

# 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

## ABSTRACT

Neural networks can be compressed to reduce memory and computational requirements, or to increase accuracy by facilitating the use of a larger base architecture. In this work, we use the information bottleneck principle to inspire a novel energy deep function for compressing deep networks. Theoretical analysis supports the compression performance while empirical results show state-of-art compression rates across an array of datasets and network architectures.

## CONTACT

Bin Dai:  
[daib13@mails.tsinghua.edu.cn](mailto:daib13@mails.tsinghua.edu.cn)  
 Chen Zhu:  
[chenzhu@umd.edu](mailto:chenzhu@umd.edu)  
 David Wipf:  
[davidwip@microsoft.com](mailto:davidwip@microsoft.com)

## Code Available



## Background

Neural networks are over-parameterized in terms of

- Number of parameters ( $r_W$ )
- Computational cost (FLOP)
- Memory footprint ( $r_N$ )

Previous compression strategies:

- Design more efficient network structures
- Quantize network weights
- Apply tensor/matrix decompositions
- Prune existing network structures, e.g.
  - Connections
  - Weight groups / activations
  - Bayesian approaches
  - Group Lasso
  - Smoothed  $l_0$ -norm approach

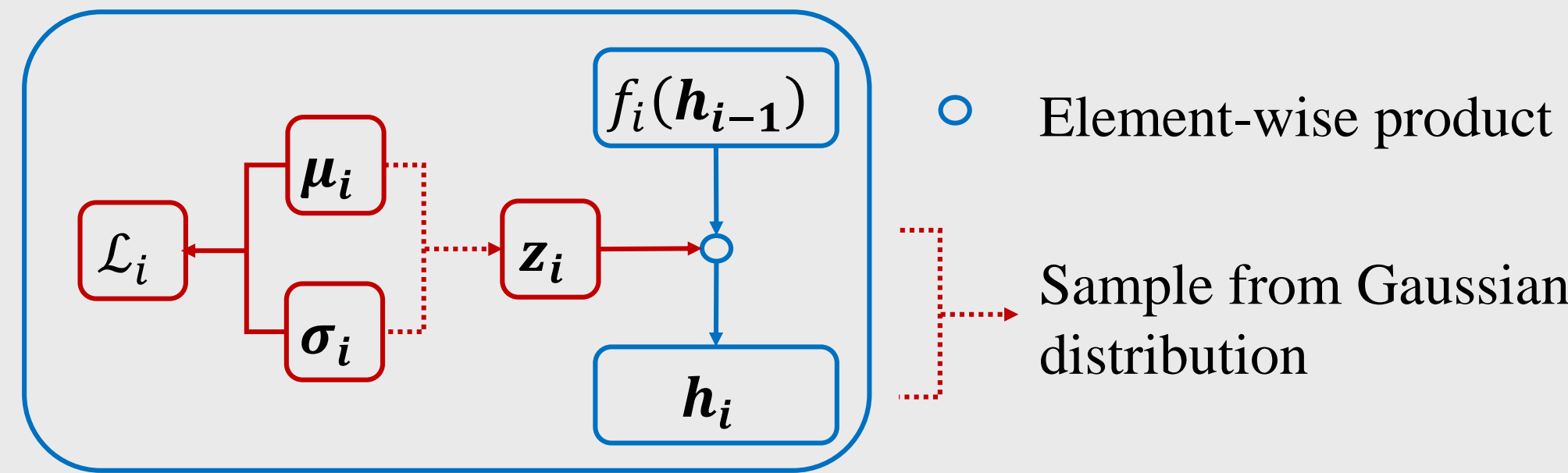


Figure 1. VIBNet layer/block structure

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

Table 1. LeNet300-100 on MNIST

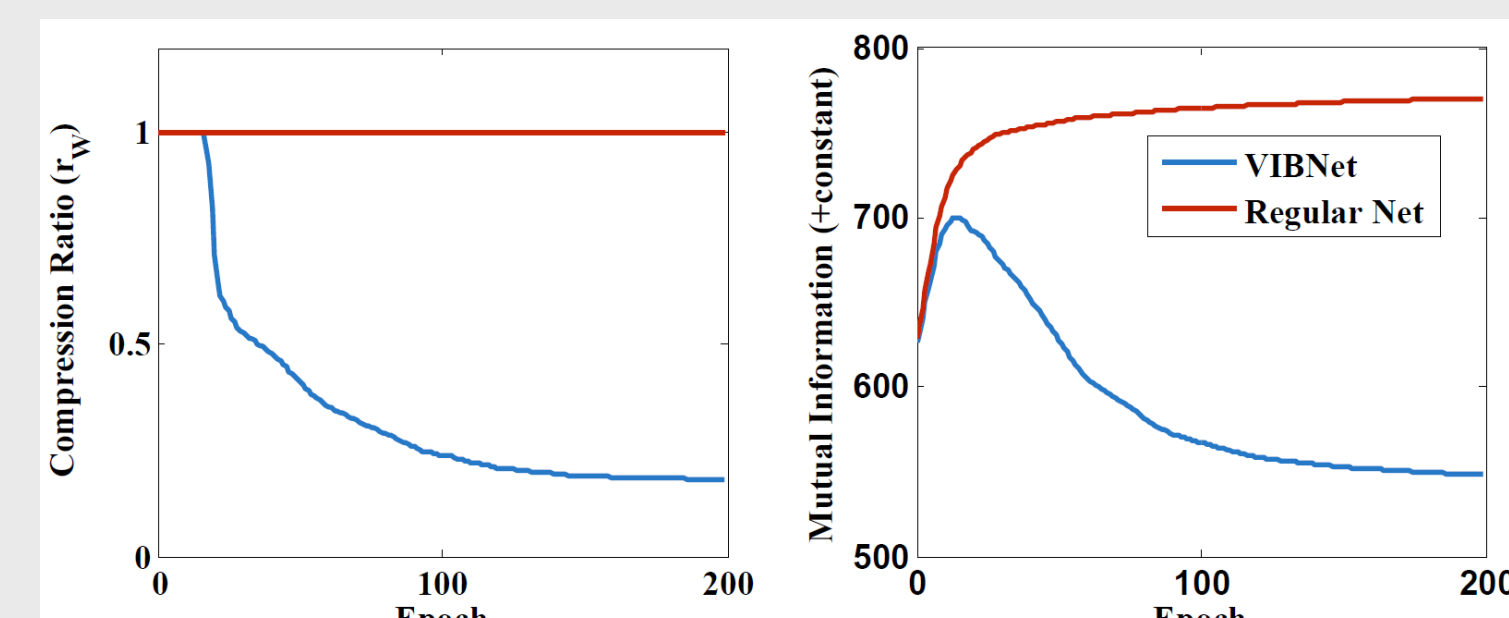


Figure 2. Compression ratio and mutual information between the first hidden layer and the input layer.

## Method

We interpret a feedforward network as a Markov chain (All the layers are stochastic):

$$y \rightarrow x \rightarrow h_1 \rightarrow \dots \rightarrow h_L \rightarrow \hat{y}$$

❖ Intuition:

Minimize the mutual information  $I(h_i; h_{i-1})$  to remove inter-layer redundancy (*information bottleneck*)  
 Maximize the mutual information  $I(h_i; y)$  to encourage accurate predictions of  $y$

(Tishby & Zaslavsky, arXiv 2015)

❖ Layer-wise objective:  $\mathcal{L}_i = \gamma_i I(h_i; h_{i-1}) - I(h_i; y)$

❖ Variational upper bound:

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

❖ Parametric forms of  $p(h_i|h_{i-1})$ ,  $q(h_i)$  and  $q(y|h_L)$

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

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

$q(y|h_L)$  is multinomial for classification and Gaussian for regression.

❖ Final variational IB-inspired objective function (VIBNet):

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

Regularization term

Data-fit term

Method	$r_W$ (%)	$r_N$ (%)	FLOP(mil)	Error(%)
GD	1.38	32.00	0.250	1.1
GL	23.69	19.35	0.201	1.0
VD	9.29	60.78	0.660	1.0
SBP	19.66	21.15	0.213	<b>0.9</b>
BC-GNJ	0.95	35.03	0.283	1.0
BC-GHS	<b>0.64</b>	22.80	0.153	1.0
L0	8.92	85.82	1.113	<b>0.9</b>
L0-sep	1.08	40.36	0.389	1.0
VIBNet	<b>0.83</b>	<b>15.55</b>	<b>0.094</b>	<b>1.0</b>

Table 2. LeNet5-Caffe on MNIST

VD: Molchanov et al., ICML 2017  
 BC-GNJ / BC-GHS: Louizos et al., NIPS 2017  
 L0 / L0-sep: Louizos et al., ICLR 2017  
 DN: Pan et al., arXiv 2016  
 GD: Srinivas & Babu, arXiv 2016  
 GL: Wen et al., NIPS 2016

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

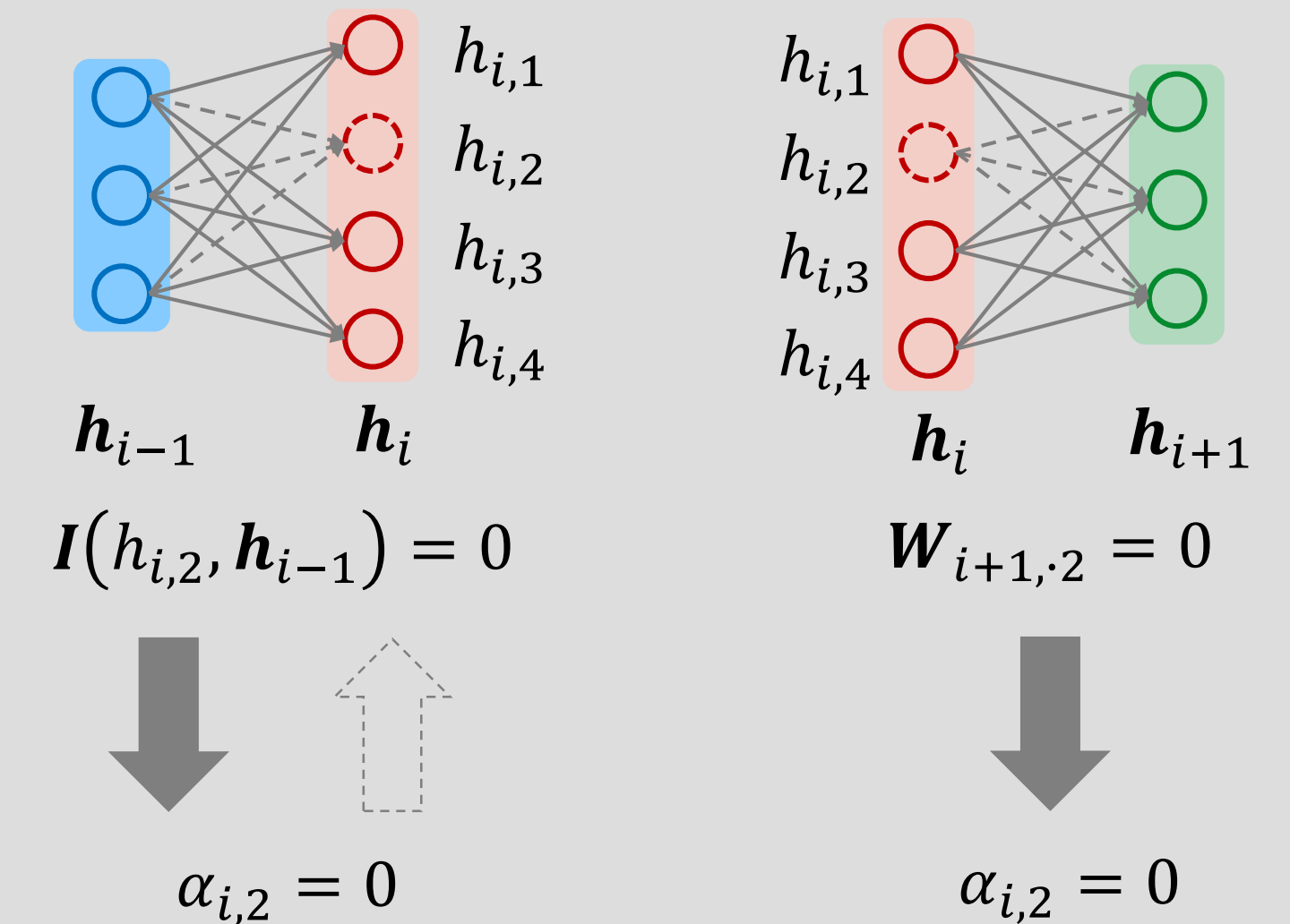
Table 3. VGG16 (modified) on CIFAR10

SBP / SBPa: Neklyu-dov et al., NIPS 2017  
 PF: Li et al., ICLR 2016  
 NS-Single / NS-Best: Liu et al., ICCV 2017  
 RNP: Lin et al., NIPS 2017

## Theoretical Results

❖ Reduced redundancy via intrinsic sparsity:

Define  $\alpha_{i,j} = \frac{\mu_{i,j}^2}{\sigma_{i,j}^2}$ , at the global minimum of (\*)



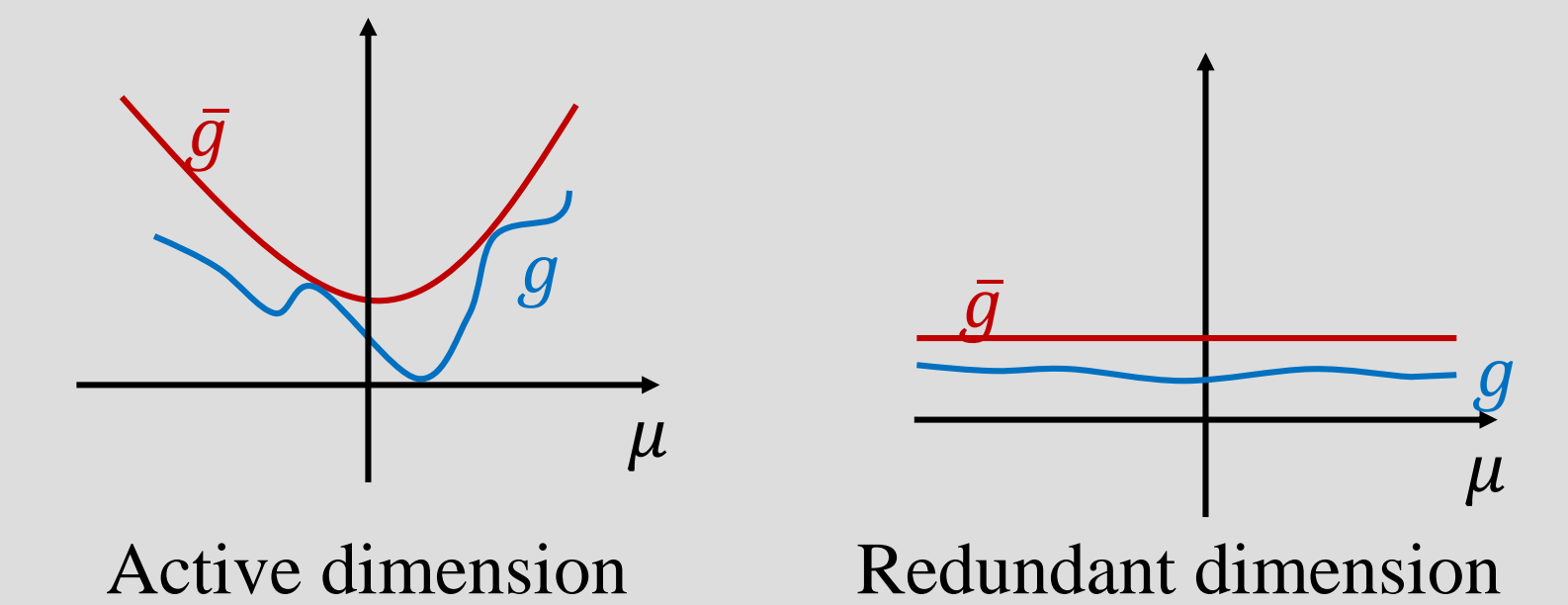
❖ Analysis of tractable upper bounds:

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

$$\int p(\epsilon) g(\epsilon; \theta, W) d\epsilon$$

For fixed network weights  $W'$ , replace  $g(\epsilon; \theta, W)$  with a quadratic upper bound

$$\bar{g}(\epsilon; \theta) \triangleq z(\epsilon; \theta)^T A^T A z(\epsilon; \theta) + b^T z(\epsilon; \theta) + c$$



$A$  becomes low-rank if the network is over-parameterized. At any local minimum of the upper bound of  $\tilde{\mathcal{L}}$ , we have

$$\|\mu^*\|_0 = \|\alpha^*\|_0 \leq \text{rank}[A] + 1 \Rightarrow \text{sparse}$$

Method	$r_W$ (%)	$r_N$ (%)	FLOP(mil)	Error(%)
RNP	-	-	160	38.0
VIBNet	<b>22.75</b>	<b>59.80</b>	<b>133.6</b>	<b>37.6 (37.4)</b>
NS-Single	24.90	-	250.5	26.5
NS-Best	20.80	-	214.8	26.0
VIBNet	<b>15.08</b>	<b>73.80</b>	<b>203.1</b>	<b>25.9 (25.1)</b>

Table 4. VGG16 (modified) on CIFAR100