

# Compressing Neural Networks using the Variational Information Bottleneck

**Bin Dai**<sup>1</sup>; Chen Zhu<sup>2</sup>; Baining Guo<sup>3</sup>; David Wipf<sup>3</sup>

<sup>1</sup>Tsinghua University, <sup>2</sup>University of Maryland, <sup>3</sup>Microsoft Research

# Neural Networks are Often Over-Parameterized

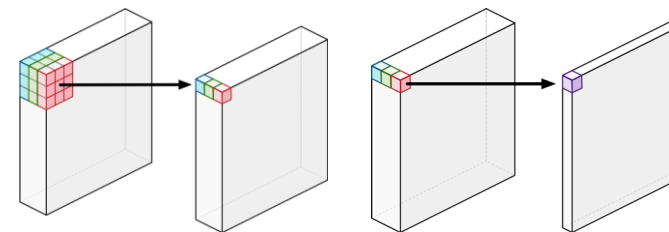
*(Denil et al., 2013)*

## Consequences:

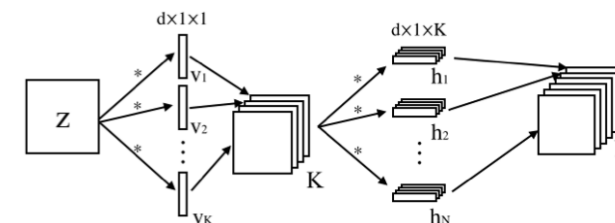
- Unnecessarily large model size
- Increased computational cost
- Extra run-time memory footprint

# Network Compression Methods

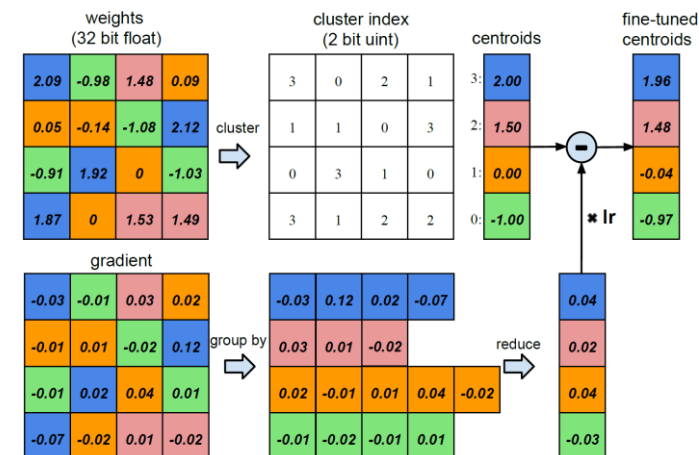
- Design more efficient network structure  
(Howard et al., 2017; Dong et al., 2017; Iandola et al., 2016)



- Tensor/matrix decomposition  
(Jaderberg et al., 2014; Zhang et al., 2016; Yu et al., 2017)

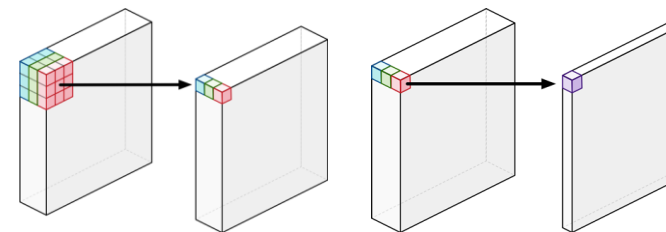


- Weight Quantization  
(Courbariaux et al., 2016; 2015; Han et al., 2015a; Mellempudi et al., 2017; Rastegari et al., 2016)

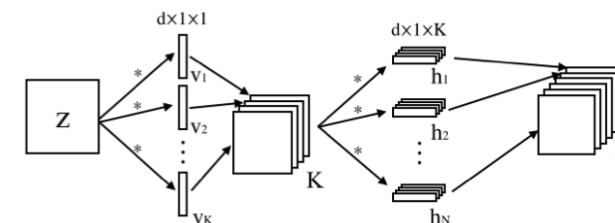


# Network Compression Methods

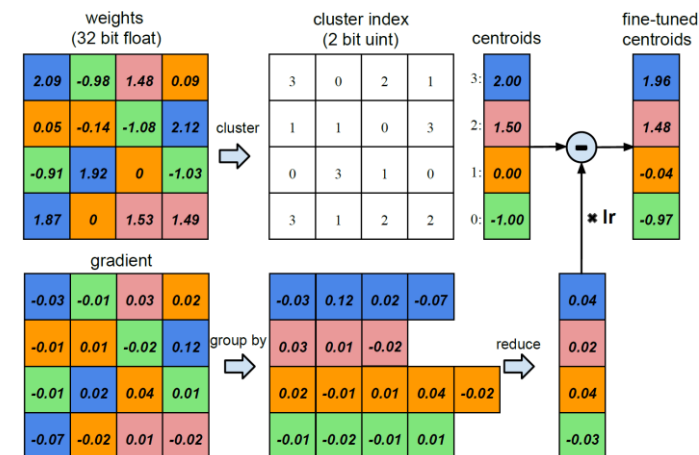
- Design more efficient network structure  
(Howard et al., 2017; Dong et al., 2017; Iandola et al., 2016)



- Tensor/matrix decomposition  
(Jaderberg et al., 2014; Zhang et al., 2016; Yu et al., 2017)



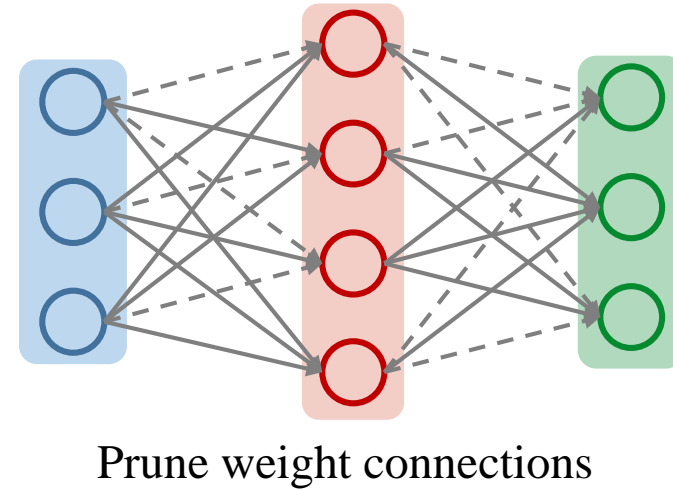
- Weight Quantization  
(Courbariaux et al., 2016; 2015; Han et al., 2015a; Mellempudi et al., 2017; Rastegari et al., 2016)



- Prune existing network structure ...

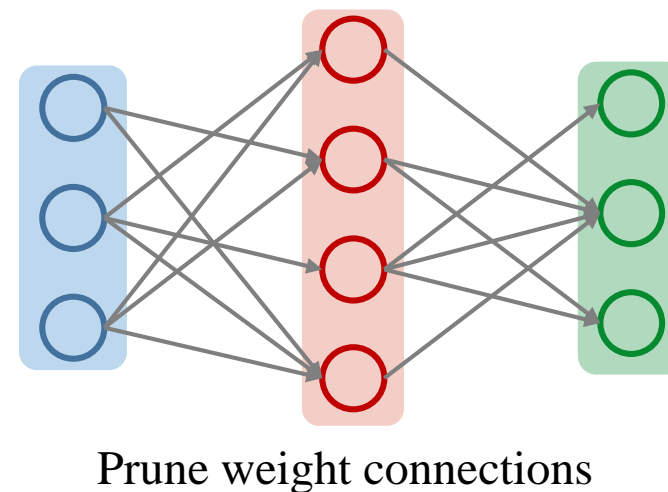
# Pruning Methods

- Prune weight connections  
(*Han et al., 2015b; Guo et al., 2016; LeCun et al., 1990*)
- Prune weight groups / activations
  - Group lasso  
(*Liu et al., 2017; Pan et al., 2016; Wen et al., 2016*)
  - Bayesian approach  
(*Louizos et al., 2017a; Neklyudov et al., 2017*)
  - Smoothed  $l_0$  approach  
(*Louizos et al., 2017b*)



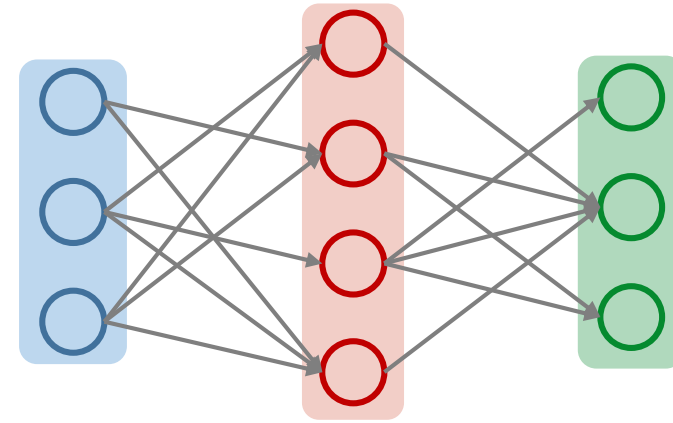
# Pruning Methods

- Prune weight connections  
(*Han et al., 2015b; Guo et al., 2016; LeCun et al., 1990*)
- Prune weight groups / activations
  - Group lasso  
(*Liu et al., 2017; Pan et al., 2016; Wen et al., 2016*)
  - Bayesian approach  
(*Louizos et al., 2017a; Neklyudov et al., 2017*)
  - Smoothed  $l_0$  approach  
(*Louizos et al., 2017b*)

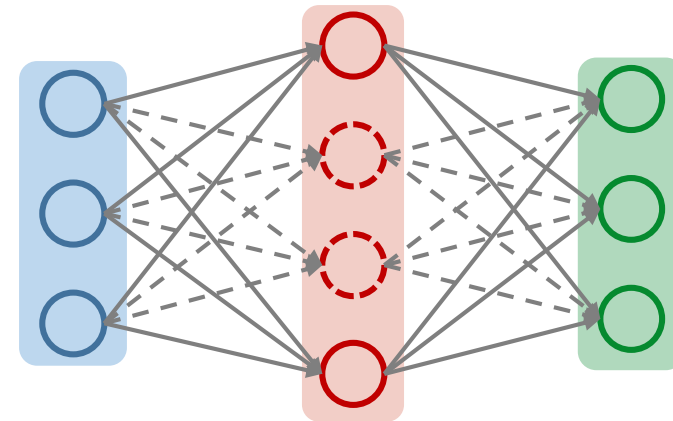


# Pruning Methods

- Prune weight connections  
(*Han et al., 2015b; Guo et al., 2016; LeCun et al., 1990*)
- Prune weight groups / activations
  - Group lasso  
(*Liu et al., 2017; Pan et al., 2016; Wen et al., 2016*)
  - Bayesian approach  
(*Louizos et al., 2017a; Neklyudov et al., 2017*)
  - Smoothed  $l_0$  approach  
(*Louizos et al., 2017b*)



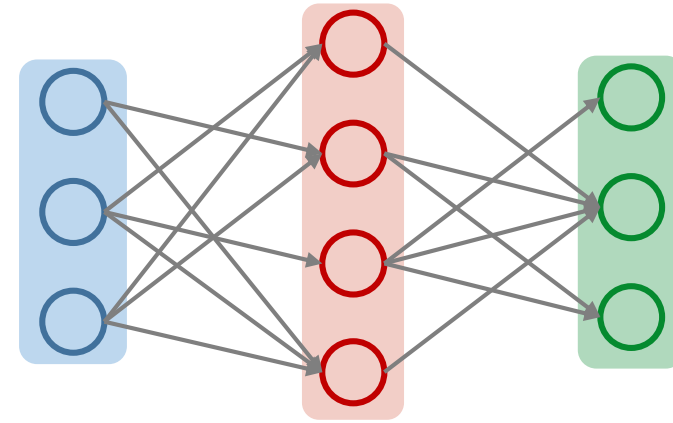
Prune weight connections



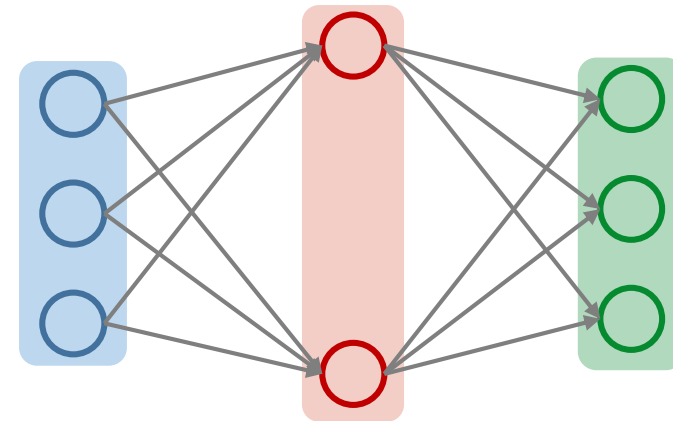
Prune weight groups / activations

# Pruning Methods

- Prune weight connections  
(*Han et al., 2015b; Guo et al., 2016; LeCun et al., 1990*)
- Prune weight groups / activations
  - Group lasso  
(*Liu et al., 2017; Pan et al., 2016; Wen et al., 2016*)
  - Bayesian approach  
(*Louizos et al., 2017a; Neklyudov et al., 2017*)
  - Smoothed  $l_0$  approach  
(*Louizos et al., 2017b*)



Prune weight connections



Prune weight groups / activations



# Markov Chain Interpretation of Network

$$y \rightarrow x(h_0) \rightarrow h_1 \rightarrow \cdots \rightarrow \underbrace{h_{i-1} \rightarrow h_i}_{\text{red oval}} \rightarrow \cdots \rightarrow h_L \rightarrow \hat{y} \quad (\text{Tishby \& Zaslavsky, 2015})$$

Defines  $p(h_i|h_{i-1})$

Degenerates to a Dirac-delta function if the network is deterministic

We consider non-degenerate distribution in our model

# Markov Chain Interpretation of Network

$$y \rightarrow x(h_0) \rightarrow h_1 \rightarrow \cdots \rightarrow h_{i-1} \rightarrow h_i \rightarrow \cdots \rightarrow h_L \rightarrow \hat{y} \quad (\text{Tishby \& Zaslavsky, 2015})$$

Approximates  $p(y|h_L)$  via some tractable alternative  $p(\hat{y} | h_L)$

# Intuition

$$y \rightarrow x(h_0) \rightarrow h_1 \rightarrow \dots \rightarrow h_{i-1} \rightarrow h_i \rightarrow \dots \rightarrow h_L \rightarrow \hat{y} \quad (\text{Tishby \& Zaslavsky, 2015})$$



- Maximize the mutual information  $I(h_i; y)$  between  $h_i$  and  $y$

For high-accuracy prediction

- Minimize the mutual information  $I(h_i; h_{i-1})$  between  $h_i$  and  $h_{i-1}$

For compression

Information  
Bottleneck

# Layer-wise Energy

- Minimize the mutual information  $I(h_i; h_{i-1})$  between  $h_i$  and  $h_{i-1}$
- Maximize the mutual information  $I(h_i; y)$  between  $h_i$  and  $y$

$$\mathcal{L}_i = \gamma_i I(h_i; h_{i-1}) - I(h_i; y)$$

$\gamma_i$  determines the strength of the bottleneck

# Upper Bound

$$\mathcal{L}_i = \gamma_i I(h_i; h_{i-1}) - I(h_i; y)$$

$$\gamma_i \mathbb{E}_{\{x,y\} \sim \mathcal{D}, h_{1:i-1} \sim p(h_{1:i-1}|x)} [\mathbb{KL}[p(h_i|h_{i-1}) || q(h_i)]] \geq \gamma_i I(h_i; h_{i-1})$$

$$- \mathbb{E}_{\{x,y\} \sim \mathcal{D}, h \sim p(h|x)} [\log q(y|h_L)] \geq -I(h_i; y)$$

$$\tilde{\mathcal{L}}_i = \gamma_i \mathbb{E}_{\{x,y\} \sim \mathcal{D}, h_{1:i-1} \sim p(h_{1:i-1}|x)} [\mathbb{KL}[p(h_i|h_{i-1}) || q(h_i)]]$$

$$- \mathbb{E}_{\{x,y\} \sim \mathcal{D}, h \sim p(h|x)} [\log q(y|h_L)]$$

$\mathcal{D}$ : data distribution

$q(h_i)$ : variational approximation of  $p(h_i)$

$h_{1:i}: \{h_j\}_{j=1}^i$

$q(y|h_L)$ : variational approximation of  $p(y|h_L)$

$h: h_{1:L}$

# Parameterization of $q(y|h_L)$

$$\begin{aligned}\tilde{\mathcal{L}}_i &= \gamma_i \mathbb{E}_{\{x,y\} \sim \mathcal{D}, h_{1:i-1} \sim p(h_{1:i-1}|x)} [\mathbb{KL}[p(h_i|h_{i-1})||q(h_i)]] \\ &\quad - \mathbb{E}_{\{x,y\} \sim \mathcal{D}, h \sim p(h|x)} [\log q(y|h_L)]\end{aligned}$$

$q(y|h_L)$ : multinomial distribution for classification task

Gaussian distribution for regression task

# Parameterization of $q(h_i)$

$$\begin{aligned}\tilde{\mathcal{L}}_i = & \gamma_i \mathbb{E}_{\{x,y\} \sim \mathcal{D}, h_{1:i-1} \sim p(h_{1:i-1}|x)} [\mathbb{KL}[p(h_i|h_{i-1}) || q(h_i)]] \\ & - \mathbb{E}_{\{x,y\} \sim \mathcal{D}, h \sim p(h|x)} [\log q(y|h_L)]\end{aligned}$$

$$q(h_i) = \mathcal{N}(h_i; 0, \text{diag}[\xi_i])$$

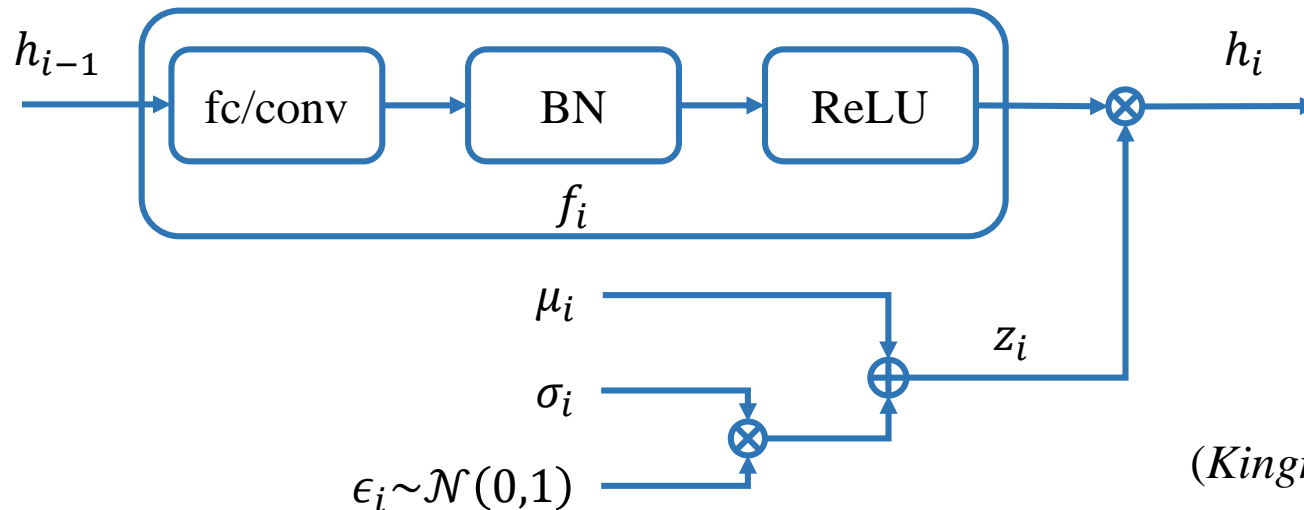
$\xi_i$ : unknown vector of variances learned from data

*(Tipping 2001)*

# Parameterization of $p(h_i|h_{i-1})$

$$\tilde{\mathcal{L}}_i = \gamma_i \mathbb{E}_{\{x,y\} \sim \mathcal{D}, h_{1:i-1} \sim p(h_{1:i-1}|x)} [\mathbb{KL}[p(h_i|h_{i-1}) || q(h_i)]] \\ - \mathbb{E}_{\{x,y\} \sim \mathcal{D}, h \sim p(h|x)} [\log q(y|h_L)]$$

$$p(h_i|h_{i-1}) = \mathcal{N}(h_i; f_i(h_{i-1}) \odot \mu_i, \text{diag}[f_i(h_{i-1})^2 \odot \sigma_i^2])$$



$$z_i = \mu_i + \epsilon_i \odot \sigma_i$$

Reparameterization trick

(Kingma and Welling, 2014; Rezende et al., 2014)



# Final Objective Function

$$\inf_{\xi_i > 0} 2 \mathbb{E}_{\{x,y\} \sim \mathcal{D}, h_{1:i-1} \sim p(h_{1:i-1}|x)} [\mathbb{KL}[p(h_i|h_{i-1})||q(h_i)]]$$
$$\equiv \sum_{j=1}^{r_i} \left[ \log \left( 1 + \frac{\mu_{i,j}^2}{\sigma_{i,j}^2} \right) + \psi_{i,j} \right]$$

$$\psi_{i,j} \triangleq \log \mathbb{E}_{h_{i-1} \sim p(h_{i-1})} [f_{i,j}(h_{i-1})^2] - \mathbb{E}_{h_{i-1} \sim p(h_{i-1})} [\log f_{i,j}(h_{i-1})^2]$$

# Final Objective Function

$$\tilde{\mathcal{L}} = \underbrace{\sum_{i=1}^L \gamma_i \sum_{j=1}^{r_i} \log \left( 1 + \frac{\mu_{i,j}^2}{\sigma_{i,j}^2} \right)}_{\text{Regularization term}} - \underbrace{L \mathbb{E}_{\{x,y\} \sim \mathcal{D}, h \sim p(h|x)} [\log q(y|h_L)]}_{\text{Data-fit term}} \quad (*)$$

Regularization term

Data-fit term

$r_i$ : Number of neurons in the  $i$ -th layer.

Define  $\alpha_{i,j} = \frac{\mu_{i,j}^2}{\sigma_{i,j}^2}$ .

# Final Objective Function

$$\tilde{\mathcal{L}} = \underbrace{\sum_{i=1}^L \gamma_i \sum_{j=1}^{r_i} \log \left( 1 + \frac{\mu_{i,j}^2}{\sigma_{i,j}^2} \right)}_{\text{Regularization term}} - \underbrace{L \mathbb{E}_{\{x,y\} \sim \mathcal{D}, h \sim p(h|x)} [\log q(y|h_L)]}_{\text{Data-fit term}} \quad (*)$$

Regularization term

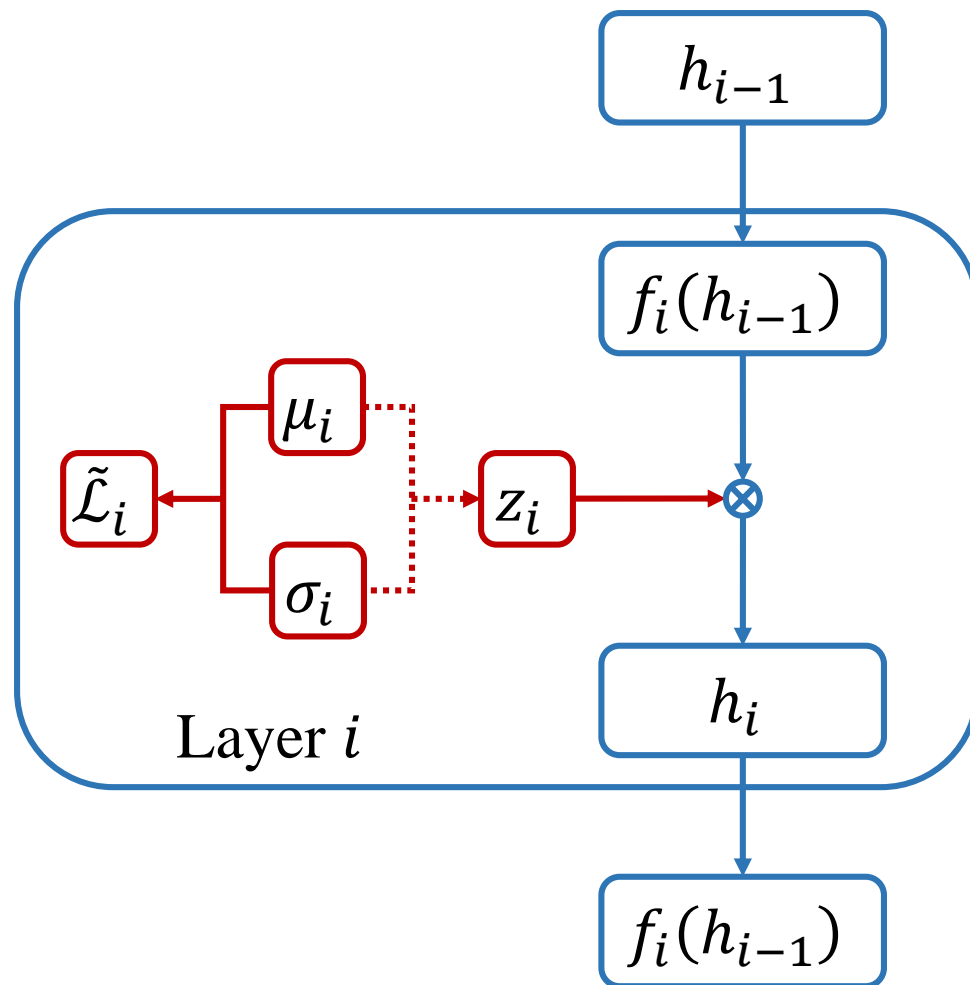
Data-fit term

$r_i$ : Number of neurons in the  $i$ -th layer.

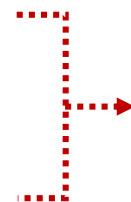
Define  $\alpha_{i,j} = \frac{\mu_{i,j}^2}{\sigma_{i,j}^2}$ .

$$\tilde{\mathcal{L}} = \sum_{i=1}^L \gamma_i \sum_{j=1}^{r_i} \log(1 + \alpha_{i,j}) - L \mathbb{E}_{\{x,y\} \sim \mathcal{D}, h \sim p(h|x)} [\log q(y|h_L)] \quad (*)$$

# Network Structure



Element-wise product



Sample from Gaussian distribution

Variational Information Bottleneck Network  
(VIBNet)

# Reduce Redundancy via Intrinsic Sparsity

$$\tilde{\mathcal{L}} = \underbrace{\sum_{i=1}^L \gamma_i \sum_{j=1}^{r_i} \log(1 + \alpha_{i,j})}_{\text{Regularization term}} - \underbrace{L \mathbb{E}_{\{x,y\} \sim \mathcal{D}, h \sim p(h|x)} [\log q(y|h_L)]}_{\text{Data-fit term}} \quad (*)$$

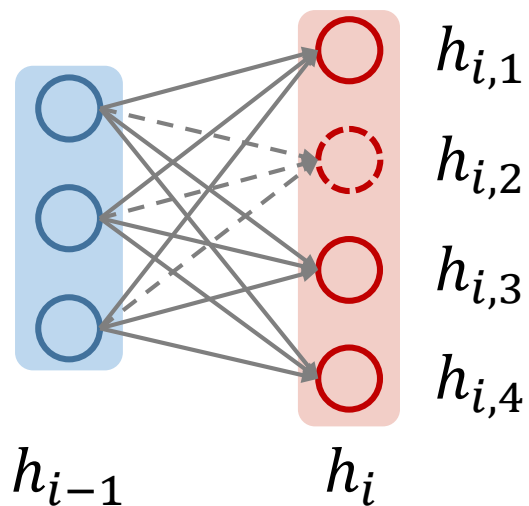
$\log(1 + \alpha_{i,j})$ : concave non-decreasing function on  $[0, \infty)$

- ✓ Push some  $\alpha_{i,j}$  to exactly 0 and leaving others mostly unchanged
- ✗ Push all  $\alpha_{i,j}$  towards smaller values

(Chen et al., 2017)

# Relationship Between $\alpha_{i,j}$ and $I(h_{i,j}, h_{i-1})$

Proposition 1. At any minimum of (\*)

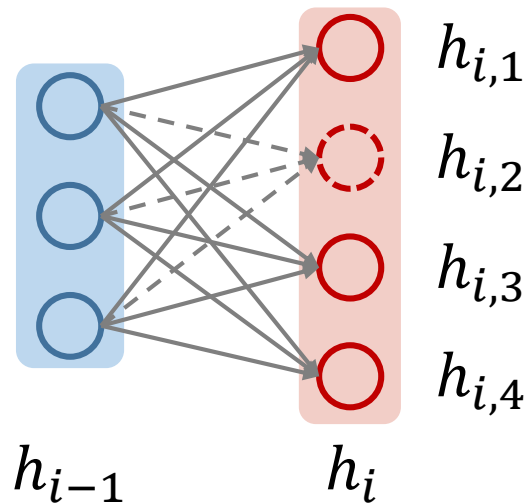
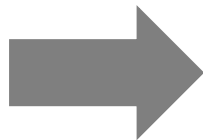


$$\alpha_{i,2} = 0$$

# Relationship Between $\alpha_{i,j}$ and $I(h_{i,j}, h_{i-1})$

Proposition 1. At any minimum of  $(*)$ ...

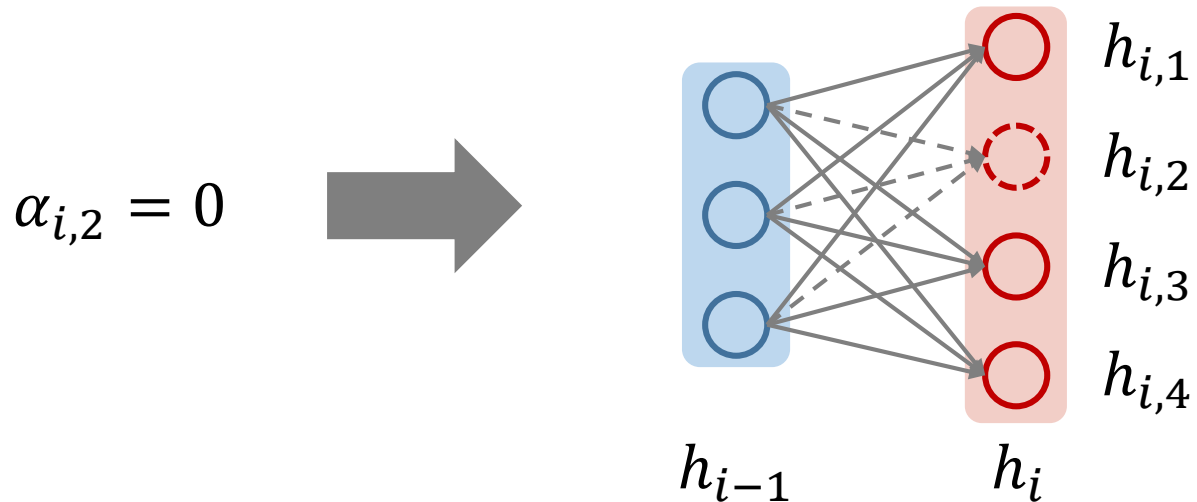
$$\alpha_{i,2} = 0$$



$$I(h_{i,2}, h_{i-1}) < \psi_{i,2}$$

# Relationship Between $\alpha_{i,j}$ and $I(h_{i,j}, h_{i-1})$

Proposition 1. At any minimum of (\*)

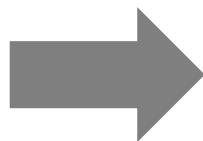
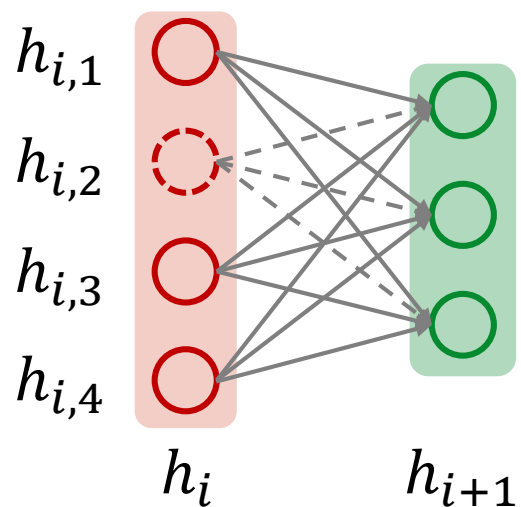


$$I(h_{i,2}, h_{i-1}) = 0 \quad \text{With added assumptions.}$$



# Relationship Between $\alpha_{i,j}$ and $W_{i+1,\cdot j}$

Proposition 2. At any minimum of (\*)

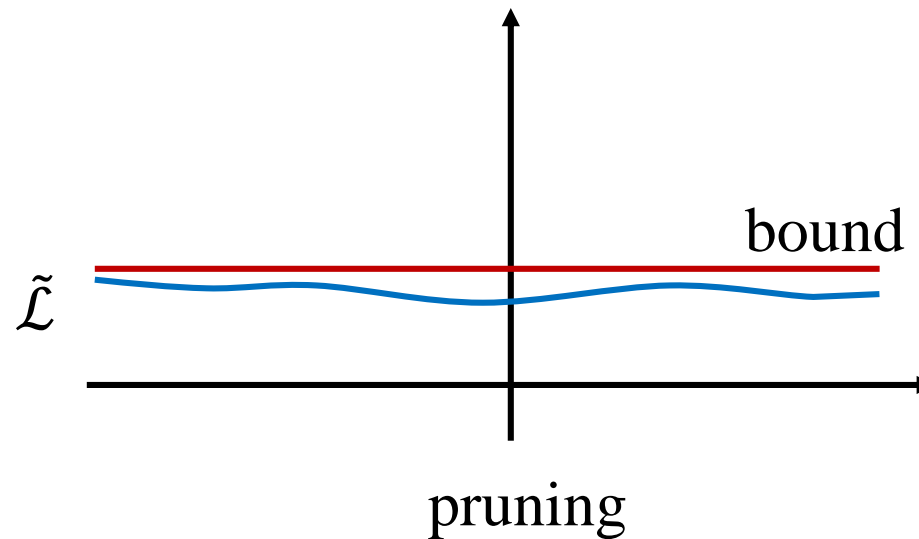
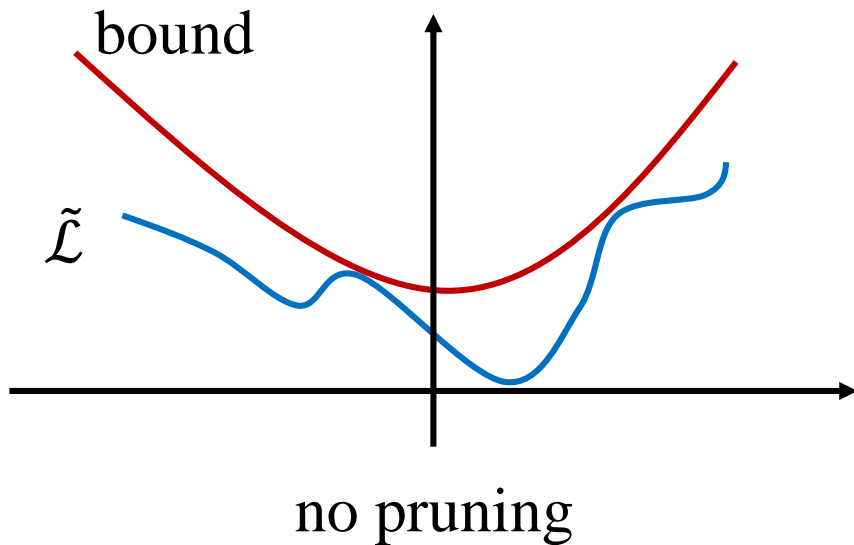


$$\alpha_{i,2} = 0$$

$$W_{i+1,\cdot 2} = 0$$

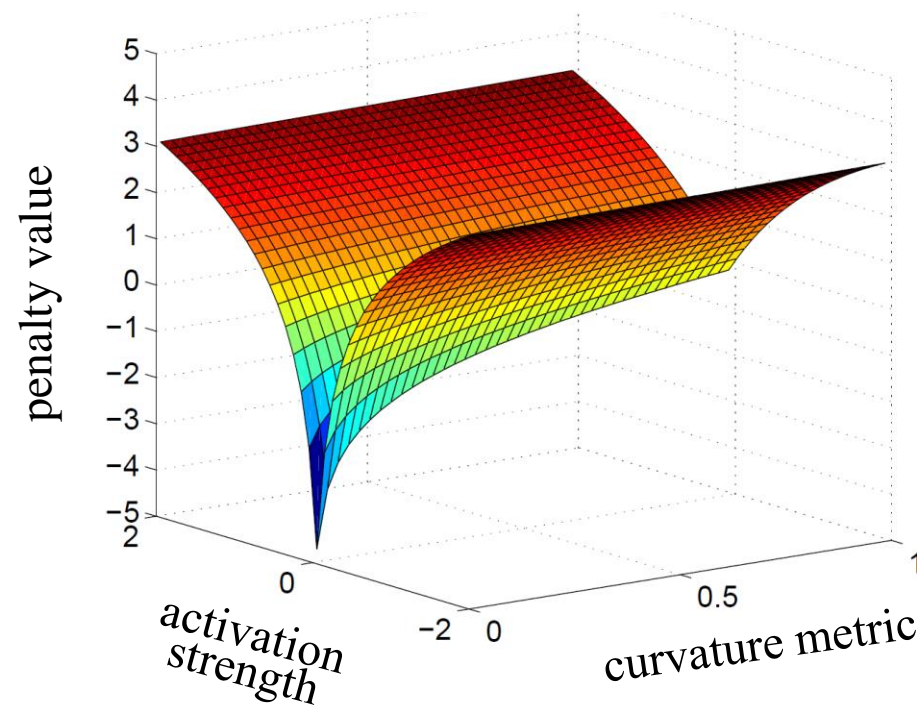
# Analysis of Tractable Upper Bound

At any minimum, redundant activations will provably be pruned



# Analysis of Tractable Upper Bound

Effective sparsity penalty shape is adaptive relative to the curvature of data fitting terms:

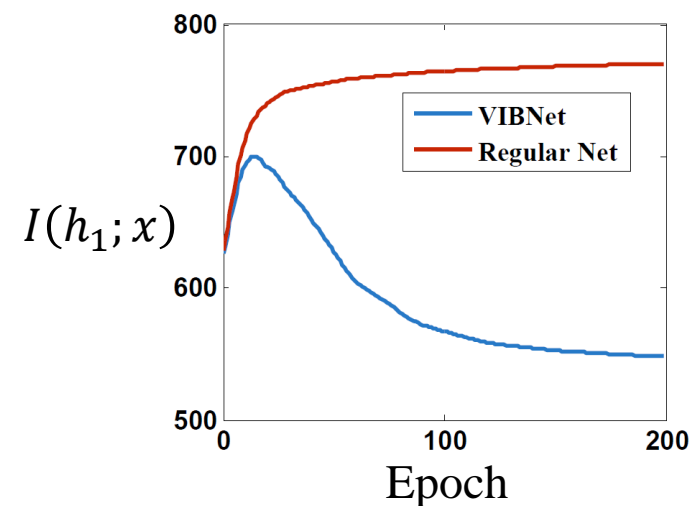
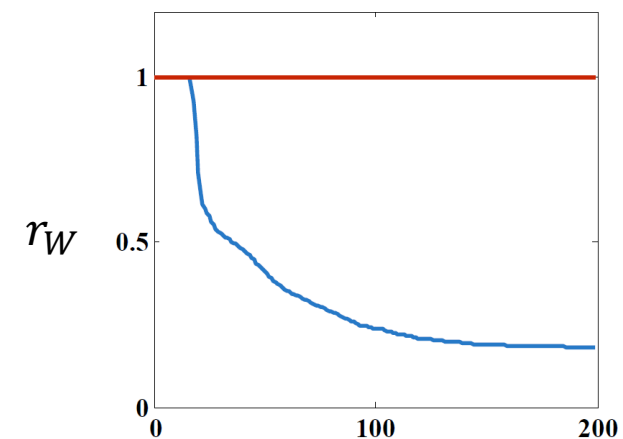


# Training and Testing

- Training
  - Bottleneck parameters set with simple heuristic to roughly match accuracy of existing methods
- Pruning
  - Small values of  $\alpha_{i,j}$  are set to zero (SGD will not push to exactly zero)
- Testing
  - Use the mean value of  $p(h_i|h_{i-1})$  rather than sampling
  - Use multiple samples can further improve the accuracy (at the cost of computation time) *(Louizos et al., 2017a)*
  - If desired, can finetune the pruned model to further boost the accuracy

# LeNet300-100 on MNIST

Method	$r_W(\%)$	$r_N(\%)$	Error(%)	Pruned Model
VD ( <i>Molchanov et al., 2017</i> )	25.28	58.95	1.8	512-114-72
BC-GNJ ( <i>Louizos et al., 2017a</i> )	10.76	32.85	1.8	278-98-13
BC-GHS ( <i>Louizos et al., 2017a</i> )	10.55	34.71	1.8	311-86-14
L0 ( <i>Louizos et al., 2017b</i> )	26.02	45.02	<b>1.4</b>	219-214-100
L0-sep ( <i>Louizos et al., 2017b</i> )	10.01	32.69	1.8	266-88-33
DN ( <i>Pan et al., 2016</i> )	23.05	57.94	1.8	542-83-61
<b>VIBNet</b>	<b>3.59</b>	<b>16.98</b>	<b>1.6</b>	<b>97-71-33</b>



$$r_W = \frac{\# \text{ params left}}{\# \text{ params total}}$$

$$r_N = \frac{\text{Memory footprint of the pruned model}}{\text{Memory footprint of the original model}}$$

# VGG16 on CIFAR10

Method	$r_W(\%)$	$r_N(\%)$	Error(%)	FLOP(Mil)
BC-GNJ ( <i>Louizos et al., 2017a</i> )	6.57	81.68	8.6	141.5
BC-GHS ( <i>Louizos et al., 2017a</i> )	5.40	74.82	9.0	121.9
VIBNet	<b>5.30</b>	<b>49.57</b>	<b>8.8 (8.5)</b>	<b>70.63</b>
PF ( <i>Li et al., 2017b</i> )	35.99	83.97	6.6	206.3
SBP ( <i>Neklyudov et al., 2017</i> )	7.01	80.72	7.5	136.0
SBPa ( <i>Neklyudov et al., 2017</i> )	5.78	66.46	9.0	99.20
VIBNet	<b>5.45</b>	<b>57.86</b>	<b>6.5 (6.1)</b>	<b>86.82</b>
NS-Single ( <i>Liu et al., 2017</i> )	11.50	-	6.2	195.5
NS-Best ( <i>Liu et al., 2017</i> )	8.60	-	5.9	147.0
VIBNet	<b>5.79</b>	<b>59.60</b>	<b>6.2 (5.8)</b>	<b>116.0</b>

Previous works used three different modifications of VGG16 on CIFAR10.

# VGG16 on CIFAR100

Method	$r_W(\%)$	$r_N(\%)$	Error(%)	FLOP(Mil)
RNP ( <i>Lin et al., 2017</i> )	-	-	38.0	160
VIBNet	<b>22.75</b>	<b>59.80</b>	37.6 (37.4)	<b>133.6</b>
NS-Single ( <i>Liu et al., 2017</i> )	24.90	-	26.5	250.5
NS-Best ( <i>Liu et al., 2017</i> )	20.80	-	26.0	214.8
VIBNet	<b>15.08</b>	<b>73.80</b>	25.9 (25.7)	<b>203.1</b>

# Conclusion

- We proposed a network compression model inspired by the information bottleneck principle
- Theoretical analysis shows that the objective tends to accumulate useful information in a sparse set of neurons
- The adaptive sparsity penalty that emerges from our model has advantages over traditional fixed sparsity penalties
- Empirical results show that our model can produce better performance than previous compression models



# Q & A



<https://github.com/zhuchen03/VIBNet>

Poster #128