

---

# AI Project 2: Multi-label Classification on YouTube-8M

---

**KAR CHUN TEONG**

518030990014

Department of Computer Science  
Shanghai Jiao Tong University  
teongkarchun@sjtu.edu.cn

**MATSUNAGA TAKEHIRO**

518030990028

Department of Computer Science  
Shanghai Jiao Tong University  
matsunagatakehiro@sjtu.edu.cn

**EDUARDO WANG ZHENG**

Department of Computer Science  
Shanghai Jiao Tong University  
eduardowangzheng@sjtu.edu.cn

## Abstract

The authors develop two models, which are Deep Fully Connected Neural Networks with Merge-and-Run Mappings(DMRNets) and long short-term memory(LSTM), to achieve multi-label classification on the YouTube-8M dataset.

## 1 Introduction

This project is the group project of CS410 Artificial Intelligence 2020 Fall in Shanghai Jiao Tong University named Multi-label Classification on YouTube-8M. The main purpose is to use a neural network model taking two video-level features(mean\_rgb and mean\_audio) as input, to output the prediction of labels of videos from the YouTube-8M dataset. The authors develop two different neural network models, which are Deep Fully Connected Neural Networks with Merge-and-Run Mappings(DMRNets) and LSTM in their attempts to achieve this goal.

## 2 Neural Network Architecture

### 2.1 Deep Fully Connected Neural Networks with Merge-and-Run Mappings(DMRNets)

#### 2.1.1 Model Overview

One of the neural network architecture we consider is the deep fully connected neural networks with merge-and-run mappings, this is actually a variation of the deep convolutional neural networks with merge-and-run mappings, which coincidentally has the same short-form as DMRNets, to differentiate these two architectures, we would use DMRNets(FCN) as the short-form of Deep fully connected neural networks with Merge-and-Run Mappings and DMRNets(CNN) as the short-form of Deep convolutional neural networks with Merge-and-Run Mappings. DMRNets(CNN) itself is a variation of the residual neural network (ResNet).

ResNet originated from Deep Convolutional Neural Networks(CNN), researchers have been pondering at the question of, does stacking more convolutional layers (which is increasing the depth) leads to better performance? They found out that stacking layers leads to a degradation problem, as they increase the depth, accuracy get saturated then degrades rapidly, and this is not caused by overfitting, the more layers they stack, the higher the training error. Some proposed a solution which is residual mapping, the basic building block of residual mapping is shown in Figure 1. Denote the

desired underlying mapping as  $H(x) = F(x) + x$ , instead of directly fitting to  $H(x)$ , fit the layer to  $F(x)$ , and fit another layer to the identity mapping which is  $x$ , at the end merge the two layers to get the desired mapping  $H(x)$ . The authors of the article argued and proved that this modification significantly reduced the degradation problem while increasing the depth of the model, for more information please refer to the original article[1].

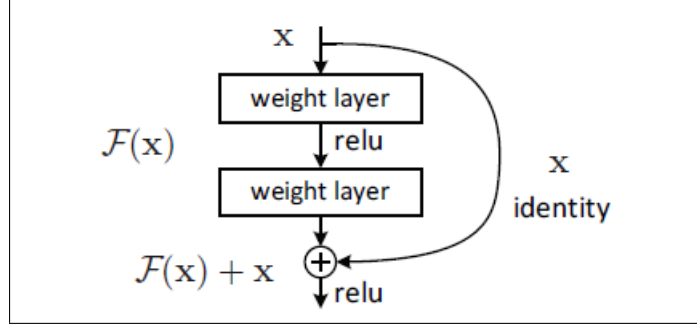


Figure 1: Residual Learning: a building block.

Next, we move on to the variation based on ResNet, the DMRNets(CNN). The idea behind this architecture is to assemble the residual branches in parallel through a merge-and-run mapping. Merge means to average the inputs of these residual branches, and by Run it means to add the average to the output of each residual branch as the input of the subsequent residual branch. To further illustrate the model, please refer to Figure 2. Figure 2(a) shows 2 residual branches assembled sequentially, Figure 2(b) shows 2 residual branches assembled in parallel and with identity mappings, and finally Figure 2(c) shows 2 residual branches assembled in parallel and with the proposed merge-and-run mappings. The authors of the article argued and proved that this modification further reduced the training difficulty of ResNet and retained the advantages of ResNet, for more information please refer to the original article[2].

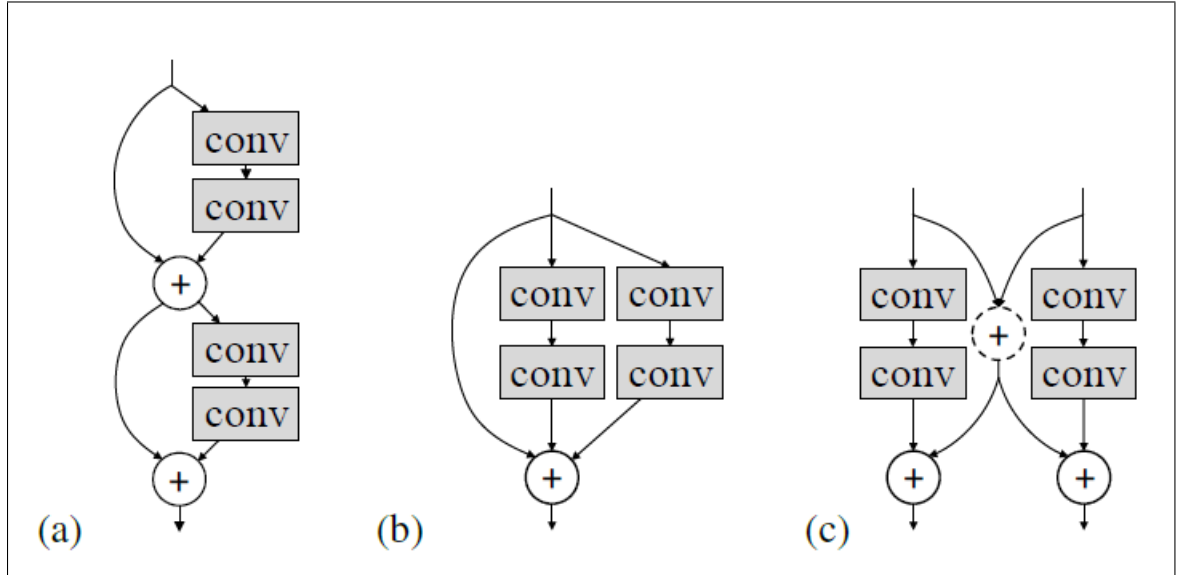


Figure 2: Comparison between different assembly blocks.

In this project, we adopt and do some variations on the DMRNets(CNN) by using the fully connected layers instead of convolutional layers as the basic layers to form the DMRNets(FCN). Part of the model structure is shown in Figure 3. First we define a fc\_block which consists of dense layer,

batch normalization layer, leaky ReLU layer, and dropout layer. Then we define the merge-and-run function, which separate input into 2 branches, branch a pass through the fc\_block while branch b pass through identity mapping, which is essentially mapping to itself, then we pass both branches into an average layer, then the output passes through another leaky ReLU layer. In Figure 3, we can see both input first go through a fc\_block, then into a merge-and-run mappings. The model is actually so deep that the rest of the structure is too large to be shown in the report, but it's essentially the same process, repeatedly pass through fc\_block and merge-and-run mappings and finally merged into a single output, and the full picture is included in the submission folder.

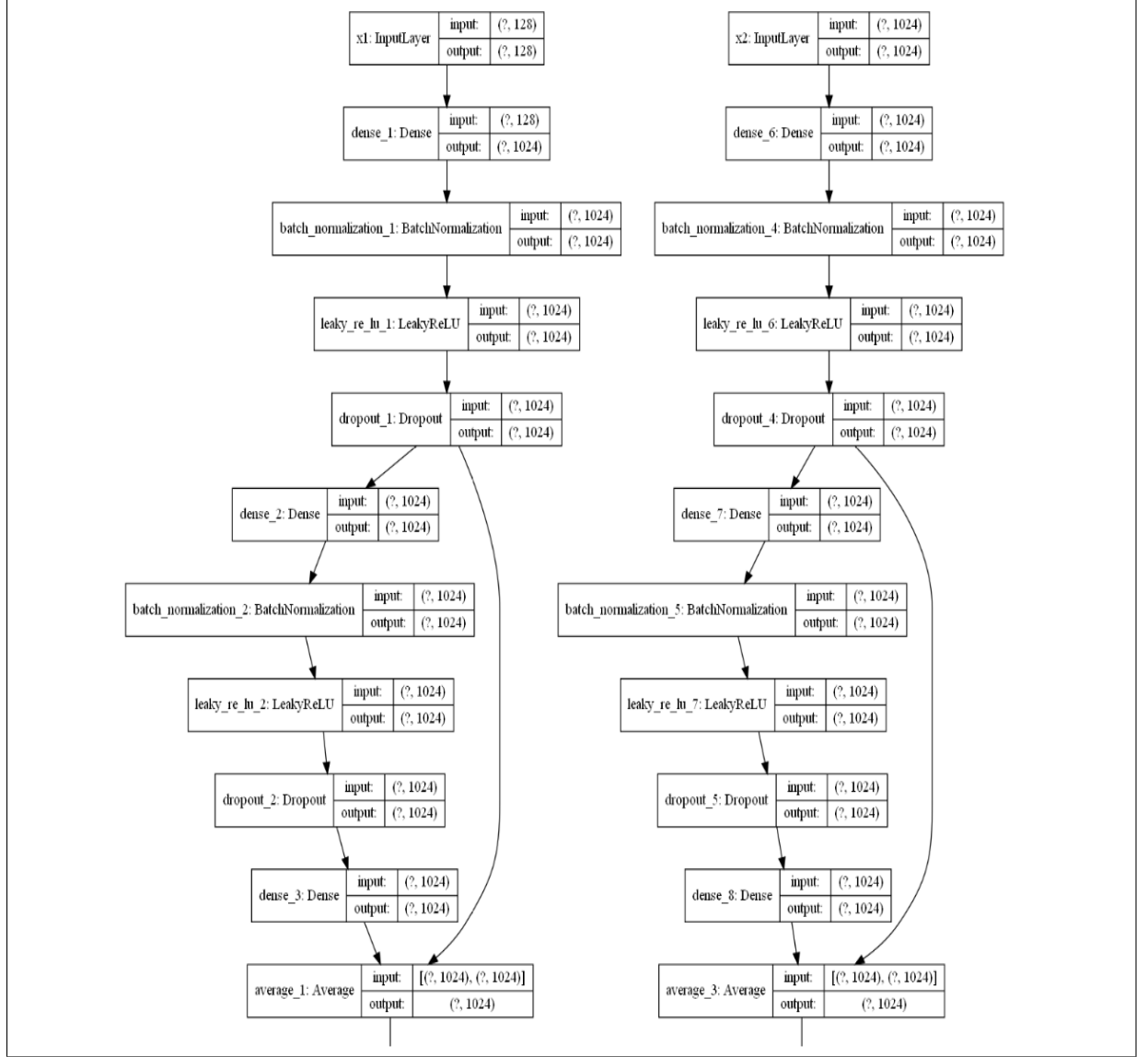


Figure 3: Part of the structure of DMRNets(FCN).

### 3 Contributions

- KAR CHUN TEONG:
  - Production management(organization and distribution of workload, arrangement of meeting schedule, etc.)
  - Analysis and implementation of DMRNets model

- Design and compilation of final report and PPT
- MATSUNAGA TAKEHIRO
  - Analysis and implementation of LSTM model
  - Training and testing of LSTM model
- EDUARDO WANG ZHENG
  - Implementation of DMRNets model
  - Training and testing of DMRNets model

## References

References follow the acknowledgments. Use unnumbered first-level heading for the references. Any choice of citation style is acceptable as long as you are consistent. It is permissible to reduce the font size to small (9 point) when listing the references. **Note that the Reference section does not count towards the eight pages of content that are allowed.**

- [1] Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun. Deep Residual Learning for Image Recognition.
- [2] Liming Zhao, Mingjie Li, Depu Meng, Xi Li, Zhaoxiang Zhang. Deep Convolutional Neural Networks with Merge-and-Run Mappings. <https://arxiv.org/abs/1611.07718>.
- [3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience* **15**(7):5249-5262.