VERY DEEP CONVOLUTIONAL NEURAL NETWORKS FOR RAW WAVEFORMS

#### Overview

- Classifies raw Audio waveform into 10 classes.
- Uses only Conv/Pool Layers.
- Uses Batch Normalization & Glorot initialization to avoid exploding or vanishing gradients.
- Builds models from scratch (no pre-training).
- Proposes five different Models.

#### Data

- The dataset consists of 8732 audio clips (7GB).
- Each audio clip is a wav file (≈ 4 seconds).
- Sample rate: 44100 HZ
- Channels: Stereo
- Total duration: 9.7 hours.

#### Data Classes

- 0) air\_conditioner
- •
- 1) Car\_horn

- 1
- 2) children\_playing



3) dog\_bark

- •
- 4) drilling

- **1**
- 5) engine\_idling



6) gun\_shot

•

7) jackhammer

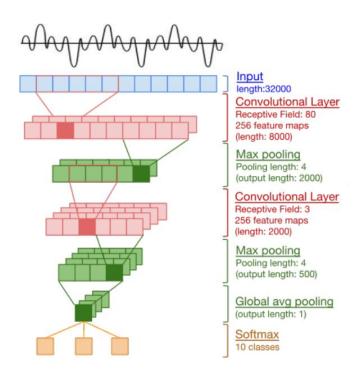
- 1
- 8) siren



9) street\_music



## Models



M3 (0.2M)	M5 (0.5M)	M11 (1.8M)	M18 (3.7M)	M34-res (4M)	
7	Input: 32	2000x1 time-doi	main waveform		
[80/4, 256]	[80/4, 128]	[80/4, 64]	[80/4, 64]	[80/4, 48]	
	Maxp	ool: 4x1 (outpu	t: $2000 \times n$ )		
[3, 256]	[3, 128]	$[3, 64] \times 2$	$[3, 64] \times 4$	$\left[\begin{array}{c}3,48\\3,48\end{array}\right]\times3$	
	Max	pool: 4x1 (outp	ut: 500×n)		
	[3, 256]	$[3, 128] \times 2$	$[3, 128] \times 4$	$\left[\begin{array}{c}3,96\\3,96\end{array}\right]\times4$	
	Maxpool: $4x1$ (output: $125 \times n$ )				
	[3, 512]	$[3, 256] \times 3$	$[3, 256] \times 4$	$\left[\begin{array}{c} 3,192\\ 3,192 \end{array}\right] \times 6$	
	Maxpool: $4x1$ (output: $32 \times n$ )				
		$[3,512] \times 2$	$[3,512] \times 4$	$\left[\begin{array}{c}3,384\\3,384\end{array}\right]\times3$	
	Global a	verage pooling	$(output: 1 \times n)$		
		Softmax			

# My Implementation

My implementation is divided into three different parts:

- 1. Preprocessing the data.
- 2. Creating the models.
- 3. Replicating the results.

# Data Preprocessing

```
class ToMono(torch.nn.Module):
   def forward(self, waveform: torch.Tensor) -> torch.Tensor:
        return torch.mean(waveform, dim=0, keepdim=True)
class Normalize(torch.nn.Module):
   def forward(self, waveform: torch.Tensor) -> torch.Tensor:
       return (waveform-waveform.mean()) / waveform.std()
class Pad(torch.nn.Module):
   def __init__(self, value: float, size: int):
       super(Pad, self).__init__()
       self.value = value
       self.size = size
   def forward(self, waveform: torch.Tensor) -> torch.Tensor:
        return torch.nn.functional.pad(waveform, (0, self.size-max(waveform.shape)), "constant", self.value)
audio_transform = torch.nn.Sequential(*[
   ToMono(), #converts audio channels to mono
    torchaudio.transforms.Resample(orig_freq=441000, new_freq=8000), # downsamples audio signal to 8000 HZ
   Normalize(), # normalize audio signal to have mean=0 & std=1
   Pad(value=0, size=32000),
```

# Models

#### Two base classes:

- 1. CNN
- 2. CCNRes

M3 (0.2M)	M5 (0.5M)	M11 (1.8M)	M18 (3.7M)	M34-res (4M)		
Input: 32000x1 time-domain waveform						
[80/4, 256]	[80/4, 128]	[80/4, 64]	[80/4, 64]	[80/4, 48]		
	Maxpe	ool: 4x1 (outpu	t: 2000 × n)			
[3, 256]	[3, 128]	$[3, 64] \times 2$	$[3, 64] \times 4$	$\left[\begin{array}{c}3,48\\3,48\end{array}\right]\times3$		
Maxpool: 4x1 (output: 500×n)						
	[3, 256]	$[3, 128] \times 2$	$[3, 128] \times 4$	$\left[\begin{array}{c}3,96\\3,96\end{array}\right]\times4$		
	Maxpool: $4x1$ (output: $125 \times n$ )					
	[3, 512]	$[3, 256] \times 3$	$[3, 256] \times 4$	$\left[\begin{array}{c}3,192\\3,192\end{array}\right]\times 6$		
	Maxpool: $4x1$ (output: $32 \times n$ )					
		$[3,512] \times 2$	$[3,512] \times 4$	$\left[\begin{array}{c}3,384\\3,384\end{array}\right]\times3$		
	Global a	verage pooling	(output: $1 \times n$ )			
		Softmax				

CNN

**CNNRes** 

#### CNN Base Class

```
def __init__(self, channels, conv_kernels, conv_strides, conv_padding, pool_padding, num_classes=10):
    assert len(conv_kernels) == len(channels) == len(conv_strides) == len(conv_padding)
    super(CNN, self).__init__()
   prev_channel = 1
    for i in range(len(channels)):
       block = []
        for j, conv_channel in enumerate(channels[i]):
           block.append( torch.nn.Conv1d(in_channels = prev_channel, out_channels = conv_channel, kernel_size = conv_kernels[i], stride = conv_strides[i], padding = conv_padding[i]) )
           prev_channel = conv_channel
           block.append( torch.nn.BatchNorm1d(prev_channel) )
           block.append( torch.nn.ReLU() )
        self.conv_blocks.append( torch.nn.Sequential(*block) )
    for i in range(len(pool_padding)):
        self.pool_blocks.append( torch.nn.MaxPoolId(kernel_size = 4, stride = 4, padding = pool_padding[i]) )
    self.linear = torch.nn.Linear(prev_channel, num_classes)
def forward(self, inwav):
    for i in range(len(self.conv_blocks)):
       inwav = self.conv_blocks[i](inwav)
        if i < len(self.pool_blocks): inwav = self.pool_blocks[i](inwav)</pre>
    out = self.global_pool(inwav).squeeze()
    out = self.linear(out)
    return out.squeeze()
```

#### CNN Derived classes

```
m3 = CNN(channels = [[256], [256]],
         conv_kernels = [80, 3],
         conv_strides = [4, 1],
         conv_padding = [38, 1],
         pool_padding = [0, 0])
m5 = CNN(channels = [[128], [128], [256], [512]],
         conv_kernels = [80, 3, 3, 3],
         conv_strides = [4, 1, 1, 1],
         conv_padding = [38, 1, 1, 1],
         pool_padding = [0, 0, 0, 2])
m11 = CNN(channels = [[64], [64]*2, [128]*2, [256]*3, [512]*2],
          conv_kernels = [80, 3, 3, 3, 3],
          conv_strides = [4, 1, 1, 1, 1],
          conv_padding = [38, 1, 1, 1, 1],
          pool_padding = [0, 0, 0, 2])
m18 = CNN(channels = [[64], [64]*4, [128]*4, [256]*4, [512]*4],
          conv_kernels = [80, 3, 3, 3, 3],
          conv_strides = [4, 1, 1, 1, 1],
          conv_padding = [38, 1, 1, 1, 1],
          pool_padding = [0, 0, 0, 2])
```

### CNNRes Base Class

```
def __init__(self, channels, conv_kernels, conv_strides, conv_padding, pool_padding, num_classes=10):
    assert len(conv_kernels) == len(channels) == len(conv_strides) == len(conv_padding)
    super(CNNRes, self).__init__()
    prev channel = 1
       torch.nn.Conv1d(in_channels = prev_channel, out_channels = channels[0][0], kernel_size = conv_kernels[0], stride = conv_strides[0], padding = conv_padding[0]).
       torch.nn.BatchNorm1d(channels[0][0]),
        torch.nn.MaxPoolld(kernel_size = 4, stride = 4, padding = pool_padding[0]),
    prev_channel = channels[0][0]
        for j, conv_channel in enumerate(channels[i]):
            block.append( ResBlock(prev_channel, conv_channel, conv_kernels[i], conv_strides[i], conv_padding[i]) )
            prev channel = conv channel
        self.res_blocks.append( torch.nn.Sequential(*block) )
    for i in range(1, len(pool_padding)):
        self.pool_blocks.append( torch.nn.MaxPoolld(kernel_size = 4, stride = 4, padding = pool_padding[i]) )
    self.linear = torch.nn.Linear(prev_channel, num_classes)
def forward(self, inwav):
    inway = self.conv_block(inway)
        inwav = self.res_blocks[i](inwav)
        if i < len(self.pool_blocks): inwav = self.pool_blocks[i](inwav)
    out = self.global_pool(inwav).squeeze()
    out = self.linear(out)
    return out.squeeze()
```

#### M32-Res

```
m32 = CNNRes(channels = [[48], [48]*3, [96]*4, [192]*6, [384]*3],

conv_kernels = [80, 3, 3, 3, 3],

conv_strides = [4, 1, 1, 1, 1],

conv_padding = [38, 1, 1, 1, 1],

pool_padding = [0, 0, 0, 2])
```

#### Results

Model	Test
M3	56.12%
M5	63.42%
M11	69.07%
M18	71.68%
M34-res	63.47%

M3: 35% (15 epochs)

M5: 45% (50 epochs)

M11: 35% (10 epochs)

M18: 32% (8 epochs)

M34-res = 37% (10 epochs)

# Code

Check the GitHub Repository:

https://github.com/Anwarvic/CNN-for-Raw-Waveforms

# Questions?