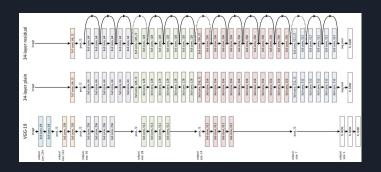
X-Vector embedding

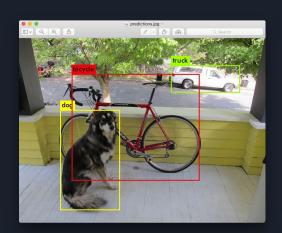
A modern approach to vectorize speaker identity.

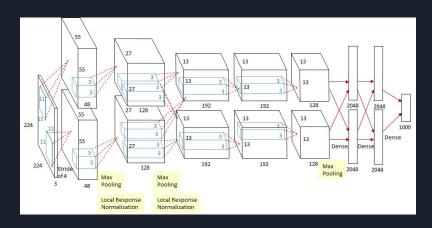


The need for embedding

The state of the art of image recognition







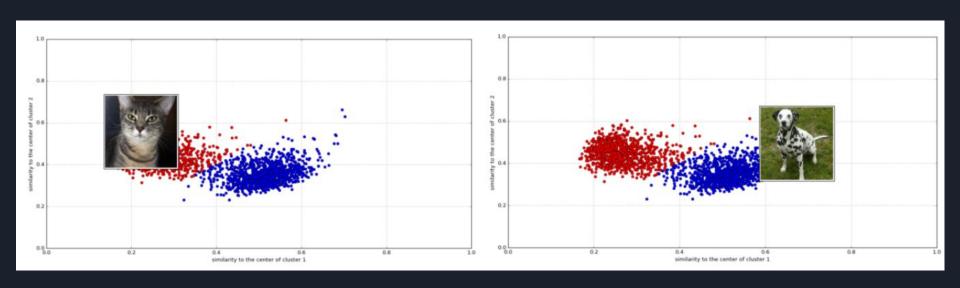
What are the input features?

Well the image itself!

and we are done....

		165	187	209	58	7
		100	107	200	90	
	14	125	233	201	98	159
253	144	120	251	41	147	204
67	100	32	241	23	165	30
209	118	124	27	59	201	79
210	236	105	169	19	218	156
35	178	199	197		14	218
115	104	34	111	19	196	
32	69	231		74		

But what if we would like to cluster the images?

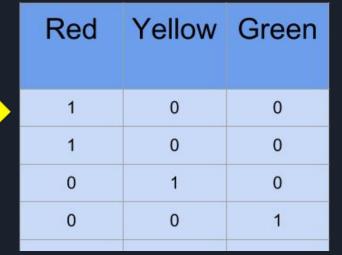


One step further....

let's say now that's our input is a text

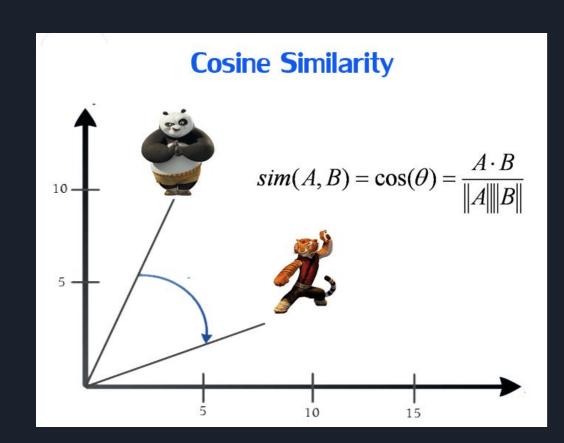
one hot encoding?

Color
Red
Red
Yellow
Green
Yellow

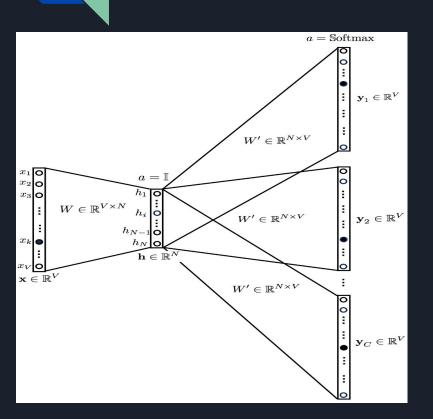


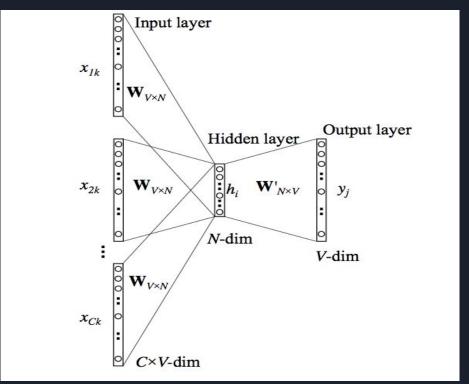
we just lost the semantic context of a word.... Our objective is to have words with similar context occupy close spatial positions.

Mathematically, the cosine of the angle between such vectors should be close to 1, i.e. angle close to 0.

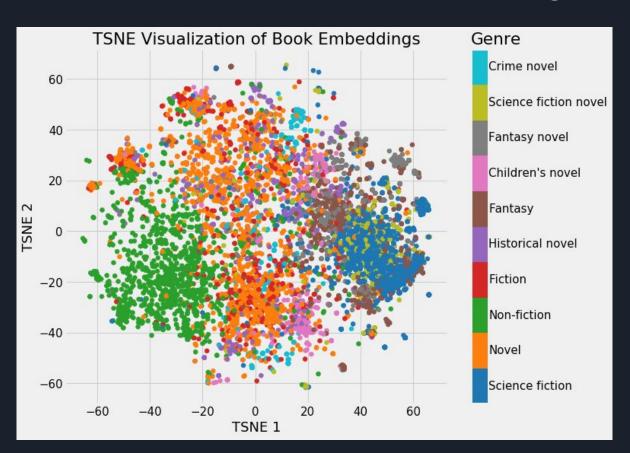


Word2vec came to help





The need for fixed size embedding

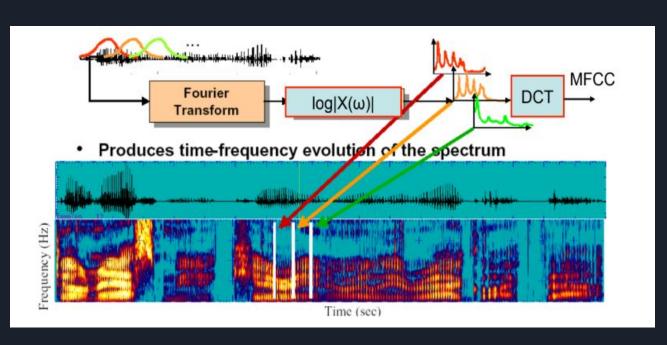


And now for speech embedding

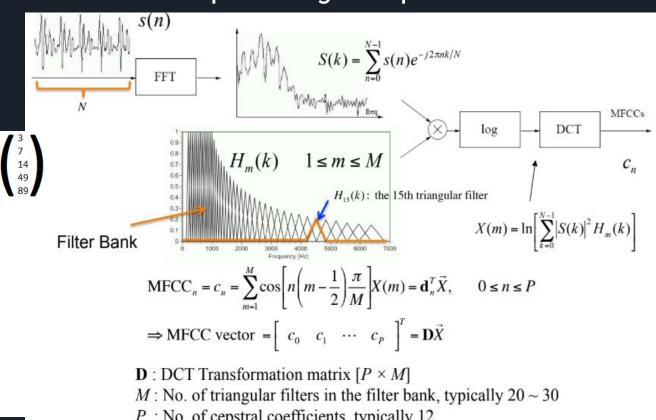
The first line of speech features

Speech is a time-varying signal conveying multiple layers of information

- Words
- Speaker
- Language
- Imotion



MFCC - Mel-frequency cepstrum



0.5 0 0.01 0.02 0.03 0.04 0.05 0.06 0.07

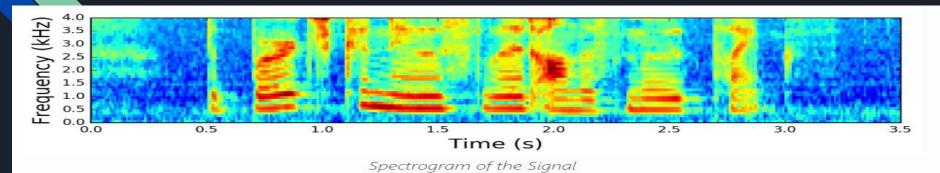
1k

2k

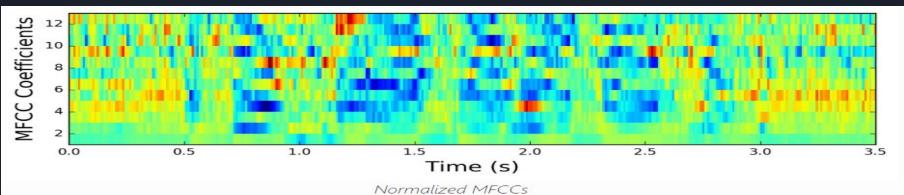
500

P: No. of cepstral coefficients, typically 12 c_0 : Logarithm of energy of the current frame

MFCC



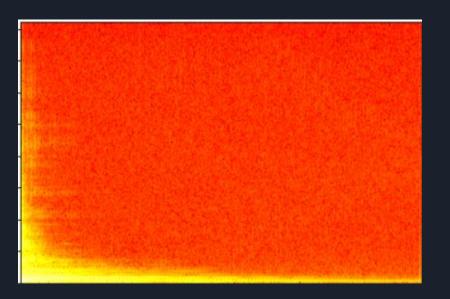




MFCC - Intuition

each musical instrument produces unique "fingerprint"

Acoustic guitar spectrogram vs MFCC

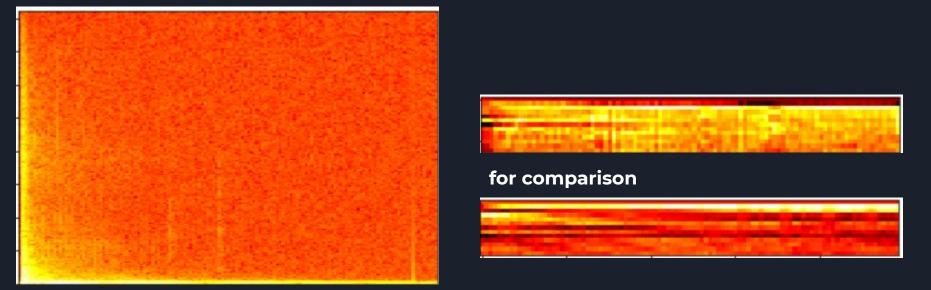




MFCC - Intuition

each musical instrument produces unique "fingerprint"

Bass drum spectrogram vs MFCC



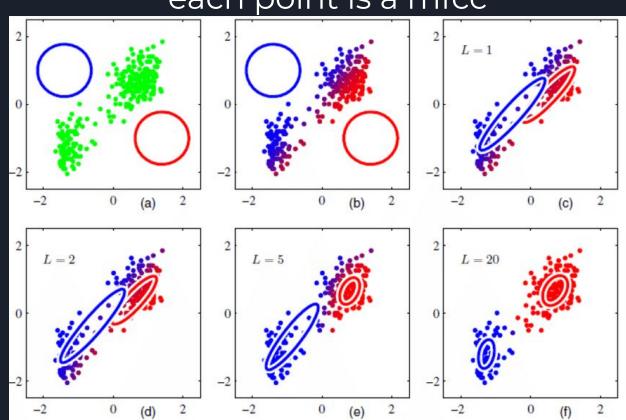
And now for the embedding part...

The pre Deep Learning era - The **I-vector**Unsupervised embedding model
each point is a mfcc

UBM - GMM – Universal Background Model

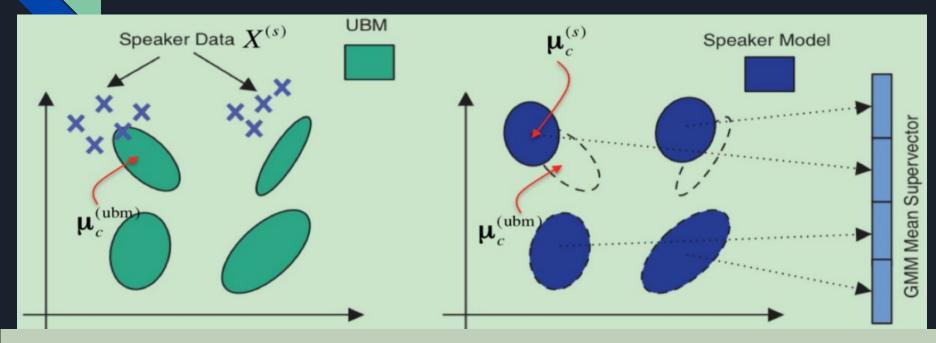
Use MFCC as features

Use 2048 (standard value)
Gaussians to describe each of the features



Adaptive GMM - Extract

Adapt the UBM to the speaker using AGMM (using EM or MAP adaptation)



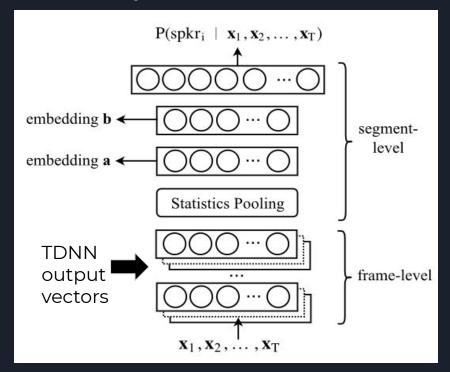
The 65k parameters vector aka Super Vector, via dimensionality reduction the commonly **I-vector** is produced (using un/supervised methods) ville Vestman, Tomi Kinnunen 2018

The X-vector

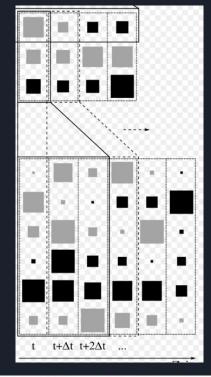
X - vector architecture

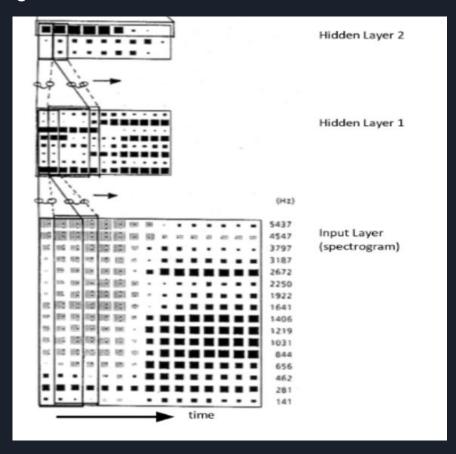
x-vector consists of three main parts:

- TDNN layers
- Stats pooling
- Fully connected layers

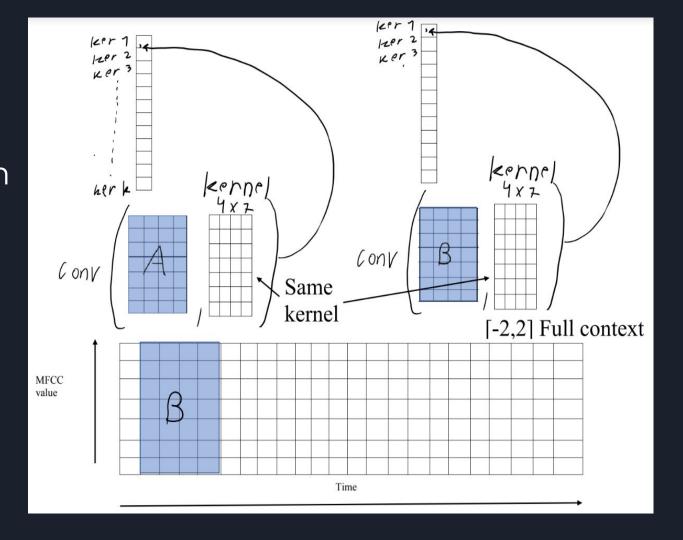


The TDNN - Time delay neural net

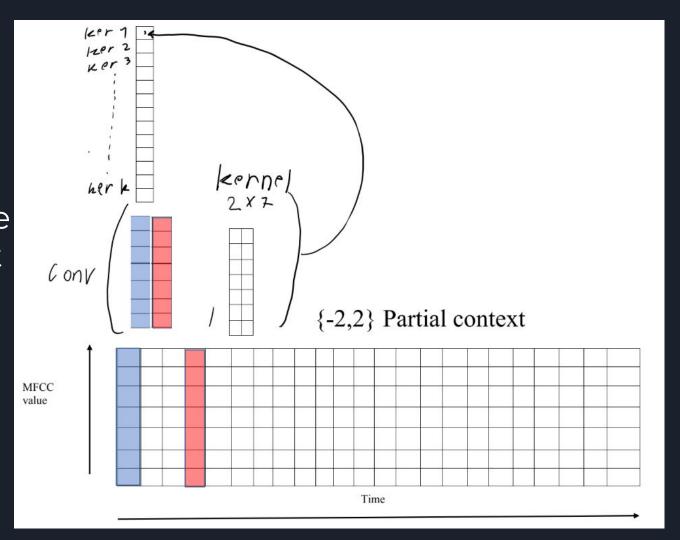




Implementation using 1d convolution layer



The TDNN introduced the partial context idea



TDNN usage beside X-vectors

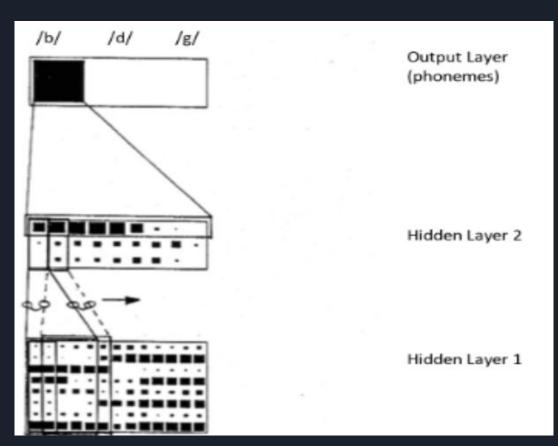
The TDNN could be used for phonemes recognition.

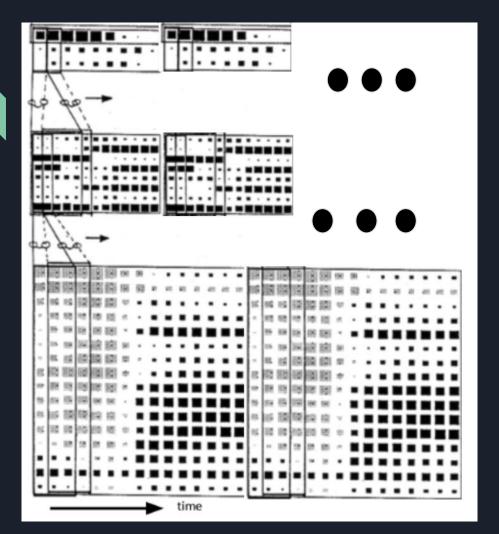
ALEXANDER WAIBEL et al.

And even for keyword spotting

Ming Sun, David Snyder et al

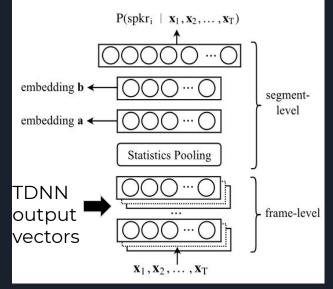
Depending on the output layer of our choice...



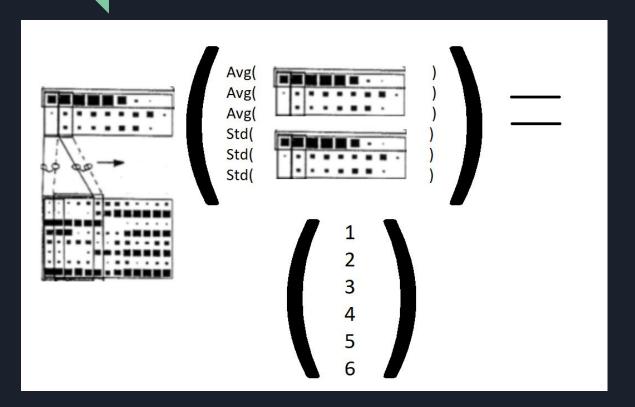


Depending on the input length we will receive **dynamic length** end layer.

Which means we still don't have the **fixed size** vector that we wanted



Statistics - pooling

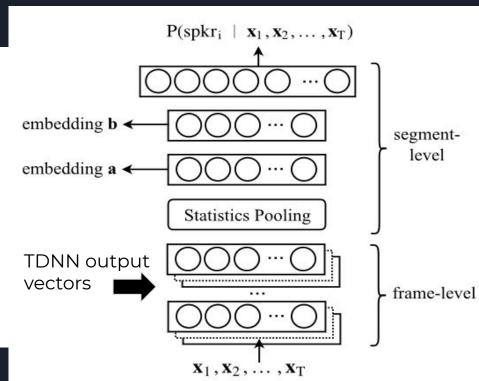


Simply take the average and the standard deviation for each channel of the final TDNN output and concatenate them

And now we got our desired **fixed length** vector

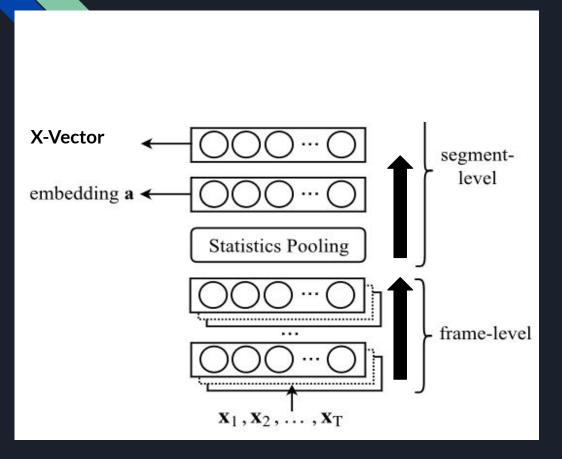
Two final steps and we have out x - vector

Layer	Layer context	Total context	Input x output
frame1	[t-2, t+2]	5	120x512
frame2	$\{t-2, t, t+2\}$	9	1536x512
frame3	$\{t-3, t, t+3\}$	15	1536x512
frame4	$\{t\}$	15	512x512
frame5	$\{t\}$	15	512x1500
stats pooling	[0,T)	T	1500Tx3000
segment6	{0}	T	3000x512
segment7	{0}	T	512x512
softmax	{0}	T	512x <i>N</i>



The stats pooling layer output is than sent to two Fully - connected layers

Chopping the network head



After learning the network parameters we chop off the softmax layer, and forward propagate an input signal

Demonstration

The most common public implementation of X-vectors is using the Kaldi framework

Until now...

Pytorch implementation of X-vectors

```
context = [-2,2]
    input dim = 24
    output dim = 512
   net1 = TDNN(context, input dim, output dim, full context=True)
   context = [-2, 1, 2]
    input dim = 512
   output dim = 512
   net2 = TDNN(context, input dim, output dim, full context=False)
10
11
   context = [-3, 1, 3]
12
    input dim = 512
    output dim = 512
13
                                                         Layer
                                                                   Layer context
                                                                                  Total context
                                                                                               Input x output
    net3 = TDNN(context, input dim, output dim, fu
                                                        frame1
                                                                    [t-2, t+2]
                                                                                                  120x512
15
                                                        frame2
                                                                  \{t-2, t, t+2\}
                                                                                       9
                                                                                                 1536x512
16
   context = [1]
                                                                  \{t-3, t, t+3\}
                                                        frame3
                                                                                      15
                                                                                                 1536x512
    input dim = 512
                                                                                      15
18
    output dim = 512
                                                         frame4
                                                                                                  512x512
   net4 = TDNN(context, input dim, output dim, fu
19
                                                         frame5
                                                                        \{t\}
                                                                                       15
                                                                                                 512x1500
20
                                                                                       T
                                                      stats pooling
                                                                       [0,T)
                                                                                                1500Tx3000
    context = [1]
                                                                        \{0\}
                                                                                                 3000x512
                                                        segment6
22
    input dim = 512
                                                                                                  512x512
                                                                        \{0\}
    output dim = 1500
                                                        segment7
   net5 = TDNN (context, input dim, output dim, fi
                                                        softmax
                                                                        \{0\}
                                                                                                  512xN
```

So how exactly is TDNN implemented?

```
# Perform the convolution with relevant input frames
for c, i in enumerate(valid_steps):
    features = torch.index_select(x, 2, context+i)
    # print ('features taken:{}'.format(features))
    xs[:_\(\lambda:\lambda\)] = F.convld(features, kernel, bias_=_bias)[:_\(\lambda:\lambda\)]
return xs
```

How the stats pooling and the final layers looks like?

```
class StatsPooling(nn.Module):
    def __init__(self):
        super(StatsPooling_self).__init__()

def forward(self_varient_length_tensor):
    mean = varient_length_tensor.mean(dim=1)
    std = varient_length_tensor.std(dim=1)
    return torch.cat((mean_std)_dim=1)
```

```
class FullyConnected(nn.Module):
    def __init__(self):
        super(FullyConnected, self).__init__()
        self.hidden1 = nn.Linear(3000_512).double()
        self.hidden2 = nn.Linear(512_512).double()

def forward(self, x):
    x = F.relu(_self.hidden1(x))
    x = F.relu(_self.hidden2(x))
    return x
```

The Full network

```
SP = StatsPooling()
FC = FullyConnected()

Final = nn.Linear(512,4).double()

net = nn.Sequential(net1, net2, net3, net4, net5, SP, FC, Final)
```

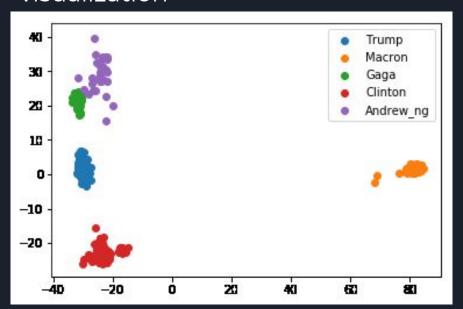
```
mel size: 24
                                              loss: 0.02339517055584004 step took 22.59490942955017
sequence length: 400
                                              output:
output shape: torch.Size([5, 396, 512])
                                              tensor([[ 6.6202, -6.6046, 0.9639, -0.9110],
mel size: 512
                                                     [0.8131, -3.1778, -2.1677, 4.8661],
sequence length: 396
                                                     [5.7950, -4.4012, -1.2745, -0.0496],
output shape: torch.Size([5, 392, 512])
                                                     [4.0205, -12.8199, 8.0885, 1.0369],
mel size: 512
                                                     [ 1.5110, -2.8568, -2.4932, 4.1227]], device='cuda:0',
sequence length: 392
                                                    dtype=torch.float64, grad fn=<AddmmBackward>)
output shape: torch.Size([5, 386, 512])
                                              labels:
mel size: 512
                                              tensor([[1, 0, 0, 0],
sequence length: 386
                                                     [0, 0, 0, 1],
output shape: torch.Size([5, 385, 512])
                                                     [1, 0, 0, 0],
mel size: 512
                                                     [0, 0, 1, 0],
sequence length: 385
                                                     [0, 0, 0, 1]], device='cuda:0', dtype=torch.uint8)
output shape: torch.Size([5, 384, 1500])
```

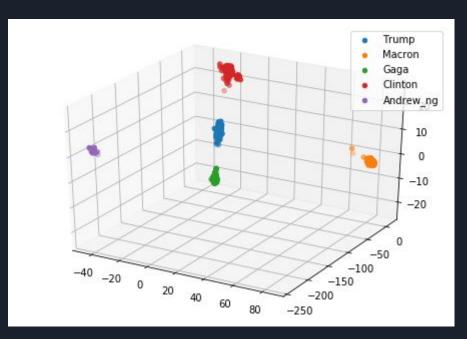
The **x-vector** embedding

```
1 net_xvec = nn.Sequential(net1, net2, net3, net4, net5, SP, FC)
```

net xvec.cuda()

LDA dimensionality reduction for visualization





https://github.com/Dannynis/xvector_pytorch

Questions?