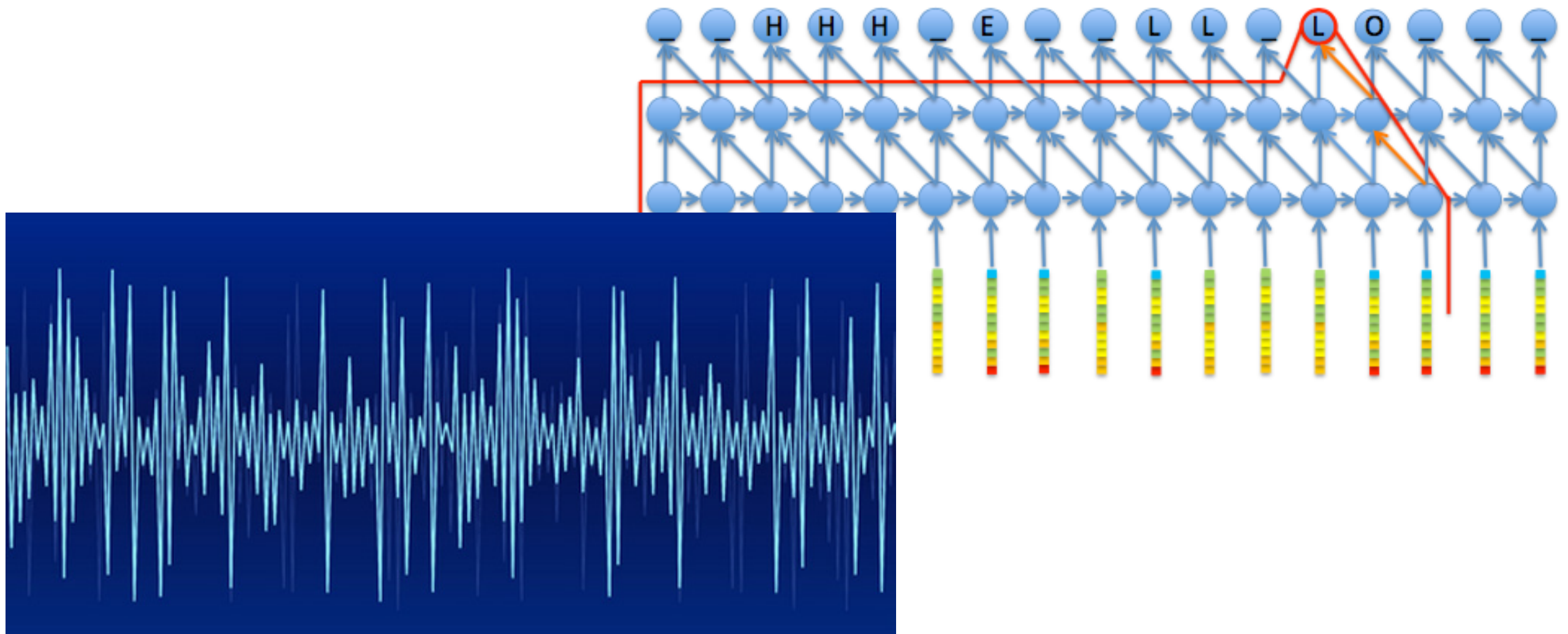


Speech Recognition with Deep Learning

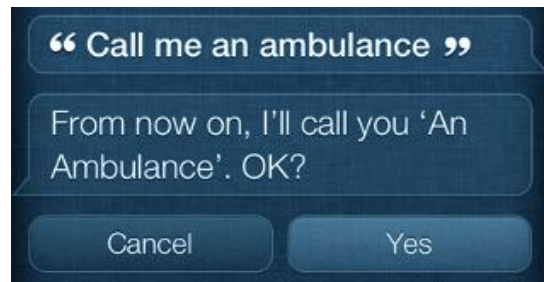
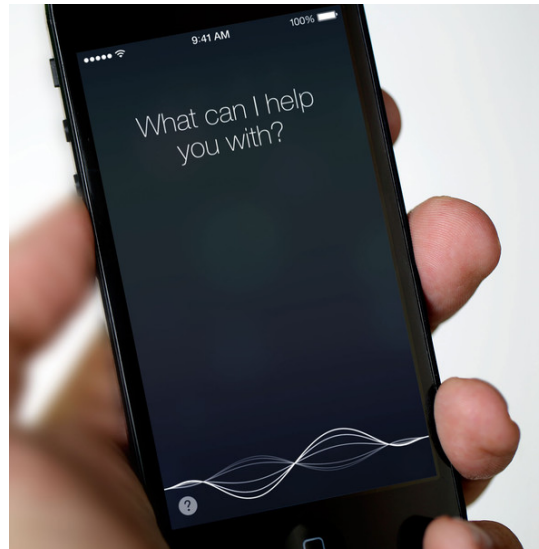
Patrick Foldenauer

based on lecture by Adam Coates at Bay Area Deep Learning
School 2016

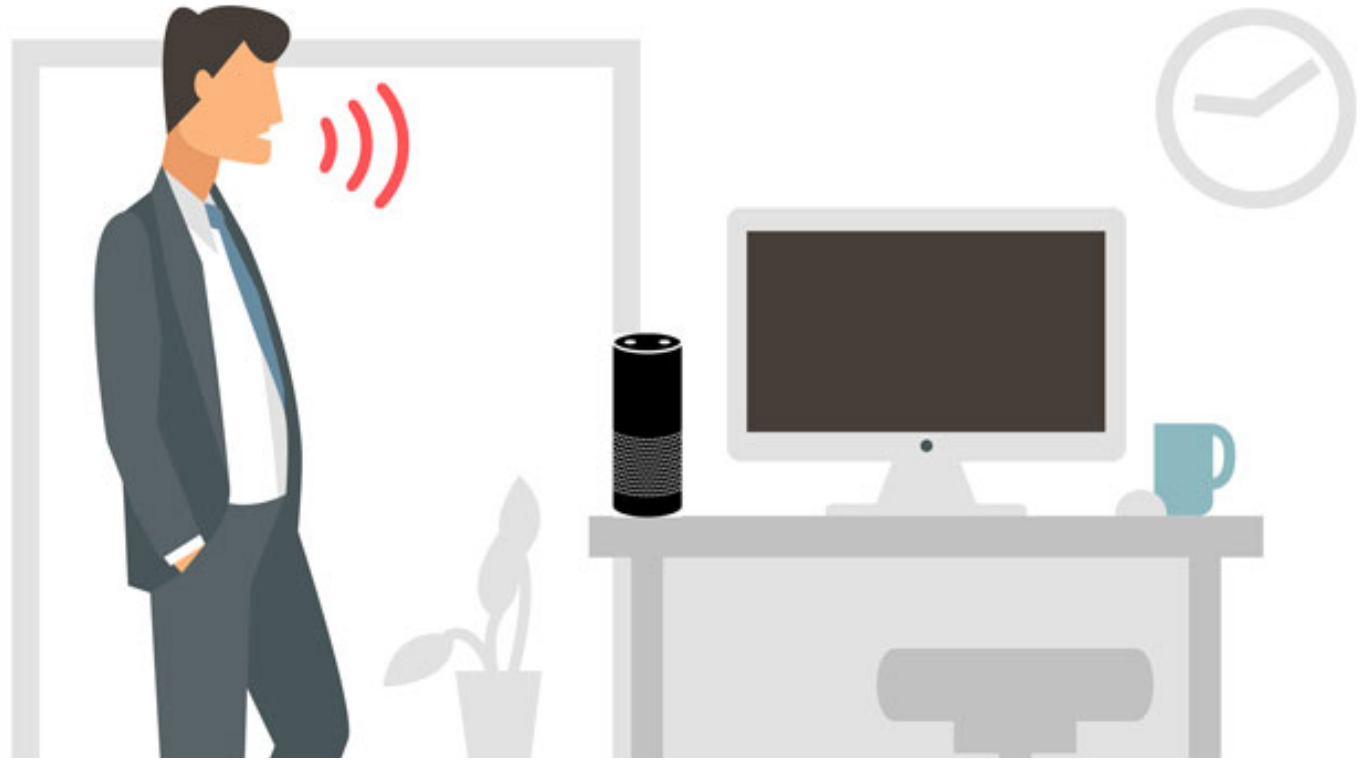


- Motivation
- ASR Workchain
- Deep Learning ASR
 - Preprocessing
 - Connectionist Temporal Classification (CTC)
 - Decoding
- Example/Summary

Speech recognition applications



Mobile devices



Hands-free interaction



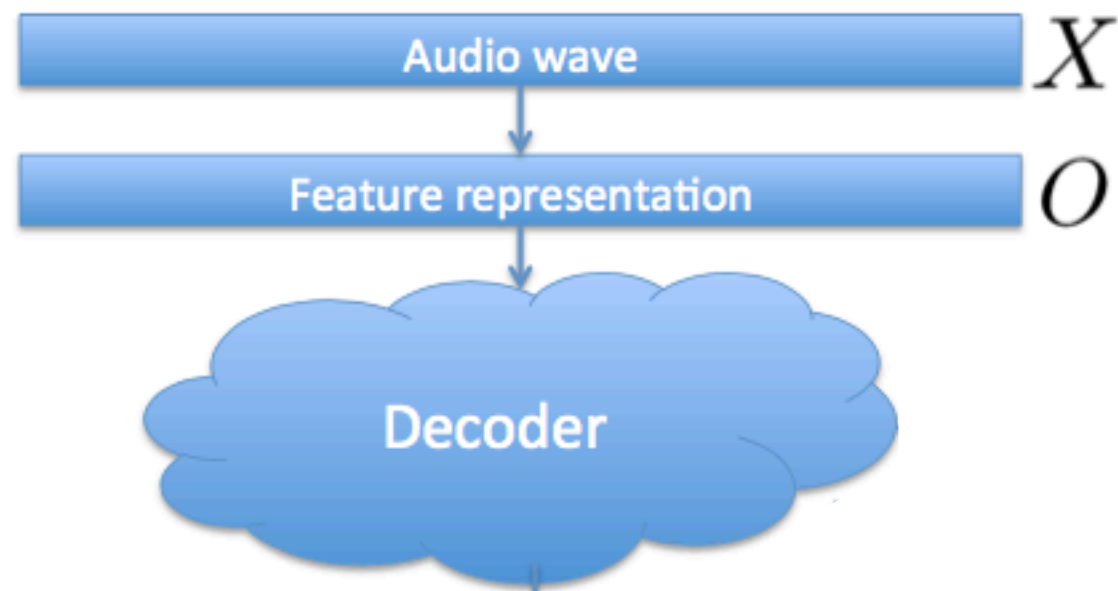
Content captioning



Voice verification

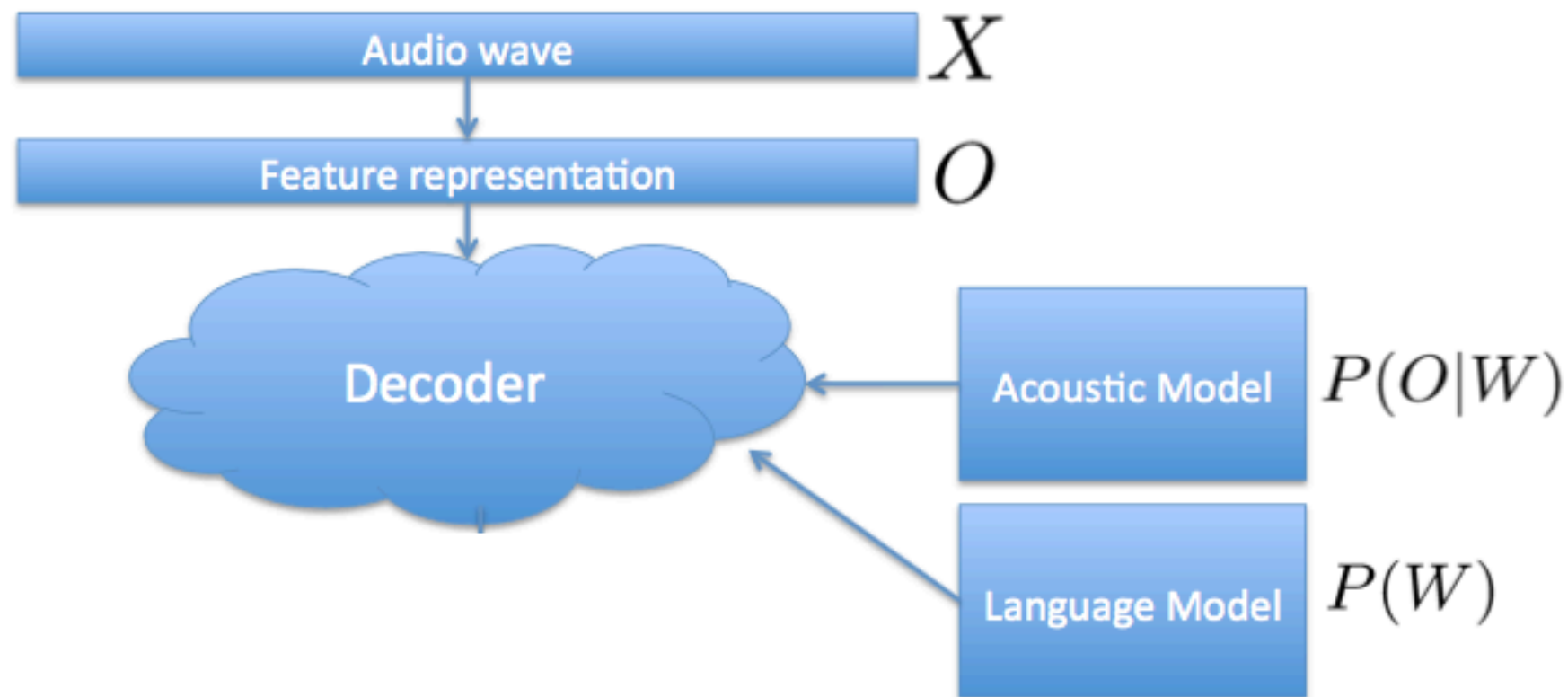
ASR Workchain

- Split up the task into different components:



ASR Workchain

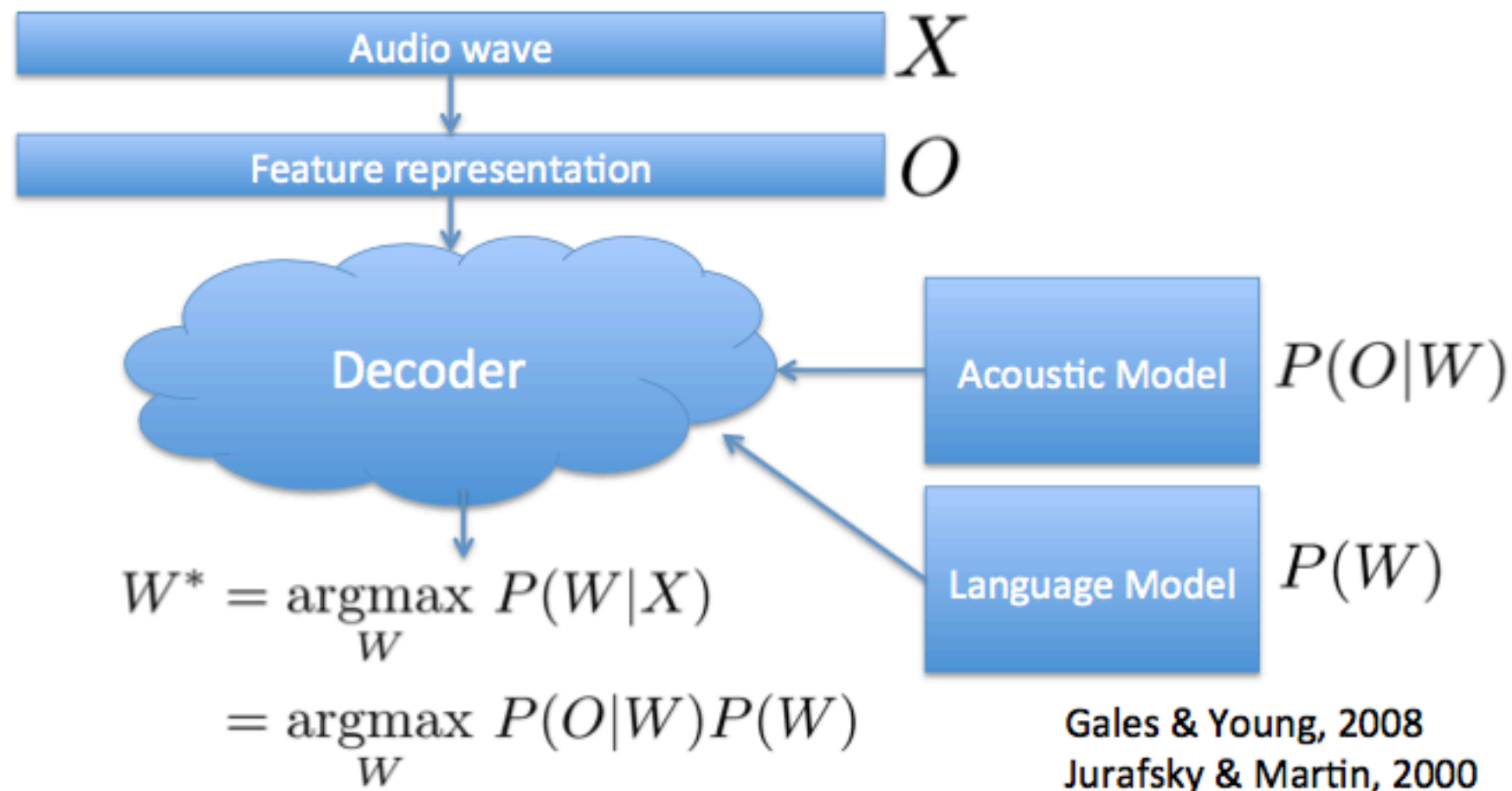
- Split up the task into different components:



Gales & Young, 2008
Jurafsky & Martin, 2000

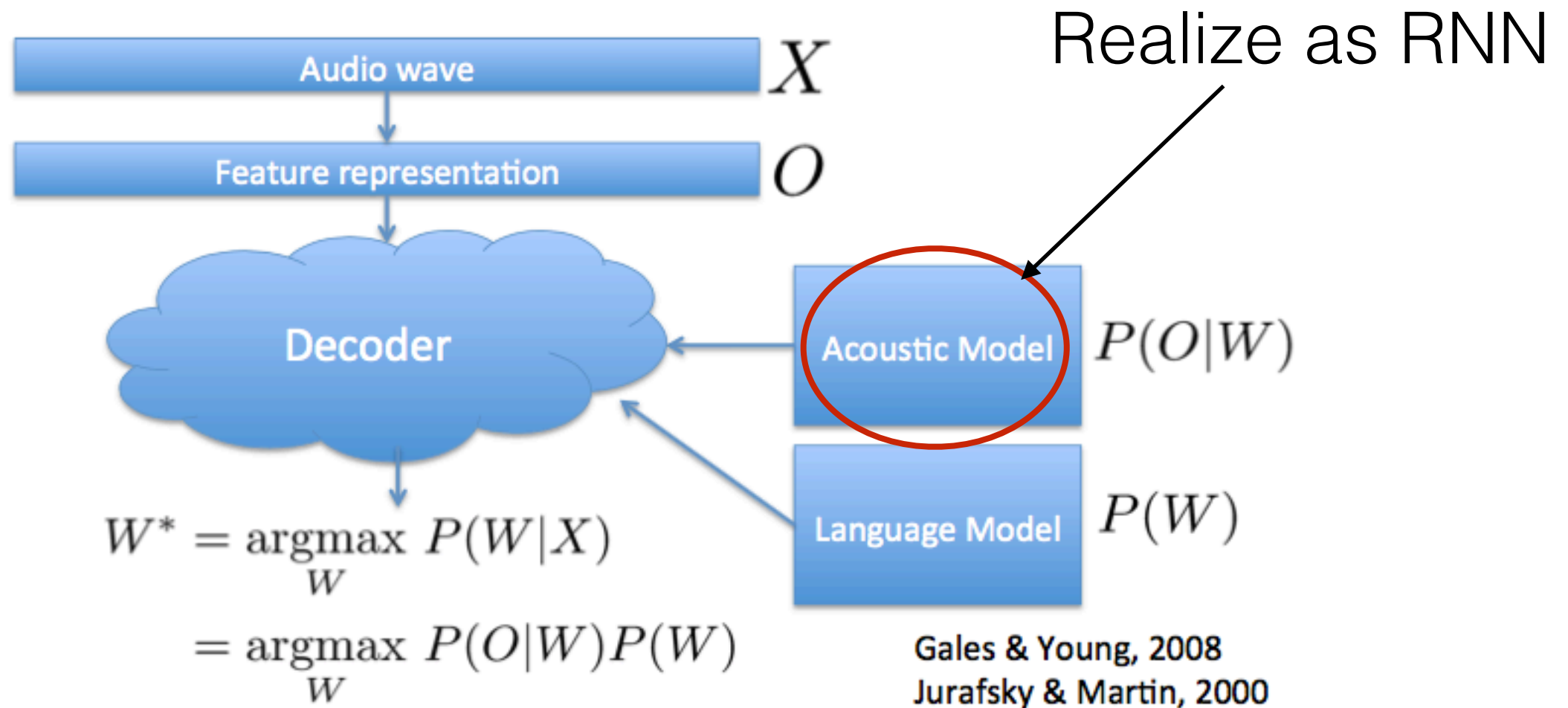
ASR Workchain

- Split up the task into different components:

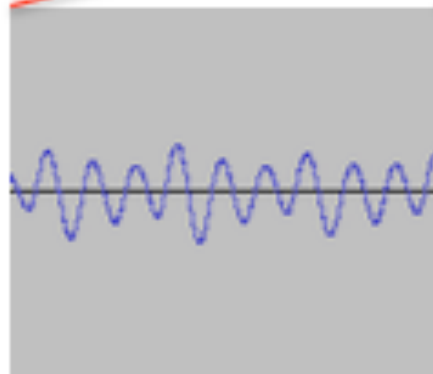
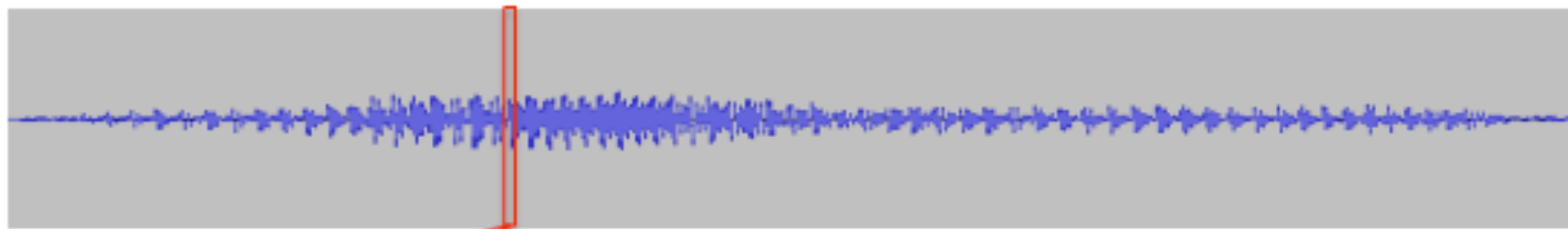


ASR Workchain

- Split up the task into different components:



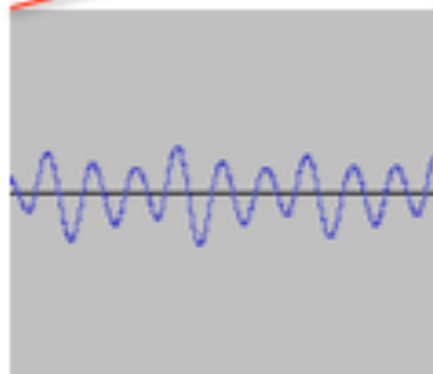
Preprocessing



20ms

Preprocessing

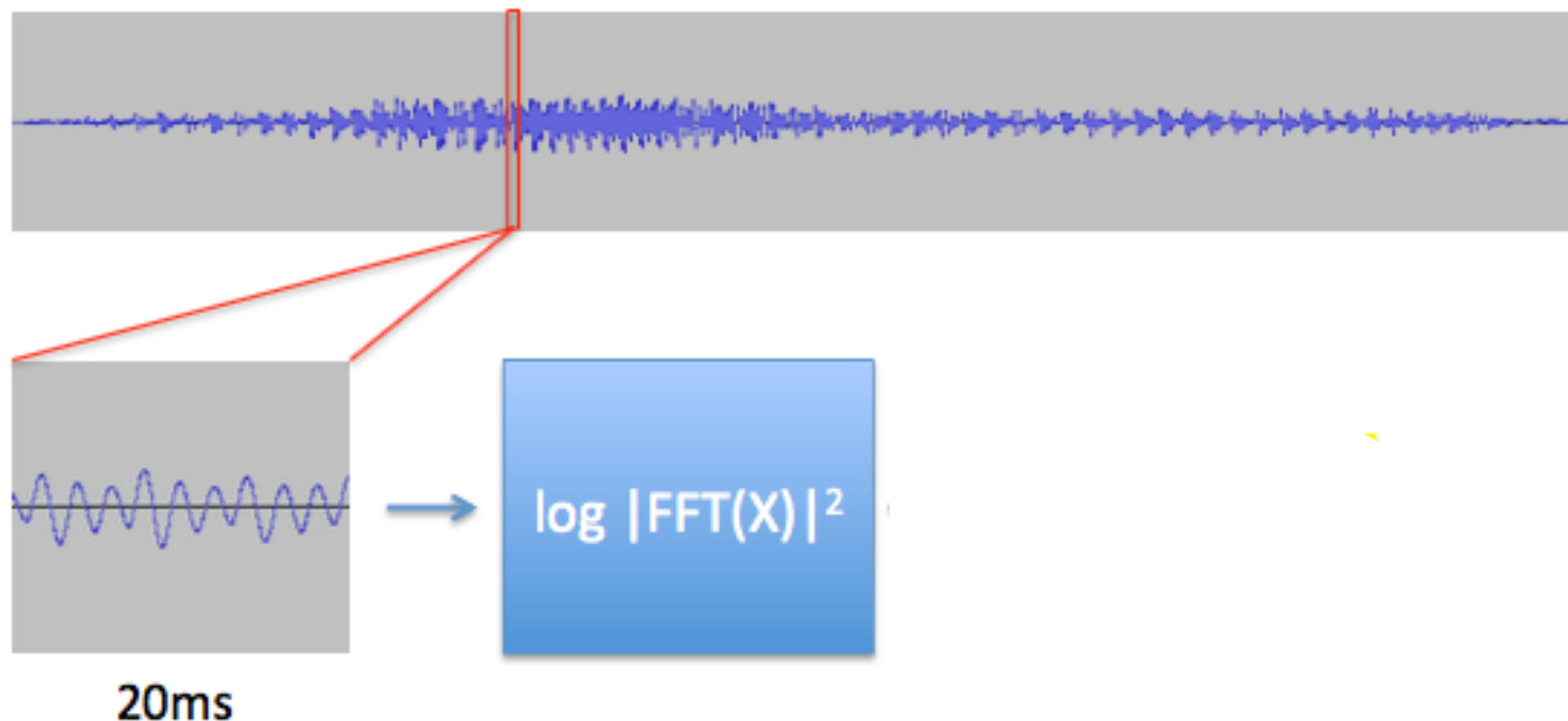
- Sample audio signal/“discretize”



20ms

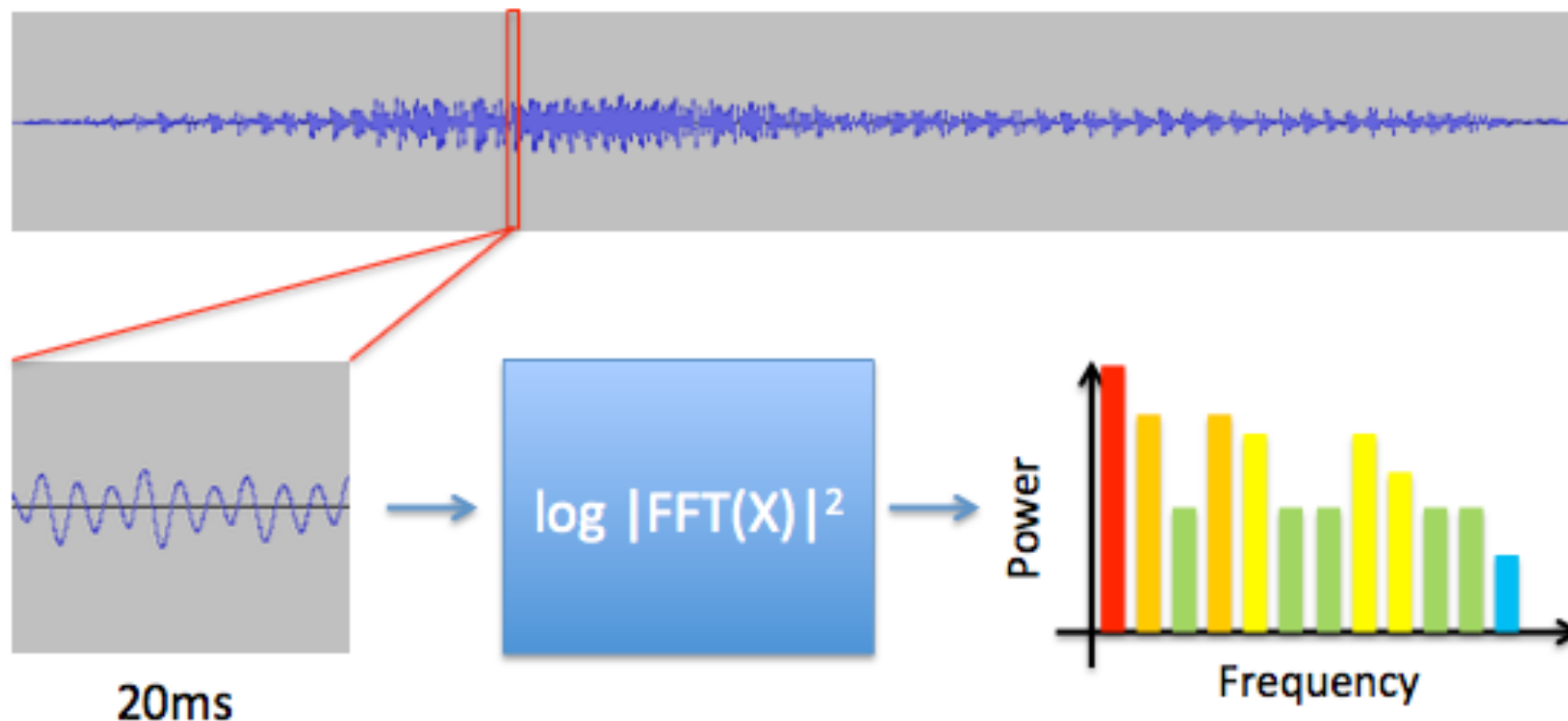
Preprocessing

- Sample audio signal/“discretize”
- Extract power spectrum (i.e. the features)



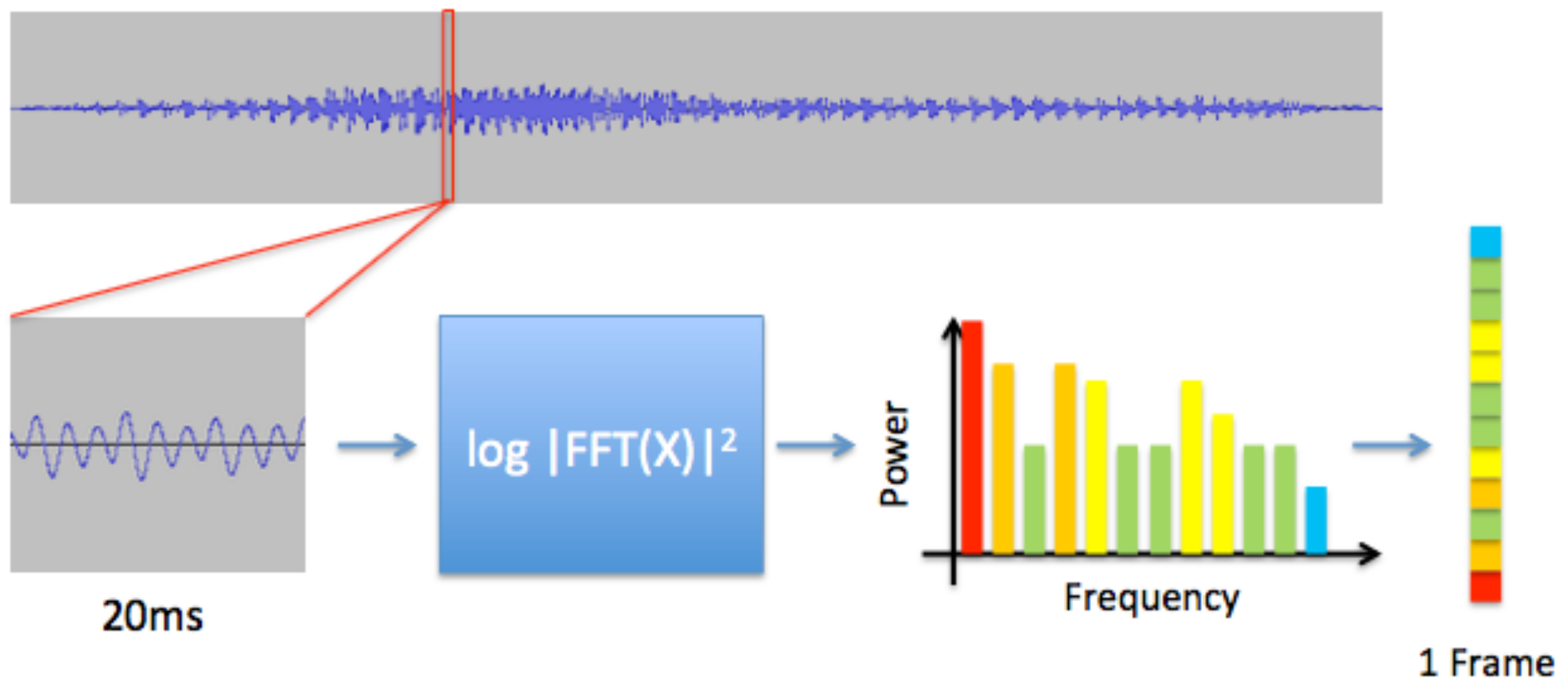
Preprocessing

- Sample audio signal/“discretize”
- Extract power spectrum (i.e. the features)



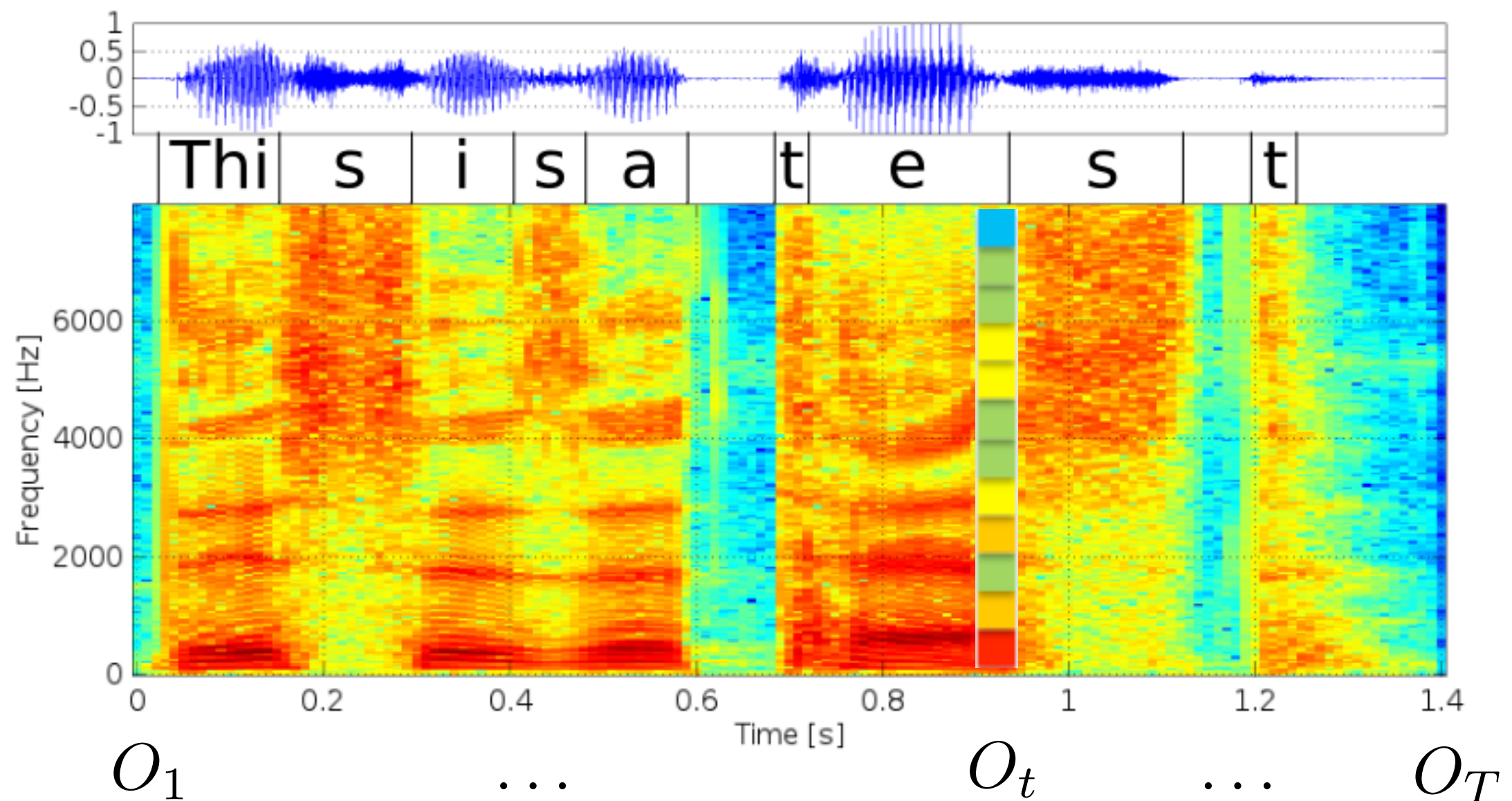
Preprocessing

- Sample audio signal/“discretize”
- Extract power spectrum (i.e. the features)
- Construct local feature vector from power spectrum



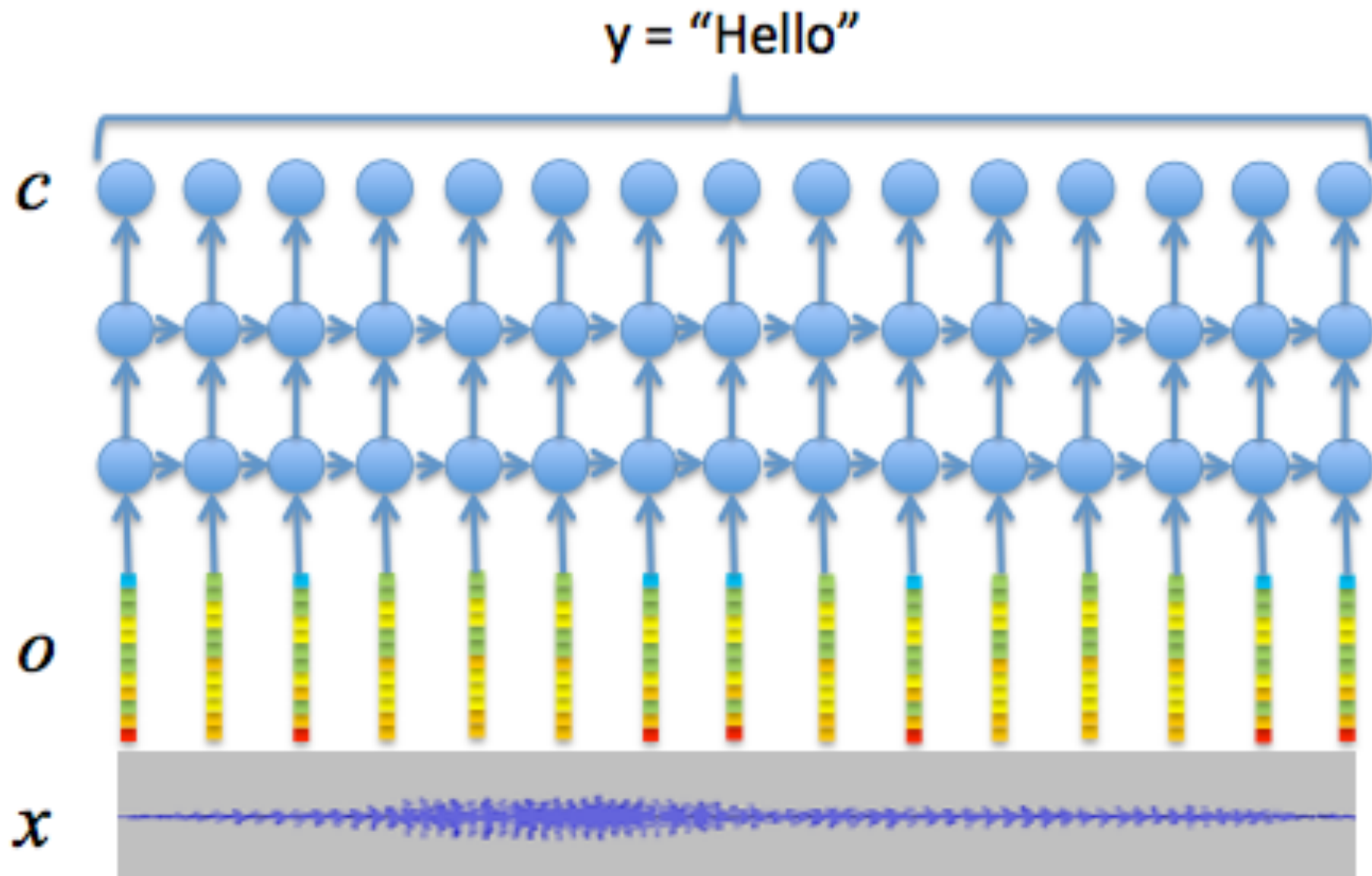
Preprocessing

Generate full spectrogram for the audio sequence



RNN Acoustic Model

Create RNN that outputs for sequence x transcription y



Connectionist Temporal Classification

Main issue: $\text{length}(y) \leq \text{length}(x)$

Connectionist Temporal Classification

Main issue: $\text{length}(y) \leq \text{length}(x)$

Addressed by CTC [Graves *et al.*, 2006]:

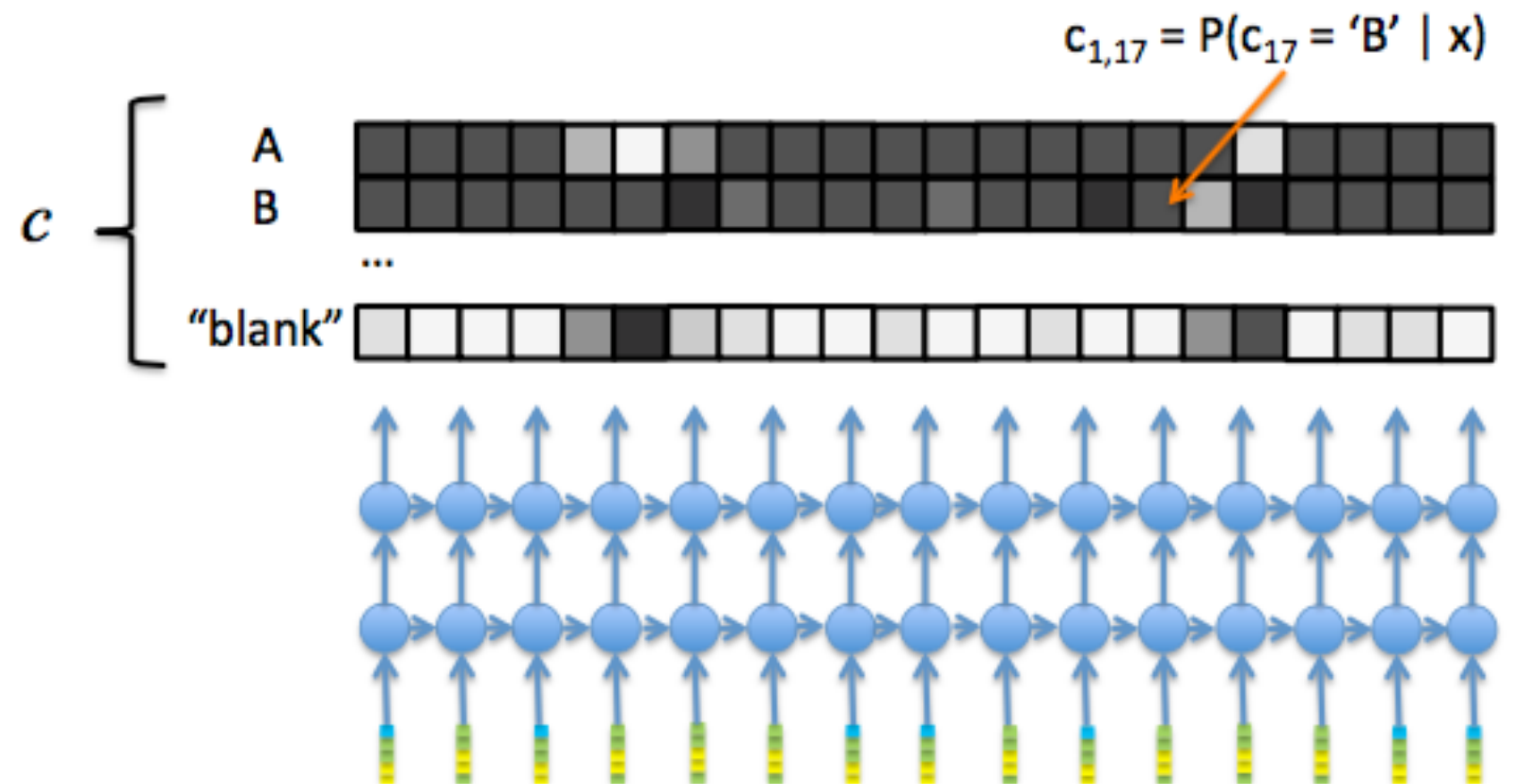
1. Output neurons c encode distribution over letters:
($\text{length}(c) = \text{length}(x)$) :

$$c \in \{A, B, C, \dots, Z, \text{blank}, \text{space}\}$$

2. Define a mapping $\beta(c) \rightarrow y$
3. Maximize the likelihood of true labels y^* for given input x under this model

Connectionist Temporal Classification

1. RNN output:

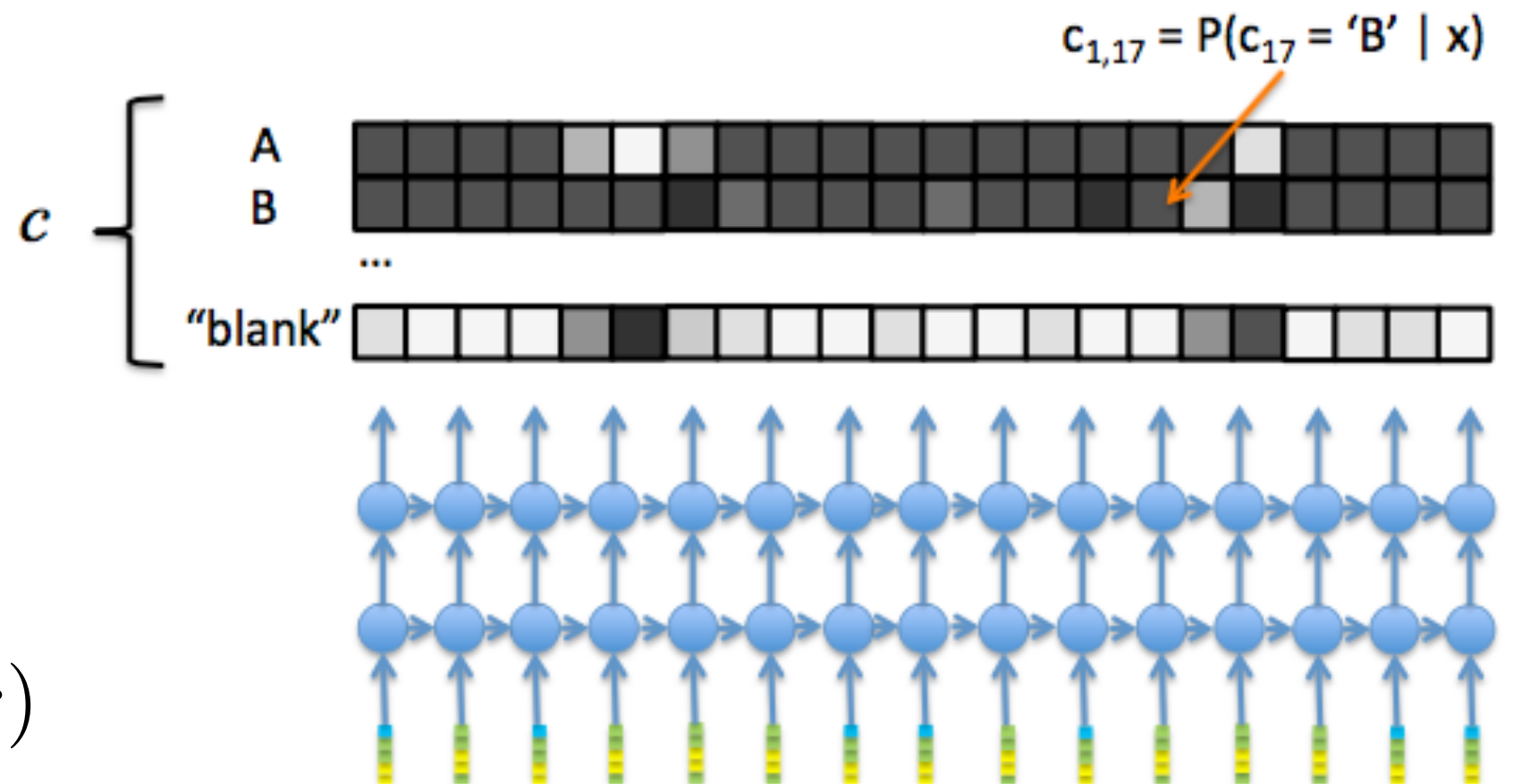


Connectionist Temporal Classification

1. RNN output:

Under assumption of independence, output defines distribution over whole sequences of characters c :

$$P(c|x) = \prod_{i=1}^N P(c_i|x)$$



$$P(c = \text{HHH_E_LL_LO_}___|x) = P(c_1 = \text{H}|x)P(c_2 = \text{H}|x) \dots P(c_{17} = \text{blank}|x)$$

Connectionist Temporal Classification

2. Define mapping:

- Eliminate duplicate characters then remove blanks:

$$\beta(c = \text{HHH_E_LL_LO_}) = \text{"HELLO"}$$

Connectionist Temporal Classification

2. Define mapping:

- Eliminate duplicate characters then remove blanks:

$$\beta(c = \text{HHH_E_LL_LO_}) = \text{"HELLO"}$$

- Mapping implies distribution over *transcriptions* y :

$$P(c|x) = \begin{cases} 0.1 & \text{HHH_E_LL_LO_} & \text{"HELLO" v.1} \\ 0.02 & \text{HH_E_LL_LO_} & \text{"HELLO" v.2} \\ 0.01 & \text{HHH_E_L_L_LOH_} & \text{"HELL OH"} \\ 0.01 & \text{HHH_EE_LL_L_O_} & \text{"HELLO" v.3} \\ \dots & \text{YY_E_LL_LO_W_} & \text{"YELLOW"} \end{cases}$$

$$P(y|x) = \sum_{c:\beta(c)=y} P(c|x)$$

$$P(\text{"HELLO"}) = 0.1 + 0.02 + 0.01 + \dots$$

Connectionist Temporal Classification

3. Maximize likelihood:

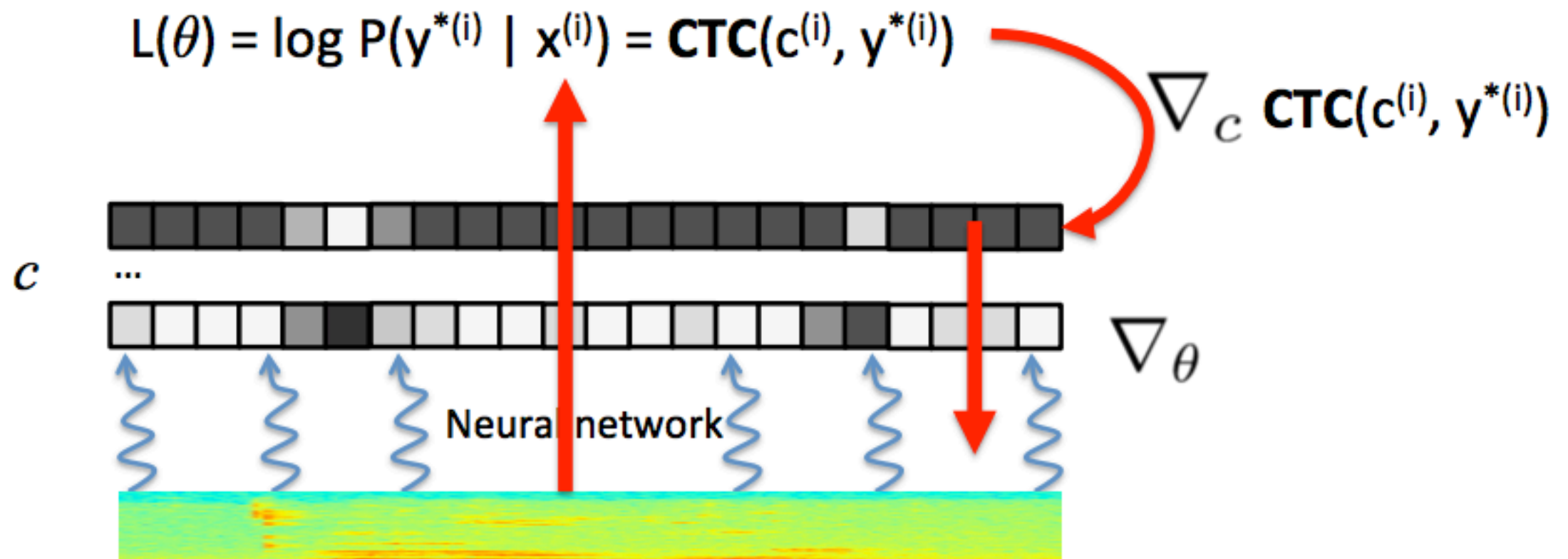
Update model parameters θ such that correct label y^* maximizes the log likelihood:

$$\begin{aligned}\theta^* &= \arg \max_{\theta} \sum_i \log P(y^{*(i)} | x^{(i)}) \\ &= \arg \max_{\theta} \sum_i \log \sum_{c: \beta(c)=y^{*(i)}} P(c | x^{(i)})\end{aligned}$$

Connectionist Temporal Classification

How do we maximize:

- Use gradient descent methods on the CTC objective function (based on log likelihood)
- Backpropagate error and adjust model parameters θ



Summary

- Deep Learning is becoming more and more important for state-of-the-art speech recognition
- Success of DL RNNs based on CTC training algorithm
- Still a lot of engineering required for high accuracy speech system (language models, huge data sets, ...)

References

- A Coates. “Lecture at Bay Area Deep Learning School 2016” <https://youtu.be/9dXiAecyJrY?t=13874>
- A Geitgey. “Machine Learning is Fun!” <https://medium.com/@ageitgey/machine-learning-is-fun-part-6-how-to-do-speech-recognition-with-deep-learning-28293c162f7a>
- A Graves, S Fernández, F Gomez, J Schmidhuber. “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks.” ICML, 2006.
- D Kriesel. “A Brief Introduction to Neural Networks” http://www.dkriesel.com/en/science/neural_networks