

# Deep Neural Network Based Polyphonic Sound Event Detection

Emre Cakir, Giambattista Parascandolo

Audio Research Group  
Tampere University of Technology, Finland

February 24, 2017

SGN-41007 Guest lecture



TAMPERE UNIVERSITY OF TECHNOLOGY

## 1 Sound Event Detection

## 2 Problem Formulation

- Sound Representation
- Sound Event Classification
- Experiments setup

## 3 Results

## 4 Conclusions



# Sound Events



TAMPERE UNIVERSITY OF TECHNOLOGY

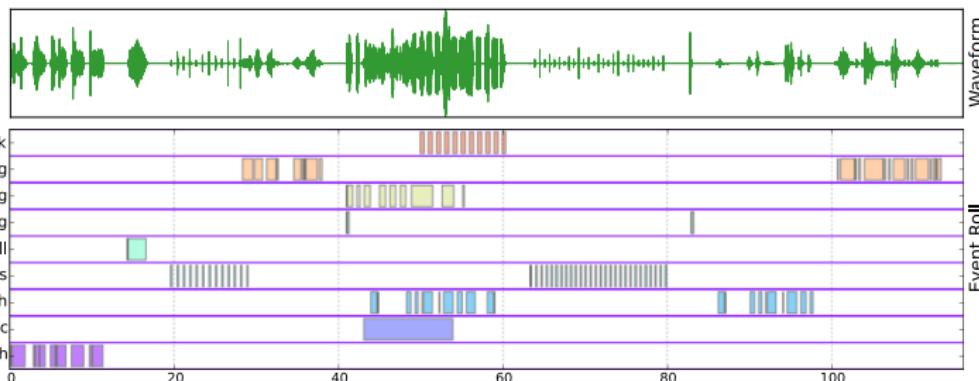
# SED Application Areas

- Audio surveillance (gun shots, scream, footsteps)
- Health-care monitoring (coughing, sneezing frequency)
- Context-aware devices (a mobile phone detecting the user is in a meeting and adjusting settings accordingly)
- Smart-home systems
- etc.



# Polyphonic Sound Event Detection

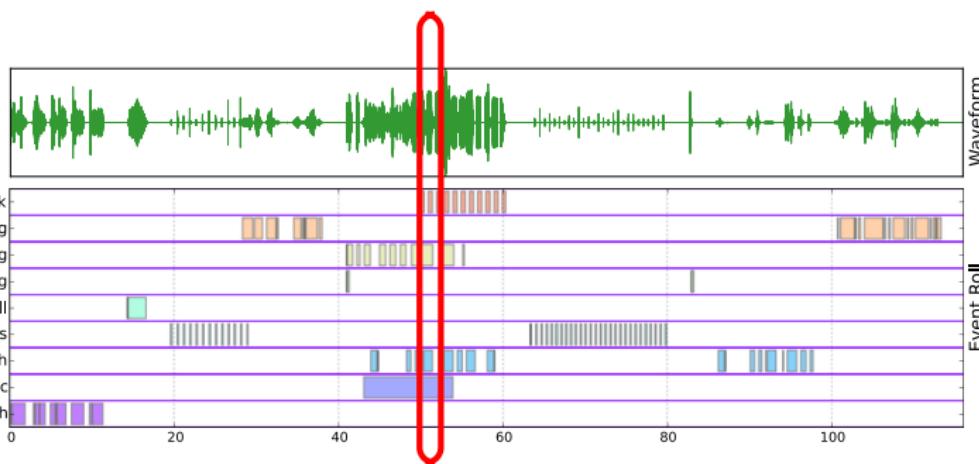
- ① Estimate start and end time of each sound event
- ② Assign a textual label for each event



TAMPERE UNIVERSITY OF TECHNOLOGY

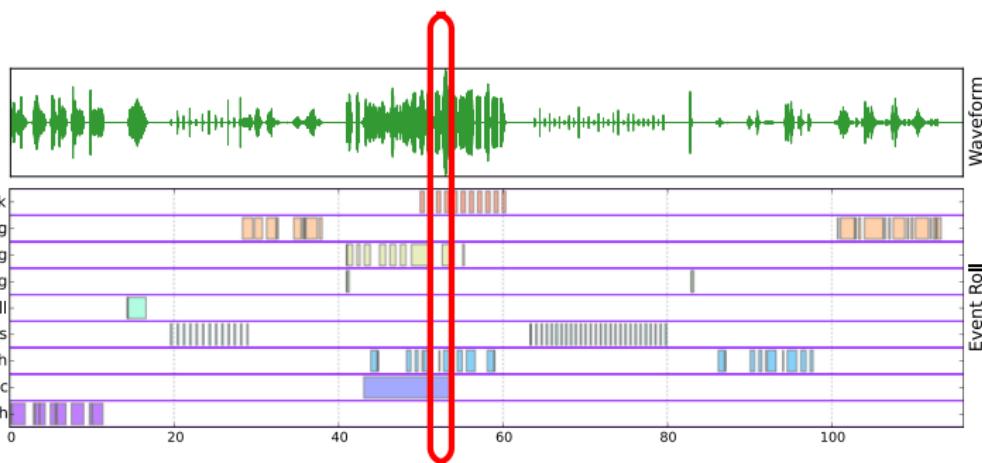
# Polyphonic Sound Event Detection

- ① Estimate start and end time of each sound event
- ② Assign a textual label for each event



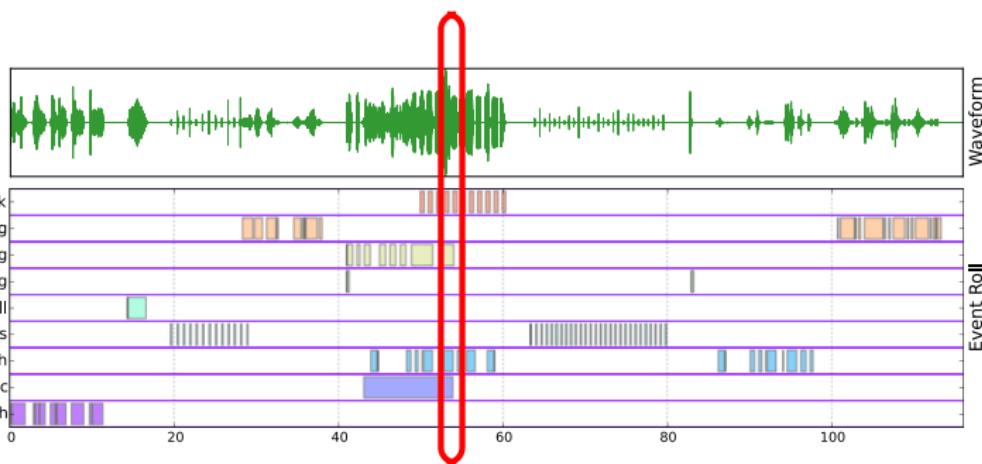
# Polyphonic Sound Event Detection

- ① Estimate start and end time of each sound event
- ② Assign a textual label for each event



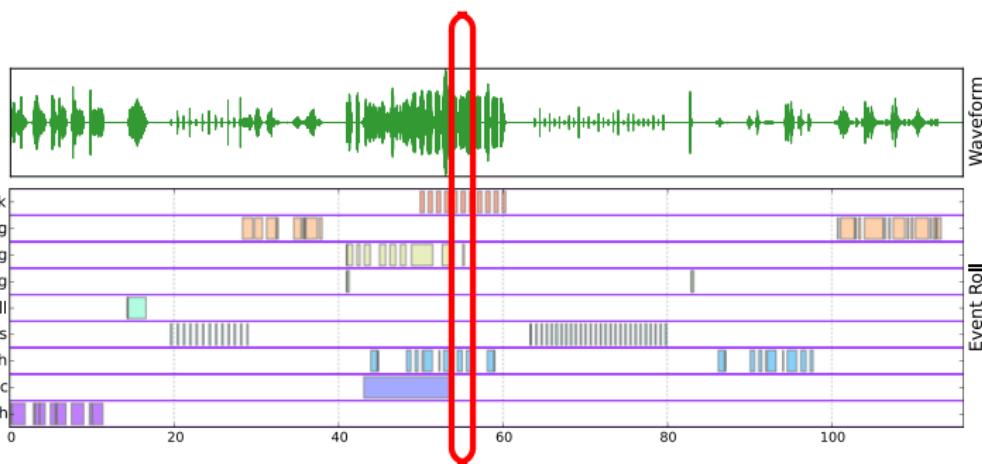
# Polyphonic Sound Event Detection

- ① Estimate start and end time of each sound event
- ② Assign a textual label for each event



# Polyphonic Sound Event Detection

- ① Estimate start and end time of each sound event
- ② Assign a textual label for each event



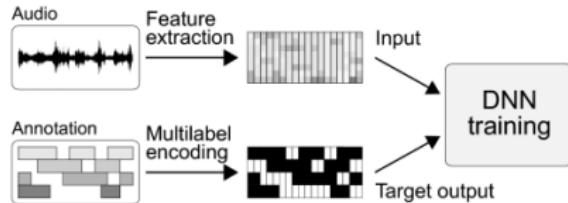
# Challenges for SED in real-life

- Environmental noise
- Overlapping sound events
- Variation among sources for the same event

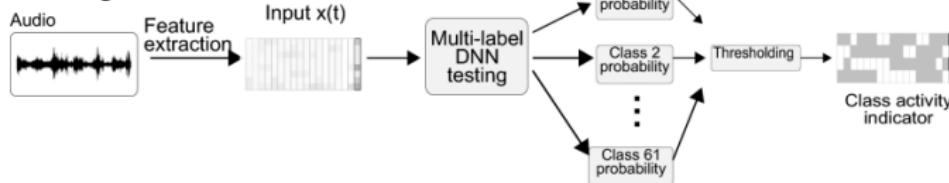


# System overview

## Training



## Testing



## Polyphonic Sound Event Detection

- ① Sound Representation (feature extraction)
- ② Sound Event Classification (DNN learning)



# Polyphonic Sound Event Detection

- ① **Sound Representation**
- ② Sound Event Classification



TAMPERE UNIVERSITY OF TECHNOLOGY

# Sound Representation

Sound representation most often in frequency domain using short audio frames

- More informative than raw audio
- 2-D input images, similar to image recognition / object detection
- Perception based features



# Sound representation steps

- Frame blocking
- Windowing
- STFT for each short frame
- Perception based post-processing (Mel-band, Bark-band etc.)
- Normalization over each feature

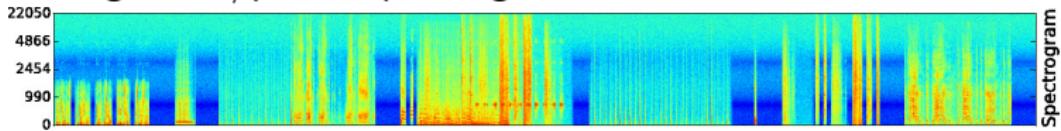


# Sound Representation

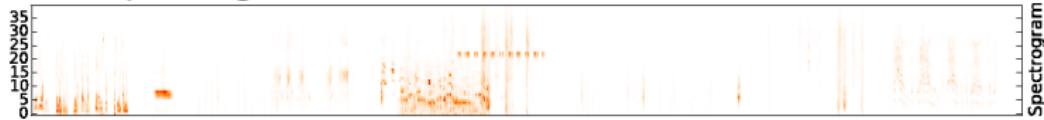
- Audio waveform



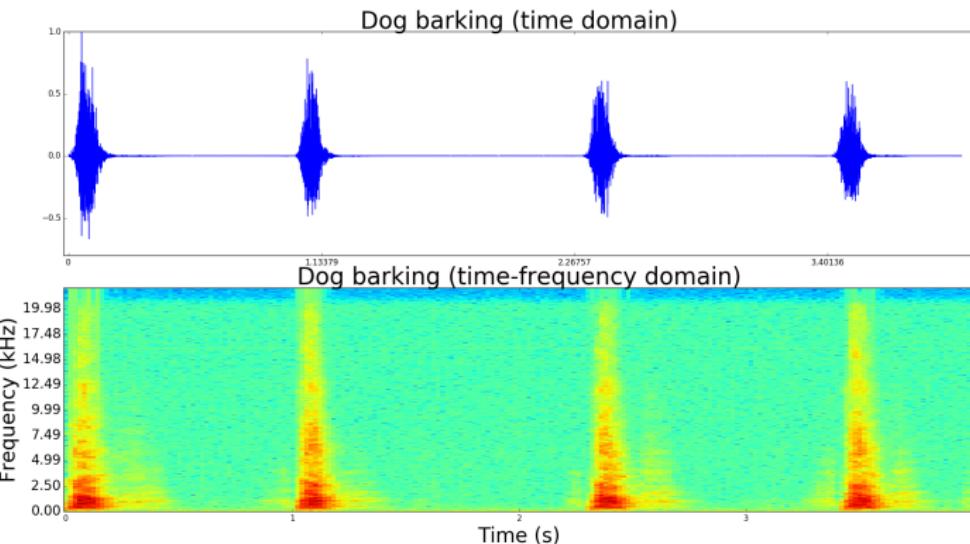
- Magnitude/power spectrogram



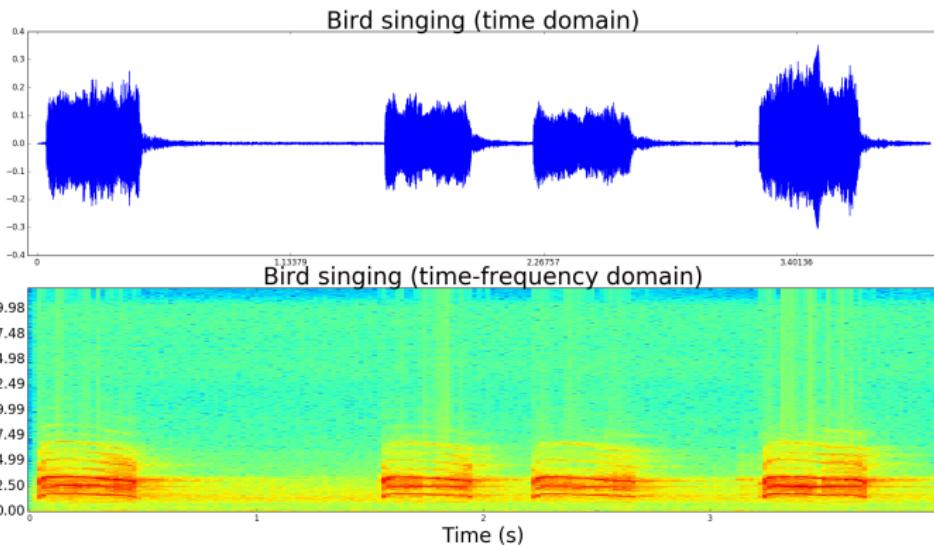
- Mel spectrogram



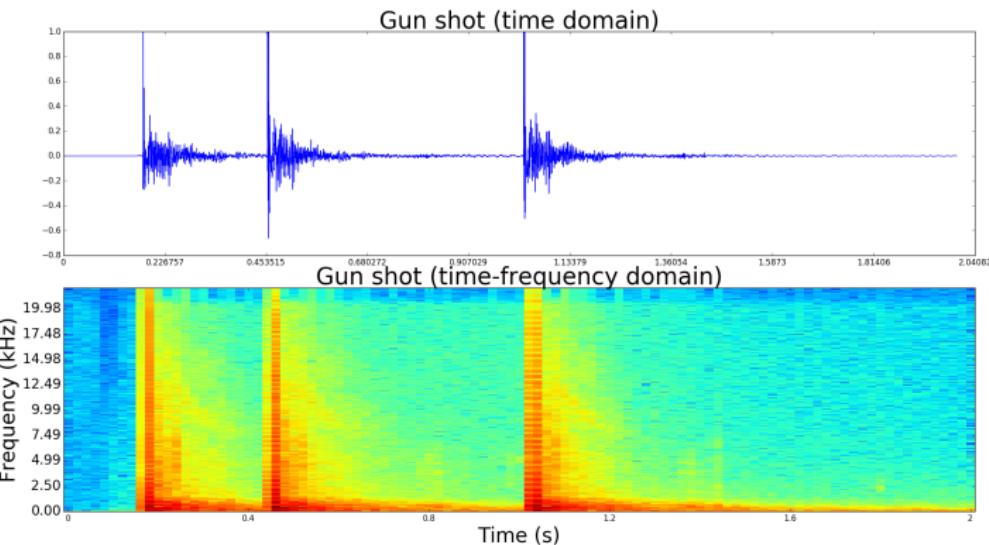
# Sound Representation



# Sound Representation

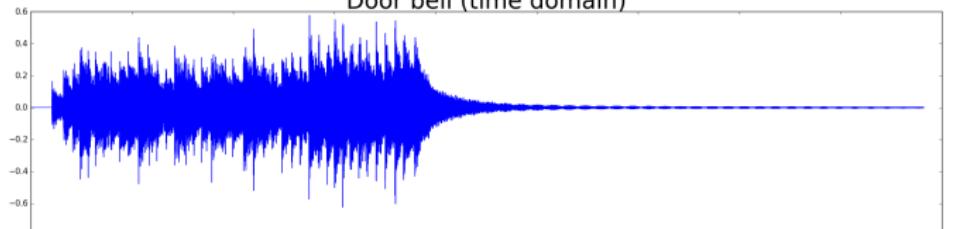


# Sound Representation

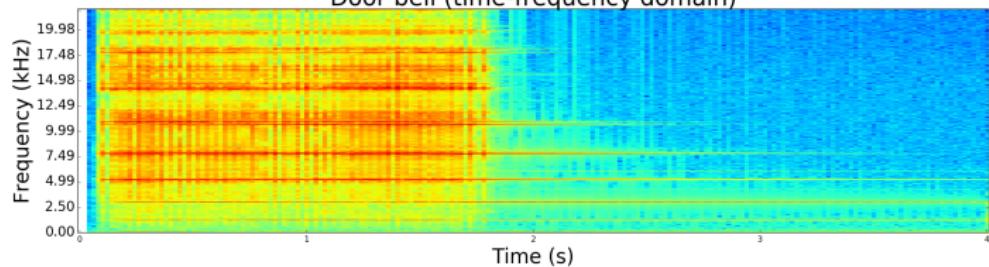


# Sound Representation

Door bell (time domain)



Door bell (time-frequency domain)



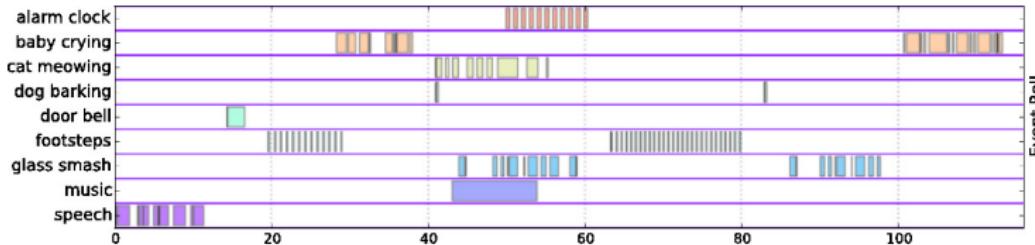
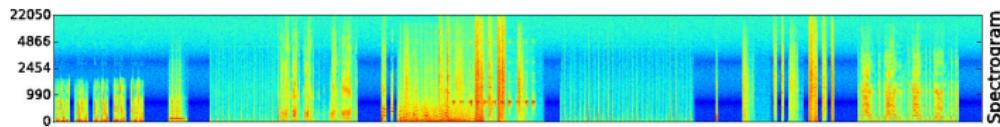
TAMPERE UNIVERSITY OF TECHNOLOGY

## Polyphonic Sound Event Detection

- ① Sound Representation
- ② **Sound Event Classification**



# Sound Event Classification

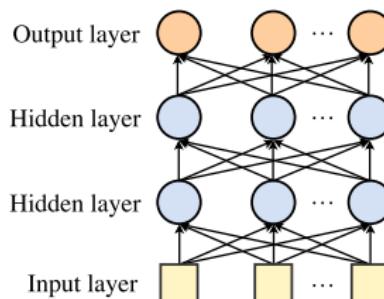


# Previous work

- Gaussian mixture model (GMM)
- Hidden Markov model (HMM)
- Feedforward neural networks (FNN)
- Convolutional neural networks (CNN)
- Recurrent neural networks (RNN)



# Previous work (FNN)



- Simple, high expressional power and accuracy compared to established methods (+)
- Context window input to FNN beneficial (+)



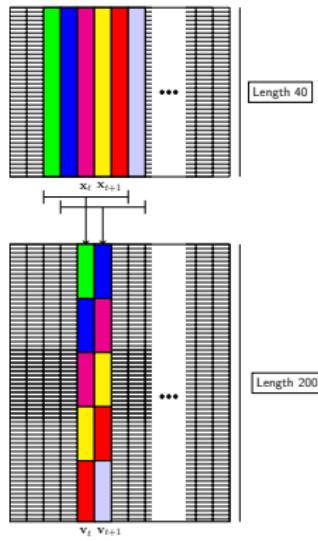
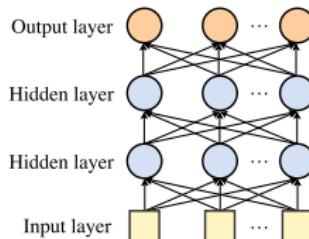


Figure: Concatenation of  $x_t$  with its 2 adjacent frames on both sides to form  $v_t$ .



# Previous work (FNN)

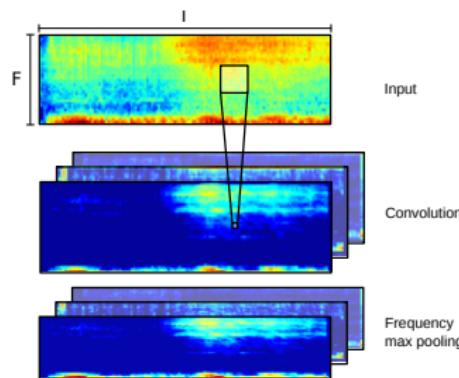


- Simple, high expressional power and accuracy compared to established methods (+)
- Context window input to FNN beneficial (+)
- Not robust for small variations in frequency content due to fixed connections (-)
- Short temporal context (-)



TAMPERE UNIVERSITY OF TECHNOLOGY

# Previous work (CNN)

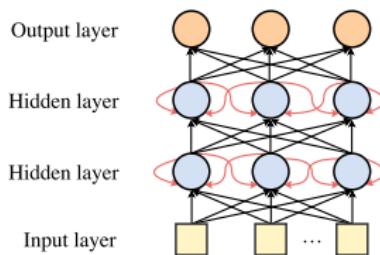


- Input frequency shift invariance due to shared weights and max pooling (+)
- Short context window (-)



TAMPERE UNIVERSITY OF TECHNOLOGY

# Previous work (RNN)



Feedback structure: neuron activations from previous inputs are used while calculating the activation for the current input

- Able to model long-term dependencies between consecutive inputs (+)
- Not robust for variations in frequency content (-)



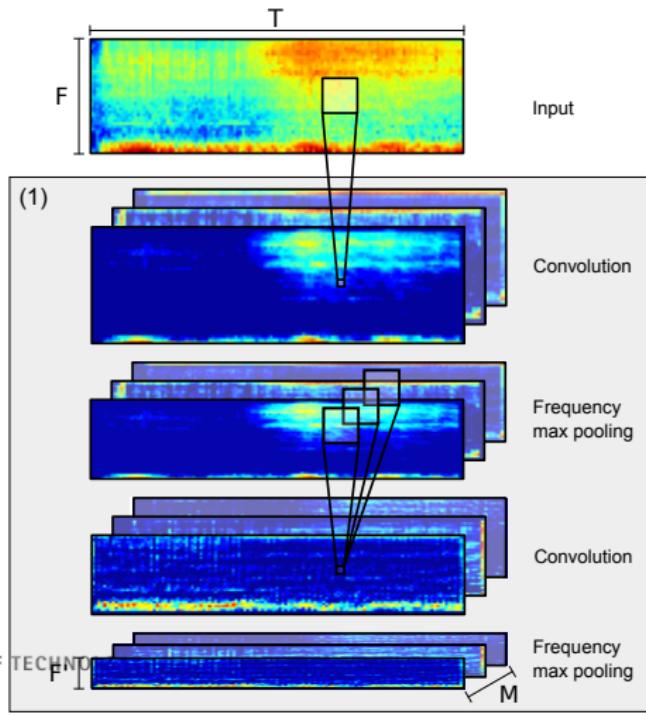
# Proposed method

- Combine CNN and RNN in a single deep learning architecture called Convolutional recurrent neural network (CRNN)
- Apply CRNN on polyphonic, scene-independent SED

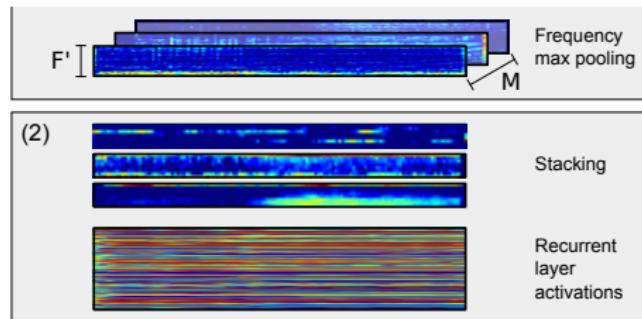
ArXiv: <https://arxiv.org/abs/1702.06286>



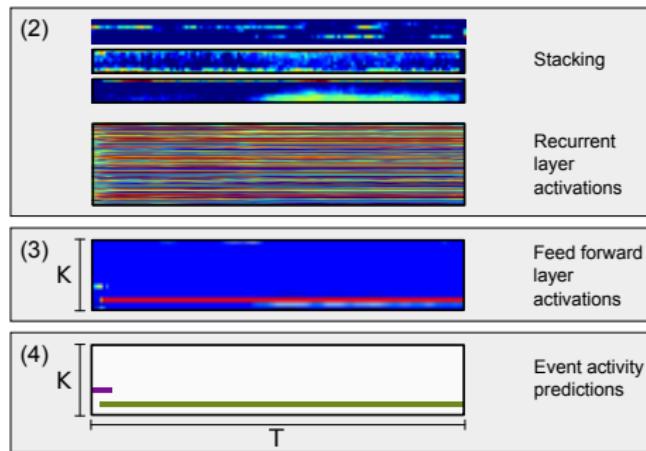
# Proposed method - Convolutional Recurrent Neural Network



# Proposed method - Convolutional Recurrent Neural Network



# Proposed method - Convolutional Recurrent Neural Network



CNN works as a feature extractor, RNN does temporal integration and language modeling.

Google and Baidu successfully used similar architectures for speech recognition.

CNNs and RNNs as subsets of CRNNs

- A CNN is a CRNN with 0 recurrent layers
- An RNN is a CRNN with 0 convolutional layers



4 datasets.

Preprocessing: 40 log mel bands, 40ms frames, 50% overlap

Grid-search for network hyper-parameters

- Hidden units for conv and recur layers {96, 256}



4 datasets.

Preprocessing: 40 log mel bands, 40ms frames, 50% overlap

Grid-search for network hyper-parameters

- Hidden units for conv and recur layers {96, 256}
- Hidden layers 1 to 3 for recur, 1 to 4 for conv



4 datasets.

Preprocessing: 40 log mel bands, 40ms frames, 50% overlap

Grid-search for network hyper-parameters

- Hidden units for conv and recur layers {96, 256}
- Hidden layers 1 to 3 for recur, 1 to 4 for conv
- Several frequency maxpooling arrangements



4 datasets.

Preprocessing: 40 log mel bands, 40ms frames, 50% overlap

Grid-search for network hyper-parameters

- Hidden units for conv and recur layers {96, 256}
- Hidden layers 1 to 3 for recur, 1 to 4 for conv
- Several frequency maxpooling arrangements
- Select the best architecture on validation to compute test results



- ReLU as hidden activations, sigmoid for the output layer. Cut-off for binary thresholding at 0.5.



- ReLU as hidden activations, sigmoid for the output layer. Cut-off for binary thresholding at 0.5.
- $5 \times 5$  kernels in conv layers



- ReLU as hidden activations, sigmoid for the output layer. Cut-off for binary thresholding at 0.5.
- $5 \times 5$  kernels in conv layers
- Gated Recurrent Units (GRU) in all recurrent layers. Simpler alternative to LSTM.



- ReLU as hidden activations, sigmoid for the output layer. Cut-off for binary thresholding at 0.5.
- $5 \times 5$  kernels in conv layers
- Gated Recurrent Units (GRU) in all recurrent layers. Simpler alternative to LSTM.
- Dropout in all layers 0.25



- ReLU as hidden activations, sigmoid for the output layer. Cut-off for binary thresholding at 0.5.
- $5 \times 5$  kernels in conv layers
- Gated Recurrent Units (GRU) in all recurrent layers. Simpler alternative to LSTM.
- Dropout in all layers 0.25
- Batch normalization in all non-recurrent layers



- ReLU as hidden activations, sigmoid for the output layer. Cut-off for binary thresholding at 0.5.
- $5 \times 5$  kernels in conv layers
- Gated Recurrent Units (GRU) in all recurrent layers. Simpler alternative to LSTM.
- Dropout in all layers 0.25
- Batch normalization in all non-recurrent layers
- Adam as optimizer, cross-entropy as loss



- ReLU as hidden activations, sigmoid for the output layer. Cut-off for binary thresholding at 0.5.
- $5 \times 5$  kernels in conv layers
- Gated Recurrent Units (GRU) in all recurrent layers. Simpler alternative to LSTM.
- Dropout in all layers 0.25
- Batch normalization in all non-recurrent layers
- Adam as optimizer, cross-entropy as loss
- Keras for Python, Theano backend, trained on Tesla GPUs



# Results

- Baseline: GMM and FNN
- CNN, RNN and CRNN with hyperparameters from grid search
- Mean and std for ten experiments with different random weight initialization



# Results - TUT-SED Synthetic 2016

**Table:** F1 score and error rate results for single frame segments ( $F1_{frm}$  and  $ER_{frm}$ ) and one second segments ( $F1_{1sec}$  and  $ER_{1sec}$ ). Bold face indicates the best performing method for the given metric. Standard deviation is presented after  $\pm$  sign.

TUT-SED Synthetic 2016				
Method	$F1_{frm}$	$ER_{frm}$	$F1_{1sec}$	$ER_{1sec}$
GMM [37]	40.5	0.78	45.3	0.72
FNN [15]	49.2 $\pm$ 0.8	0.68 $\pm$ 0.02	50.2 $\pm$ 1.4	1.1 $\pm$ 0.1
CNN	59.8 $\pm$ 0.9	0.56 $\pm$ 0.01	59.9 $\pm$ 1.2	0.78 $\pm$ 0.08
RNN	52.8 $\pm$ 1.5	0.6 $\pm$ 0.02	57.1 $\pm$ 0.9	0.64 $\pm$ 0.01
<b>CRNN</b>	<b>66.4<math>\pm</math>0.6</b>	<b>0.48<math>\pm</math>0.01</b>	<b>68.7<math>\pm</math>0.7</b>	<b>0.47<math>\pm</math>0.01</b>



# Effect of convolutional filter shape

**Table:**  $F1_{frm}$  for accuracy vs. convolution filter shape for TUT-SED Synthetic 2016 dataset.  $(*, *)$  represents filter lengths in frequency and time axis, respectively.

Filter shape	(3,3)	(5,5)	(11,11)	(1,5)	(5,1)	(3,11)	(11,3)
$F1_{frm}$	67.2	<b>68.3</b>	62.6	28.5	60.6	67.4	61.2

Thumb rule for CNNs: use  $3 \times 3$  kernels, strided convolutions instead of pooling, global average pooling at the end, many CONV - (L)RELU - BATCHNORM blocks or residual connections.



# Number of parameters vs. accuracy

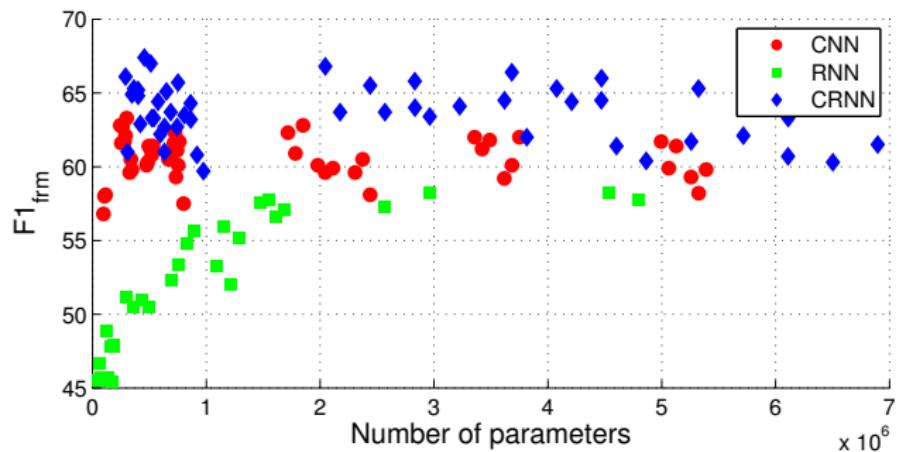


Figure: Number of parameters vs. accuracy for CNN, RNN and CRNN.



# Frequency shift invariance

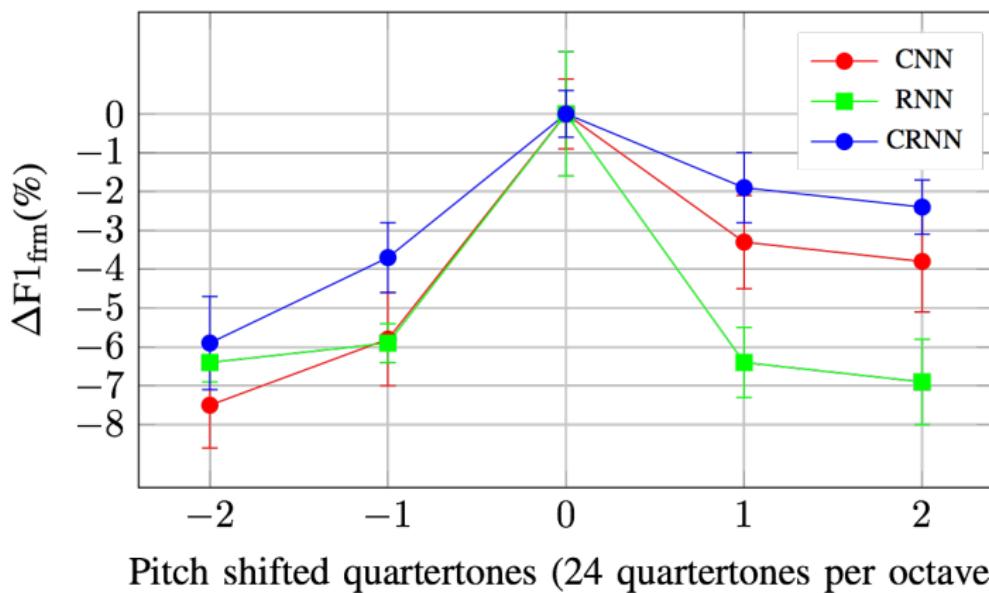
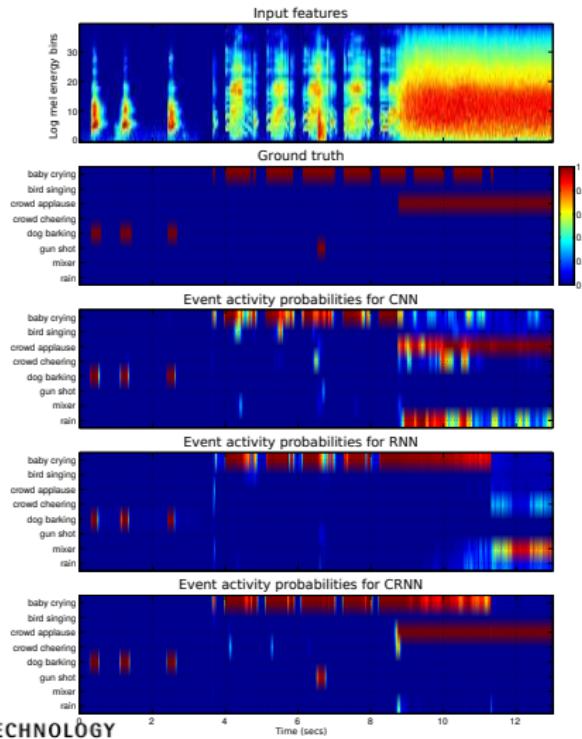


Figure: Accuracy vs. pitch-shifting over  $\pm 2$  quartertones for CNN, RNN and

# Closer look on network outputs



TAMPERE UNIVERSITY OF TECHNOLOGY

Figure: Input features, ground truth and event activity probabilities for CNN, RNN.

# Results - TUT-SED 2009 (CASA)

TUT-SED 2009 - CASA	
Method	$F1_{frm}$
GMM [37]	33.0
FNN [15]	$60.9 \pm 0.4$
CNN	$64.8 \pm 0.2$
RNN	$62.4 \pm 1.0$
<b>CRNN</b>	<b><math>69.7 \pm 0.4</math></b>

**Table:** For a cleaner view here we only show the results using the main metric. The results for the other metrics are consistent with these, and can be found in the paper.



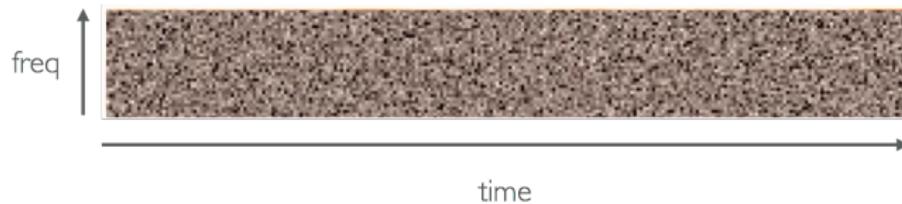
# Results - TUT-SED 2009 (CASA)

On the 2 smaller datasets there was little improvement :/



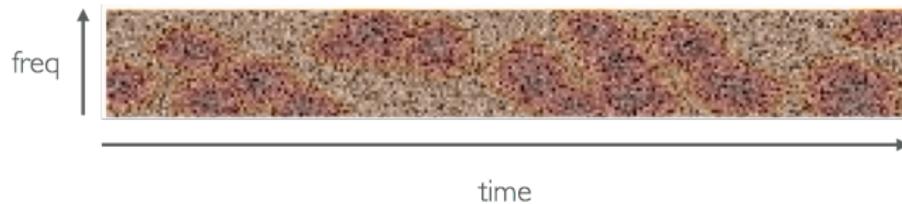
# What did our neurons learn?

Let's feed a random noise input, compute gradients to maximize the activation of a specific neuron and use them to modify the input and make it more likable to the neuron! (a cool new technique presented a few years ago)



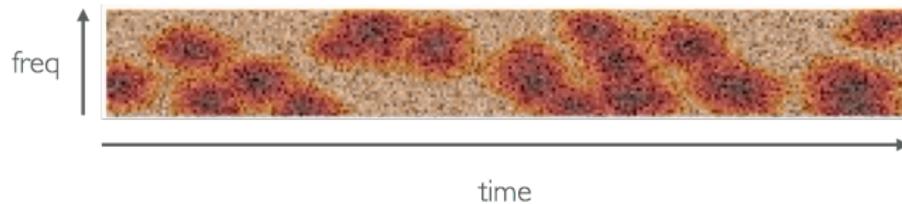
# What did our neurons learn?

Let's feed a random noise input, compute gradients to maximize the activation of a specific neuron and use them to modify the input and make it more likable to the neuron! (a cool new technique presented a few years ago)



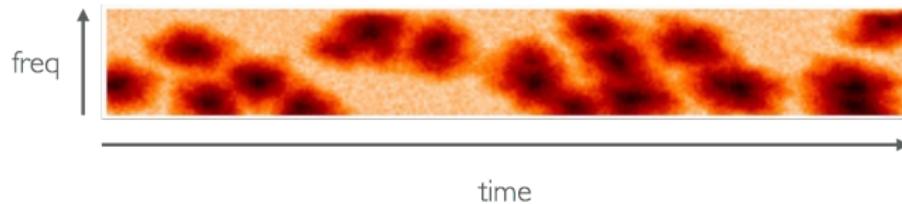
# What did our neurons learn?

Let's feed a random noise input, compute gradients to maximize the activation of a specific neuron and use them to modify the input and make it more likable to the neuron! (a cool new technique presented a few years ago)

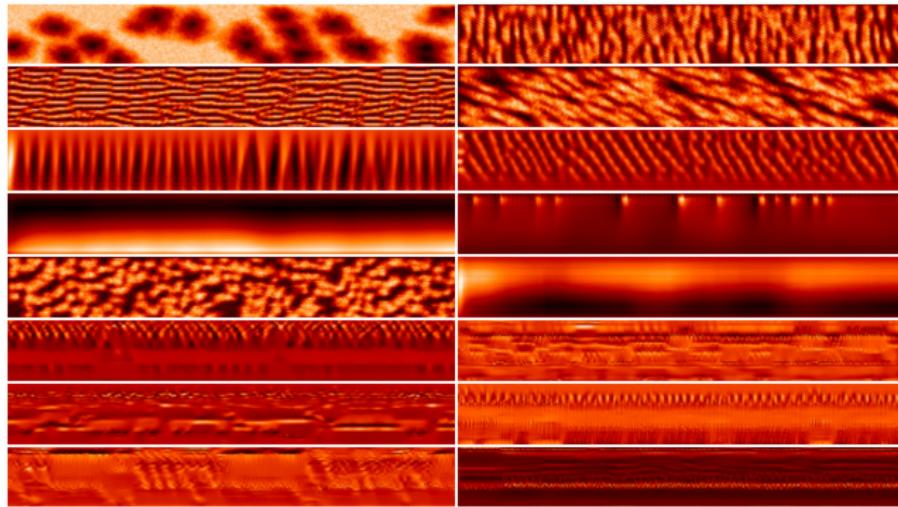


# What did our neurons learn?

Let's feed a random noise input, compute gradients to maximize the activation of a specific neuron and use them to modify the input and make it more likable to the neuron! (a cool new technique presented a few years ago)



# Visualization of convolutional layers



**Figure:** Two columns of crops from input patterns that would strongly activate certain neurons from different layers of the CRNN. On the horizontal axis is time, on the vertical axis mel bands. On both columns the rows 1 and 2 are from neurons in the first convolutional layer, columns 3 to 5 from the second, and columns from 6 to 8 from the third.



# Conclusions

- CRNNs considerably improve the results on SED compared to CNN and RNN alone



# Conclusions

- CRNNs considerably improve the results on SED compared to CNN and RNN alone
- Small datasets make training very difficult



# Possible ways to improve?

- ?
- ?
- ?
- ?



# Possible ways to improve

- Data augmentation



TAMPERE UNIVERSITY OF TECHNOLOGY

# Possible ways to improve

- Data augmentation
- Transfer learning (e.g. train the network first on a large dataset, then fine-tune on the smaller one)



# Possible ways to improve

- Data augmentation
- Transfer learning (e.g. train the network first on a large dataset, then fine-tune on the smaller one)
- Residual connections (Search for ResNets)



# Possible ways to improve

- Data augmentation
- Transfer learning (e.g. train the network first on a large dataset, then fine-tune on the smaller one)
- Residual connections (Search for ResNets)
- More data

