



# YOHO model for Audio Segmentation and Sound Event Detection

Davide Capone [SM3500601]    Enrico Stefanel [SM3500554]  
`{davide.capone, enrico.stefanel}@studenti.units.it`

Data Science and Scientific Computing Master's Course  
Department of Mathematics and Geosciences  
University of Trieste

A.Y. 2023–2024

# Contents

Introduction

YOHO model

Work done and results

Conclusions



# Audio Segmentation and Sound Event Detection

The goal of automatic sound event detection (SED) methods is to recognize what is happening in an audio signal and when it is happening<sup>1</sup>. In practice, the goal is to recognize at what temporal instances different sounds are active within an audio signal. An example of sound event detection is presented below.

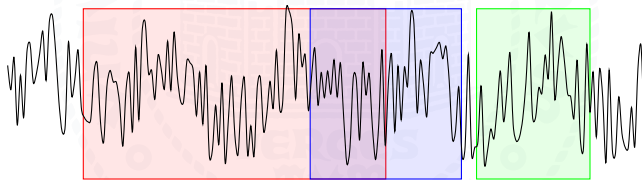


Figure 1: Event Detection in an audio track.

---

<sup>1</sup>Annamaria Mesaros et al. “Sound Event Detection: A tutorial”. In: *IEEE Signal Processing Magazine* 38 (2021), pp. 67–83. URL: <https://api.semanticscholar.org/CorpusID:235795366>.

## Datasets

Common datasets for Audio Segmentation and Sound Event Detection problems are:

- **TUT Sound Event Detection:** primarily consists of street recordings with traffic and other activity, with audio examples of 2.56 s and a total size of approximately 1.5 h. It has six unique audio classes – Brakes Squeaking, Car, Children, Large Vehicle, People Speaking, and People Walking;
- **Urban-SED:** purely synthetic dataset, with audio example of 10 s and a total size of about 30 h. It has ten unique audio classes – Air Conditioner, Car Horn, Children Playing, Dog Bark, Drilling, Engine Idling, Gun Shot, Jackhammer, Siren, and Street Music.

The first dataset is too small to train a Deep Neural Network model and requires use of augmentation techniques (like **SpecAugment**<sup>2</sup>).

<sup>2</sup>Daniel S. Park et al. “SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition”. In: *Interspeech 2019*. 2019, pp. 2613–2617. DOI: [10.21437/Interspeech.2019-2680](https://doi.org/10.21437/Interspeech.2019-2680).

## Metrics I

A popular toolbox for Sound Event Detection models evaluation is **SED Eval**<sup>3</sup>. The **F<sub>1</sub>-score** for a SED model prediction is defined as:

$$P = \frac{TP}{TP + FP}$$

$$R = \frac{TP}{TP + FN}$$

$$F_1\text{-score} = 2 \times \frac{P \times R}{P + R}$$

where TP, FP and FN are respectively (for each audio segment):

- the ground truth and system output both indicate an event to be active;
- the ground truth indicates an event to be inactive, but the system as active;
- the ground truth indicates an event to be active, but the system as inactive.

## Metrics II

The **Error rate** keeps in consideration *Substitutions*, *Insertions*, and *Deletions*:

$$S(k) = \min(\text{FN}(k), \text{FP}(k))$$

$$D(k) = \max(0, \text{FN}(k) - \text{FP}(k))$$

$$I(k) = \max(0, \text{FP}(k) - \text{FN}(k))$$

$$\text{ER} = \frac{\sum_{k=1}^K S(k) + \sum_{k=1}^K D(k) + \sum_{k=1}^K I(k)}{\sum_{k=1}^K N(k)}$$

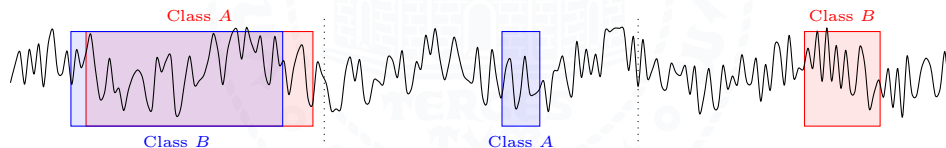


Figure 2: Example of **ground truth** and model **predicted labels** in three segments.

<sup>3</sup> Annamaria Mesaros, Toni Heittola, and Tuomas Virtanen. “Metrics for Polyphonic Sound Event Detection”. In: *Applied Sciences* 6.6 (2016). ISSN: 2076-3417. DOI: 10.3390/app6060162. URL: <https://www.mdpi.com/2076-3417/6/6/162>.

# YOHO model

Presented in 2021, **YOHO**<sup>4</sup> is a novel and lightweight real-time algorithm for *audio segmentation* and *sound event detection*:

- it aims to detect acoustic classes and their temporal boundaries by treating the problem as a **regression task**;
- inspired by *YOLO* algorithm for machine vision.

---

<sup>4</sup>Satvik Venkatesh, David Moffat, and Eduardo Reck Miranda. “You Only Hear Once: A YOLO-like Algorithm for Audio Segmentation and Sound Event Detection”. In: *Applied Sciences* 12.7 (Mar. 2022), p. 3293. ISSN: 2076-3417. DOI: 10.3390/app12073293. URL: <http://dx.doi.org/10.3390/app12073293>.

## Input shape

The network accepts as input a mel-spectrogram: a compressed time-frequency representation of audio signal that allows the model to better capture relevant audio patterns.

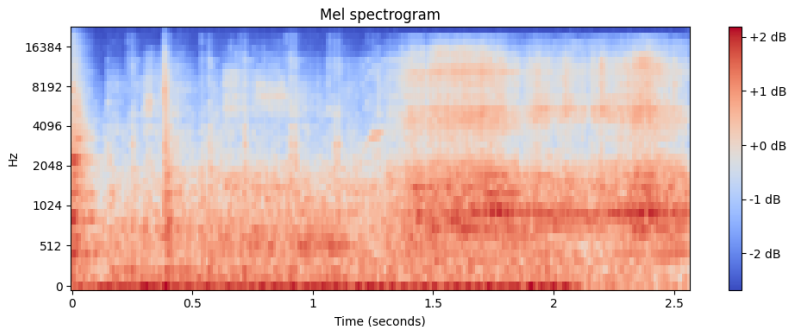


Figure 3: An example of mel-spectrogram.



# Network Architecture

The YOHO model is composed of three main parts:

- **Input layer:** mel-spectrogram derived from the audio signal;
- **MobileNet Backbone:** It consists of alternating depthwise and pointwise convolution layers, which significantly reduce the number of parameters and computational complexity. This efficient design makes it ideal for real-time applications where both speed and low resource consumption are critical;
- **Additional layers:** these are designed to refine the feature maps extracted by the backbone, adjusting the representation to better suit the sound event detection task;
- **Output layer:** a sigmoid function outputs the probability of each class being present in each time advancement with the indication of start and stop of the event.

The total number of trainable parameters is 3 930 590.

## Output shape

The output shape of the YOHO model is a tensor of shape (batch size, time steps, audio clas

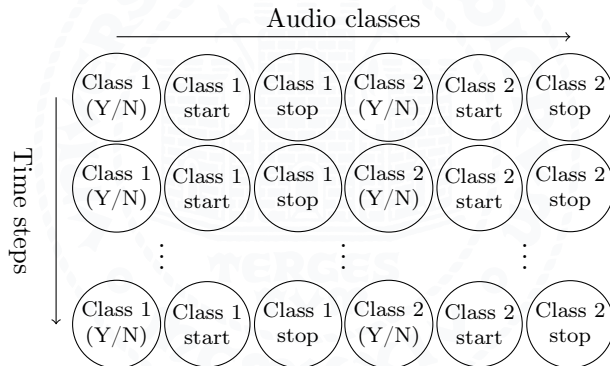


Figure 4: The YOHO output shape.

## Loss Function

$$\mathcal{L}_c(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + \\ (\hat{y}_2 - y_2)^2 + (\hat{y}_3 - y_3)^2 & \text{if } y_1 = 1 \\ (\hat{y}_1 - y_1)^2, & \text{if } y_1 = 0 \end{cases}$$

where  $y$  and  $\hat{y}$  are the ground-truth and predictions respectively.  $y_1 = 1$  if the acoustic class is present and  $y_1 = 0$  if the class is absent.  $y_2$  and  $y_3$ , which are the start and endpoints for each acoustic class are considered only if  $y = 1$ . In other words,  $(\hat{y}_1 - y_1)^2$  corresponds to **the classification loss** and  $(\hat{y}_2 - y_2)^2 + (\hat{y}_3 - y_3)^2$  corresponds to **the regression loss**.

## Other Details

All the **hyperparameters** used in the model training process were replicated from the original paper to **ensure consistency**:

- Learning Rate: 0.001;
- Optimizer: Adam with a weight decay of 0.01;
- Dropout Rate: 10 % (to prevent overfitting).

## Our 2 cents

We introduced a couple of improvements to the model, in order to make the training more efficient:

- use of **cosine annealing**<sup>5</sup> to dynamically adjust the learning rate during the training phase;
- performance improving while maintaining accuracy with **autocast**, that automatically chooses the precision for GPU operations;
- take advantage of GPUs high-performances using a **bigger batch size** (128 vs 32).

---

<sup>5</sup>Ilya Loshchilov and Frank Hutter. “SGDR: Stochastic Gradient Descent with Restarts”. In: *CoRR* abs/1608.03983 (2016). arXiv: 1608.03983. URL: <http://arxiv.org/abs/1608.03983>.

## Used technologies

Starting from the original paper, we implemented the system using PyTorch<sup>6</sup>, writing the code keeping in mind that it had to be **clear** and permit **reproducible tests**.

```
$ python3 -m yoho.train --help
usage: train.py [-h] [--name NAME] [--epochs EPOCHS] [--batch-size BATCH_SIZE] [--cosine-annealing]
               [--autocast] [--spec-augment]

options:
  -h, --help                show this help message and exit
  --name NAME                The name of the model
  --epochs EPOCHS           The number of epochs to train the model
  --batch-size BATCH_SIZE   The batch size for training the model
  --cosine-annealing         Use cosine annealing learning rate scheduler
  --autocast                Use autocast to reduce memory usage
  --spec-augment             Augment the training data using SpecAugment
```

Listing 1: Training script parameters

We used ORFEO<sup>7</sup> computational resources for the trainings of the models.

<sup>6</sup>All the code is available at <https://github.com/enstit/YOH024>.

<sup>7</sup><https://www.areasciencepark.it/piattaforme-tecnologiche/data-center-orfeo/>

## Implementation challenges

- A major limitation is that the original paper lacks detailed metrics or loss curves, preventing direct comparison and raising concerns about the **reproducibility** and **transparency** of their research.
- The code was poorly documented and didn't allow us to replicate their methods, hindering a clear understanding of their approach.
- One of the dataset used by the authors (used to test the algorithm on music–speech detection problem) was not publicly available due to copyrights (was manually recorder and labelled by the authors).

# Training results I

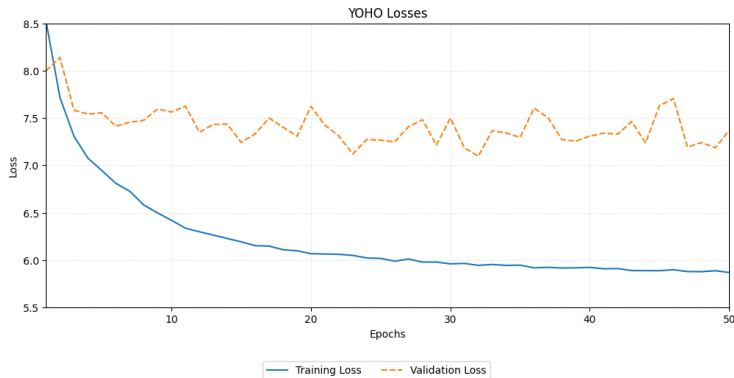


Figure 5: Training and validation loss for YOHO model on UrbanSED dataset.



## Training results II

We obtained an  $F_1$ -score of 54.23 and an error rate of 62.49 %.

Algorithm	$F_1$ -score
CRNN with envelope estimation	64.70
YOHO	59.50
CNN	56.88
CRNN	55.96
Our YOHO	54.23

Table 1: Segment-based overall  $F_1$ -score on the UrbanSED dataset.

We did bad, but not that bad after all...

# Project workflow

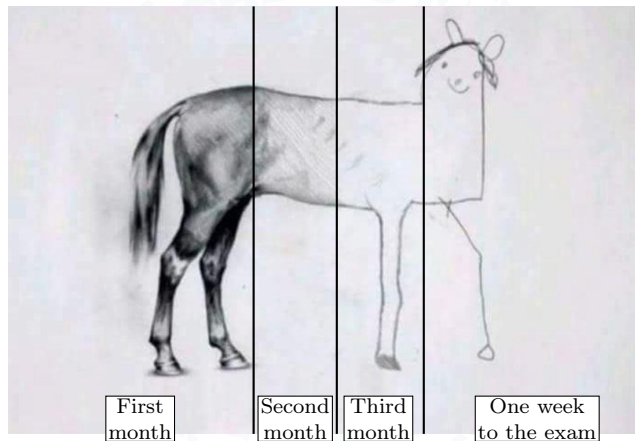


Figure 6: The roadmap of our journey.

## Lessons learned I

We learned many things during this project:

- ① how to exploit computer vision approaches on different input types (e.g., audio files);
- ② to hyper-tune High Performance Computer resources in a Deep Learning scenario to make the best we can with what we have in our hands;
- ③ not all papers found in the literature meet the necessary requirements for replicability. We faced several challenges during the implementation of the algorithm, but scientific researches should be accessible to anyone that has a basic understanding of what is being discussed;

## Lessons learned II

And, most important, we learned that...

*It's all about the journey, not the destination.*

Thank you for your attention.