

Speech Processing with Neural Networks

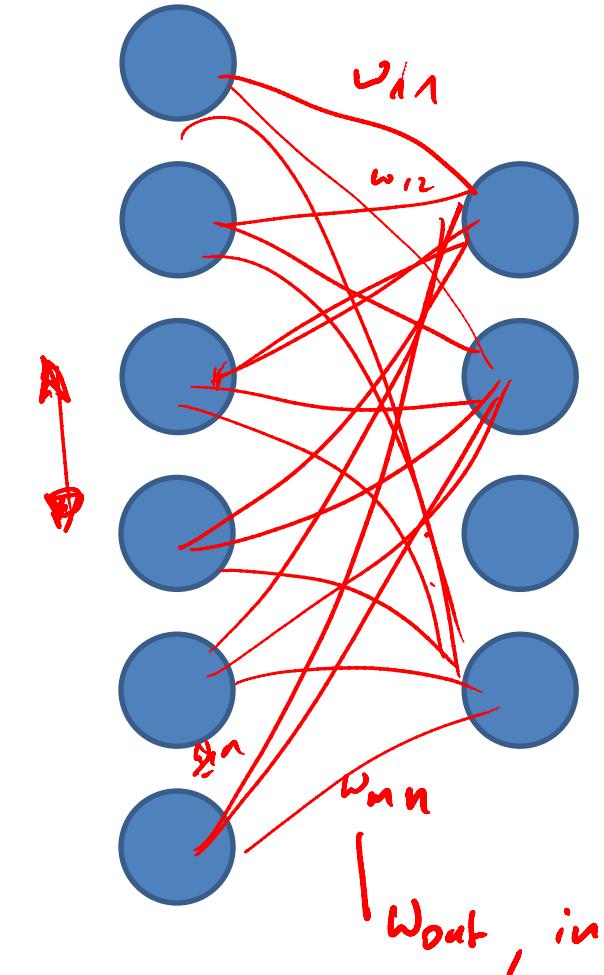
Jan Chorowski

RECAP: DIFFERENT NEURON TYPES

Different Layers for Different Problems

Dense layers:

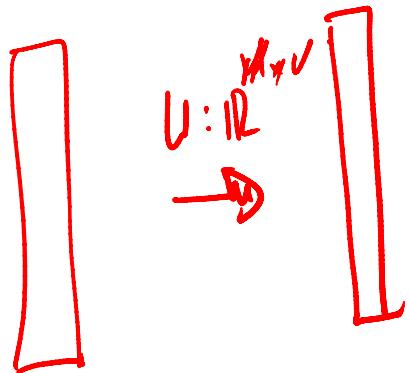
- Global connectivity:
All outputs use all inputs
- No weight sharing
 $O(N^2)$ weights
- Trains the same when
inputs/outputs permuted
- Any combination of input/output
dimension possible



Dense Layers

Dense!

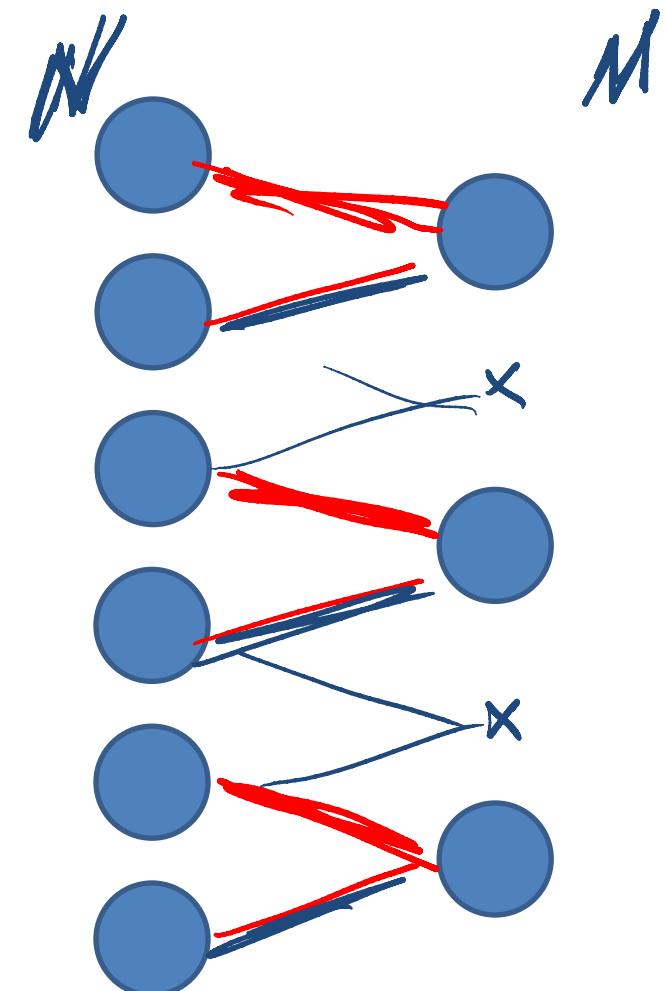
$$\mathbb{R}^N \rightarrow \mathbb{R}^M$$



Different Layers for Different Problems

Convolutional layers

- Local connectivity – each output depends on small set of inputs
- Weight sharing – number of weights independent of input size
- Depends on sensible input/output ordering (e.g. can't shuffle pixels in an image)
- Number of output proportional to number of inputs (give-or-take padding)

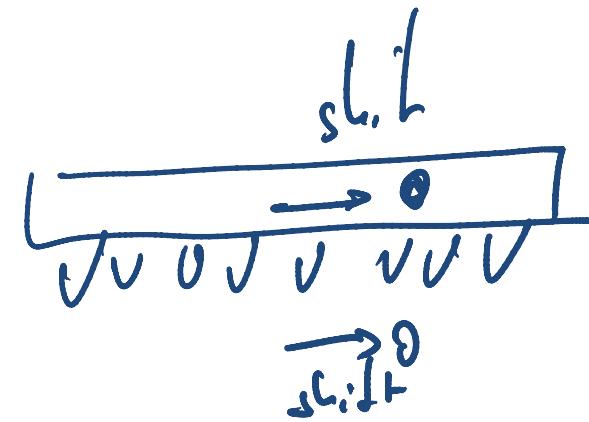
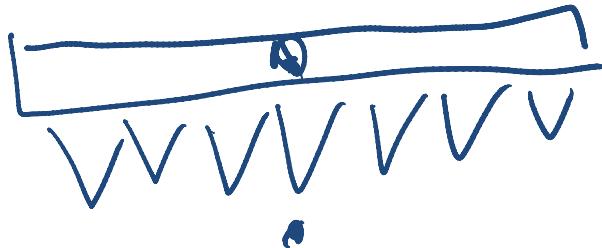


Conv layers

1D

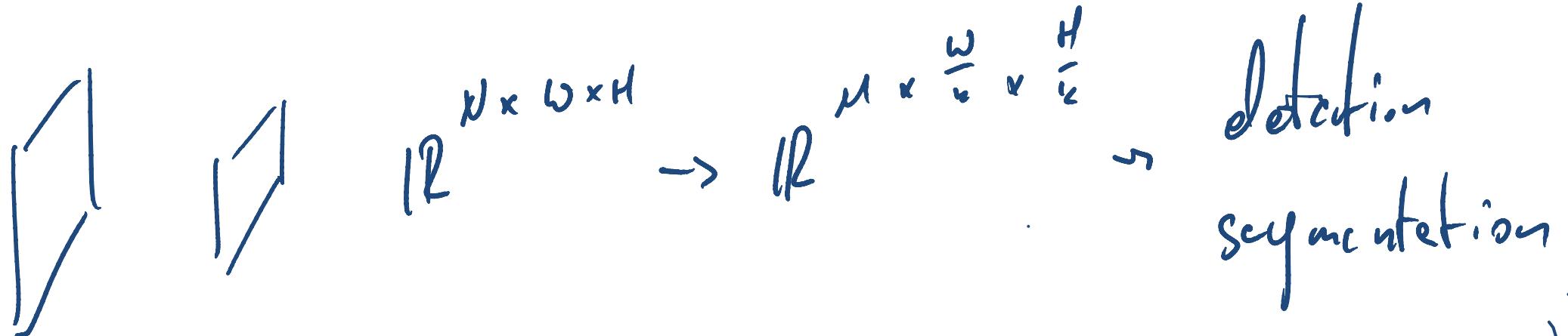
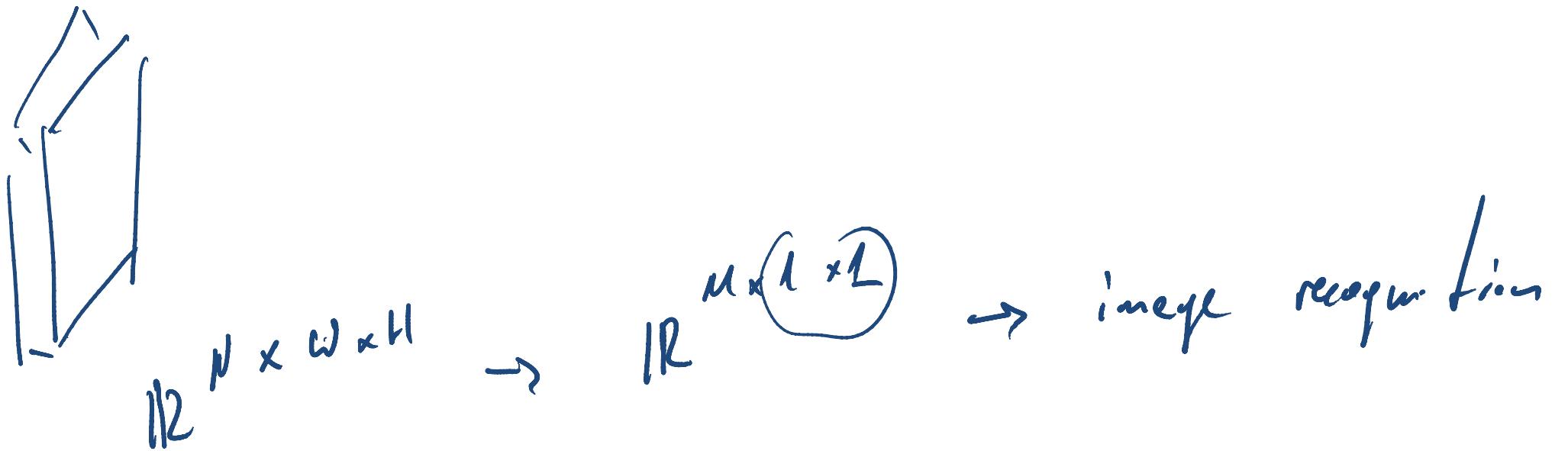
input $\mathbb{R}^{N \times T}$
 ↑
 decay
 time axis

out $\mathbb{R}^{M \times \frac{T}{k}}$
 ↑
 k



$$k : nL$$

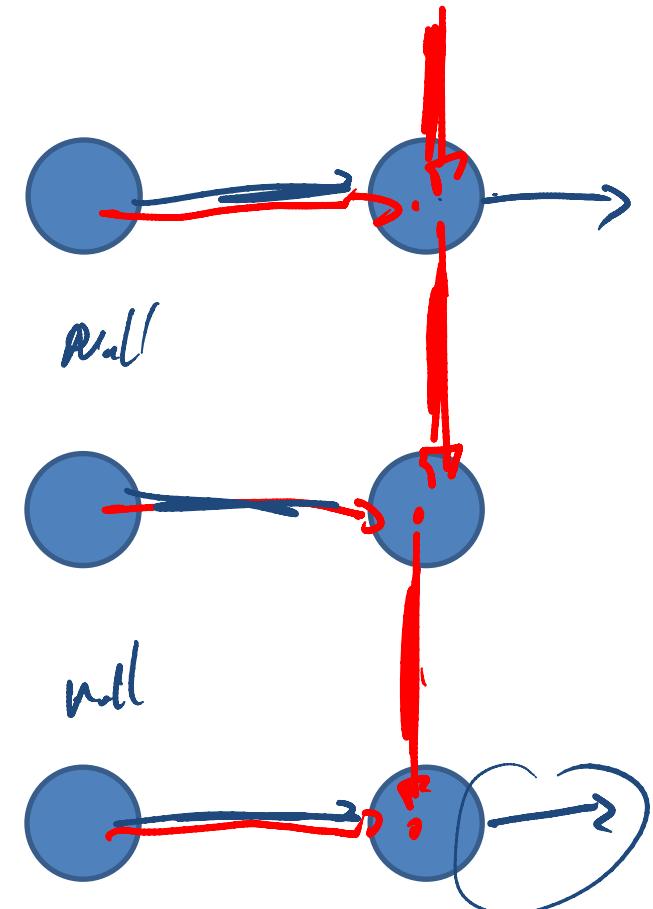
Conv nets



Different Layers for Different Problems

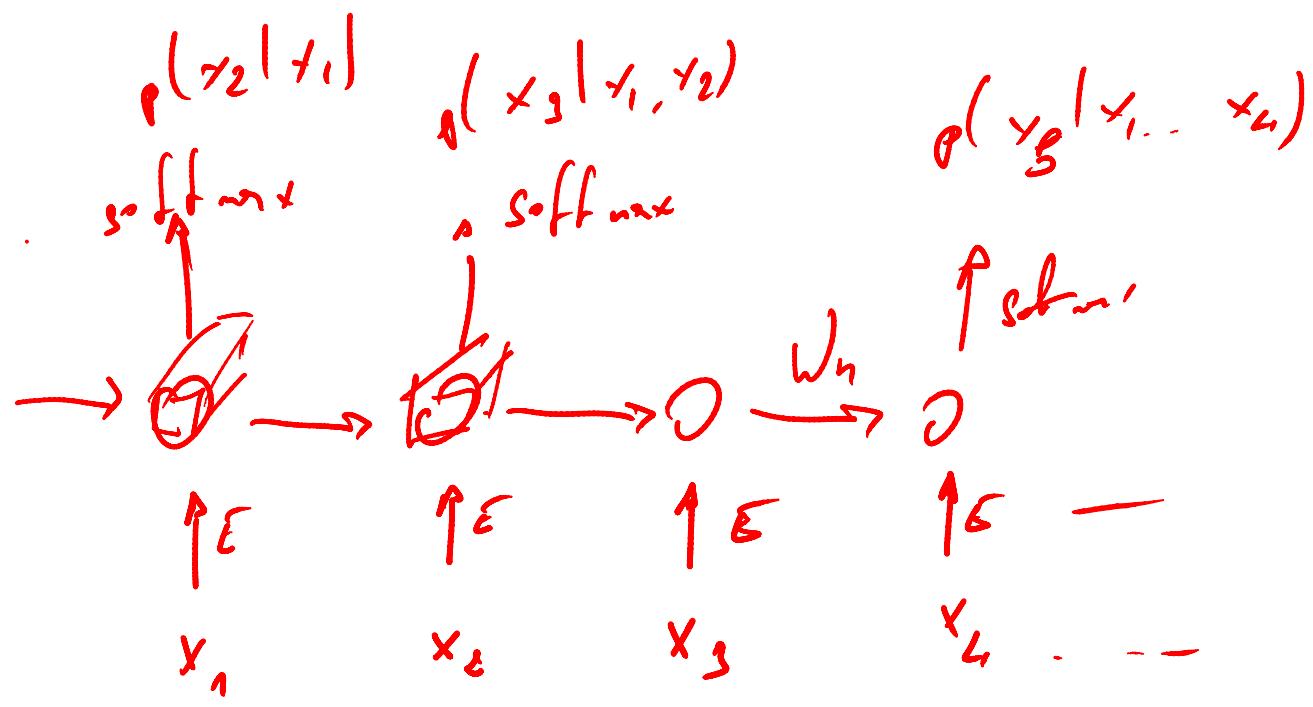
Recurrent Layers

- Local connectivity
- Weight sharing
all steps use same function
- Order dependent
sequential computation
- Number of output proportional to
number of inputs
or forces only one input or one output



Ep. Recurrent Language Model

Task: predict $p(y_n | x_1, \dots, x_m)$

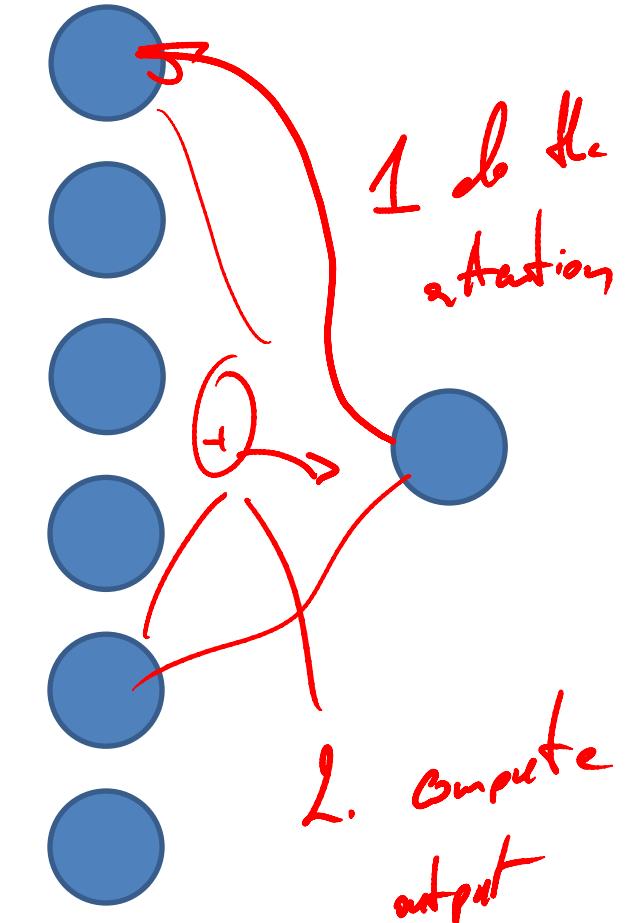


W_h - recurrent weights
(e.g. for an LSTM cell)
 E - word embed
matrix (same for
all timesteps)

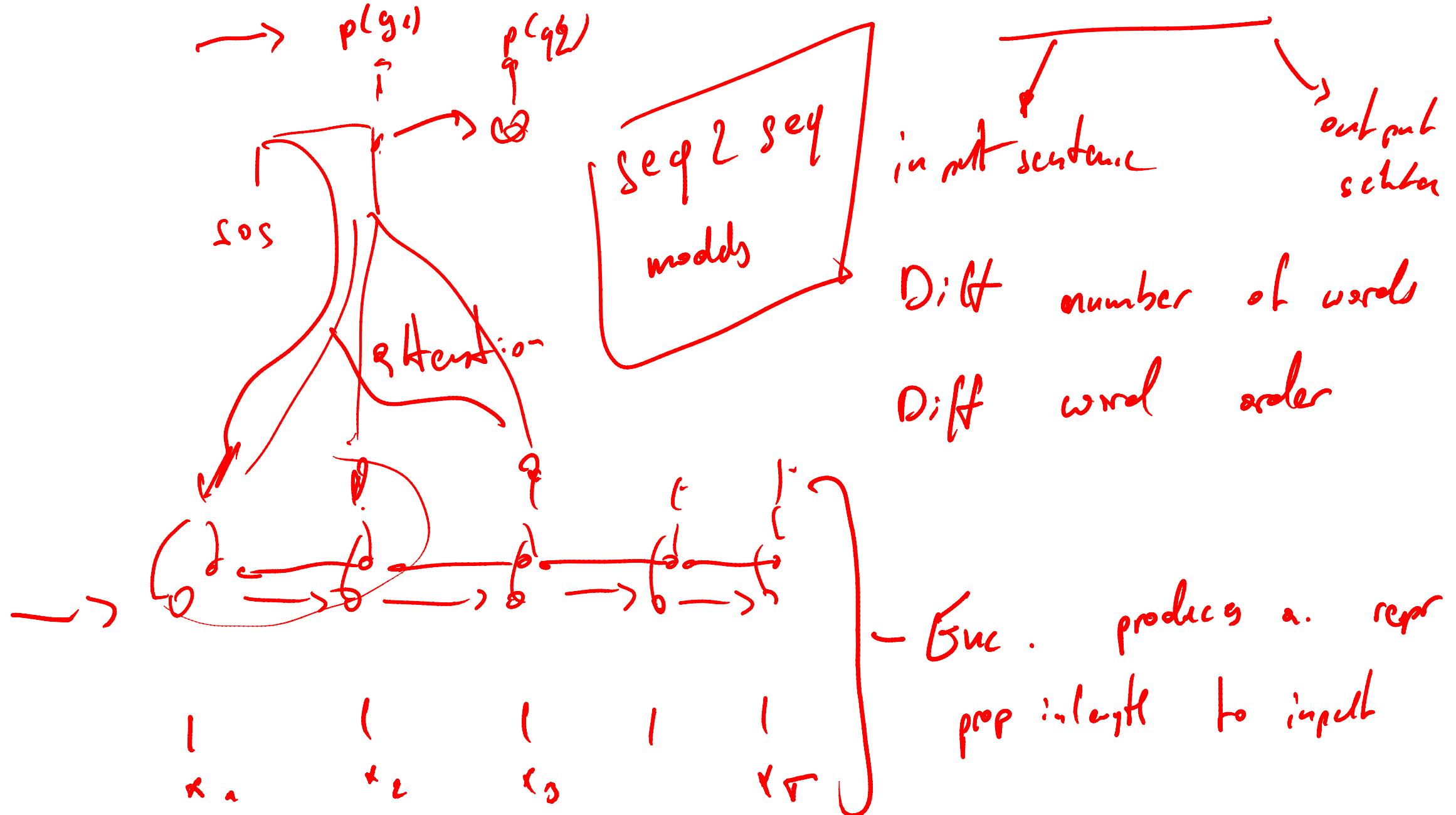
Different Layers for Different Problems

Attention neuron

- Global connectivity:
The neuron selects some or all input units
- Works the same when inputs permuted
- Works for dynamically changing number of inputs



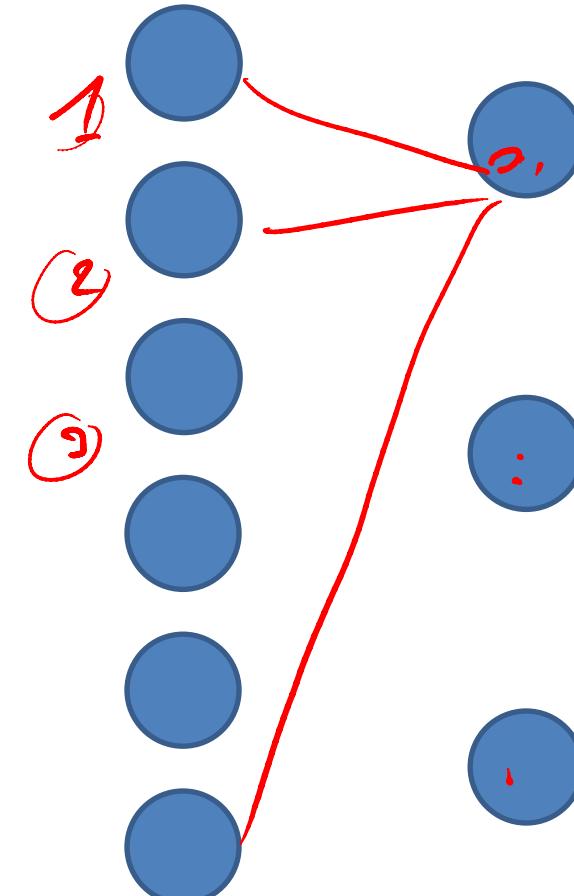
Recurrent Attention-Based Translation



Different Layers for Different Problems

Transformer (self-attention) layers

- Global connectivity – each output accesses all inputs using attention
- Weight sharing – all outputs use the same weights (but choose different subsets of inputs)
- Optionally can depend on input/output ordering (can use positional embeddings)
- Number of outputs independent from number of inputs



Transformer (self-attention) Translation

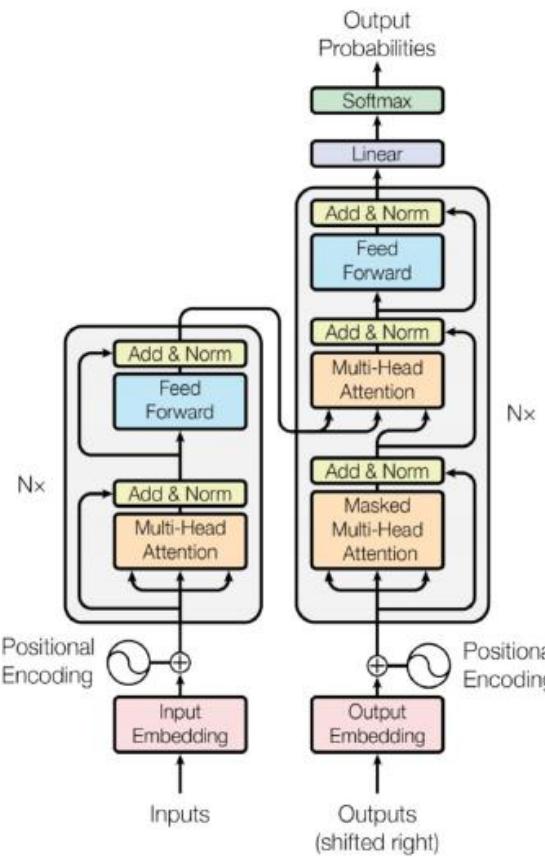
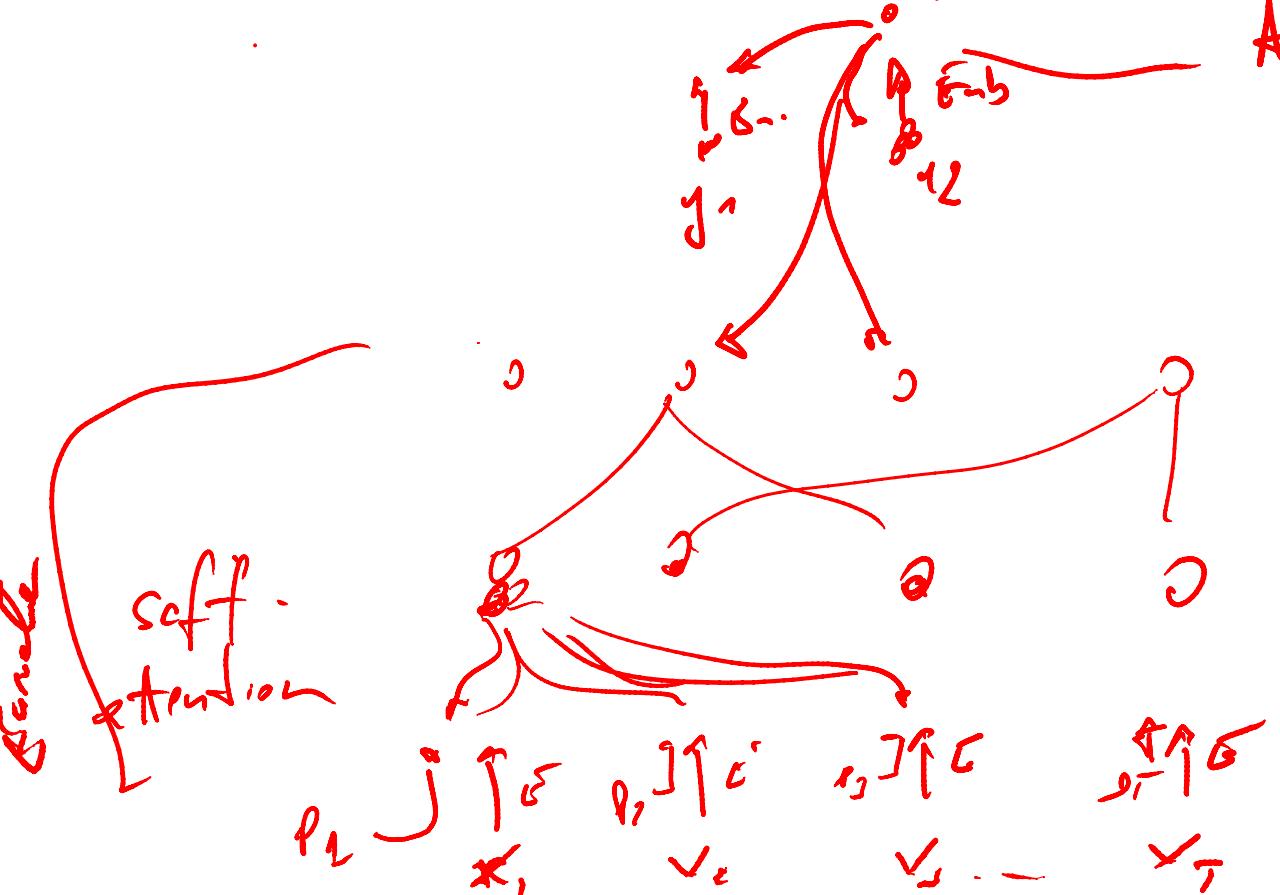


Figure 1: The Transformer - model architecture.

~~Transformer~~ (self-attention) Translation

input : word sequence

output : word sequence



Attention over : all encoder's states
+ all PAST decoder states

self-attention

self-attention : same weights
for all neurons, parallel compute

E_i Embed matrix, same for all words

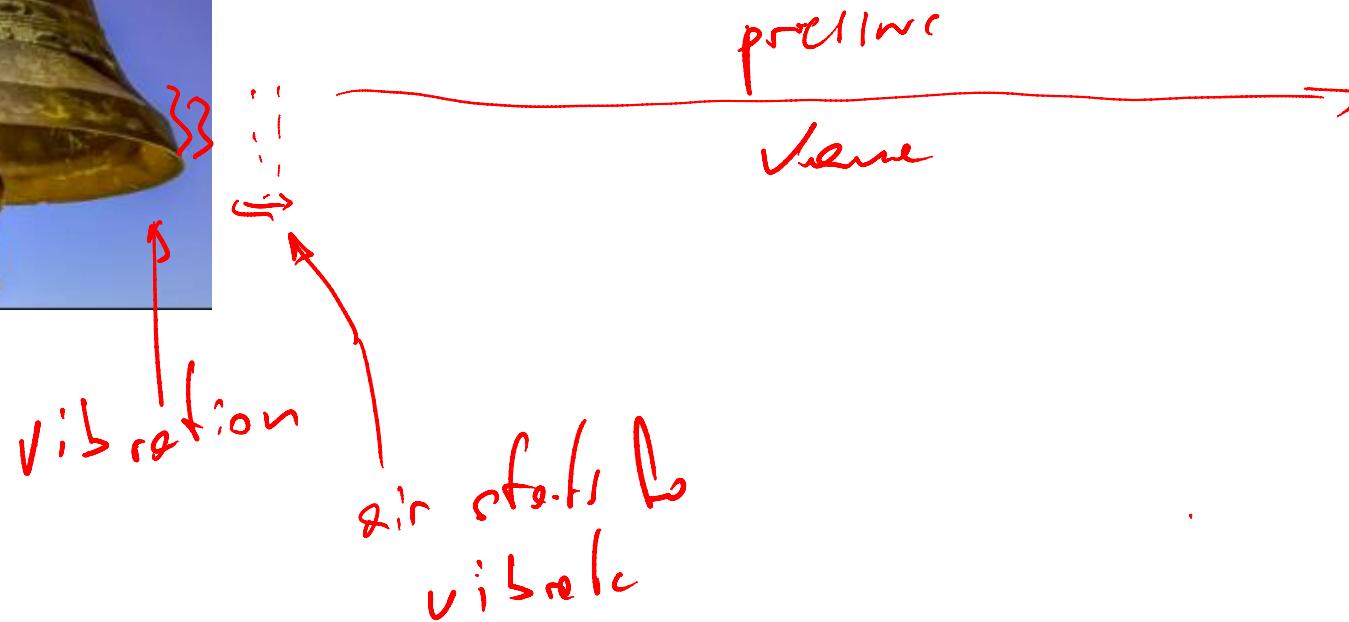
Architecture summary

Neural architectures exploit problem domain properties:

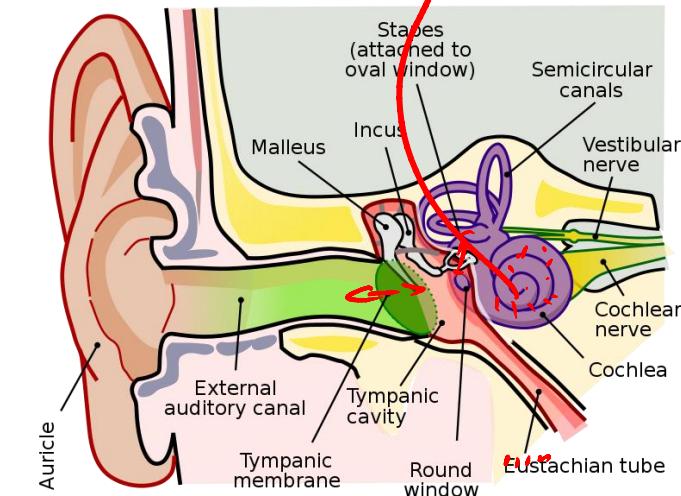
- Local vs global connectivity, match data layout
1D (words, audio), 2D (images), 3D (movies, tomography)
- Weight sharing:
 - All image location use same set of filters
 - All words use same embedding matrix

SPEECH AND SOUND PROCESSING

Sounds

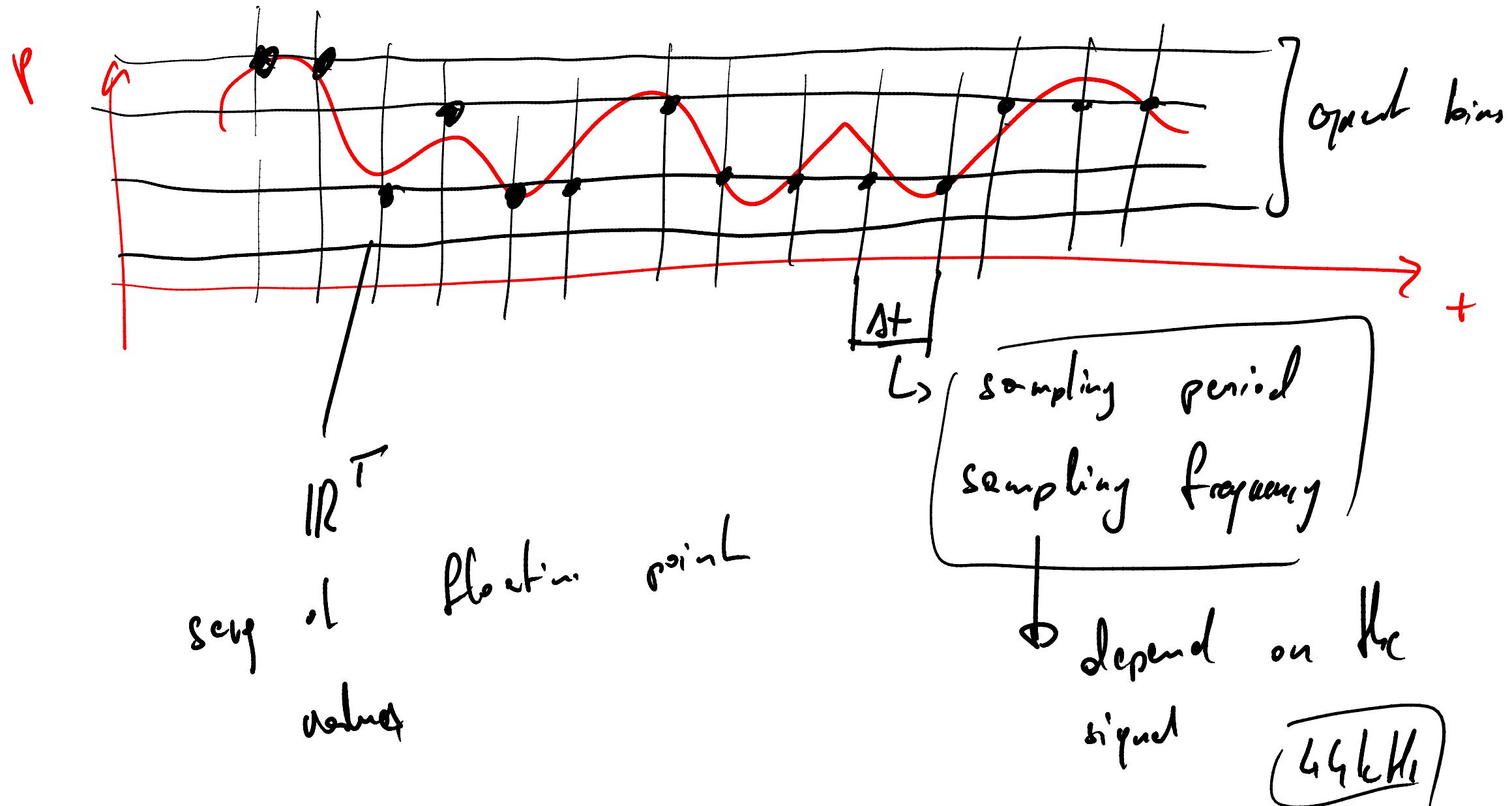


we feel the motion
of these tiny hairs



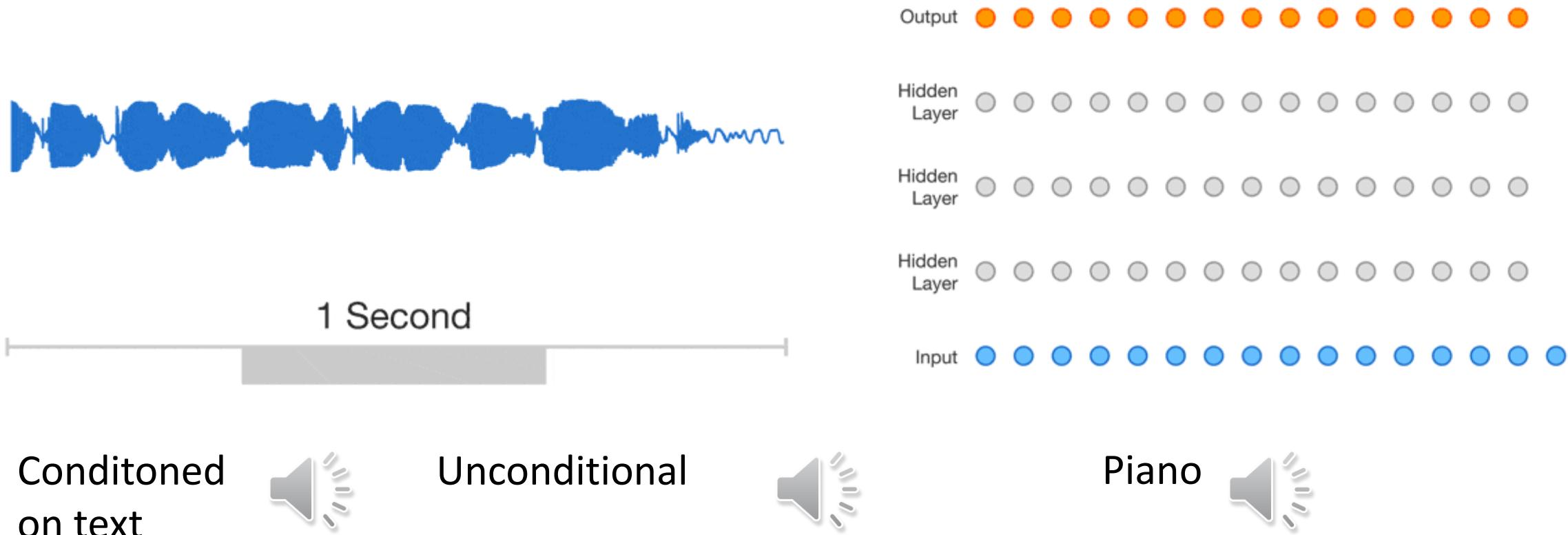
modulating
an electrical
property

Digital sound: discrete in time and magnitude



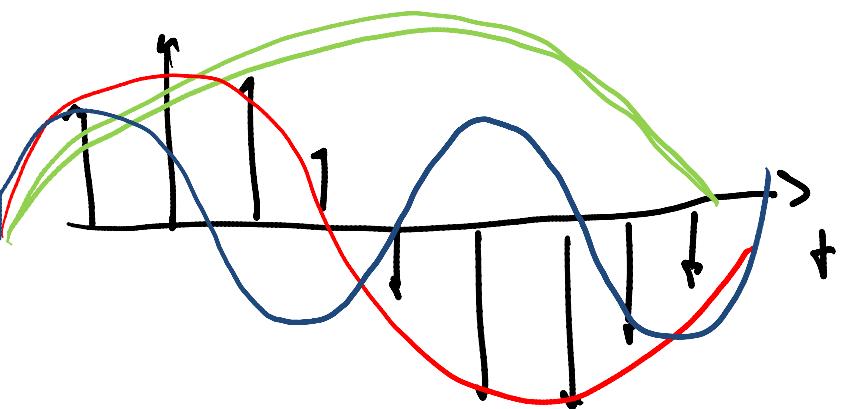
Nnets for direct sound sample processing

WaveNet: a convnet for raw audio



Fourier Transform: Time \leftrightarrow Frequency

Sampled signal



sequence of flats

correlate with different frequencies

 Grado A empiricidade

Fourier

A hand-drawn diagram illustrating a complex number in the complex plane. The horizontal axis is labeled F and the vertical axis is unlabeled. A vector originates from the origin. A green tick mark on the vertical axis is labeled 1 . A blue tick mark on the horizontal axis is labeled 1 . A red tick mark on the vector is labeled α . A black arrow points along the vector, starting from the origin. A curved arrow at the top left points towards the vector, labeled "Ampplitude". A curved arrow at the top right points towards the angle between the vector and the positive F -axis, labeled "phase". To the right of the angle label, the text "- complex number" is written.

↑ we will only look at
amplitude and drop the phase

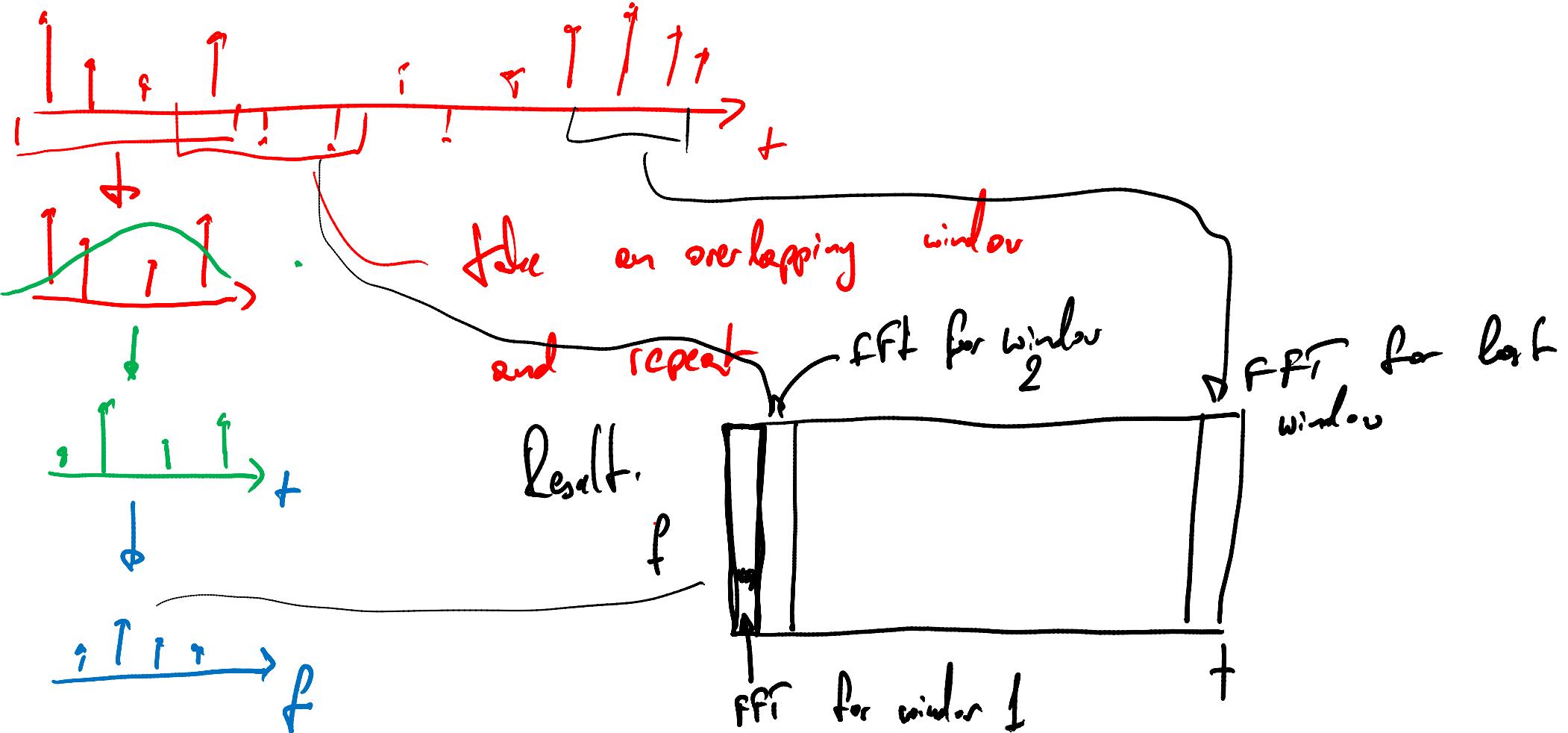
Short Time Fourier Transform STFT

Intuition: time signal \leftrightarrow time x frequency representation

Input

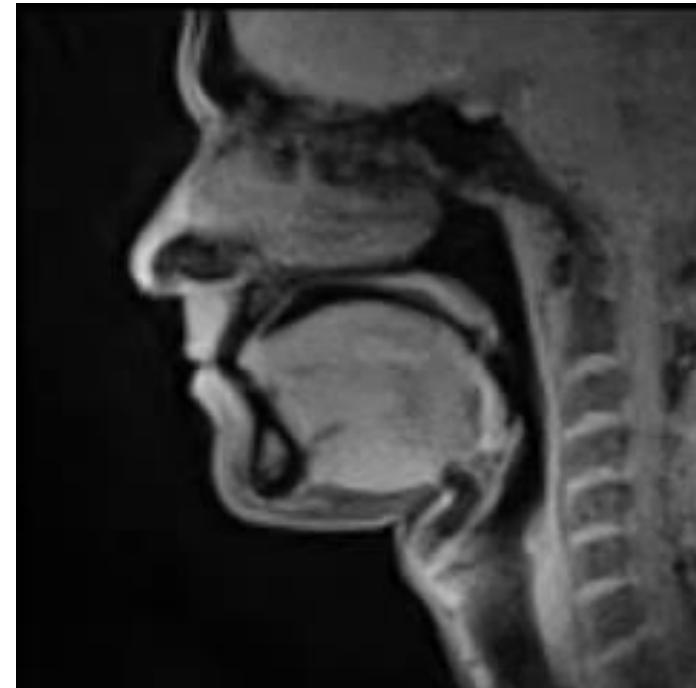
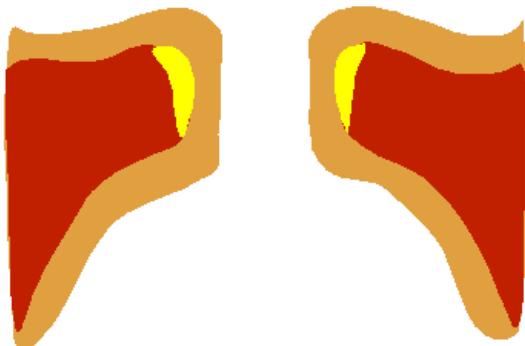
1: take a
"window",
scale down
on edges

Fourier
Transform

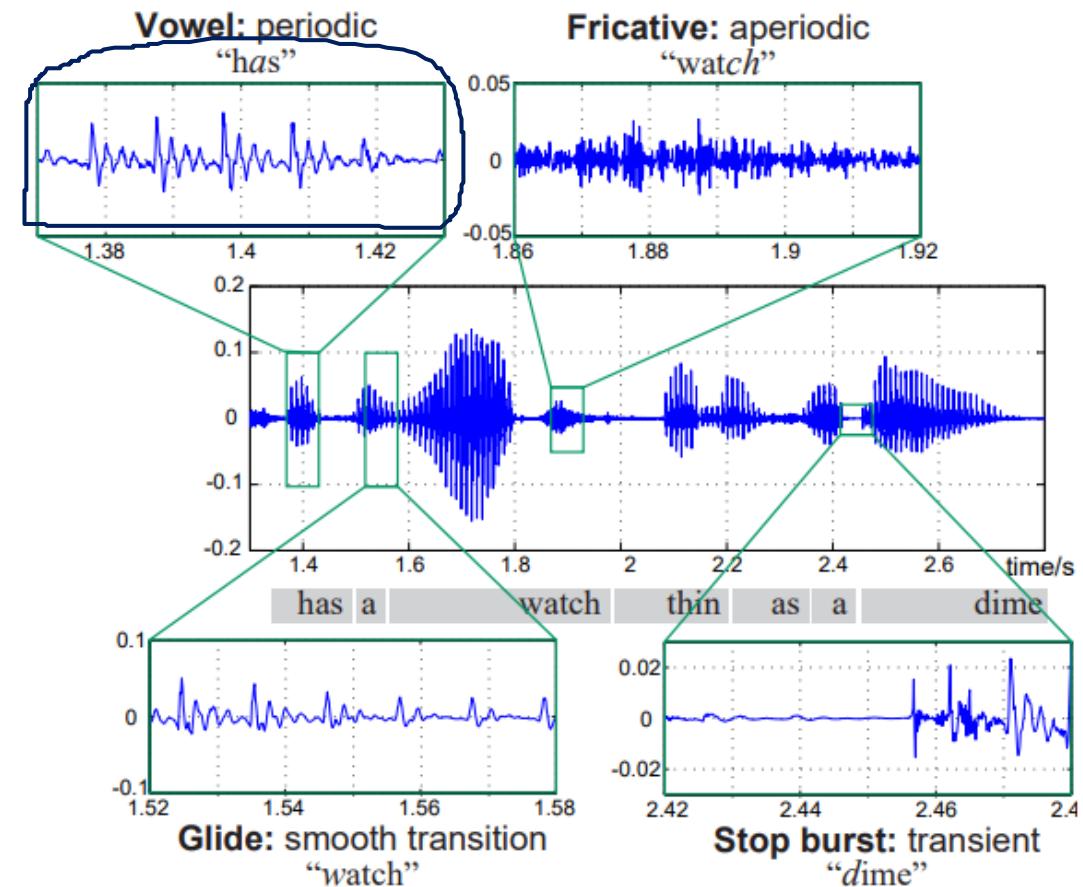


Speech

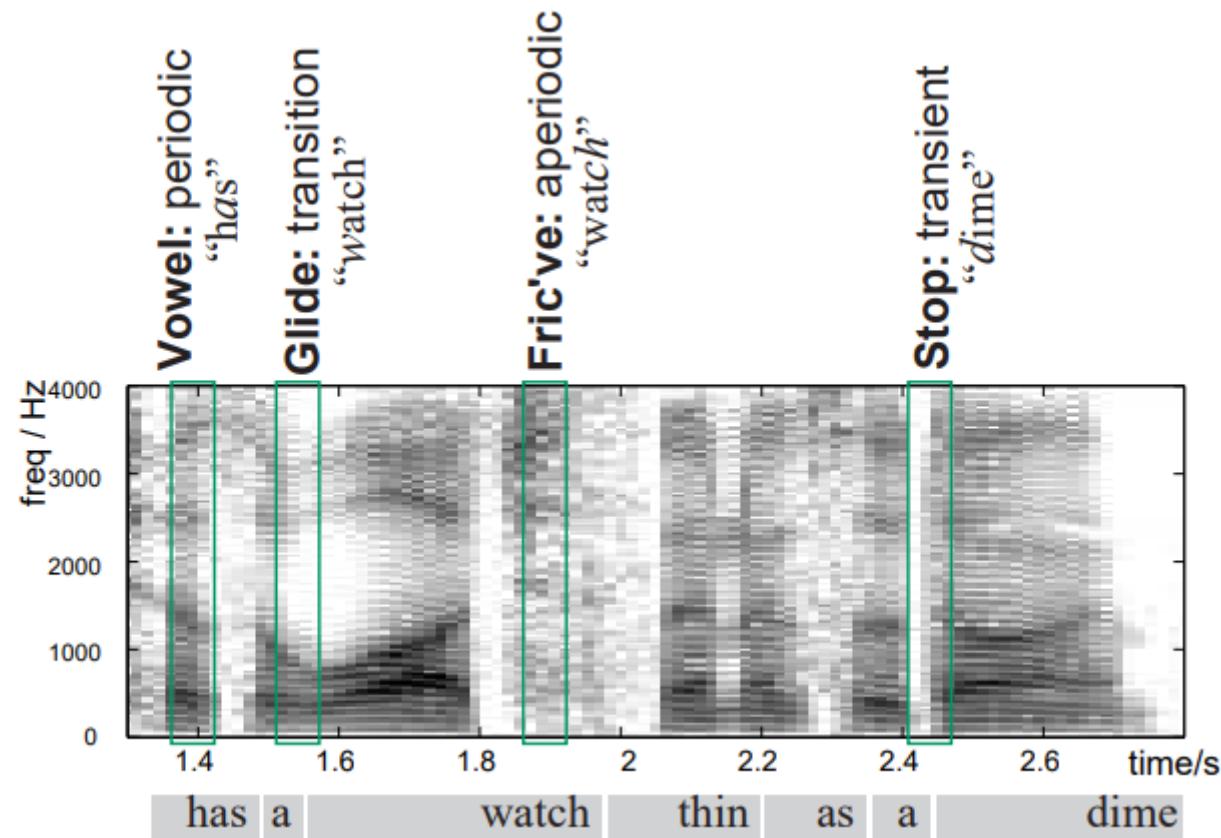
- Vocal cords vibrate (think of saying aaaa)
- The throat, mouth and tongue filter out some frequencies,



Speech time domain signal



Speech STFT



SPEECH RECOGNITION

Speech recognition: prepare a transcription

“Speech to text”:

find a word sequence that matches the recorded utterance

Error measure: Word Error Rate (WER)

Reference: It is a sunny day

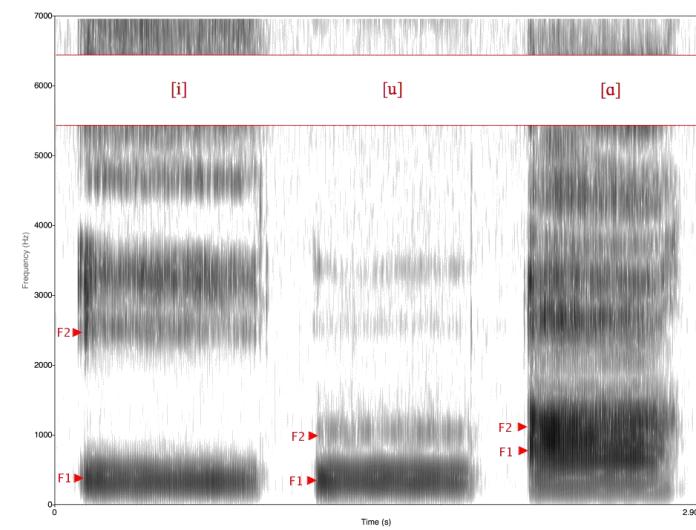
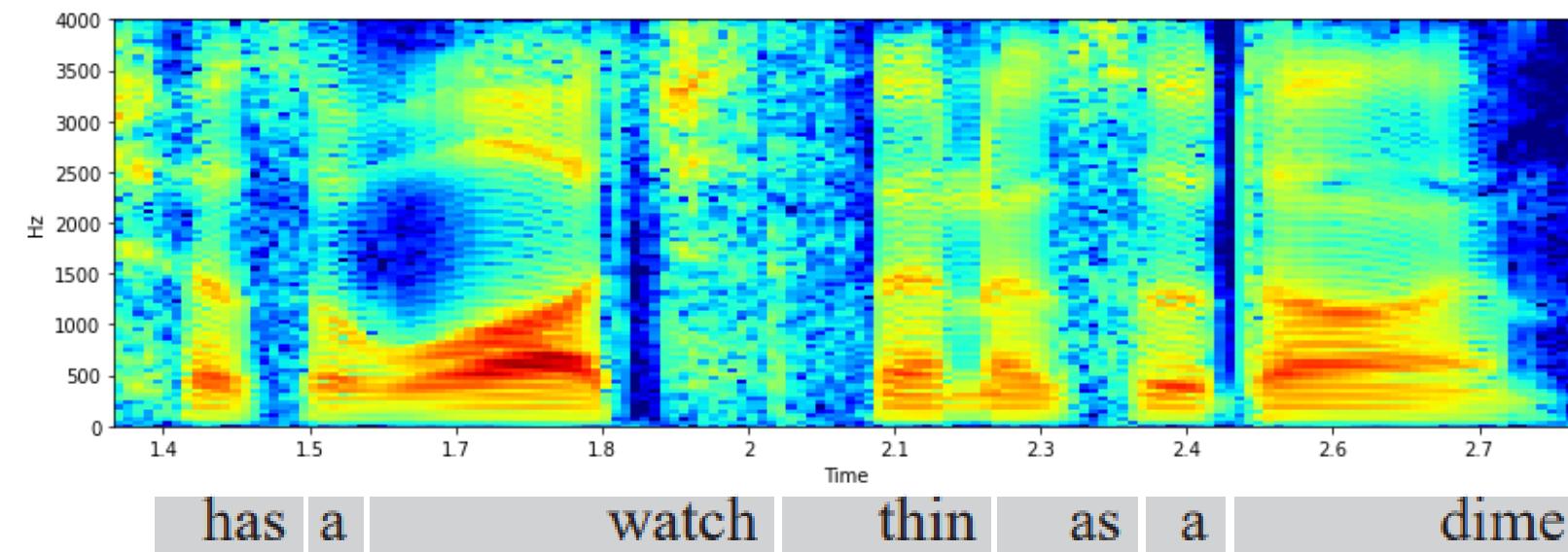
Recognized: It was sunny all day

Word status: OK S D OK | OK

$$\text{WER} = (\text{S} + \text{D} + \text{I}) / \text{len}(\text{ref})$$

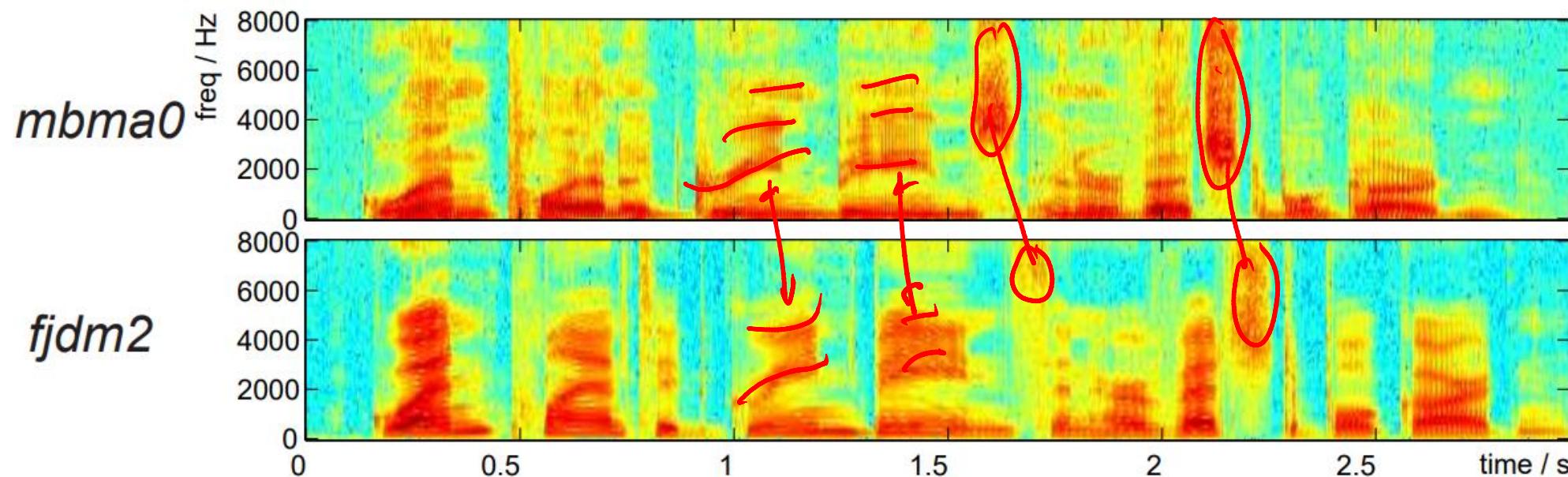
What the net should discard/recognize?

- Discard: vocal cord frequency (f0)
- Keep: formant frequency (set position of the mouth, tongue...)

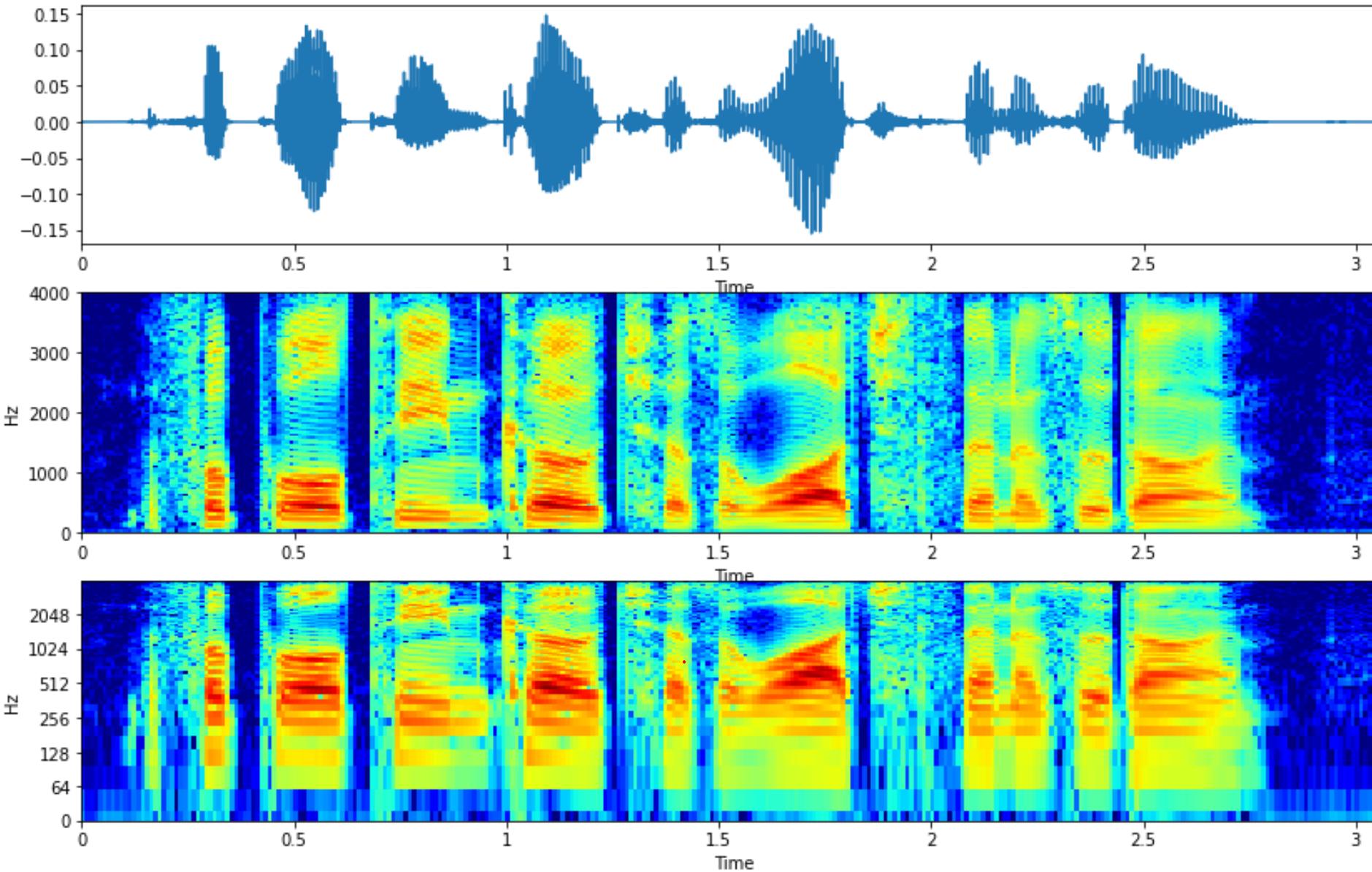


Speech Differences Between Individuals

- Differences between individuals:
 - f0 (base vocal cord frequency)
 - Speed, prosody
 - accents



Speech Features for Neural Networks



- raw

STFT

Log cepstrum

Frey expts

MEL scale

Speech Features: Intuitions

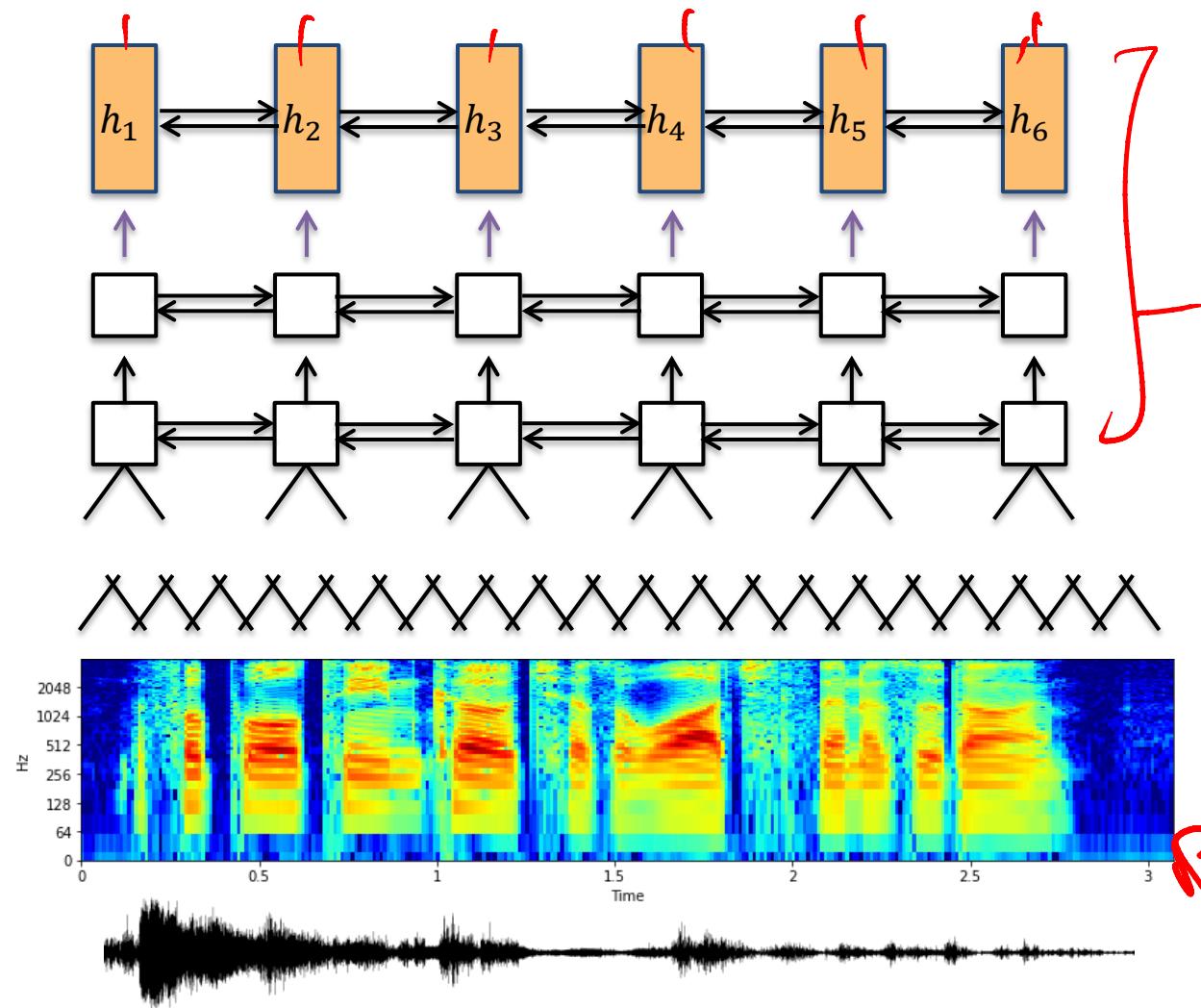
An utterance is represented as a $Time \times Frequency$ array of amplitudes.

Length along time axis varies between utterances.

The number of frequencies is fixed:

- Treat as a 1D signal with $|F|$ channels, apply 1D convolutions.
- Treat as a 2D signal (image) and apply 2D convolutions.

A typical speech processing architecture



$(Bi)RNNs$

2D convs, per-freq bias

Rets: NBL blk h

Speech – Time Alignment

Need alignments between:

- transcript (sequence of words, or phonemes, or graphemes)
- audio (sequence of frames, sampled uniformly in time)

Some possible solutions:

- Attention, treat as sequence-to-sequence task
- Design a loss function that handles the alignment (CTC)