



KYUSHU
UNIVERSITY

FUJITSU



ACT-∞

Layer-wise Interpretation of Deep Neural Networks Using Identity Initialization



Shohei Kubota
(Kyushu Univ.)



Hideaki Hayashi
(Kyushu Univ.)

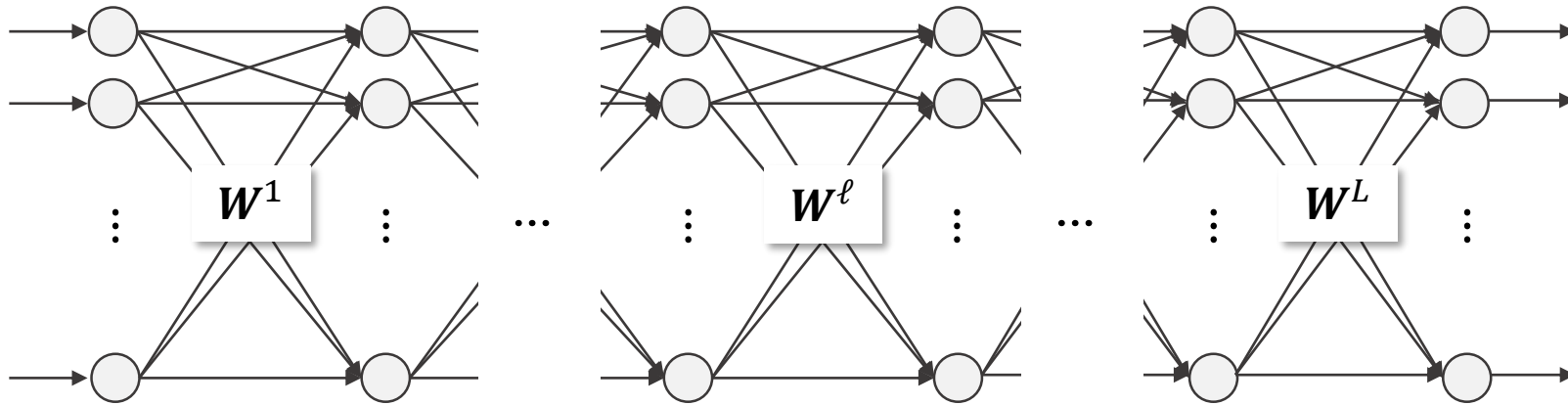


Tomohiro Hayase
(Fujitsu Lab.)



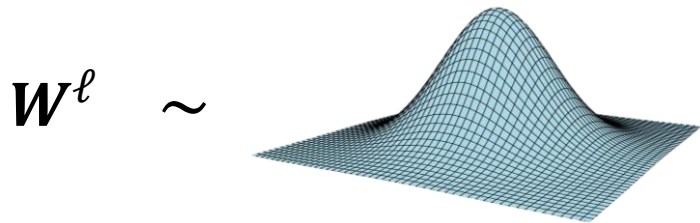
Seiichi Uchida
(Kyushu Univ.)

Can identity initialized deep neural nets be trained?



Standard initialization

- Random matrix



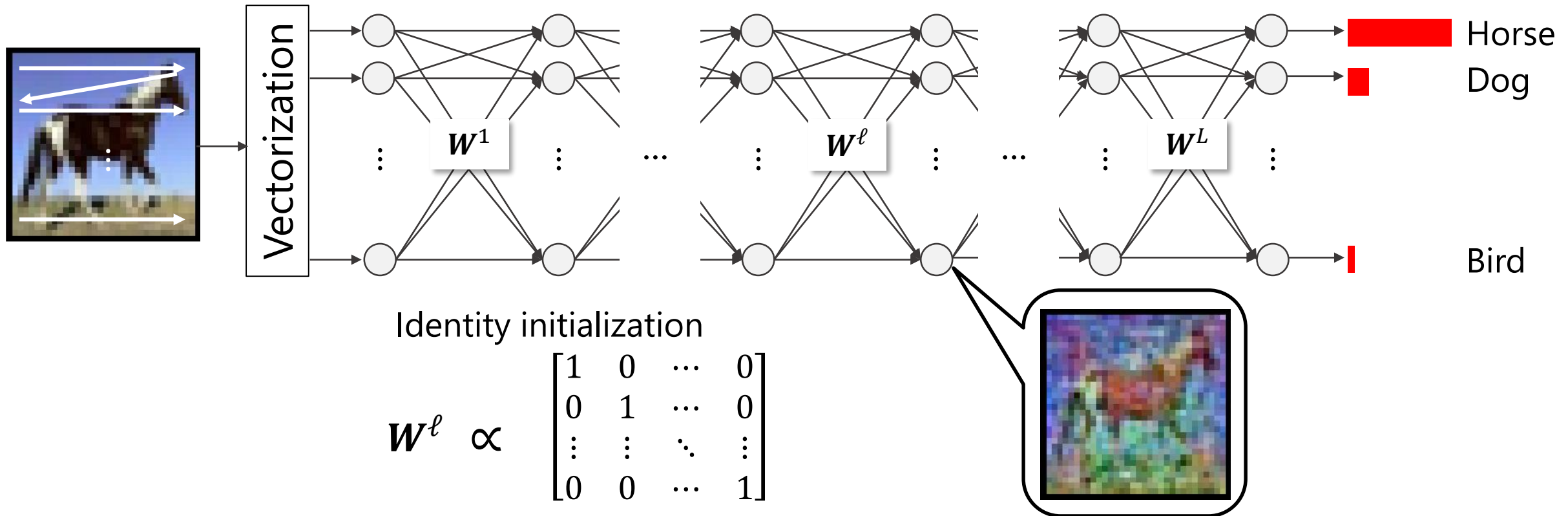
Initialization in this study

- Identity matrix

$$W^l \propto \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix}$$

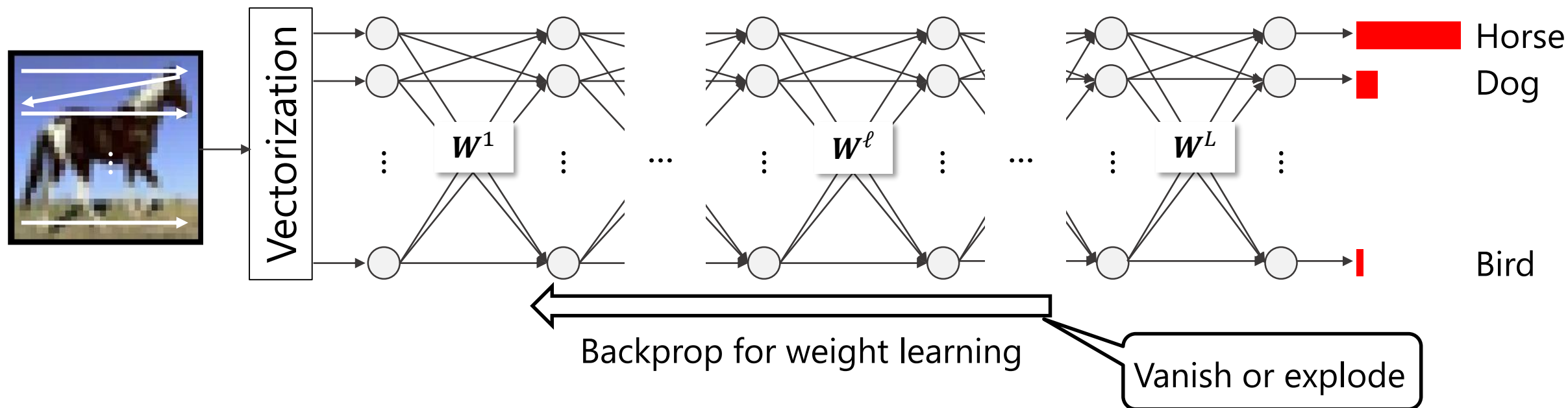
Interpretability of the internal process

- Learned weights are close to the identity matrix
- Intermediate features preserve the input structure



Gradient vanishing and exploding

- Identity-initialized deep neural networks suffer from the gradient vanishing/exploding problems



Investigate the potential of identity initialization

- **Theory:**

We show a condition under which the identity-initialized deep multilayer perceptron (MLP) prevents gradient vanishing/exploding.

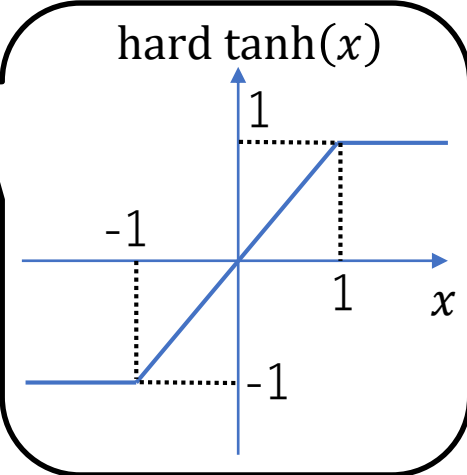
- **Application:**

We propose an interpretable MLP structure using identity initialization.

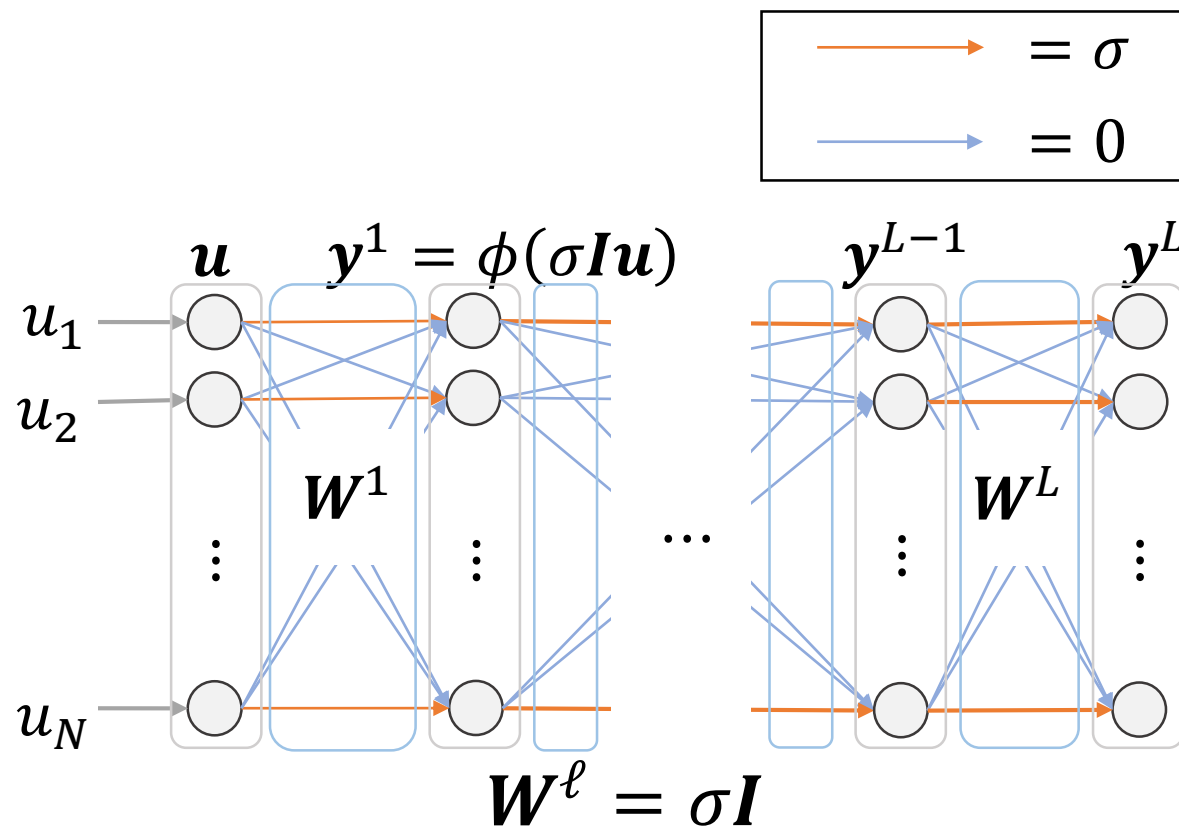
Identity initialization

- **MLP of L layers with width N and $W^\ell \in \mathbb{R}^{N \times N}$**

- **Initialization: $W^\ell = \sigma I$**

- **Activation ϕ**


- **Input u follows $N(0, q^0)$**



Dynamical isometry

- **Singular values of Jacobian J (or eigenvalues of $J^T J$) concentrate around 1**

$$J = \frac{\partial \mathbf{y}^L}{\partial \mathbf{u}} = \mathbf{D}^1 \mathbf{W}^1 \mathbf{D}^2 \mathbf{W}^2 \dots \mathbf{D}^L \mathbf{W}^L$$

Output (red arrow pointing to \mathbf{y}^L)

Input (red underline under \mathbf{u})

Weight of the first layer (purple underline under \mathbf{W}^1)

Derivative of the activation function at the L -th layer (blue underline under \mathbf{D}^L)

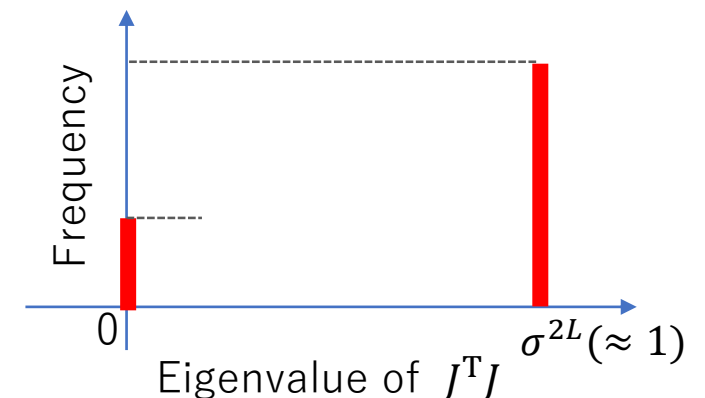
Theoretical result

- **The eigenvalue distribution is a Bernoulli function**

$$\mu_{JJ^T}(\lambda) = \begin{cases} \alpha^L \delta(\lambda - \sigma^{2L}) - (1 - \alpha^L) \delta(\lambda) & (\sigma > 1) \\ \alpha^1 \delta(\lambda - \sigma^{2L}) - (1 - \alpha^1) \delta(\lambda) & (\sigma \leq 1) \end{cases}$$

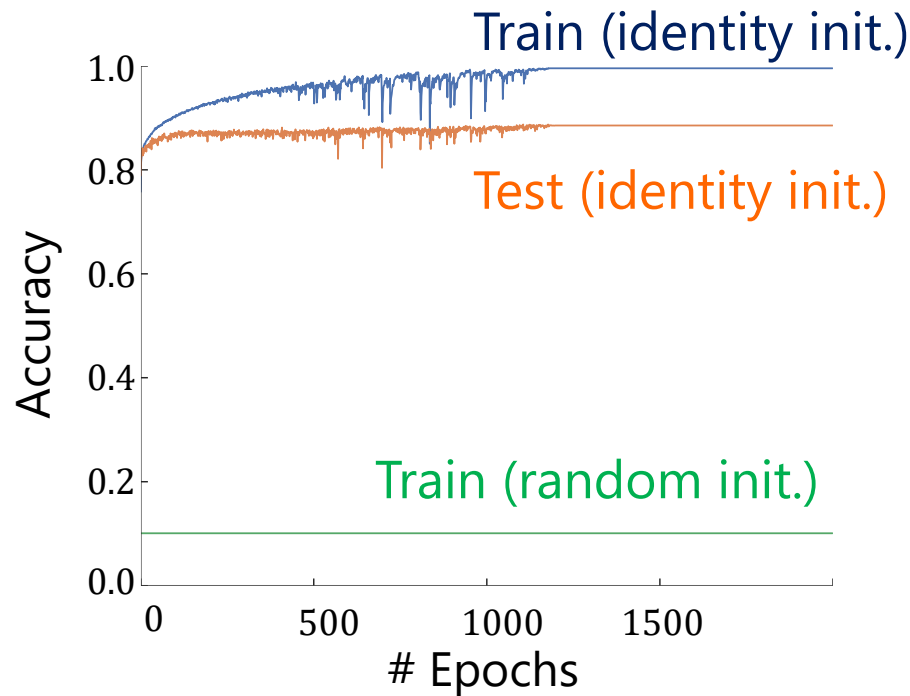
α^ℓ : Constant depending on σ (constant) and q^0 (input variance)

- **We can prevent gradient vanishing by setting σ and q^0 appropriately**

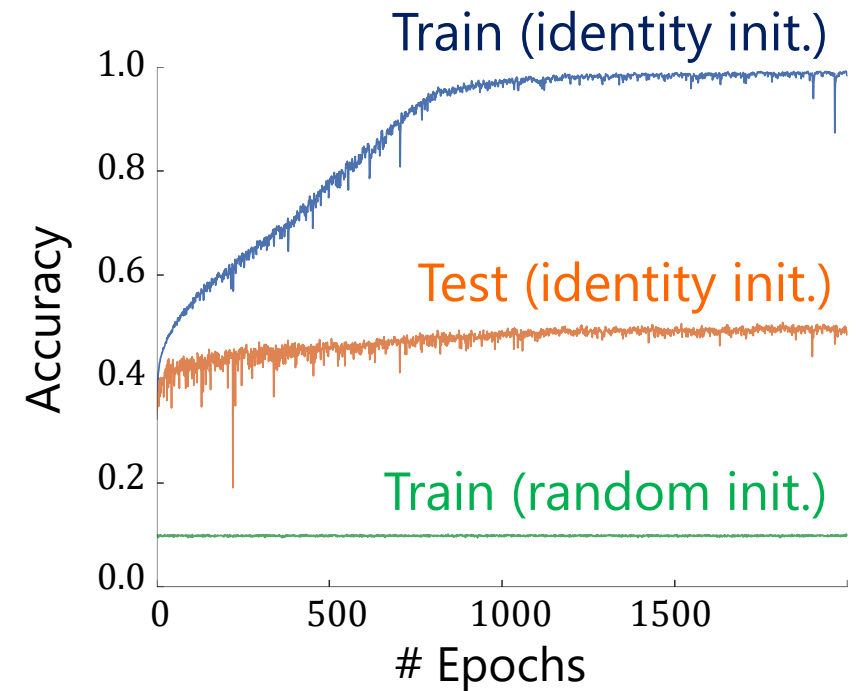


Identity-initialized networks can be trained

- **MLP with 100 layers**
- **Dataset: Fashion-MNIST and CIFAR-10 (10 classes)**



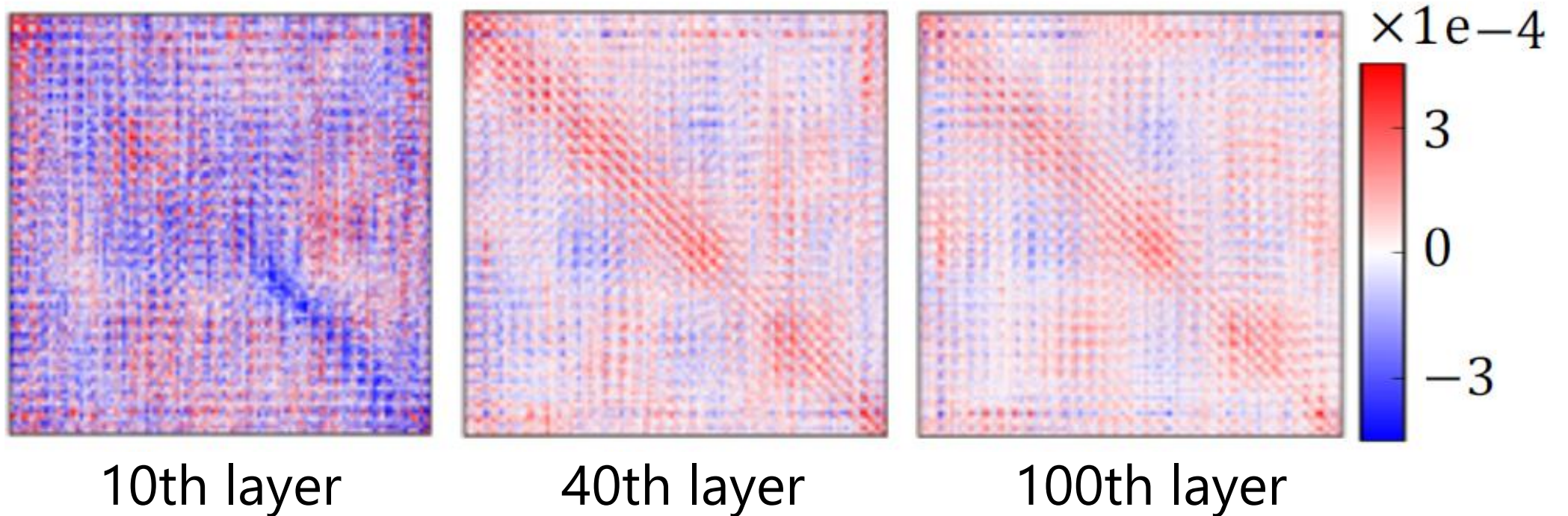
(a) Fashion-MNIST



(b) CIFAR-10

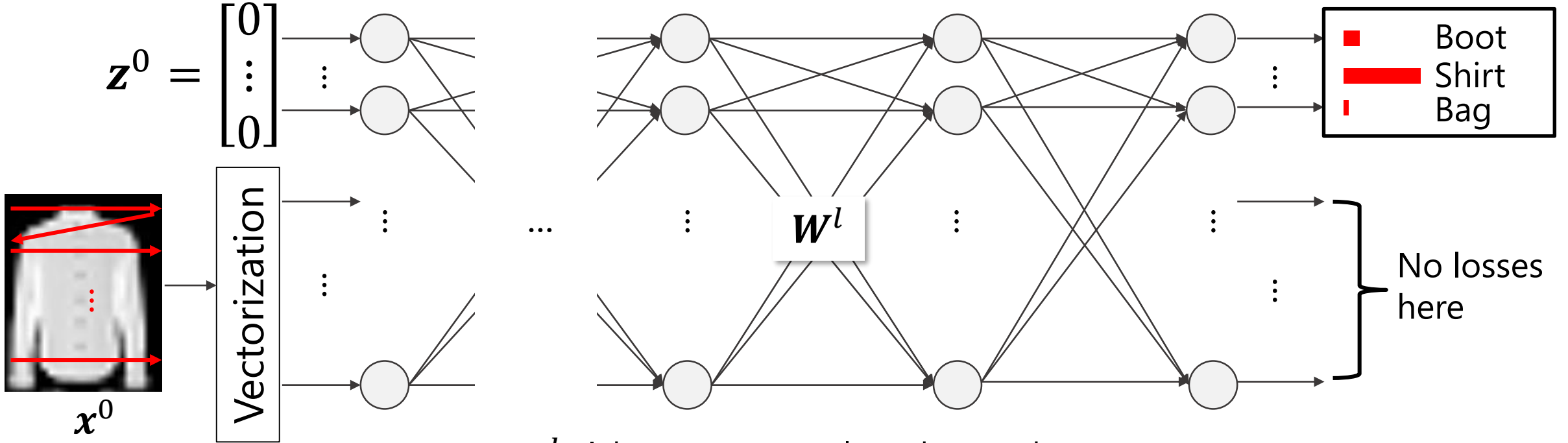
Close to the identity matrix

- **Difference of the weight matrix from the identity matrix after training**



Network structure that enhances interpretability

Same dimension as the number of classes

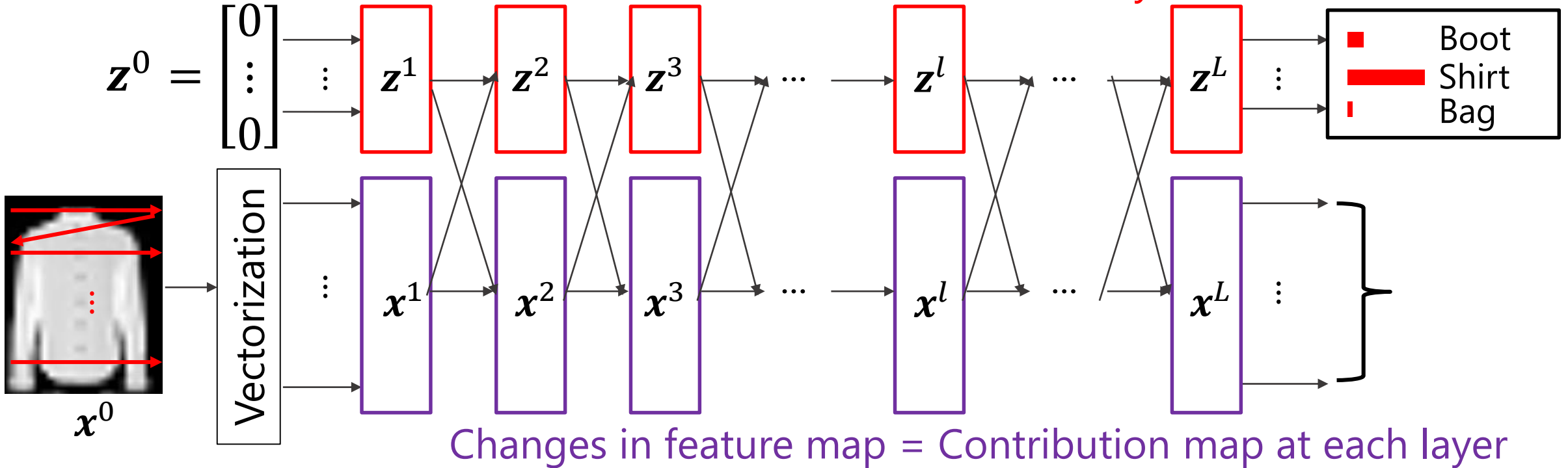


W^l : Identity-initialized weight

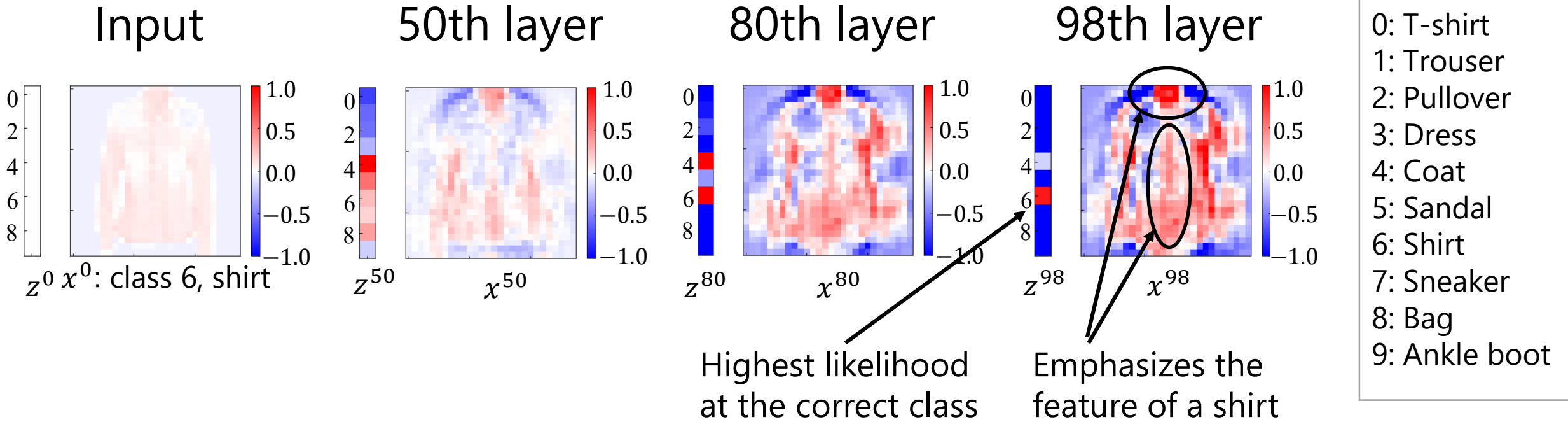
Network structure that enhances interpretability

Same dimension as the number of classes

Likelihood at each layer



Emphasize the areas important for classification



z^l : Class likelihood at l -th layer

x^l : Feature map at l -th layer

Conclusion

- **Investigate the potential of identity initialization**
 - Condition for preventing the vanishing/exploding gradients
 - Network structure enhancing interpretability
- **Future work**
 - Analysis of changes in feature values during learning
 - Application to other structures such as convolutional neural networks



KYUSHU
UNIVERSITY

FUJITSU



ACT-

Layer-wise Interpretation of Deep Neural Networks Using Identity Initialization

Shohei Kubota, Hideaki Hayashi, Tomohiro Hayase, and Seiichi Uchida

Poster Session

MLSP-42: Neural Network Pruning

Friday, 11 June from 11:30 to 12:15 (EDT)