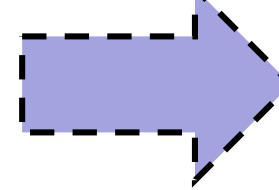


## Motivation:

### Deep Neural Networks

- Excellent performance
- Slow training speed:
  - VGG [Simonyan and Zisserman] on ILSVRC2014 with 16 conv layers takes approximately one month using 4 GPUs
- Vanishing gradient problem in back propagation (using the chain rule)
- Can we train DNNs faster?
  - Our network: 64 fully connected layers, 1024 nodes per layer, 1 GPU, training time ~3 hours
  - State-of-the-art for supervised hashing.

SOL



## Optimization:

- Introduce auxiliary variables:

$$\min_{\theta, \tilde{\theta}, W, \{z_{i,m}\}} \Omega(\theta, W) + \sum_i \ell(W^T z_{i,M}, y_i)$$

$$s. t. z_{i,0} = x_i, \tilde{\theta}_i^{(m)} = \theta^{(m)}, 1 \leq m \leq M, \forall i$$

$$z_{i,m} = f_m(z_{i,m-1}; \theta^{(m)}), 1 \leq m \leq M, \forall i$$

- Augmented Lagrangian (or Adaptive Direction

### Method of Multipliers):

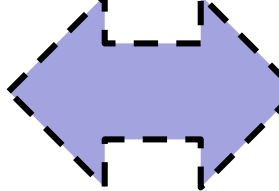
$$\min_{\theta, \tilde{\theta}, W, \{z_{i,m}\}, \{u_{i,m}\}} \Omega(\theta, W) + \sum_i \ell(W^T z_{i,M}, y_i)$$

$$+ \frac{\beta}{2} \sum_{i,m=1}^M \|z_{i,m} - f_m(z_{i,m-1}; \tilde{\theta}_i^{(m)}) + u_{i,m}\|_2^2$$

$$+ \frac{\gamma}{2} \sum_{i,m=1}^M \|\theta^{(m)} - \tilde{\theta}_i^{(m)} + v_{i,m}\|_2^2$$

$$s. t. z_{i,0} = x_i, \forall i$$

OPT



- Gradients w.r.t. variables

$$\frac{\partial g}{\partial \tilde{\theta}_i^{(m)}} = \frac{\partial \Omega}{\partial \tilde{\theta}_i^{(m)}} - \beta \sum_i (z_{i,m} - f_m(z_{i,m-1}; \tilde{\theta}_i^{(m)}) + u_{i,m}) \frac{\partial f_m}{\partial \tilde{\theta}_i^{(m)}} - \gamma \sum_i (\theta^{(m)} - \tilde{\theta}_i^{(m)} + v_{i,m})$$

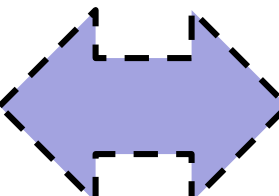
$$\frac{\partial g}{\partial z_{i,m}} = \frac{\partial \ell}{\partial z_{i,m}} + \beta (z_{i,m} - f_m(z_{i,m-1}; \tilde{\theta}_i^{(m)}) + u_{i,m})$$

$$\frac{\partial g}{\partial z_{i,m}} = -\beta \frac{\partial f_{m+1}}{\partial z_{i,m}} (z_{i,m+1} - f_{m+1}(z_{i,m}; \tilde{\theta}_i^{(m+1)}) + u_{i,m+1}) + \beta (z_{i,m} - f_m(z_{i,m-1}; \tilde{\theta}_i^{(m)}) + u_{i,m})$$

- **Supervised hashing**

- Goal: learn hash codes for linear classification [Shen et al.]
- General problem:

$$\min_{\theta, W, B} \Omega(\theta, W) + \sum_i \ell(W^T b_i, y_i), s. t., b_i = \text{sgn}(F(x_i, \theta)), \forall i$$



## Method

### Train DNN with optimization problem:

$$\min_{\theta, W} \Omega(\theta, W) + \sum_i \ell(W^T F_M(x_i), y_i)$$

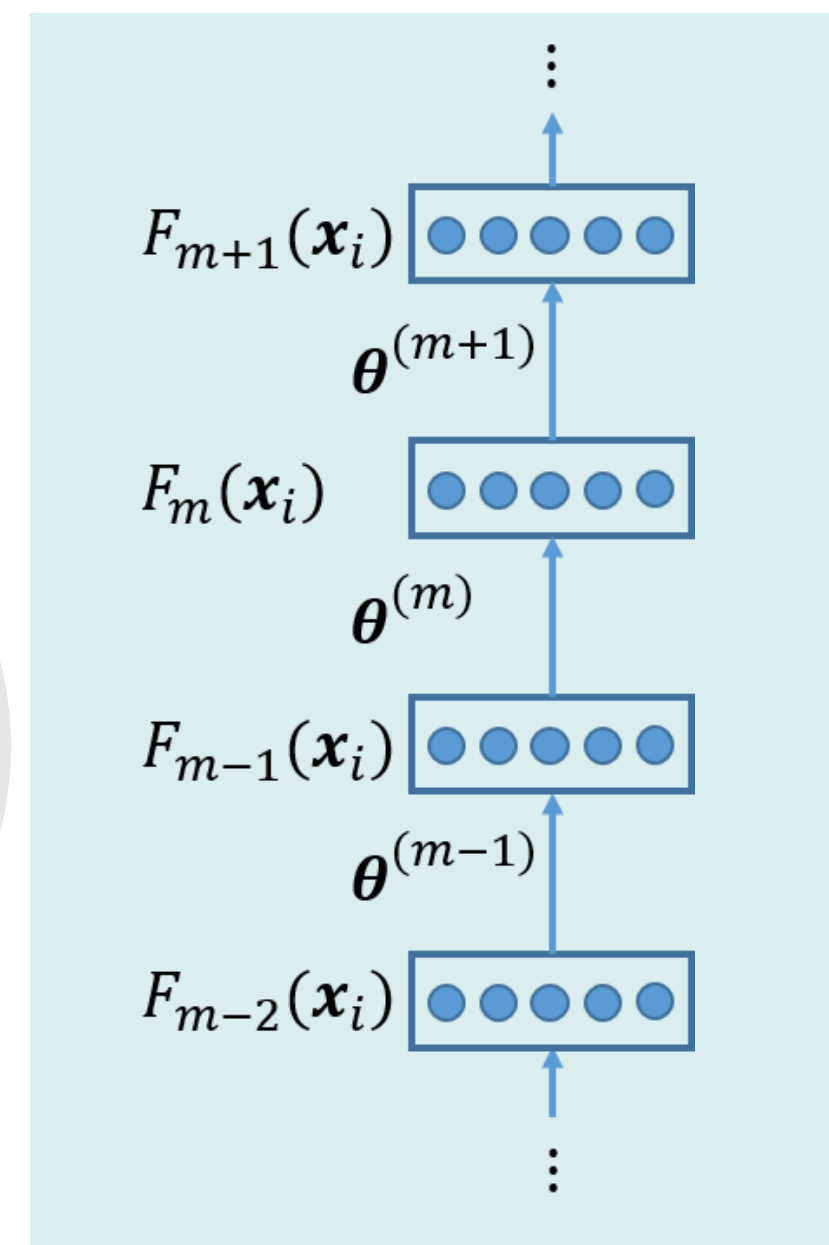
$$s. t. F_0(x_i) = x_i, F_m(x_i) = f_m(F_{m-1}(x_i); \theta^{(m)}), 1 \leq m \leq M, \forall i$$

$\theta = \{\theta^{(m)}\}_{m=1, \dots, M}$ : network weights

$\theta^{(m)} \in \mathbb{R}^{D_m \times D_{m-1}}$ : layer weights

$f_m: \mathbb{R}^{D_{m-1}} \rightarrow \mathbb{R}^{D_m}$ : activation function

e.g.  $f_m(x_i; \theta^{(m)}) = \max\{0, \theta^{(m)} x_i\}$



Decomposes the training process

- Break long term dependency:  $z_{i,m} = F_m(x_i), \forall i, m$

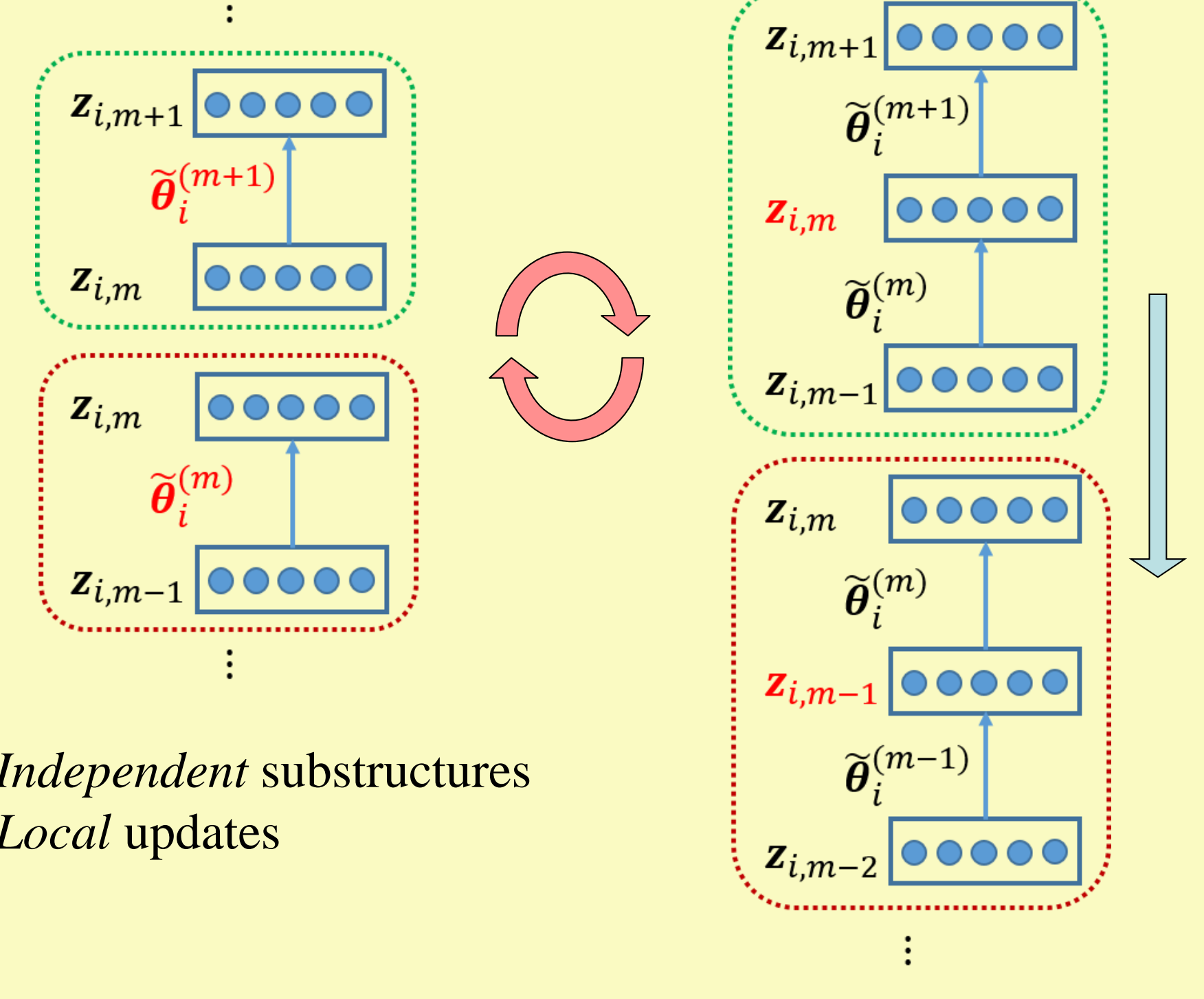
- Dependency between loss  $\ell$  and regularizer  $\Omega$ :

$$\tilde{\theta}_i^{(m)} = \theta^{(m)}, \forall i, \forall m$$

$\theta$  updates

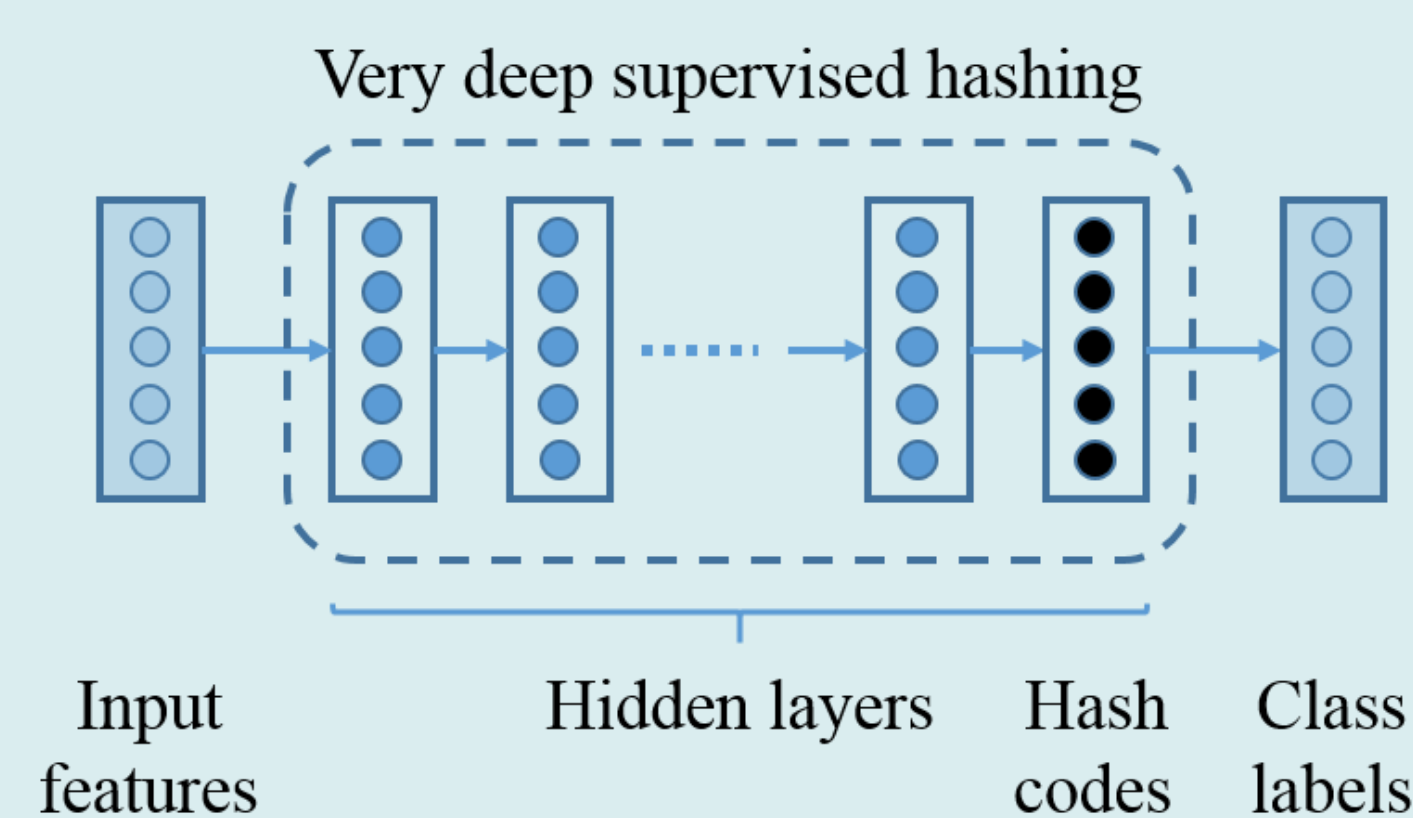
$z$  updates

### Alternate minimization



- ✓ Independent substructures
- ✓ Local updates

### Application: Hashing



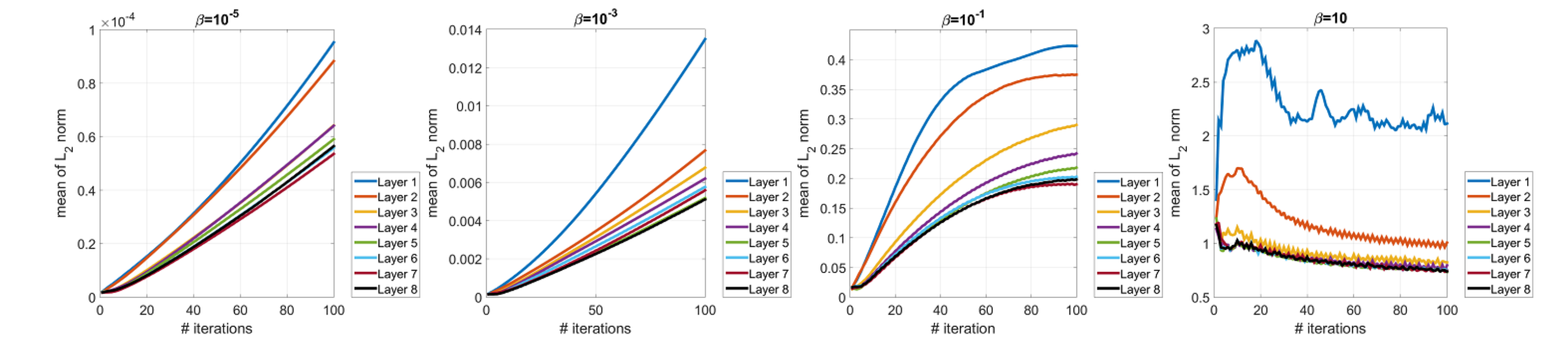
ANLYS



## Toy Example

- MNIST dataset: 8 layers, 64 nodes per layer

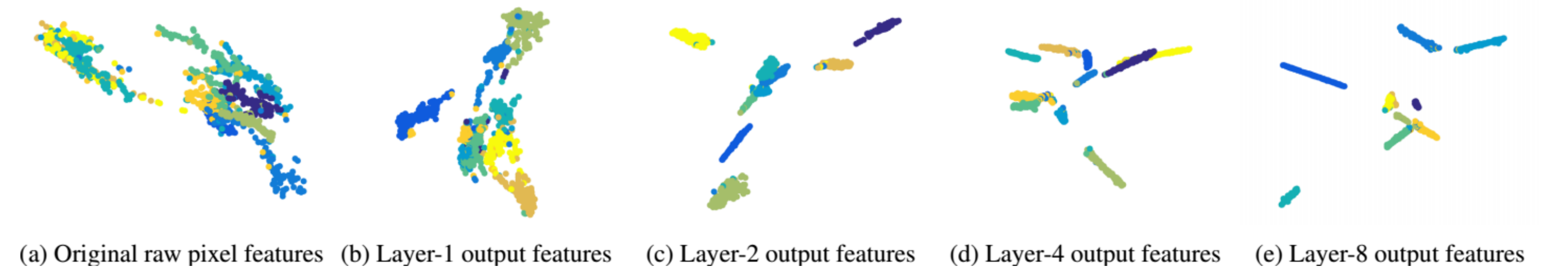
### Empirical convergence



### Computational complexity

- For two-layer substructure:  $O(D_m D_{m+1})$
- Total complexity for the network:  $O(\sum_{m=0}^M D_m D_{m+1} N)$

### Per layer visualization



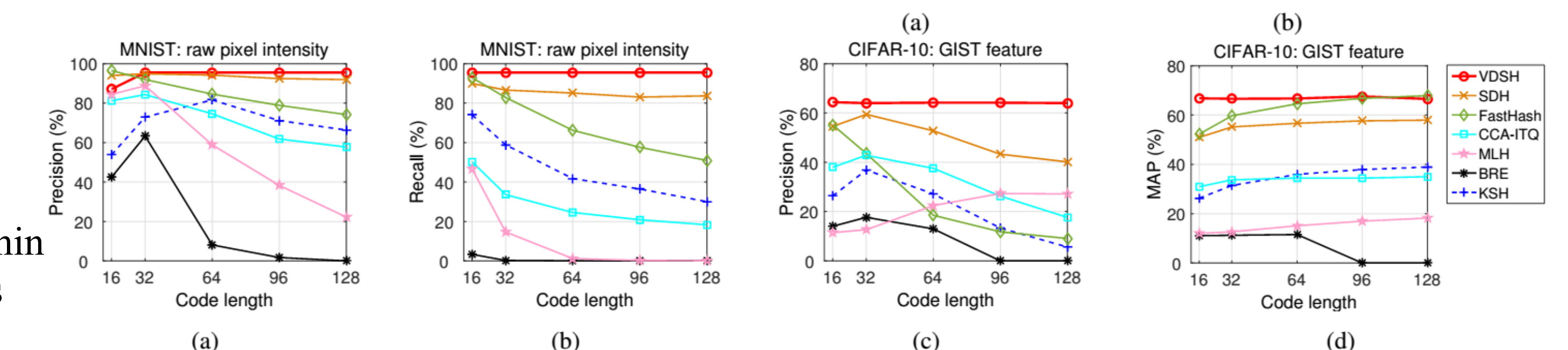
## Experiments on Image Retrieval

- Dataset

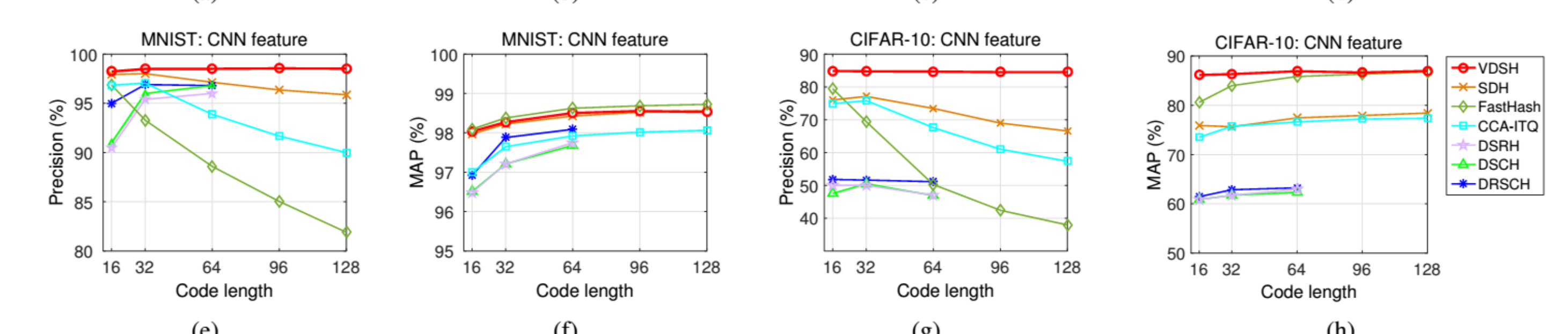
	Content	# Classes	# Images	No. Query
CIFAR-10	Natural	10	60k	100 * 10
MNIST	Digits	10	70k	100 * 10
NUS-WIDE	Natural	81	~270k	100 * 21

- Retrieval results

MNIST:  
# layers: 48  
# nodes: 256  
Training: 15 min  
Query: 6.6 ms



CIFAR-10:  
# layers: 16  
# nodes: 1024  
Training: 1 h  
Query: 4ms



- Hash code visualization (CIFAR-10, 64bit codes)

