



An Experimental Study on Light Speech Features for Small-Footprint Keyword Spotting

Iván López-Espejo¹, Zheng-Hua Tan¹ and Jesper Jensen^{1,2}

¹Department of Electronic Systems, Aalborg University, Denmark

²Oticon A/S, Denmark

{ivl,zt,jje}@es.aau.dk, jesj@demant.com



CASPR

Centre for Acoustic Signal Processing Research

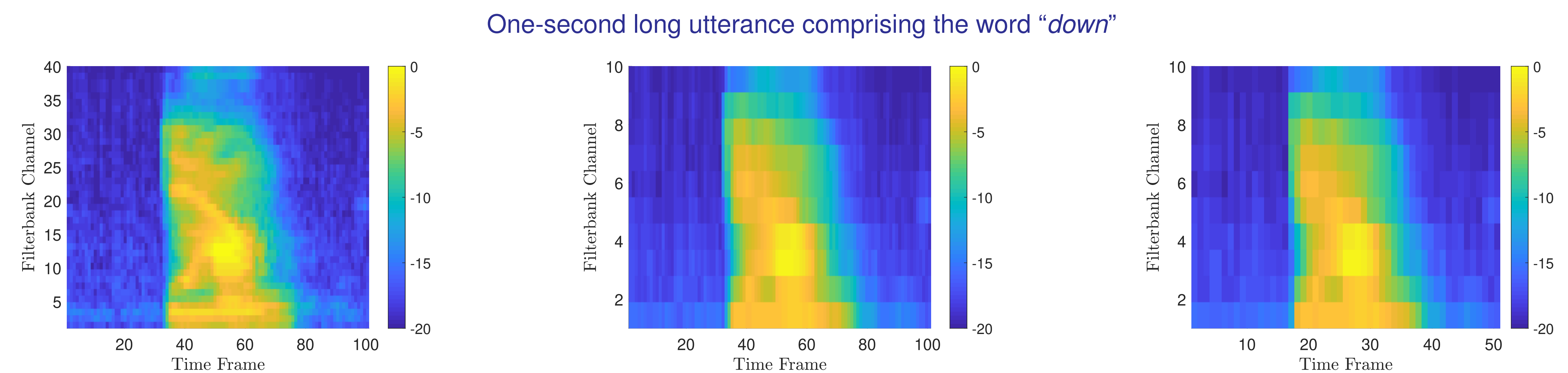
AALBORG UNIVERSITY
DENMARK

Introduction

- ▶ **Keyword spotting (KWS)** enables voice interaction with small devices like smartphones, tablets and smartwatches
- ▶ *Due to computational and energy constraints*, embedding typical always-on KWS technology can pose a challenge
 1. Acoustic model parameter **quantization**
 2. **Reduction of the number of parameters and/or multiplications** of the acoustic model
- ▶ There is much **redundant information** in speech features fed into modern KWS acoustic models
- ▶ **Feature matrix size reduction (2.)** → Remarkable **multiplication number reduction** in state-of-the-art acoustic models (*CNNs + residual connections*)

Methodology

- ▶ **Log-Mel and Mel-frequency cepstral coefficients (MFCCs)**
- ▶ Smaller no. of filterbank channels/cepstral coefficients (**features**) and time resolution reduction by increasing the analysis window hop size



Results and Conclusions

The hop size/number of time frames is 10 ms/101:

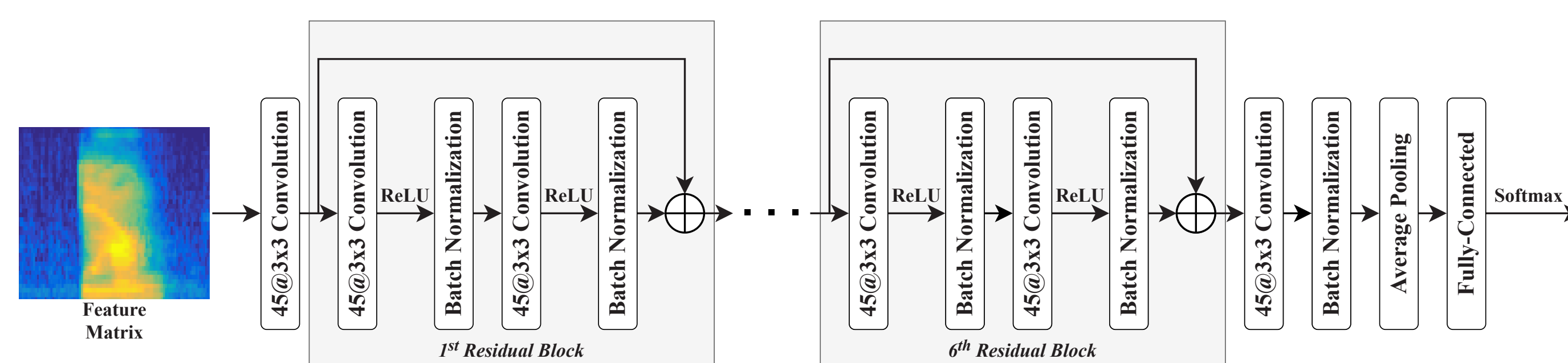
No. of Features	No. of Mult.	Log-Mel			MFCC		
		Accuracy (%)	Training Time per Epoch (s)	Inference Time (μ s)	Accuracy (%)	Training Time per Epoch (s)	Inference Time (μ s)
40	895M	95.33 \pm 0.28	98.7 \pm 0.4	980 \pm 7	95.24 \pm 0.96	96.5 \pm 0.4	980 \pm 12
20	424M	95.70 \pm 0.58	58.3 \pm 0.7	590 \pm 11	95.55 \pm 0.65	55.2 \pm 0.4	595 \pm 7
10	188M	95.34 \pm 0.76	36.2 \pm 0.3	390 \pm 21	95.24 \pm 0.64	36.3 \pm 0.3	385 \pm 12
5	71M	93.00 \pm 0.36	27.3 \pm 0.5	292 \pm 16	92.60 \pm 0.87	26.3 \pm 0.5	289 \pm 10

The number of features is 10:

Hop Size (ms) / No. of Frames	No. of Mult.	Log-Mel			MFCC		
		Accuracy (%)	Training Time per Epoch (s)	Inference Time (μ s)	Accuracy (%)	Training Time per Epoch (s)	Inference Time (μ s)
10 / 101	188M	95.34 \pm 0.76	36.2 \pm 0.3	390 \pm 21	95.24 \pm 0.64	36.3 \pm 0.3	385 \pm 12
20 / 51	93M	94.63 \pm 0.65	24.2 \pm 0.5	265 \pm 11	94.61 \pm 0.89	24.1 \pm 0.3	265 \pm 9
30 / 34	61M	94.53 \pm 0.47	19.2 \pm 0.4	216 \pm 10	93.50 \pm 0.83	19.0 \pm 0.4	216 \pm 14
40 / 26	46M	93.24 \pm 0.50	15.3 \pm 0.4	197 \pm 9	92.36 \pm 0.55	15.2 \pm 0.3	201 \pm 8

Experimental Framework

- ▶ We use the popular **Google Speech Commands Dataset** defining **10 keywords**: "yes", "no", "up", "down", "left", "right", "on", "off", "stop" and "go"
- ▶ **Skip connections** make the number of multiplications of the acoustic model heavily depend on the feature matrix size, since the successive feature maps have to preserve it for addition throughout the residual blocks



- ▶ **The final fully-connected layer has 11 nodes** corresponding to the 10 different keywords plus the non-keyword/filler class

- ▶ **Strong linear relationship** ($R^2 = 0.9641$, $p = 0.0001$) between the no. of multiplications of the acoustic model and its energy consumption
- ▶ The standard size of both **log-Mel and MFCC matrices can be reduced by a factor of 8** while essentially maintaining KWS performance
 1. 9.6 \times number of multiplications/energy consumption reduction
 2. 4.0 \times training time reduction
 3. 3.7 \times inference time reduction
- ▶ This is a finding to bear in mind when designing *light and compact KWS systems* intended to be embedded on low-resource devices