

End-To-End Framework for Keyword Spotting

Stefano Ivancich

Department of
Information Engineering
University of Padova, Italy

Email: stefano.ivancich@studenti.unipd.it

Luca Masiero

Department of
Information Engineering
University of Padova, Italy

Email: luca.masiero.8@studenti.unipd.it

Abstract—In this paper, we present an end-to-end approach for the keyword spotting task, consisting of a sliding window of one second, a Voice Activity Detection module, a feature extraction module, a neural network and a fusion rule.

We explore the application of a 1D Convolutional Neural Network (CNN) that directly learns a representation from the audio signal, a Depthwise Convolution (DSConv) and an ensemble of the best models. The proposed models achieve a maximum accuracy of 96.6% with a small footprint of 127K parameters on the Google Speech Commands Dataset V2 (for the 10-commands and 21-commands recognition task).

We also compare the performances of our networks on different feature extraction methods such as 80 Mel, 40 MFCCs and 40MFCCs + deltas.

Index Terms—Human voice, Command recognition, Deep learning, Keyword spotting.

I. INTRODUCTION

The goal of keyword spotting is to detect a relatively small set of predefined keywords in a stream of user utterances, usually in the context of an intelligent agent on a mobile phone or a consumer “smart home” device.

Since the cloud-based interpretation of speech input requires transferring audio recordings from the user’s device, significant privacy implications should be kept in mind. Therefore, on-device keyword spotting has two main objectives: the first one concerns the recognition of common commands such as “on” and “off” (as well as other frequent words such as “yes” and “no”): this recognition can be accomplished directly on the user’s device; secondly, keyword spotting can be used to detect “command triggers”, such as “hey Siri”, that provides explicit hints for interactions directed at the device itself. It is additionally desirable that such models have a small footprint (for example, measured in the number of model parameters) so they can be deployed on low-power and performance-limited devices.

In the last years, neural networks have been shown to provide effective solutions to the small-footprint keyword spotting problem. Research typically focuses on a trade-off between achieving a high detection accuracy, with, on the other hand, a small footprint. Compact models are usually variants derived from a full model that sacrifices accuracy for a smaller model footprint.

In this work, we focus on **convolutional neural networks** (CNNs), a class of models that has been successfully applied to small-footprint keyword spotting; in particular, we explore

the use of **Depthwise Separable Convolutions (DSConv)** and **1DCNN**.

In this paper we will:

- present an end-to-end architecture for the keyword spotting task;
- try different audio features (raw data, 80 Mels, 40 MFCC and 40 MFCC with Deltas and Delta-Deltas);
- train different type of neural networks (1D CNN on raw data and DSConv CNN on features);
- build an ensemble of the best models;
- compare prediction speed and the number of parameters between the models.

The report is structured as follows. In Section II we present the current state of the art in the speech recognition field, in Section III we show our approach in order to tackle the problem and in Section IV we explain the pre-processing techniques used. In Section V we describe the various architectures we used, in Section VI we report their results and in Section VII we make some extra considerations on future developments and improvements.

II. RELATED WORK

The first system similar to a modern ASR was built in the mid-1950s by researchers at the Bell laboratories and was able to recognize numerical digits from speech using formants of the input audio. These formants are a concentration of the acoustic energy around a particular frequency in the input file wave.

For the next thirty years, various researchers developed devices capable of recognize vowels and consonants using different types of features (like phonemes), until the introduction, in the mid 1980s, of the Hidden Markov Models (HMMs). This approach, described in [1], represented a significant shift from simple pattern recognition methods (based on templates and a spectral distance measure) to a statistical method for speech processing and was possible thank to the incredible advancements in the computer computational power during those years.

In recent times, however, the HMMs were challenged by the introduction of Deep Learning and several architectures that work very well with these type of problems: Convolutional Neural Networks take advantage of weight-sharing and the convolution operation (which is shift-invariant in the data rep-

resentation domain) while Recurrent Neural Networks (RNN) are able to store a huge quantity of information.

Our two neural networks are inspired by [2] and [3]. In [2] the authors apply 1D CNN on raw signals for environment sound classification; we slightly modified the structure to adapt it to our task. In [3] the authors made some improvement on the MobileNet for Mobile Vision Applications using DSConv.

III. PROCESSING PIPELINE

The aim of the end-to-end architecture we propose, represented in **Figure 1**, is to handle audio signals of variable lengths and learn from the audio signal received.

A. Variable Audio Length

One of the main challenges of neural networks in audio processing is that the length of the input sample should be fixed, but the sound captured may have various duration. For this reason, it is necessary to adapt the model making it able to deal with audio signals of different lengths. Moreover, a model should be used for continuous prediction of input audio signals. One way to avoid this constraint, imposed by the input layer of the neural network, is to split the audio signal into different frames of fixed length using a sliding window of 1-second width.

For long contiguous audio recordings, instead of increasing the input dimension of the network (and consequently the number of parameters and its complexity), it is preferable to split the audio waveform into shorter frames: in this way, we keep the network compact so it can process audio waveforms of any length.

However, letting the model process each frame is too expensive (from the computational point of view). For this reason, we can use a VAD (Voice Activity Detection) or a Silence Filter to let the model process just the frames containing voice. This component should be very compact, not so expensive and fast; it could be part of a software or even a little piece of hardware (like some home device that we use nowadays). Only a batch of subsequent windows will be processed. For each of these windows a Feature Extraction module will extract 80 Mels, 40 MFCC, 40MFCC+deltas or no extraction at all and then this will be the input of our model.

B. Aggregation of Audio Frames

During the classification we need to aggregate the predictions to come up to a decision, as illustrated in **Figure 1**; for this reason, different fusion rules can be used to reach a final decision, such as the *majority vote* or the *sum rule*.

IV. SIGNALS AND FEATURES

As stated in Section I, we used the Google Speech Dataset V2. This dataset contains 105,829 audio files in .wav format (duration = one second) divided in thirty classes, plus an additional class containing five different type of noise of variable duration (about ten seconds for each file).

The core words are "Yes", "No", "Up", "Down", "Left", "Right", "On", "Off", "Stop", "Go", "Zero", "One", "Two",

"Three", "Four", "Five", "Six", "Seven", "Eight", and "Nine". To help distinguish unrecognized words, there are also ten auxiliary words, that most speakers only said once: "Bed", "Bird", "Cat", "Dog", "Happy", "House", "Marvin", "Sheila", "Tree", and "Wow".

We decided to create two different datasets:

- the first one is made of ten classes ("Yes", "No", "Up", "Down", "Left", "Right", "On", "Off", "Stop", "Go"), where each class contains approximately 3,600 files;
- the second is made of twenty-one classes ("Yes", "No", "Up", "Down", "Left", "Right", "On", "Off", "Stop", "Go", "Zero", "One", "Two", "Three", "Four", "Five", "Six", "Seven", "Eight", "Nine" and "Unknown"), where each class contains approximately 3,600 files, except for the "Unknown" class that contains all the remaining auxiliary words.

A. Dataset Partitioning

We divided both datasets into train-validation-test sets as suggested by Google Dataset's README file and we used a split of 30K-3K-3K for the ten classes dataset and 84K-9K-11K for the twenty-one classes dataset. The training set is used to train the network, while the validation one is used to evaluate the network performance during the training and to allow early-stopping to save the model with the lowest validation error, preventing overfitting. The test set is used to evaluate the network's final score that reaches the best performance.

B. Features extraction

We decided to study four different types of features in order to compare them and tried to see which was the best able to solve our problem:

- **No feature extraction:** we trained a specific model (1D CNN) directly on the raw signal (a 16,000 elements vector);
- **80 Mel spectrogram:** we took the Fourier Transform of (a windowed excerpt of) the raw signal and mapped the powers of the spectrum obtained into the mel scale (using triangular overlapping windows);
- **40 Mel-frequency cepstral¹ coefficients (MFCC):** those coefficients are obtained after the computation of the Discrete Cosine Transform (DCT) on the logarithm of the Mel spectrogram;
- **40 MFCC + Delta + Delta-Delta:** matrix containing the MFCCs previously computed, their first derivative (Delta) and second derivative (Delta-Delta).

We used the *Librosa Python Library* to compute those features, in particular we used the following functions:

- `librosa.feature.melspectrogram(...)`
- `librosa.power_to_db(...)`
- `librosa.feature.mfcc(...)`
- `librosa.feature.delta(...)`

¹The cepstrum is the result of the Fourier Transform applied to the signal spectrum, in decibel.

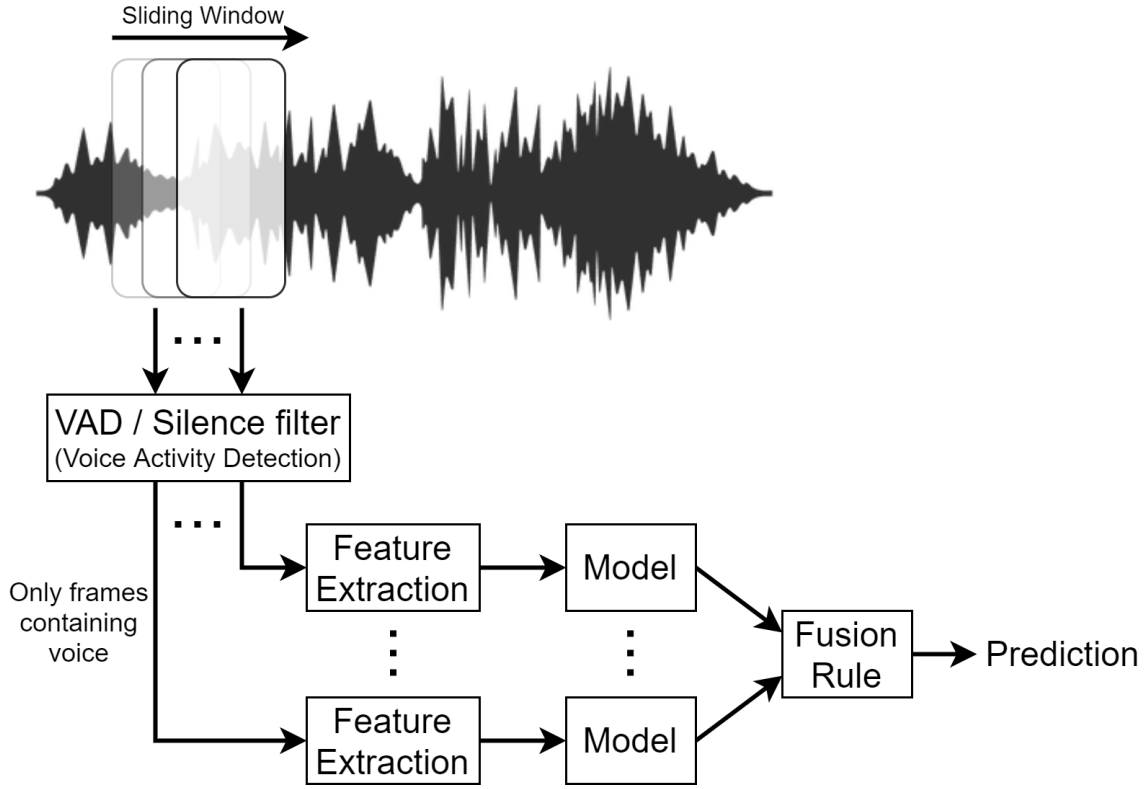


Figure 1: The end-to-end framework proposed consists of framing the input audio signal into several frames with appropriate overlapping percentage, a Voice Activity Detection module that filters out frames that do not contain voice, a feature extraction module, a neural network and a fusion rule.

V. LEARNING FRAMEWORK

A. 1D CNN on raw data

The proposed architecture (shown in **Figure 2** at page 4) is made of four *convolutional layers* interlaced with max pooling layers and followed by two *fully connected layers* and an output layer.

The input consists of an array of 16,000 dimensions, which represents a 1-second audio sample at 16kHz. The proposed 1D CNN has large receptive fields in the first convolutional layers since it is assumed that the first layers should have a more “global view” of the audio signal. The output of the last pooling layer, for all feature maps, is flattened and represents the input for a fully connected layer. In order to reduce overfitting, batch normalization is applied after the activation function of each convolution layer. After the last pooling layer, there are two fully connected layers with 128 and 64 neurons respectively; dropout is applied with a probability of 0.25 for both layers.

Keeping in mind the architecture shown in **Figure 2**, it is possible to omit a signal processing module because the network is strong enough to extract relevant low-level and high-level information from the audio waveform.

In the case of multiclass classification, the number of neurons of the output layer is the number of classes. Using the *softmax* as activation function for the output layer, each

output neuron indicates the membership degree of the input samples for each class.

B. DSConv Model

This model uses Depthwise Separable Convolutions. A separable convolution uses less parameters, less memory and less computations than regular convolutional layers, and in general it even performs better. There are two types of Separable convolution: Spatial and Depthwise.

Spatial Separable Convolutions divide a kernel into two smaller kernels (e.g. division of a 3×3 kernel into a 3×1 and 1×3 kernel). Instead of considering one convolution with 9 multiplications, we handled two convolutions with 3 multiplications each (6 in total) to achieve the same result. With a lower number of multiplications, the computational complexity significantly decreases, and the network is able to run and train faster.

Depthwise Separable Convolutions use kernels that cannot be “split” into two smaller kernels that apply two different convolutions: the *depthwise convolution*, which result will be the input for the *pointwise* (1×1) one.

In **Figure 3** (at page 4) you can see the difference between a normal 3D Convolution and a Depthwise one, and the final Depthwise Separable Convolutions block we used.

In our work we adopted a Depthwise Separable CNN inspired by [3]. This CNN is made of a 2D CNNs followed by batch

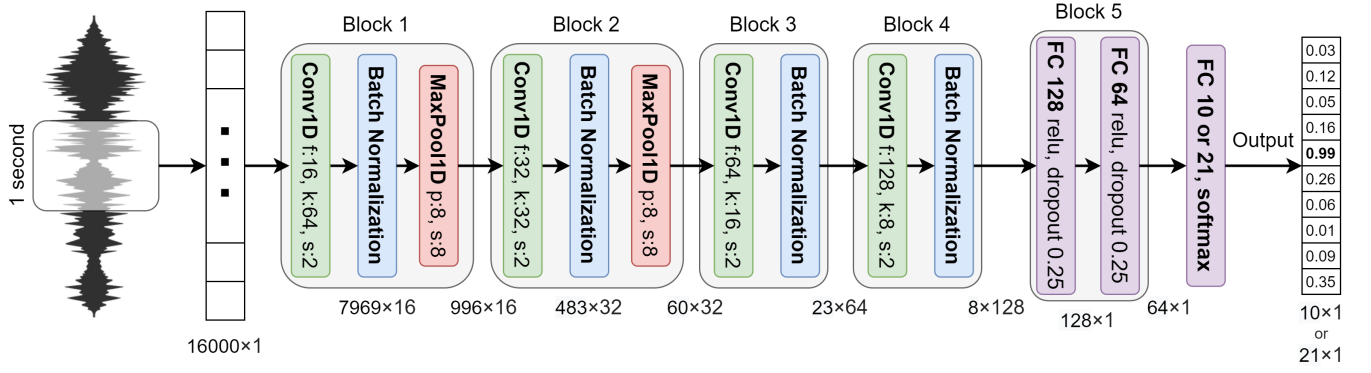


Figure 2: The architecture of the proposed 1D CNN model.

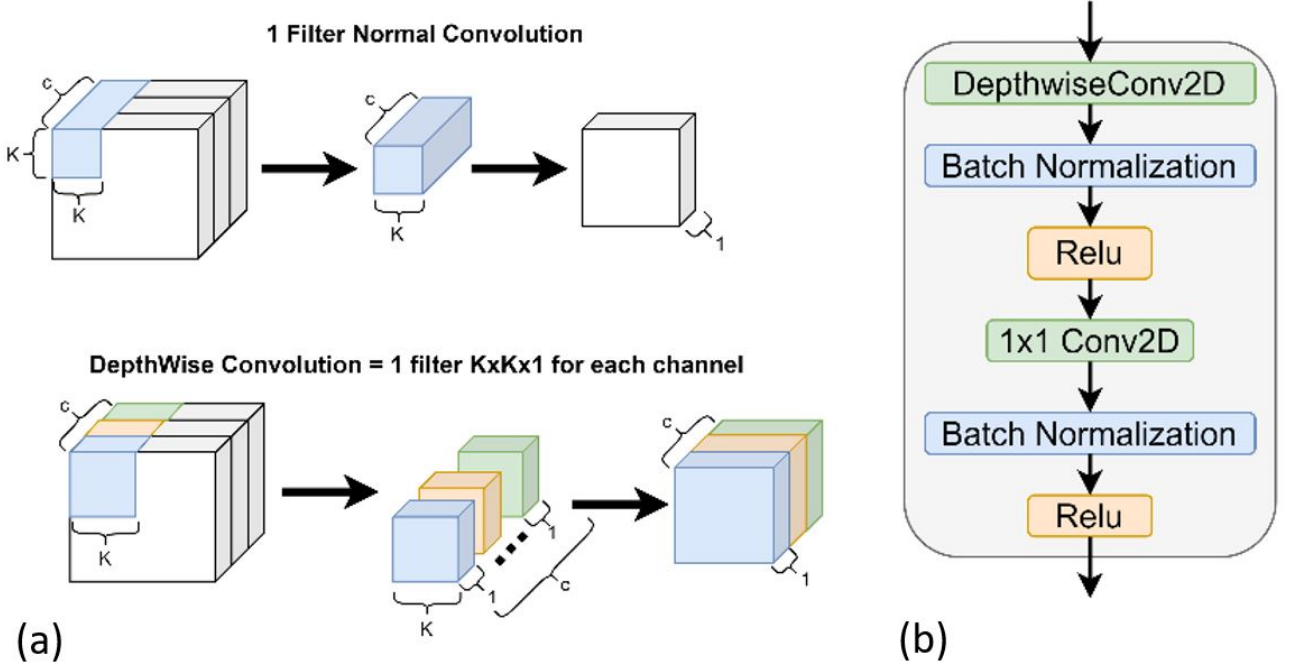


Figure 3a: Difference between normal convolution and depthwise convolution. Figure 3b: DSConv layer.

normalization and four or five DSConv2D layers. An average pooling, with padding 2×2 and stride 2×2 , followed by a fully-connected layer, is used in the end to provide global interaction and reduce the total number of parameters in the final layer. We implemented three different variants: *Large*, *Medium* and *Small*, described in the Tables at pages 5 and 6, where we tried to reduce the number of parameters still maintaining a good accuracy.

C. Ensemble

The ensemble simply took the two best models and computed the average of their final prediction, in this way we reached the best accuracy paying a high price: a large network.

All models are trained with the **Google Tensorflow 2.1** Framework with **Keras** using the standard *Sparse categorical cross entropy* and *Nadam optimizer*. With a batch size of 32, the models were trained for 100 epochs with initial learning

rate of 10^{-4} . An early stopping with a patience of 10 epochs was applied to avoid overfitting.

VI. RESULTS

We trained each DSConv network for all the features we described in Section IV to see if there was any difference between each set. As stated before, we used two different datasets, one made of ten class and another one of twenty-one.

We found that the number of convolutional layers played a key role in detecting high-level concepts, both for the DSConv models and the 1DCNN one. The number of convolutional layers in our models was determined in an experiment using the audio files of the Google Speech Commands V2 Dataset.

We faced the overfitting problem, recognizing the importance of the early-stopping procedure during the training comparing the test scores of the last epochs of the models with

Model Size		Layers					
		Conv2D	DSCConv2D 1	DSCConv2D 2	DSCConv2D 3	DSCConv2D 4	DSCConv2D 5
Small	# Filters	64	64	64	64	64	
	Filter size	10×4	3×3	3×3	3×3	3×3	NA
	Stride	2×2	1×1	1×1	1×1	1×1	
Medium	#Filters	172	172	172	172	172	
	Filter size	10×4	3×3	3×3	3×3	3×3	NA
	Stride	1×2	2×2	1×1	1×1	1×1	
Large	# Filters	276	276	276	276	276	276
	Filter size	10×4	3×3	3×3	3×3	3×3	3×3
	Stride	1×2	2×2	1×1	1×1	1×1	1×1

TABLE I: The configurations of the DSConv model layers considering different model sizes.

the best ones found by the `ModelCheckpoint` function of Keras (using the validation set). It took some time, for us, to choose the best metrics to use to evaluate our models, and at the end we decided *accuracy* was the best choice. We also measured the prediction speed of each model with and without Feature Extraction.

We also measured the prediction speed of each model with and without Feature Extraction. We should not consider multiplication as metric because it is an indirect alternative for a direct metric such as latency.

Furthermore, we noticed that predicting one example or a batch made of ten examples required almost the same amount of time: this happened because TensorFlow makes tensors operations. In the final implementation this fact should be taken into account in order to minimize the computational cost.

In Tables II, III, IV, V, VI you can see the accuracy and the speed prediction of our models: 40MFCCs + Deltas does not perform better, on the contrary the models have worse accuracy, a larger number of parameters and they are slower. There is not so much difference between using 80 Mels or 40 MFCCs, calculating MFCC requires slightly more time since they are calculated after the Mels but MFCC requires almost the half of parameters. The fastest network is the 1DCNN on raw data since it has not to calculate any feature extraction, but is the worst in terms of accuracy. The best accuracy is obtained by the ensemble network (the slower and larger model). The smallest model is the DSConvSmall + 40MFCCs with 127K parameters, 92.9% accuracy and 38.23 ms of prediction speed.

VII. CONCLUDING REMARKS

In this paper, an end-to-end architecture for the keyword spotting task has been proposed. We built two neural networks and an ensemble between the best models.

The 1DCNN consists of five convolutional layers. Instead of using handcrafted static filterbanks (such as the ones used to extract MFCC features), the proposed 1DCNN learns the filters directly from the audio waveform. The other network receives

	10 commands (30K-3K-3K)	21 commands (84K-9K-11K)
Accuracy %	93.0	89.1
# parameters	257,018	257,733
Speed (ms)	28.71	28.25

TABLE II: 1DCNN performances on 10 and 21-commands datasets.

as input features 80Mels, 40MFCCs or 40MFCCs + Deltas and consists of four or five DSConv layers. Furthermore, the 1DCNN is faster than the other models, but unfortunately does not reach the same accuracy.

As we can see from the results we obtained, all the features work well but 40MFCCs seems slightly better than the others, the 1DCNN on raw waveform performed the worst way. We then reduced the number of parameters of the DSConv model by lowering the number of the filters for each convolution. We proposed three version of the DSConv model: Large, Medium and Small. Our models achieve excellent results but we did not beat the state of the art accuracy of 97.4% reached in [4]. You can find a summary of our model compared with the state of the art in TABLE VII at page 7. We noticed that if we pick more frames per second, as described in the end-to-end architecture, the model has many “points of view” of the word, so in the real world it actually performs better than the test accuracy measured. In various tests we made, the networks (even the 1DCNN, the one with the worst accuracy) captured all the words we were saying thank to the multiple windows – fusion rule technique.

Further improvement could be exploring the SincConv, attention mechanism and LSTMs.

You can find the code on github [X]

	10 commands (30K-3K-3K)			21 commands (84K-9K-11K)		
	<i>80 Mels</i>	<i>40 MFCC</i>	<i>40 + deltas MFCC</i>	<i>80 Mels</i>	<i>40 MFCC</i>	<i>40 + deltas MFCC</i>
Accuracy %	96.0	95.3	93.9	93.4	boh	boh
# parameters	874,930	571,330	1,178,530	1,375,881	boh	boh
Speed (ms)	33.39	30.24	boh	33.79	boh	boh
Speed with FE	41.44	45.32	boh	41.62	boh	boh

TABLE III: DSConv Large performances on 10 and 21-commands datasets.

	10 commands (30K-3K-3K)			21 commands (84K-9K-11K)		
	<i>80 Mels</i>	<i>40 MFCC</i>	<i>40 + Δ MFCC</i>	<i>80 Mels</i>	<i>40 MFCC</i>	<i>40 + Δ MFCC</i>
Accuracy %	94.3	95.0	94.8	92.7	92.2	boh
# parameters	469,398	262,998	675,798	832,673	399,233	boh
Speed (ms)	30.75	29.55	boh	32.25	30.01	boh
Speed with FE	38.49	38.23	boh	39.72	46.76	boh

TABLE IV: DSConv Medium performances on 10 and 21-commands datasets.

	10 commands (30K-3K-3K)			21 commands (84K-9K-11K)		
	<i>80 Mels</i>	<i>40 MFCC</i>	<i>40 + Δ MFCC</i>	<i>80 Mels</i>	<i>40 MFCC</i>	<i>40 + Δ MFCC</i>
Accuracy %	92.5	92.9	boh	90.0	89.2	boh
# parameters	300,618	127,818	boh	604,757	241,877	boh
Speed (ms)	31.00	29.27	boh	32.86	29.48	boh
Speed with FE	38.31	38.23	boh	37.97	48.15	boh

TABLE V: DSConv Small performances on 10 and 21-commands datasets.

	10 commands (30K-3K-3K)			21 commands (84K-9K-11K)		
	<i>DSConv Large 80 Mels + DSConv Large 40 MFCC</i>	<i>boh</i>	<i>boh</i>	<i>DSConv Large 80 Mels + DSConv Large 40 MFCC</i>	<i>boh</i>	<i>boh</i>
Accuracy %	96.6	92.9	boh	90.0	89.2	boh
# parameters	1,446,260	127,818	boh	604,757	241,877	boh
Speed (ms)	89.28	boh	boh	boh	boh	boh

TABLE VI: Ensemble performances on 10 and 21-commands datasets.

WHAT WE LEARNT AS STUDENTS

We cannot say we faced an easy problem. With this project we learnt the meaning of the word "teamwork". It took a long time, for us, to read, understand, analyze and synthesize all the articles we read; we never did something like this before. We co-operated writing the code (run on an ASUS ROG GL703GM *Scar Edition*, with processor Intel Core i7-8750H Hexa-Core, Nvidia GeForce GTX 1060 and 16GB RAM memory), and in the meanwhile we started to think to what to write to present a well-written paper.

We learnt how to respect deadlines given the one to the other, **we learnt the one from the other.**

REFERENCES

- [1] L. R. Rabiner, "A tutorial on hidden markov models and selected applications in speech recognition," *Proceedings of the IEEE*, vol. 77, pp. 257–286, Feb 1989.
- [2] Sajjad Abdoli, Patrick Cardinal and Alessandro Lameiras Koerich. *End-to-End Environmental Sound Classification using a 1D Convolutional Neural Network*. arXiv:1904.08990v1 [cs.SD] 18 Apr 2019
- [3] Yundong Zhang, Naveen Suda, Liangzhen Lai and Vikas Chandra. *Hello Edge: Keyword Spotting on Microcontrollers*. arXiv:1711.07128v3 [cs.SD] 14 Feb 2018.
- [4] Simon Mittermaier, Ludwig Kurzinger, Bernd Waschneck and Gerhard Rigoll. *SMALL-FOOTPRINT KEYWORD SPOTTING ON RAW AUDIO DATA WITH SINC-CONVOLUTIONS*. arXiv:1911.02086v1 [eess.AS] 5 Nov 2019