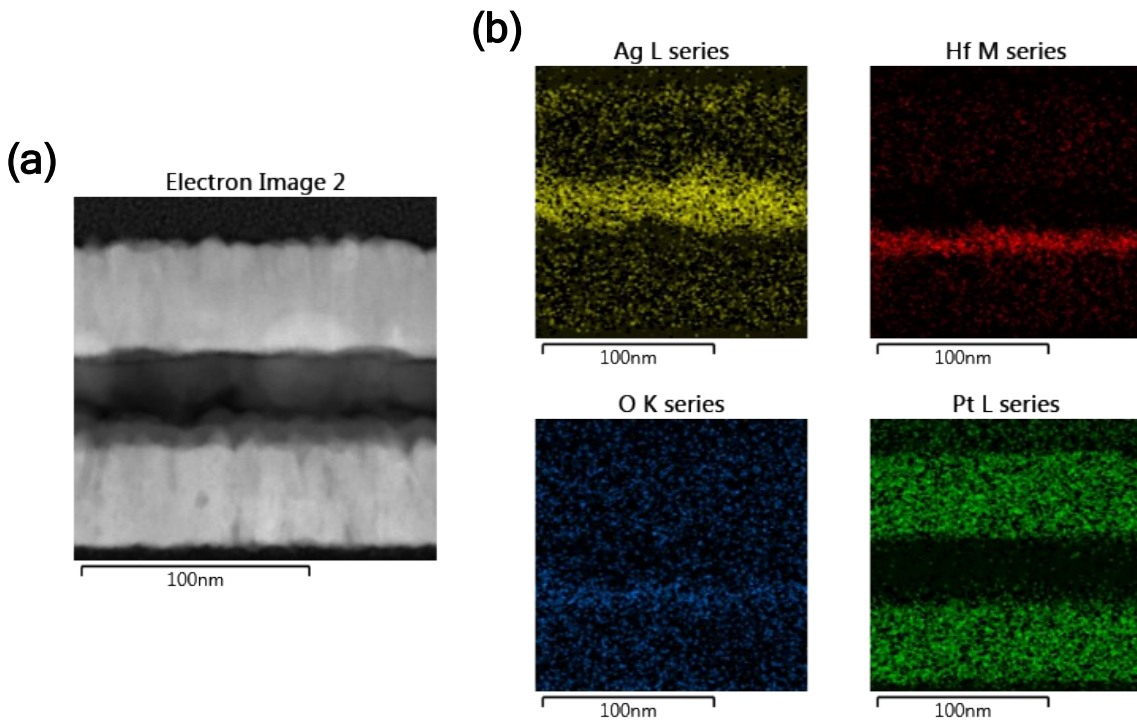


**Fig. S5.** (a) Cross-sectional transmission electron microscopy (TEM) image and (b) electrical dispersive spectroscopy (EDS) results.
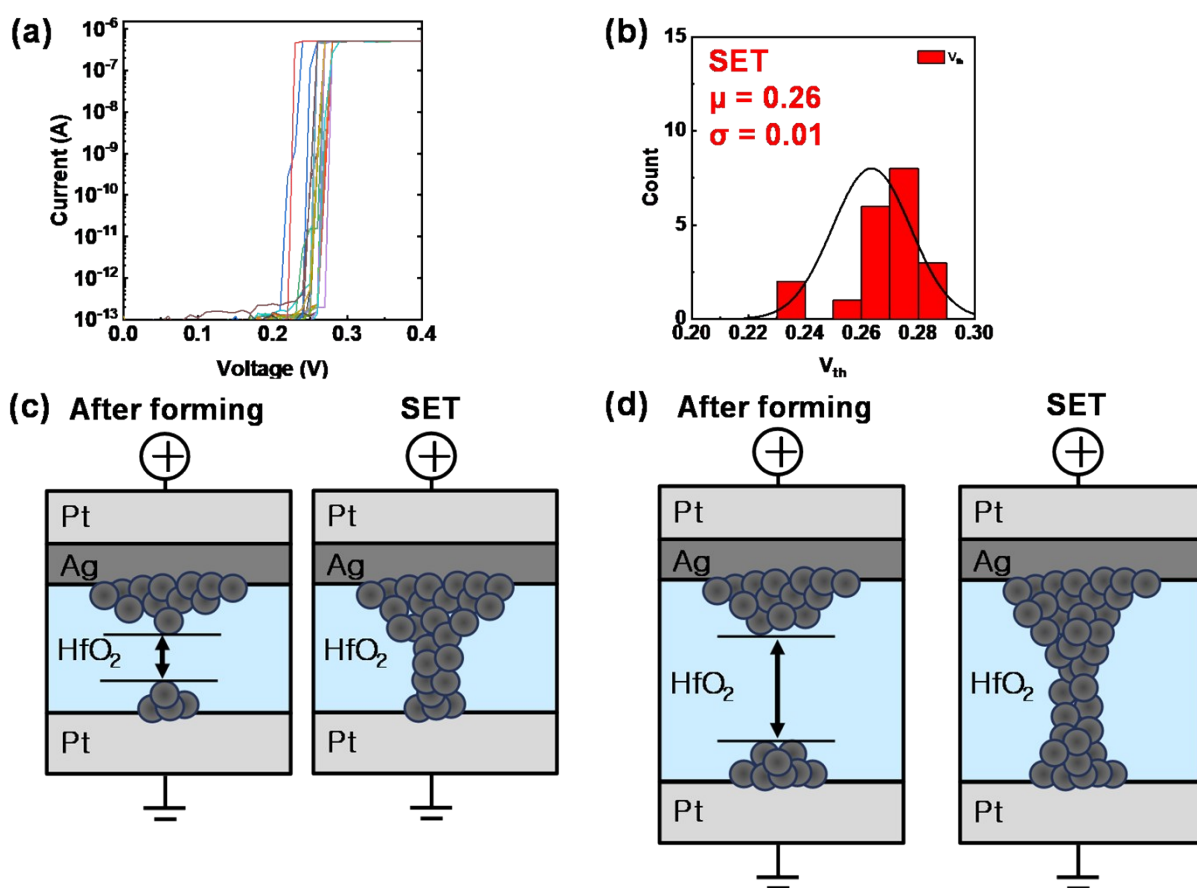
**Fig. S6.** (a) *I-V* characteristics and (b) statistical distribution of $V_{th}$ of the AHP device with an $HfO_2$ thickness of 3 nm. Schematic illustrations of the threshold-switching mechanism of the AHP devices with (c) 3 nm and (d) 8 nm $HfO_2$ layers.
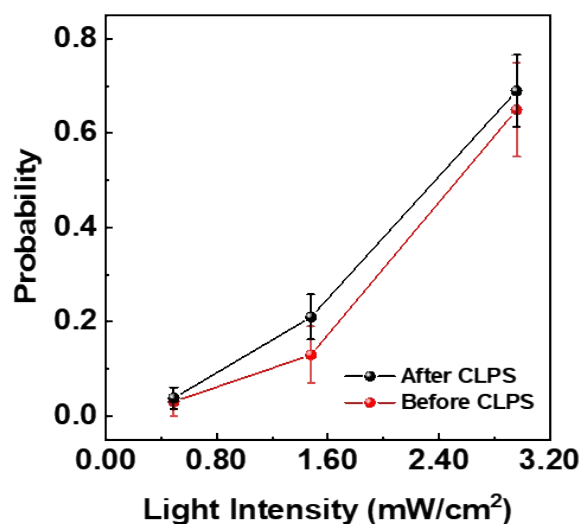
**Fig. S7.** Switching probability of photo-responsive neurons observed at light intensities of 0.49 mW/cm², 1.48 mW/cm², and 2.96 mW/cm² before and after the closed-loop pulse switching (CLPS) endurance test. The error bar shows the variability in 10 measurements.

In this work, the pulse switching endurance measurement was conducted by the CLPS method with a custom-made board, as detailed in previous work.[8] This method involves applying incremental pulses until the desired resistance is achieved. Specifically, during the SET process, positive incremental pulses were applied from 0.05 V with a step size of 0.05 V until the device reached the LRS resistance, which was set as 2 MΩ for this work. After each pulse, the resistance value was measured with a read pulse of 0.2 V. Once the device reached the LRS resistance, negative pulses of increasing magnitude (in absolute value) were immediately applied, starting from -0.2 V with a step size of -0.05 V, to initiate the RESET process until the target HRS resistance of 50 MΩ was attained. It should be noted that the AHP threshold-switching memristor exhibited very high resistance, which exceeded the measurement system's resolution limit. Consequently, the CLPS measurement resulted in the HRS region containing numerous data points that resembled a uniform resistance profile in the HRS region.
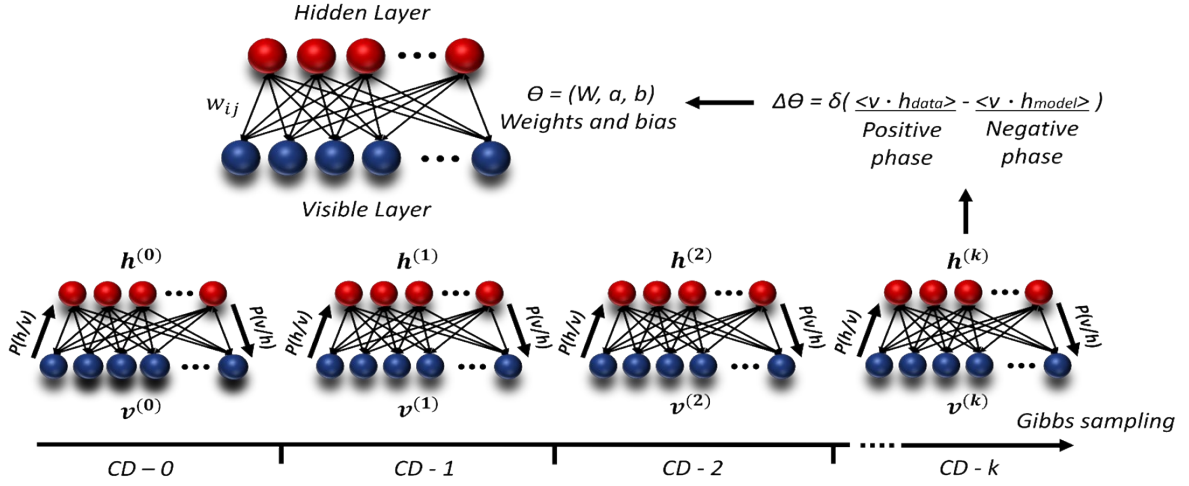
**Fig. S8.** Detailed learning process of a restricted Boltzmann machine (RBM). $v^{(0)}$ represents the state of visible neurons determined by light, and $h^{(0)}$ denotes the state of hidden neurons determined based on a given switching probability. $v^{(k)}$ and $h^{(k)}$ represent the states determined after $k$ iterations of sampling for the visible layer neurons and hidden layer neurons, respectively. The term $p(h/v)$ defines the probability of hidden neurons becoming active given the state of visible neurons, expressed in a sigmoid form. $\vartheta$ represents the set of parameters that are updated through the contrastive divergence (CD) algorithm.

The training process of an RBM involves updating the weights using the contrastive divergence (CD) algorithm followed by Gibbs sampling. The CD algorithm is divided into two main phases: the positive phase and the negative phase. In the positive phase, the goal is to capture the characteristics of the input data based on the actual data distribution. During this phase, the input data is presented to the visible layer of the RBM, from which the activation probabilities of the hidden layer are calculated to derive the expectation concerning the input data. The negative phase then generates new data samples through Gibbs sampling, based on the current weights, to learn about the distribution of the model. This phase calculates the expectation for the actual data.

The weight update is performed by calculating the difference between the expectation of the data (from the positive phase) and the model's expectation (from the negative phase), multiplying this difference by the learning rate, and adding it to the previous weights. By iteratively adjusting the weights through this process, the hidden layer of the RBM learns to capture hidden features that can accurately reconstruct the input data. This process allows the RBM to approximate the complex probability distribution of the input data more effectively over time.

In Gibbs sampling, the activation probabilities of the hidden neurons are determined based on the inputs from the visible layer. Following this, the activation probabilities of the input neurons in the visible layer are calculated based on the activation states of the hidden neurons, resulting in a reconstructed output. This process includes calculating the sum of the weighted inputs and the bias term for each neuron. Specifically, the neurons probabilistically determine their state of either 0 or 1 using the following equations.

$$p(h_j = 1|v) = \frac{1}{1 + e^{-(\sum_i v_i w_{ij} + b_j)}} \quad p(v_i = 1|h) = \frac{1}{1 + e^{-(\sum_j h_j w_{ij} + a_i)}}$$

In this equation, $v$ represents the state of the input neurons, $h$ represents the state of the hidden neurons, and $w$ indicates the weights between the two layers. Also, $v_i$ and $h_j$ represent the state of the $i$-th input neuron and the $j$-th hidden neuron, respectively.

Using the features learned by the RBM, a fully connected neural network finely tunes the weights through backpropagation for classification. This approach leverages the efficiency of the RBM in extracting features, allowing for precise classification.
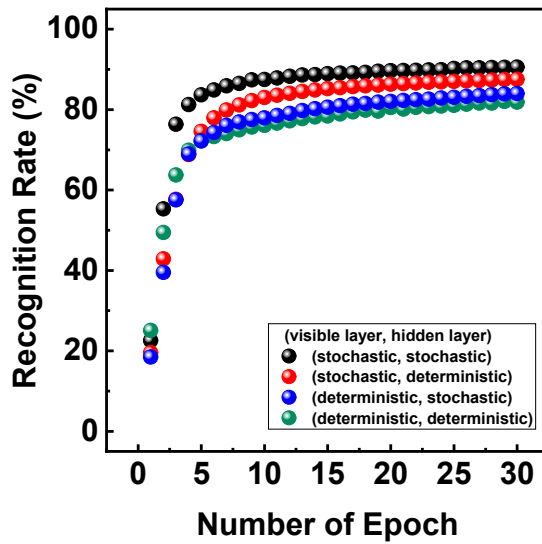
**Fig. S9.** The recognition rate of MNIST handwritten digits depends on whether the input neurons in the visible layer and the hidden neurons in the hidden layer are stochastic or deterministic.
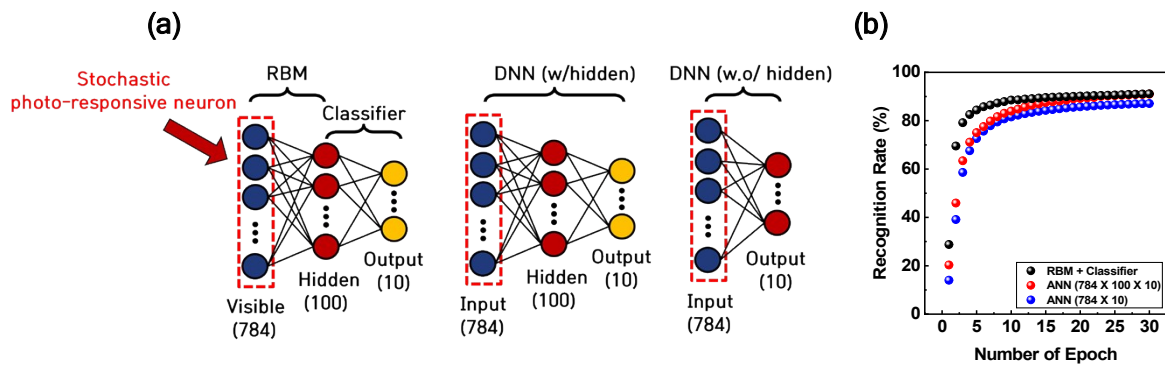


**Fig. S10.** (a) Schematics of RBM network, deep neural network with hidden layer, and deep neural network without hidden layer. (b) Recognition rate of MNIST handwritten digits according to the epochs in each case.
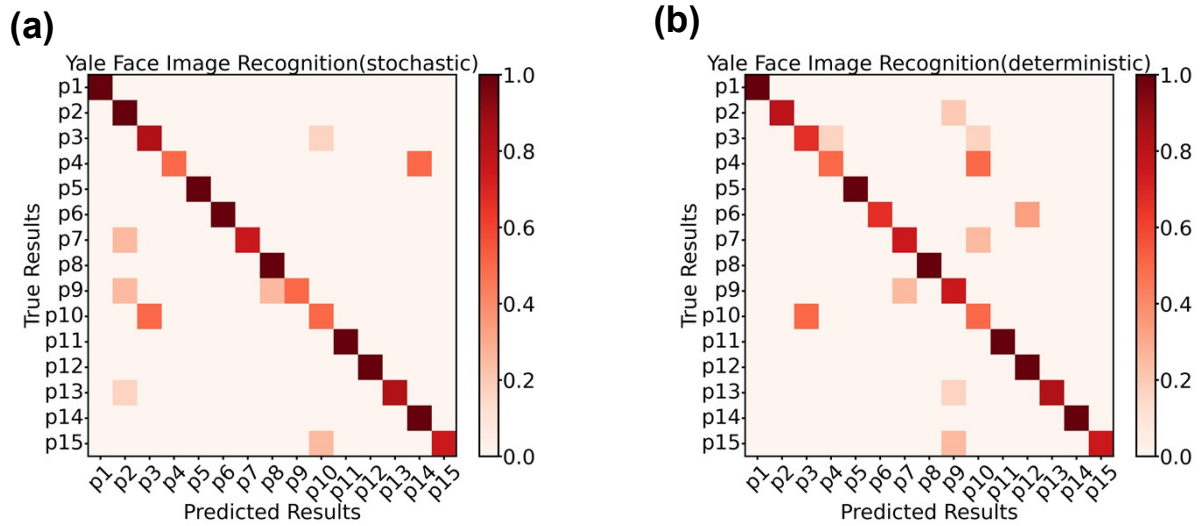
**Fig. S11.** Confusion matrix of Yale face image recognition using (a) the stochastic neurons and (b) the deterministic neurons.
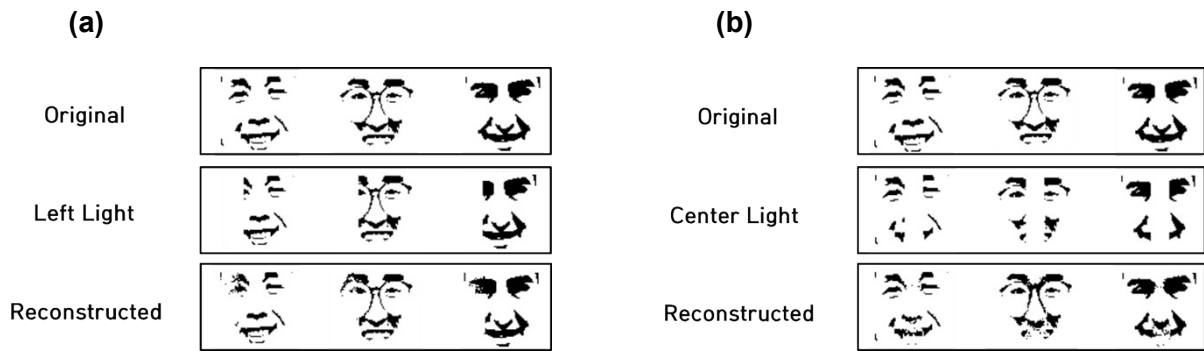


**Fig. S12.** Image reconstruction of (a) left-side light-shone face images and (b) center light-shone face images using the in-sensor RBM.

**Table S1.** Comparison of the MNIST recognition task accuracy using the restricted Boltzmann machine (RBM) between this study and previously reported research.

| Reference | Accuracy (%) | Device (Neuron) | Network size | Image reconstruction | Light responsive system |
|---|---|---|---|---|---|
| Software | 91.7 | - | 784 x 100 x 10 | - | - |
| Mao et al., *Adv. Electron. Mater.*, 2022, **8**, 2100918.[9] | 91.2 | Ag/IGZO/ITO | 784 x 500 x 10 | No | No |
| Heo et al., *Adv. Sci.*, 2024, 2405768.[10] | 90.63 | W/ZnTe/W | 784 x 500 x 10 | No | No |
| Li et al., *Nano Lett.*, 2024, **24**, 5420–5428.[11] | 93.0 | SOT - MTJ | 25 x 2 | Yes | No |
| This study | 90.9 | Ag/HfO$_2$/Pt | 784 x 100 x 10 | Yes | Yes |

# References

1    P. Wu, T. He, H. Zhu, Y. Wang, Q. Li, Z. Wang, X. Fu, F. Wang, P. Wang, C. Shan, Z. Fan, L. Liao, P. Zhou and W. Hu, InfoMat, 2022, **4**, 12275.

2    J. K. Han, Y. W. Chung, J. Sim, J. M. Yu, G. B. Lee, S. H. Kim and Y. K. Choi, Sci. Rep., 2022, **12**, 1818.

3    D. Kimovski, N. Saurabh, M. Jansen, A. Aral, A. Al-Dulaimy, A. B. Bondi, A. Galletta, A. V. Papadopoulos, A. Iosup and R. Prodan, IEEE Internet Comput., 2023, **28**, 6-16.

4    C. Yang, B. Sun, G. Zhou, T. Guo, C. Ke, Y. Chen, J. Shao, Y. Zhao and H. Wang, ACS Mater. Lett., 2023, **5**, 504–526.

5    K. Roy, A. Jaiswal and P. Panda, Nature, 2019, **575**, 607–617.

6    E. Neftci, S. Das, B. Pedroni, K. Kreutz-Delgado and G. Cauwenberghs, *Front. Neurosci.,* 2014, **7**, 272.

7    M. Jerry, A. Parihar, B. Grisafe, A. Raychowdhury, and S. Datta, in 2017 Symposium on VLSI Technology (VLSIT), IEEE, Kyoto, Japan 2017, pp. T186-T187.

8    H. J. Kim, T. H. Park, K. J. Yoon, W. M. Seong, J. W. Jeon, Y. J. Kwon, Y. Kim, D. E. Kwon, G. S. Kim, T. J. Ha, S. G. Kim, J. H. Yoon and C. S. Hwang, Adv. Funct. Mater., 2019, **29**, 1806278.

9    H. Mao, Y. He, C. Chen, L. Zhu, Y. Zhu, Y. Zhu, S. Ke, X. Wang, C. Wan and Q. Wan, Adv. Electron. Mater., 2022, **8**, 2100918.

10   J. Heo, S. Kim, S. Kim and M. H. Kim, Adv. Sci., 2024, 2405768.

11   X. Li, C. Wan, R. Zhang, M. Zhao, S. Xiong, D. Kong, X. Luo, B. He, S. Liu, J. Xia, G. Yu and X. Han, *Nano Lett.*, 2024, **24**, 5420–5428.