

Exploring Methods of Signal Processing for Keywords Spotting in Hearing Devices

Cypress Payne¹, Jhyv Philor², Tadesse Ghirmai³, Kaibao Nie³

¹ Seattle Pacific University, ² University of Florida, ³ University of Washington Bothell



INTRODUCTION

Problem: Changing settings in hearing aids and cochlear implants can be inconvenient by hand. Using keywords spotting (KWS) allows parameter adjustment to be done by voice.

Current feature extraction methods for KWS:

- o Bark spectrum, Mel spectrum, Mel frequency cepstral coefficients (MFCC)

Current deep learning techniques for KWS:

- o Recurrent Neural Network (RNN), Convolutional Neural Networks (CNN), Long-short term Memory (LSTM) [1]

Proposed Method: Applying amplitude modulation on an audio signal to extract information to train a CNN for keywords spotting. Amplitude modulation is vital for speech recognition because at certain frequencies useful speech information can be extracted [2] [3].

We compared the accuracy of a CNN using amplitude modulated signals with different bands versus Mel spectrum-based signal processing.

METHODS

Google speech dataset

- o 65000 different utterances of words were split into training, validation, and testing sets [4].
- o 11 keywords, consisting of “up”, “down”, and “one” through “nine” were selected for volume control

Preprocessing

- o An audio feature extractor in MATLAB was used to convert each audio signal into a Mel spectrogram.
- o Amplitude modulation was obtained by passing the signal through multiple contiguous bandpass filters and performing a Hilbert Transform to extract temporal envelopes (Figure 1).

Neural Network

- o A CNN with 3 hidden layers was setup in MATLAB (Figure 1).

EXTRACTION

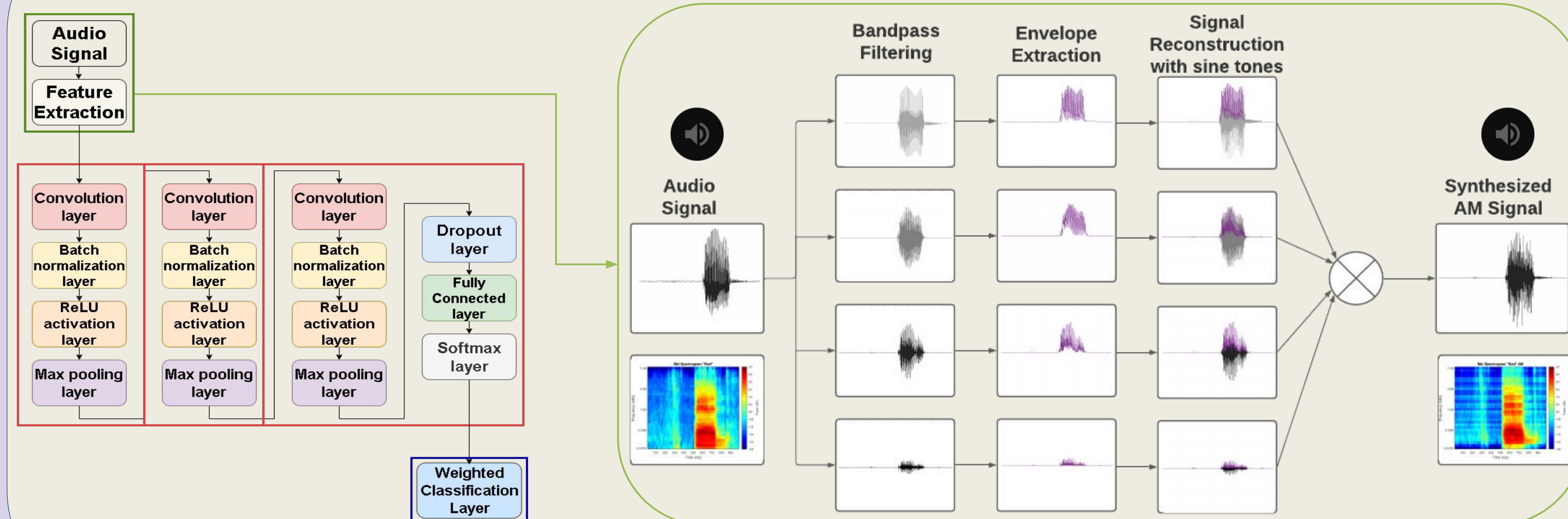


Figure 1: The left of the diagram outlines the entire convolutional neural network, with input, hidden layers, and output. The right displays the feature extraction algorithm that performs amplitude modulation on the signal before it is fed to the neural network.

RESULTS

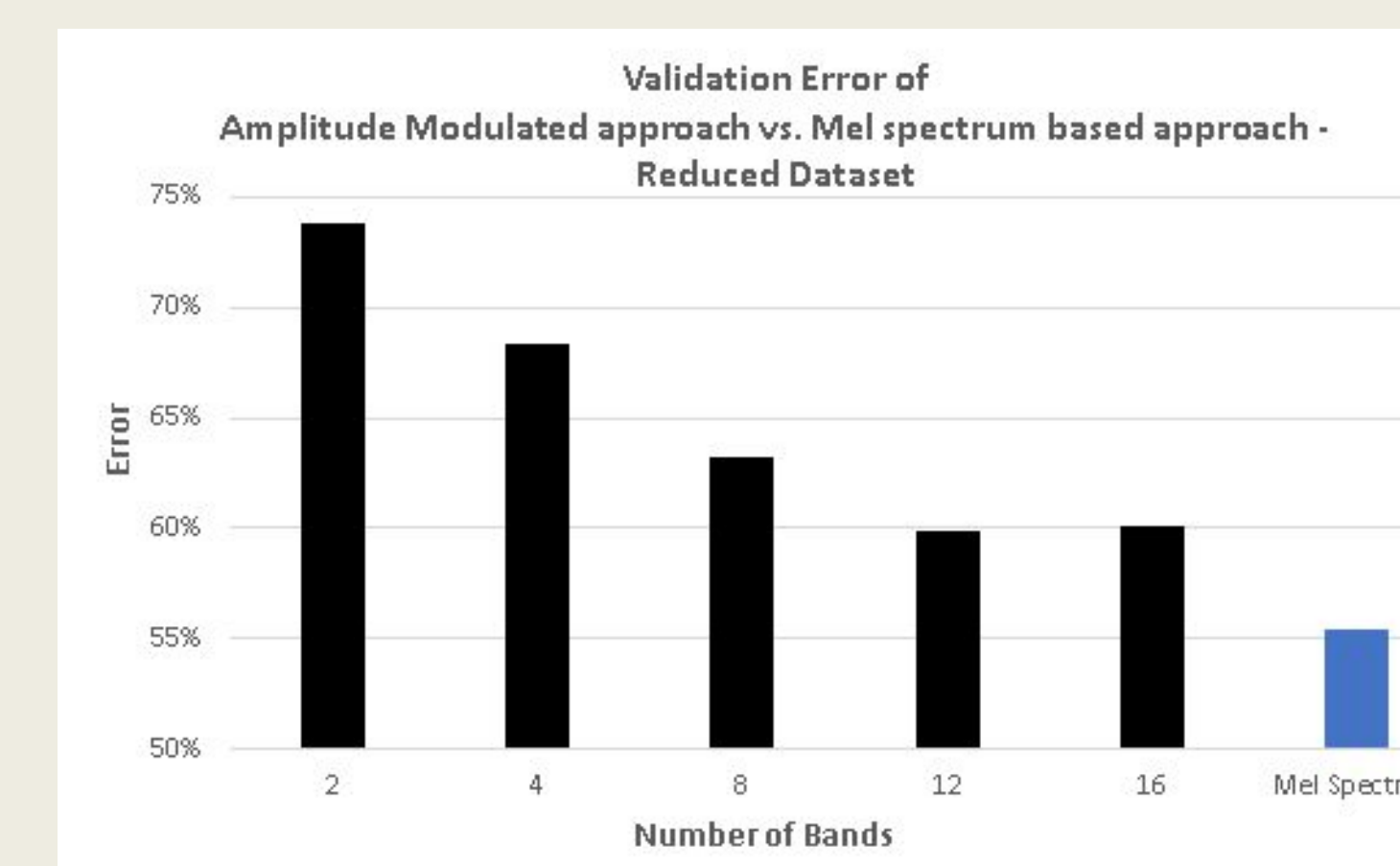


Fig. 2: Accuracy of AM signals with different numbers of bands compared to Mel spectrum accuracy using a reduced dataset

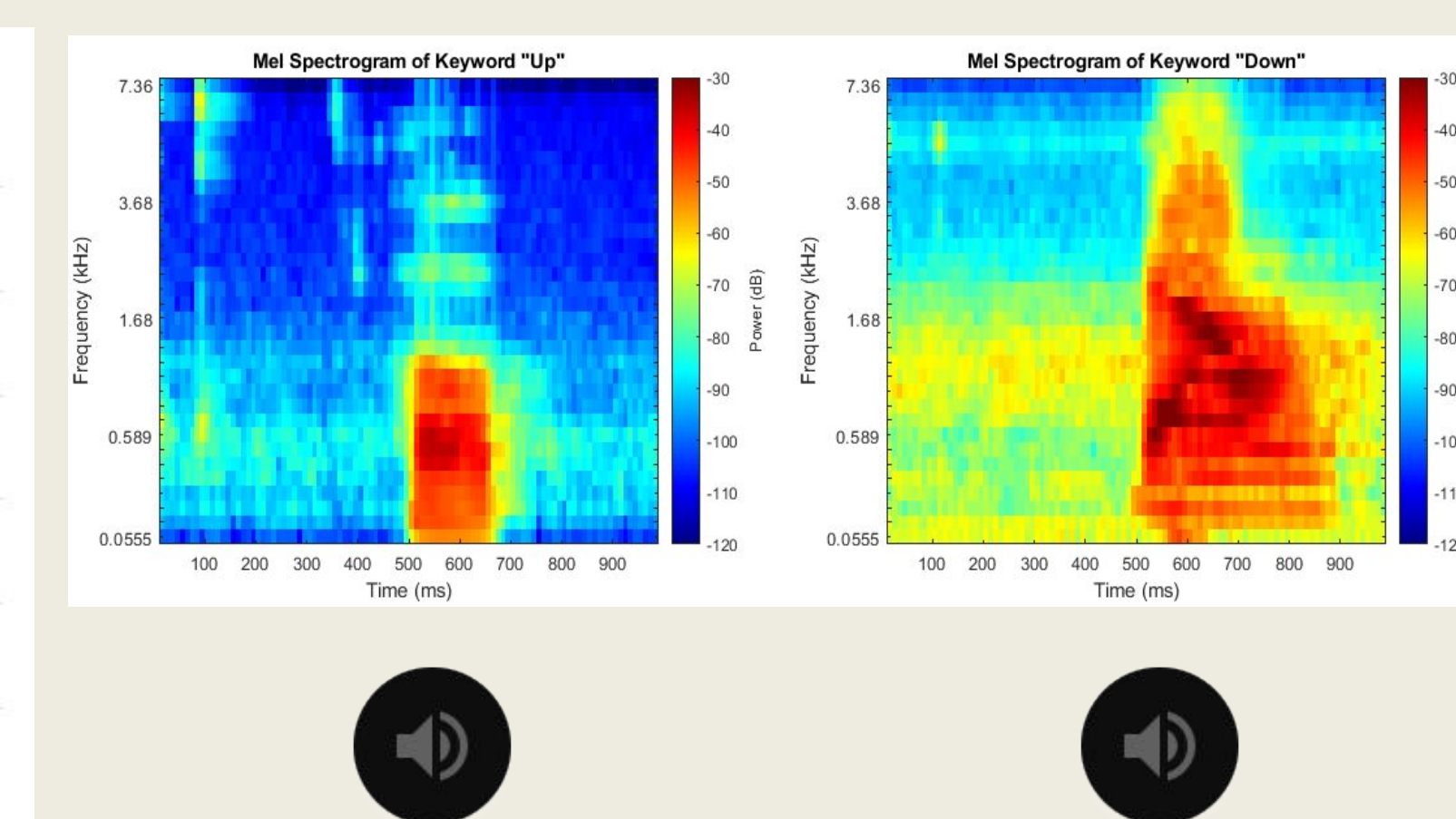


Fig. 3: The left is a Mel spectrogram and recording of the word “up” and the right is a Mel spectrogram and recording of the word “down”.

Confusion Matrix for Validation Data												
True Class	up	down	one	two	three	four	five	six	seven	eight	nine	unknown
up	379											74
down	1	367										81
one		1	339	1		6	1	1				52
two		1	2	350	1	4						59
three			2	4	353	4						59
four			5	1	1	395	4	4	2	1	1	43
five			3	1	3	2	1	361	1	1	1	63
six					3	2		1	409	1	9	20
seven			2	3		1	1	2	5	1	376	51
eight					7	6				1	2	40
nine			2	1	4				2		1	358
unknown	127	90	82	65	104	34	70	36	51	48	94	597
	72.7%	78.9%	77.9%	81.0%	79.1%	86.4%	80.8%	88.9%	85.5%	82.8%	73.5%	91.0%
	27.3%	21.1%	22.1%	19.0%	24.9%	11.6%	19.2%	11.1%	14.5%	17.2%	26.5%	9.0%
	up	down	one	two	three	four	five	six	seven	eight	nine	unknown

Fig. 4: Confusion matrix of a CNN trained with the full dataset of Mel spectrograms identifying keywords from the validation set Validation Error: 13.6105%

Confusion Matrix for Validation Data												
True Class	up	down	one	two	three	four	five	six	seven	eight	nine	unknown
up	387											61
down	5	364										75
one		2	336									55
two			7	360	4	1						44
three			1	7	366	1						45
four			3	1	7	1	389	4	2	5	2	43
five			4	3	1	5		373	1	1	11	50
six				3	1	4		2	408	1	2	21
seven			1	2	2	1	2	5	1	386	1	42
eight					8	8	1	3	4	1	346	2
nine					1	3	1	2	1	1	348	57
unknown	124	74	90	94	108	37	91	38	50	48	85	5949
	71.9%	82.4%	76.7%	74.8%	72.9%	86.8%	77.1%	86.5%	84.8%	80.7%	74.8%	91.8%
	28.1%	17.6%	23.3%	25.2%	27.1%	11.2%	22.9%	11.5%	15.2%	19.3%	25.2%	8.2%
	up	down	one	two	three	four	five	six	seven	eight	nine	unknown

Fig. 5: Confusion matrix of a CNN trained with the full dataset of Mel spectrograms of both original and amplitude modulated signals identifying keywords from the validation set Validation Error: 13.8085%

CONCLUSIONS

- It is feasible to perform keywords spotting using a convolutional neural network for volume adjustments of hearing devices.
- Mel-based spectral information is vital to train a neural network for keywords spotting, but combining this with amplitude information is a viable method
- Results are consistent with previous studies which show that with 8 or more bands, amplitude modulated audio can be recognized by human subjects.
- More research into amplitude modulation is necessary to determine if solely amplitude modulated signal processing is feasible for keywords spotting systems

REFERENCES

- [1] Y. Zhang, N. Suda, L. Lai, and V. Chandra, “Hello Edge: Keyword Spotting on Microcontrollers,” *ArXiv*, vol. abs/1711.07128, 2017.
- [2] J. H. Won, C. Lorenzi, K. Nie, X. Li, E. M. Jameyson, W. R. Drennan, and J. T. Rubinstein, “The ability of cochlear implant users to use temporal envelope cues recovered from speech frequency modulation,” *The Journal of the Acoustical Society of America*, vol. 132, no. 2, pp. 1113–1119, 2012.
- [3] E. Scheirer and M. Slaney, “Construction and evaluation of a robust multifeature speech/music discriminator,” 1997 IEEE International Conference on Acoustics, Speech, and Signal Processing, 1997, pp. 1331–1334.
- [4] P. Warden, “Launching the Speech Commands Dataset,” Google AI Blog, 24-Aug-2017. [Online]. Available: <https://ai.googleblog.com/2017/08/launching-speech-commands-dataset.html>.

ACKNOWLEDGEMENTS

This work was supported by the National Science Foundation REU Award #1757395 at the University of Washington Bothell.

INTRODUCTION

Problem: Changing settings in hearing aids and cochlear implants can be inconvenient by hand. Using keywords spotting (KWS) allows parameter adjustment to be done by voice.

Current feature extraction methods for KWS:

- o Bark spectrum, Mel spectrum, Mel frequency cepstral coefficients (MFCC)

Current deep learning techniques for KWS:

- o Recurrent Neural Network (RNN), Convolutional Neural Networks (CNN), Long-short term Memory (LSTM) [1]

Proposed Method: Applying amplitude modulation on an audio signal to extract information to train a CNN for keywords spotting. Amplitude modulation is vital for speech recognition because at certain frequencies useful speech information can be extracted [2] [3].

We compared the accuracy of a CNN using amplitude modulated signals with different bands versus Mel spectrum-based signal processing.

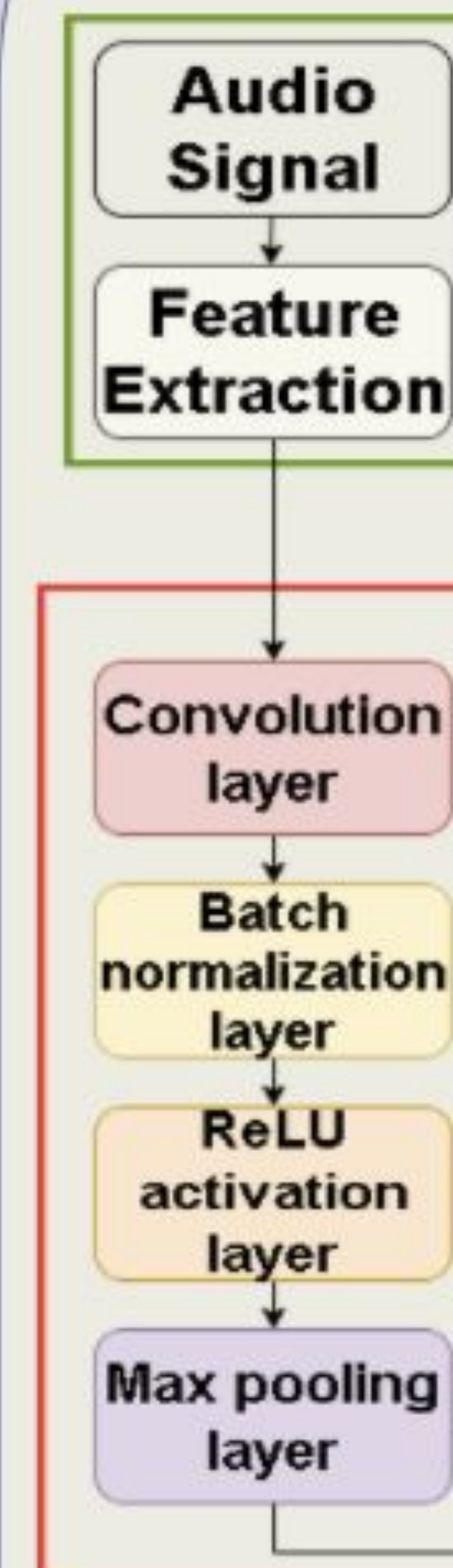


Figure
that po

We compared the accuracy of a CNN using amplitude modulated signals with different bands versus Mel spectrum-based signal processing.

METHODS

Google speech dataset

- 65000 different utterances of words were split into training, validation, and testing sets [4].
- 11 keywords, consisting of “up”, “down”, and “one” through “nine” were selected for volume control

Preprocessing

- An audio feature extractor in MATLAB was used to convert each audio signal into a Mel spectrogram.
- Amplitude modulation was obtained by passing the signal through multiple contiguous bandpass filters and performing a Hilbert Transform to extract temporal envelopes (Figure 1).

Neural Network

- A CNN with 3 hidden layers was setup in MATLAB (Figure 1).

RESULTS

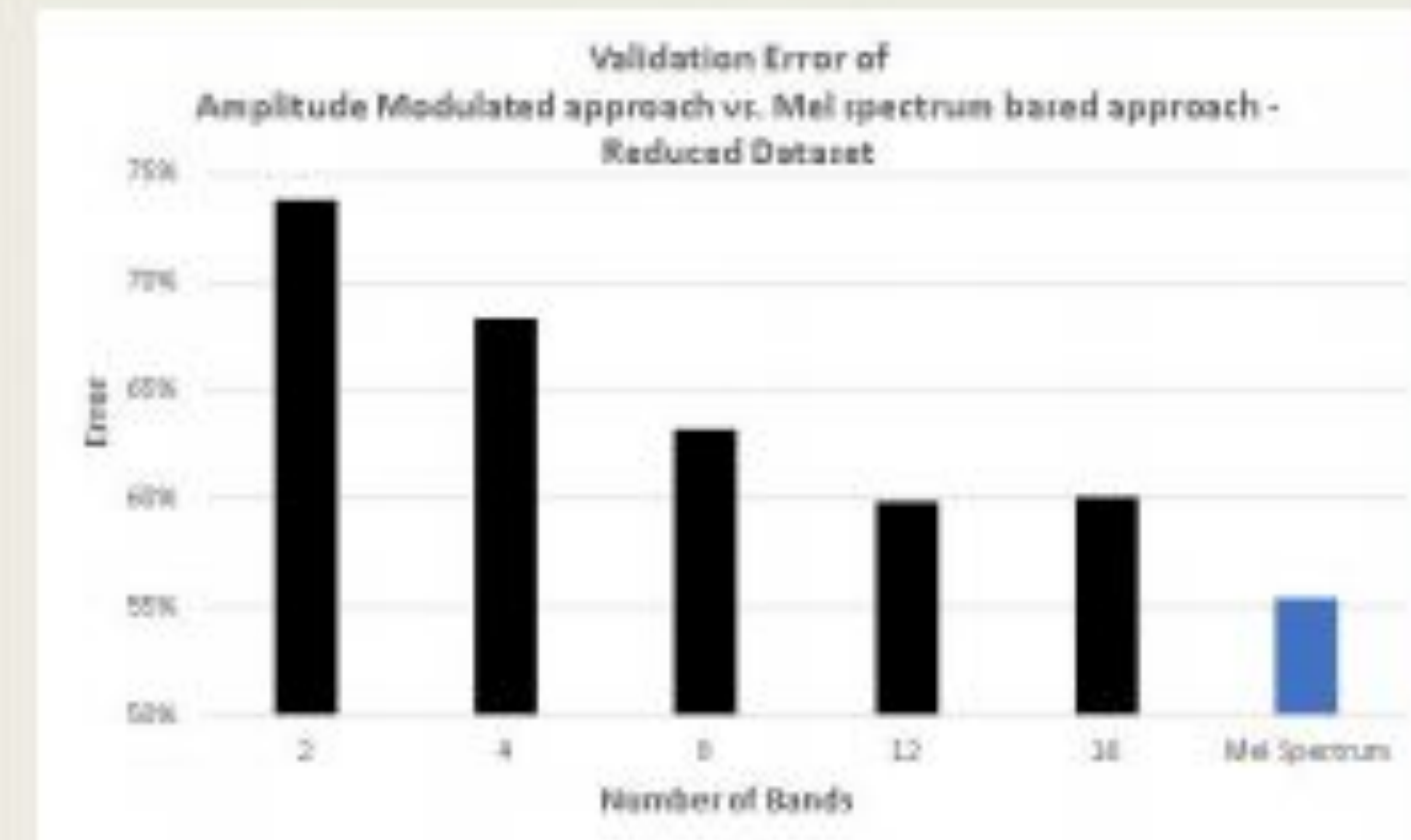


Fig. 2: Accuracy of AM signals with different numbers of bands compared to Mel spectrum accuracy using a reduced dataset

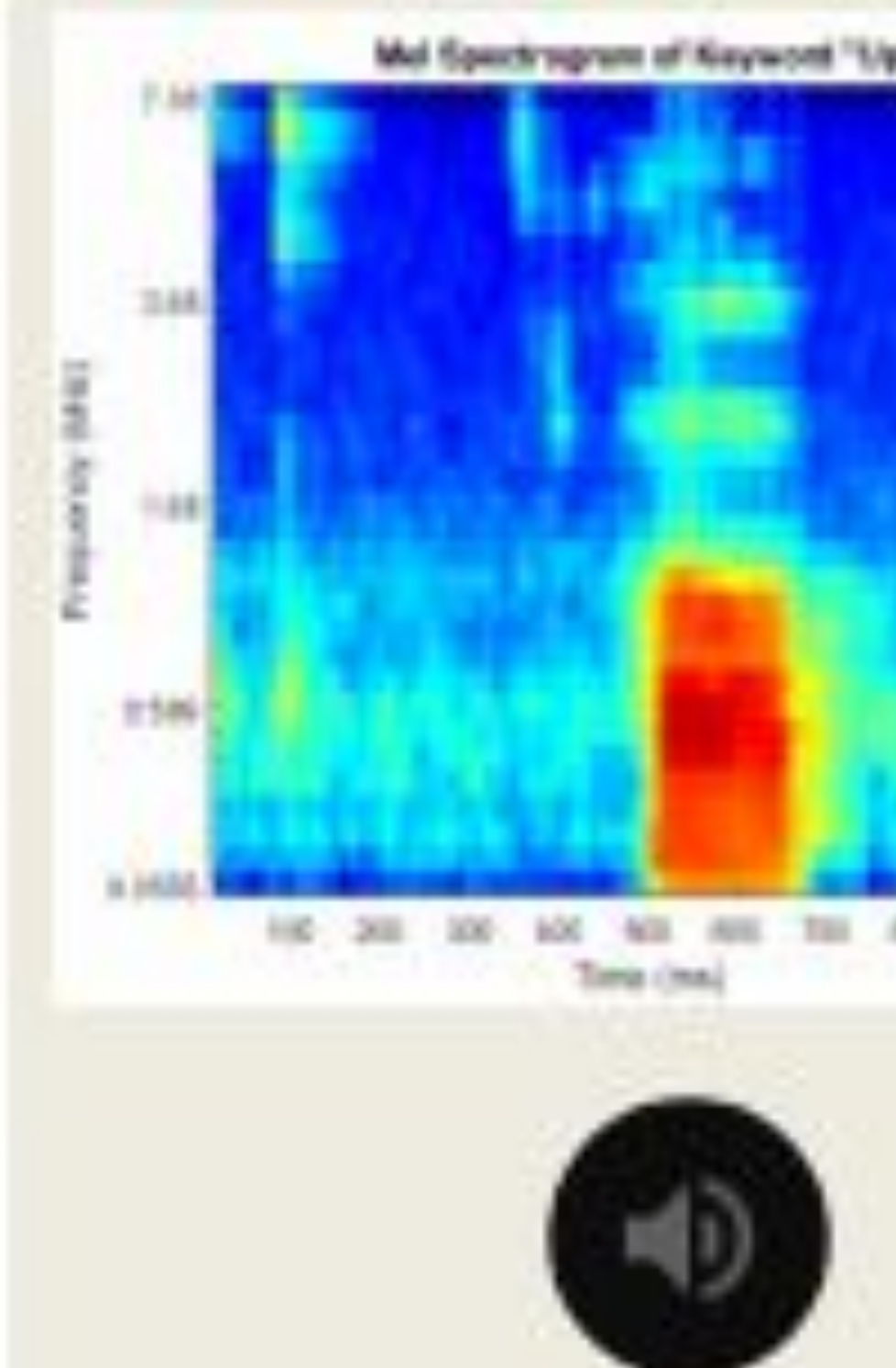


Fig. 3: The left is a Mel spectrogram of the keyword 'up' and the right is a Mel spectrogram of the keyword 'down'.



Fig. 4: Confusion matrix of a CNN trained with the full dataset of Mel spectrograms identifying keywords from the validation set. Validation Error: 13.6105%



Fig. 5: Confusion matrix of a CNN trained with the full dataset of Mel spectrograms identifying keywords from the validation set. Validation Error: 13.8085%

EXTRACTION

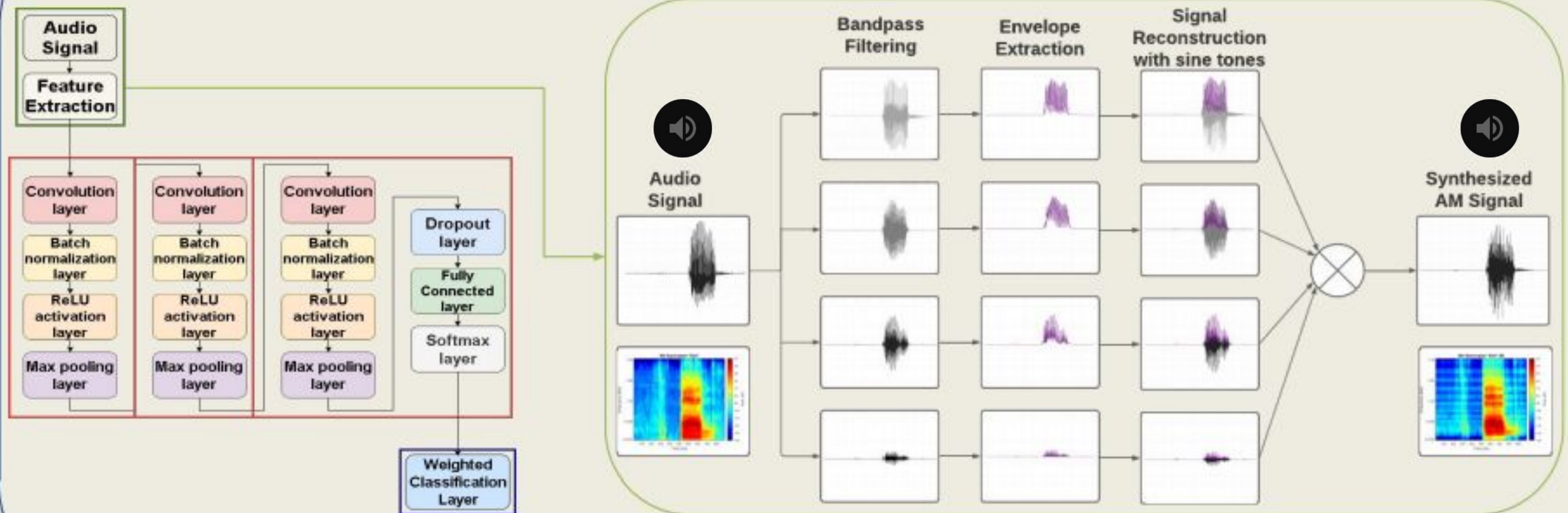
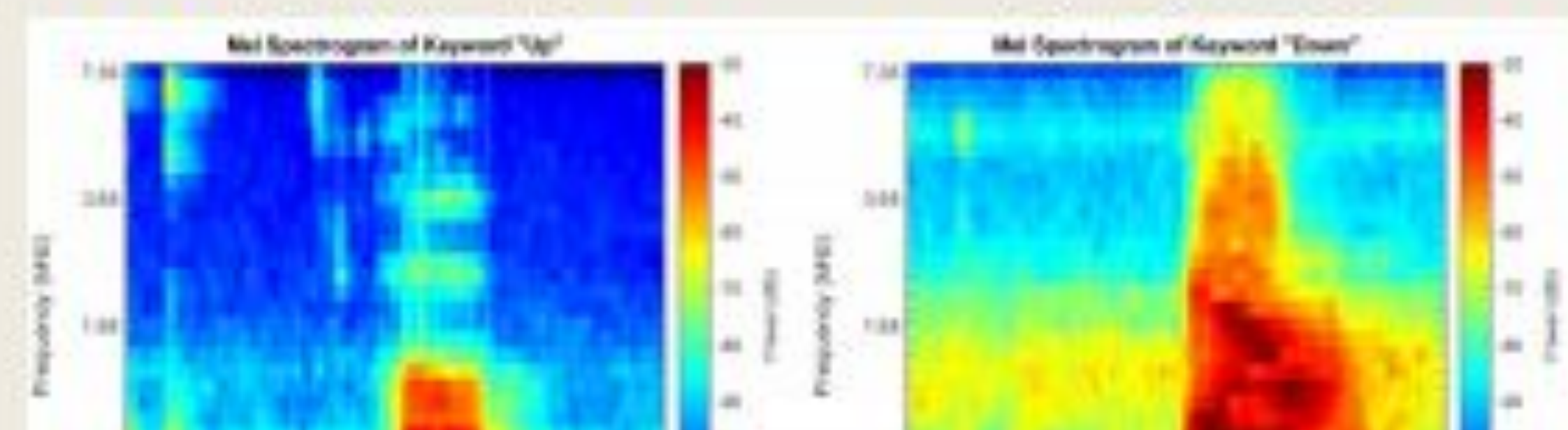


Figure 1: The left of the diagram outlines the entire convolutional neural network, with input, hidden layers, and output. The right displays the feature extraction algorithm that performs amplitude modulation on the signal before it is fed to the neural network.

RESULTS



CONCLUSIONS

- It is feasible to perform keywords spotting using a convolutional neural network for volume adjustments of hearing devices.
- Mel-based spectral information is vital to train a neural network for keywords spotting, but combining this with amplitude information is a viable method
- Results show that with 8 or more bands, amplitude modulated

RESULTS

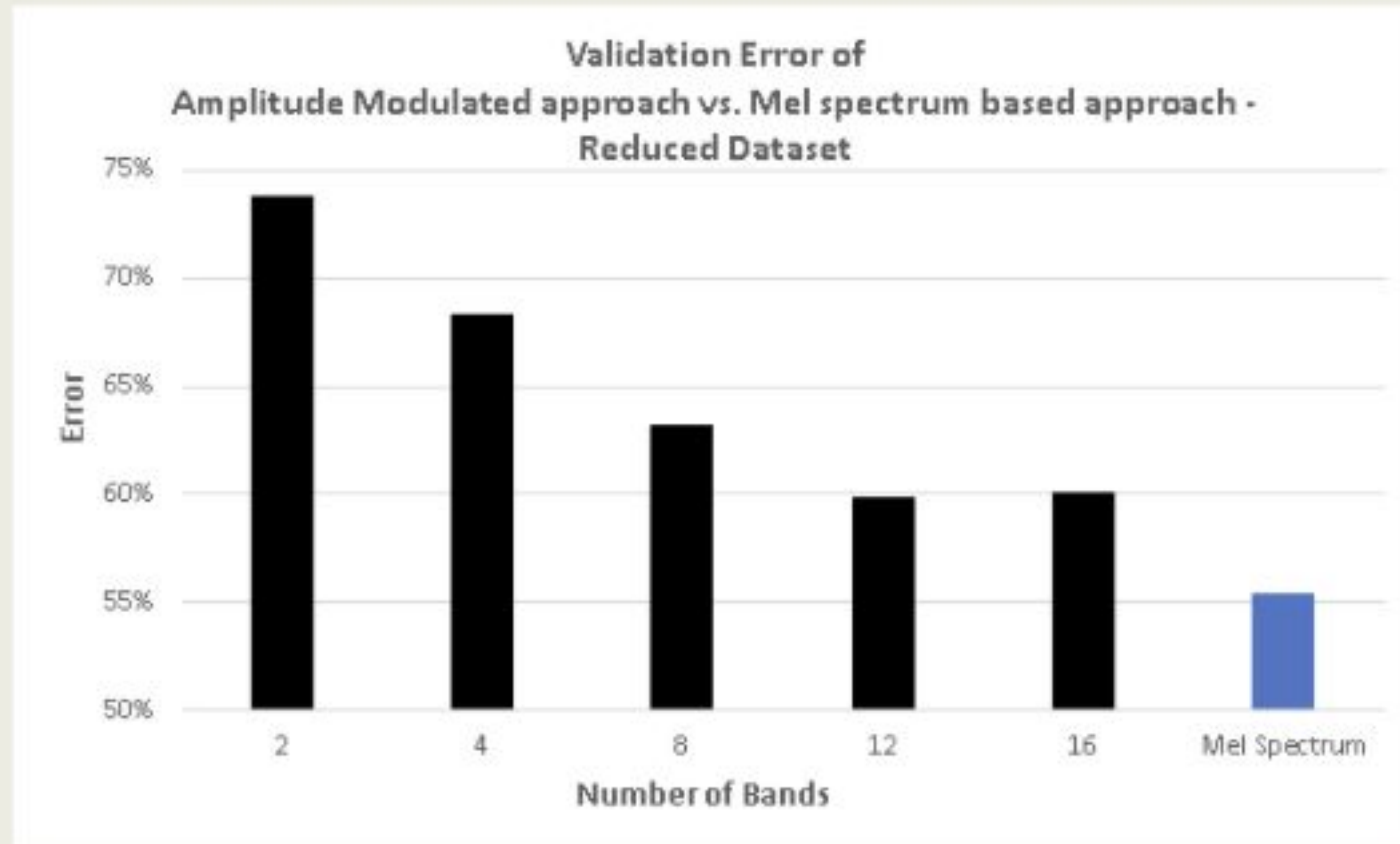


Fig. 2: Accuracy of AM signals with different numbers of bands compared to Mel spectrum accuracy using a reduced dataset

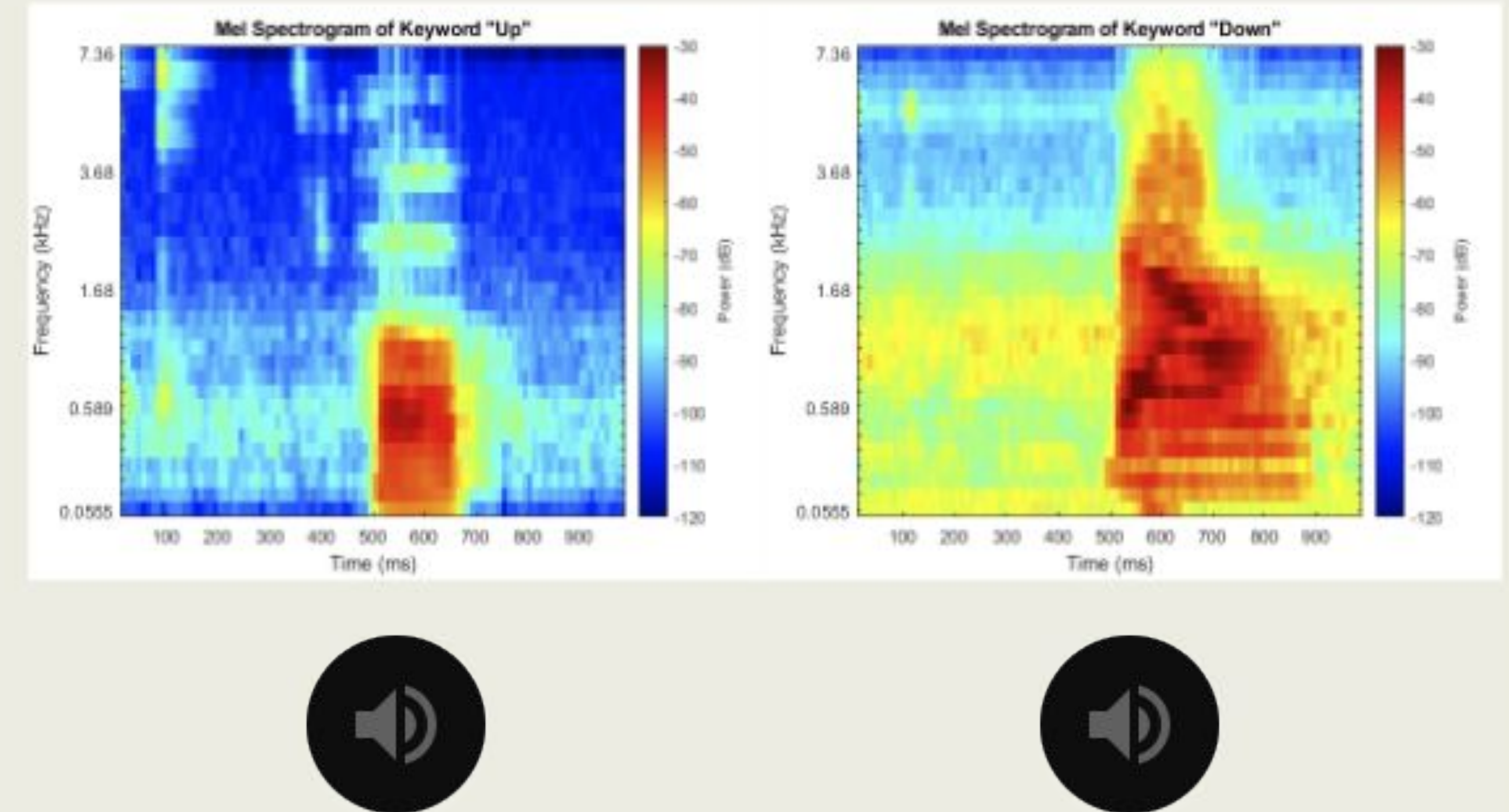


Fig. 3: The left is a Mel spectrogram and recording of the word “up” and the right is a Mel spectrogram and recording of the word “down”.

Confusion Matrix for Validation Data												
True Class	up	down	one	two	three	four	five	six	seven	eight	nine	unknown
up	379	1	1	1	1	1	1	1	1	1	1	74
down	1	367	1	1	1	1	1	1	1	1	1	81
one	1	1	339	1	1	1	1	1	1	1	1	52
two	1	2	1	350	1	1	1	1	1	1	1	59
three	1	1	1	1	353	1	1	1	1	1	1	59
four	2	1	1	1	1	395	1	1	1	1	1	43
five	3	1	1	1	1	1	361	1	1	1	1	63
six	2	1	1	1	1	1	1	409	1	1	1	20
seven	2	3	1	1	1	1	1	1	376	1	1	51
eight	2	1	1	1	1	1	1	1	1	361	1	40
nine	2	1	1	1	1	1	1	1	1	1	358	47
unknown	127	90	82	65	104	34	70	36	51	48	94	5987

82.6%	17.4%
80.5%	19.5%
82.7%	17.3%
82.7%	17.3%
79.5%	20.5%
86.1%	13.9%
80.4%	19.6%
91.1%	8.9%
84.5%	15.5%
86.2%	13.8%
86.3%	13.7%
88.2%	11.8%

72.7%	78.9%	77.9%	81.0%	75.1%	88.4%	80.8%	88.9%	85.5%	82.8%	73.5%	91.0%
27.3%	21.1%	22.1%	19.0%	24.9%	11.6%	19.2%	11.1%	14.5%	17.2%	26.5%	9.0%
up	down	one	two	three	four	five	six	seven	eight	nine	unknown

Fig. 4: Confusion matrix of a CNN trained with the full dataset of Mel spectrograms identifying keywords from the validation set
Validation Error: 13.6105%

Confusion Matrix for Validation Data												
True Class	up	down	one	two	three	four	five	six	seven	eight	nine	unknown
up	387	3	2	1	2	1	1	1	1	1	1	61
down	5	364	3	1	2	1	2	2	1	3	3	75
one	2	1	336	2	4	2	1	1	1	7	7	55
two	2	1	1	360	4	1	1	1	5	6	6	44
three	1	1	1	1	366	1	2	1	16	5	5	45
four	3	1	7	1	1	389	4	2	5	2	1	43
five	4	1	3	1	5	1	373	1	1	11	50	50
six	3	1	1	1	4	2	1	408	2	6	2	21
seven	1	2	2	2	1	2	5	1	386	1	42	42
eight	5	1	1	8	8	1	3	4	1	346	2	41
nine	1	1	1	1	3	1	2	1	1	1	348	57
unknown	124	74	90	94	108	37	91	38	50	48	85	5949

84.3%	15.7%
79.8%	20.2%
82.0%	18.0%
85.1%	14.9%
82.4%	17.6%
84.7%	15.3%
83.1%	16.9%
90.9%	9.1%
86.7%	13.3%
82.6%	17.4%
83.9%	16.1%
87.6%	12.4%

71.9%	82.4%	76.7%	74.8%	72.9%	88.8%	77.1%	88.5%	84.8%	80.7%	74.8%	91.8%
28.1%	17.6%	23.3%	25.2%	27.1%	11.2%	22.9%	11.5%	15.2%	19.3%	25.2%	8.2%
up	down	one	two	three	four	five	six	seven	eight	nine	unknown

Fig. 5: Confusion matrix of a CNN trained with the full dataset of Mel spectrograms of both original and amplitude modulated signals identifying keywords from the validation set
Validation Error: 13.8085%

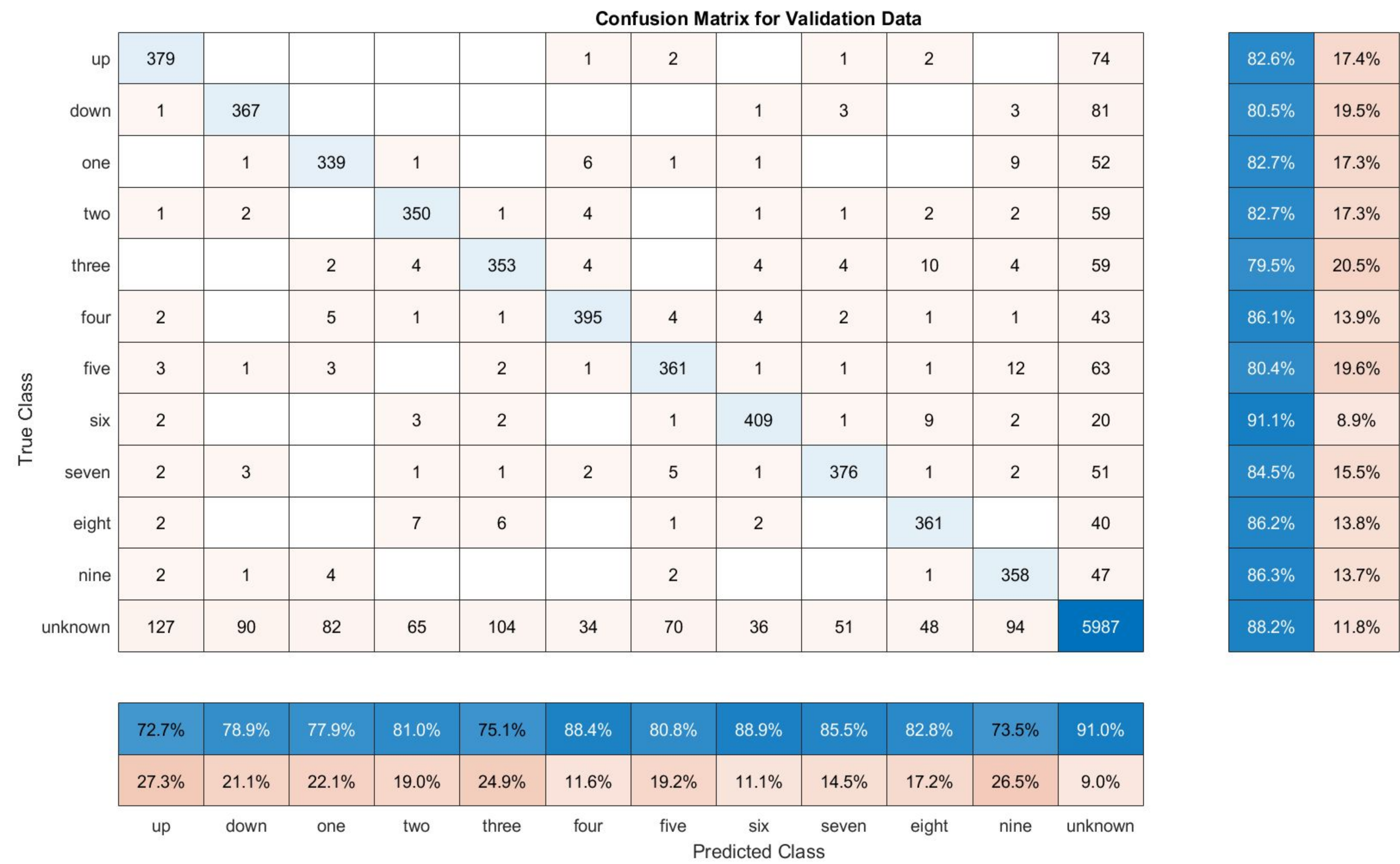


Fig. 4: Confusion matrix of a CNN trained with the full dataset of Mel spectrograms identifying keywords from the validation set
Validation Error: 13.6105%

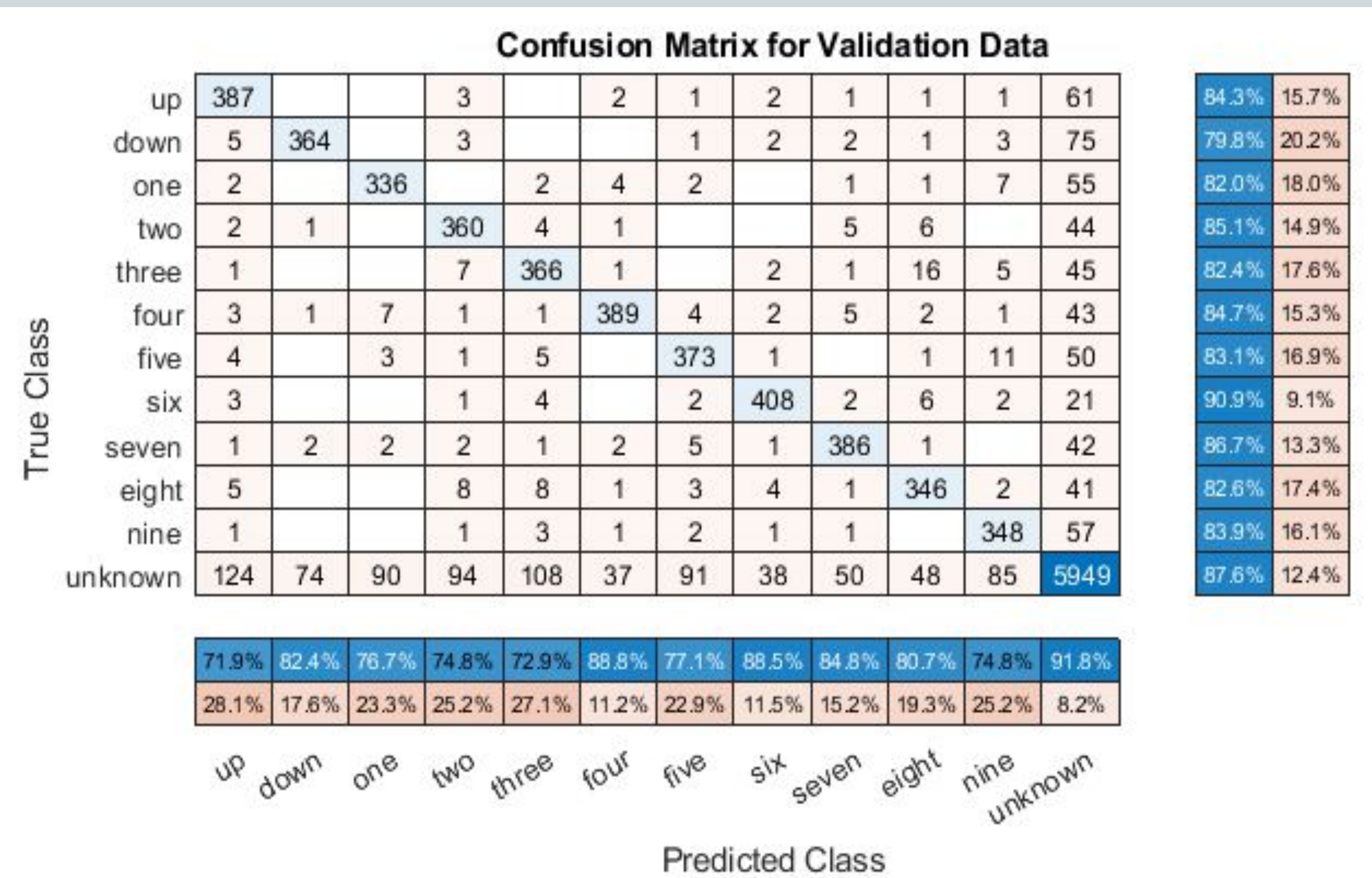


Fig. 5: Confusion matrix of a CNN trained with the full dataset of Mel spectrograms of both original and amplitude modulated signals identifying keywords from the validation set
Validation Error: 13.8085%

RESULTS

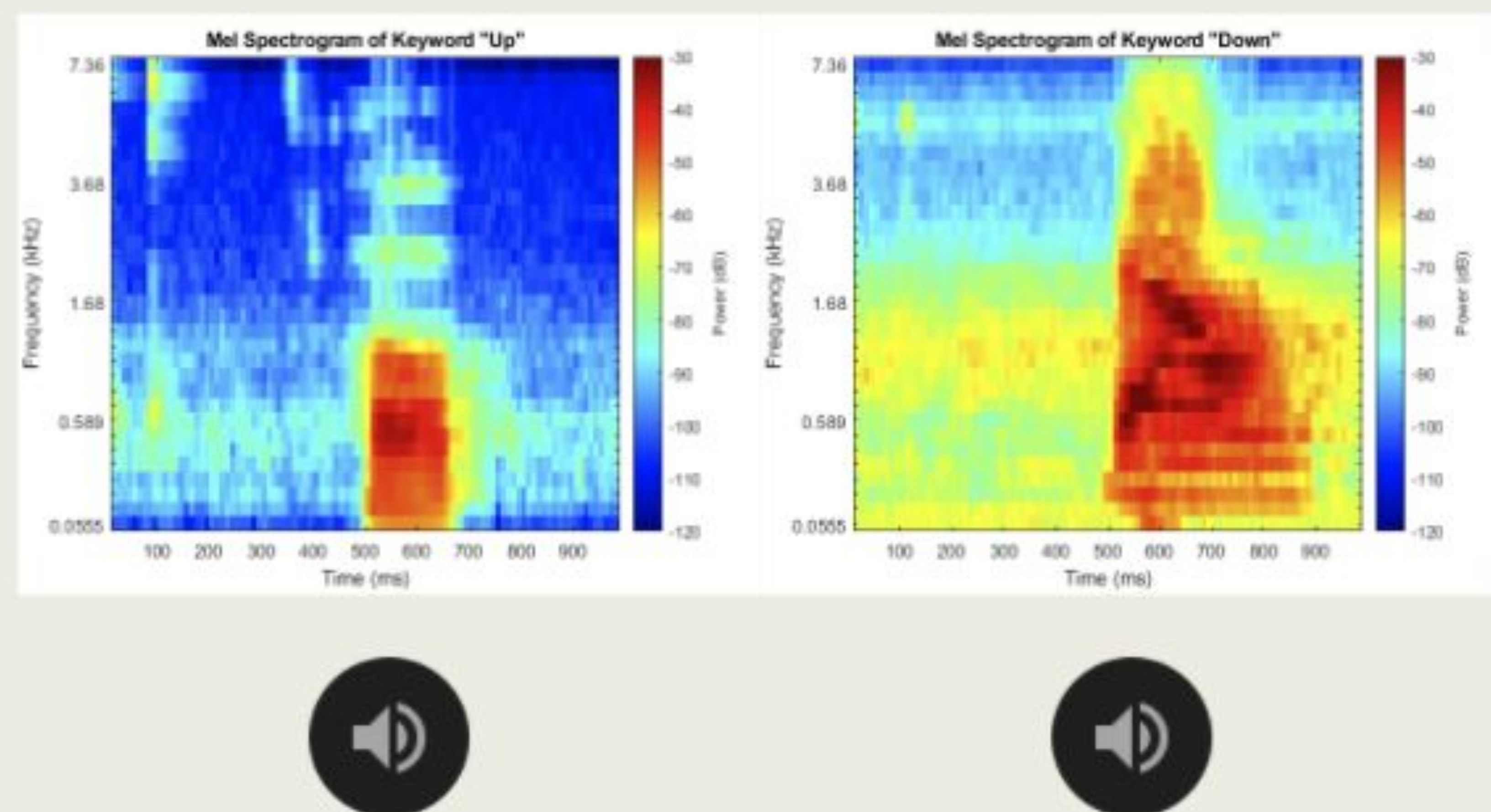


Fig. 3: The left is a Mel spectrogram and recording of the word “up” and the right is a Mel spectrogram and recording of the word “down”.

Confusion Matrix for Validation Data												
True Class	up	down	one	two	three	four	five	six	seven	eight	nine	unknown
up	387			3		2	1	2	1	1	1	61
down	5	364		3			1	2	2	1	3	75
one	2		336		2	4	2		1	1	7	55
two	2	1		360	4	1			5	6		44
three	1			7	366	1		2	1	16	5	45
four	3	1	7	1	1	389	4	2	5	2	1	43
five	4		3	1	5		373	1		1	11	50
six	3			1	4		2	408	2	6	2	21
seven	1	2	2	2	1	2	5	1	386	1		42
eight	5			8	8	1	3	4	1	346	2	41
nine	1			1	3	1	2	1	1		348	57
unknown	124	74	90	94	108	37	91	38	50	48	85	5949
	71.9%	82.4%	76.7%	74.8%	72.9%	86.8%	77.1%	88.5%	84.8%	80.7%	74.8%	91.8%
	28.1%	17.6%	23.3%	25.2%	27.1%	11.2%	22.9%	11.5%	15.2%	19.3%	25.2%	8.2%
	up	down	one	two	three	four	five	six	seven	eight	nine	unknown
Predicted Class												

Fig. 5: Confusion matrix of a CNN trained with the full dataset of Mel spectrograms of both original and amplitude modulated signals identifying keywords from the validation set

Validation Error: 13.8085%

CONCLUSIONS

- It is feasible to perform keywords spotting using a convolutional neural network for volume adjustments of hearing devices.
- Mel-based spectral information is vital to train a neural network for keywords spotting, but combining this with amplitude information is a viable method
- Results are consistent with previous studies which show that with 8 or more bands, amplitude modulated audio can be recognized by human subjects.
- More research into amplitude modulation is necessary to determine if solely amplitude modulated signal processing is feasible for keywords spotting systems

REFERENCES

- [1] Y. Zhang, N. Suda, L. Lai, and V. Chandra, “Hello Edge: Keyword Spotting on Microcontrollers,” *ArXiv*, vol. abs/1711.07128, 2017.
- [2] J. H. Won, C. Lorenzi, K. Nie, X. Li, E. M. Jarneyson, W. R. Drennan, and J. T. Rubinstein, “The ability of cochlear implant users to use temporal envelope cues recovered from speech frequency modulation,” *The Journal of the Acoustical Society of America*, vol. 132, no. 2, pp. 1113–1119, 2012.
- [3] E. Scheirer and M. Slaney, “Construction and evaluation of a robust multifeature speech/music discriminator,” 1997 IEEE International Conference on Acoustics, Speech, and Signal Processing, 1997, pp. 1331–1334.
- [4] P. Warden, “Launching the Speech Commands Dataset,” Google AI Blog, 24-Aug-2017. [Online]. Available: <https://ai.googleblog.com/2017/08/launching-speech-commands-dataset.html>.

ACKNOWLEDGEMENTS

This work was supported by the National Science Foundation REU Award #1757395 at the University of Washington Bothell.