



# Neuron Campaign for Initialization Guided by Information Bottleneck Theory

Haitao Mao<sup>1,2\*</sup> Xu Chen<sup>1,3\*</sup> Qiang Fu<sup>1\*</sup> Lun Du<sup>1\*</sup> Shi Han<sup>1</sup> Dongmei Zhang<sup>1</sup>

<sup>1</sup> Microsoft Research Asia, China,

<sup>2</sup> University of Electronic Science and Technology of China, China

<sup>3</sup> Peking University, China,

\* These authors contribute equally in this work.

## 1. Introduction

Initialization plays a critical role in the training of deep neural networks (DNN). Existing initialization strategies mainly focus on stabilizing the training process to mitigate gradient vanish/explosion problems. However, these initialization methods are lacking in consideration about how to enhance generalization ability. The **Information Bottleneck (IB) theory** is a well-known understanding framework to provide an explanation about the generalization of DNN. Guided by the insights provided by IB theory, we design two criteria for better initializing DNN. And we further design a neuron campaign initialization algorithm to efficiently select a good initialization for a neural network on a given dataset. The experiments on MNIST dataset show that our method can lead to a better generalization performance with faster convergence.

## 2. motivation & target

- ▶ The purpose of DNN's training procedure: **find a good local minima** near where the loss landscape is flat.
- ▶ Effect of initialization strategy: inappropriate initialization may leads to **stuck in a bad local minima**.
- ▶ A proper initialization: help to faster convergence to the desired solution with good generalization

## 3. Challenges

This work try to address the following challenges:

- ▶ How to enhance the generalization ability during the initialization phase?
- ▶ What is the relationship between model's primary discriminative ability and the generalization ability?
- ▶ How to promote primary discriminative ability of an initial model?

## 4. Information Bottleneck Theory

The information bottleneck (IB) theory provides insights for DNN's generalization ability from: (1) **Input information maintenance**: in the early training phase, DNN increases mutual information between input  $X$  and the latent representation  $Z$  as  $\max I(X, Z)$ , while keeps compressing the input information by minimization of  $I(X, Z)$ . (2) **Target-related information enhancement**: the target-related information keeps growing during the entire training procedure through maximizing the mutual information between the model target  $Y$  and latent representation  $Z$  as  $\max I(Z, Y)$ .

## 7. Conclusion

In this work, we explore the initialization guided with the Information Bottleneck Theory and propose IBCI with efficient neuron campaign algorithm. Two practical criteria, **Input information maintenance** and **Target-related information enhancement** are proposed to enhance the model's primary discriminative ability while initializing. Comparison with widely-used initialization strategies (Xavier and He) demonstrate that our IBCI leads to better generalization ability with faster convergence.

## 5. Proposed Method: IBCI

### 1.IB based Initialization Principle

The two principles guided by IB can be formulated as follows:

$$\max_W \sum_{d=1}^D \alpha_i I(\mathbf{X}; \mathbf{Z}_d) + (1 - \alpha_i) I(\mathbf{Z}_d; \mathbf{Y}) \quad (1)$$

$$\text{s.t. } \forall j, \left\| \mathbf{W}_{:,j}^{(i)} \right\|_2 < \varepsilon, \quad \mathbf{W} = \{\mathbf{W}^{(i)} | i = 1, 2, \dots, D\},$$

where  $\alpha_i$  controlling the importance between input information maintenance and target-related information enhancement.

### 2. Mutual Information Simplification

The input information maintenance criterion can be simplified as

$$\text{tr}(\Sigma_i), \quad (2)$$

where  $\Sigma_i$  is the covariance matrix of  $Z_i$ .

The target-related maintenance criterion is approximated as:

$$\text{tr}(\hat{\mathbf{H}}\hat{\mathbf{H}}^T) - \frac{1}{N} \sum_{j=0}^N \text{tr}(\hat{\mathbf{Z}}_j^T \Pi^j \hat{\mathbf{Z}}_j) \quad (3)$$

More details of the notations can be found in our paper. Intuitively, the

first term denotes the inter-class variance and the second term represents the intra-class variance.

### 3. Neuron Campaign

We utilize a greedy algorithm to find the optimal neuron subset from a large set of randomly initialized neurons and then integrate them as the initialization weight.

#### Algorithm 1 Neuron Campaign initialization algorithm

**Input:** Candidate weight matrix  $\mathbf{W}$

- 1: **for**  $t=1$  to  $T$  **do**
- 2: Update generalized orthonormalization matrix at  $t$  steps:  $\mathbf{A}_t = (\mathbf{A}_{t-1}, \mathbf{a}_t^T)^T$
- 3: Calculate the null space projection by  $\mathbf{P}_t = \mathbf{P}_{t-1} - \mathbf{a}_t \mathbf{a}_t^T \mathbf{W}$
- 4: Select optimal neuron whose index is chosen by  $i = \max_j s_i \frac{\|\mathbf{p}_t'\|}{\|\mathbf{W}_i\|}$
- 5: Update  $\mathbf{w}^* = \mathbf{W}_{:,i}$
- 6: Normalize basis of the generalized orthonormalization matrix as  $\mathbf{a}_{t+1} = \mathbf{p}_t^i / \|\mathbf{p}_t^i\|$
- 7: **end for**

**Output:** Winning neurons formed weight matrix  $\mathbf{W}'$

## 6. Experiments

We conduct experiments on MNIST compared with widely-used initialization strategies, Xavier, He and LSUV. Three different MLPs with ReLU as the activation function are selected as the basic model. The results are organized on Table 1. We also study the effect of **Target Information Enhancement (TIE)** and **Input Information Maintenance (IIE)** on Table 2.

Strategy	Layers	Vanilla	LSUV	IBCI
Xavier	2	2.04 ± 0.03 (75)	2.05 ± 0.06 (51)	<b>1.93 ± 0.06 (60)</b>
	3	1.82 ± 0.05 (52)	1.80 ± 0.07 (63)	<b>1.71 ± 0.09 (36)</b>
	5	2.83 ± 0.16 (98)	3.13 ± 0.17 (69)	<b>2.53 ± 0.09 (78)</b>
He	2	2.03 ± 0.03 (65)	2.00 ± 0.04 (70)	<b>1.93 ± 0.07 (57)</b>
	3	1.83 ± 0.05 (54)	1.86 ± 0.07 (71)	<b>1.73 ± 0.04 (35)</b>
	5	2.76 ± 0.07 (80)	2.90 ± 0.12 (77)	<b>2.62 ± 0.08 (73)</b>

**Table 1:** minimal error rate and corresponding epoch comparison of IBCI with baseline methods on MNIST.

Strategy	Layers	TIE	IIM
Xavier	2	<b>1.93 ± 0.06 (60)</b>	2.04 ± 0.07 (58)
	3	<b>1.71 ± 0.09 (36)</b>	1.82 ± 0.03 (43)
	5	<b>2.53 ± 0.09 (78)</b>	2.68 ± 0.05 (82)
He	2	<b>1.93 ± 0.07 (57)</b>	2.07 ± 0.06 (59)
	3	<b>1.73 ± 0.04 (35)</b>	1.83 ± 0.07 (42)
	5	<b>2.62 ± 0.08 (73)</b>	2.89 ± 0.11 (74)

**Table 2:** minimal error rate and corresponding epoch comparison of IBCI with methods with only one criterion.