

## Automatic Speech Recognition

Well, the task at hand is straightforward:

Take a waveform looking like this



This process is termed as "Speech-to-Text"

If map it to words being uttered → It is raining!

But if all this was, then we would not be here (still 9 pages to go)

The raw waveform as such does not tell us anything useful. Put differently, we need features that we can input to a model so that it can spit out what was being said.

But we have Amplitude-vs-time information, so why do we need something else and can't we use waveform as an input to a Deep learning architecture??

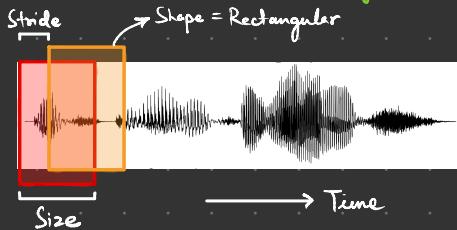
→ The 2 questions are separate and should not be confused. Let's address the 2<sup>nd</sup> one to begin with.  
A raw waveform depicts sound on a basic level, which is nothing but vibrations propagating through a medium. But to make a machine/model understand, we need to convert it from this physical notion to digital signal.

Now a sound disturbance carries with itself → frequency, amplitude, energy, etc. These are the properties of interest that we need to extract, which gets fed into a model.

The waveforms are just one-dimensional. What we need is 'Energy' information across various 'Frequency Bands' at a given 'Time'. Therefore mel spectrograms are what we use as an input to a speech recognition model.

## Converting sound into Mel Spectrogram

- First step is to digitize the signal  $\rightarrow$  Sampling & Quantization  
# Samples per second  $\rightarrow$  Storing Amplitude as 8 bit / 16 bit integers
- Next we do Windowing  $\rightarrow$  the part of signal that is responsible for a particular fundamental phoneme will be extracted through windowing.



$$\rightarrow \text{The value post windowing} \Rightarrow w(t) = f(t) \cdot s(t)$$

↓  
window function

$$e^{i\theta} = \cos \theta + i \sin \theta$$

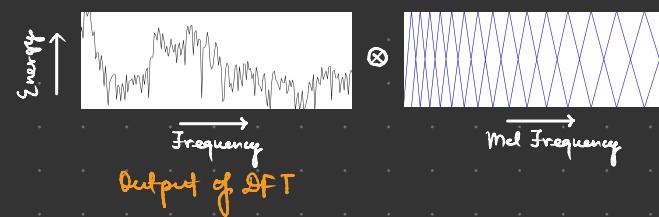
$\rightarrow$  Euler's Formula

- To extract energy from the different frequency bands we need Discrete Fourier Transform

$$\text{The magnitude & phase of a } k^{\text{th}} \text{ Frequency component is given by} \rightarrow E(k) = \sum_{t=0}^N w(t) \exp(-i \frac{2\pi}{N} kt)$$

$\rightarrow$  Total number of Frequency Bands

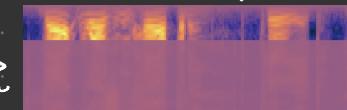
- Finally to aid human recognition, we need to increase the weight of lower frequencies.  
So we transform frequency to Mel scale that lets us collect energies on a 'log scale'



$$\text{Mel Scale} \rightarrow \text{mel}(k) = 1127 \ln(1 + f/700)$$



Feature Extraction



Phoneme Classification  
(J)

## Working with Mel Spectograms

- The huge data contains
  - Utterances
  - Phoneme state labels
- or  $\Rightarrow$  Mel Spectral vectors

Each utterance contains in its first dimension # Frames while # Frequency bands in the second dimension

The individual frame is of size  $x \times m$  and any adjacent frames will be  $y \times m$  apart.

Frame size

Frame stride

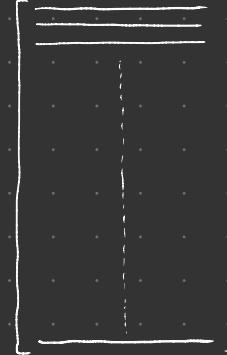
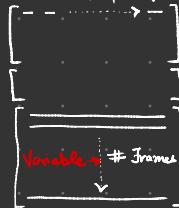
Against each utterance there is a phoneme states vector, which contains a phoneme state for every frame in that utterance. Entries in the array correspond to Energy values (log scale) in different frequency bands that vary across time.

The job is to build a model that learns for every frame a phoneme state out of the 138 classes.

- First naive model will be that takes each frame independently & learns the phoneme label against it. We'll test both MLP & ConvID against this structure of data

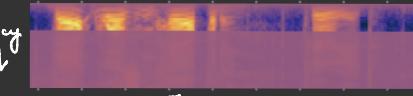


Fixed  $\rightarrow$  # frequency bands



An Utterance

Frequency  
 $\downarrow$



Frames  $\rightarrow$

Now given we are working with Temporal Data, taking each frame on its own will only take us to a best  $\sim 25\%$  on the accuracy scale. This learning is worse than a model giving random output ( $\sim 50\%$  accuracy).

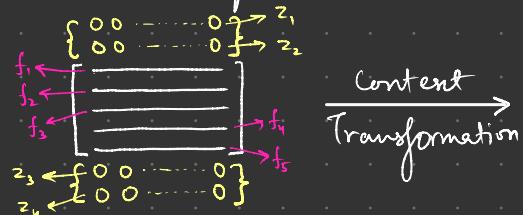
The flaw in this plan is that the model only sees at the current frame to learn the current phoneme class. What if our model got access to energy values in the frequency bands of previous & next frames with respect to the current frame and then learn the phoneme label {given that we need to work with MLP & Conv only} !!

So we concatenate the frames from previous & next time steps to the current frame. The thing to take care of is to not leak this context between utterances. Therefore, we need to do this step of augmentation at an utterance by utterance basis.

Let's say we choose `context size=2`, then the first 2 & last 2 frames of each utterance won't have enough features post data augmentation.

To solve this issue, we pad each utterance with an array, of constant = 0, of shape (`context size, # frequencies`) on each side & then perform the transformation.

Let's take an example utterance



Input to model

$$\begin{bmatrix} z_1 \oplus z_2 \oplus f_1 \oplus f_2 \oplus f_3 \\ z_2 \oplus f_1 \oplus f_2 \oplus f_3 \oplus f_4 \\ f_1 \oplus f_2 \oplus f_3 \oplus f_4 \oplus f_5 \\ f_2 \oplus f_3 \oplus f_4 \oplus f_5 \oplus z_3 \\ f_3 \oplus f_4 \oplus f_5 \oplus z_4 \oplus z_1 \end{bmatrix}$$

Output Phoneme states

$$\begin{bmatrix} - \\ - \\ - \\ - \\ - \end{bmatrix}$$

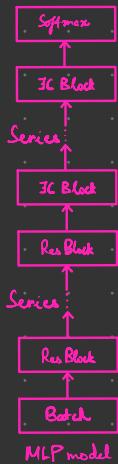
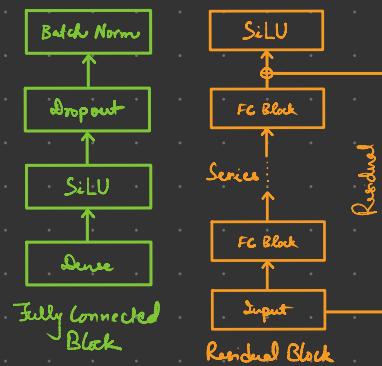
We can apply MLP architecture on this type of data

Number of Rows remain same

but the size of each row increases to  $\# \text{Features} * (2 * \text{content size} + 1)$

Another issue that we face is if the total number of frames of all utterances combined is already large then creating new features by increasing context size will increase the size of data beyond a normal system can handle. Therefore we will use TensorFlow's generator to create batches on the fly which will be called only when we do `model.fit`, so the whole data won't be loaded into the RAM while training. We'll also modify the iterator to augment the data based on the `content size` as a hyper-parameter.

## Dense Residual Architecture



- Fully Connected Block is made of 'Dense', 'Activation', 'Dropout' & 'Batch Norm' Layers
- Residual Block offers a 'skip connection' by joining the output of last JC Block & input to pass as the argument of next Res Block post applying a non-linear activation
- Collection of the 'Res Blocks' & then 'JC Blocks' followed by 'Softmax Layer' makes a deep MLP

can be 'concatenation' or 'summation'

SiLU is short for 'Sigmoid-weighted Linear Unit'  $\text{silu}(x, \beta) = x \cdot \text{sigmoid}(\beta x)$

↳ like a smoothed version of ReLU

→ shape = (batch\_size, [2 \* context size + 1] \* # frequencies)

$$= \frac{x}{1 + \exp(-\beta x)}$$



What is Dropout ??

→ Consider an image of a human face & car. The task is to train the model to detect a face.



Imagine 3 units of hidden layer & parameters randomly initialized



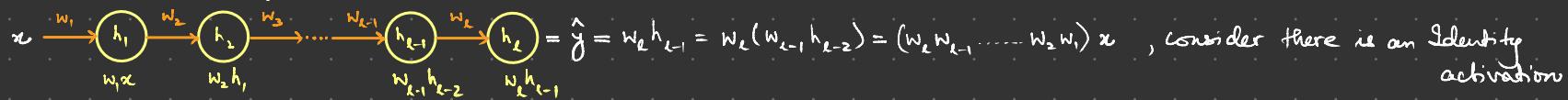
Each unit learns a feature about the input image. But more importantly, let's say the 3rd unit learns about the mouth in a face.

But what if there was a test image with the mouth not visible.

The task Dropout does is it will mask the 3rd unit and force the model to think that it has not learnt about a mouth. So in the next iteration, the model makes some other unit learn about mouth or a new feature to detect face, like nose or ear. The model is then made to learn new features or call it regularization.

What is Batch Normalization ??

To set the record straight BN is not a regularization method. It is in fact an adaptive reparametrization technique to improve the process of optimisation in very deep nets. Lets understand this through a toy example



After one step of Backpropagation  $\rightarrow \hat{y} = (w_l - \alpha g_l) \cdot (w_{l-1} - \alpha g_{l-1}) \dots (w_2 - \alpha g_2) (w_1 - \alpha g_1) x \rightarrow$  each weight follows an update

The prediction is of the form  $\rightarrow \hat{y} \sim (\alpha g^{W^{l-1}} + \alpha^2 g^2 W^{l-2} + \dots + \alpha^l g^l W^0) x + \underbrace{w^l x}_{w_i' = w_i - \alpha g_i}$ ,

Assume the gradient to be 1  $\rightarrow$  fixed by the cost function Before BP where  $g_i = \text{gradient}$

Now In an optimization process, we want the difference to be small  $\Rightarrow$  Small  $\Delta \hat{y}$

$\Delta \hat{y} \sim \alpha w^{l-1} + \alpha^2 w^{l-2} + \dots + \alpha^l w^0 \rightarrow$  The weights of each layer when multiplied can either be very small or very large

Therefore the choice of learning rate becomes difficult to get stable updates in a deep network

Let's make few changes to the network

what BatchNorm does is that it normalizes the hidden layer at first. In doing we are making units standardized

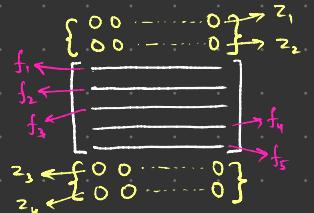
This helps reduce the problem of simultaneous updates across layers of the deep network. But in doing the network becomes vanilla. To restore the balance, we let the layers learn its own mean & standard deviation



$\tilde{\mu} = (\mu, \sigma, \lambda, \beta)$   
Fixed for learnable every hidden layer

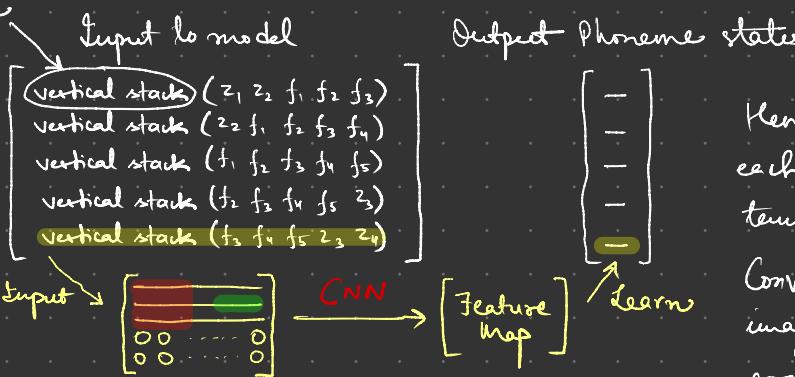
$$BN = \begin{cases} h_i' = \frac{h_i - \mu}{\sigma} \\ h_i'' = \lambda h_i' + \beta \end{cases}$$

Instead of flattening post transformation



Content  
Transformation

Input to model

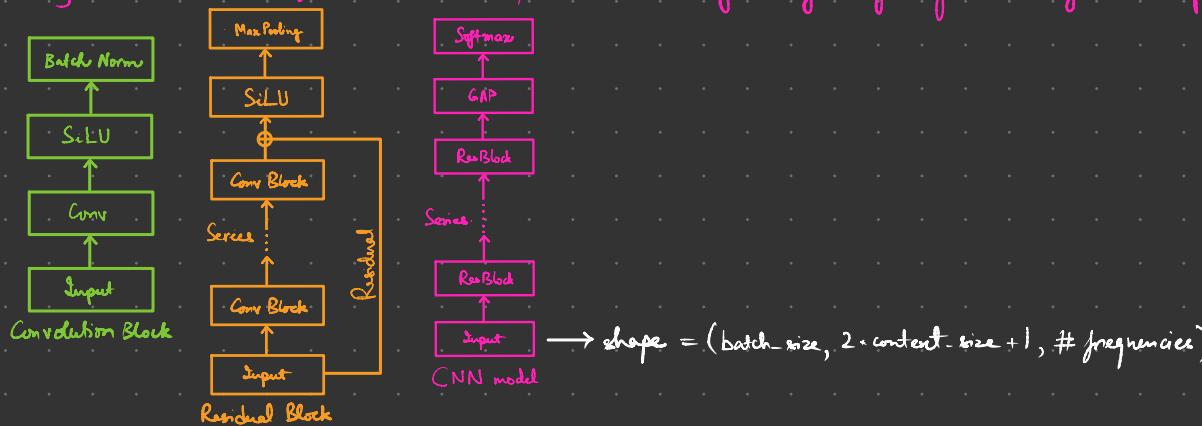


Here we can consider each stack as collection of temporal vectors and apply ConvID or see it as an image and see that how good Conv2D performs

What the custom CNN model looks like??

The architecture will be built using 'Sub-classing' & 'Functional' approaches

- There will be an elementary 'Convolution block', comprising of a 'Conv' layer, a 'Batch Normalization' layer & a 'SiLU' layer
- Using this a recurrent part, we will build a 'Residual block' containing series of {Convolution blocks & a 'Skip Connection'} followed by a 'MaxPooling' layer.
- Lastly, we combine series of Residual blocks by a 'Global Average Pooling' layer followed by the output 'Dense' layer



$\rightarrow \text{shape} = (\text{batch\_size}, 2 \cdot \text{content\_size} + 1, \# \text{ frequencies})$