

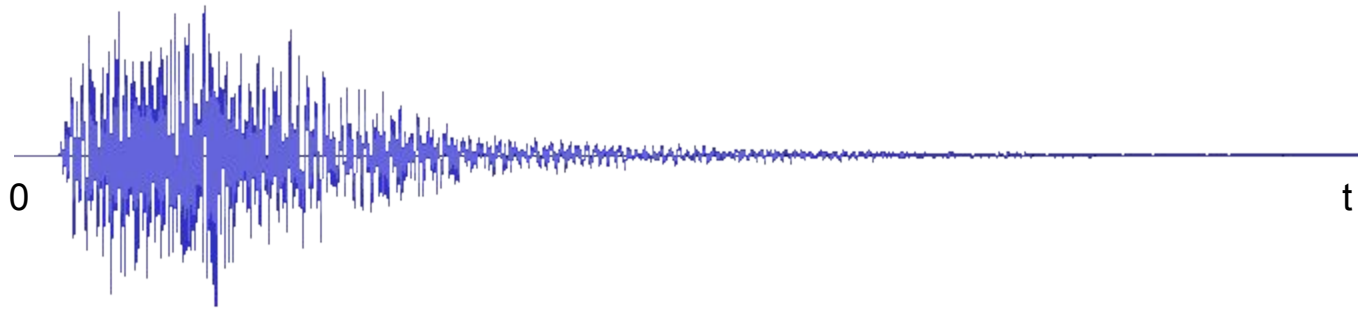
# Sound event detection using weakly labeled dataset with convolutional and recurrent neural network

Sharath Adavanne, Tuomas Virtanen

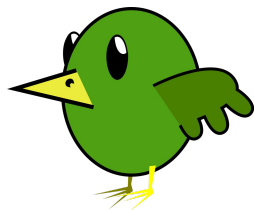
Laboratory of Signal Processing, Tampere University of Technology, Finland

# Outline

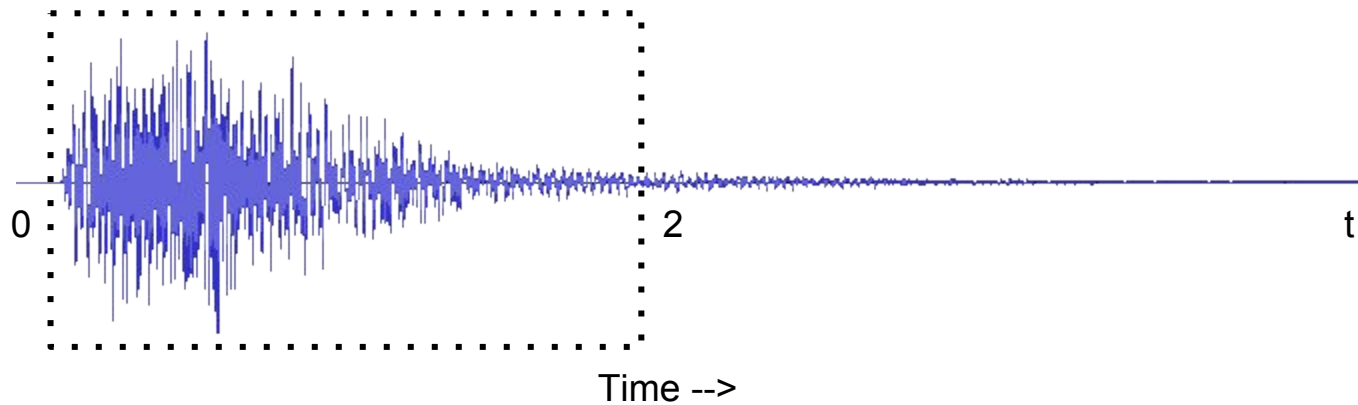
- Introduction
  - Sound tagging / Weak labels
  - Sound event detection / Strong labels
- Dataset
- Proposed neural network
- Results

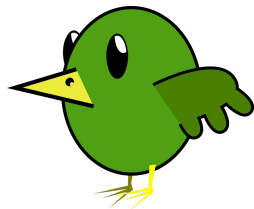


Time -->

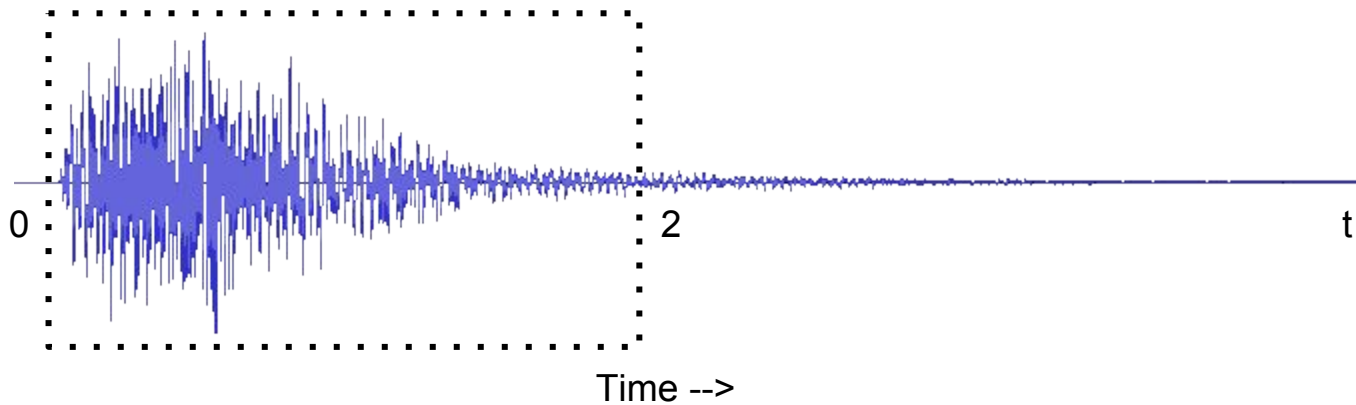


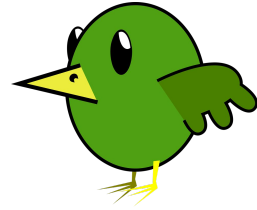
## Sound event detection



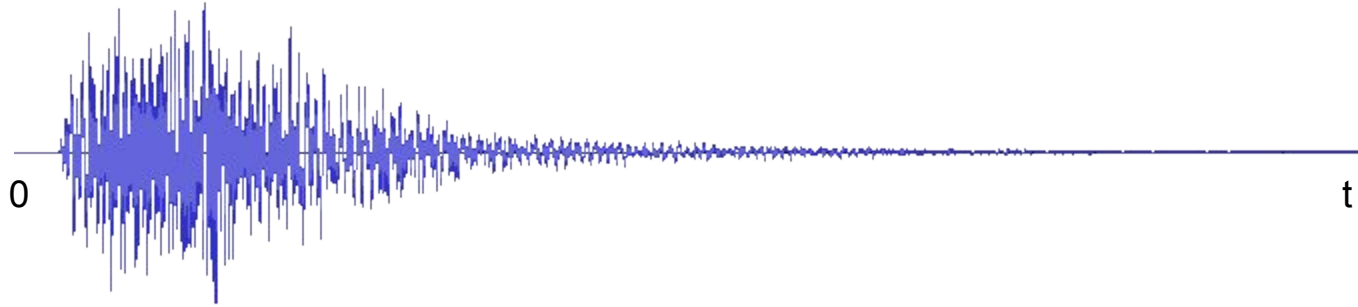


## Strong label Sound event detection

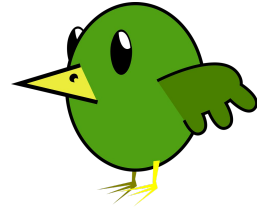




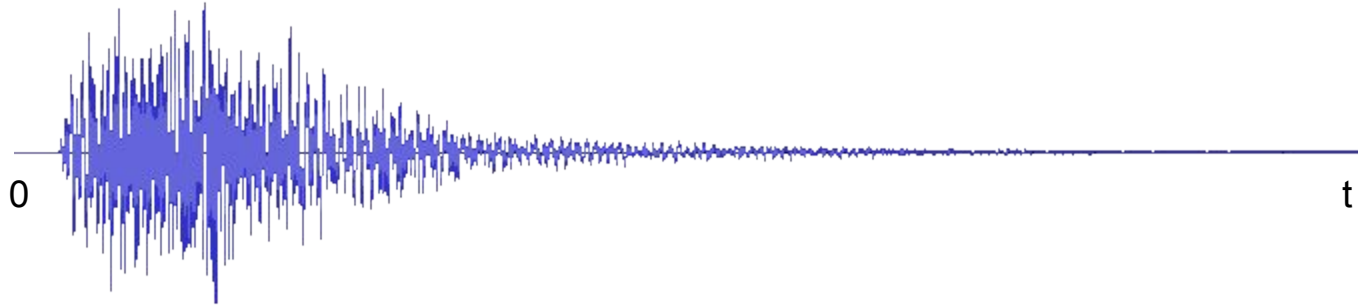
**Weak label**



Time -->



## Weak label Sound tagging



Time -->

Sound event detection using weakly labeled dataset  
with convolutional and recurrent neural network



Sound event detection using **weakly labeled** dataset  
with convolutional and recurrent neural network

Sound event detection using **sound tagging** dataset  
with convolutional and recurrent neural network

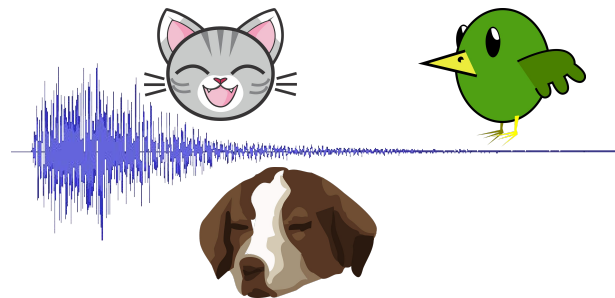
# Dataset

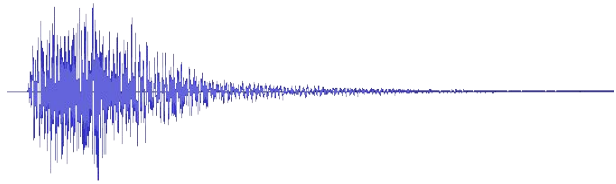
- Subset of Google's AudioSet
- 51,172 for training, 488 for testing
- 10s clips (padded with zeros, if not 10s)



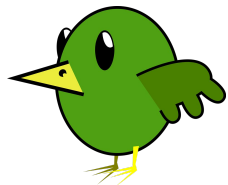
# Dataset

- Subset of Google's AudioSet
  - 51,172 for training, 488 for testing
  - 10s clips (padded with zeros, if not 10s)
- 
- Weak labeled
  - 17 classes - Car, Bus, Train, Truck etc..
  - Single recording can have more than one sound source



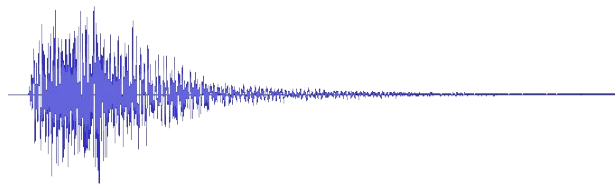
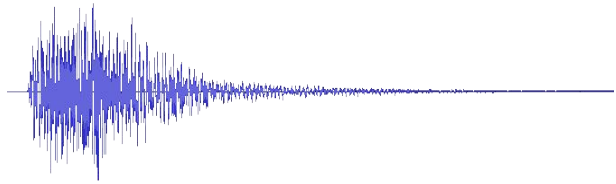


Neural network



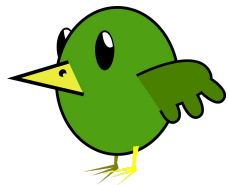
**Weak label  
Sound Tagging**

**Training procedure**

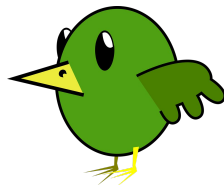


Training

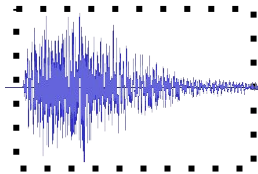
Testing



**Weak label  
Sound Tagging**



**Sound tagging**



**Sound event detection**

Input

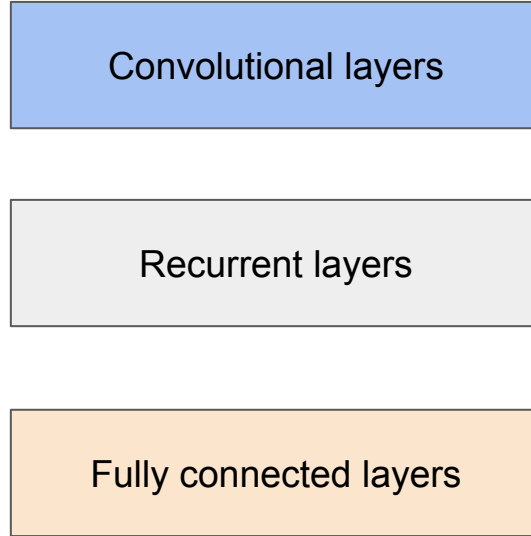
Convolutional layers

Recurrent layers

Fully connected layers

Output

Convolutional and  
recurrent neural network  
(CRNN) architecture



Input

Convolutional and  
recurrent neural network  
(CRNN) architecture

Convolutional layers

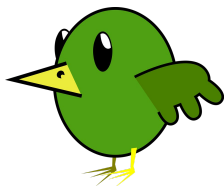
Recurrent layers

Fully connected layers

Why?

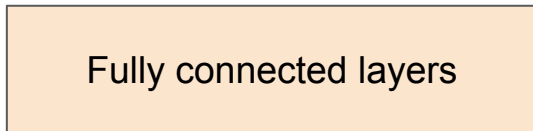
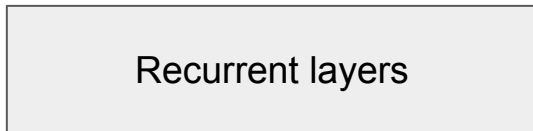
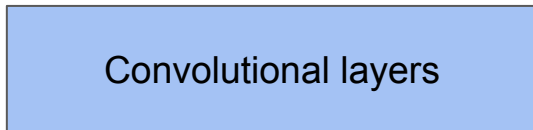
Output

Convolutional and recurrent neural network (CRNN) architecture

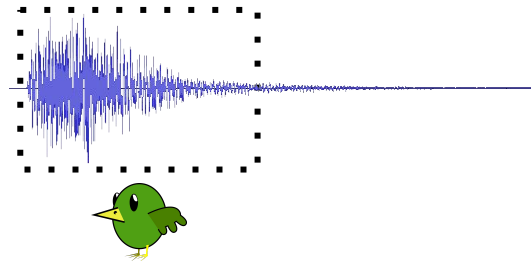


**Sound tagging**

Input



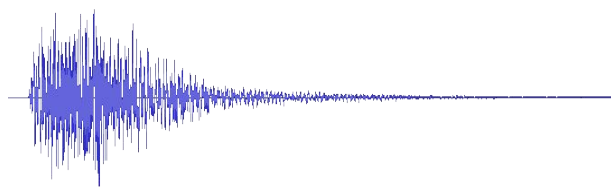
Output



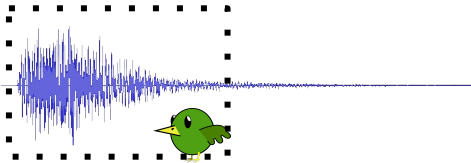
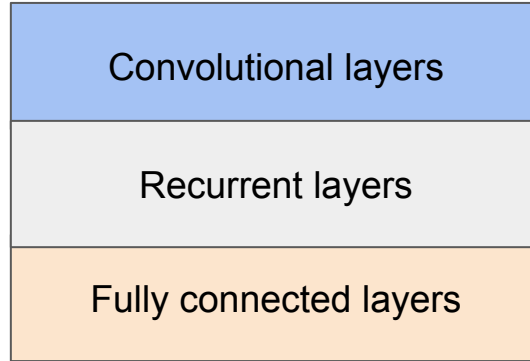
**Sound event detection**



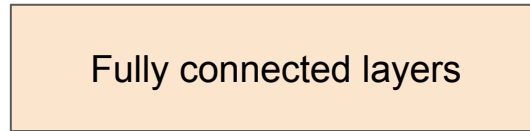
Convolutional and recurrent neural network (CRNN) architecture for sound event detection from weak labels



Input

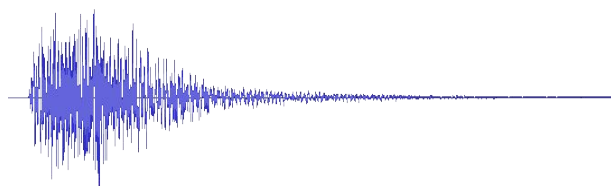


Output 1 : Sound event detection

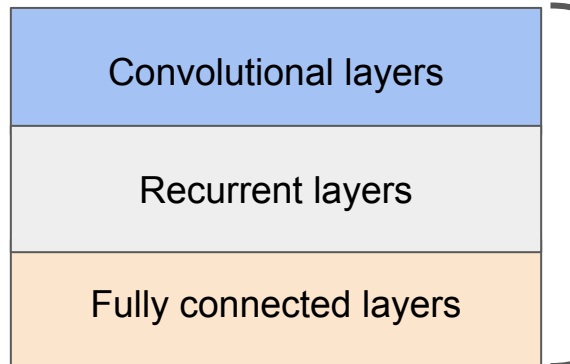


Output 2 : Sound tagging

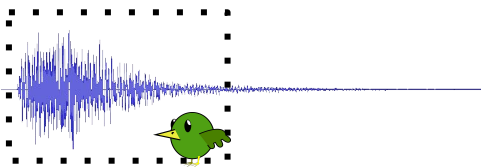
Convolutional and recurrent neural network (CRNN) architecture for sound event detection from weak labels



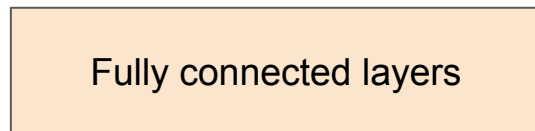
Input  
Log mel-band energies of complete 10 second audio



**Time-distributed layers:** Sequence input and sequence output



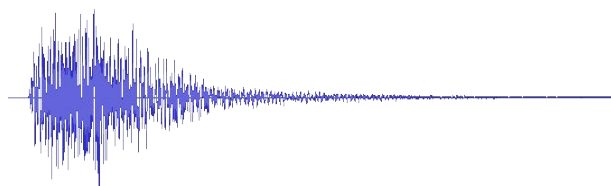
Output 1 : Sound event detection



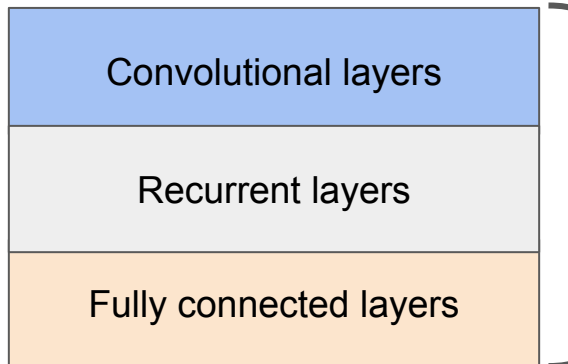
Sequence input and number of classes output

Output 2 : Sound tagging

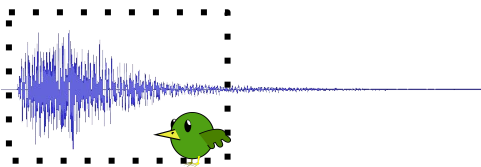
Convolutional and recurrent neural network (CRNN) architecture for sound event detection from weak labels



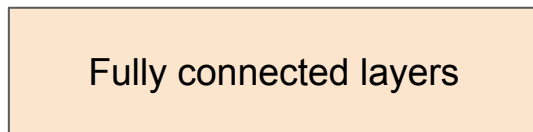
Input  
Log mel-band energies of complete 10 second audio



**Time-distributed layers:** Sequence input and sequence output



Output 1 : Sound event detection

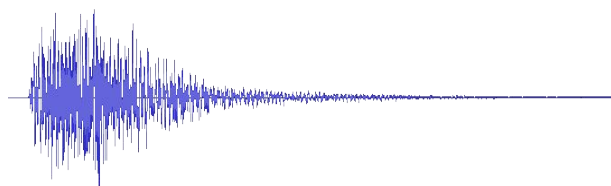
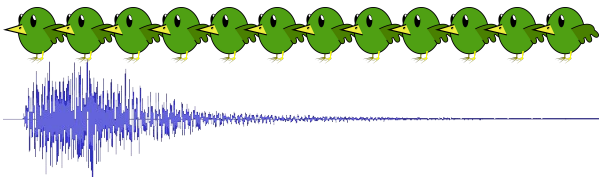


Sequence input and number of classes output

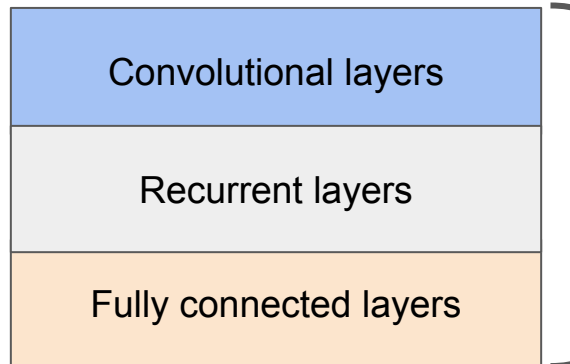


Output 2 : Sound tagging

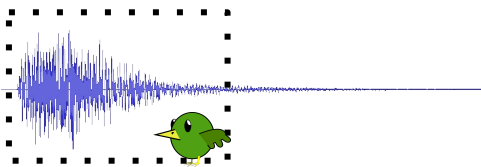
Convolutional and recurrent neural network (CRNN) architecture for sound event detection from weak labels



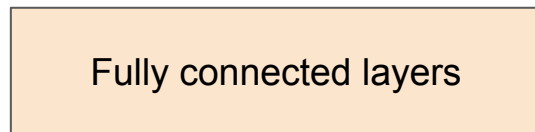
Input  
Log mel-band energies of complete 10 second audio



**Time-distributed layers:** Sequence input and sequence output



Output 1 : Sound event detection

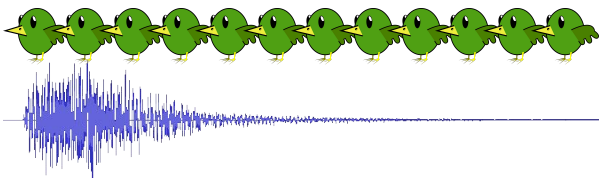


Sequence input and number of classes output

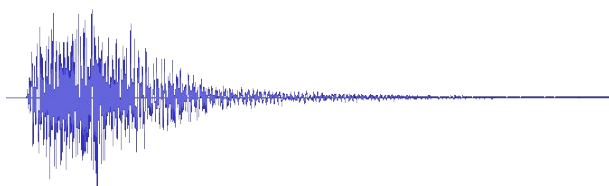


Output 2 : Sound tagging

Convolutional and recurrent neural network (CRNN) architecture for sound event detection from weak labels

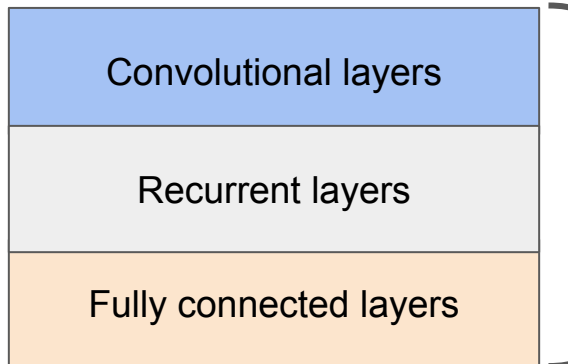


Loss 1 : strong label loss

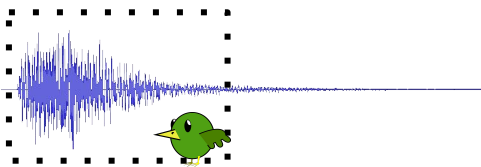


Input

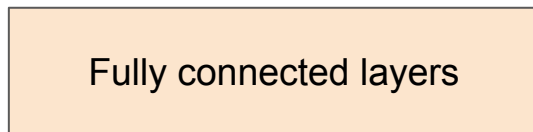
Log mel-band energies of complete 10 second audio



**Time-distributed layers:** Sequence input and sequence output



Output 1 : Sound event detection



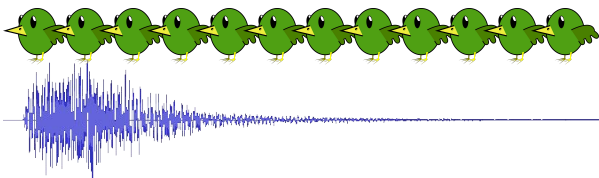
Sequence input and number of classes output

Loss 2 : weak label loss

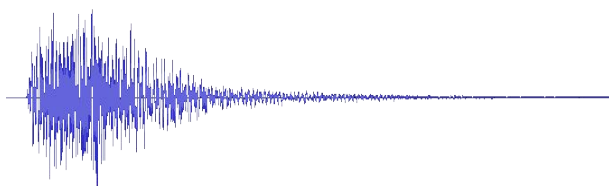


Output 2 : Sound tagging

Convolutional and recurrent neural network (CRNN) architecture for sound event detection from weak labels

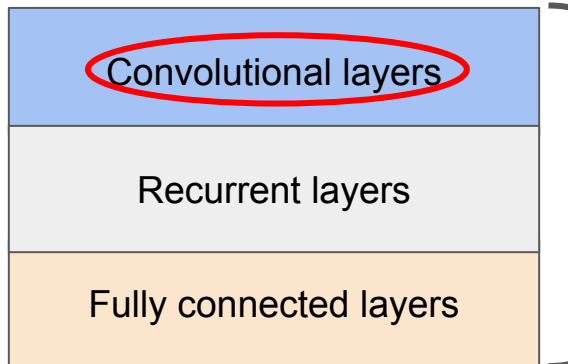


Loss 1 : strong label loss

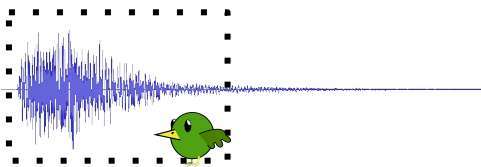


Input

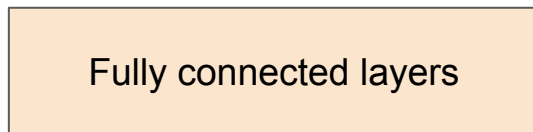
Log mel-band energies of complete 10 second audio



**Time-distributed layers:** Sequence input and sequence output



Output 1 : Sound event detection



Sequence input and number of classes output

Loss 2 : weak label loss



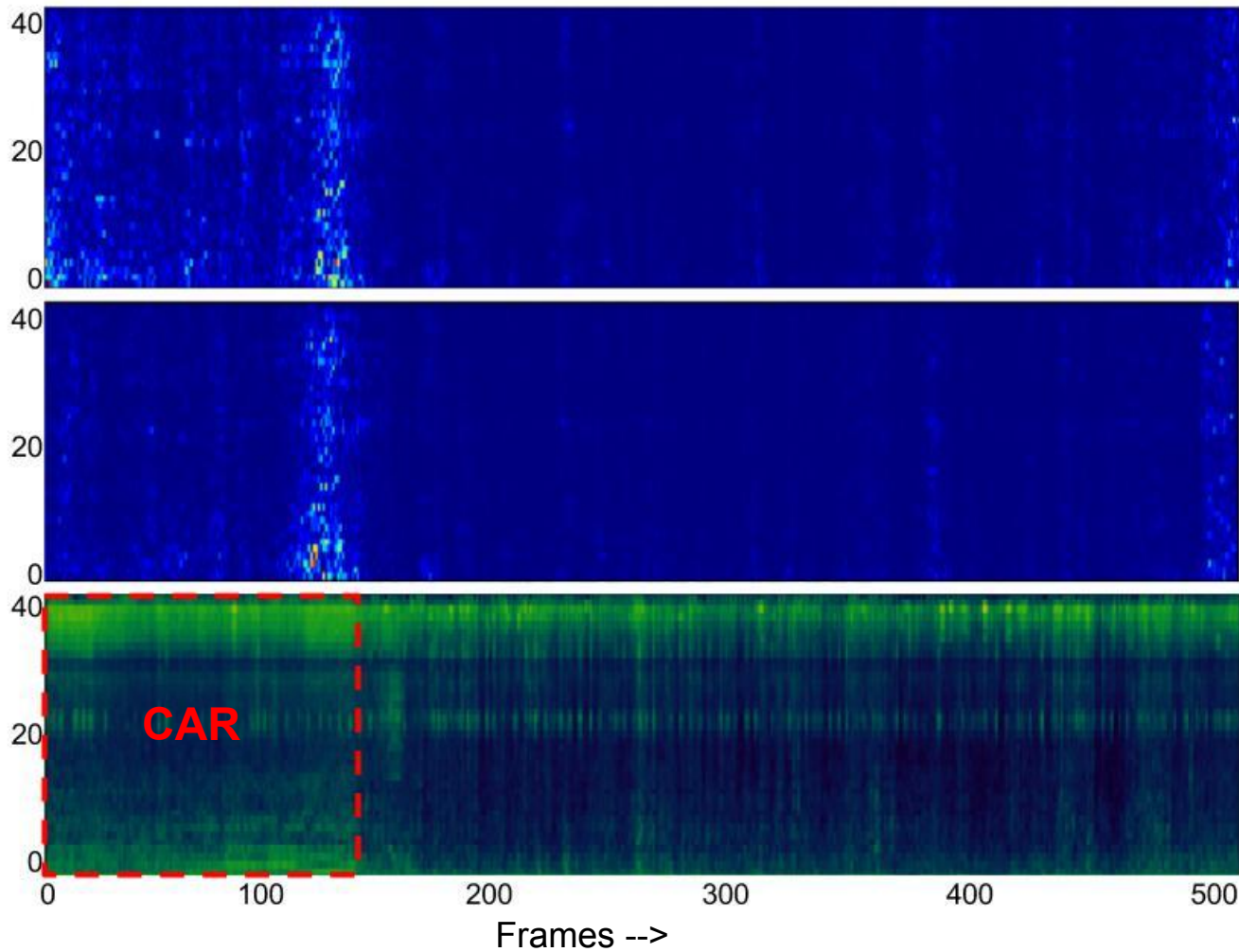
Output 2 : Sound tagging

**Saliency map for  
-jc0NAxK8M\_30.000\_40.000  
recording**

Gradients for strong label

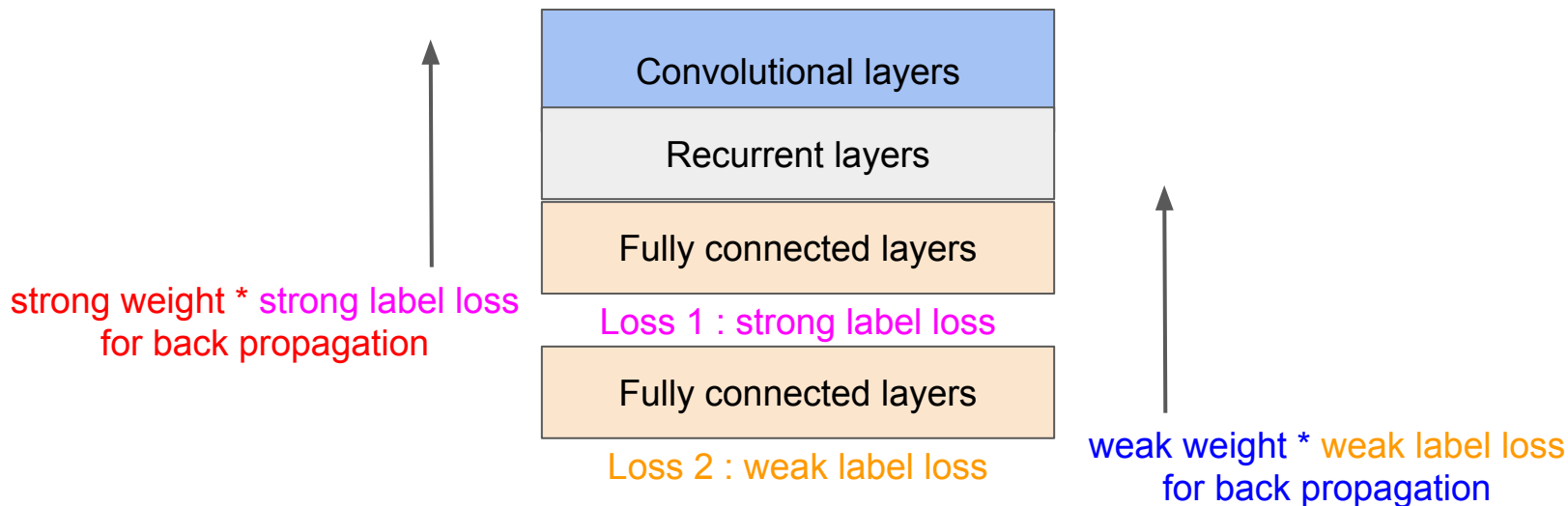
Gradients for weak label

Log mel-band energy



# Results

Strong Weight	Weak Weight	Sound tagging				Sound Event detection		
		Precision	Recall	F-score		Error rate	F-score	
0.002	1	44.9	37.0	40.5		1.38	10.9	





# Results

Strong Weight	Weak Weight	Sound tagging				Sound Event detection		
		Precision	Recall	F-score		Error rate	F-score	
0.002	1	44.9	37.0	40.5		1.38	10.9	
0.02	1	44.2	36.5	40.0		1.13	17.0	
0.2	1	<b>47.5</b>	39.6	43.2		0.84	38.1	
1	1	<b>47.5</b>	<b>39.7</b>	<b>43.3</b>		0.84	38.8	
1	0.2	47.3	39.5	43.0		0.84	38.6	
1	0.02	25.5	20.6	22.8		<b>0.81</b>	41.1	
1	0.002	20.5	16.5	18.3		<b>0.81</b>	<b>42.4</b>	

# Results

Strong Weight	Weak Weight	Sound tagging				Sound Event detection		
		Precision	Recall	F-score		Error rate	F-score	
0.002	1	44.9	37.0	40.5		1.38	10.9	
0.02	1	44.2	36.5	40.0		1.13	17.0	
0.2	1	<b>47.5</b>	39.6	43.2		0.84	38.1	
<b>1</b>	<b>1</b>	<b>47.5</b>	<b>39.7</b>	<b>43.3</b>		0.84	<b>38.8</b>	
1	0.2	47.3	39.5	43.0		0.84	38.6	
1	0.02	25.5	20.6	22.8		<b>0.81</b>	41.1	
1	0.002	20.5	16.5	18.3		<b>0.81</b>	<b>42.4</b>	

# Results

Strong Weight	Weak Weight	Sound tagging				Sound Event detection		
		Precision	Recall	F-score		Error rate	F-score	
0.002	1	44.9	37.0	40.5		1.38	10.9	
0.02	1	44.2	36.5	40.0		1.13	17.0	
0.2	1	<b>47.5</b>	39.6	43.2		0.84	38.1	
1	1	<b>47.5</b>	<b>39.7</b>	<b>43.3</b>		0.84	38.8	
1	0.2	47.3	39.5	43.0		0.84	38.6	
1	0.02	25.5	20.6	22.8		<b>0.81</b>	41.1	
1	0.002	20.5	16.5	18.3		<b>0.81</b>	<b>42.4</b>	

# Results

Strong Weight	Weak Weight	Sound tagging				Sound Event detection		
		Precision	Recall	F-score	F-score chng	Error rate	F-score	Error rate chng
0.002	1	44.9	37.0	40.5		1.38	10.9	
0.02	1	44.2	36.5	40.0		1.13	17.0	
0.2	1	<b>47.5</b>	39.6	43.2	<b>46.6</b>	0.84	38.1	0.80
1	1	<b>47.5</b>	<b>39.7</b>	<b>43.3</b>	45.5	0.84	38.8	0.81
1	0.2	47.3	39.5	43.0	44.5	0.84	38.6	0.82
1	0.02	25.5	20.6	22.8		<b>0.81</b>	41.1	
1	0.002	20.5	16.5	18.3	26.3	<b>0.81</b>	<b>42.4</b>	<b>0.79</b>

# Results

Strong Weight	Weak Weight	Sound tagging				Sound Event detection		
		Precision	Recall	F-score	F-score chng	Error rate	F-score	Error rate chng
0.002	1	44.9	37.0	40.5		1.38	10.9	
0.02	1	44.2	36.5	40.0		1.13	17.0	
0.2	1	<b>47.5</b>	39.6	43.2	<b>46.6</b>	0.84	38.1	0.80
1	1	<b>47.5</b>	<b>39.7</b>	<b>43.3</b>	45.5	0.84	38.8	0.81
1	0.2	47.3	39.5	43.0	44.5	0.84	38.6	0.82
1	0.02	25.5	20.6	22.8		<b>0.81</b>	41.1	
1	0.002	20.5	16.5	18.3	26.3	<b>0.81</b>	<b>42.4</b>	<b>0.79</b>
					55.6			0.66

# Outline

- **Introduction**
  - Sound tagging / Weak labels
  - Sound event detection / Strong labels

# Outline

- **Introduction**
  - Sound tagging / Weak labels
  - Sound event detection / Strong labels
- **Dataset**
  - Weak labels
  - More than one label for each recording

# Outline

- **Introduction**

- Sound tagging / Weak labels
- Sound event detection / Strong labels

- **Dataset**

- Weak labels
- More than one label for each recording

- **Proposed neural network**

- Convolutional recurrent neural network (CRNN)
- Two sequential outputs - strong label followed by weak



# Outline

- **Introduction**

- Sound tagging / Weak labels
- Sound event detection / Strong labels

- **Dataset**

- Weak labels
- More than one label for each recording

- **Proposed neural network**

- Convolutional recurrent neural network (CRNN)
- Two sequential outputs - strong label followed by weak

- **Results**

- CRNN was shown to learn temporal information, given just weak labels
- Best result was against our intuition: equal scaling of strong and weak loss

# Outline

- **Introduction**

- Sound tagging / Weak labels
- Sound event detection / Strong labels

- **Dataset**

- Weak labels
- More than one label for each recording

- **Proposed neural network**

- Convolutional recurrent neural network (CRNN)
- Two sequential outputs - strong label followed by weak

- **Results**

- CRNN was shown to learn temporal information, given just weak labels
- Best result was against our intuition: equal scaling of strong and weak loss

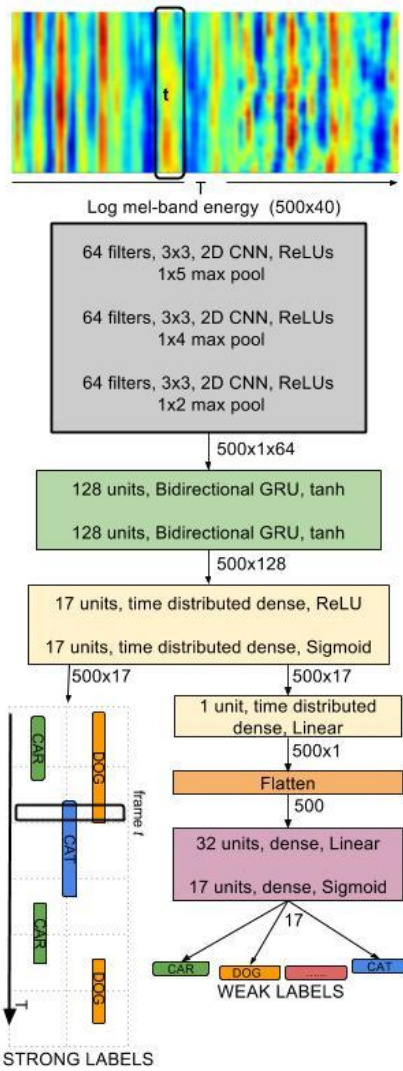
- **Future work**

- More fine tuning
- Strong labels for high energy regions only
- Attention layers

Thank you.

Sharath Adavanne, Tuomas Virtanen

Laboratory of Signal Processing, Tampere University of Technology, Finland



218,000 trainable weights