

# **YouTube-8M: A Large-Scale Video Classification Benchmark**

**S Abu-El-Haija (2016)**

**발표자: 김사라  
2017.04.24**

# YouTube-8M Dataset

- **A large-scale benchmark dataset for general multi-label video classification.**
- **Here, they treat the task of video classification as that of producing labels that are relevant to a video given its frame.**
- **YouTube-8M is not restricted to action classes alone.**

# Vertical

Arts & Entertainment ▾

# Filter

Guitar|

# Entities

- Acoustic guitar
- Cort Guitars
- Electric guitar
- Flamenco guitar
- Guitar**
- Guitar Center
- Guitar Hero
- Guitar Hero III: Legends of Rock
- Guitar amplifier
- Lead guitar
- PRS Guitars
- Pedal steel guitar
- Resonator guitar
- Steel guitar
- Steel-string acoustic guitar
- Twelve-string guitar
- Washburn Guitars

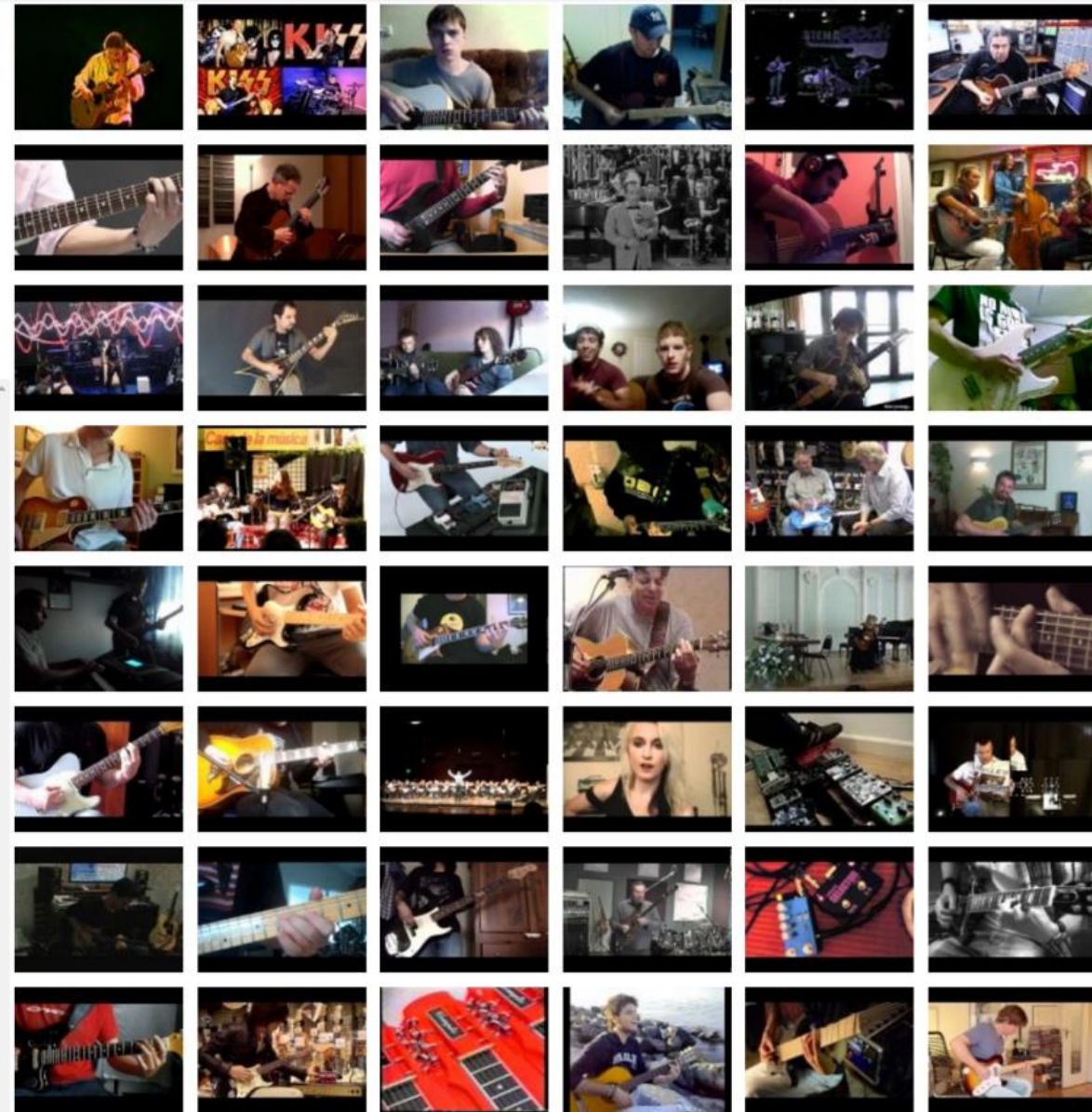
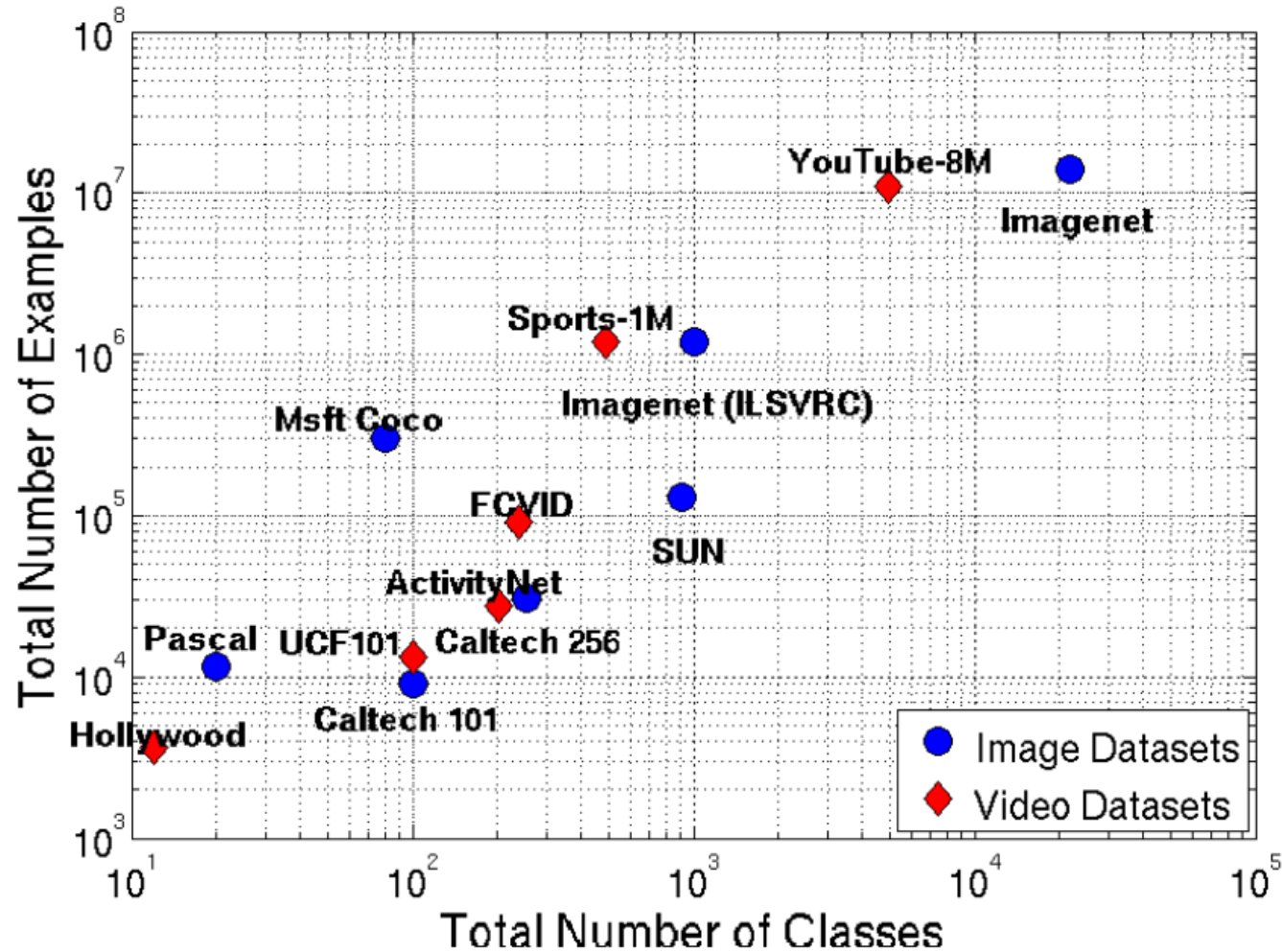


Figure1: This screenshot of a dataset explorer depicts a subset of videos in the dataset annotated with the entity “Guitar”.

# YouTube-8M Dataset

- YouTube-8M contains more than 8 million videos from 4,800 classes.



# YouTube-8M Dataset

- We use the YouTube video annotation system to obtain topic annotations for a video, and to retrieve videos for a given topics.
- The annotations are provided in the form of Knowledge Graph entities. We use Knowledge Graph entities to succinctly describe the main themes of a video, e.g., a video of biking on dirt roads and cliffs would have a central topic/theme of Mountain Biking, not Dirt, Road, Person, Sky, and so on.
- The aim of the dataset is not only to understand what is present in each frame of the video, but also identify the few key topics that best describe what the video is about.

## 1. Vocabulary Construction

- every label in the dataset should be distinguishable using visual information alone,
- each label should have sufficient number of videos for training models and for computing reliable metrics on the test set.

Entity Name	Entity URL	Entity Description
Thunderstorm	<a href="http://www.freebase.com/m/0jb2l">http://www.freebase.com/m/0jb2l</a>	A thunderstorm, also known as an electrical storm, a lightning storm, or a thundershower, is a type of storm characterized by the presence of lightning and its acoustic effect on the Earth's atmosphere known as thunder. The meteorologically assigned cloud type associated with the thunderstorm is the cumulonimbus. Thunderstorms are usually accompanied by strong winds, heavy rain and sometimes snow, sleet, hail, or no precipitation at all...

How difficult is it to identify this entity in images or videos (without audio, titles, comments, etc)?

- 1. Any layperson could
- 2. Any layperson after studying examples, wikipedia, etc could
- 3. Experts in some field can
- 4. Not possible without non-visual knowledge
- 5. Non-visual

(a) Screenshot of the question displayed to human raters.

## 2. Collecting Videos

- have at least 1,000 views.
- We exclude too short ( $< 120$  secs) or too long ( $> 500$  secs) videos.
- Randomly sample 10 million videos among them.
- Obtained all entities for the sampled 10 million videos using the YouTube video annotation system.
- Filtered out entities with less than 200 videos, and videos with no remaining entities. This reduced the size of our data to 8,264,650 videos.

## 3. Features

- We pre-process the videos and extract frame-level features using a state-of-art deep model: the publicly available Inception network trained on ImageNet.
- We decode each video at 1 frame-per-second up to the first 360 seconds, feed the decoded frames into the Inception network, and fetch the ReLu activation of the last hidden layer, before the classification layer.
- The feature vector is 2048-dimensional per second of video. Afterwards, we apply PCA (+ whitening) to reduce feature dimensions to 1024.



# YouTube-8M Dataset

- ✓ **While this removes motion information from the videos, recent work shows diminishing returns from motion features as the size and diversity of the video data increase.**
- ✓ **We hope to extend the dataset with audio and motion features in the futures.**

# YouTube-8M Dataset

- The YouTube-8M dataset contains 4,800 classes and a total of 8,264,650 videos.
- The average length of a video in the dataset is 229.6 seconds, which amounts to ~1.9 billion frames across the dataset.

Top-level Category	1 <sup>st</sup> Entity	2 <sup>nd</sup> Entity	3 <sup>rd</sup> Entity	4 <sup>th</sup> Entity	5 <sup>th</sup> Entity	6 <sup>th</sup> Entity	7 <sup>th</sup> Entity
<b>Arts &amp; Entertainment</b>	Concert	Animation	Music video	Dance	Guitar	Disc jockey	Trailer
<b>Autos &amp; Vehicles</b>	Vehicle	Car	Motorcycle	Bicycle	Aircraft	Truck	Boat
<b>Beauty &amp; Fitness</b>	Fashion	Hair	Cosmetics	Weight training	Hairstyle	Nail	Mascara
<b>Books &amp; Literature</b>	Book	Harry Potter	The Bible	Writing	Magazine	Alice	E-book
<b>Business &amp; Industrial</b>	Train	Model aircraft	Fish	Water	Tractor pulling	Advertising	Landing
<b>Computers &amp; Electronics</b>	Personal computer	Video game console	iPhone	PlayStation 3	Tablet computer	Xbox 360	Microsoft Windows
<b>Finance</b>	Money	Bank	Foreign Exchange	Euro	United States Dollar	Credit card	Cash
<b>Food &amp; Drink</b>	Food	Cooking	Recipe	Cake	Chocolate	Egg	Eating
<b>Games</b>	Video game	Minecraft	Action-adventure game	Strategy video game	Sports game	Call of Duty	Grand Theft Auto V
<b>Health</b>	Medicine	Raw food	Ear	Glasses	Injury	Dietary supplement	Dental braces
<b>Hobbies &amp; Leisure</b>	Fishing	Outdoor recreation	Radio-controlled model	Wedding	Christmas	Hunting	Diving
<b>Home &amp; Garden</b>	Gardening	Home improvement	House	Kitchen	Garden	Door	Swimming pool
<b>Internet &amp; Telecom</b>	Mobile phone	Smartphone	Telephone	Website	Sony Xperia	Google Nexus	World Wide Web
<b>Jobs &amp; Education</b>	School	University	High school	Teacher	Kindergarten	Campus	Classroom
<b>Law &amp; Government</b>	Tank	Firefighter	President of the U.S.A.	Soldier	President	Police officer	Fighter aircraft
<b>News</b>	Weather	Snow	Rain	News broadcasting	Newspaper	Mattel	Hail
<b>People &amp; Society</b>	Prayer	Family	Play-Doh	Human	Dragon	Angel	Tarot
<b>Pets &amp; Animals</b>	Animal	Dog	Horse	Cat	Bird	Aquarium	Puppy
<b>Real Estate</b>	House	Apartment	Condominium	Dormitory	Mansion	Skyscraper	Loft
<b>Reference</b>	Vampire	Bus	River	City	Mermaid	Village	Samurai
<b>Science</b>	Nature	Robot	Eye	Ice	Biology	Skin	Light
<b>Shopping</b>	Toy	LEGO	Sledding	Doll	Shoe	My Little Pony	Nike; Inc.
<b>Sports</b>	Motorsport	Football	Winter sport	Cycling	Basketball	Gymnastics	Wrestling
<b>Travel</b>	Amusement park	Hotel	Airport	Beach	Roller coaster	Lake	Resort
<b>Full vocabulary</b>	<b>Vehicle</b>	<b>Concert</b>	<b>Animation</b>	<b>Music video</b>	<b>Video game</b>	<b>Motorsport</b>	<b>Football</b>

Table 1: Most frequent entities for each of the top-level categories.

# YouTube-8M Dataset

- Frame-level data (total size of 1.71 TB)

Each video has

- a. “video\_id”: unique id for the video,
- b. “label”: list of labels of that video,
- c. Each frame has “rgb”: float array of length 1024,
- d. Each frame has “audio”: float array of length 128.

A video  $v$  is given by a sequence of frame-level features (rgb)  $x_{1:F_v}^v$ ,  
where  $x_j^v$  is the feature of the  $j$ -th frame from video  $v$ .

# YouTube-8M Dataset

- Video-level data (total size of 31GB)

Each video has

- a. “video\_id”
- b. “labels”
- c. “mean\_rgb”: float array of length 1028
- d. “mean\_audio”: float array of length 128

Formally, a video-level feature  $\phi(x_{1:F_v}^v)$  is a fixed-length representation at the video-level. We use a simple aggregation technique for getting these video-level representation:

1. First, second order and ordinal statistics
2. Feature normalization

- **Video-level data (total size of 31GB)**

1. **First, second order and ordinal statistics:**

we extract the mean  $\mu^v \in R^{1024}$  and the standard-deviation  $\sigma^v \in R^{1024}$ , and the top 5 ordinal statistics  $Top_5(x_{1:F_v}^v)$  for each dimension

$$\phi(x_{1:F_v}^v) = \begin{bmatrix} \mu(x_{1:F_v}^v) \\ \sigma(x_{1:F_v}^v) \\ Top_5(x_{1:F_v}^v) \end{bmatrix}$$

2. **Feature normalization:**

we subtract the mean  $\phi(\cdot)$  then use PCA to decorrelate and whiten the features.

# Evaluation Metric

- For each entity  $e$ , we consider  $(x_i, g_i^e)_{i=1, \dots, N}$  for a binary classifier, where  $N$  is the number of videos,  $g_i^e \in \{0, 1\}$  is the ground-truth which is 1 if labels of  $i$ -th video contains entity  $e$ , and zero otherwise.
- Mean Average Precision (mAP)
  - Let  $y_i$  be the annotation scores for  $i$ -th video.
  - For each entity  $e$ , we sort all the non-zero annotations according to the model score. (denoted as  $y_i^e$ )
  - At a given threshold  $\tau$ , the precision  $P^e(\tau)$  and recall  $R^e(\tau)$  are given by

$$P^e(\tau) = \frac{\sum_i I(y_i^e \geq \tau) g_i^e}{\sum_i I(y_i^e \geq \tau)},$$
$$R^e(\tau) = \frac{\sum_i I(y_i^e \geq \tau) g_i^e}{\sum_i g_i^e}.$$

- **Mean Average Precision (mAP)**
  - The average precision, approximating the area under the precision-recall curve, can then be computed as

$$AP^e = \sum_{j=1}^{10000} P(\tau_j) [R(\tau_j) - R(\tau_{j+1})],$$

where  $\tau_j = j/10000$ .

- The mean average precision is computed as the unweighted mean of all the per-class average precisions:

$$mAP = \frac{\sum_e AP^e}{\# \text{ of entities}}.$$