# Large-Scale Video Classification with Cloud Computing

Tyler Kohan, Murat Turkeli, Yunsheng Bai

# Introduction

- Kaggle contest: Google Cloud & YouTube-8M Video Understanding Challenge



$100,000 prize money

- Hosted by Google / Youtube
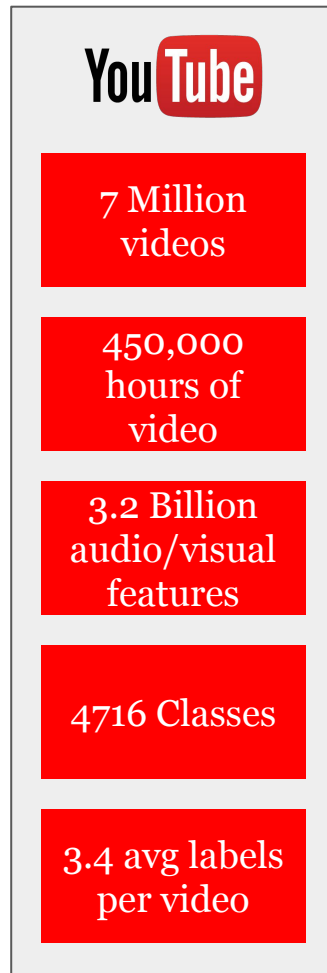


- Classify a video into multiple classes

# Youtube-8M v2 Video Dataset

- Training dataset: with labels, lots of videos
- Testing dataset: no labels, less videos
- Each video **:**
  - **rgb** + **audio** features
  - Multiple labels for each video



| 7 Million data point | | |
|---|---|---|
| Training (70%) | Eval (20%) | Test (10%) |

Train our models          Submission csv file

YouTube

7 Million videos

450,000 hours of video

3.2 Billion audio/visual features

4716 Classes

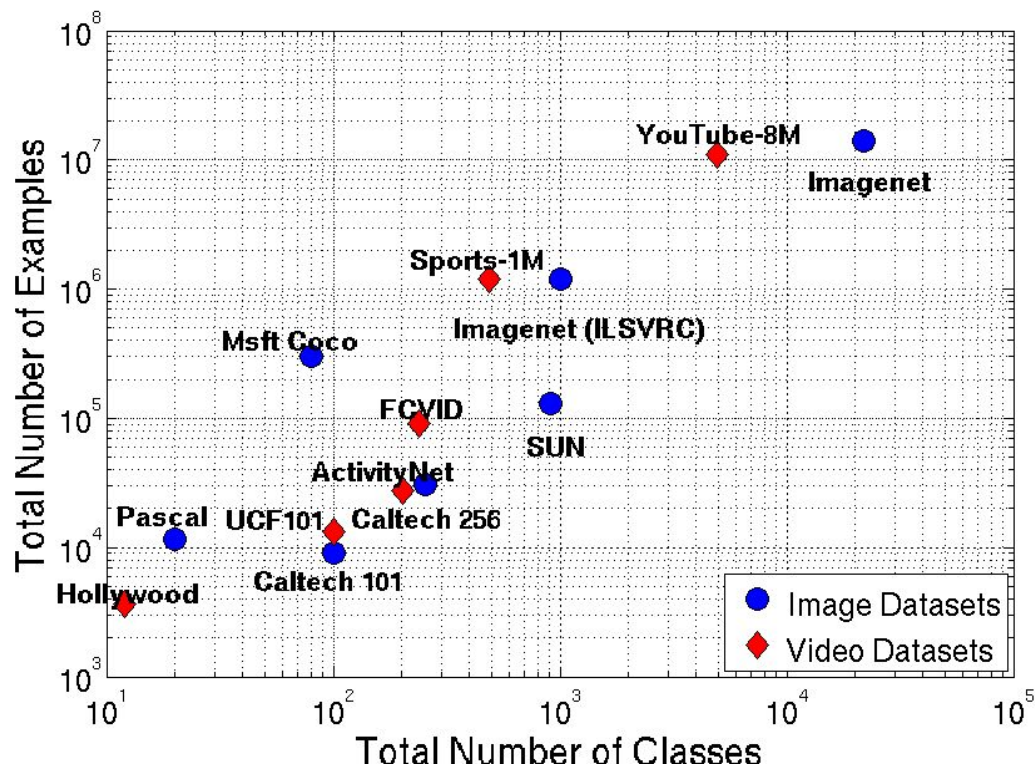3.4 avg labels per video

**~2 TB of data!**

Large datasets are good for advances in both image and video understanding tasks.

# Frame-level features dataset (1.71 TB)

**8 million**

```
video_id:   a0zzNorfSIw

 labels:   [48 10 71]

    rgb:   [[0 1 2 3 4 … 1023]

           [1 1 2 3 4 … 1023]

           ...

           [299 1 2 3 4 … 1023]]

  audio:   [[0 1 2 … 127]

           [1 1 2 … 127]

           ...

           [299 1 2 … 127]]
```

# Video-level features dataset (30 GB)

**8 million**

```
   video_id:   a0zzNorfSIw

     labels:   [48 10 71]

   mean_rgb:   [0 1 2 3 4 … 1024]

 mean_audio:   [0 1 2 .. 127]
```

**video_id:** `a0zzNorfSIw`

**labels:** `[48 10 71]`

**rgb:**
```
[[0 1 2 3 4 … 1023]
 [1 1 2 3 4 … 1023]
 ...
 [299 1 2 3 4 … 1023]]
```

**audio:**
```
[[0 1 2 … 127]
 [1 1 2 … 127]
 ...
 [299 1 2 … 127]]
```

ReLU activation of the last hidden layer

Inception network trained on ImageNet

Raw pixels

# Frame-level features dataset

# Video-level features dataset

```
video_id:  a0zzNorfSIw

  labels:  [48 10 71]

     rgb:  [        frame        ]
          [1 1 2 3 4 … 1023]

          ...

          299 1 2 3 4 … 1023]]

   audio:  [ 1 2 … 127]
          [ 1 2 … 127]
          .
          [ 1 2 … 127]]
```

300 frames

mean

Take the mean
(or std) across
frames

```
   video_id:  a0zzNorfSIw

     labels:  [48 10 71]

   mean_rgb:  [0 1 2 3 4 … 1024]

 mean_audio:  [ 1 2 .. 127]
```

# TensorFlow and Google Cloud ML

- **TensorFlow:** open source software library that makes it easy for us to perform complex machine learning concepts with limited knowledge
- **Google Cloud Machine Learning:** neural-net ML service that includes a platform to create our own training models





Used over $2500 of Google Cloud ML credit!

# Evaluation

| Video ID | Label Confidence Pairs |
|----------|------------------------|
| 100011194 | 1 0.983786 4 0.900343 297 0.891204 2292 0.792589 933 0.688224 ... |
| 100635497 | 92 0.716859 1 0.714576 926 0.422048 202 0.387686 4 0.254472 ... |
| ... | ... |

submission.csv

$$GAP = \sum_{i=1}^{N} p(i)\Delta r(i)$$

p(i) is precision, r(i) is the recall, and N is the number of videos

relevant elements

false negatives     true negatives

true positives     false positives

selected elements

How many selected items are relevant?     How many relevant items are selected?

Precision =     Recall =

# Global Average Precision

- Based on **precision** and **recall**
- Precision: out of the the labels we predicted, how many are correct
- Recall: out of all the actual labels, how many did we find
- Consider only up to 20 label/confidence pair per video



$$GAP = \sum_{i=1}^{N} p(i)\Delta r(i)$$

Corresponds to the **area** under the **curve**



relevant elements

false negatives | true negatives

true positives | false positives

selected elements

How many selected items are relevant?

How many relevant items are selected?

Precision =

Recall =

# Methods

## Models we identified as promising

| CNN | MOE | LSTM | Adaboost |
|-----|-----|------|----------|
| ~60 % | **~80%** | ~78% | ? |

### Mixture of Experts (MOE)



## Model development cycle

# Paper Reading and Knowledge Extraction

| # | Title | Summary | Extraction |
|---|-------|---------|------------|
| 1 | Learning Temporal Regularity in Video Sequences | Use an autoencoder to learn regularity in video sequences. | Fully connected neural networks could work. |
| 2 | Object Detection from Video Tubelets | Propose a special temporal convolutional neural network to incorporate | Object detection could assist classification. |
| 35 | Long-term Recurrent Convolutional Networks for Visual Recognition and | Propose a novel recurrent convolutional architecture based on RNN. The idea is to learn on sequential data sets to investigate events that are recurrent or temporally deep. | Interesting as it shows that learning sequentially (sequential video frames), we can improve on other methods. Also says the novel LRCN can be easily integrated in existing models. |

https://docs.google.com/document/d/1zYSWcDrX38v2glNRXyKmdLDINnkPCHrp01gyqaet2fs/edit?usp=sharing

The input data matrix has dimension (batch_size, 1152).

video

| 1 | 1 | 2 … 1023 | 0 | 1 | 2 … 127 |
|---|---|---|---|---|---|
| … | | | | | |
| 99 | 1 | 2 … 1023 | 0 | 1 | 2 … 127 |

mean_rgb    mean_audio

Expert 1
Expert 2
Expert 3
Gating Network

$\Sigma$

```
0.1     0.2
0.2     0.7
…
0.4715  0.1
0.1     0.8
0.1     0.5
…
0.4715  0.4
…
0.99    0.7
0.1     0.3
…
0.4715  0.1
```

```
0.01       0.18
0.04       0.56
…
0.4715²   0.05285
0.01       0.72
0.01       0.45
…
0.4715²   0.5285*0.4
…
0.99²      0.007
0.01       0.27
…
0.4715²   0.05285
```

reduce
_sum

```
0.19
0.6
…
0.52437
0.73
0.46
…
0.68292
…
0.999
0.28
…
0.52437
```

*

```
0   1   2 … 4715   0   1   2 … 4715
1   1   2 … 4715   0   1   2 … 4715
..
99  1   2 … 4715   0   1   2 … 4715
```

reshape

```
0     0
1     1
…
4715  4715
…
1     1
1     1
…
4715  4715
..
99    99
1     1
…
4715  44715
```

```
0.1      0.9
0.2      0.8
…
0.4715   0.5285
0.1      0.9
0.1      0.9
…
0.4715   0.5285
…
0.99     0.01
0.1      0.9
…
0.4715   0.5285
```

reshape

```
0.19   0.6   … 0.52437
0.73   0.46  … 0.68292
...
0.999  0.28  … 0.52437
```

Input:
batch_size by (1024+128=1152)

| 0 | 1 | 2 ... 1023 | 0 | 1 | 2 ... 127 |
| 1 | 1 | 2 ... 1023 | 0 | 1 | 2 ... 127 |
| ... | | | | | |
| 99 | 1 | 2 ... 1023 | 0 | 1 | 2 ... 127 |

mean_rgb      mean_audio

Each neuron has
A weight vector
$[w_0\ w_1\ \ldots\ w_{1151}]$
and a bias:
b
and an activation f
(ReLU)

| 0.1 | 0.2 |
| 0.2 | 0.7 |
| ... | |
| 0.4715 | 0.1 |
| 0.1 | 0.8 |
| 0.1 | 0.5 |
| ... | |
| 0.4715 | 0.4 |
| ... | |
| 0.99 | 0.7 |
| 0.1 | 0.3 |
| ... | |
| 0.4715 | 0.1 |

| 0 | 1 | 2 ... 4715 | 0 | 1 | 2 ... 4715 |
| 1 | 1 | 2 ... 4715 | 0 | 1 | 2 ... 4715 |
| ... | | | | | |
| 99 | 1 | 2 ... 4715 | 0 | 1 | 2 ... 4715 |

reshape

| 0 | 0 |
| 1 | 1 |
| ... | |
| 4715 | 4715 |
| ... | |
| 1 | 1 |
| 1 | 1 |
| ... | |
| 4715 | 4715 |
| ... | |
| 99 | 99 |
| 1 | 1 |
| ... | |
| 4715 | 44715 |

| 0.1 | 0.9 |
| 0.2 | 0.8 |
| ... | |
| 0.4715 | 0.5285 |
| 0.1 | 0.9 |
| 0.1 | 0.9 |
| ... | |
| 0.4715 | 0.5285 |
| ... | |
| 0.99 | 0.01 |
| 0.1 | 0.9 |
| ... | |
| 0.4715 | 0.5285 |

*

| 0.01 | 0.18 |
| 0.04 | 0.56 |
| ... | |
| $0.4715^2$ | 0.05285 |
| 0.01 | 0.72 |
| 0.01 | 0.45 |
| ... | |
| $0.4715^2$ | 0.5285*0.4 |
| ... | |
| $0.99^2$ | 0.007 |
| 0.01 | 0.27 |
| ... | |
| $0.4715^2$ | 0.05285 |

reduce
_sum

| 0.19 |
| 0.6 |
| ... |
| 0.52437 |
| 0.73 |
| 0.46 |
| ... |
| 0.68292 |
| ... |
| 0.999 |
| 0.28 |
| ... |
| 0.52437 |

reshape

| 0.19 | 0.6 | ... | 0.52437 |
| 0.73 | 0.46 | ... | 0.68292 |
| ... | | | |
| 0.999 | 0.28 | ... | 0.52437 |

Input:
batch_size by (1024+128=1152)

feature vector

| 1 | 1 | 2 | ... | 1023 | 0 | 1 | 2 | ... | 127 |
|---|---|---|-----|------|---|---|---|-----|-----|
| ... | | | | | | | | | |
| 99 | 1 | 2 | ... | 1023 | 0 | 1 | 2 | ... | 127 |

mean_rgb   mean_audio

neuron output

Each neuron has
A weight vector

weight vector

and a bias:

bias

and an activation f
(ReLU)

0.1    0.2
0.2    0.7
...
0.4715  0.1
0.1    0.8
0.1    0.5
...
0.4715  0.4
...
0.99   0.7
0.1    0.3
...
0.4715  0.1

0.1    0.9
0.2    0.8

0.01    0.18
0.04    0.56
...
$0.4715^2$  0.05285
0.01    0.72
0.01    0.45
...
$0.4715^2$  0.5285*0.4

*

reduce
_sum

0.19
0.6
...
0.52437
0.73
0.46
...
0.68292

| | 1 | 2 | ... | 4715 | 0 | 1 | 2 | ... | 4715 |
|---|---|---|-----|------|---|---|---|-----|------|
| 1 | 1 | 2 | ... | 4715 | 0 | 1 | 2 | ... | 4715 |
| ... | | | | | | | | | |
| 99 | 1 | 2 | ... | 4715 | 0 | 1 | 2 | ... | 4715 |

reshape

f(   feature vector   *   weight vector   +   bias   )=   neuron output

4715  44715

Input:
batch_size by (1024+128=1152)

feature vector

| 1 | 1 | 2 | ... | 1023 | 0 | 1 | 2 | ... | 127 |
| 99 | 1 | 2 | ... | 1023 | 0 | 1 | 2 | ... | 127 |

mean_rgb     mean_audio

neuron 0

Class "Dog"

Class "Bicycle"

Class "Gift"

Neuron 4716

Class "Dog"

Class "Bicycle"

Class "Gift"

Each neuron corresponds to a class. For example, neuron 1 corresponds to "Dog." Neurons 0 - neuron 4715 are expert 1; Neurons 4716 - neuron 9431 are expert 2. Each neuron produces one number.

| 0.1 | 0.2 |
| 0.3 | 0.7 |

0.4715 0.1

| 1 | 2 | ... | 4715 | 0 | 2 | ... | 4715 |
| 1 | 1 | 2 | ... | 4715 | 0 | 1 | 2 | ... | 4715 |
| 99 | 1 | 2 | ... | 4715 | 0 | 1 | 2 | ... | 4715 |

reshape

| 0 | 0 |
| 1 | 1 |
| ... | |
| 4715 | 4715 |
| ... | |
| 1 | 1 |
| 1 | 1 |
| ... | |
| 4715 | 4715 |
| .. | |
| 99 | 99 |
| 1 | 1 |
| ... | |
| 4715 | 44715 |

| 0.1 | 0.9 |
| 0.2 | 0.8 |
| ... | |
| 0.4715 | 0.5285 |
| 0.1 | 0.9 |
| 0.1 | 0.9 |
| ... | |
| 0.4715 | 0.5285 |
| 0.99 | 0.01 |
| 0.1 | 0.9 |
| ... | |
| 0.4715 | 0.5285 |

*

| 0.01 | 0.72 |
| 0.01 | 0.45 |
| ... | |
| $0.4715^2$ | $0.5285*0.4$ |
| ... | |
| $0.99^2$ | 0.007 |
| 0.01 | 0.27 |
| ... | |
| $0.4715^2$ | 0.05285 |

reduce _sum

0.19
0.6
...
0.52437
0.73
0.46
...
0.68292
...
0.999
0.28
...
0.52437

| 0.19 | 0.6 | ... | 0.52437 |
| 0.73 | 0.46 | ... | 0.68292 |
| ... | | | |
| 0.999 | 0.28 | ... | 0.52437 |

reshape

Input:
batch_size by (1024+128=1152)

video

| 1 1 2 ... 1023 | 0 1 2 ... 127 |
| ... | |
| 99 1 2 ... 1023 | 0 1 2 ... 127 |

mean_rgb    mean_audio

4716 neuros: Expert 1

4716 neuros: Expert 2

video

Expert 1    Expert 2

reshape

```
0.1    0.2
0.2    0.7
...
0.4715  0.1
0.1    0.8
0.1    0.5
...
0.4715  0.4
...
0.99    0.7
0.1    0.3
...
0.4715  0.1
```
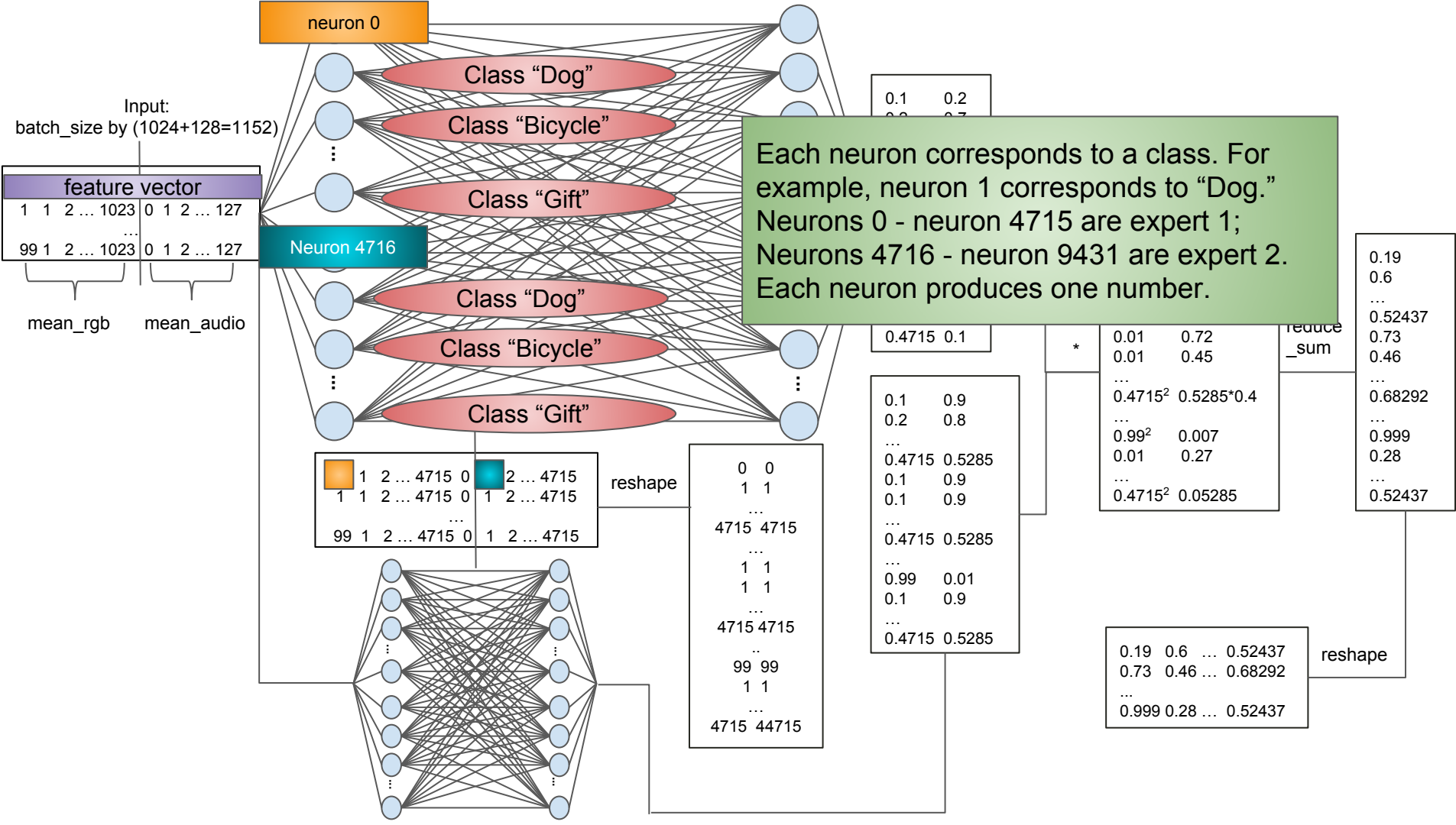
```
0.1    0.9
0.2    0.8
...
0.4715  0.5285
0.1    0.9
0.1    0.9
...
0.4715  0.5285
...
0.99    0.01
0.1    0.9
...
0.4715  0.5285
```

```
0    0
1    1
...
4715  4715
...
1    1
1    1
...
4715 4715
..
99   99
1    1
...
4715  44715
```

*

```
0.01    0.18
0.04    0.56
...
0.4715²  0.05285
0.01    0.72
0.01    0.45
...
0.4715²  0.5285*0.4
...
0.99²    0.007
0.01    0.27
...
0.4715²  0.05285
```

reduce_sum

```
0.19
0.6
...
0.52437
0.73
0.46
...
0.68292
...
0.999
0.28
...
0.52437
```

```
0.19  0.6  ... 0.52437
0.73  0.46 ... 0.68292
...
0.999 0.28 ... 0.52437
```

reshape

Input

| Expert 1 | o1 |
| Expert 2 | o2 |
| Expert 3 | o3 |

x    g1
x    g2
x    g3

Σ    o

Gating Network

Input:
batch_size by (1024+128=1152)

video 0
video 1
...
video 99

mean_rgb    mean_audio

|  | Exp 1 Prob | Exp 2 Prob |
|---|---|---|
| Class 0 | 0.1 | 0.2 |
| Class 1 | 0.2 | 0.7 |
| ... | | |
| Class 4715 | 0.4715 | 0.1 |
| Class 0 | 0.1 | 0.8 |
| Class 1 | 0.1 | 0.5 |
| ... | | |
| Class 4715 | 0.4715 | 0.4 |
| ... | | |
| | 0.99 | 0.7 |
| Class 0 | 0.1 | 0.3 |
| Class 1 | ... | |
| Class 4715 | 0.4715 | 0.1 |

video 0
video 1
...
video 1

video 0
video 1
video 99

0.19
0.6
...
0.52437
0.73
0.46
...
0.68292
...
0.999
0.28
...
0.52437

Expert 1    Expert 2

| 1 | 1 | 2 | ... | 4715 | 0 | 1 | 2 | ... | 4715 |
|---|---|---|---|---|---|---|---|---|---|
| 99 | 1 | 2 | ... | 4715 | 0 | 1 | 2 | ... | 4715 |

99 99
1 1
...
4715 44715

0.73  0.46 ... 0.68292
...
0.999 0.28 ... 0.52437

reshape

Each expert gives a probability.
How to combine them so that we have ONE
probability for each class and each video?
Instead of giving each expert 50% weight,
introduce **another neural network** to get a
weight matrix (called "**gating network**")!

Input:
batch_size by (1024+128=1152)

| 0 | 1 | 2 | ... | 1023 | 0 | 1 | 2 | ... | 127 |
| 1 | 1 | 2 | ... | 1023 | 0 | 1 | 2 | ... | 127 |
| ... |
| 99 | 1 | 2 | ... | 1023 | 0 | 1 | 2 | ... | 127 |

mean_rgb          mean_audio

This **expert network** has 2 experts.

Sigmoid Layer

The only difference between the expert and gating networks: the gating network uses softmax function:

$$\sigma(x_j) = \frac{e^{x_j}}{\sum_i e^{x_i}}$$

So that each row adds up to 1.

| Expert 1 | Expert 2 |
| 1 | 1 | 2 | ... | 4715 | 0 | 1 | 2 | ... | 4715 |
| ... |
| 99 | 1 | 2 | ... | 4715 | 0 | 1 | 2 | ... | 4715 |

reshape

Expert 1   Expert 2

| 4715 | 5 |
| ... |
| 1 | 1 |
| 1 | 1 |
| ... |
| 4715 | 4715 |
| ... |
| 99 | 99 |
| 1 | 1 |
| ... |
| 4715 | 44715 |

This **gating network** produces a gating matrix.

Softmax Layer

| **0.1** | **0.9** |
| 0.2 | 0.8 |
| ... |
| 0.4715 | 0.5285 |
| 0.1 | 0.9 |
| 0.1 | 0.9 |
| ... |
| 0.4715 | 0.5285 |
| ... |
| 0.99 | 0.01 |
| 0.1 | 0.9 |
| ... |
| 0.4715 | 0.5285 |

Gating matrix

| ... | |
| 0.4715² | 0.5285*0.4 |
| ... | |
| 0.4715² | 0.05285 |

For this video and class, assign **90%** weight to expert 2.

For this video and class, assign **10%** weight to expert 1.

| ... | |
| 0.68292 | |
| ... | |
| .999 | |
| .28 | |
| ... | |
| 0.52437 | |

| 0.19 | 0.6 | ... | 0.52437 |
| 0.73 | 0.46 | ... | 0.68292 |
| ... |
| 0.999 | 0.28 | ... | 0.52437 |

reshape

Input:
batch_size by (1024+128=1152)

| 0 | 1 | 2 | ... | 1023 | 0 | 1 | 2 | ... | 127 |
| 1 | 1 | 2 | ... | 1023 | 0 | 1 | 2 | ... | 127 |
| ... | | | | | | | | | |
| 99 | 1 | 2 | ... | 1023 | 0 | 1 | 2 | ... | 127 |

mean_rgb        mean_audio

This **expert network** has 2 experts.

Sigmoid Layer

**Expert Network**

Gating Network

**Exp 1 Prob**    **Exp 2 Prob**

| 0.1 | 0.2 |
| 0.2 | 0.7 |
| ... | |
| 0.4715 | 0.1 |
| 0.1 | 0.8 |
| 0.1 | 0.5 |
| ... | |
| 0.4715 | 0.4 |
| | |
| 0.99 | 0.7 |
| 0.1 | 0.3 |

**Exp 1 Weig**  15  **Exp 2 Weig**

| 0.1 | 0.9 |
| 0.2 | 0.8 |
| ... | |
| 0.4715 | 0.5285 |
| 0.1 | 0.9 |
| 0.1 | 0.9 |
| ... | |
| 0.4715 | 0.5285 |
| | |
| 0.99 | 0.01 |
| 0.1 | 0.9 |
| ... | |
| 0.4715 | 0.5285 |

Gating matrix

This **gating network** produces a gating matrix.

Softmax Layer

| Expert 1 | 5 | 0 | Expert 2 | |
| 1 | 1 | 2 | ... | 4715 | 0 | 1 | 2 | ... | 4715 |
| ... | | | | | | | | | |
| 99 | 1 | 2 | ... | 4715 | 0 | 1 | 2 | ... | 4715 |

reshape

Expert 1    Expert 2

| 4715 | 5 |
| ... | |
| 1 | 1 |
| 1 | 1 |
| ... | |
| 4715 | 4715 |
| .. | |
| 99 | 99 |
| 1 | 1 |
| ... | |
| 4715 | 44715 |

*

| 0.01 | 0.18 |
| 0.04 | 0.56 |
| ... | |
| $0.4715^2$ | 0.05285 |
| 0.01 | 0.72 |
| 0.01 | 0.45 |
| ... | |
| $0.4715^2$ | 0.5285*0.4 |
| ... | |
| $0.99^2$ | 0.007 |
| 0.01 | 0.27 |
| ... | |
| $0.4715^2$ | 0.05285 |

**reduce _sum**

| 0.19 |
| 0.6 |
| ... |
| 0.52437 |
| 0.73 |
| 0.46 |
| ... |
| 0.68292 |
| ... |
| 0.999 |
| 0.28 |
| ... |
| 0.52437 |

**Element-wise matrix multiplication** followed by **row summation** produces ONE probability for each class and each video.

Input:
batch_size by (1024+128=1152)

| 0 | 1 | 2 ... 1023 | 0 | 1 | 2 ... 127 |
| 1 | 1 | 2 ... 1023 | 0 | 1 | 2 ... 127 |
| ... | | | | | |
| 99 | 1 | 2 ... 1023 | 0 | 1 | 2 ... 127 |

mean_rgb    mean_audio

This **expert network** has 2 experts.

Sigmoid Layer

This **gating network** produces a gating matrix.

Softmax Layer

Expert 1    Expert 2

| 1 | 1 | 2 ... 4715 | 0 | 1 | 2 ... 4715 |
| ... | | | | | |
| 99 | 1 | 2 ... 4715 | 0 | 1 | 2 ... 4715 |

reshape

Expert 1    Expert 2

| 4715 | 5 |
| ... | |
| 1 | 1 |
| 1 | 1 |
| ... | |
| 4715 | 4715 |
| ... | |
| 99 | 99 |
| 1 | 1 |
| ... | |
| 4715 | 44715 |

Gating matrix

Exp 1 Prob    Exp 2 Prob

| **0.1** | **0.2** |
| 0.2 | 0.7 |
| ... | |
| 0.4715 | 0.1 |
| 0.1 | 0.8 |
| 0.1 | 0.5 |
| ... | |
| 0.4715 | 0.4 |
| ... | |
| 0.99 | 0.7 |
| 0.1 | 0.3 |

Exp 1 Weig    15    Exp 2 Weig

| **0.1** | **0.9** |
| 0.2 | 0.8 |
| ... | |
| 0.4715 | 0.5285 |
| 0.1 | 0.9 |
| 0.1 | 0.9 |
| ... | |
| 0.4715 | 0.5285 |
| ... | |
| 0.99 | 0.01 |
| 0.1 | 0.9 |
| ... | |
| 0.4715 | 0.5285 |

*

| **0.01** | **0.18** |
| 0.04 | 0.56 |
| ... | |
| $0.4715^2$ | 0.05285 |
| 0.01 | 0.72 |
| 0.01 | 0.45 |
| ... | |
| $0.4715^2$ | 0.5285*0.4 |
| ... | |
| $0.99^2$ | 0.007 |
| 0.01 | 0.27 |
| ... | |
| $0.4715^2$ | 0.05285 |

reduce _sum

| **0.19** |
| 0.6 |
| ... |
| 0.52437 |
| 0.73 |
| 0.46 |
| ... |
| 0.68292 |
| ... |
| 0.999 |
| 0.28 |
| ... |
| 0.52437 |

**Expert Network**

Gating Network

**0.1*0.1+0.2*0.9 =0.01+0.18=0.19**. Expert 1 and 2 **together** think that the **weighted average probability** that video 0 belongs to class 0 is pretty low.

# Mixture of Expert (MOE) model gives our current best result.

# Results



```
tensorboard --logdir=tensorboard
--logdir=gs://eecs3511_yt8m_train
_bucket/MoeModel_std_1_2_3_4_5_20
exp_lm/ --port=8080
```

# Links

- Dataset: https://research.google.com/youtube8m/
- Kaggle: https://www.kaggle.com/c/youtube8m
- Github: https://github.com/yunshengb/youtube-8m
- Paper: https://arxiv.org/pdf/1609.08675.pdf
- Paper Reading: https://docs.google.com/document/d/1zYSWcDrX38v2glNRXyKmdLDINnkPCHrp01gyqaet2fs/edit?usp=sharing