

# Layer-wise Interpretation of Deep Neural Networks Using Identity Initialization

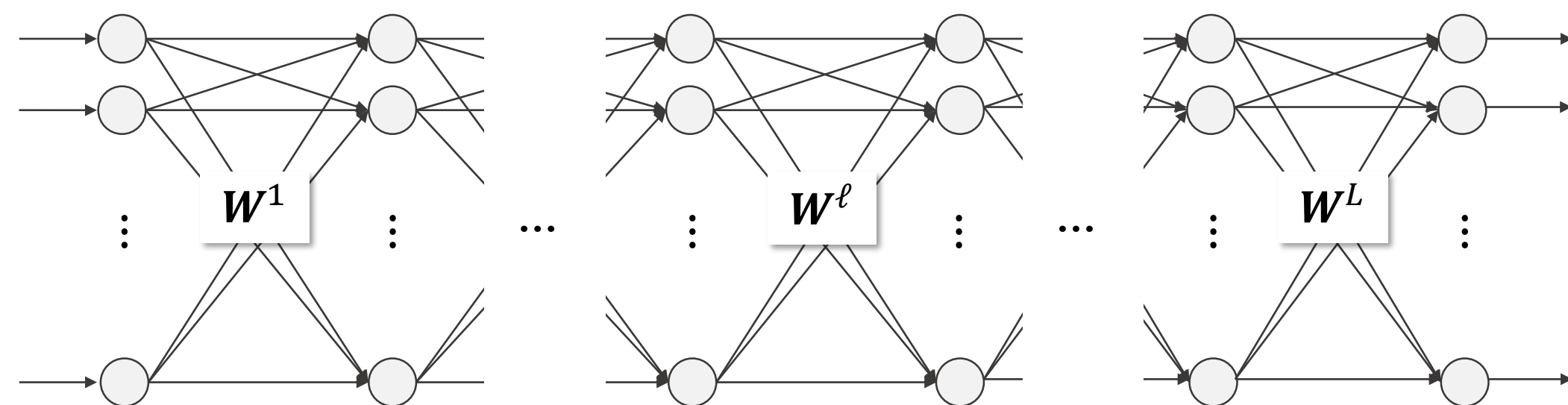
Shohei Kubota<sup>1</sup>, Hideaki Hayashi<sup>1</sup>, Tomohiro Hayase<sup>2</sup>, and Seiichi Uchida<sup>1</sup>

<sup>1</sup>Kyushu University, Japan, <sup>2</sup>Fujitsu lab., Japan



## Motivation

Can identity initialized deep neural networks be trained?



### Standard initialization

- Random matrix

$$W^\ell \sim \text{random matrix}$$

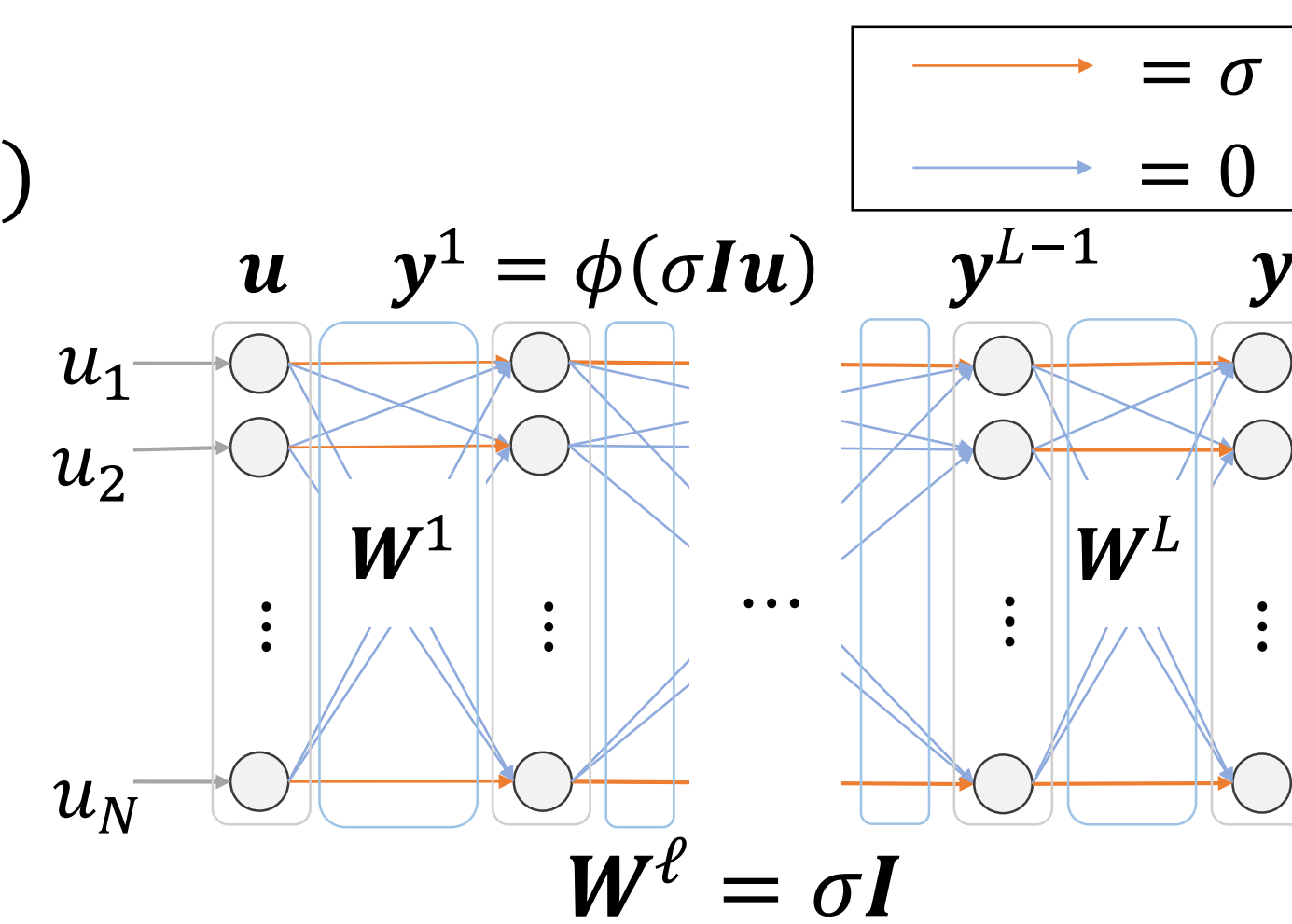
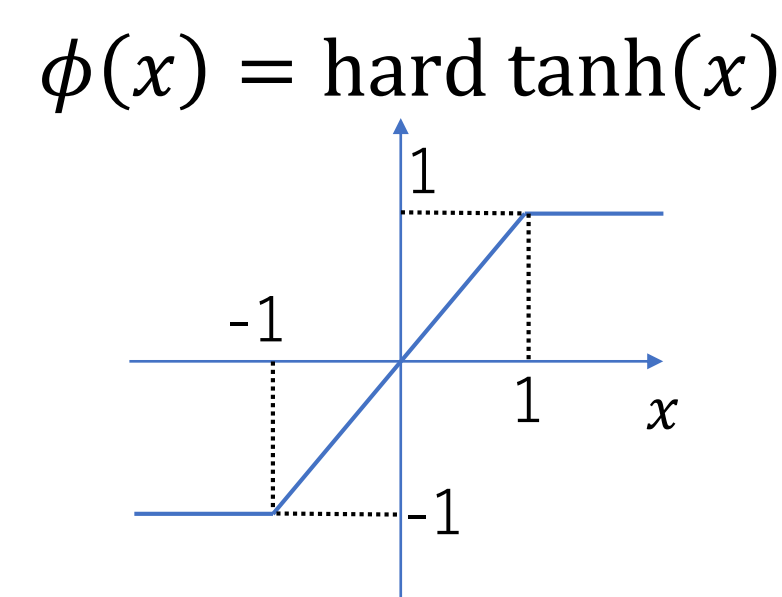
### Initialization in this study

- Identity matrix

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

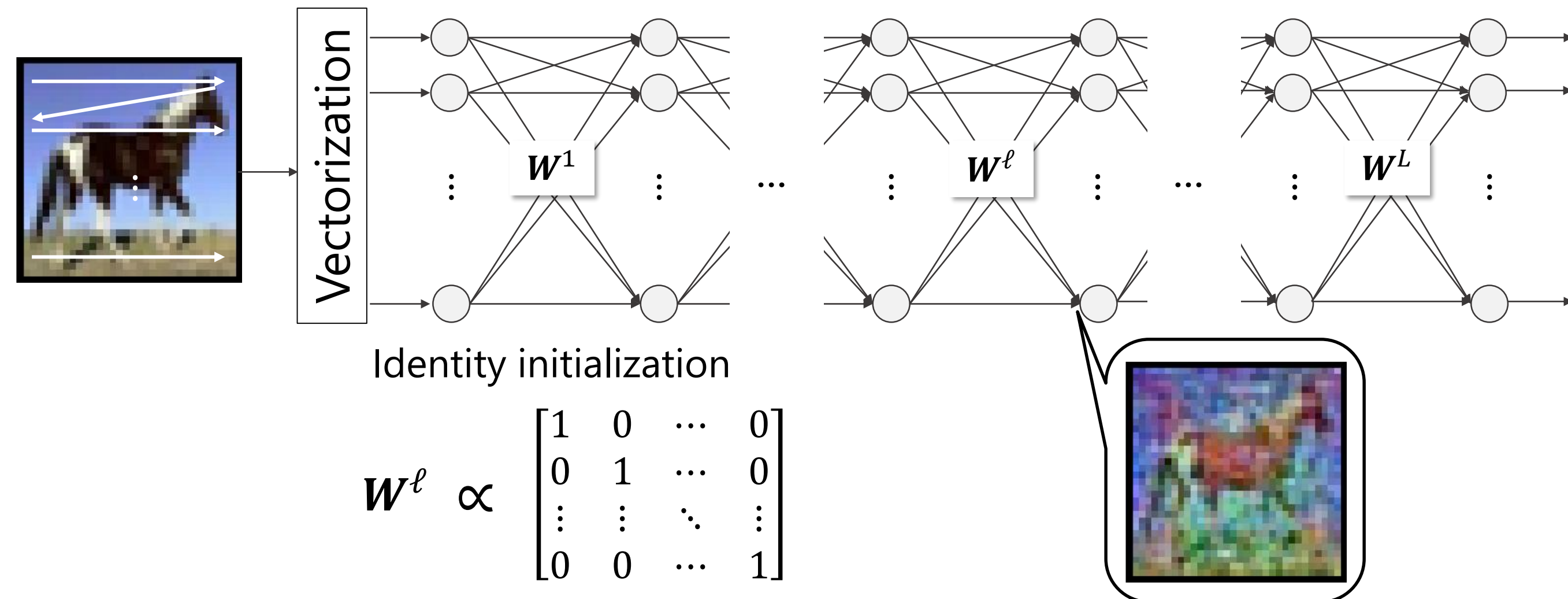
## Identity initialization

- MLP of  $L$  layers with width  $N$  and  $W^\ell \in \mathbb{R}^{N \times N}$
- Initialization:  $W^\ell = \sigma I$
- Input  $u$  follows  $N(0, q^0)$
- Activation  $\phi$



## Merit | Interpretability

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



## Difficulty | Gradient vanishing/exploding

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

## Contribution

- **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.

## Condition to prevent gradient vanishing

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

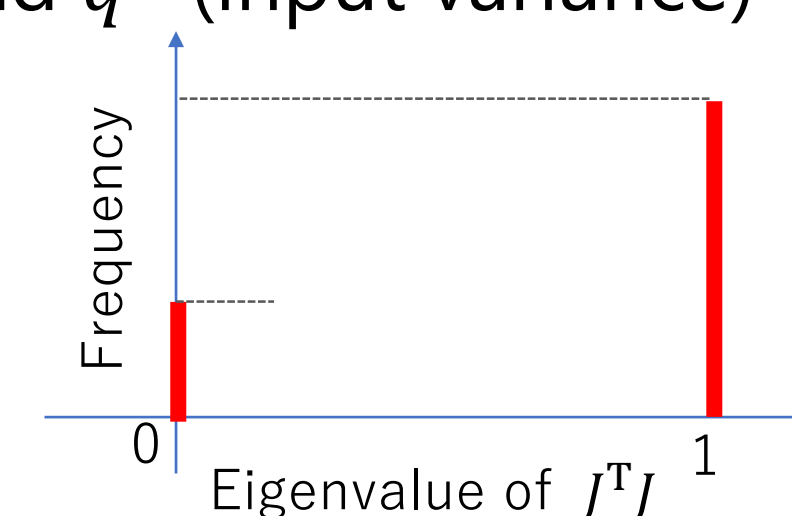
$$J = \frac{\partial y^L}{\partial u} = D^1 W^1 D^2 W^2 \dots D^L W^L$$

$D^l$ : Derivative of the activation function at the  $l$ -th layer  
 $W^l$ : Weight of the first layer

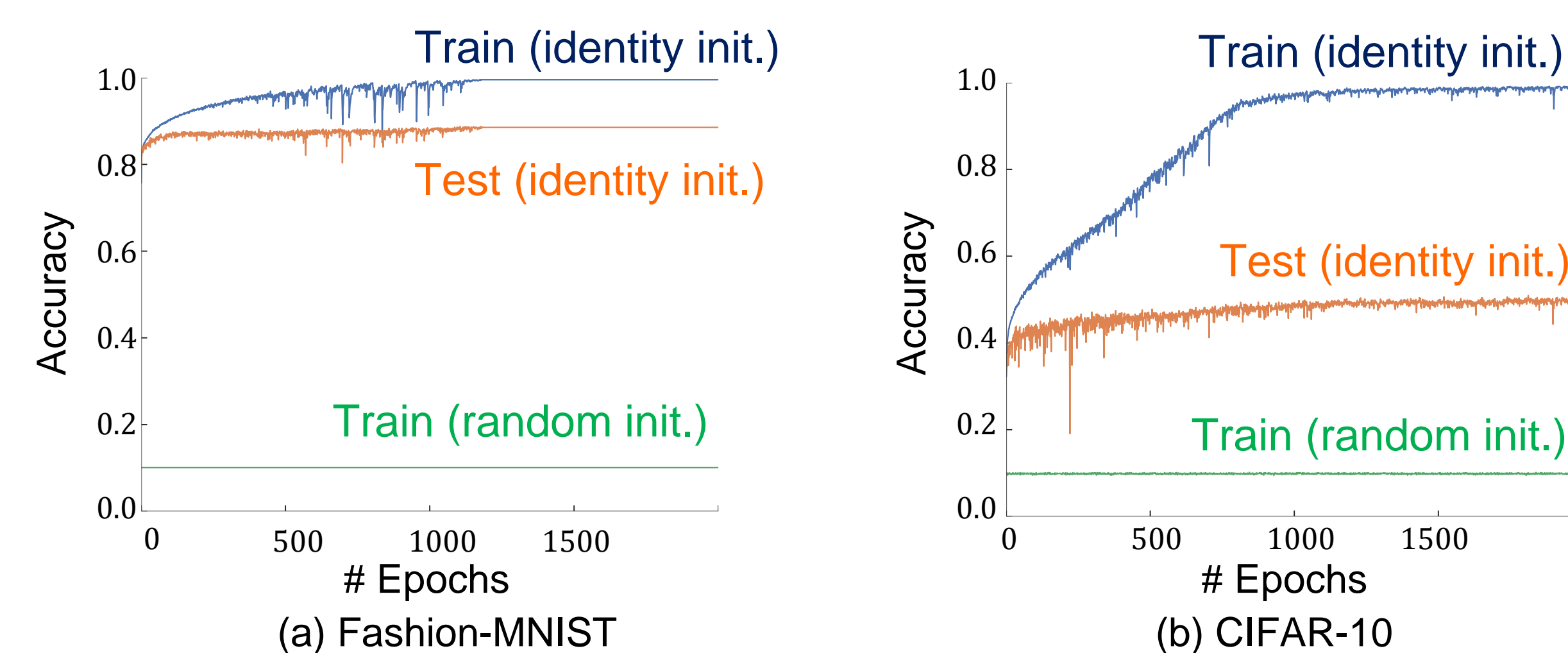
## Theoretical result

- The eigenvalue distribution is a Bernoulli function
- $$\mu_{JJ^T}(\lambda) = \begin{cases} \alpha^L \delta(\lambda - \sigma^{2L}) - (1 - \alpha^L) \delta(\lambda) & (\sigma > 1) \\ \alpha^1 \delta(\lambda - \sigma^{2L}) - (1 - \alpha^1) \delta(\lambda) & (\sigma \leq 1) \end{cases}$$
- $\alpha^\ell$ : Constant depending on  $\sigma$  (constant) and  $q^0$  (input variance)

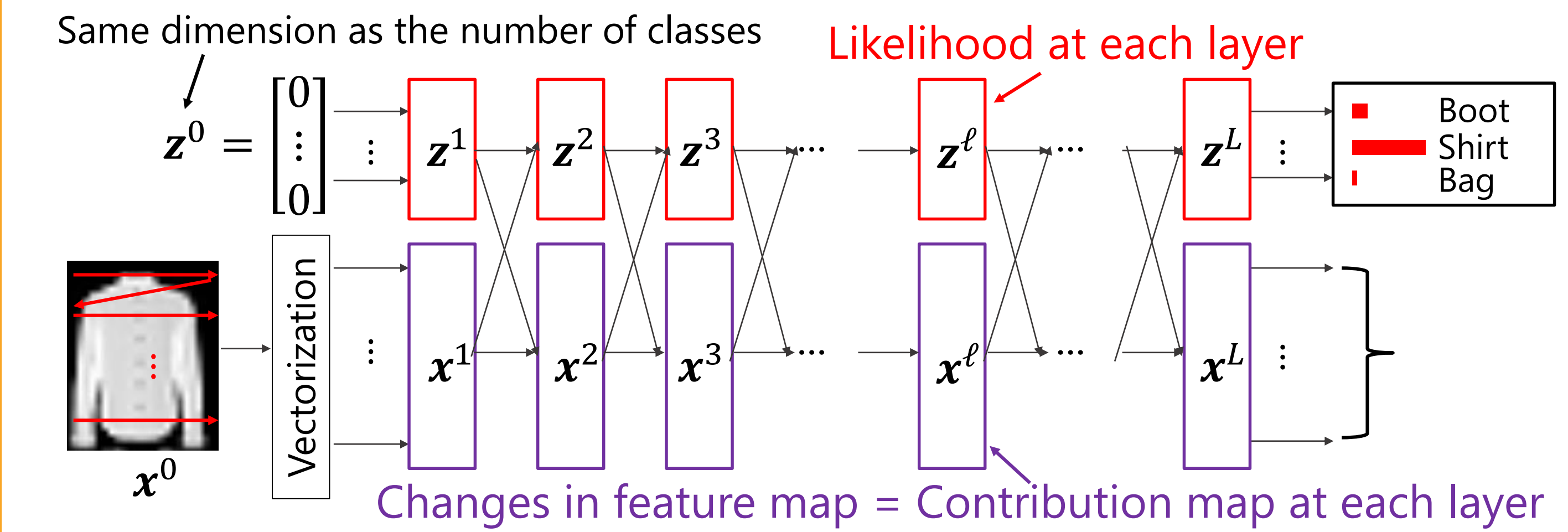
- We can prevent gradient vanishing by setting  $\sigma$  and  $q^0$  appropriately



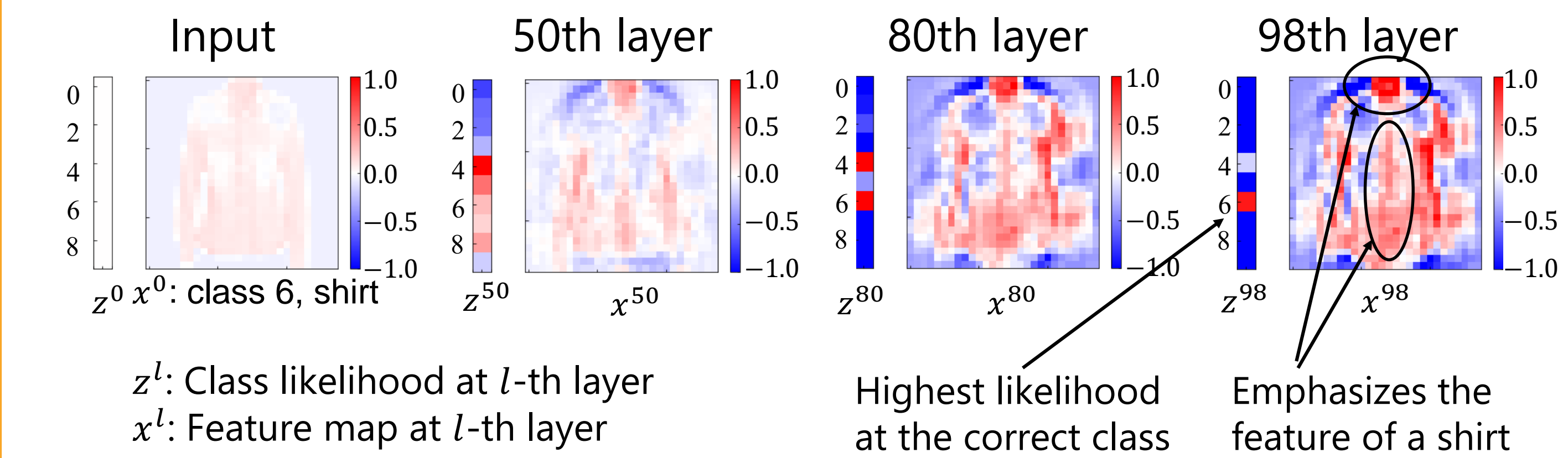
## Experimental Results



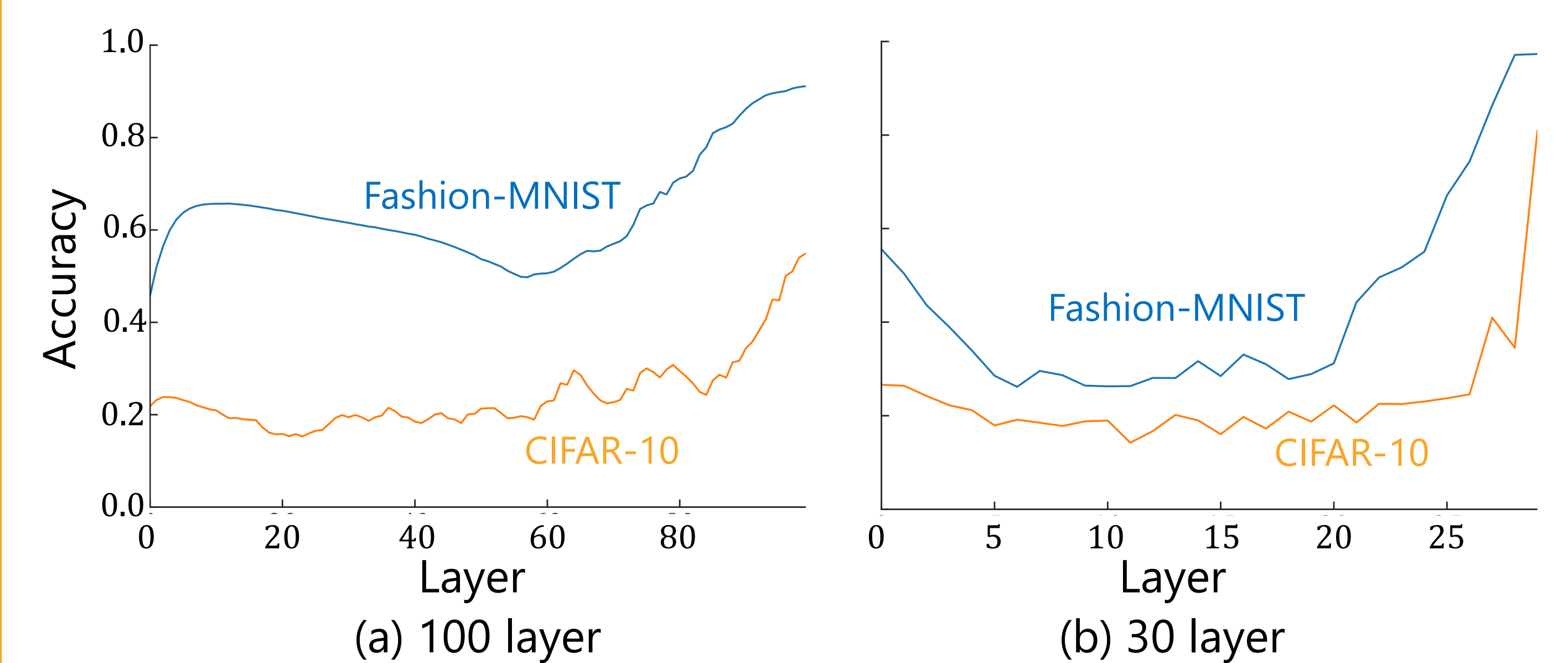
## Network structure that enhances interpretability



## Changes in feature map



## Classification accuracy for each 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