

Optical neural network

**What's the difference between computational imaging and
optical computing?**

Outline

- **Analogy of ONN**
- ONN implementation
- A glimpse
- A optics-inspired design

History

inline hologram

- Inline: The Reference wave is **co-axis** with object wave.
- Hologram: Two wave interfere beyond the object and the interference pattern recorded **at some distance**.
- Plane wave v.s. spherical wave (Gabor hologram)

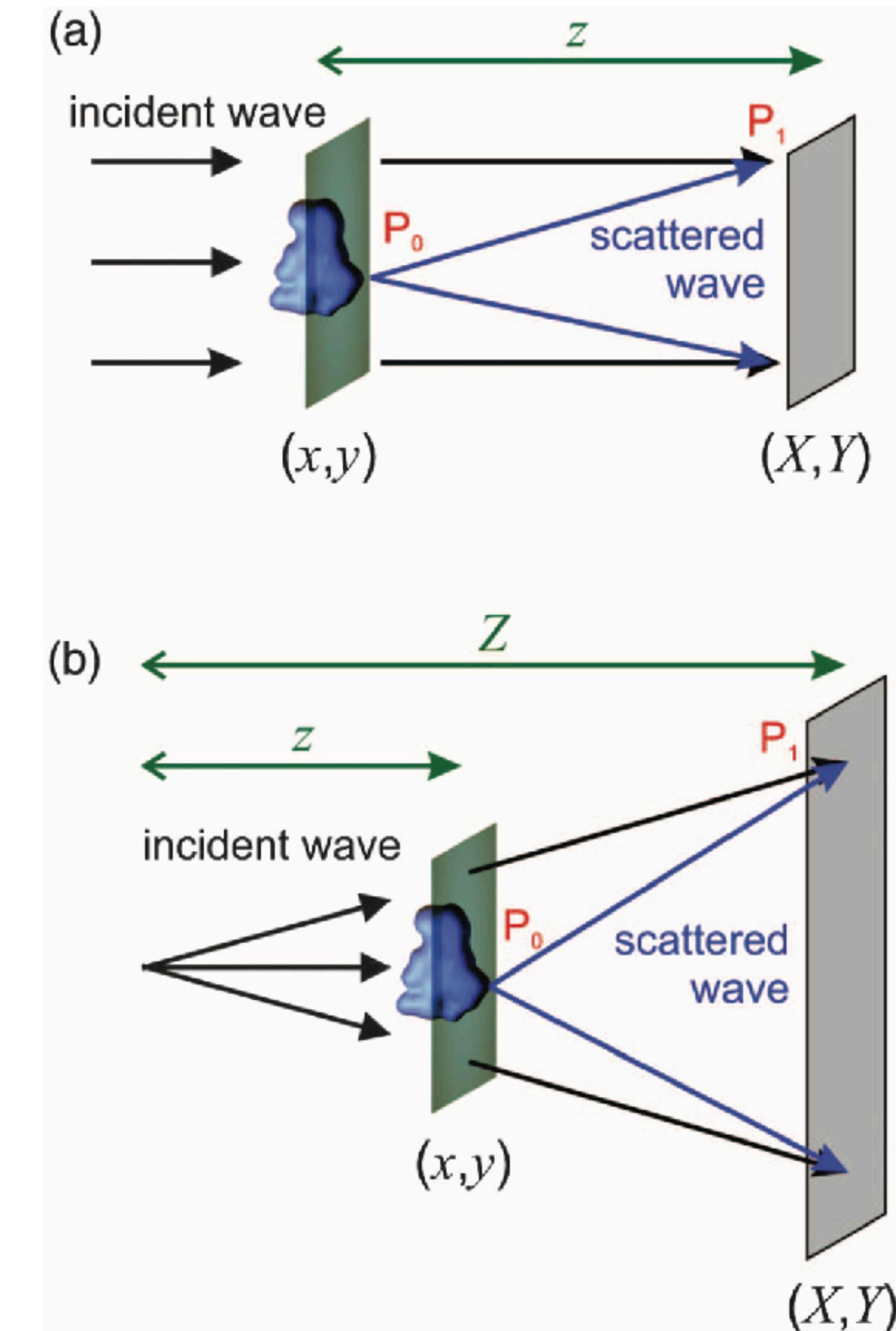


Fig. 1: In-line holography schemes realized with a **(a)** plane wave **(b)** a spherical wave

History

inline hologram

- Forwarding
- Back-propagation

object plane's incident optical wave & exit wave :
 $U_{\text{exit wave}}(x,y) = U_{\text{incident}}(x,y) \cdot t(x,y)$ (1)

where $t(x,y)$ is transmission function that is
 $t(x,y) = e^{-\alpha x,y} + i\phi(x,y)$
 absorption phase distribution,
 $= |t(x,y)|$. (2)

(2) \rightarrow (1) \Rightarrow
 $U_{\text{exit wave}}(x,y) - U_{\text{incident}}(x,y) = U_{\text{incident}} \cdot f(x,y)$.

The recorded wave of image plane :

$$\begin{aligned} U_{\text{image}}(X,Y) &= \frac{1}{\lambda i} \iint U_{\text{incident}}(x,y) \cdot t(x,y) \cdot \frac{e^{(ik\vec{x}\vec{r})}}{|z|} dx dy \\ &= K \iint |t(x,y)| \cdot \frac{e^{i\frac{\lambda z}{\lambda z}[(x-X)^2 + (y-Y)^2]}}{|(x-X)^2 + (y-Y)^2 + z^2|^{\frac{1}{2}}} dx dy \\ |z| \gg \rho^2/\lambda &= K \iint |t(x,y)| \cdot \frac{e^{i\frac{\lambda z}{\lambda z}[(x-X)^2 + (y-Y)^2]}}{|z|} dx dy \\ &= \frac{K}{z} \iint t(x,y) \cdot e^{i\frac{\lambda z}{\lambda z}[(x-X)^2 + (y-Y)^2]} dx dy \\ &= \boxed{\frac{K}{z} \cdot t(X,Y) \otimes S(X,Y)} \quad \text{Conv.} \end{aligned}$$

The recorded hologram $\Rightarrow H(X,Y) = |U_{\text{image}}(X,Y)|^2$

$$= |R|^2 + |O|^2 + R^*O + RO^*$$

normalization $\Rightarrow H_{\text{norm}} = 1 + \frac{|O|^2}{|R|^2} + \frac{R^*O}{|R|^2} + \frac{RO^*}{|R|^2}$

$$\Rightarrow H_0 = H_{\text{norm}} - 1 = \frac{R^*O + RO^*}{|R|^2} \quad (\text{bg} = 0)$$

* H_0 · hanning to reduce hologram edge effect

also Fresnel-Kirchoff $\Rightarrow U(X,Y) = K \iint R(X,Y) H_0(X,Y) \cdot e^{i\frac{\lambda k\vec{x}\vec{r}}{|z|}} dx dy$.

$R=1$ & far-field condition \Rightarrow

$$\begin{aligned} &= \frac{K}{z} \iint |H_0| \cdot e^{-\frac{|k\vec{x}\vec{r}|}{|z|}} dx dy \\ &= \boxed{\frac{K}{z} \cdot H_0 \otimes S^*} \end{aligned}$$

$$+ \frac{U}{R} \Rightarrow f + I = t = e^{-\alpha} e^{-i\phi}$$

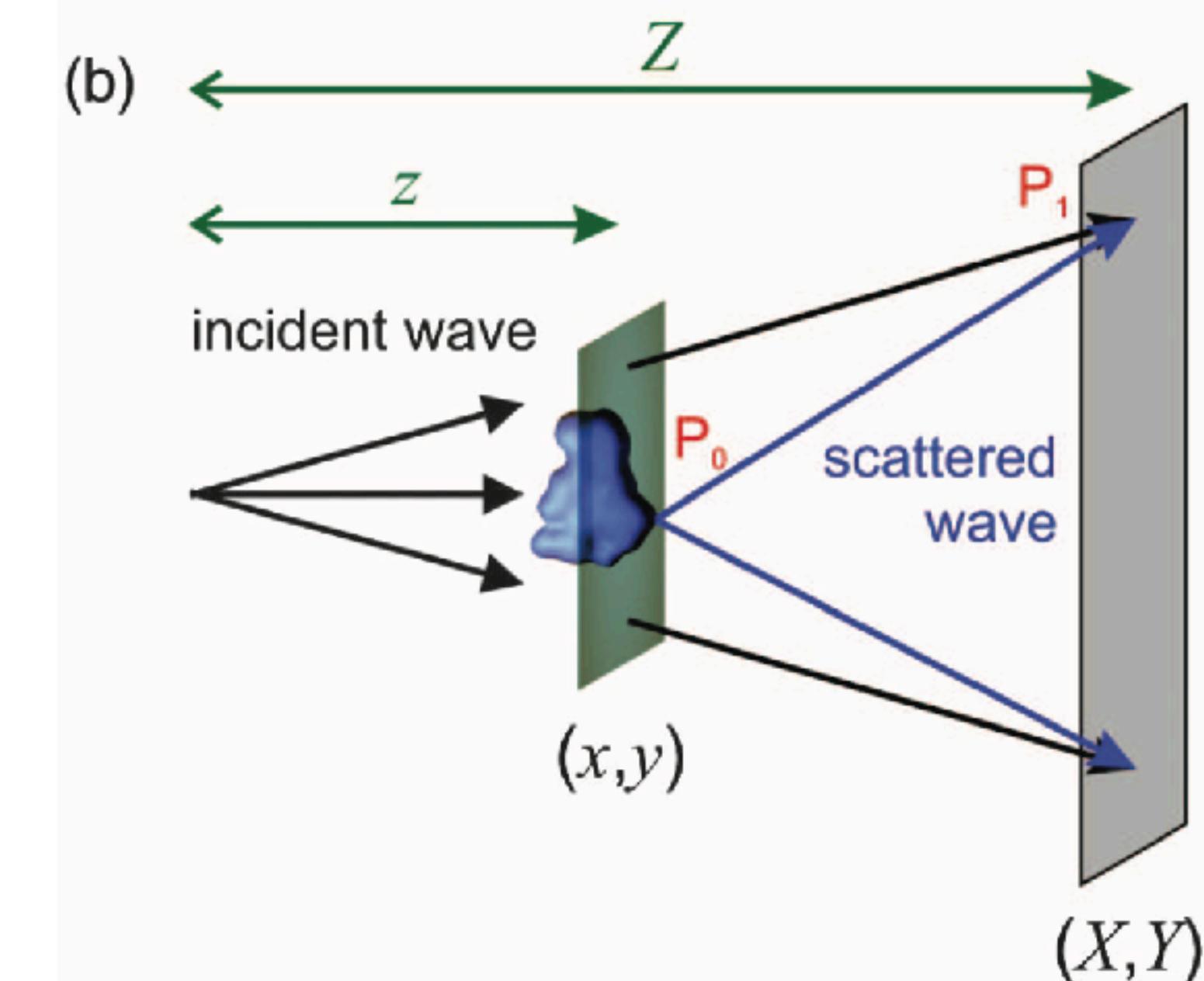
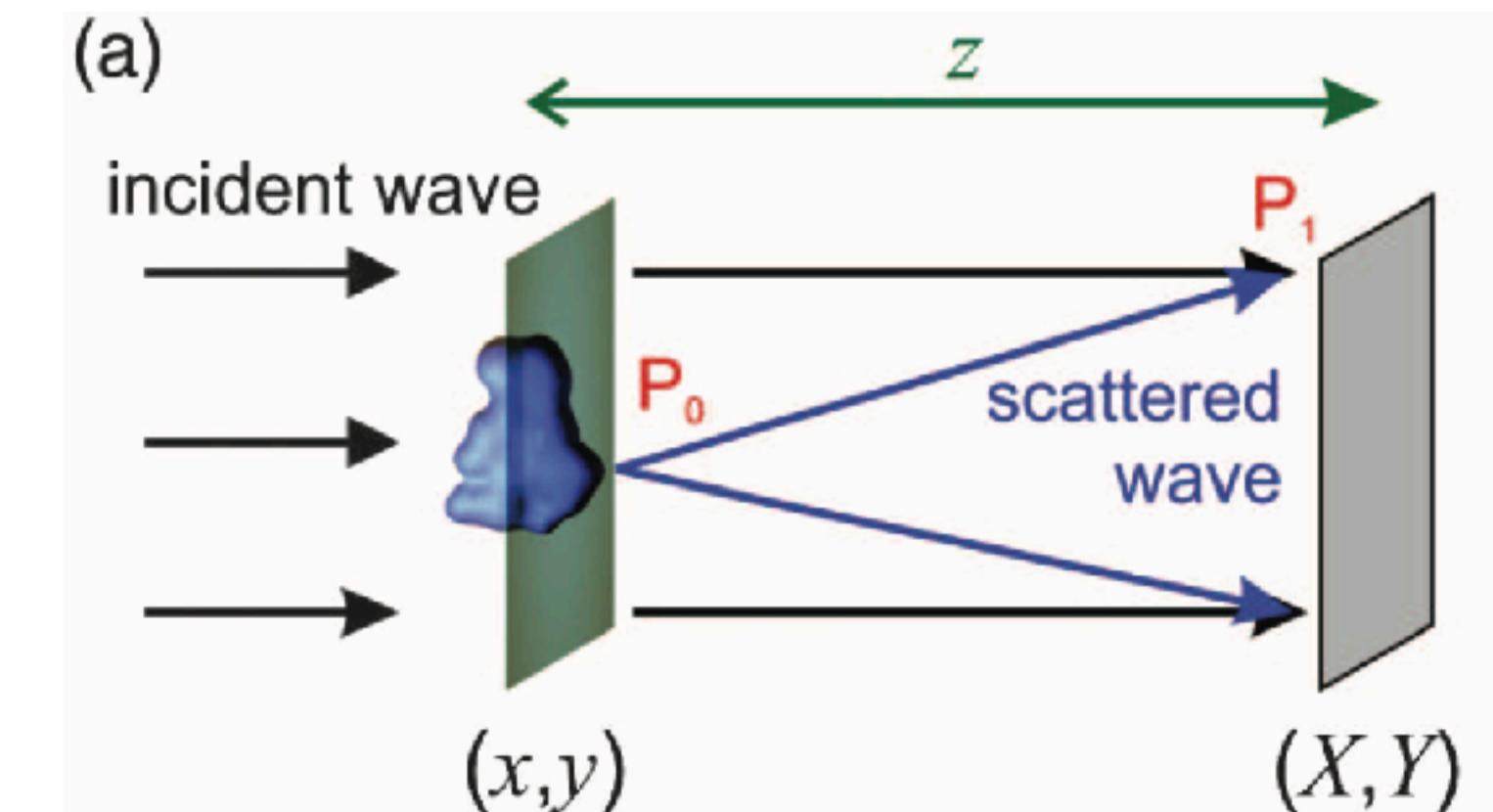


Fig. 1: In-line holography schemes realized with a **(a)** plane wave **(b)** a spherical wave

History

Optical neural network

- A optical implementation of an ANN
- Hopfield neural netowrk
- Dispersive medium for optical analog of the neural network Applied Optics, Marcus S. Cohen
- Diffractive neural network

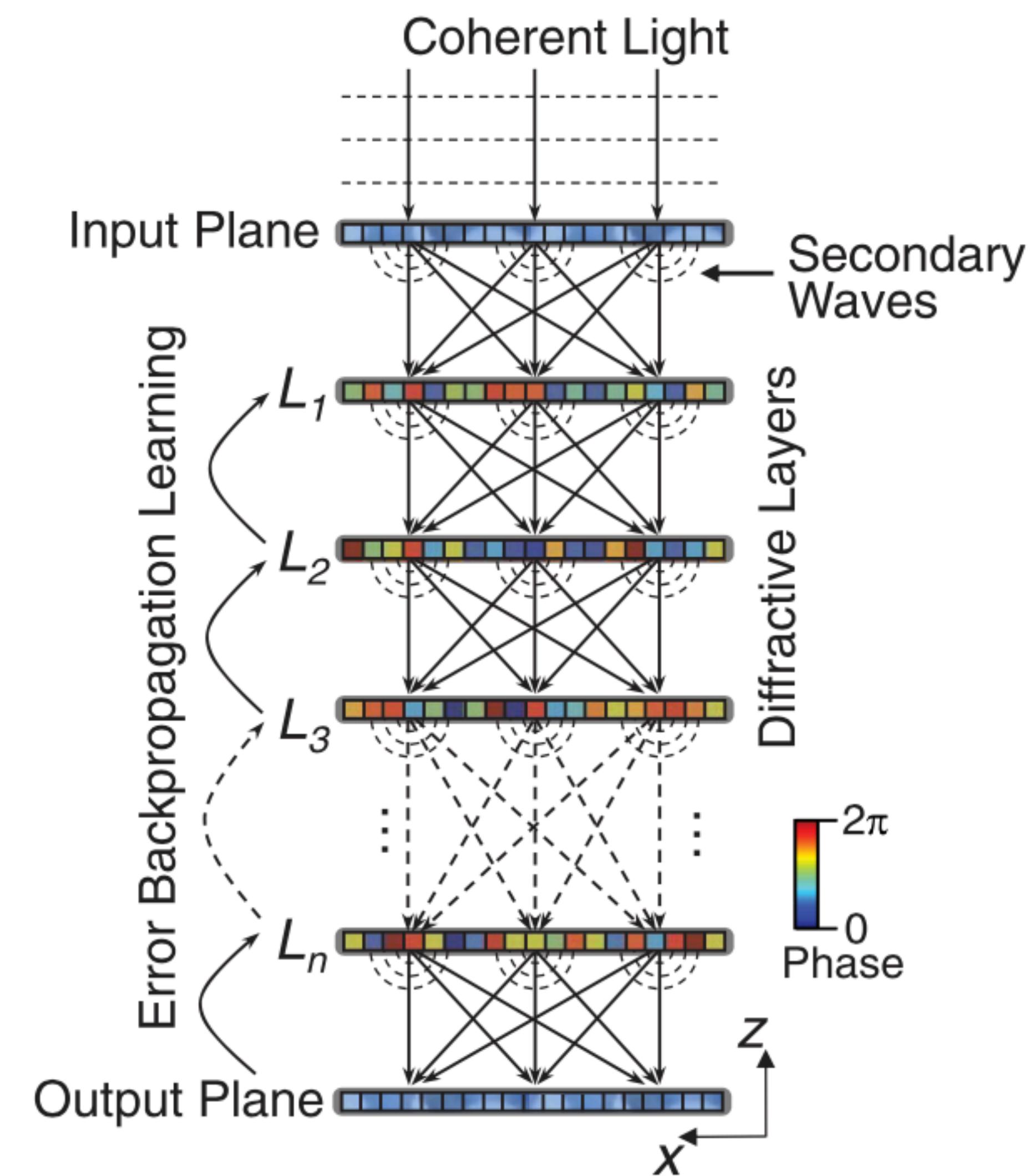
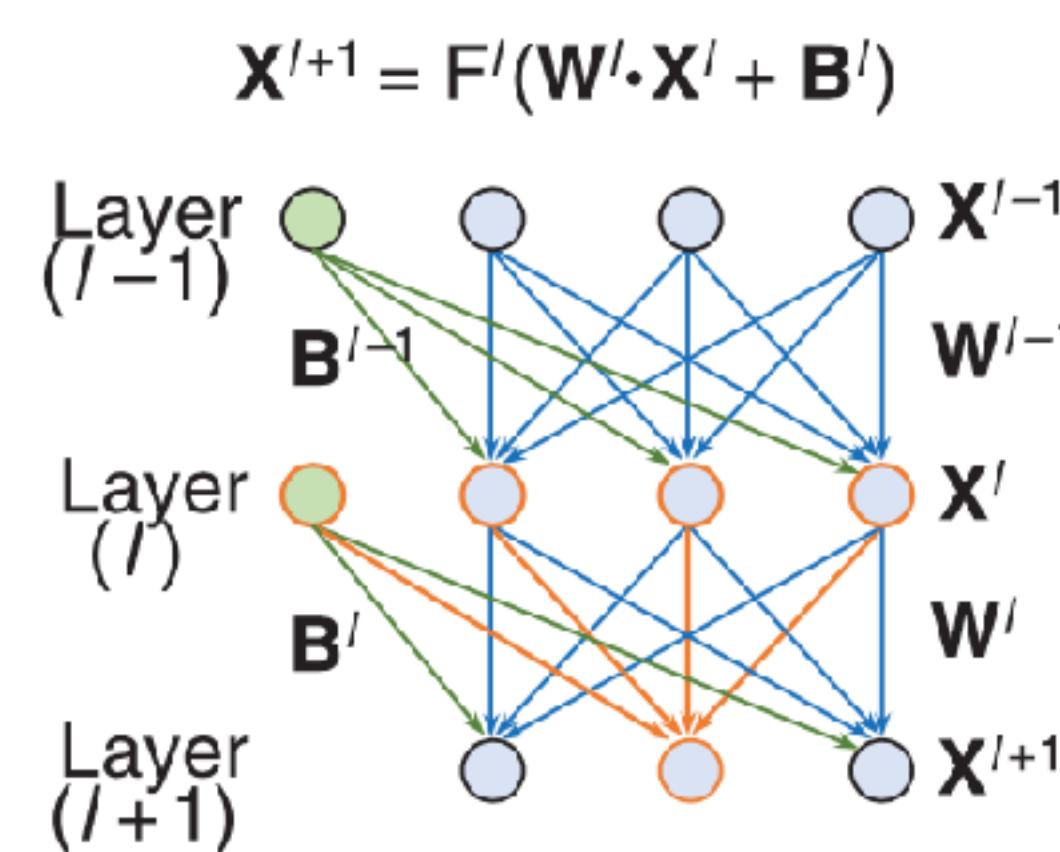
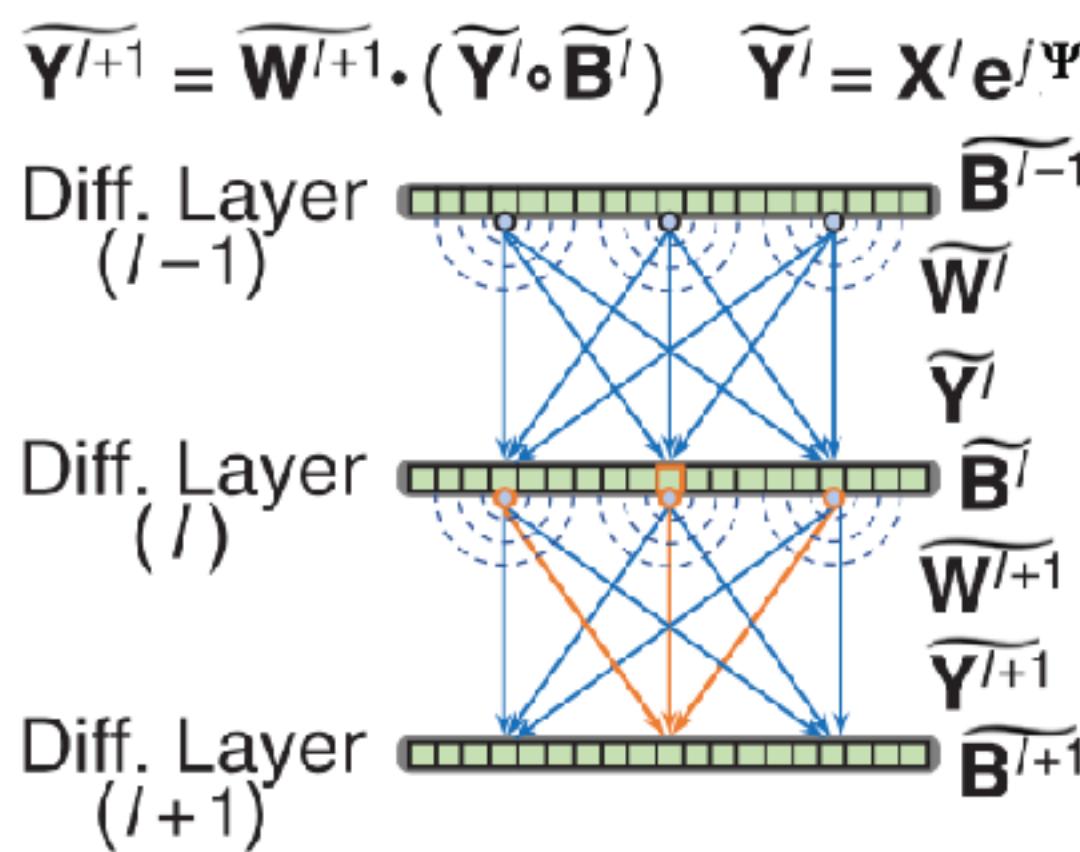


Fig. 2: Diffractive deep neural network (D2NNs)

Schematic

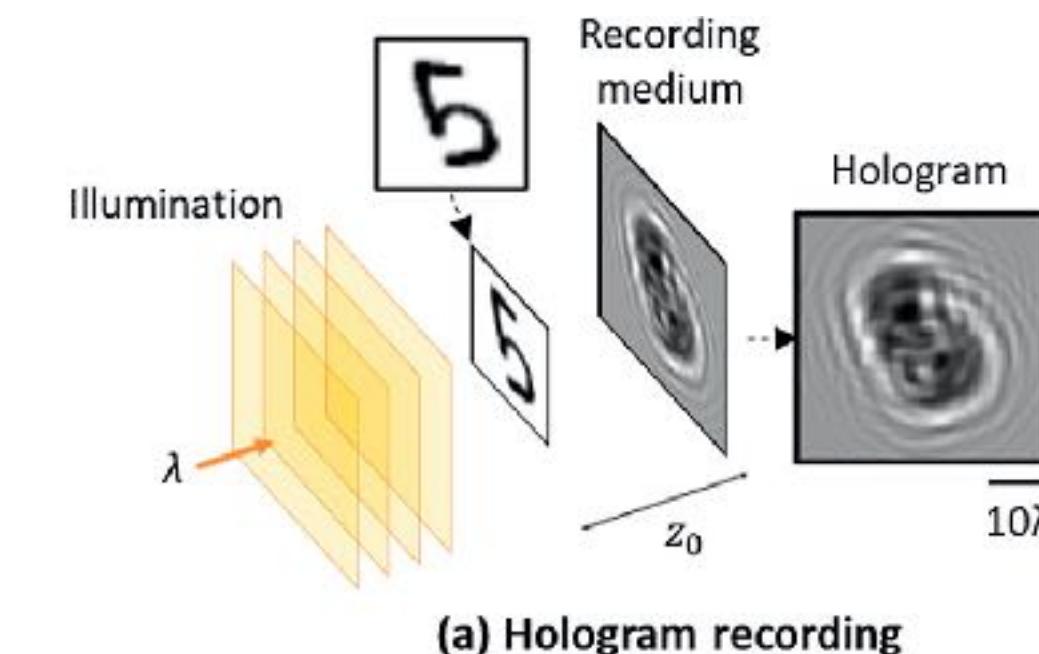
All-optical reconstruction

1. A plane wave as illuminating source.
 2. Diffraction pattern is recorded.
 3. Again illustrate hologram followed by diffractive neural network.
- Rayleigh-Sommerfeld formulation of diffraction

$$u_k^l = \sum_m \frac{z_l}{(r_{mk}^l)^2} \left(\frac{1}{2\pi r_{mk}^l} + \frac{1}{j\lambda} \right) \exp\left(j\frac{2\pi r_{mk}^l}{\lambda}\right) v_m^{l-1}$$

$$r_{mk}^l = \sqrt{(x_k - x_m)^2 + (y_k - y_m)^2 + z_l^2}$$

λ is the wavelength of the optical wave;
 z is the axial distance between layers $l-1$ and l ;
 (x_m, y_m) and (x_k, y_k) are the transverse coordinates of feature m of layer $l-1$ and feature k of layer l ;
 v^0 and u^{L+1} are the optical fields at the input and the output fields-of-view of the diffractive network.



(a) Hologram recording

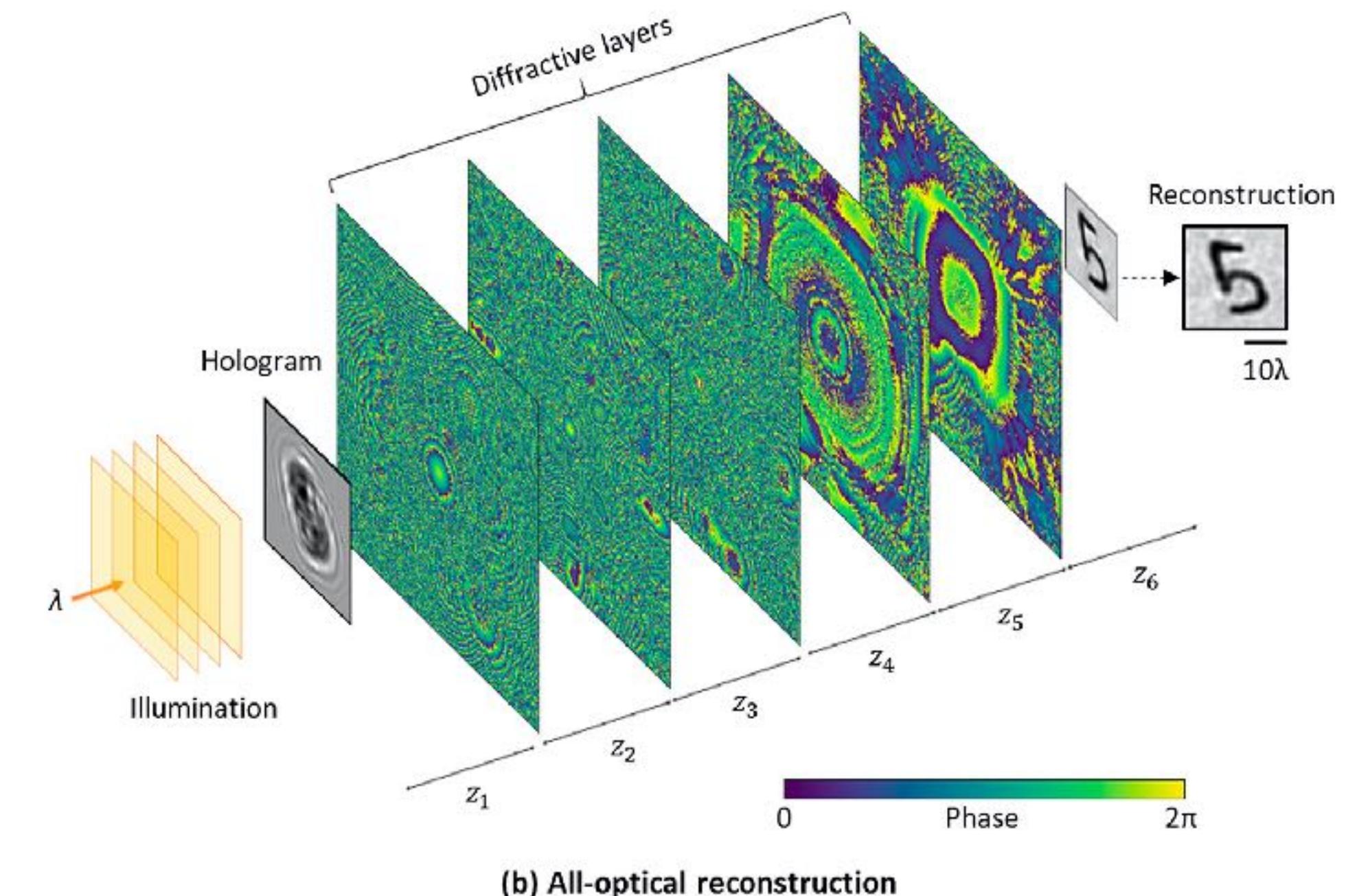


Fig. 3: Hologram recording and reconstruction scheme for the proposed all-optical holographic reconstruction method

Schematic

All-optical reconstruction

- Complex-valued transmittance of pixel k of layer l is denoted with t_k^l

- A thin diffractive layer

- Thus ==> $v_k^l = u_k^l \cdot t_k^l = \{ \sum_m w_{km}^l v_m^{l-1} \} \cdot t_k^l$

- where

$$w_{km}^l = \frac{z_l}{(r_{mk}^l)^2} \left(\frac{1}{2\pi r_{mk}^l} + \frac{1}{j\lambda} \right) \exp \left(j \frac{2\pi r_{mk}^l}{\lambda} \right)$$

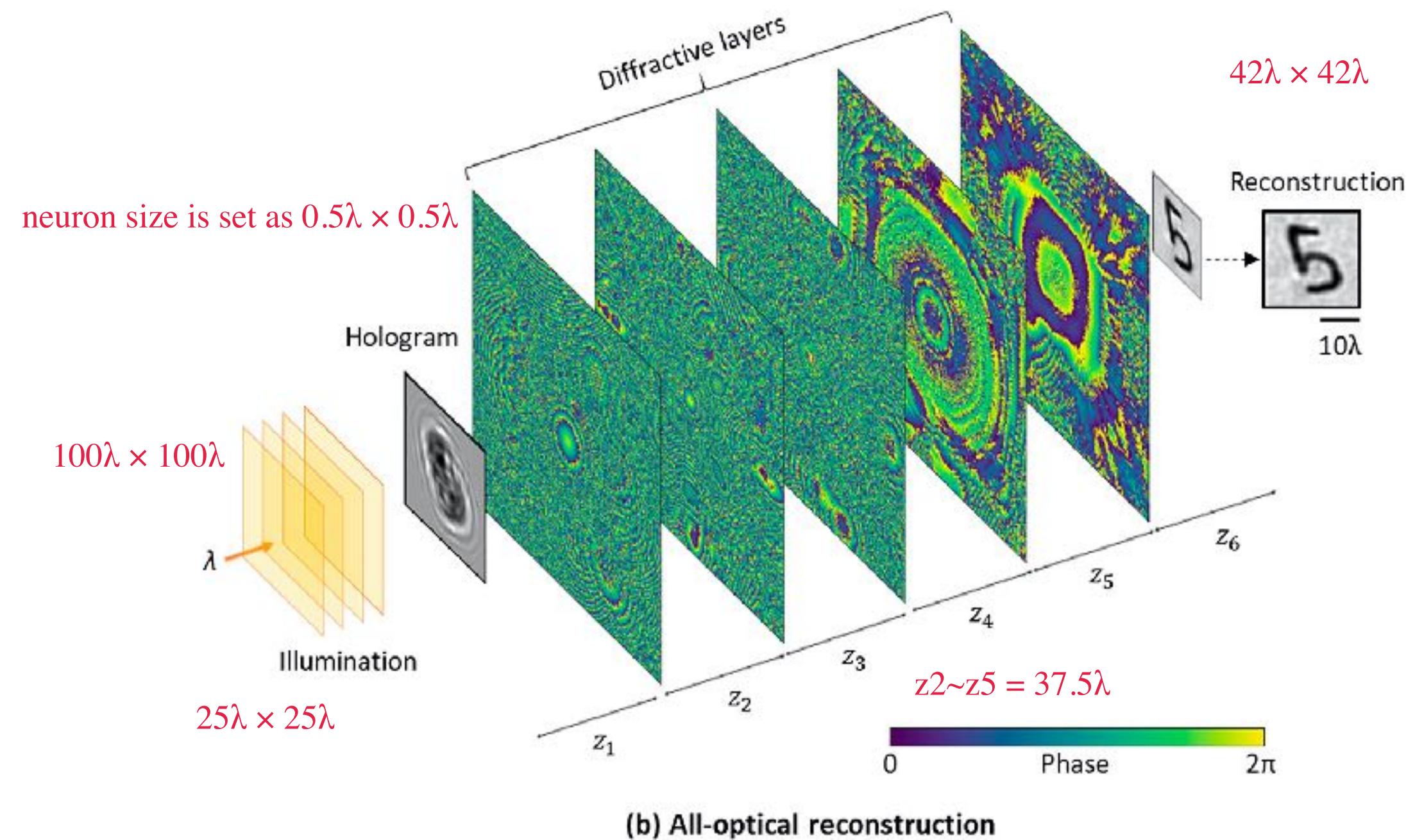
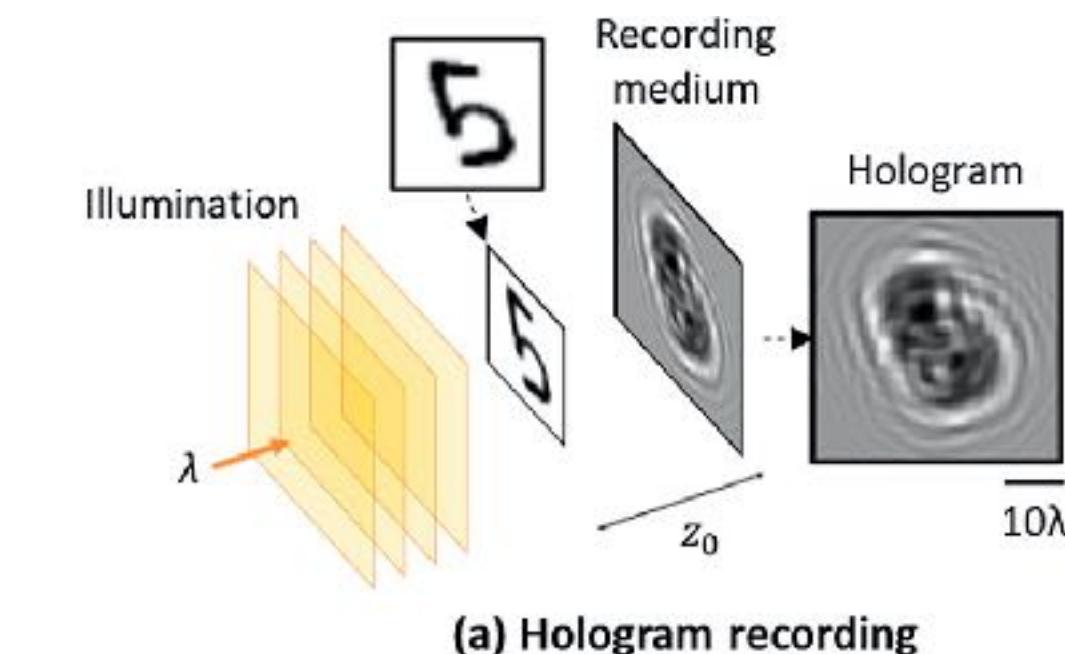


Fig. 3: Hologram recording and reconstruction scheme for the proposed all-optical holographic reconstruction method ($\lambda = 600$ nm)

Results

reconstruction with trained model

- Loss function:

$$L = L_{\text{pixel}} + 1000L_{\text{fourier}} + \eta L_{\text{efficiency}}$$

$$L_{\text{pixel}} = \frac{1}{N_p} \sum_{p=1}^{N_p} |y_p - \hat{y}_p|$$

$$L_{\text{fourier}} = \frac{1}{N_p} \sum_{p=1}^{N_p} |F\{y\}_p - F\{\hat{y}\}_p|^2$$

$$L_{\text{efficiency}} = 1 - \frac{P_{\text{out}}}{P_{\text{illum}}}$$

Where L_{pixel} is MSE; L_{fourier} is Fourier domain's MSE; $L_{\text{efficiency}}$ is efficiency of light energy.

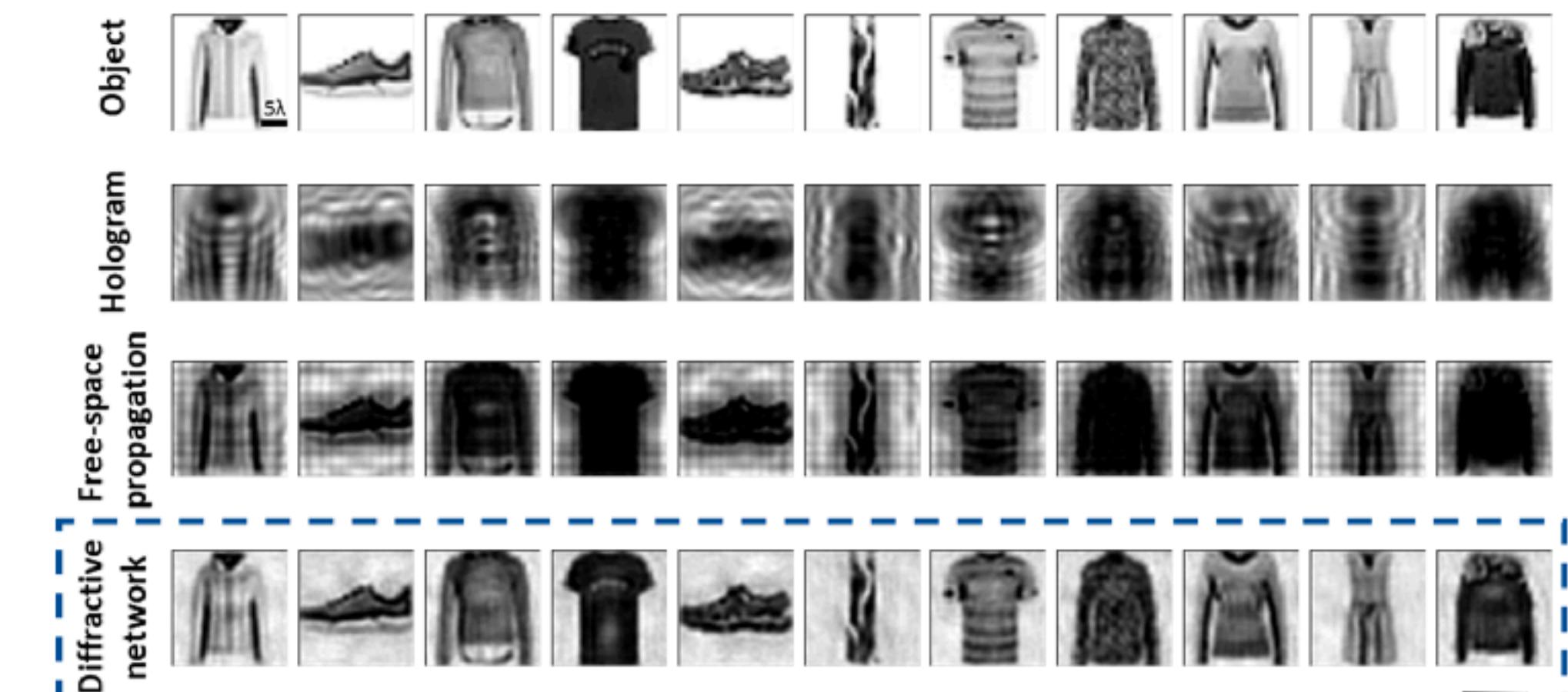
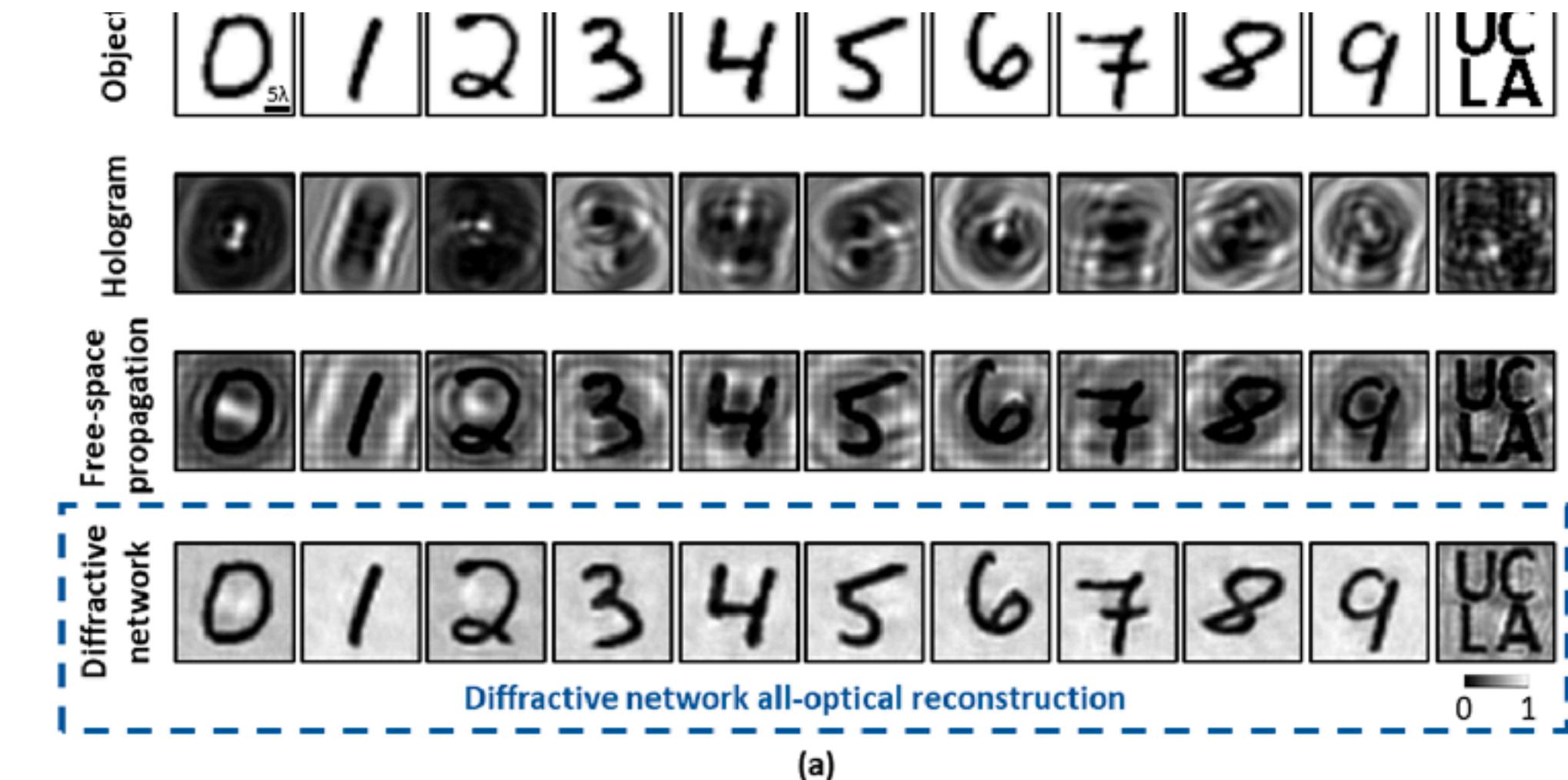


Fig. 4: Performance of all-optical hologram reconstruction using diffractive neural network ($z_6=30\lambda$)

Results

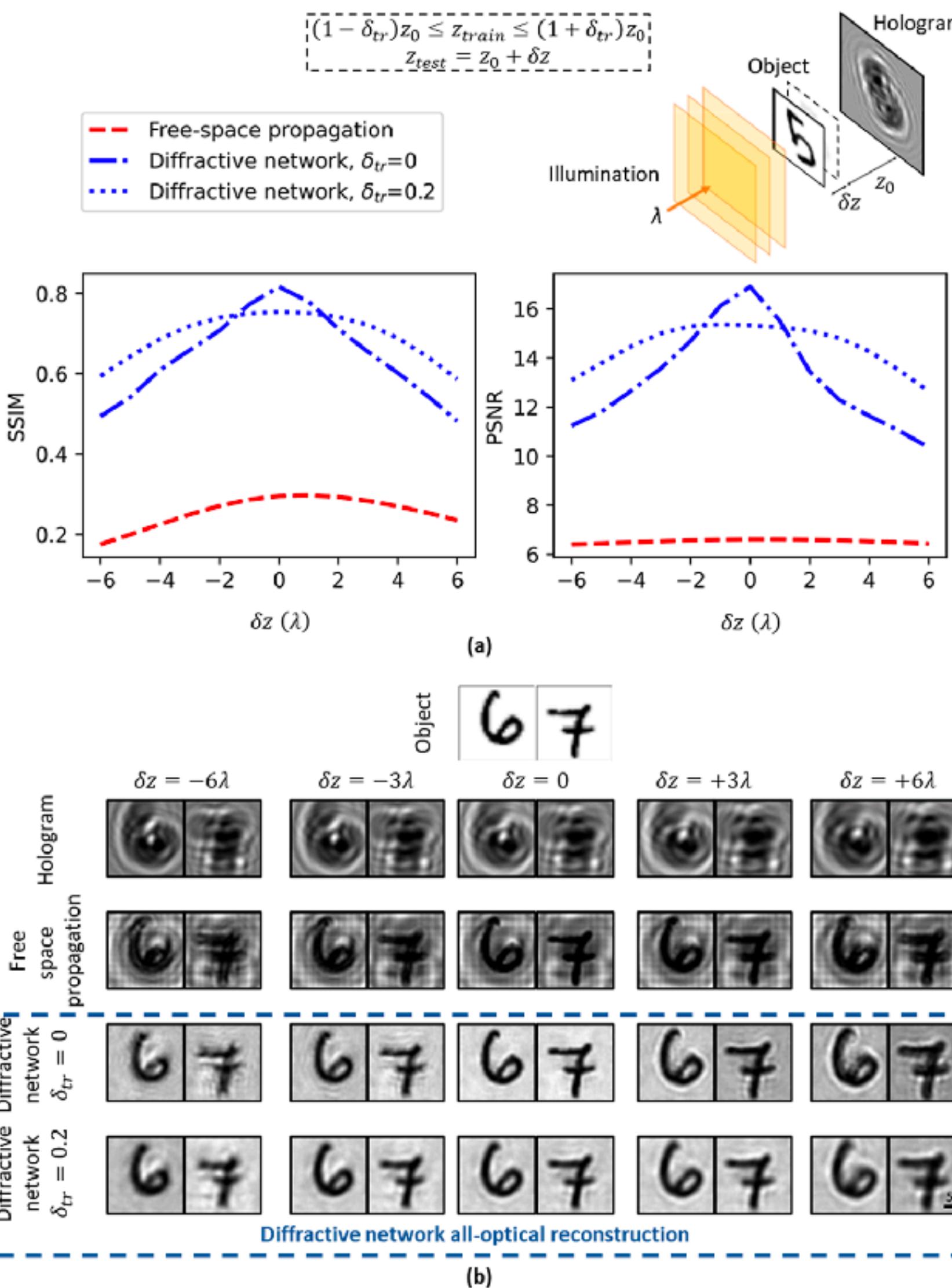


Fig. 5: Robustness of all-optical diffractive network reconstruction against **the hologram recording distance variance**.

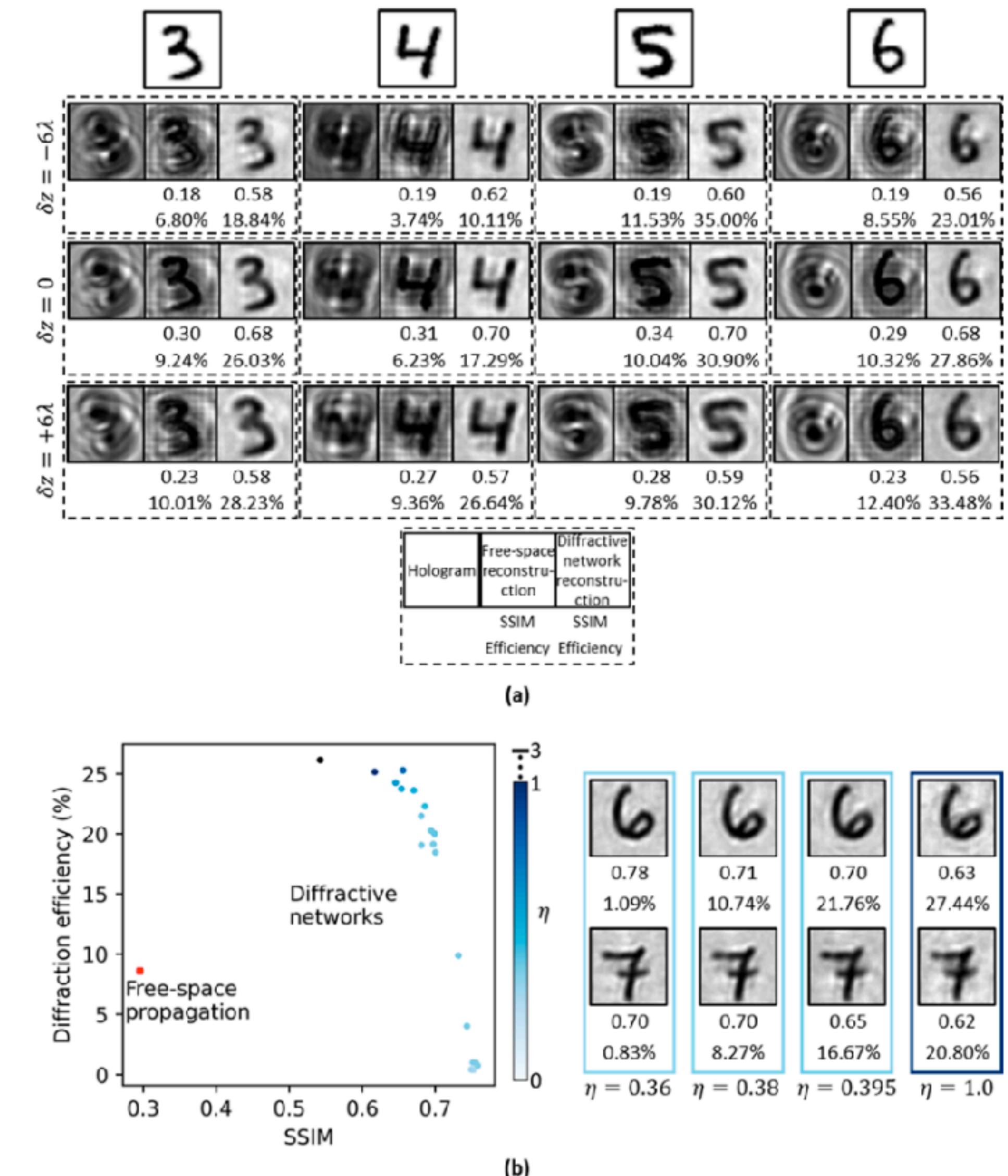


Fig. 6: Diffraction efficiency improvement of all-optical holographic reconstruction performed by diffractive network.

Conclusion

- Phase recovery & twin-image elimination based on all-optical processor (High speed).
- Diffractive layers can perform as neural (physics based connection layer).
- Method that vaccine the model the changing of hologram recording distance.
- Rayleigh–Sommerfeld diffraction integral was computed using the Angular Spectrum method
- Tensorflow + Adam + mixed Loss GTX1080Ti + Win 10

Outline

- Analogy of ONN
- **ONN implementation**
- A glimpse
- A optics-inspired design

Schematic

The optical path of ONN

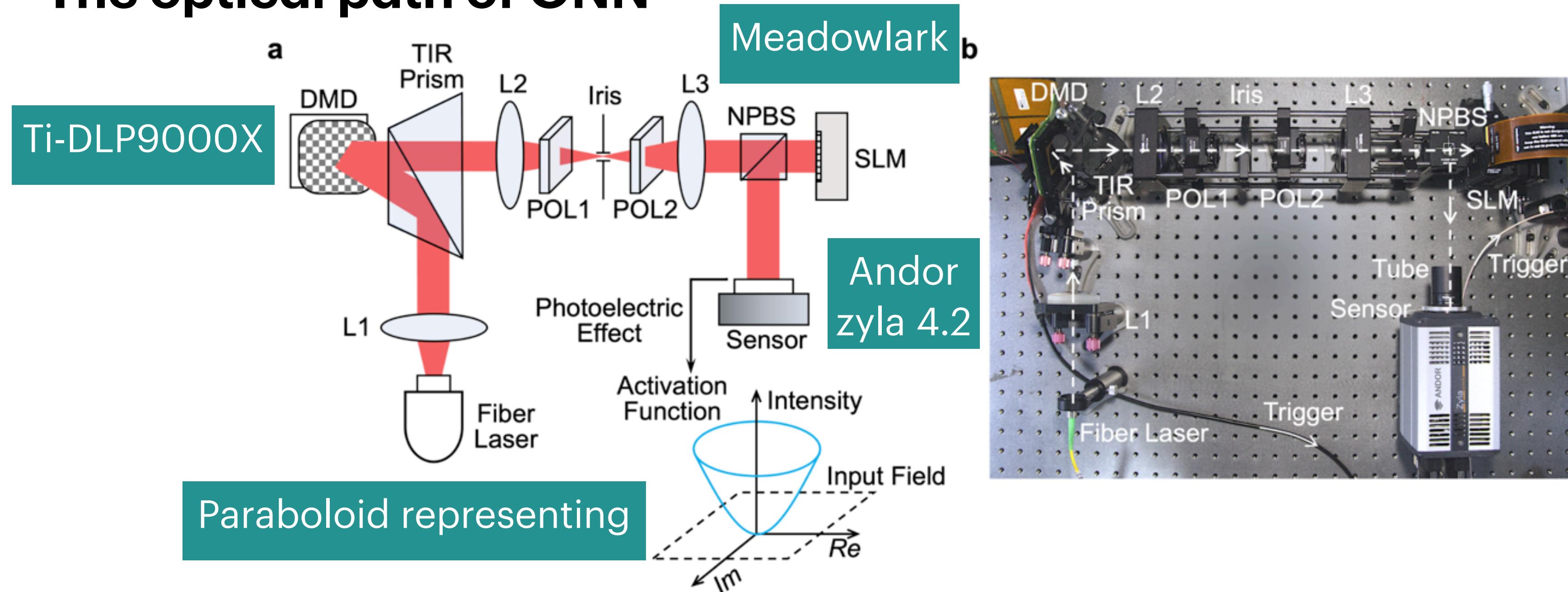


Fig. 2-1 Schematic and photo of the experimental DPU prototype

- **NPBS:** non-polarized beam splitter.
- Laser: 698nm; DMD: 149989Hz(2560 × 1600 , 1-Bit) + 7.6 μm
- SLM: 422.4Hz (1,920 × 1,152, 8-Bits) + 7.6 μm + efficiency 88% + update time of ~2.4 ms
- sCMOS: 100 f.p.s. (2048 × 2048, 16-Bits) + 6.5 μm

Dataflow

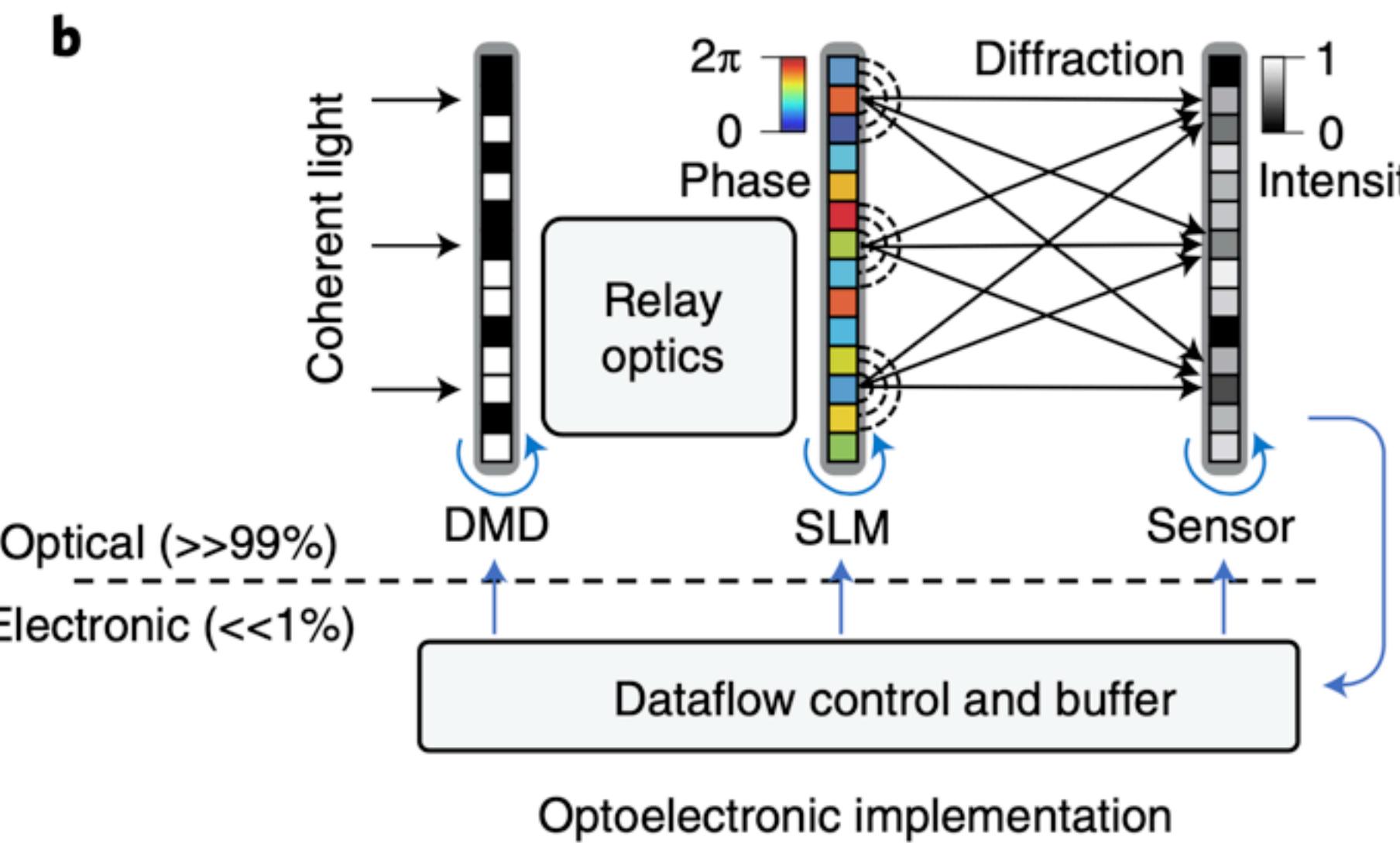
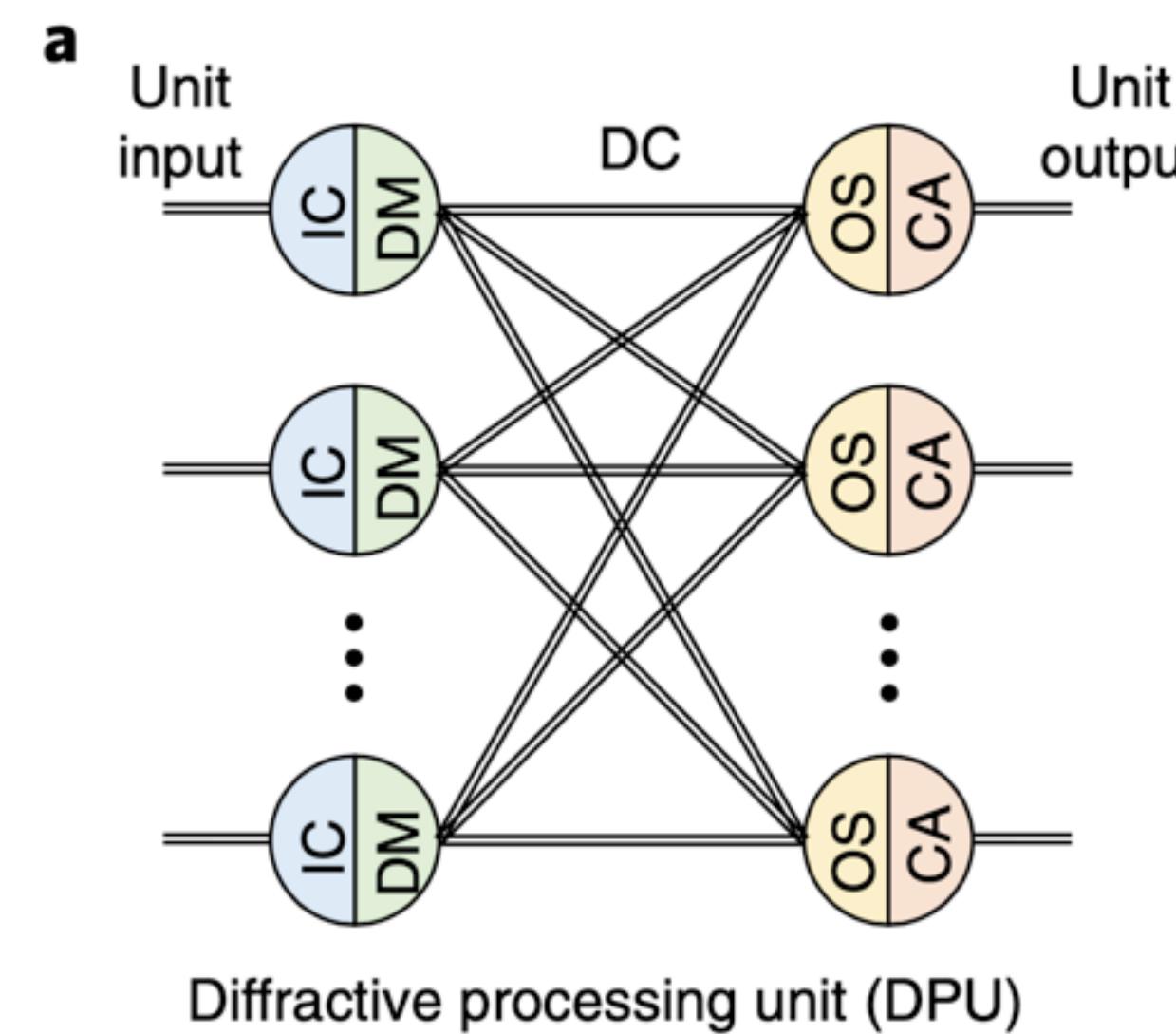
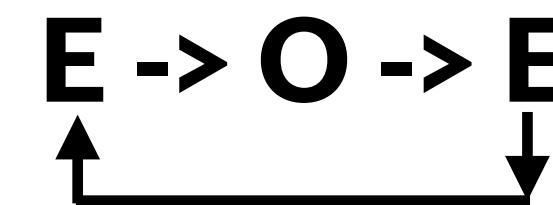


Fig. 2-2 Dataflow of perceptron-like optoelectronic computing building block

- IC: information coding; DM: Diffractive modulation; DC: Diffractive connection
- OS: optical field summation; CA: complex activation;

Network structure

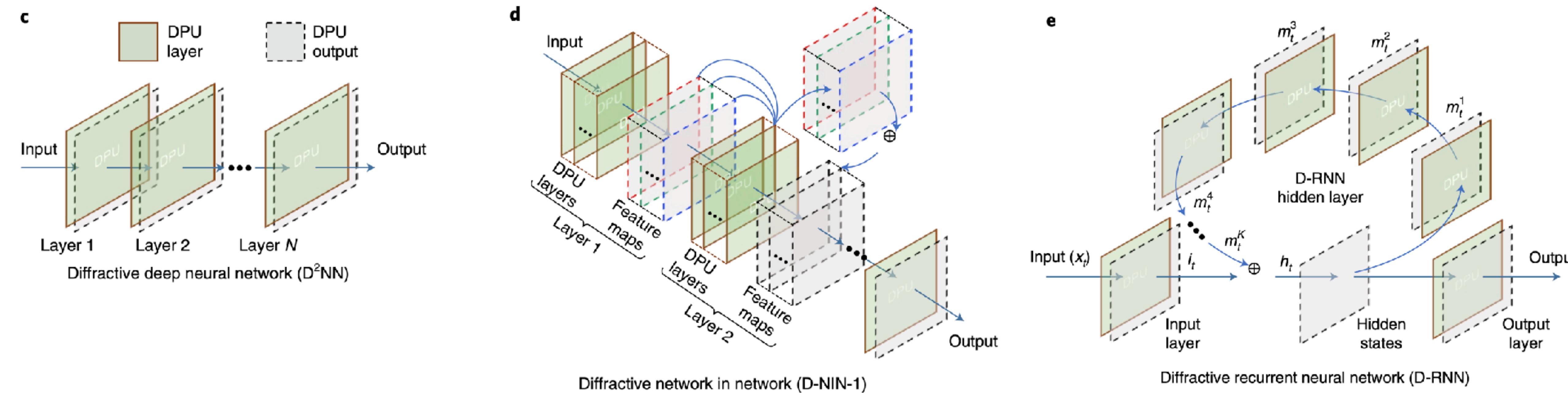


Fig. 2-3 Different DNN, that are D2NN, DNIN, D-RNN

- Single layer, multi-layers, and recurrent layer

$$\text{out}(N_i, C_{\text{out}_j}) = \text{bias}(C_{\text{out}_j}) + \sum_{k=0}^{C_{\text{in}}-1} \text{weight}(C_{\text{out}_j}, k) * \text{input}(N_i, k)$$

- , where (N,Cin,H,W) is the output value of the layer with input size and (N,Cout,Hout,Wout) refers to its output.

Adaptive training inter-layer fine-tuning

- Each layer contains transfer errors distort the wavefront connections of neurons.
- in situ v.s. in solico
- Robustness for few-shot

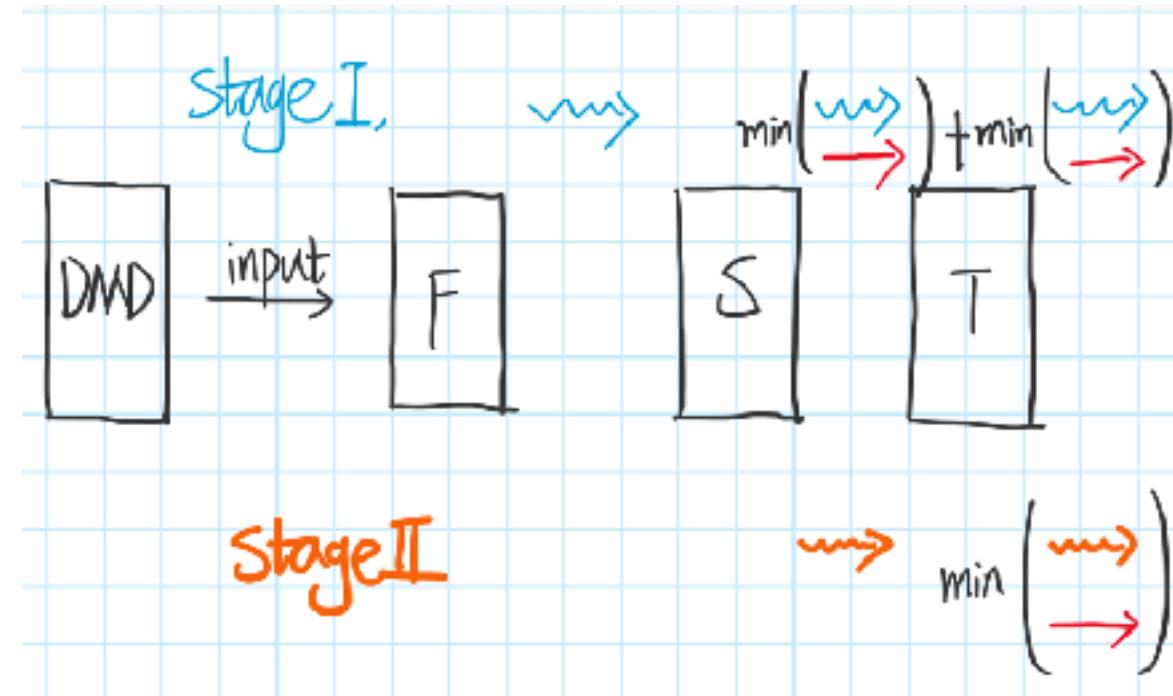


Fig. 2-5 Illustration of adaptive train & results of phase diff

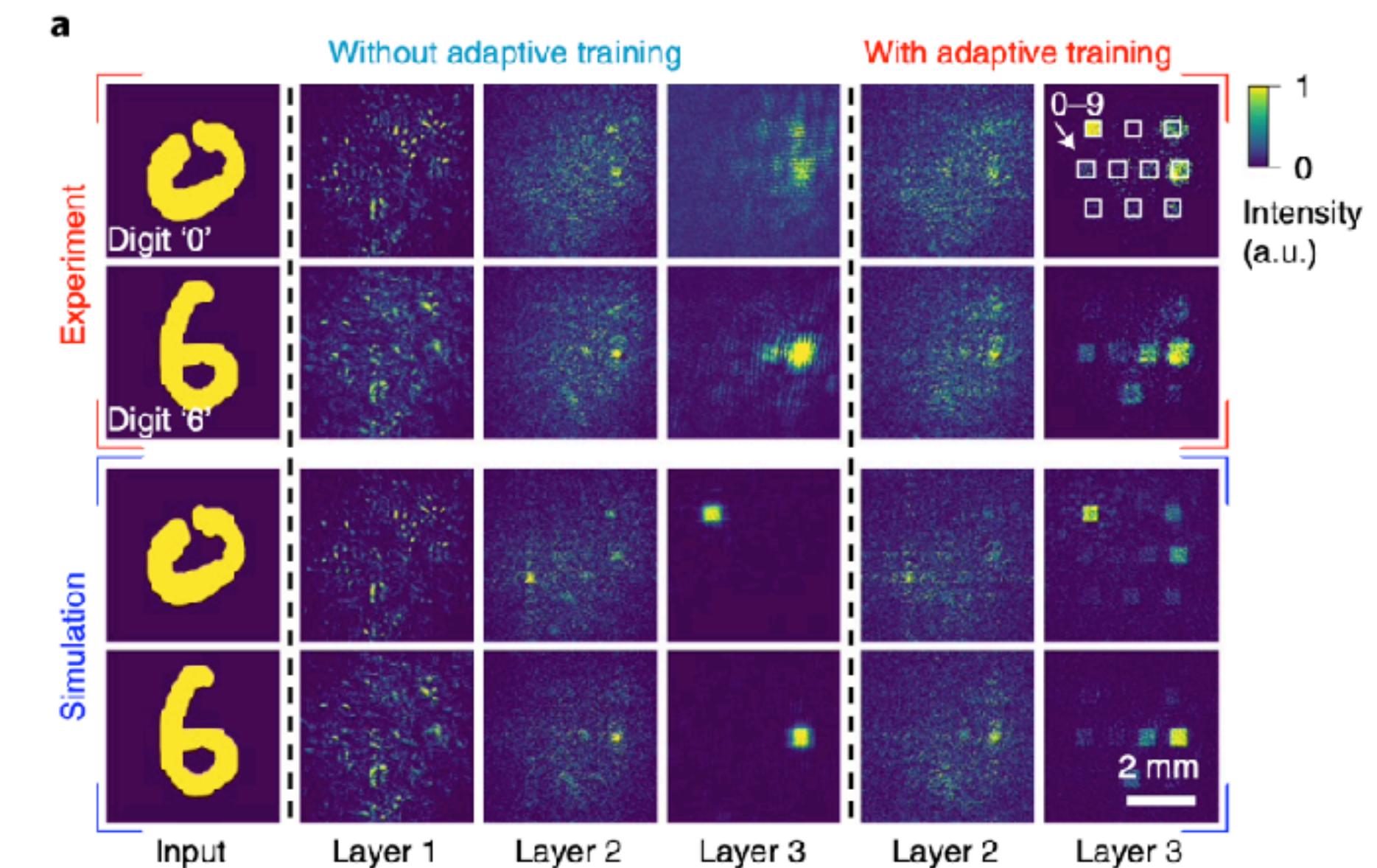


Fig. 2-4 Experimental DPU outputs of a 3-layers OE D2NN

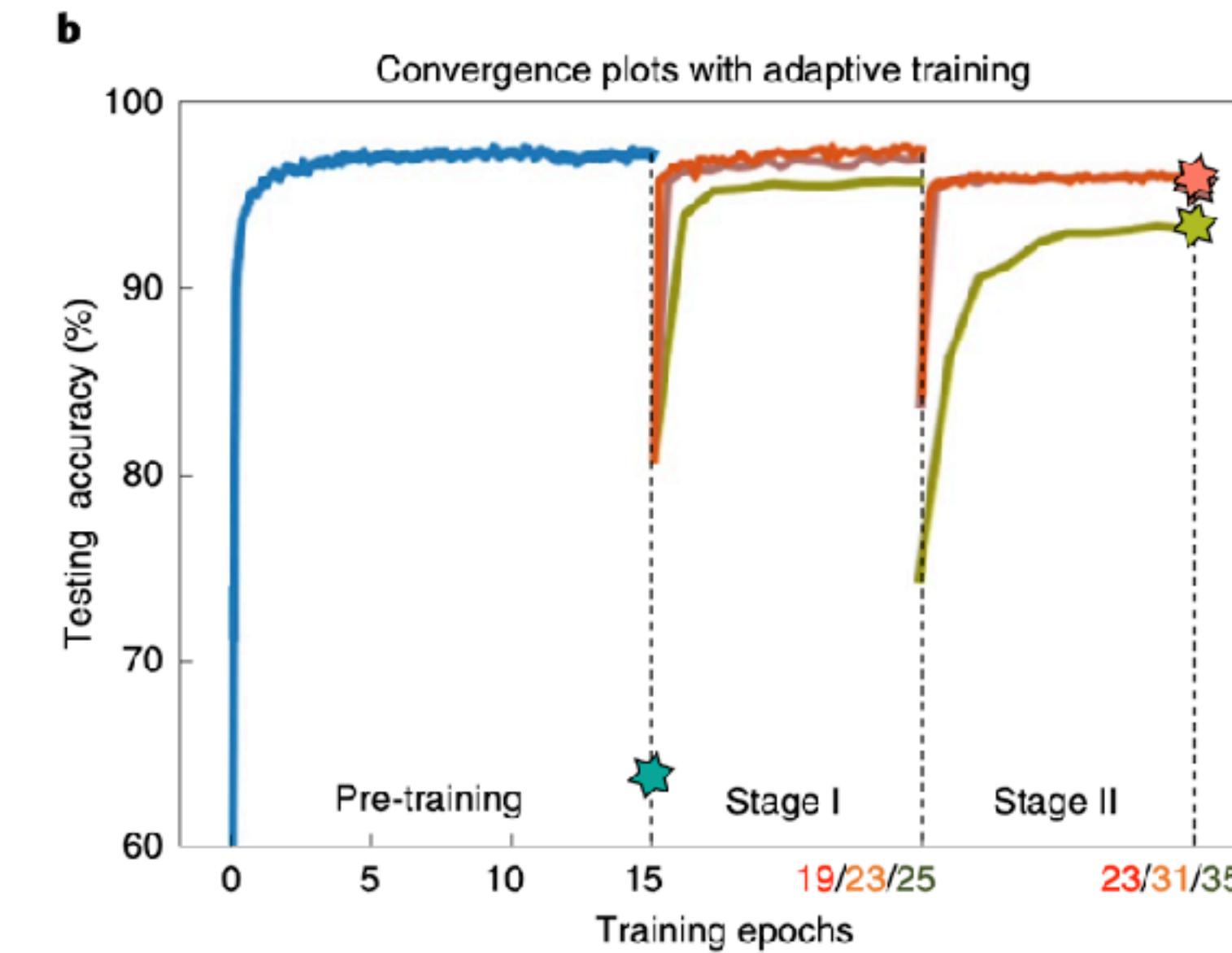
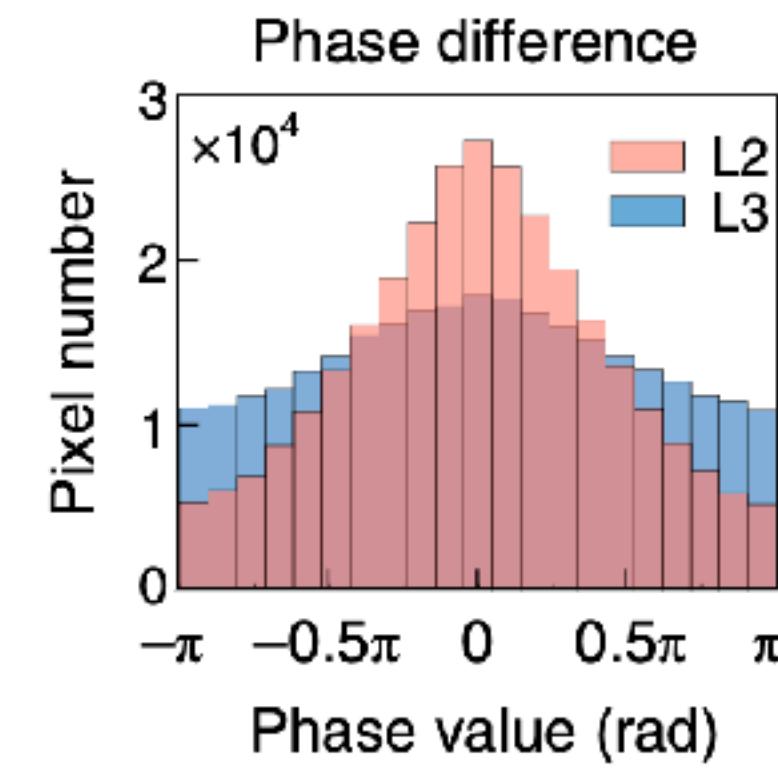


Fig. 2-6 Experimental DPU outputs of a 3-layers OE D2NN



Results

MINIST classification

- Simulating the multi-feature map of CNN

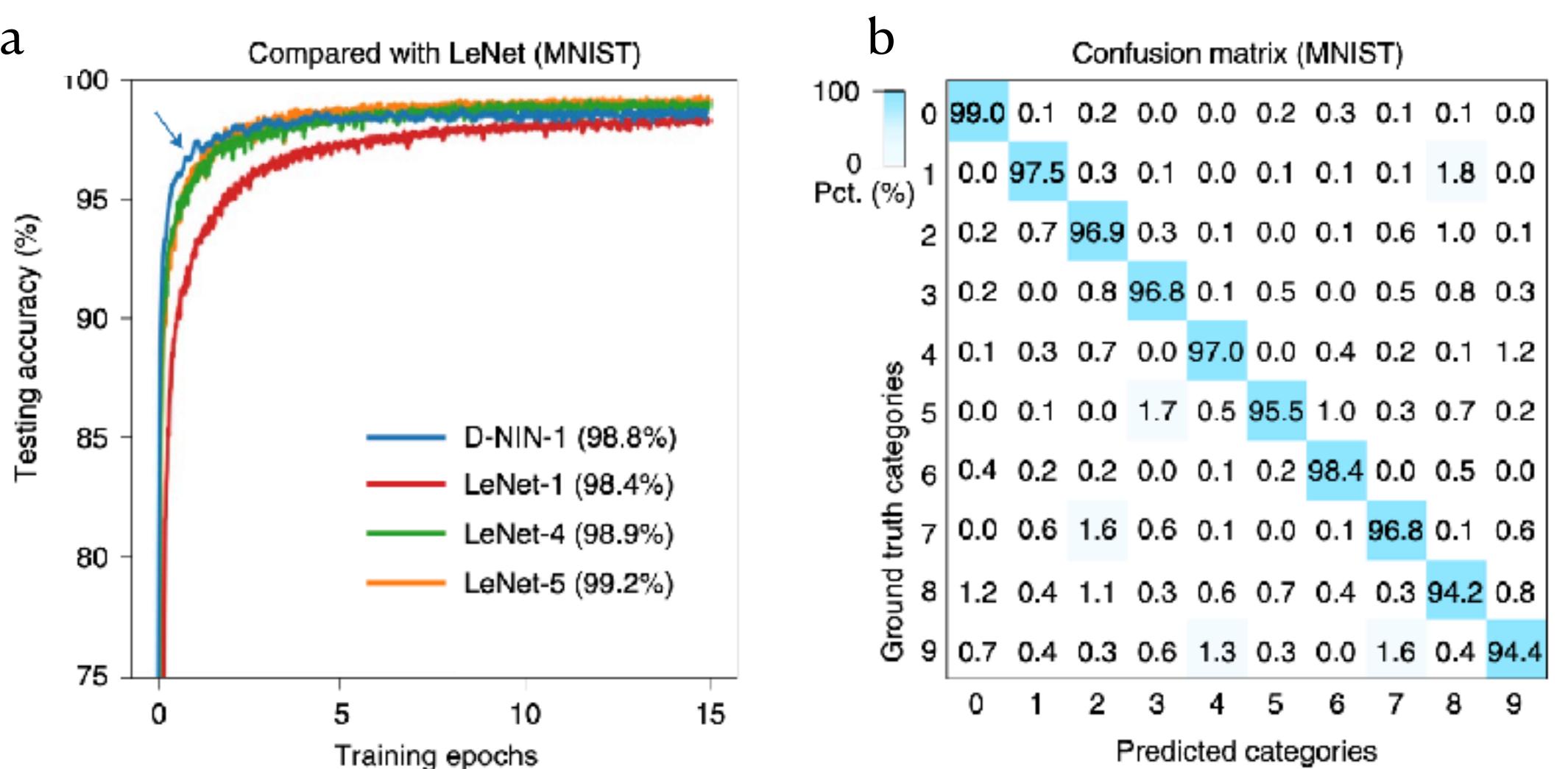


Fig. 2-6 Test results of D-NIN-1

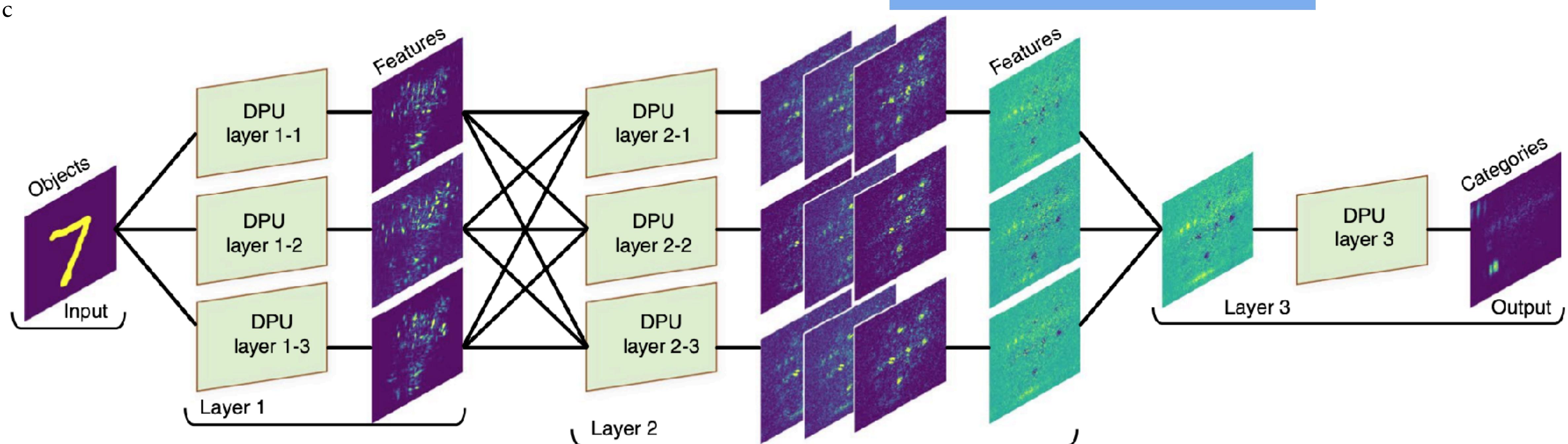


Fig. 2-7 Tensor flow detail of training with D-NIN-1

Results

optical version RNN

- Weizmann: 180*140 ppi; deinterlaced 50fps; run, walk, skip, shortly jack, jump-forward-on-two-legs, pjump, gallopsideways, wave2, wave1, or bend; 340MB.
- KTH: 160*120 ppi2; walking, jogging, running, boxing, handwaving, handclapping; 766.37MB

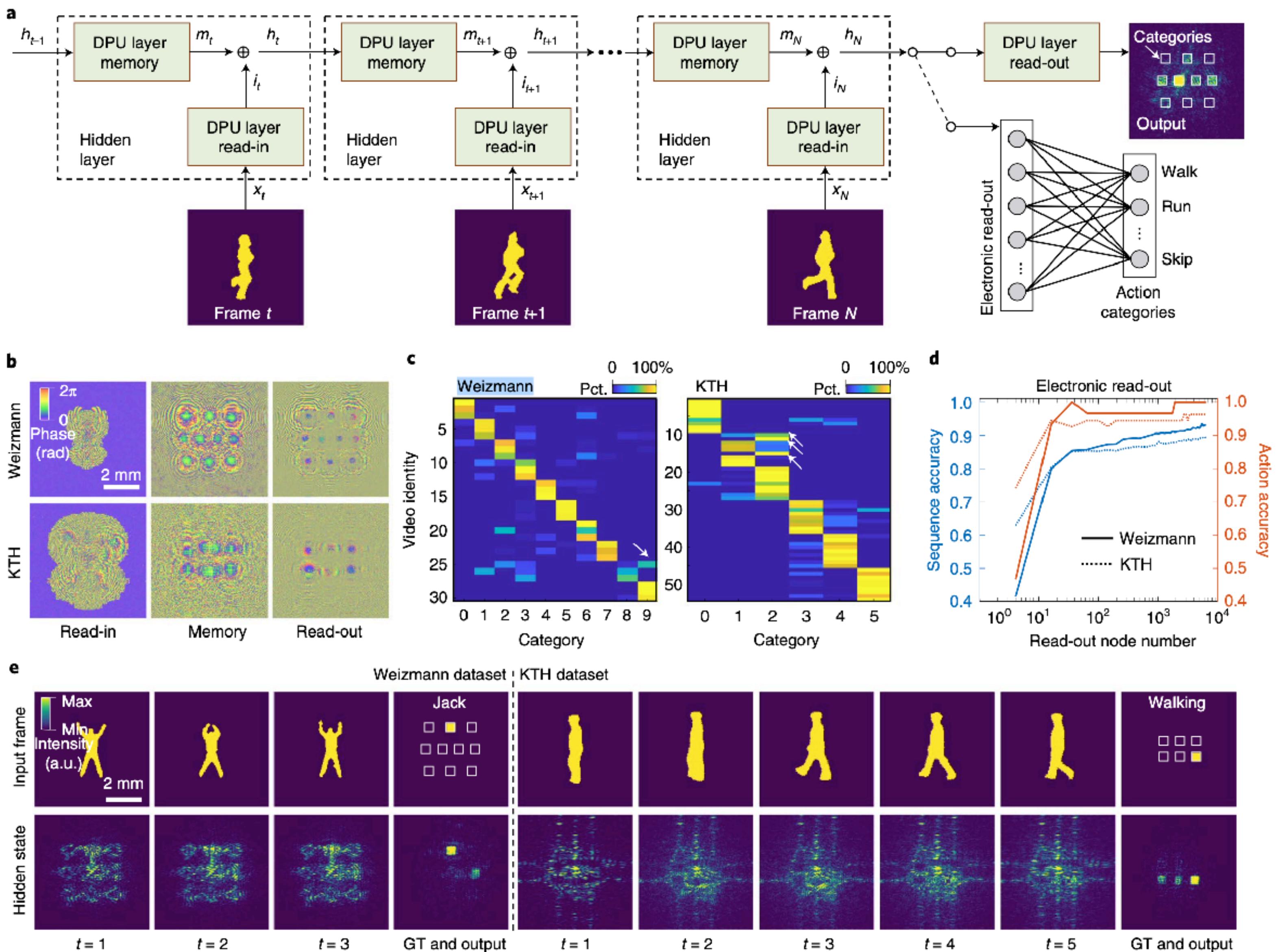


Fig. 2-8 D-RNN for human actions recognition

Conclusion

- Still a **data-driven** training.
- Speed of training is determined by **speed of DMD and sCMOS and SLM**
- Data preprocessing is important (interpolation & featuring)
- Good inducing may come from **multi-task learning** & **widely varying dataset**.
- Classification is benefited from **average pooling** of the network output.
- softmax cross-entropy := softmax loss = $E = - \sum_{j=1}^T y_i \log\left(\frac{e^{a_i}}{\sum_{k=1}^T e^{a_k}}\right)$, where y_i is ith class in array, e^{a_i} is value of **exp** (neuron output).

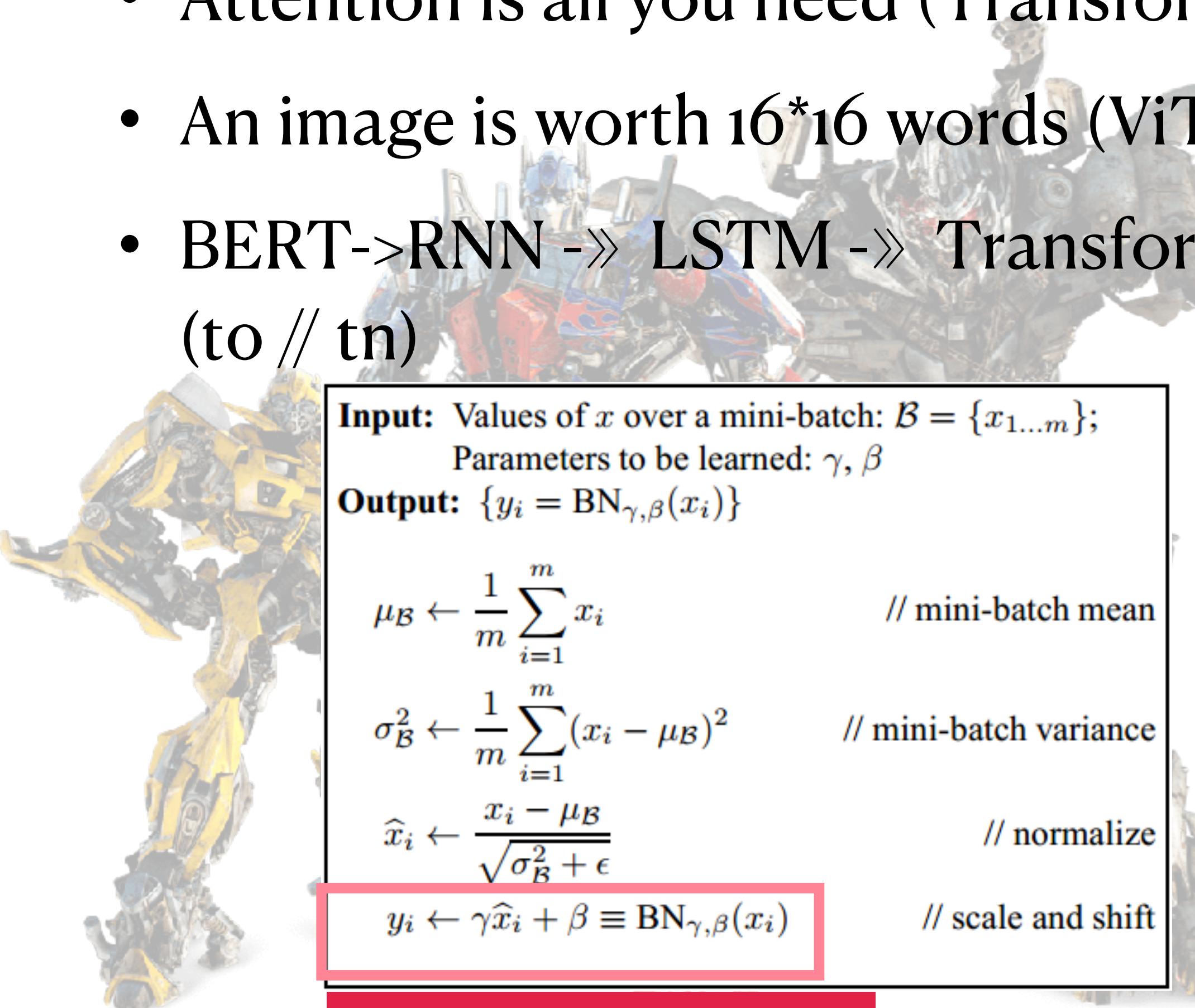
Outline

- Analogy of ONN
- ONN implementation
- A glimpse
- A optics-inspired design

Prior knowledge

What is Vision Transformer (ViT)?

- Attention is all you need (Transformer)
- An image is worth 16×16 words (ViT)
- BERT->RNN -> LSTM -> Transformer
(to // tn)



where x as input & y is output

An image is worth 16×16 Words: transformer for image recognition at scale

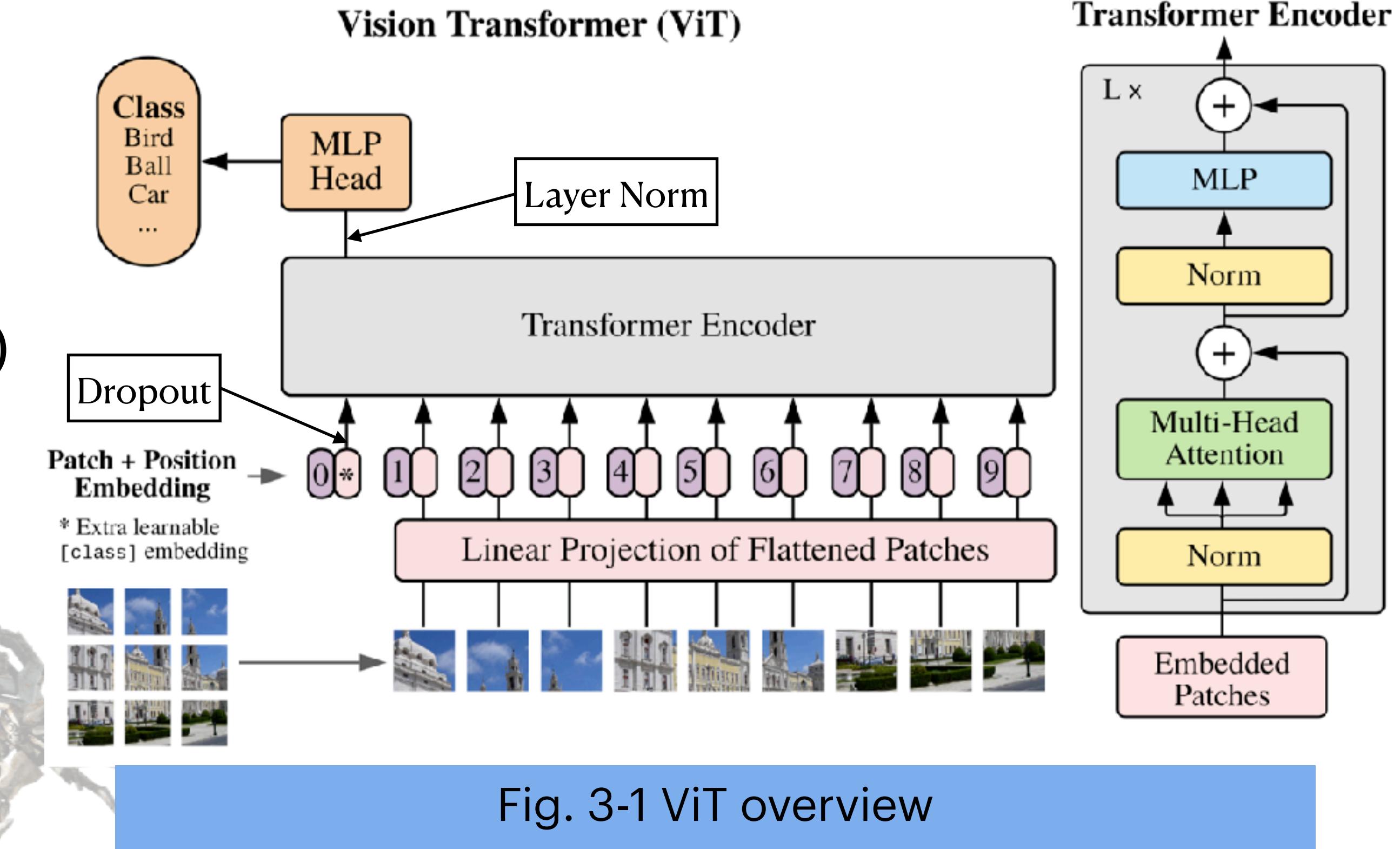


Fig. 3-1 ViT overview

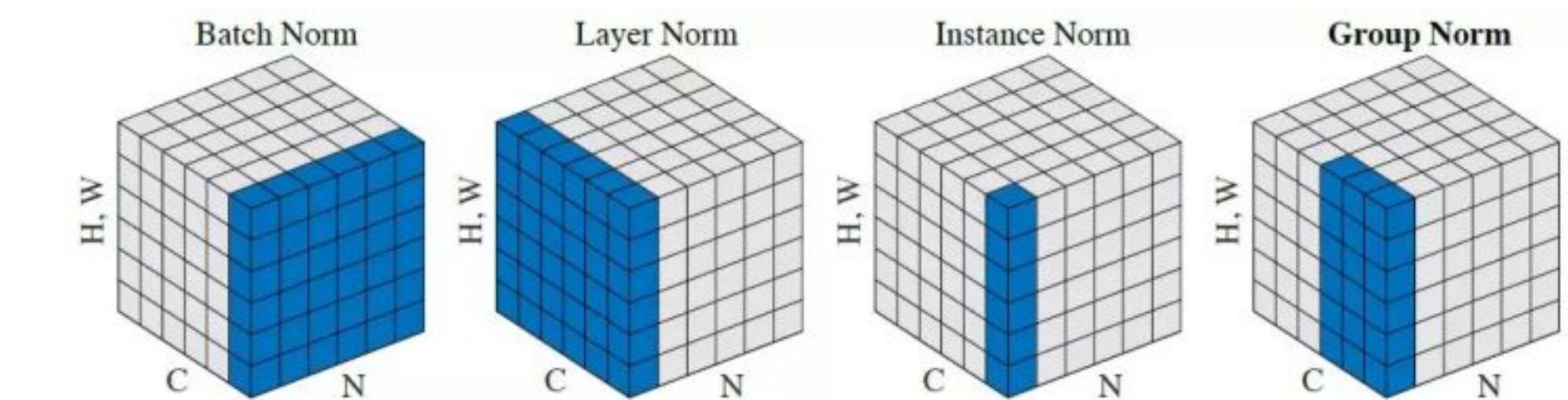


Fig. 3-2 Different norm methods

Prior knowledge

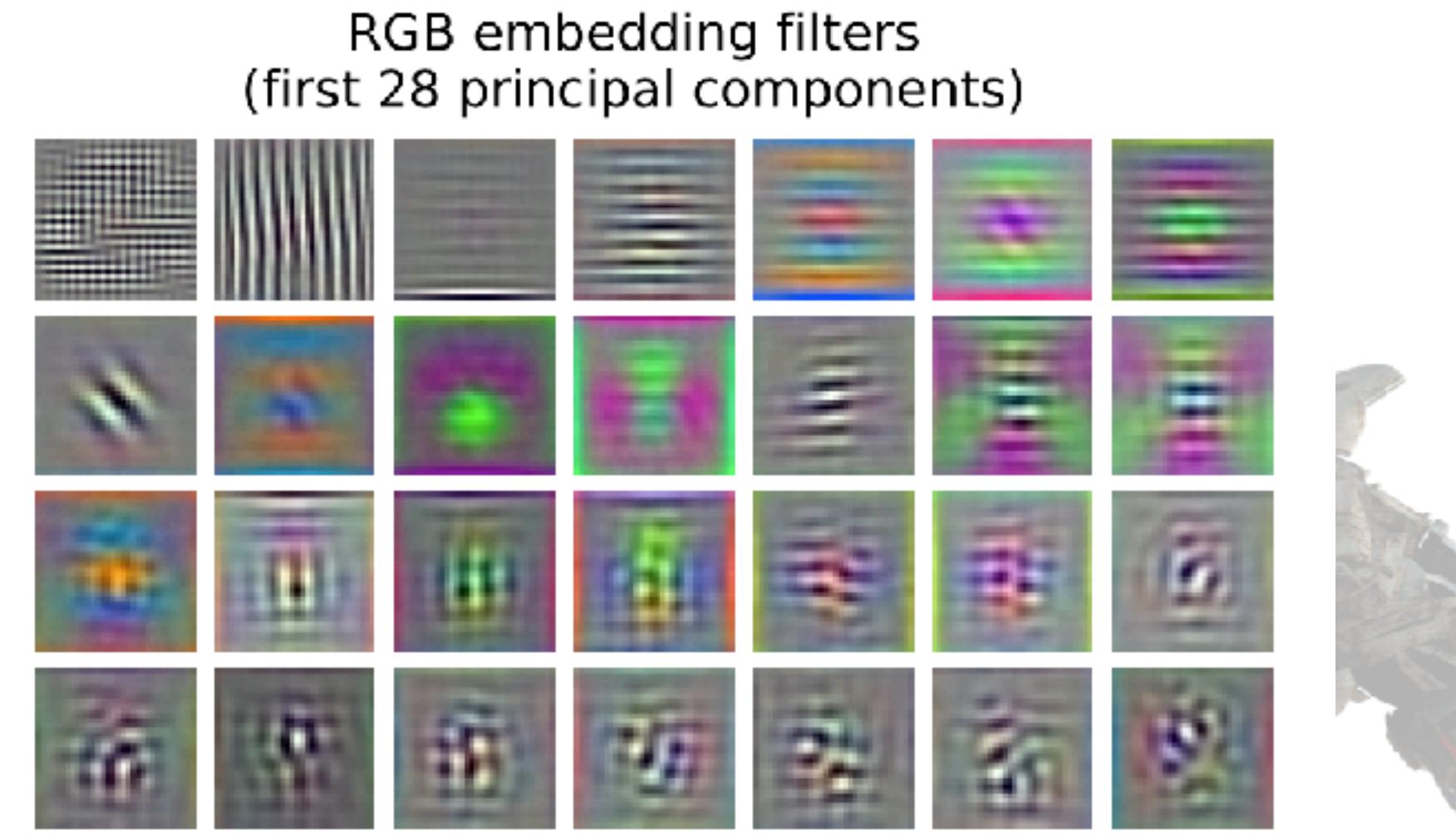


Fig. 3-3 Filters of the init linear embedding of RGB values of ViT-L/32

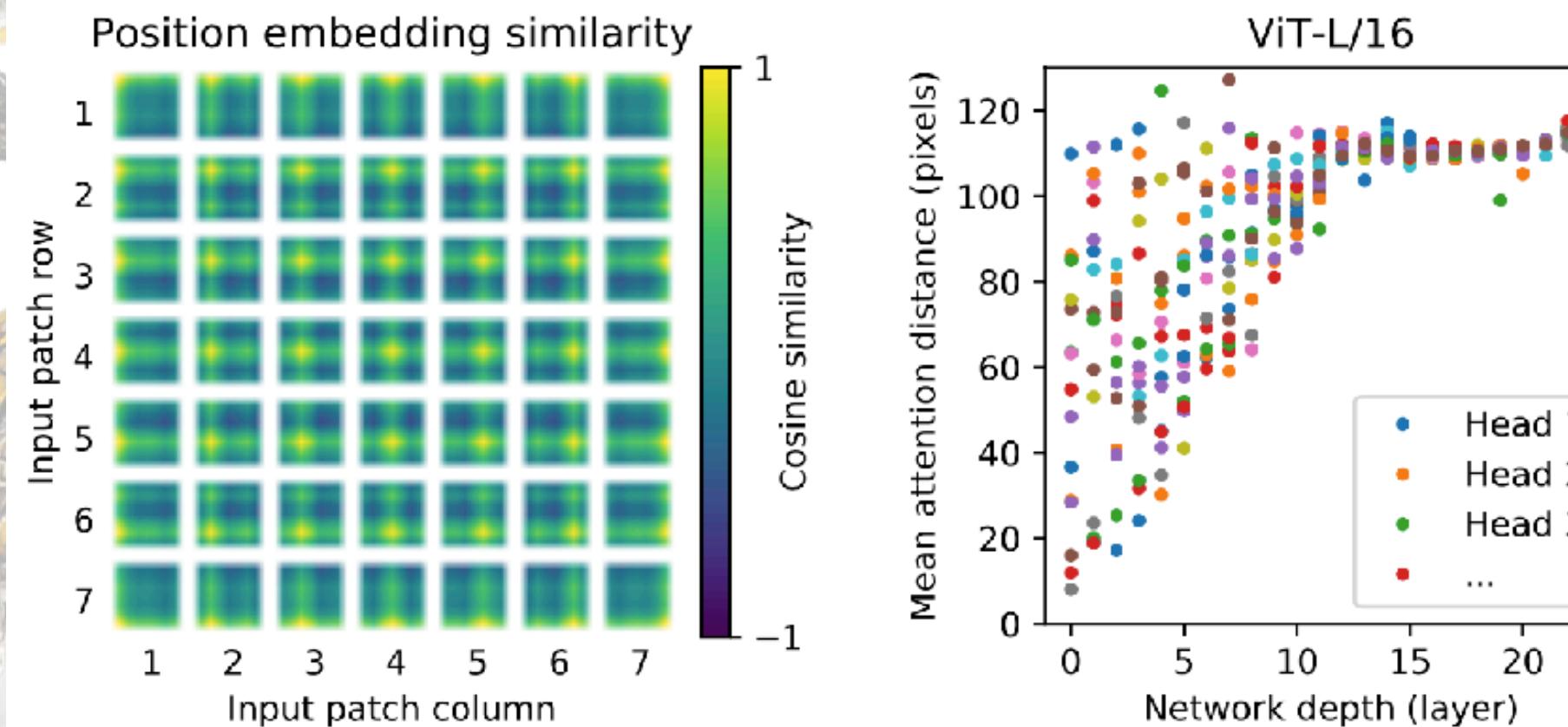


Fig. 3-4 Similarity of position embeddings of ViT-L/32

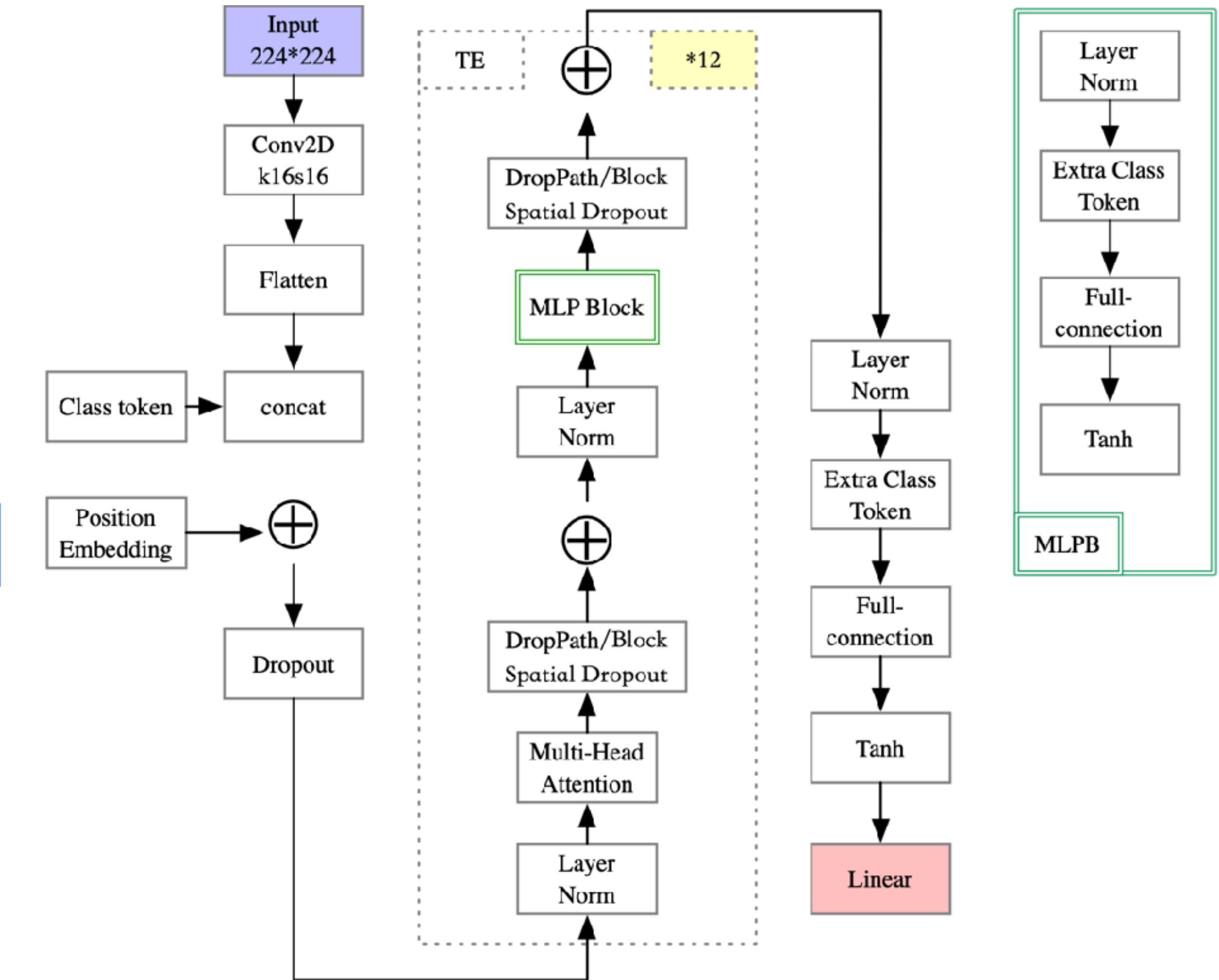


Fig. 3-5 The architecture of ViT

Masked Autoencoder

Architecture

- **Image** → masked image → flatten → encoding → embedded image + mask token → decoding → **Image**



Fig. 3-6 Architecture of MAE

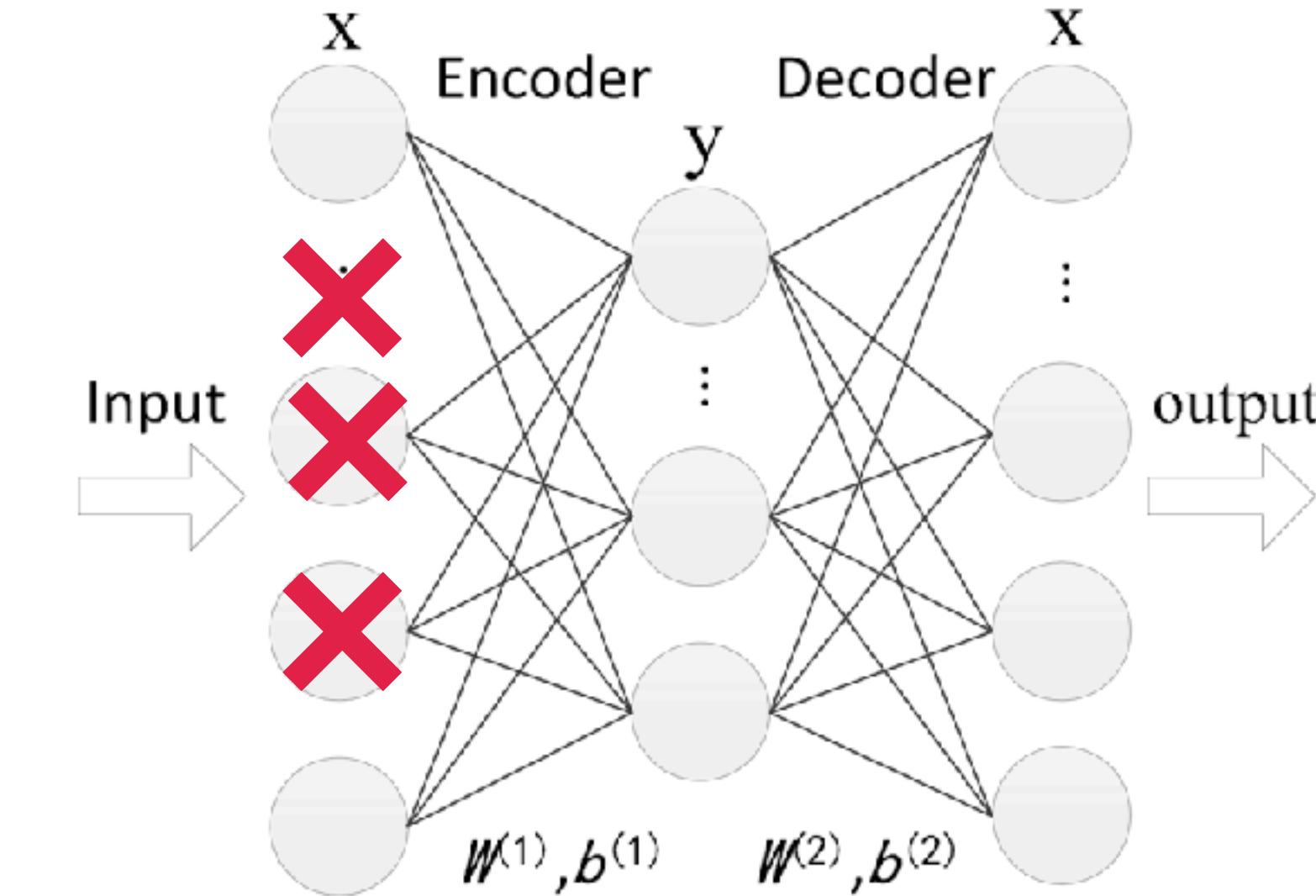


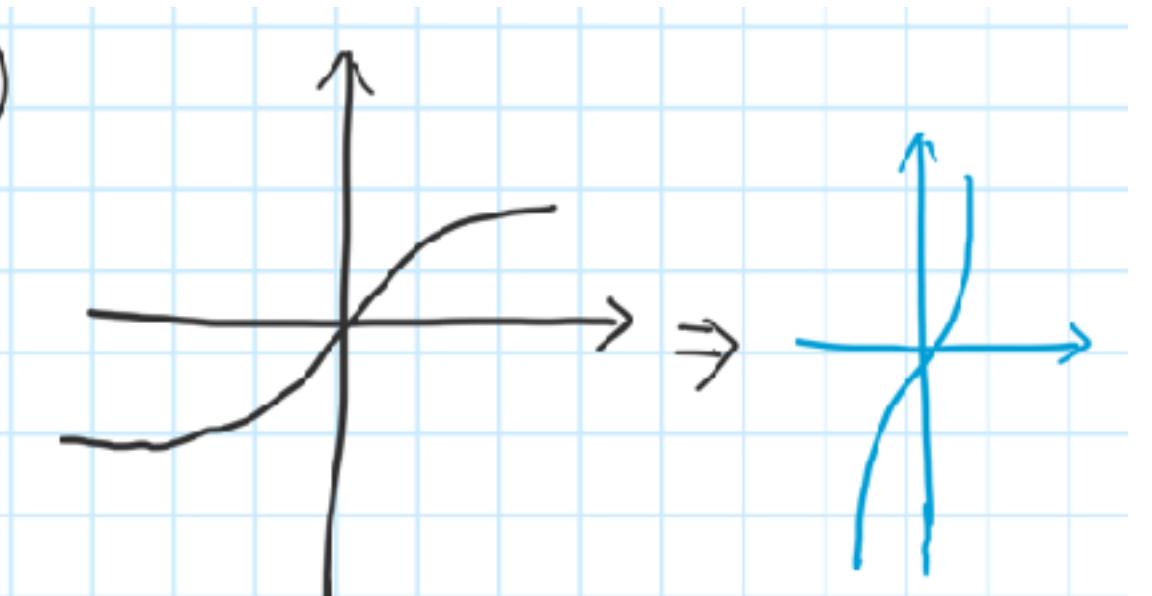
Fig. 3-7 Autoencoder (auto-associator) & Denoise autoencoder

$$y = f(x) = \tanh(w_1 x + b_1)$$

$$x = (\tanh^{-1} y - b_1) \cdot \frac{1}{w_1}$$

$$= \text{act}(W_2 y + b_2)$$

$$\boxed{\text{Loss} = \text{MSE} + L_p}$$



Masked Autoencoders Are Scalable Vision Learners

Masked Autoencoder

- Fine-tune is better.
- Random sampling works better.
- Masked autoencoder has better generalization.

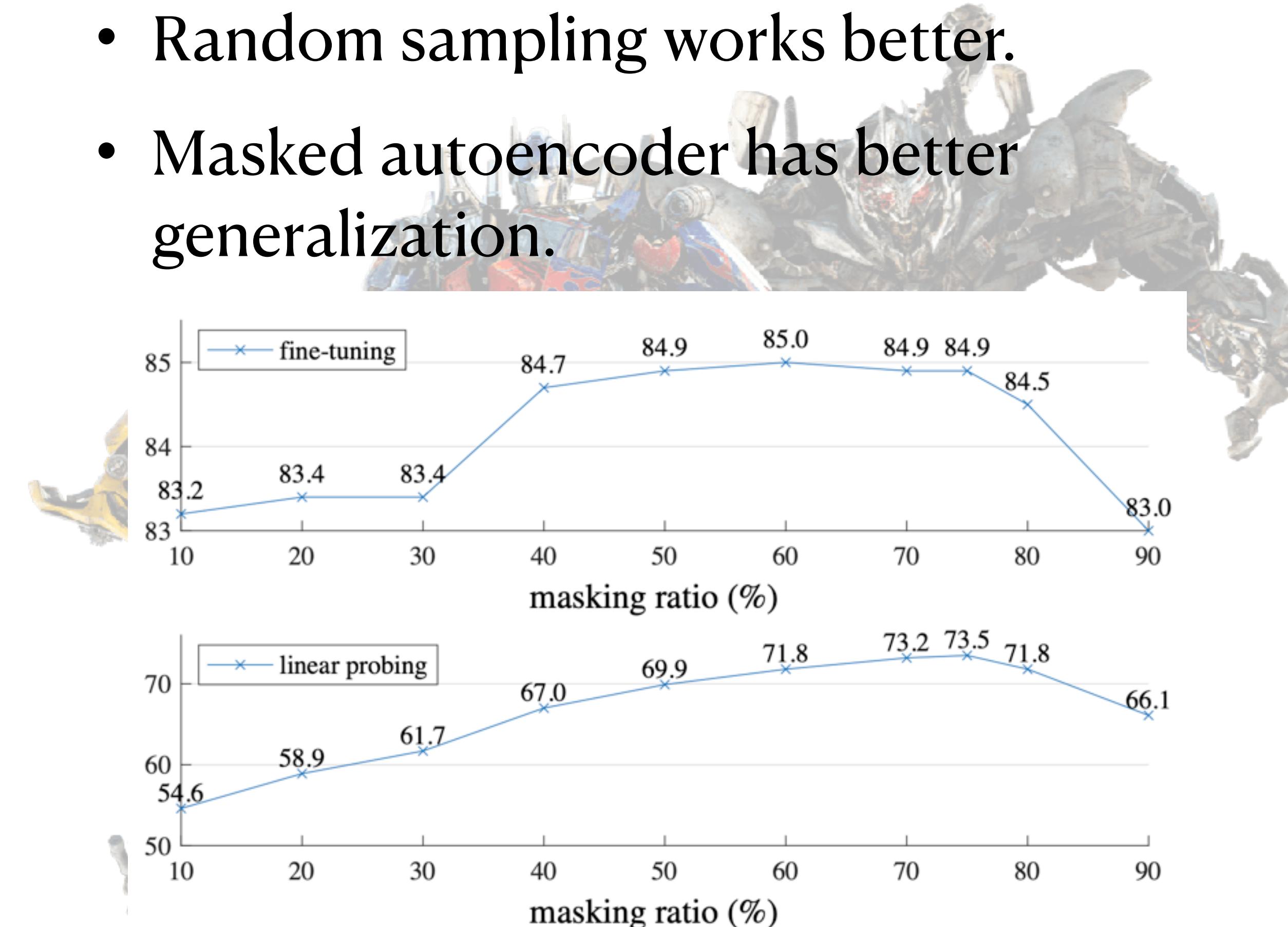


Fig. 3-8 Masking ratio v.s. validation accuracy (Image-1K)

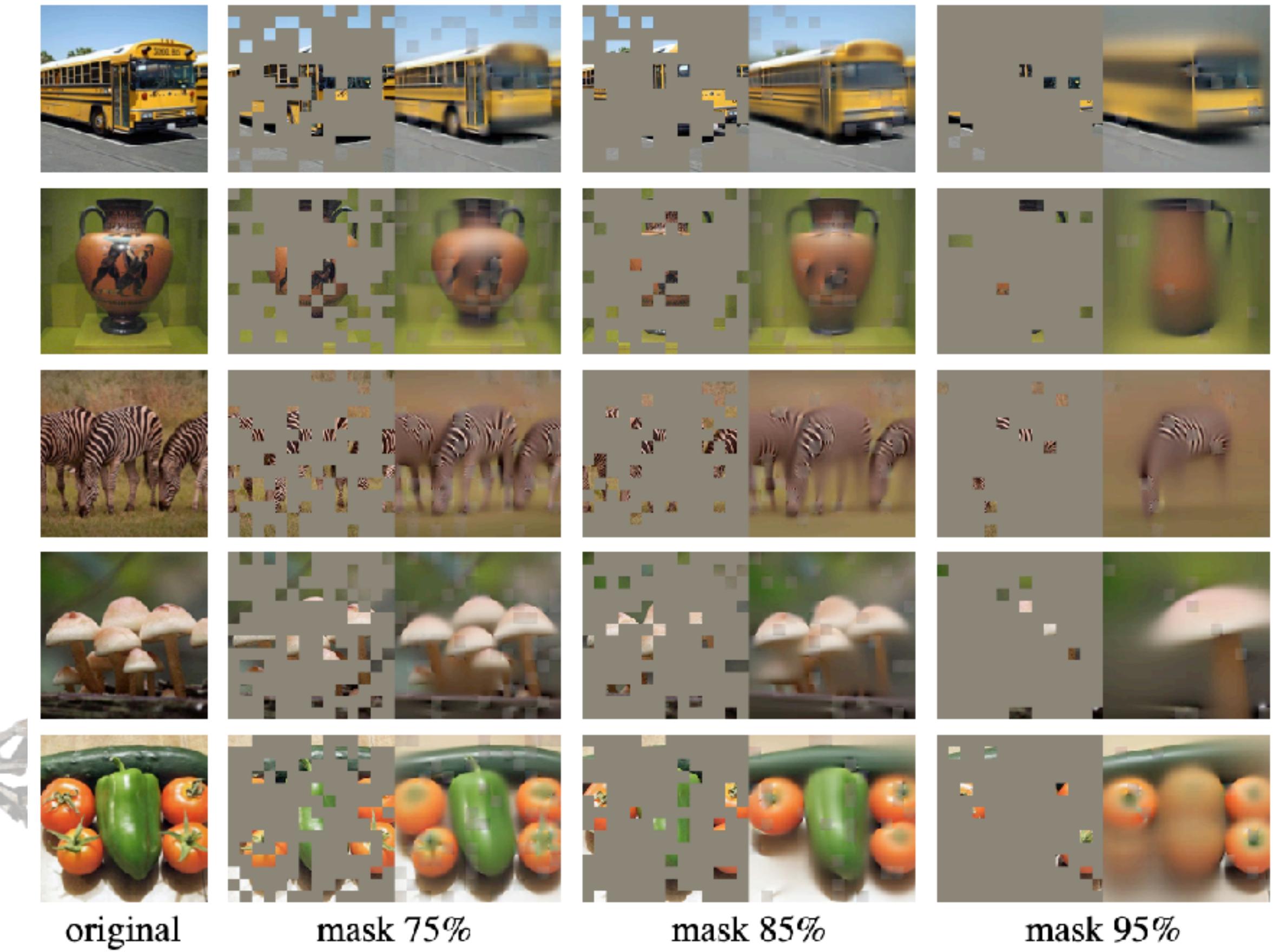


Fig. 3-9 MAE pretrained with 75% masking ratio but applied on inputs with higher masking ratio.

case	ft	lin	FLOPs	case	ratio	ft	lin
				random	75	84.9	73.5
encoder w/ [M]	84.2	59.6	3.3×	block	50	83.9	72.3
encoder w/o [M]	84.9	73.5	1×	block	75	82.8	63.9
				grid	75	84.0	66.0

Tab. 3-1 The comparison between 'with' or 'without' mask tokens

Tab. 3-2 Random sampling works best

Masked Autoencoder

Architecture

- Skipping the mask token in the encoder.
- a smaller decoder (1-block), a larger encoder (ViT-H)
- Partial Fine-tune

encoder	dec. depth	ft acc	hours	speedup
ViT-L, w/ [M]	8	84.2	42.4	-
ViT-L	8	84.9	15.4	2.8×
ViT-L	1	84.8	11.6	3.7×
ViT-H, w/ [M]	8	-	119.6 [†]	-
ViT-H	8	85.8	34.5	3.5×
ViT-H	1	85.9	29.3	4.1×



Tab. 3-3 Wall-clock time of MAE 800 epochs training in 128 TPU-V3 cores with TensorFlow. Decoder width=512. Mask ratio = 75%

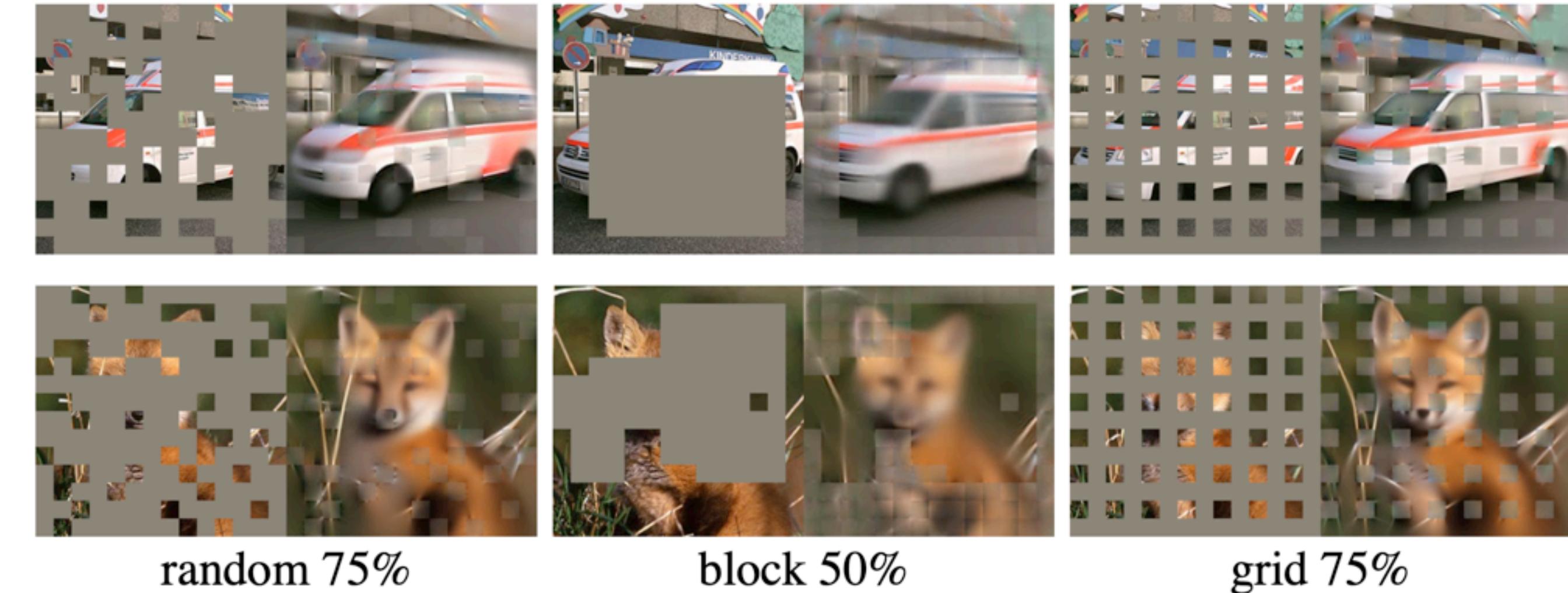


Fig. 3-10 Mask sampling strategies determine the pretext task difficulty, influencing reconstruction quality and representations

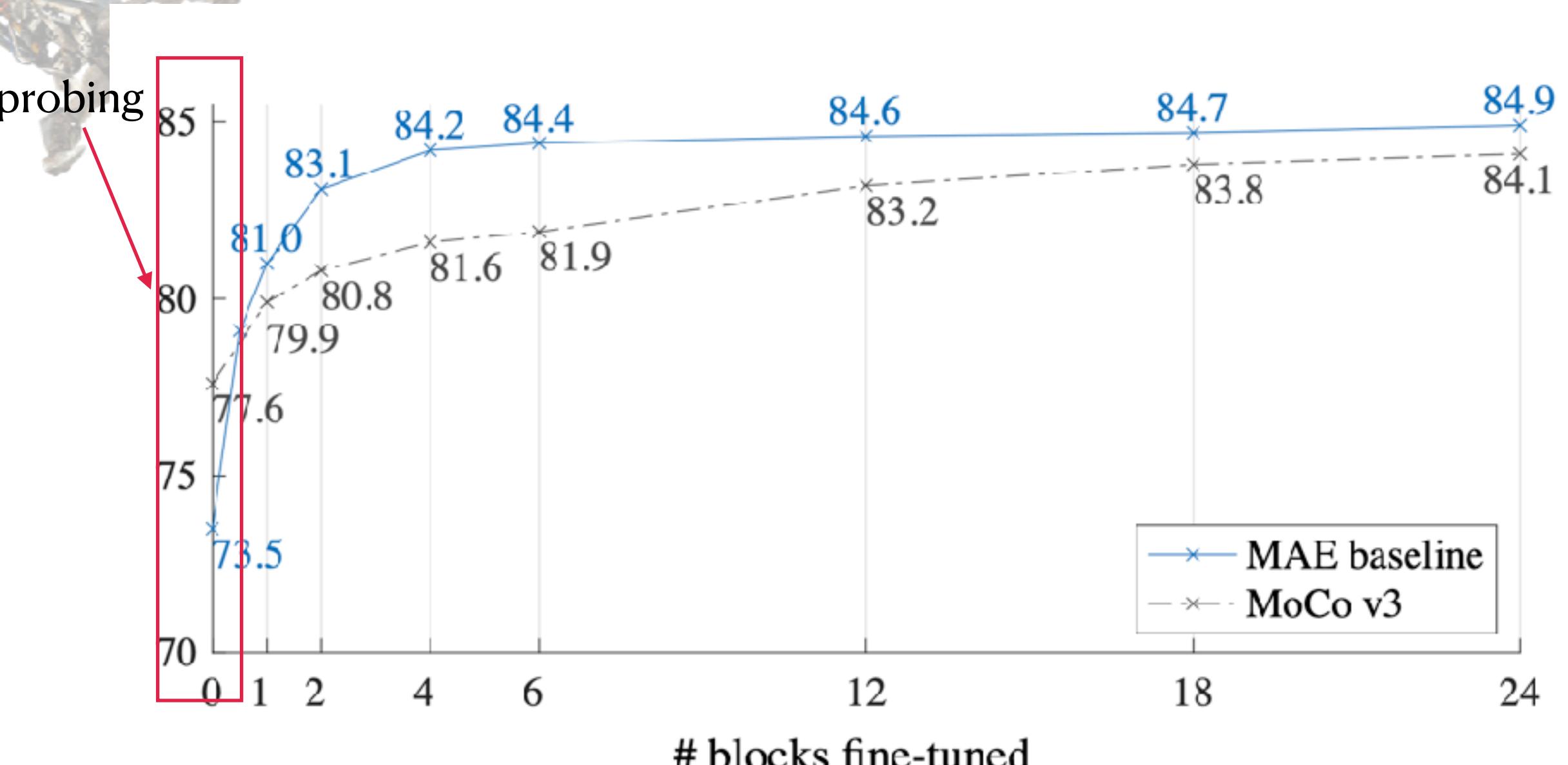


Fig. 3-11 Partial fine tuning, ViT-L using default setting

Conclusion

- CV is mimicking humans unconsciously.
- How to define optics?

Optics is the branch of physics that studies the behaviour and properties of light, including its interactions with matter and the construction of instruments that use or detect it.^[1] Optics usually describes the behaviour of visible, ultraviolet, and infrared light. Because light is an electromagnetic wave, other forms of electromagnetic radiation such as X-rays, microwaves, and radio waves exhibit similar properties.^[1]



- CI—» CV «—ONN

config	value
optimizer	AdamW [34]
base learning rate	1.5e-4
weight decay	0.05
optimizer momentum	$\beta_1, \beta_2=0.9, 0.95$ [6]
batch size	4096
learning rate schedule	cosine decay [33]
warmup epochs [19]	40
augmentation	RandomResizedCrop

Table 7. Pre-training setting.

config	value
optimizer	AdamW
base learning rate	1e-3
weight decay	0.05
optimizer momentum	$\beta_1, \beta_2=0.9, 0.999$
layer-wise lr decay [10, 2]	0.75
batch size	1024
learning rate schedule	cosine decay
warmup epochs	5
training epochs	100 (B), 50 (L/H)
augmentation	RandAug (9, 0.5) [12]
label smoothing [45]	0.1
mixup [58]	0.8
cutmix [57]	1.0
drop path [26]	0.1 (B/L) 0.2 (H)

Table 8. End-to-end fine-tuning setting.

config	value
optimizer	LARS [55]
base learning rate	0.1
weight decay	0
optimizer momentum	0.9
batch size	16384
learning rate schedule	cosine decay
warmup epochs	10
training epochs	90
augmentation	RandomResizedCrop

Table 9. Linear probing setting. We use LARS with a large batch for faster training; SGD works similarly with a 4096 batch size.

Outline

- **Analogy of ONN**
- ONN implementation
- A glimpse
- A optics inspired design

Preliminary

Gerchberg-Saxton algorithm

- To find a phase that conforms to the measurements of I & O.

$$\begin{aligned} \bullet \quad \mathbf{U}(x, y; d) &= \mathcal{F}^{-1}(\mathcal{F}(\mathbf{U}(x, y; 0)) \odot \mathbf{G}(f_x, f_y)) \\ &= \mathbf{H}(\mathbf{U}(x, y; 0)) \end{aligned} \quad (1)$$

$$\mathbf{G}(f_x, f_y) = \exp \left[i k d \sqrt{1 - (\lambda f_x)^2 - (\lambda f_y)^2} \right] \quad (2)$$

$$\mathbf{Y}(x, y; d) = |\mathbf{U}(x, y; d)|^2 = |\mathbf{H}(\mathbf{U}(x, y; 0))|^2 \quad (3)$$

$$\hat{\mathbf{U}}(x, y, 0) = \arg \min_{\mathbf{U}(x, y; 0)} \|\mathbf{Y}(x, y; d) - |\mathbf{H}(\mathbf{U}(x, y; 0))|^2\|_2^2 \quad (4)$$

$$\mathbf{I}^{(k)}(x, y; d) = \sqrt{\mathbf{Y}(x, y; d)} \cdot \exp[i \cdot \text{angle}(\mathbf{U}^{(k)}(x, y; d))]$$

$$\mathbf{O}^{(k)}(x, y; 0) = \mathbf{H}^{-1}(\mathbf{I}^{(k)}(x, y; d))$$

$$\mathbf{U}^{(k+1)}(x, y; 0) = \mathbf{P}(x, y; 0) \cdot \exp[i \cdot \text{angle}(\mathbf{O}^{(k)}(x, y; 0))]$$

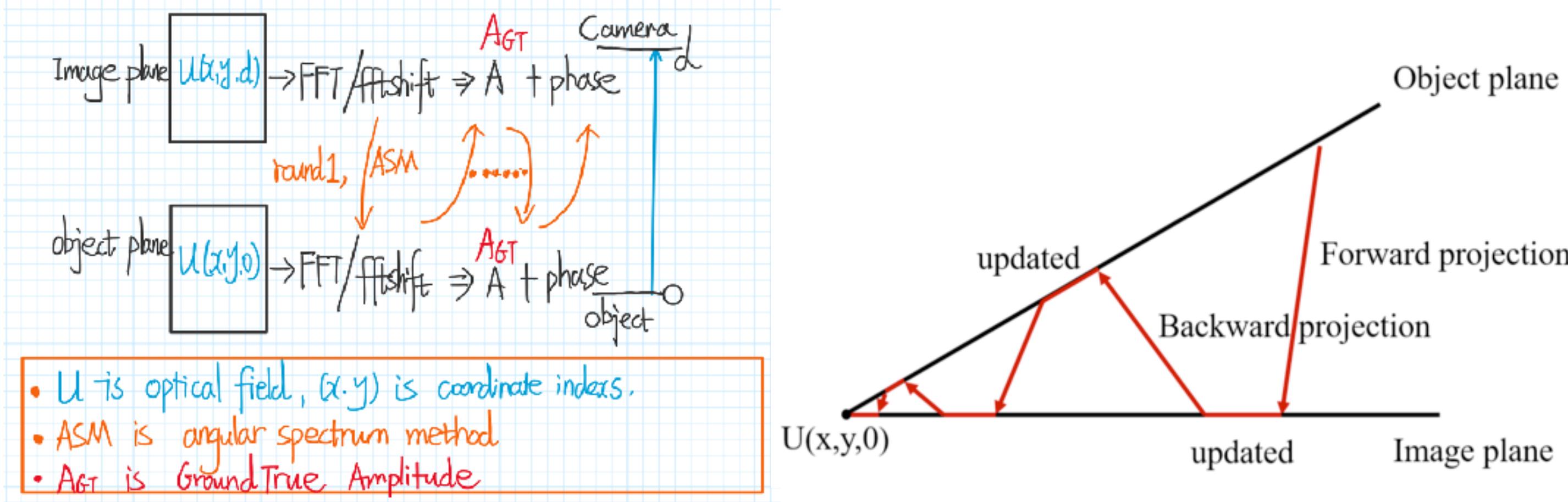


Fig. 1 The process of Gerchberg-Saxton algorithm

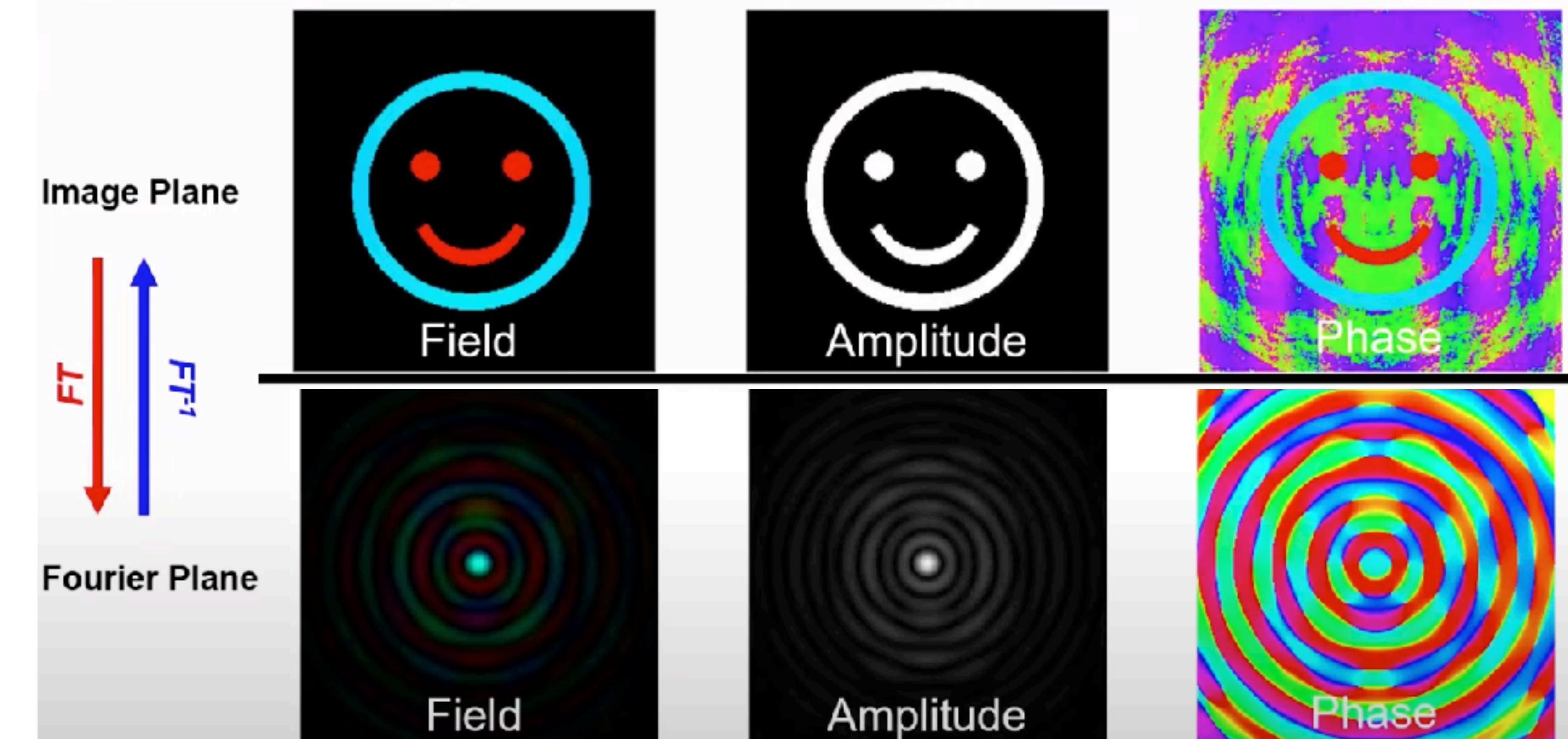


Fig. 2 Illustration of GS algorithm

Preliminary

Angular spectrum method

- **Angular spectrum method:**
 1. Decompose $U(x, y, 0)$ into 2d plane waves;
 2. Propagate each plane wave separately to the plane d ;
 3. Add the propagated plane waves together ($U(x, y, d)$)
- Where the U is optical field, k is wave vector, f is frequency, and F is fourier transformation

$$\begin{aligned}
 U(x, y, 0) &= \iint \hat{U}_0(f_x, f_y) \cdot e^{i2\pi(f_x x + f_y y)} df_x df_y \\
 &\xrightarrow{\text{Fourier i.e. } \hat{U}_0(f_x, f_y) = \frac{1}{2\pi} \iint U(x, y, 0) e^{-i2\pi(f_x x + f_y y)} dx dy} \\
 &= \iint \hat{U}_0(f_x, f_y) \cdot e^{i2\pi(f_x x + f_y y + f_z \cdot 0)} df_x df_y \\
 &\text{As: } |\vec{k}| = \frac{2\pi}{\lambda} = \sqrt{k_x^2 + k_y^2 + k_z^2} \Rightarrow f_z = \pm \frac{1}{2\pi} \sqrt{\left(\frac{2\pi}{\lambda}\right)^2 - k_x^2 - k_y^2} \\
 &\quad = \pm \sqrt{\left(\frac{1}{\lambda}\right)^2 - f_x^2 - f_y^2} \\
 \text{则 } V(x, y, z_0) &= \iint \hat{U}_0(f_x, f_y) \cdot e^{i2\pi(f_x x + f_y y + z_0 \sqrt{\left(\frac{1}{\lambda}\right)^2 - f_x^2 - f_y^2})} df_x df_y \\
 &= \iint F \{ U(x, y, 0) \} \cdot e^{i2\pi(f_x x + f_y y + z_0 \sqrt{\left(\frac{1}{\lambda}\right)^2 - f_x^2 - f_y^2})} df_x df_y
 \end{aligned}$$

Fig. 3 The process of Gerchberg-Saxton algorithm

Preliminary compressed sensing

- Conditions: **Sparsity** and **incoherence**

$$\min_X \|\nabla X \bullet d\|_1 + \frac{\lambda}{2} \|Y - \Phi X\|_2^2$$
- where the $d = (dh, dv)$ define the horizontal and vertical estimates of d
- $\Delta X = \text{grad } X$, X is Object & Y is measurement

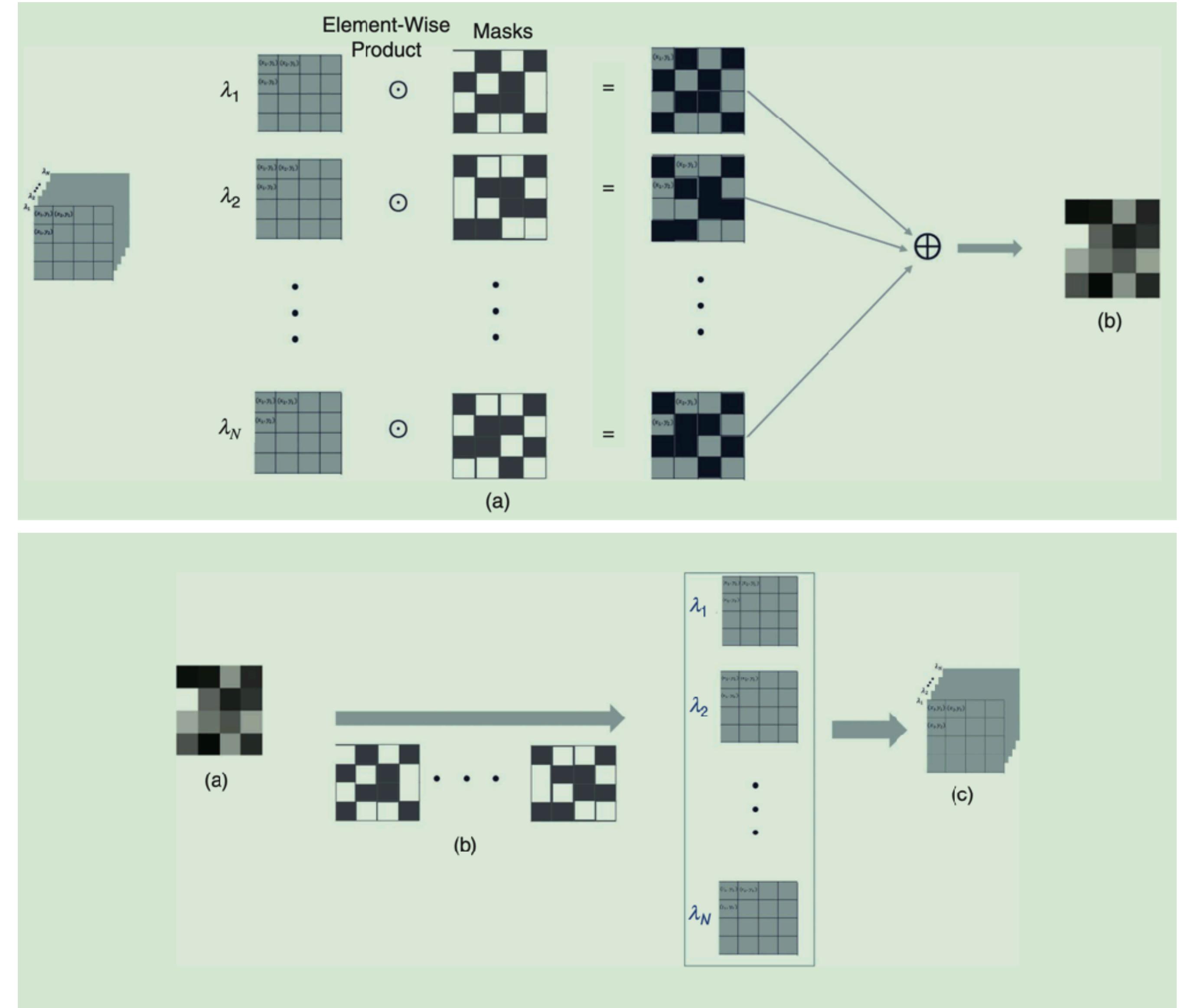


Fig. 4 (x,y,lambda) compressed sensing

PR in Lensless microscopy imaging

Current challenges

1. Limited dataset
2. Difficult to Supervise
3. Additional hardware

Schematic

- a is setup of lensless microscopy imaging
- b is structure diagram of the proposed NN
- c is complex-valued U-net (2D)

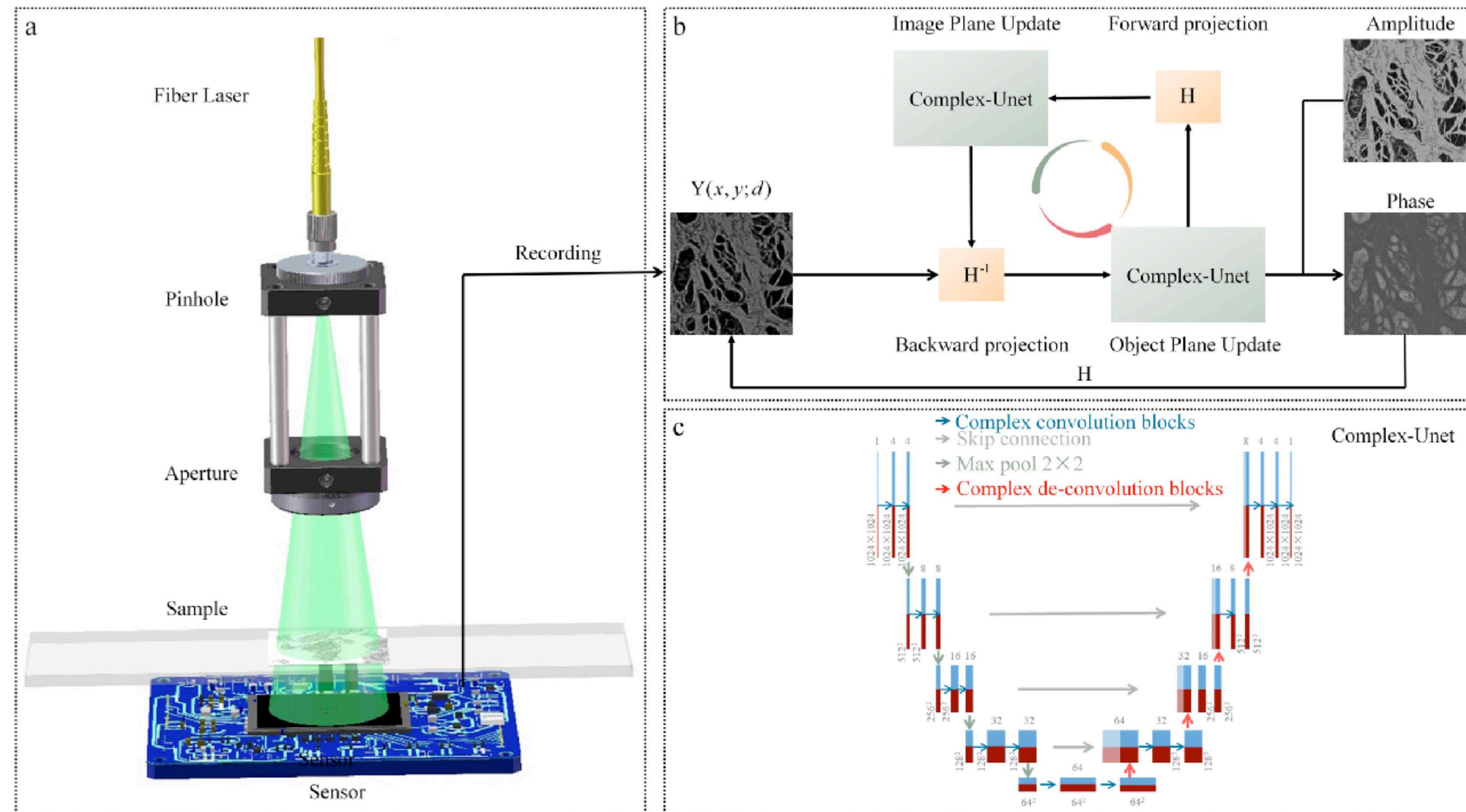


Fig. 5 Schematic diagram of the imaging process

PR in Lensless microscopy imaging

- Backward → forward → Backward
- Constraint : ASM and I from CMOS
- Complex-valued UNet:

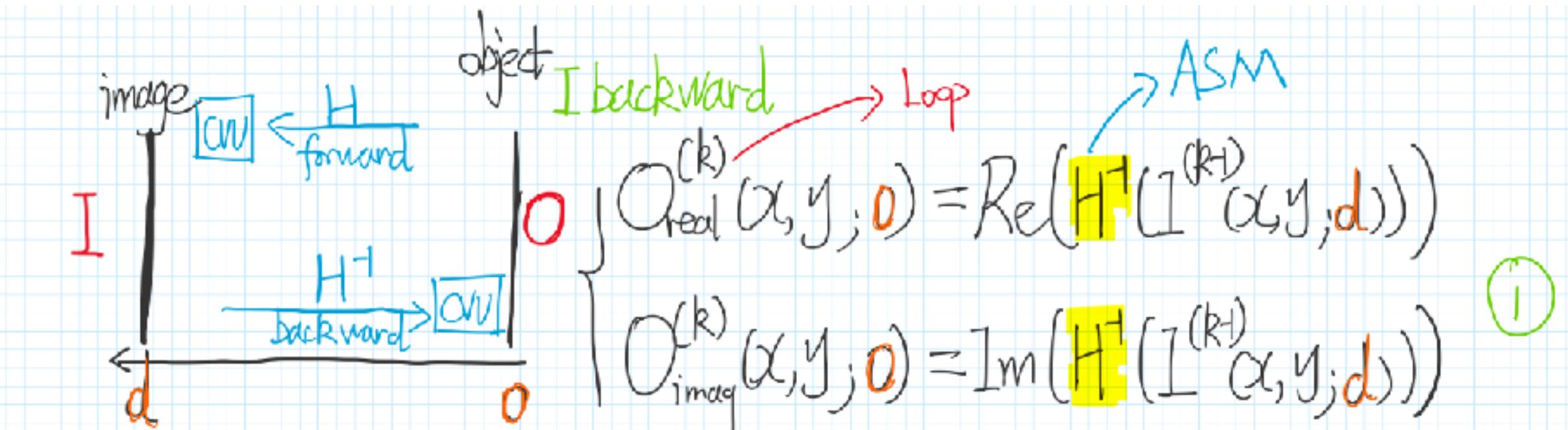
Convolution

$$\begin{aligned} \mathbf{F} * \mathbf{h} &= (\mathbf{F}_r + i\mathbf{F}_i) * (\mathbf{h}_r + i\mathbf{h}_i) \\ &= (\mathbf{F}_r * \mathbf{h}_r - \mathbf{F}_i * \mathbf{h}_i) + i(\mathbf{F}_r * \mathbf{h}_i + \mathbf{F}_i * \mathbf{h}_r) \end{aligned}$$

LeakReLU

$$\mathbb{C}LeReLU(z) = LeReLU(z_r) + iLeReLU(z_i)$$

Why LeakReLU?



$$\text{Signal restoration} \Rightarrow U^{(k)}(x, y; 0) = CV_{\text{Net}}(O^{(k)}_{\text{CV}}(x, y; 0), \theta) \quad \textcircled{2}$$

[happen on O plane]

$$\begin{aligned} \text{II forward} \Rightarrow V^{(k)}_{\text{real}}(x, y; d) &= Re(H(V^{(k)}(x, y; 0))) \\ V^{(k)}_{\text{imag}}(x, y; d) &= Im(H(V^{(k)}(x, y; 0))) \end{aligned} \quad \textcircled{3}$$

$$\text{Signal recovery} \Rightarrow I^{(k)}(x, y; d) = CV_{\text{Net}}(V^{(k)}(x, y; d), \phi) \quad \textcircled{4}$$

III backward as $\textcircled{1} \otimes \textcircled{2}$

then

$$W^* = \arg \min_W \| H(N_w(Y)) \|^2 - \| Y \|^2$$

PR in Lensless microscopy imaging

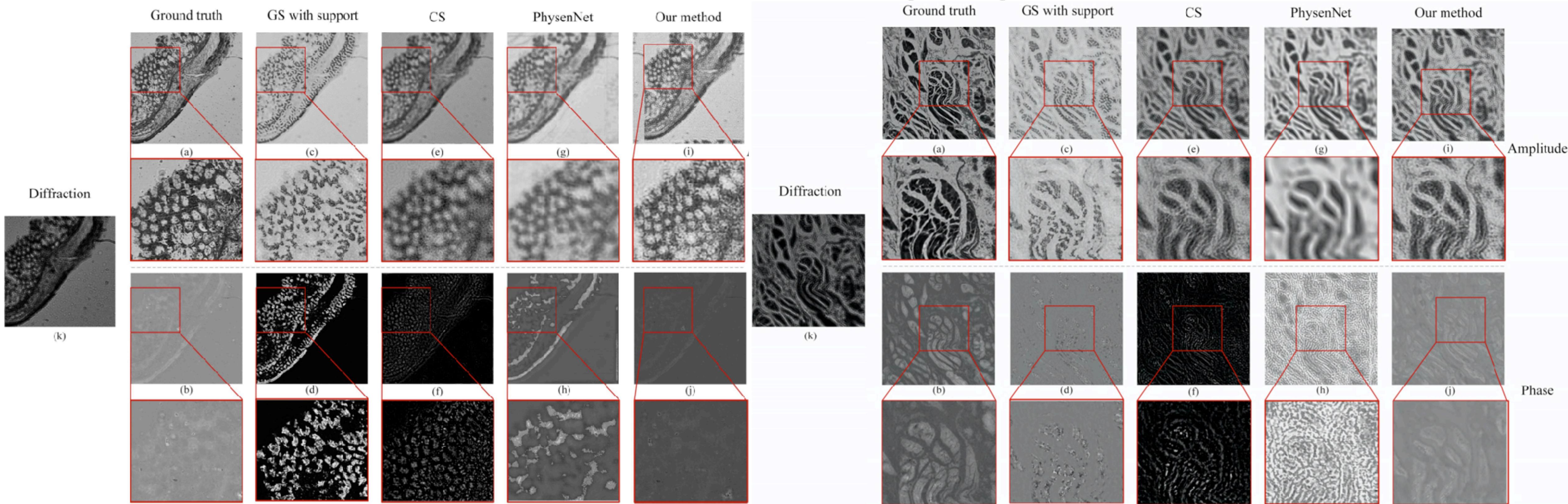
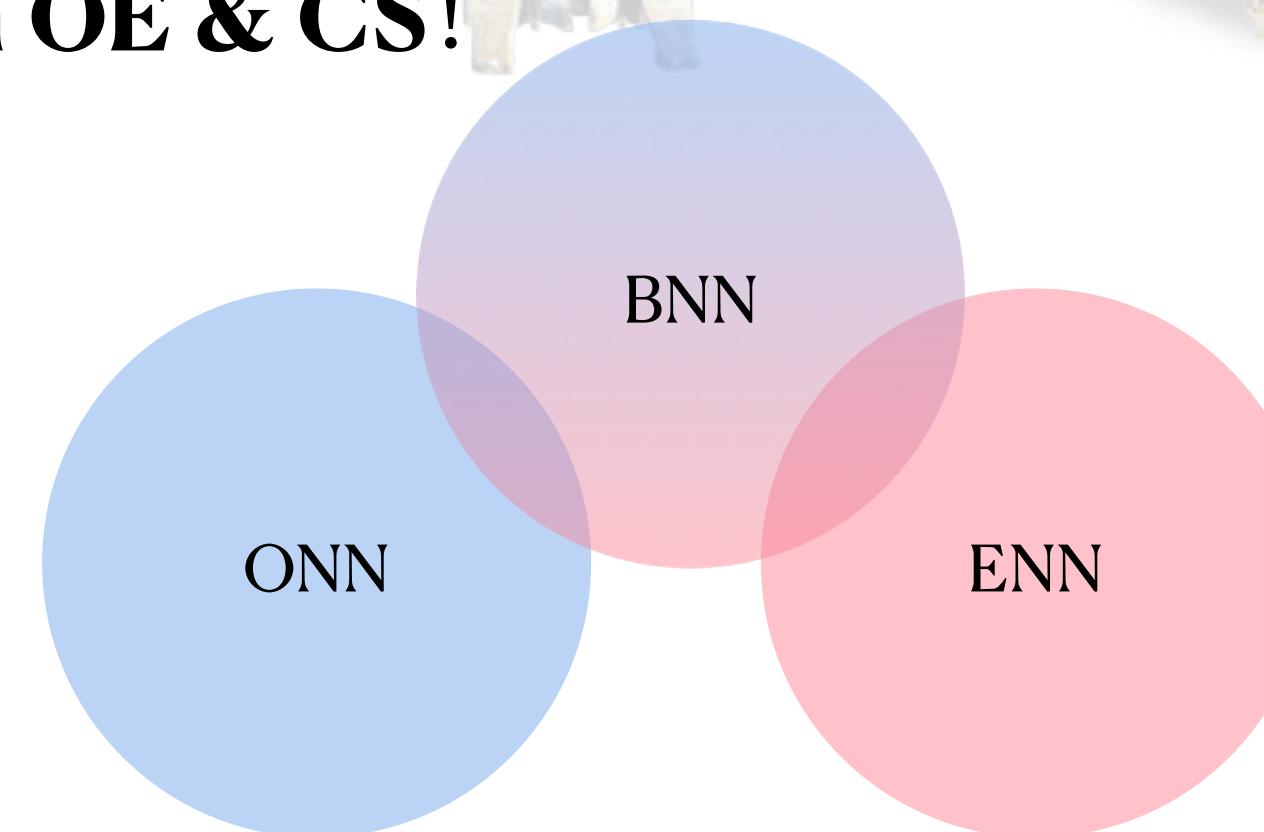


Fig. 6 Visual comparisons on the H&E stained pathological slides of rat intestine

Fig. 6 Visual comparisons on the H&E stained pathological slides of human esophagus cancer cell

Conclusion

- OC: or photonic computing, using photons (from correlation source) for computation
- CI: process of indirectly imaging from measurements using algos that rely on computing.
- ONN: a physical implementation of an **artificial neural network** with **optical components**.
- BNNs function on an electrochemical basis. ONNs use electromagnetic waves. ENNs are based on **simulation**.
- We can interchange techiques in OE & CS!



Outline

- **Analogy of ONN (E₂O | imaging)**
- ONN implementation (E₂OE | reasoning)
- A glimpse (B₂E | reasoning)
- A optics inspired design (O₂E | imaging)