

# Encoding Video and Label Priors for Multi-label Video Classification on YouTube-8M dataset

Seil Na  
Seoul National University  
seil.na@vision.snu.ac.kr

YoungJae Yu  
Seoul National University  
yj.yu@vision.snu.ac.kr

Sangho Lee  
Seoul National University  
sangho.lee@vision.snu.ac.kr

Jisung Kim  
SK Telecom Video Tech. Lab  
joyful.kim@sk.com

Gunhee Kim  
Seoul National University  
gunhee@snu.ac.kr

## Abstract

*YouTube-8M is the largest video dataset for multi-label video classification. In order to tackle the multi-label classification on this challenging dataset, it is necessary to solve several issues such as temporal modeling of videos, label imbalances, and correlations between labels. We develop a deep neural network model, which consists of four components: Video Pooling Layer, Classification Layer, Label Processing Layer, and Loss Function. We introduce our newly proposed methods and discusses how existing models operate in the YouTube-8M Classification Task, what insights they have, and why they succeed (or fail) to achieve good performance. Most of the models we proposed are very high compared to the baseline models, and the ensemble of the models we used is 8th in the Kaggle Competition.*

## 1. Introduction

Many challenging problems have been studied in computer vision research toward video understanding, such as video classification [9, 13], video captioning [19, 20], video QA [21], and MovieQA [16], to name a few. YouTube-8M [1] is the largest video dataset for multi-label video classification. Its main problem is to predict the most relevant labels for a given video out of 4,716 predefined classes. Therefore, it requires jointly solving two important problems; video classification and multi-label classification.

From the view of the video classification, YouTube-8M is challenging in that it covers more general classes like soccer, game, vehicle, and food, while existing video classification datasets focus on more specific class groups, such as sports in Sports-1M [9], and actions in UCF-101 [13]. Therefore, unlike the importance of modeling motion fea-

tures in UCF-101 [13] or Sports-1M [9], it is important to capture more generic video information(e.g. temporal encoding method for video, audio feature modeling) in YouTube-8M.

From the view of multi-label classification, the key issues to solve in YouTube-8M are label imbalances and correlations between labels. YouTube-8M involves 4,716 class labels, and the number of videos belonging to each class is significantly different, which causes a label imbalance issue that the classifier fits to the biased data. At the same time, many classes are closely related one another, such as {Football, Kick, Penalty kick, Indoor soccer} or {Super Mario Bros, Super Mario World, Super Mario bros 3, Mario Kart, Mario Kart 8}. It is also challenging to resolve the correlations between labels to decide final prediction.

Based on the challenges of the multi-label video classification task on YouTube-8M described the above, we focus on addressing i) temporal encoding for video, ii) relieving the label imbalance problem, and iii) utilizing the correlated label information. Our model consists of four components: i) video pooling layer, ii) classification layer, iii) label processing layer, and iv) loss function. The proposed components indeed show significant performance improvement over the baseline models of YouTube-8M [1], and finally our ensemble model is ranked 8th in the Google Cloud & YouTube-8M Video Understanding Challenge<sup>1</sup>.

## 2. The Model

Figure 1 shows the overall pipeline of our model. We first present video features we used, and then explain its four key components in the following sections.

<sup>1</sup><https://www.kaggle.com/c/youtube8m/leaderboard>. (Team name: SNUVL X SKT)

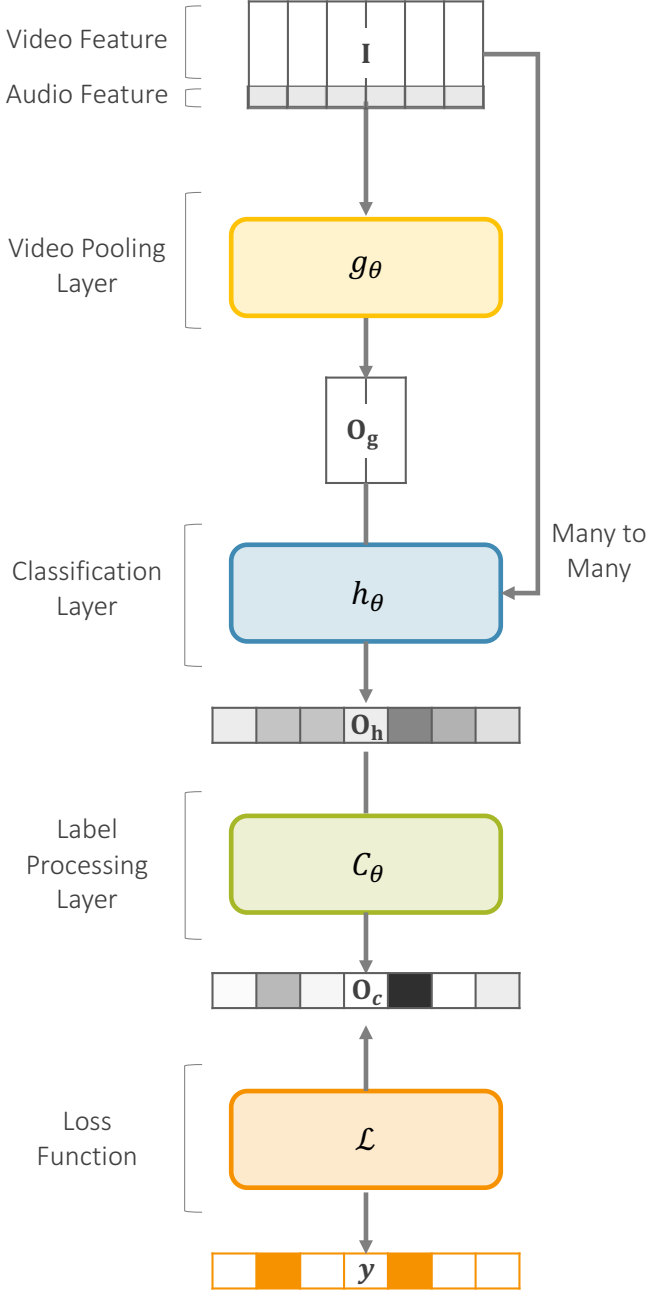


Figure 1. A complete pipeline illustration of our model for YouTube-8M classification. The input of the model concatenates the frame feature  $\mathbf{I}^f$  and audio feature  $\mathbf{I}^a$ , which is denoted by  $\mathbf{I}$ . We divided the model into four components (Video Pooling Layer, Classification Layer, Label Processing Layer, Loss Function) and experimented with the newly proposed model or variants of previously proposed model in each component. Each of the four components is denoted by  $(g_\theta, h_\theta, C_\theta, \mathcal{L})$ , and each output of first 3 layer is denoted by  $(\mathbf{O}_g, \mathbf{O}_h, \mathbf{O}_c)$ .

## 2.1. Video Features

The inputs of the model are frame features and audio features of a video clip. The frame features are obtained by sampling a clip at 1-second interval, and extracting 2,048-dimensional vector from every frame through Inception-V3 Network [15] pretrained to ImageNet [4]. Then the feature is reduced to an 1,024-dimension via PCA (+ whitening), then quantized, and finally L2-normalized. As a result, for a video of  $T$  seconds long, its frame is  $\mathbf{I}^f \in \mathbb{R}^{T \times 1,024}$ . The audio features re extracted using the VGG-inspired acoustic model [6] followed by L2-normalization, which is denoted by  $\mathbf{I}^a \in \mathbb{R}^{T \times 128}$ .

As the input of our model, we concatenate the frame feature  $\mathbf{I}^f$  and the audio feature  $\mathbf{I}^a$  at every time step, denoted by  $\mathbf{I} \in \mathbb{R}^{T \times 1,152}$ . We test the compact bilinear pooling [5] with various dimensions between  $\mathbf{I}^f$  and  $\mathbf{I}^a$ , but all of them have significantly lower performance than the simple concatenation. From now on, we use  $\mathbf{I}$  to denote the input features of a video over all frames, and  $\mathbf{I}_t \in \mathbb{R}^{1,152}$  as  $t$ -th frame vector.

## 2.2. Video Pooling Layer

The Video Pooling Layer  $g_\theta : \mathbb{R}^{T \times 1,152} \rightarrow \mathbb{R}^d$  is defined as a parametric function that encodes a sequence of  $T$  feature vectors  $\mathbf{I}$  into a  $d$ -dimensional embedding vector. We test five different encoding structures as follows.

### 2.2.1 A Variant of LSTM

The Long Short Term Memory (LSTM) model [7] is one of the most popular frameworks for modeling sequence data. We use a variant of the LSTM as follows:

$$\mathbf{i}_t = \sigma(\mathbf{I}_t \mathbf{U}^i + \mathbf{s}_{t-1} \mathbf{W}^i + \mathbf{b}_i) \quad (1)$$

$$\mathbf{f}_t = \sigma(\mathbf{I}_t \mathbf{U}^f + \mathbf{s}_{t-1} \mathbf{W}^f + \mathbf{b}_f) \quad (2)$$

$$\mathbf{o}_t = \sigma(\mathbf{I}_t \mathbf{U}^o + \mathbf{s}_{t-1} \mathbf{W}^o + \mathbf{b}_o) \quad (3)$$

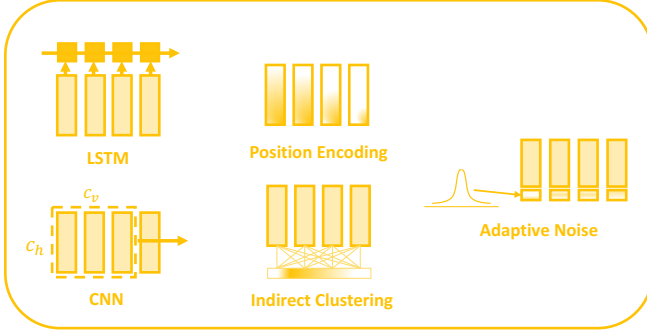
$$\mathbf{g}_t = \tanh(\mathbf{I}_t \mathbf{U}^g + \mathbf{s}_{t-1} \mathbf{W}^g + \mathbf{b}_g) \quad (4)$$

$$\mathbf{c}_t = \mathbf{c}_{t-1} \circ \mathbf{f}_t + \mathbf{g}_t \circ \mathbf{i}_t \quad (5)$$

$$\mathbf{s}_t = \tanh(\mathbf{c}_t) \circ \mathbf{o}_t \quad (6)$$

where  $t$  denotes each time step,  $\mathbf{i}, \mathbf{f}, \mathbf{o}$  are the input, forget, output gate,  $\mathbf{c}_t, \mathbf{s}_t$  are long-term and short-term memory respectively.

The baseline model uses only the final hidden states of the LSTM (*i.e.*  $\mathbf{c}_T, \mathbf{s}_T$ ), but we additionally exploit the following two states in order to extract as much information as possible from the LSTM: i) the state  $\mathbf{M}_l = \sum_{t=1}^T \mathbf{I}_t$ : the summation of the input feature  $\mathbf{I}_t$  of each time step, and ii) the state  $\mathbf{O}_l = \sum_{t=1}^T \mathbf{g}_t$ : the summation of the output of each time step of the LSTM. We concatenate  $\mathbf{M}_l$  and  $\mathbf{O}_l$  with  $\mathbf{c}_T$  and  $\mathbf{s}_T$ . That is, if the cell size of the LSTM is  $d$ , the output  $\mathbf{g}_\theta$  of the baseline that uses  $\mathbf{c}_T, \mathbf{s}_T$  becomes



(a) Video Pooling Layer  $g_\theta$

Figure 2. The first component Video Pooling Layer  $g_\theta$  of the entire model. It takes  $T$  frame vector  $\mathbf{I}$  as input and outputs  $d$ -dimensional vector using different encoding methods respectively.

a  $(2 \times d)$ -dimensional vector, whereas our model that additionally uses  $\mathbf{M}_t$  and  $\mathbf{O}_t$  has the output of a  $(4 \times d)$ -dimensional vector. Experimentally, we choose the LSTM cell size  $d = 1,152$ . We also apply the layer normalization [2] to each layer of the LSTM for fast convergence, and also use the dropout with dropout rate=0.8 to increase the generalization capacity.

### 2.2.2 CNNs

Convolutional neural networks (CNNs) are often used to jointly capture spatial information from images or video in many computer vision tasks. That is, the convolution kernels generate output signals considering all the elements in the window together, and thus they effectively work with spatially or temporally sequential information (*e.g.* images, consecutive characters in NLP, and audio understanding). As the second candidate of Video Pooling Layer, we use the CNN to capture temporal information of video as proposed in [10]:

$$\mathbf{o}_c = \text{ReLU}(\text{conv}(\mathbf{I}, \mathbf{w}_{conv}, \mathbf{b}_{conv})) \quad (7)$$

where  $\text{conv}(\text{input}, \text{filter}, \text{bias})$  indicates convolution layer with stride 1, and  $\text{ReLU}$  indicates the element-wise ReLU activation [12].  $\mathbf{w}_{conv} \in \mathbb{R}^{c_v \times c_h \times 1 \times d}$  is a convolution filter with the vertical and horizontal filter size of  $c_v = 5, c_h = 1152$ , and  $\mathbf{b}_{conv} \in \mathbb{R}^d$  is a bias.  $\mathbf{o}_c \in \mathbb{R}^{(T-c_v+1) \times d}$  indicates the output of the convolution layer. Finally, we apply max-pooling over time for  $\mathbf{o}_c$ , obtaining the  $d$ -dimensional encoding of  $g_\theta$ .

### 2.2.3 Position Encoding

We also test the Position Encoding scheme [14] that assigns different weights to each frame vector. That is, we define

the matrix  $\mathbf{L} \in \mathbb{R}^{T \times 1152}$ , Position Encoding is simply defined as follows.

$$\mathbf{I} \leftarrow \mathbf{I} \odot \mathbf{L} \quad (8)$$

$$\mathbf{O}_g = \sum_{t=1}^T \mathbf{I}_t \quad (9)$$

where  $L_{ij} = (1 - i/T) - (j/1152)(1 - 2 \times i/T)$ ,  $\odot$  means element-wise multiplication. After applying Position Encoding, we used summation of each vector as output of  $g_\theta$ , that means  $\mathbf{O}_g$  is  $d$ -dimensional vector where  $d$  is 1,152.

### 2.2.4 Self-Attention: Indirect Clustering

The YouTube-8M dataset deals with general topics (*e.g.* soccer, game, car, animation) rather than relatively focused labels like Sports-1M or UCF-101. We here test the following hypothesis: since the topic is highly general, it may be more advantageous to focus on the most dominant parts of the video rather than temporal/motion information of individual frames. Therefore, we suggest an indirect clustering model using the self-attention mechanism as follows.

We perform a clustering on the video features over all frames  $\mathbf{I} \in \mathbb{R}^{T \times 1,152}$ , and find the cluster with the largest size (*i.e.* the largest number of elements in the cluster). Then, the vectors in this cluster may represent the main scene of the video. However, since it takes very long time to perform clustering on each video, we propose a self-attention mechanism that acts like clustering as follows.

$$p_t = \text{softmax}\left(\sum_{i=1}^T \mathbf{I}_t \mathbf{I}_i\right) \quad (10)$$

where  $p_t$  is a scalar value that indicates the soft attention to a frame  $t$  while the softmax is applied over  $p_{1..T}$ . That is, the more similar the frame vector is to the other vectors, the higher its  $p_t$  value is (*i.e.* it is more likely to be the main scene of the video). Finally, the frame encoding is simply obtained by a weighted sum of frame vectors by the attention values:  $\mathbf{g}_\theta = \sum_{t=1}^T p_t \mathbf{I}_t$ .

### 2.2.5 Adaptive Noise

Each of 4,716 classes in the YouTube-8M dataset has a different number of video examples. For example, the *Car* class has about 800,000 examples, but the *Air Gear* class has only 101 examples. Let  $n$  to be the number of labels associated with a video  $\mathbf{I}$ ; we introduce the adaptive noise structure to relieve the label imbalance problem as follows.

$$\mathbf{I}_t \leftarrow \mathbf{I}_t + \gamma \cdot \mathbf{Z} \quad (11)$$

$$\text{where } \mathbf{Z} \sim \mathcal{N}(0, I), \quad \gamma = \frac{1}{n} \cdot \sum_{i=1}^n \frac{1}{S(y_i)} \quad (12)$$

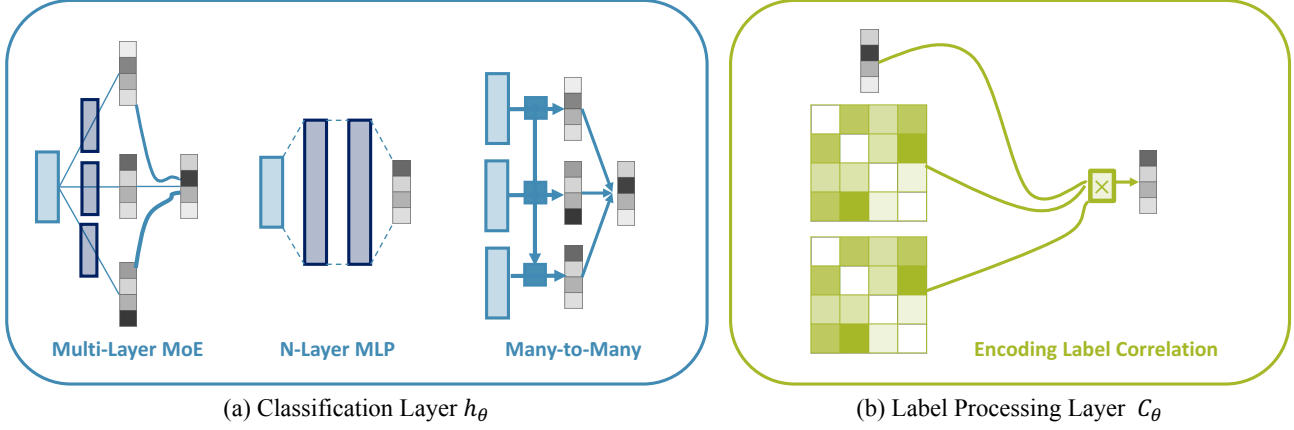


Figure 3. (a)The second component Classification Layer  $h_\theta$  of the overall model. It take the  $d$ -dimensional vector as input  $O_g$  and have score  $O_h$  for 4,716 classes as output. However, input of Many-to-Many model is  $I$  rather than  $O_g$ . (b)the 3rd component Label Processing Layer  $C_\theta$  of the entire model. It uses label correlation prior to update the score for each class.

where noise  $Z$  is sampled from normal distribution,  $y_i$  denotes the  $i$ -th label for video, and  $S(y_i)$  is the number of video examples that have label  $y_i$ . It means that we increase the generalization for small classes by adding more noise to their frame vectors. We then make summation of the vectors of all frames as the output of  $\mathbf{g}_\theta$ , as done in the position encoding.

### 2.3. Classification Layer

Classification Layer  $h_\theta$  is defined as follows.

$$\begin{cases} h_\theta : \mathbb{R}^{T \times 1152} \rightarrow \mathbb{R}^{4,716}, & \text{if } h_\theta \text{ is Many-to-Many} \\ h_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^{4,716}, & \text{otherwise} \end{cases} \quad (13)$$

That is, by default,  $h_\theta$  takes frame encoding  $O_g$  as an input and outputs score  $O_h$  for 4,716 classes. Exceptionally, as shown in figure 1, the Many-to-Many model has  $I$  as its input, where  $I$  is untouched by video pooling layer  $g_\theta$ . We have conducted experiments on the following three structures in this component. (see figure 3(a))

#### 2.3.1 Many-to-Many

Unlike other models, Many-to-Many model has a frame vector  $I$  that is not touched by video pooling layer  $g_\theta$ . It uses LSTM similar to that of 2.2.1, but it calculates the score by attaching fully-connected layer to each output of each step in LSTM and average them, which is used as output  $O_h$ . Since this model averages out the score in each frame, it has the temporal encoding ability of the RNN as well as the fact that the scores drawn in the more frequently appearing frames are reflected more. As a result, more effective video encoding could be performed.

#### 2.3.2 Variants of Mixture of Experts

The Mixtures of Experts [8] model is a binary classifier that adaptively takes into account the scores of several experts corresponding to a class. For one class, each expert  $e_i$  has a probability value between 0 and 1, and gate  $g_i$  represents the weight for each expert and is defined as follows.

$$e_i = \sigma(\mathbf{w}_e^T O_g + \mathbf{b}_e) \quad (14)$$

$$g_i = \text{softmax}(\mathbf{w}_g^T O_g + \mathbf{b}_g) \quad (15)$$

where  $\mathbf{w}_e, \mathbf{w}_g$  are  $d$ -dimensional vectors, scalars  $\mathbf{b}_e, \mathbf{b}_g$  are biases, and softmax is performed for  $\{g_1, \dots, g_E\}$ . ( $E$  is number of experts.) We extend this MoE model to multi-layer, and construct multiple fully-connected of defining probability and gate distribution of each expert. For example, the gate distribution and the expert distribution of the 2-layer MoE model are defined as follows.

$$e_i = \sigma(\mathbf{v}_e^T (\mathbf{W}_e O_g + \mathbf{b}'_e) + \mathbf{b}_e) \quad (16)$$

$$g_i = \text{softmax}(\mathbf{v}_g^T (\mathbf{W}_g O_g + \mathbf{b}'_g) + \mathbf{b}_g) \quad (17)$$

where each weight matrix  $\mathbf{W}_e, \mathbf{W}_g$  is denoted by  $\mathbf{W}_e \in \mathbb{R}^{d' \times d}, \mathbf{W}_g \in \mathbb{R}^{d' \times d}$  and each  $\mathbf{v}_e^T, \mathbf{v}_g^T$  means  $d'$ -dimensional vector.  $\mathbf{b}_e, \mathbf{b}_g, \mathbf{b}'_e, \mathbf{b}'_g$  are biases. Finally, the score  $O_h[i]$  for  $i$ -th class is determined by the weighted sum of each expert distribution;  $O_h[i] = \sum_{i=1}^k e_i g_i$ .

#### 2.3.3 Multi Layer Perceptron

The Multi Layer Perceptron model is stack of Fully-Connected Layers, one of the most basic Neural Network

structures. We experimentally set the number of layers to 3 and apply Layer Normalization [2] to each layer.

## 2.4. Label Processing Layer

Label Processing Layer  $C_\theta$  is defined as follows.

$$C_\theta : \mathbb{R}^{4,716} \rightarrow \mathbb{R}^{4,716} \quad (18)$$

This component is designed to reflect the correlation between the labels into the model. For example, YouTube-8M, which includes  $\{\textit{Soccer}, \textit{Football}, \textit{Kick}, \textit{Indoor soccer}\}$  and  $\{\textit{Super Mario Bros}, \textit{Super Mario World}, \textit{Super Mario bros 3}, \textit{Mario Kart}, \textit{Mario Kart 8}\}$ . In order to take advantage of this property, we set up the label correlation matrix  $M_c$  as follows by counting all the videos in the training set. (See Figure 3(b))

$$M_c \in \mathbb{R}^{4,716 \times 4,716} \quad (19)$$

where  $M_c$  is correlation matrix and the correlation value  $M_c[i, j]$  between the  $i$ -th label and the  $j$ -th label is calculated higher as  $i$ -th label and  $j$ -th label appear together more in the same video. Then, a new score is defined as follows to better reflect the correlation between labels through simple linear combination of matrix-vector multiplication.

$$O_c = \alpha \cdot O_h + \beta \cdot M_c O_h + \gamma \cdot M'_c O_h \quad (20)$$

Here,  $M_c$  is used as a fixed value, and  $M'_c$  is a trainable parameter initialized to the same value as  $M_c$ . Scalar values  $\alpha, \beta, \gamma$  are hyperparameters for model.

## 2.5. Loss Function

### 2.5.1 Center Loss

Center Loss was first proposed for face recognition task [18] and expanded to other field because of its effectiveness for making discriminative embedding feature [17]. The purpose of the center loss is to minimize intra-class variations while maximizing inter-class variations using the joint supervision of cross-entropy loss and center loss. The original center loss was used for the single label classification problem and it is hard to exploit it in a multi-label classification problem.

If we convert the multi label problem into a single label problem like [3], increment of centers according to the combination of labels is the simple expansion of center loss to a multi-label classification problem. However, this simple expansion is not suitable for YouTube-8M, because the number of combination for labels is too big to calculate. Therefore, we modified the center loss to suit the multi-label classification problem as follow.

$$\mathcal{L}_c = \frac{1}{N} \sum_{i=1}^N \|e_i - c_k\|_2^2, k = y_i \quad (21)$$

$$\mathcal{L} = \mathcal{L}_s + \lambda \mathcal{L}_c. \quad (22)$$

Where  $N$  denotes the number of labels in one video.  $e_i$  denotes a embedding vector from penultimate layer,  $c_k$  denotes the  $y_i$ -th corresponding class center for each  $e_i$ . The  $\mathcal{L}_s$  is the cross-entropy loss,  $\mathcal{L}_c$  is the center loss. A scalar  $\lambda$  is a hyperparameter for balancing the two loss functions.

### 2.5.2 Pseudo-Huber Loss

The Huber Loss is a combination of L2 Loss and L1 Loss, which allows the model to be trained more robustly on noise instances. In the case of the YouTube-8M, there are classes with very few instances due to the label imbalance problem, and Huber Loss is designed to better learn instances belonging to these classes. For a simple, differentiable form, we use the Pseudo-Huber Loss function, a smooth approximation of the Huber Loss, as follows.

$$\mathcal{L}_c = \delta^2 (\sqrt{1 + (\mathcal{L}_{CE}/\delta)^2} - 1) \quad (23)$$

where  $\mathcal{L}_{CE}$  means that Cross-Entropy Loss between our prediction  $O_c$  and ground-truth label  $y$ ,  $\delta$  means hyperparameter for model.

## 2.6. Training

For training our model, we choose the Adam [11] optimizer using a batch size of 128, with learning rate = 0.0006,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 1e - 8$ . We also apply learning rate decay with the rate = 0.95 for every 1.5M iterations. We train the model for 5 epochs with no early stopping. We use both official training and validation data that YouTube-8M public dataset provides for training our models.

## 3. Experiments

We use the test data from the Kaggle competition: Google Cloud & YouTube-8M Video Understanding Challenge to measure the performance of the model. The source for our model is publicly available<sup>2</sup>.

### 3.1. Experimental Setting

One of the greatest features of our model is that we have three issues to solve the YouTube-8M classification task; i)temporal encoding for video, ii)label imbalance problem, iii)correlation between labels, and we have tried several variations on each component by dividing the model pipeline into four components; i)Frame Encoding, ii)Classification Layer, iii)Label Processing Layer, iv)Loss Function. In addition, except in the Many-to-Many model in our pipeline, each component is completely independent of one another, so it is a very good structure for experimenting with a number of variant model combinations. However, it is impossible to do brute-force experiments because the number of trials is too great to test all combinations of

<sup>2</sup><https://github.com/seilna/youtube-8m>

Method	GAP@20
LSTM	0.811
LSTM-M	0.815
LSTM-M-O	<b>0.820</b>
LSTM-M-O-LN	0.815
CNN-64	0.704
CNN-256	0.753
CNN-1024	-
Position Encoding	0.782
Indirect Clustering	0.801
Adaptive Noise	0.782

Table 1. The results on various transformation structures in Video Pooling Layer  $g_\theta$  component.

variations that each component can try. Thus, we took the greedy approach, fixed the remaining 3 components in each component, and experimented the several structures only for that component, choosing the best-performing structure for each component. We used Google Average Precision<sup>3</sup> (GAP)@20 as a metric to measure the performance of the model.

## 3.2. Quantitative Results

### 3.2.1 Video Pooling Layer

There are five transform structures in Video Pooling Layer  $g_\theta$ ; i) variants of LSTM, ii) CNN, iii) Position Encoding, iv) Indirect Clustering, v) Adaptive Noise. To select the most suitable structure for  $g_\theta$ , we fixed the rest of the component  $h_\theta$  to the MoE-2 model, not  $C_\theta$  used, and  $\mathcal{L}_c$  used the cross entropy loss. Based on these settings, the results for each of the five structures are shown in table 1.

LSTM simply uses the last hidden state, LSTM-M concatenates  $M_l$ , LSTM-M-O concatenates both  $M_l$  and  $O_l$ , and LSTM-M-O-LN is a model that applies Layer Normalization for each layer of LSTM. CNN-64, CNN-256, and CNN-1024 refer to models with 64, 256, and 1024 as the output channels of CNN, respectively.

The results show that the LSTM method has the best performance for encoding frames. Within LSTM, the higher the utilization of the internal information, the higher the performance, indicating that higher performance can be expected if more information can be extracted from other methods such as skip connection in LSTM. On the other hand, unlike the expectation, the LSTM with Layer Normalization has a lower performance than that which is not. Of course, Layer Normalization had the advantage of stable and fast convergence even when 20-30 times learning rate was applied, but when compared to the final performance alone, performance was not good.

CNN showed a very poor performance unexpectedly. As

<sup>3</sup><https://www.kaggle.com/c/youtube8m#evaluation>

Method	GAP@20
Many-to-Many	0.791
2 Layer MoE-2	0.424
2 Layer MoE-16	0.421
3 Layer MLP-4096	0.802
3 Layer MLP-4096-LN	<b>0.809</b>

Table 2. The results on various transformation structures in Classification Layer  $h_\theta$  component.

CNN’s output channel increased, performance was on the rise, but channels larger than 256 were not experimented because of the memory limitations of the GPU. While we can expect that CNN will perform well for larger channels, another problem with CNN is the dramatic increase in computation cost as channels increase. However, if CNNs with different hyperparameters record similar performance to LSTM, CNN is likely to be used as good pooling method because it has the advantage that the convolution operation is fully parallelizable.

The Position Encoding model showed lower performance than LSTM, and one of the possible reasons is that the sequence modeling power of the model is weaker than LSTM.

The Indirect Clustering model has lower performance than LSTM, but it performs better than the Position Encoding model. This suggests that the assumptions we have made (the importance of the main scene in video classification) are not entirely wrong, suggesting that we need a model that can cover this issue more delicately. Also, temporal encoding of LSTM is as important as considering main scene in video classification.

Adaptive Noise do not lead to a significant performance improvement, indicating that a more sophisticated approach to the label imbalance problem is required.

### 3.2.2 Classification Layer

The Classification Layer  $h_\theta$  has three variants: i)Many-to-Many, ii)Multi-Layer MoE, iii)MLP. To select the best method for  $h_\theta$ , we fix the remaining component  $g_\theta$  with Indirect Clustering,  $C_\theta$  is not used, and  $L_c$  uses cross entropy loss. (However, in the many-to-many model,  $g_\theta$  is not used according to the definition.) Based on these settings, the results for each of the three structures are shown in Table 3.

2 Layer MoE-2 refers to a 2 layer MoE model with 2 experts, and 2 Layer MoE-16 refers to a 2 layer MoE model with 16 experts. 3 layer MLP-4096 is MLP structure with 4096 dimension of each hidden layer and 3 hidden layers, and 3 layer MLP-4096-LN applies layer normalization to each layer.

In constructing the classification layer, the MLP structure showed the best performance among the three meth-

Method	GAP@20
MoE - (1.0, 0.3, 0.0)	0.784
MoE - (1.0, 0.1, 0.0)	0.787
MoE - (1.0, 0.0, 0.1)	0.788
MoE - (1.0, 0.01, 0.0)	<b>0.790</b>
MoE - (1.0, 0.0, 0.01)	<b>0.790</b>
MoE - (1.0, 0.01, 0.01)	0.788

Table 3. The results on various hyperparameters in Label Processing Layer  $C_\theta$  component.

ods. The Many-to-Many model is an LSTM-based model, but has a lower performance than the basic LSTM structure in Table 2, indicating that the Many-to-Many framework in video classification is not always a good choice.

On the other hand, disappointingly, the Multi Layer MoE model showed severe overfitting in only two layers. As expected, this requires a lot of parameters to create an intermediate level of embedding, which is also needed for each experts, resulting in overfitting for many parameters. We experimented with different number of layers and hidden layer dimensions for the MLP structure, but 4096 dimensions of 3 layers showed the highest performance. Unusual is that, when layer normalization was applied, the MLP model showed improved performance unlike LSTM, and did not require more delicate hyperparameter tuning.

### 3.2.3 Label Processing Layer

To select the best method for  $C_\theta$ , we fix the remaining component  $g_\theta$  with Indirect Clustering,  $h_\theta$  using the MoE-16 model and  $\mathcal{L}_c$  using the cross entropy loss. As defined in Equation 20, the Label Processing Layer  $C_\theta$  uses pre-computed correlation matrix  $M_c$ .  $\alpha, \beta$ , and  $\gamma$  values can be used to control the degree to which the correlation matrix affects the class score. An experiment is shown in Table 4.

As a result, performance dropped for all  $\beta$  and  $\gamma$  values greater than zero. In other words, when using the label correlation matrix for the purpose of score update, it showed rather low classification performance.

The possible reason for the result may be that our model is too naive to reflect the correlation label prior and that label correlation in the YouTube-8M Dataset is not strong enough to improve the performance of the classification task.

### 3.2.4 Loss Function

To select the most suitable loss for  $\mathcal{L}_c$ , we fix the remaining component  $g_\theta$  with Indirect Clustering,  $h_\theta$  uses the MoE-16 model, and  $C_\theta$  is not used.

Table 5 shows the performance depending on whether or not the center loss and the pseudo-huber loss are used.

Method	GAP@20
$\mathcal{L}_{CE}$	0.798
$\mathcal{L}_{CE} + \mathcal{L}_c(\lambda = 0.001)$	0.799
Huber $_{CE}(\delta = 0.5)$	<b>0.803</b>
Huber $_{CE}(\delta = 1.0)$	0.801
Huber $_{CE}(\delta = 2.0)$	0.798
Huber $_{CE}(\delta = 3.0)$	0.794

Table 4. The results on variants in loss function  $L_c$  component.

$\mathcal{L}_{CE}$  is the cross entropy loss,  $\mathcal{L}_c$  is the center loss, and  $Huber_{CE}$  is the Pseudo-Huber loss for cross entropy term.

The results show that, firstly, using the center loss term  $\mathcal{L}_c$  gives little performance gain. Considering that we apply the center loss term to the multi-label, the performance should be improved if the labels are correlated with each other. Here we get another insight for the correlation label prior, which is that the correlation does not exist as strongly as to achieve performance improvement. Second, Huber Loss proved its ability to cover noise data somewhat from the YouTube video annotation system<sup>4</sup> by recording a relatively clear performance improvement.

### 3.3 Ensemble Model

Based on the experiments we conducted in section 3.2, we recorded test performance of 0.820 as a single model combining LSTM-M-O, MoE-2, and HuberLoss.(Due to GPU memory limitation, we could not apply MoE-16 or Multi Layer MLP model to LSTM based model) In addition, our ensemble model, which is a simple average of the scores of several models, showed a test performance of **0.839** and ranked 8th in the kaggle challenge. The interesting thing is that when we assemble several models, we did not get a significant performance increase if we combine the better single models. Rather, for increasing ensemble model’s performance, the models incorporated should be as diverse as possible.

## 4. Conclusion

We defined three issues that need to be covered in order to solve the YouTube-8M Video Classification task, and we divided the model pipeline into four components and experimented with the various structures to solve issues in each component. As a result, almost all of the deformed structures tried to perform better than the baseline, and their ensemble model recorded **0.839** in the test performance and 8th in the kaggle challenge. We also provided insights on the structures we tried on each component, on what roles each structure plays, and why they work well or poor. Based on this insight, we will explore ways of better frame encoding or use more elegant label correlation priorities as our

<sup>4</sup>[https://www.youtube.com/watch?v=wf\\_77z1H-vQ](https://www.youtube.com/watch?v=wf_77z1H-vQ)

future work.

## References

- [1] S. Abu-El-Haija, N. Kothari, J. Lee, P. Natsev, G. Toderici, B. Varadarajan, and S. Vijayanarasimhan. Youtube-8m: A large-scale video classification benchmark. *arXiv preprint arXiv:1609.08675*, 2016.
- [2] J. L. Ba, J. R. Kiros, and G. E. Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [3] A. C. de Carvalho and A. A. Freitas. A tutorial on multi-label classification techniques. In *Foundations of Computational Intelligence*. Springer, 2009.
- [4] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [5] A. Fukui, D. H. Park, D. Yang, A. Rohrbach, T. Darrell, and M. Rohrbach. Multimodal compact bilinear pooling for visual question answering and visual grounding. *arXiv preprint arXiv:1606.01847*, 2016.
- [6] S. Hershey, S. Chaudhuri, D. P. Ellis, J. F. Gemmeke, A. Jansen, R. C. Moore, M. Plakal, D. Platt, R. A. Saurous, B. Seybold, et al. Cnn architectures for large-scale audio classification. *arXiv preprint arXiv:1609.09430*, 2016.
- [7] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 1997.
- [8] M. I. Jordan and R. A. Jacobs. Hierarchical mixtures of experts and the em algorithm. *Neural computation*, 1994.
- [9] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification with convolutional neural networks. In *CVPR*, 2014.
- [10] Y. Kim. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*, 2014.
- [11] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [12] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, 2010.
- [13] K. Soomro, A. R. Zamir, and M. Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
- [14] S. Sukhbaatar, J. Weston, R. Fergus, et al. End-to-end memory networks. In *NIPS*, 2015.
- [15] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016.
- [16] M. Tapaswi, Y. Zhu, R. Stiefelhagen, A. Torralba, R. Urta-sun, and S. Fidler. Movieqa: Understanding stories in movies through question-answering. In *CVPR*, 2016.
- [17] H. Wang, Z. Li, X. Ji, and Y. Wang. Face r-cnn. *arXiv preprint arXiv:1706.01061*, 2017.
- [18] Y. Wen, K. Zhang, Z. Li, and Y. Qiao. A discriminative feature learning approach for deep face recognition. In *ECCV*, 2016.
- [19] J. Xu, T. Mei, T. Yao, and Y. Rui. Msr-vtt: A large video description dataset for bridging video and language. In *CVPR*, 2016.
- [20] Y. Yu, H. Ko, J. Choi, and G. Kim. Video captioning and retrieval models with semantic attention. *arXiv preprint arXiv:1610.02947*, 2016.
- [21] L. Zhu, Z. Xu, Y. Yang, and A. G. Hauptmann. Uncovering temporal context for video question and answering. *arXiv preprint arXiv:1511.04670*, 2015.