

# ON THE INFORMATION OF FEATURE MAPS AND PRUNING OF DEEP NEURAL NETWORKS

---

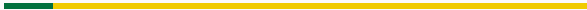
Mohammadreza Soltani

Duke University

# TABLE OF CONTENTS

1. Introduction
2. Estimating the Mutual Information
3. Proposed Pruning Algorithm
4. Experimental Results

# INTRODUCTION



# MOTIVATION AND PROBLEM STATEMENT

- Energy consumption for limited-resource embedded systems
- Very large Memory for saving the weights of a model
- Huge amount of computation for product operation
- For example, under 45nm CMOS technology,
  - A 32bit floating point add consumes 0.9pJ
  - A 32bit SRAM cache access requires 5pJ
  - A 32bit DRAM memory access takes 640pJ
- Running a 1 billion connection neural network, for example, at 20 fps needs almost 13W power just for DRAM!!!
- Need to compress model for deployment and fast inference running-time

# AN IMPORTANT OBSERVATION

- Robustness of deep architectures with *skip-connection* against coarse pruning
  - Removing a random layer doesn't hurt the performance.
  - Removing the models without *skip-connection* drops the performance dramatically.
- Our focus is to investigate this phenomena in more depth
- Studying two prominent examples of models with *skip-connection*: Resnet and DenseNet

- A *skip-unit* is defined as a set of layers, and each layer consists of sequential operations including *Conv*, *Pooling*, *ReLU*, *BN*, *Dropout*, etc.
- A *skip-units* is mathematically defined as

$$U_\ell = \Psi(T_\ell, U_{1:\ell-1}, \alpha_\ell), \quad \ell = 1, 2, \dots, L,$$

- $U_{1:\ell-1}$ , the input of  $\ell$ -th unit
- $T_\ell = f_\ell(U_{\ell-1})$ , the output in the skip-unit
- $f_\ell$ , the composition of aforementioned operations
- $\alpha_\ell$ 's are binary variables and  $\Psi$  denotes an operation that combines  $T_\ell$  and  $U_{1:\ell-1}$ .

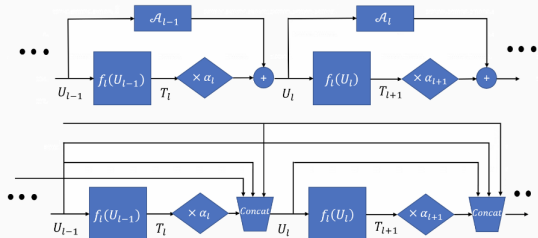
# RESNET AND DENSENET

- ResNet architecture  $\Psi_{res}$  and DenseNet architecture  $\Psi_{den}$  are respectively given by:

$$U_\ell = \Psi_{res}(T_\ell, U_{\ell:\ell-1}, \alpha_\ell) = \alpha_\ell T_\ell + \mathcal{A}_{\ell-1}U_{\ell-1},$$

$$U_\ell = \Psi_{den}(T_\ell, U_{\ell:\ell-1}, \alpha_\ell) = \text{Concat}(\alpha_\ell T_\ell, U_{1:\ell-1}),$$

- $\text{Concat}$  is the concatenation operation.
- $\mathcal{A}_{\ell-1}$  is an identity or a convolution operator.



**Figure 1:** Two consecutive skip-units in a ResNet (top) and DesNet (bottom) family, respectively.

**Compressing skip-units models:** Pruning the model by removing the redundant skip-units based on their learned information.

1. How to study the learned features?
2. How to capture the information in the learned features?
3. How to quantify the redundant the skip-units?



**Mutual Information** as a measure of information of skip-units:

- Measuring the mutual information between the skip-units and the output of the model
- Clustering the units based on their mutual information
- Keeping only the cluster heads ( $\alpha = 1$ )
- Removing the other units in each cluster from the graph of the model ( $\alpha = 0$ )

# ESTIMATING THE MUTUAL INFORMATION

---

## UPPER BOUND ON THE MUTUAL INFORMATION

- Following the method of Kolchinsky, et al., 2017, an upper bound is used for estimating the mutual information.
- Modeling the underlying probability distribution with a Gaussian mixture model with the number of components equals to the training samples.

$$\begin{aligned} I(T; Y) &= H(T) - H(T|Y) \\ &\leq -\frac{1}{n} \sum_{i=1}^n \ln \frac{1}{n} \sum_{j=1}^n \exp \left( -\frac{\|\mu_j - \mu_i\|_2^2}{2\sigma^2} \right) \\ &\quad - \sum_{k=0}^{l-1} p_k \left( -\frac{1}{n_k} \sum_{\substack{i=1 \\ y_i=k}}^n \ln \frac{1}{n_k} \sum_{\substack{j=1 \\ y_j=k}}^{n_k} \exp \left( -\frac{1}{4} \frac{\|\mu_j - \mu_i\|_2^2}{2\sigma^2} \right) \right), \end{aligned}$$

- $p_k = \frac{n_k}{n}$  denotes the probability of class  $k$  and  $n_k = \sum_{i=1}^n \mathbb{I}(y_i = k)$ .
- $\mu_i$  is the mean of each Gaussian component and  $\sigma$  is set to a small number.

# PROPOSED PRUNING ALGORITHM

---

# MULTI-STAGE PRUNING WITH INFORMATION CLUSTERING (MSPIC)

## Algorithm 1

INPUT:

$\text{DNN}^0$ : Pre-trained Deep Neural Network

$S^0$ : The index set of units

$T_l$ : Feature maps,  $l = 1, 2, \dots, |S^0|$

$N$ : Number of stages

$R^t$ : Resolution vector,  $t = 0, 1, \dots, N - 1$

$\epsilon^t$ : Error threshold,  $t = 0, 1, \dots, N - 1$

for  $t = 0, 1, \dots, N - 1$  do

  Compute  $l(T_l^t; Y)$ ,  $l = 1, \dots, |S^t|$  using  $\text{DNN}^t$

  Construct  $\mathbf{l}^t = [l(T_1^t; Y), l(T_2^t; Y), \dots, l(T_{|S^t|}^t; Y)]$

$\{ \text{Cluster}_1^{\text{policy}_1}, \dots, \text{Cluster}_{M_{\text{policy}_1}^t}^{\text{policy}_1}, \dots, \text{Cluster}_{M_{\text{policy}_{|R^t|}}^t}^{\text{policy}_{|R^t|}} \} = \text{Cluster}(R^t, \mathbf{l}^t, S^t)$

  for policy in  $R^t$  do

    for  $j = 1, 2, \dots, M_{\text{policy}}^t$  do

$a_l = 1$ ,  $l = \min_{k \in \text{Cluster}_j^{\text{policy}}}$

$a_u = 0$ ,  $\forall u \in \text{Cluster}_j^{\text{policy}} \setminus l$

    end for

    Compute  $\text{Test}_{\text{error}}^t$  for given policy

  end for

  Select one policy with  $\text{Test}_{\text{error}}^t < \epsilon^t$

$S^{t+1} \leftarrow S^t \setminus \{v : a_v = 0, v \in S\}$

$\text{DNN}^{t+1} \leftarrow$  Re-train the resulted sub-network with the trained weights in stage  $t$

end for

Return pruned model with  $|S^{N-1}|$  units

## EXPERIMENTAL RESULTS

---

# DATASETS AND MODELS

## 1. Datasets.

Dataset	Train data	Test data	Image Size	Classes
CIFAR-10	50000	10000	$32 \times 32 \times 3$	10
CIFAR-100	50000	10000	$32 \times 32 \times 3$	100
Tiny ImageNet	100000	10000	$64 \times 64 \times 3$	200

## 2. Model architectures:

Model	Units	Layers	Param. (M)	FLOPs (M)
ResNet-56	[9, 9, 9]	56	0.85	126.54
ResNet-164	[18, 18, 18]	164	1.70	254.94
DenseNet-100	[6, 12, 24, 16]	100	0.80	305.10

# EXPERIMENT OF PRUNING DNNs ON CIFAR-10 DATASET

Model	Test Accuracy	Param. (M)	FLOPs (M)	Red.(%)
ResNet-56 (full)	0.9334	0.85	126.55	-
ResNet-56	0.9128	0.23	42.30	72.72
ResNet-164 (full)	0.9569	1.70	254.94	-
ResNet164 (t=2)	0.9207	0.99	143.90	41.53
ResNet164 (t=7)	0.9173	0.47	112.50	72.02
DenseNet-100-k12 (full)	0.9531	0.77	293.55	-
DenseNet100-k12 (t=3)	0.9352	0.34	224.90	56.27
DenseNet100-k12 (t=13)	0.9437	0.29	173.20	62.66

**Table 1:** The results of pruning various DNNs on CIFAR-10 data.



# EXPERIMENT OF DNNs ON CIFAR-100 DATASET

Model	Test Accuracy	Param. (M)	FLOPs (M)	Red.(%)
ResNet-164 (full)	0.7799	1.73	254.96	-
ResNet164	0.7459	0.94	130.97	45.30
DenseNet-100-k12 (full)	0.7793	0.80	304.10	-
DenseNet100-k12	0.7408	0.38	203.35	52.37

**Table 2:** The results of pruning various DNNs on CIFAR-100 data.