

# End-to-End Framework for Keyword Spotting

## Students:

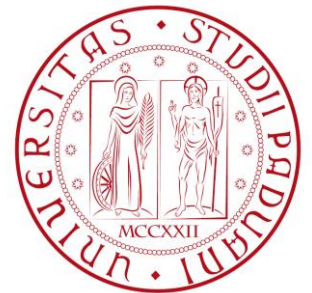
Ivancich Stefano

Masiero Luca

## Supervisors:

Prof. Rossi Michele

Meneghello Francesca



22 September 2020

- The Problem
- The Solution
- What we tried
- Architecture 1: 1DCNN on raw data
- Architecture 2: DSConv
- Architecture 3: Ensemble
- Performance comparisons vs. other papers
- Conclusions

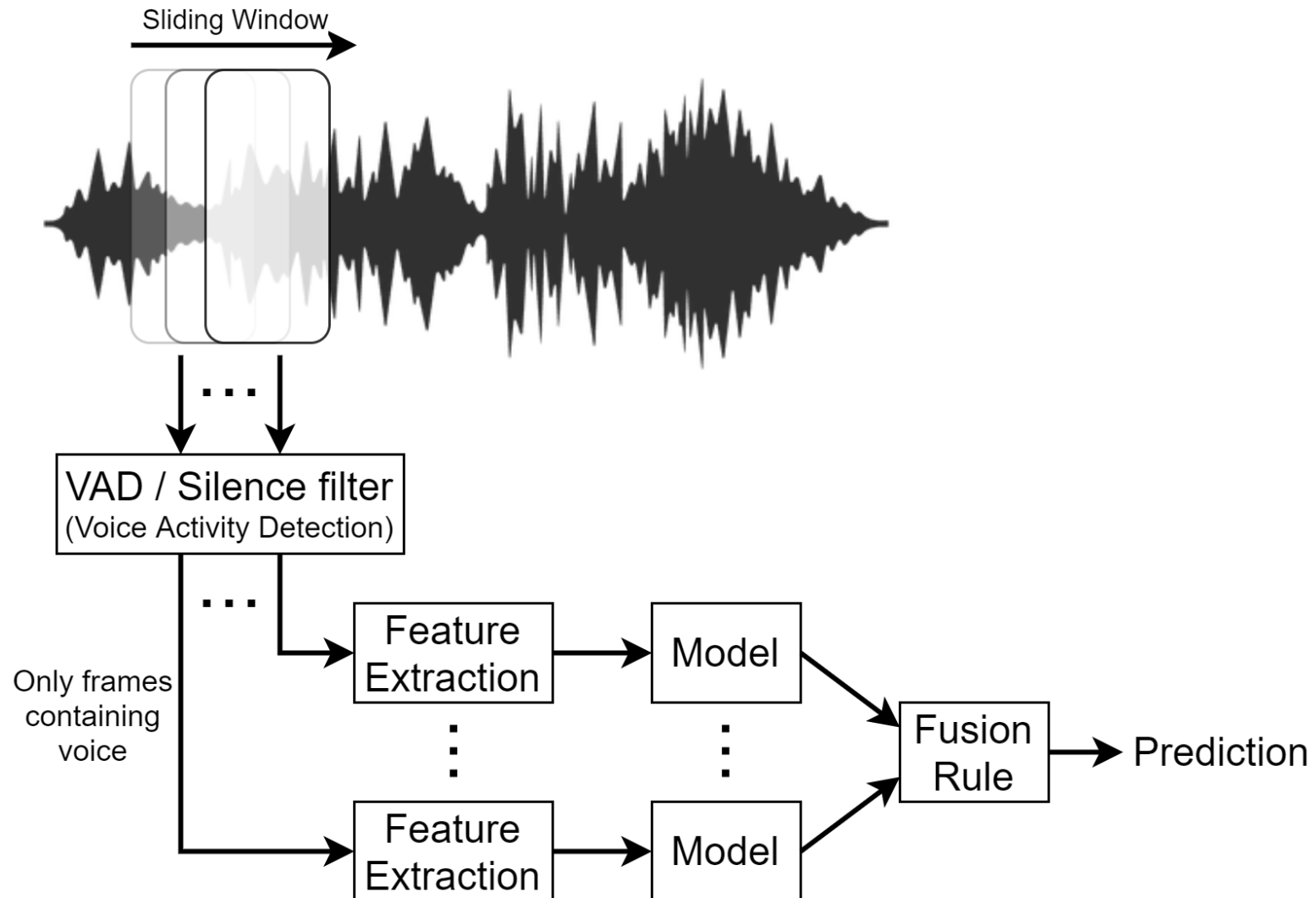
## Keyword Spotting

- Detect a relatively small set of predefined keywords (10 or 21) in a stream of user utterances.
- **Application:** Mobile phone, smart home device, consumer and robotics.
- **Constraints:** small footprint and fast (Real Time).

## Metrics

- Accuracy
- Number of parameters
- Prediction speed (milliseconds)

# The Solution



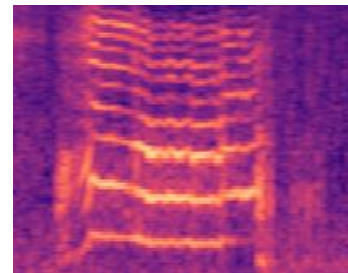
## Datasets (Google Speech Dataset V2)

- **10-commands** (“yes”, “no”, “up”, “down”, “left”, “right”, “on”, “off”, “stop”, “go”);
- **21-commands** ([...], “zero”, “one”, “two”, “three”, “four”, “five”, “six”, “seven”, “eight”, “nine”, unknown).

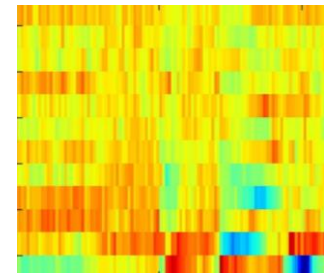
## Data Preprocessing techniques

- No preprocessing (Raw Waveform)
- 80 Mel spectrogram
- 40MFCC
- 40MFCC +  $40\Delta$  +  $40\Delta_s$  (=120)

80 Mels



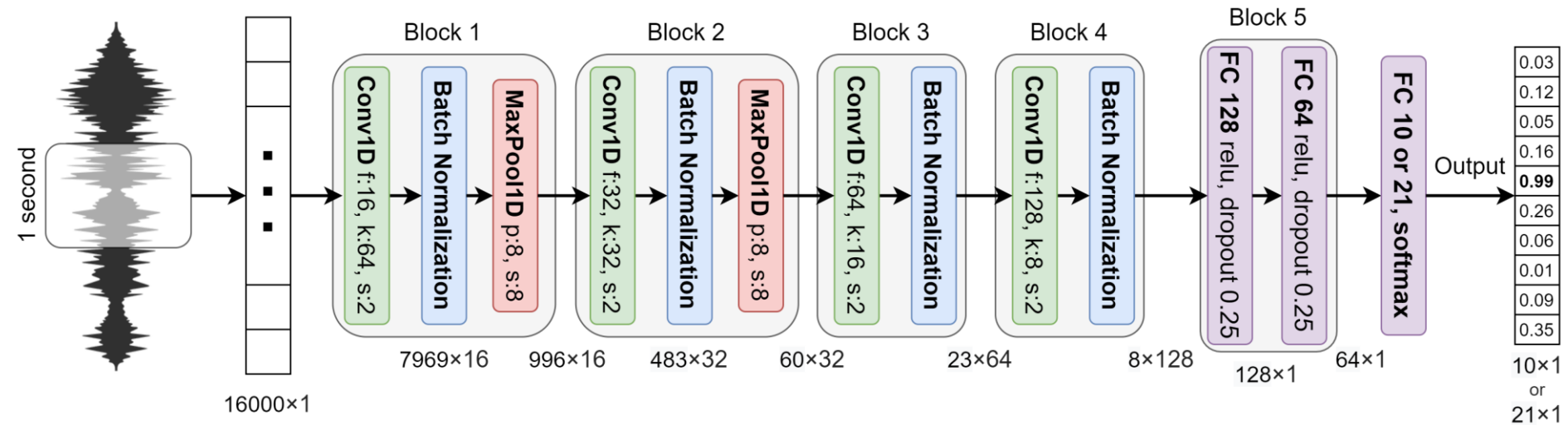
40 MFCC



## Learning Architectures

- 1DCNN on raw data
- DSConv (Small - Medium - Large)
- Ensemble (Small - Medium - Large)

# 1D CNN on raw data



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

# Separable Convolution

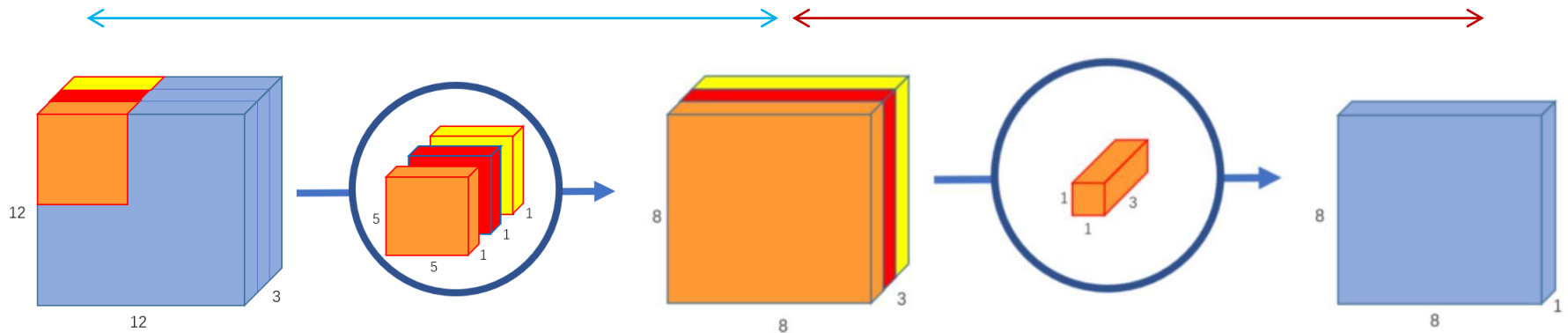
**Separable convolution** performs better than regular convolutional layers.

Two types of Separable convolution: Spatial and Depthwise.

**Depthwise Separable convolution** uses kernels that cannot be “factored” into two smaller kernels.

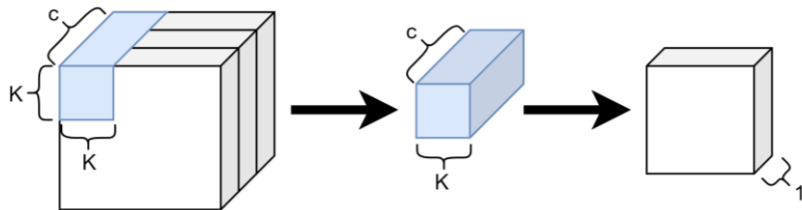
It splits a kernel into two separate kernels that do two convolutions:

- the **depthwise convolution**;
- the **pointwise (1x1) convolution**.



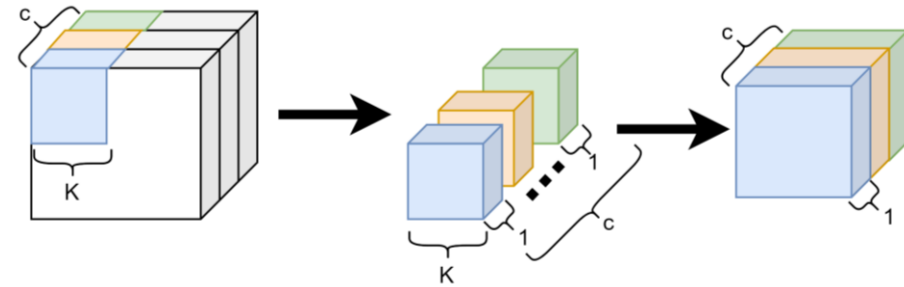
# Depthwise Separable Convolution

1 Filter Normal Convolution

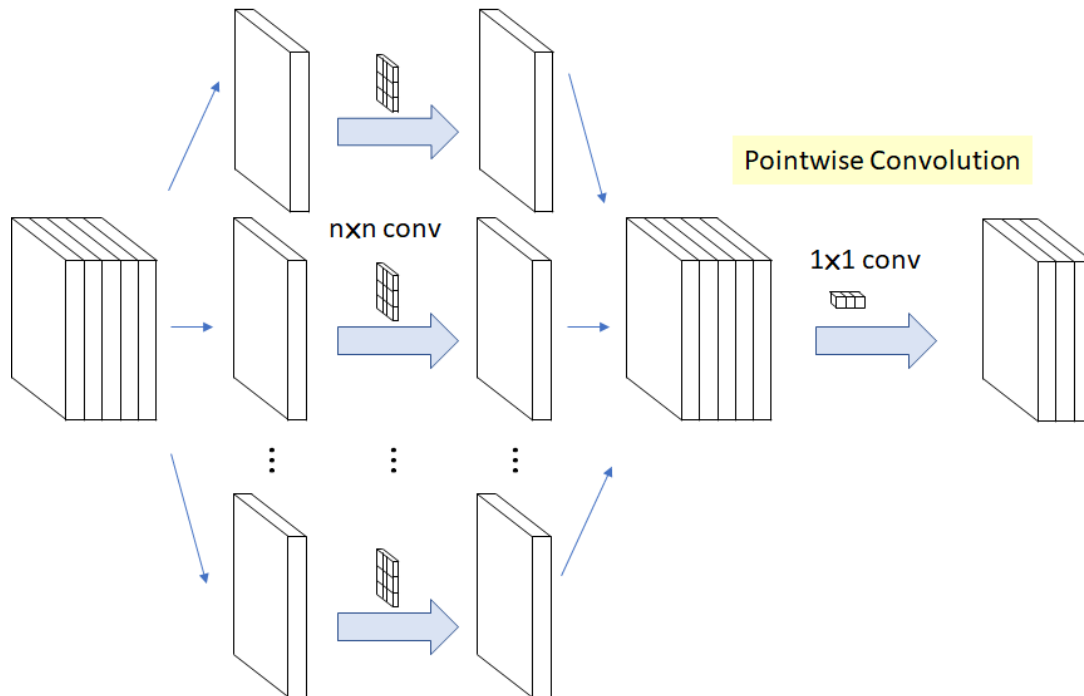


VS

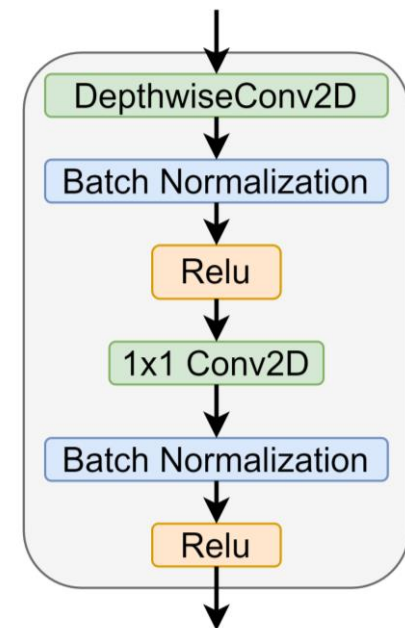
Depthwise Convolution =  $1 \times K \times K$  filter for each convolution



Depthwise Convolution

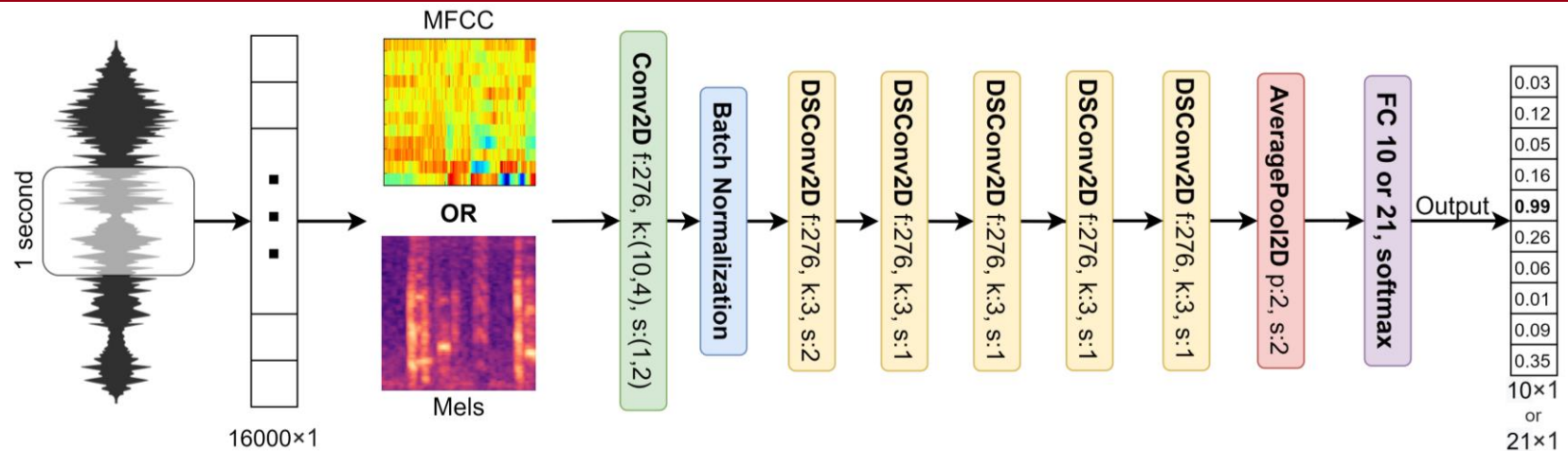


Separable Depthwise Convolution



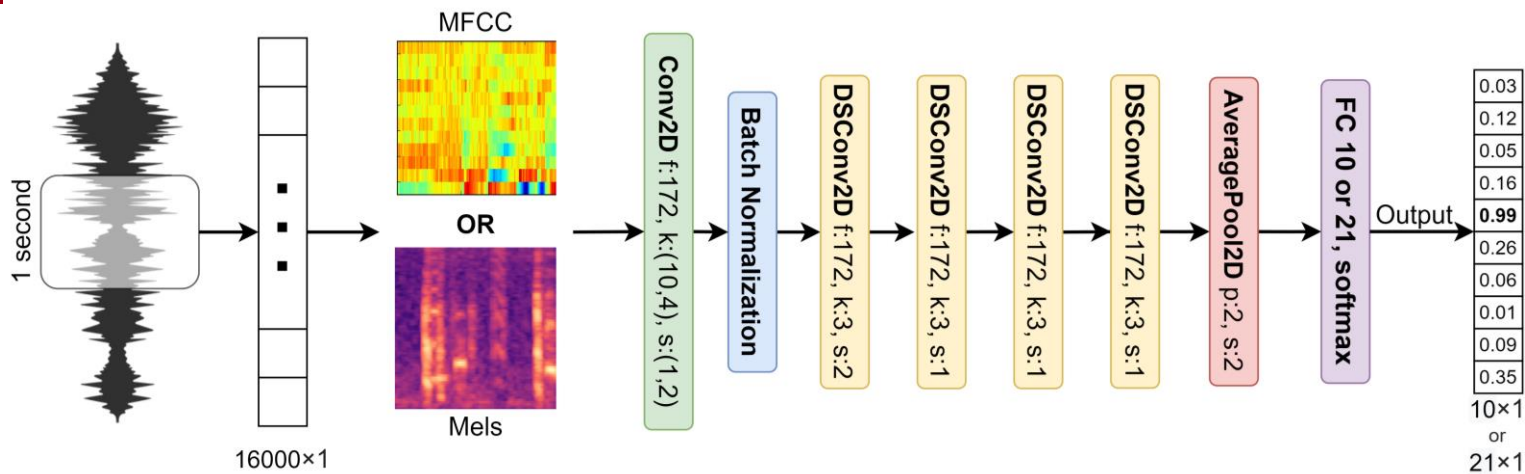


# DSCConv Model - Large



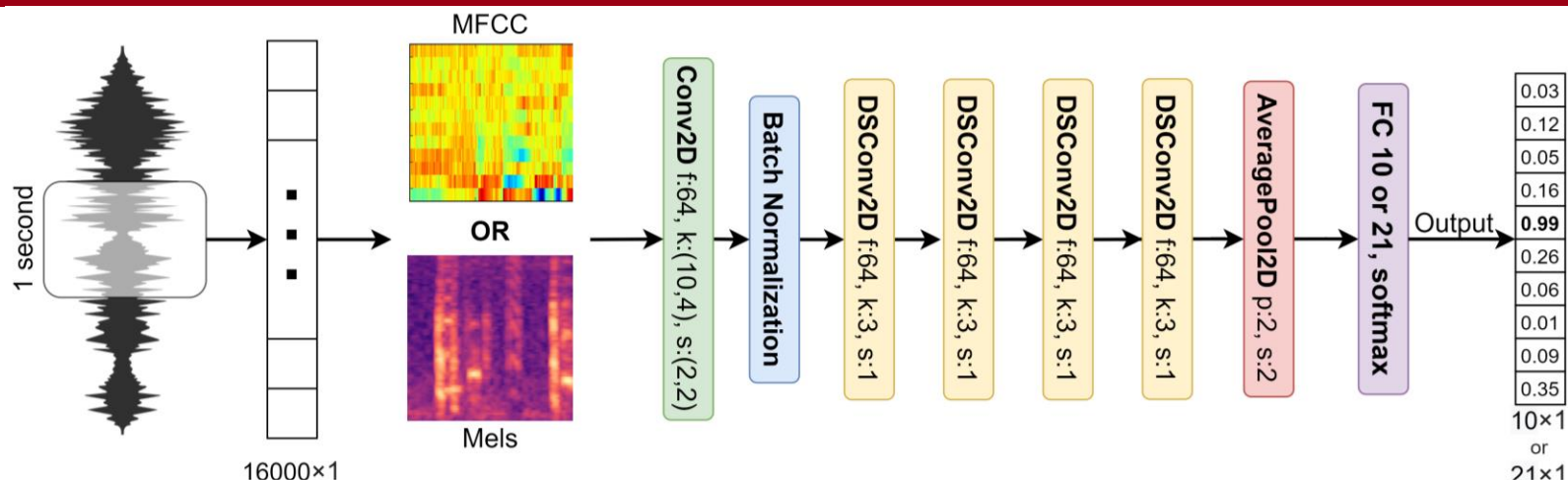
	10-commands (30K-3K-3K)			21-commands (84K-9K-11K)		
	80 Mels	40MFCC	40 + $\Delta_s$ MFCC	80 Mels	40MFCC	40 + $\Delta_s$ MFCC
Accuracy	<b>96.0</b>	95.3	95.5	93.4	<b>93.7</b>	92.7
# params	874,930	<b>571,330</b>	1,178,530	1,375,881	<b>738,321</b>	2,013,441
Speed (ms)	<b>33.39</b>	30.24	33.42	33.79	<b>30.87</b>	33.13
With FE	<b>41.44</b>	45.32	44.25	41.62	<b>39.66</b>	44.62

# DSCConv Model - Medium



	10-commands (30K-3K-3K)			21-commands (84K-9K-11K)		
	80 Mels	40MFCC	40 + $\Delta_s$ MFCC	80 Mels	40MFCC	40 + $\Delta_s$ MFCC
Accuracy	94.3	<b>95.0</b>	94.8	<b>92.7</b>	92.2	91.7
# params	469,398	<b>262,998</b>	675,798	832,673	<b>399,233</b>	1,266,113
Speed (ms)	30.75	<b>29.55</b>	30.76	<b>32.25</b>	30.01	33.03
With FE	38.49	<b>38.23</b>	41.82	<b>39.72</b>	46.76	47.08

# DSCConv Model - *Small*



	10-commands (30K-3K-3K)			21-commands (84K-9K-11K)		
	80 Mels	40MFCC	40 + $\Delta_s$ MFCC	80 Mels	40MFCC	40 + $\Delta_s$ MFCC
Accuracy	92.5	<b>92.9</b>	92.5	<b>90.0</b>	89.2	86.5
# params	300,618	<b>127,818</b>	473,418	604,757	<b>241,877</b>	967,637
Speed (ms)	31.00	<b>29.27</b>	30.54	<b>32.86</b>	29.48	30.31
With FE	38.31	<b>38.23</b>	40.23	<b>37.97</b>	48.15	40.99

# Ensemble: 10-commands

Ensemble between the best models.

	10-commands (30K-3K-3K)		
	Large	Medium	Small
	<i>DSCnv L 80 Mels</i> + <i>DSCnv L 40 MFCC</i> + <i>DSCnv L 40 MFCC <math>\Delta</math></i>	<i>DSCnv L 40 MFCC</i> + <i>DSCnv M 80 Mels</i> + <i>DSCnv M 40 MFCC <math>\Delta</math></i>	<i>DSCnv M 40 MFCC</i> + <i>DSCnv S 80 Mels</i> + <i>DSCnv S 40 MFCC</i>
<b>Accuracy</b>	96.8	96.4	95.6
<b># params</b>	2,624,790	1,303,726	691,463
<b>Speed (ms)</b>	131.01	122.04	106.09

# Ensemble: 21-commands

Ensemble between the best models.

	21-commands (84K-9K-11K)		
	Large	Medium	Small
	<i>DSCnv L 80 Mels</i> + <i>DSCnv L 40 MFCC</i> + <i>DSCnv L 40 MFCC <math>\Delta</math></i>	<i>DSCnv L 40 MFCC</i> + <i>DSCnv M 80 Mels</i> + <i>DSCnv M 40 MFCC</i>	<i>DSCnv L 40 MFCC</i> + <i>DSCnv M 40 MFCC</i> + <i>DSCnv S 40 MFCC</i>
<b>Accuracy</b>	95.2	95.0	94.2
<b># params</b>	2,498,019	1,970,227	1,379,431
<b>Speed (ms)</b>	125.9	114.21	109.03

# Performances: 10-commands

	Accuracy %	# Parameters	Speed (ms)
SincConv [Mittermaier et al.]	97.4	162K	40.35
<b>Our Ensemble Large</b>	<b>96.8</b>	2,624,790	131.01
Our Ensemble Medium	96.4	1,303,726	122.04
Our DSConvLarge + 80Mels	96.0	874,930	41,44
Our Ensemble Small	95.6	691,463	106.09
Our DSConvMedium + 40MFCC	95.0	262,998	38,23
<b>Our 1DCNN on raw data</b>	<b>93,0</b>	257,018	<b>28,71</b>
<b>Our DSConvSmall + 40MFCC</b>	<b>92,9</b>	<b>127,818</b>	<b>38,23</b>

- **Best:** Ensemble Large
- **Smaller:** DSConvSmall + 40MFCC
- **Fastest:** 1DCNN on raw data

# Performances: 21-commands

	Accuracy %	# Parameters	Speed (ms)
SincConv [Mittermaier et al.]	97.4	162K	40.35
<b>Our Ensemble Large</b>	<b>95.2</b>	2,498,019	125.9
Our Ensemble Medium	95.0	1,970,227	114.21
Our Ensemble Small	94.2	1,379,431	109.03
Our DSConvLarge + 40MFCC	93.7	738,321	39.66
Our DSConvMedium + 80Mels	92,7	832,673	39,72
Our DSConvSmall + 80Mels	90,0	604,757	37,97
<b>Our 1DCNN on raw data</b>	<b>89,1</b>	<b>257,733</b>	<b>28,25</b>

- **Best:** Ensemble Large
- **Smaller:** 1DCNN on raw data
- **Fastest:** 1DCNN on raw data

## Conclusions

- Tests: our models are **very good at classifying keywords**;
- We did not beat the state-of-the-art models;
- We found that the number of convolutional layers played a key role in detecting high-level concepts;
- No difference between 80 Mels or 40 MFCCs;
- Different model sizes in order to fit different devices.

## Future Work

- Try different hyperparameters during training;
- Change the **structure** of the **network** using:
  - SincConv;
  - GDSCov.