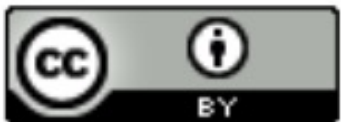


Deep Learning for Speech Processing

Danna Gurari

University of Colorado Boulder

Spring 2022



Review

- Last week:
 - Efficient learning: curriculum learning
 - Efficient learning: active learning
 - Reinforcement learning
- Assignments (Canvas):
 - Final project outline due Friday
 - Final project video due in two weeks
- Questions?

Today's Topics

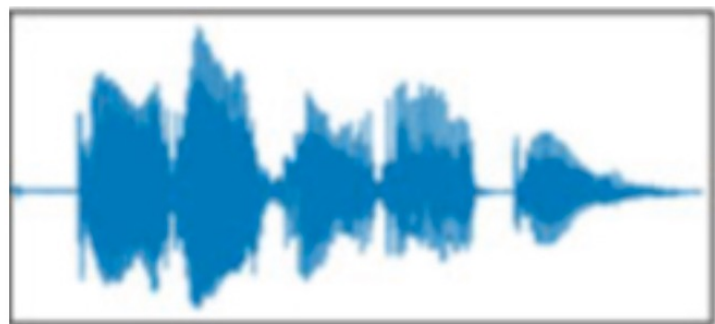
- Problem
- Applications
- Speech recognition evaluation
- Speech recognition models
- Video making tutorial

Today's Topics

- Problem
- Applications
- Speech recognition evaluation
- Speech recognition models
- Video making tutorial

Problem Definition

Input: spoken language



Raw Speech Signal

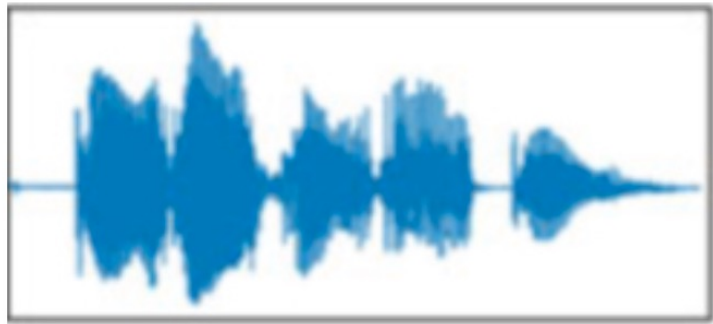


Output: machine readable text

Do you understand me

Transcription

What Is Speech?



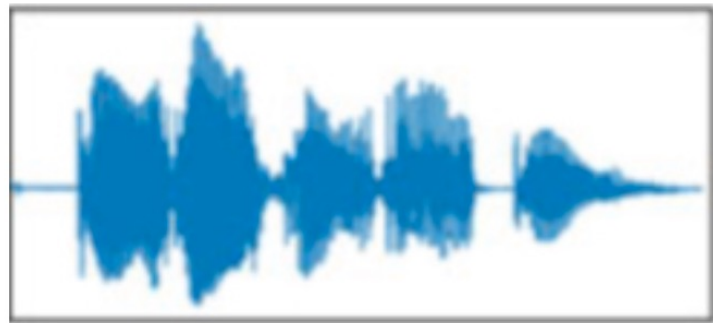
Raw Speech Signal

Compression waves created by pushing air from one's lungs and modulating it using one's tongue, teeth, and lips

Why Is Speech Processing Challenging?

Input can be diverse including different accents, volumes, pace, and cadence

Temporal data needs to be **segmented** into distinct words



Raw Speech Signal



Do you understand me

Transcription

Technology can result in many artifacts including varying quality, echos, and background noise

Today's Topics

- Problem
- **Applications**
- Speech representation
- Speech recognition models and evaluation
- Video making tutorial

Voice Typing on Mobile Devices



Voice Typing for Productivity Applications



Demo starting at 2:00: <https://www.youtube.com/watch?v=5UK4vLzU9co&t=76s>

Virtual Assistant



e.g., Amazon's Echo with Alexa

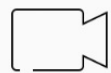


e.g., Google Home



e.g., Baidu DuerOS

Virtual Assistant



Entertainment
Video

Music, movie, television shows, variety show, short clip, audio book, and broadcasting, etc.



Information
Inquiry

Consultation, weather, stocks, flight, sports (NBA), FAQ, cookbook, images, etc.



Lifestyle
Services

Food, movie, take-out, hotel, shopping, taxi, cleaning service, travel, relaxation, and other O2O services.



Travel
Conditions

Map, route, road condition, traffic restriction, endorsement, and surrounding environment query, etc.



Utility Tools

Translation, time, calculation, exchange rate, and unit conversion, etc.



Personal
Assistant

Schedule management, alarm clock, reminder, memo and notepad, etc.



Learning

Encyclopedia, story, nursery rhyme, idiom, parenting, poetry and library, etc.



Chat and Relax

Chat, joke, poetry, idioms, and games, etc.

Audio Transcription (e.g., for Analysis & Situational/Permanent Hearing Impairments)



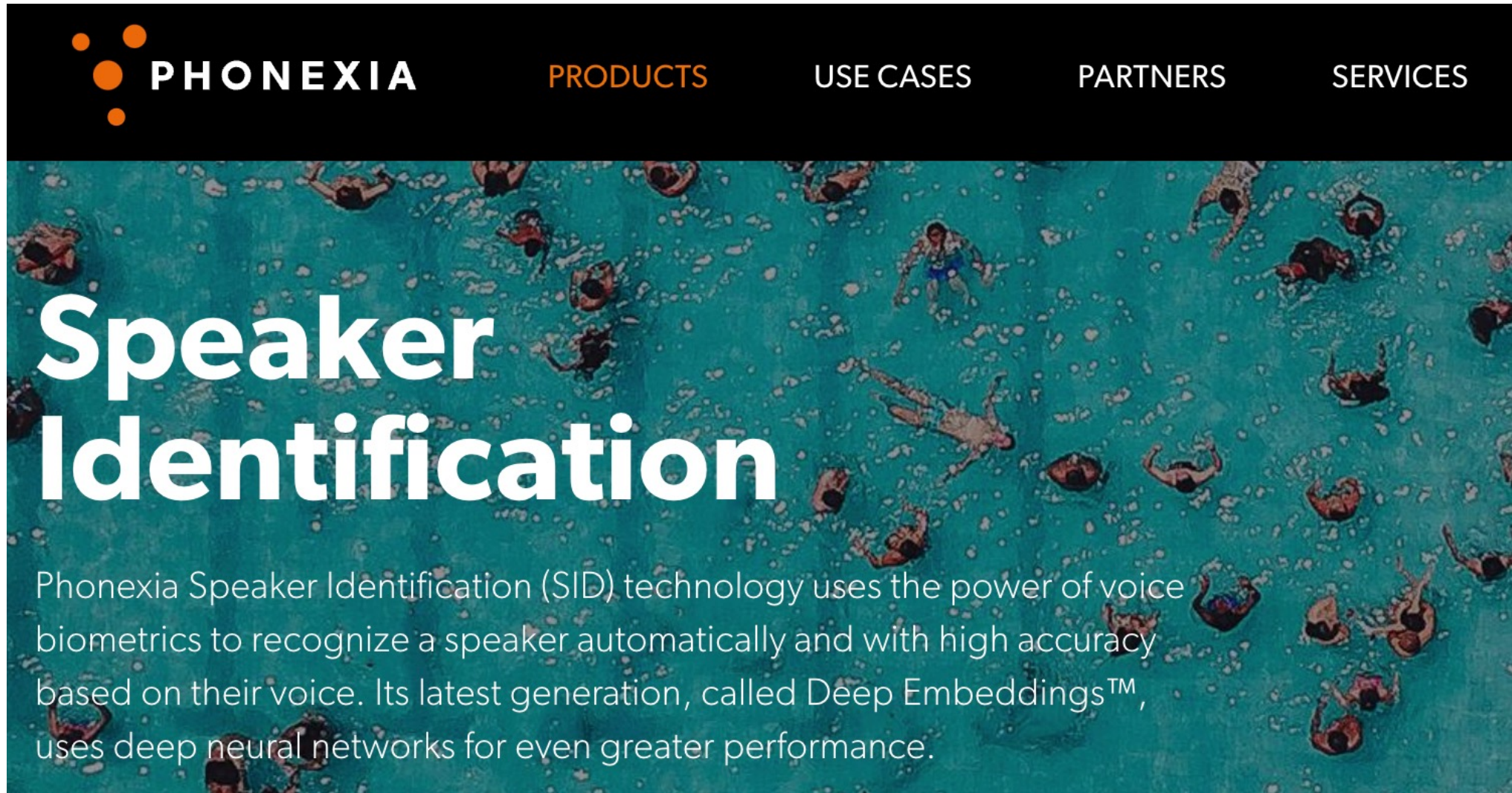
Video/Movie Captioning (e.g., for Translation & Situational/Permanent Hearing Impairments)



Speech Emotion Recognition (e.g., for Help Desks and Negotiators)



Speaker Identification (e.g., for Security)

The image shows a screenshot of the Phonexia website. At the top, there is a black navigation bar with the Phonexia logo (four orange dots) and the word 'PHONEXIA' in white. To the right of the logo are five menu items: 'PRODUCTS', 'USE CASES', 'PARTNERS', and 'SERVICES', all in white uppercase letters. Below the navigation bar is a large hero section with a teal background featuring an aerial view of many people swimming in a pool. The text 'Speaker Identification' is written in large, bold, white letters on the left side of the hero section. Below this, a paragraph of white text describes the technology: 'Phonexia Speaker Identification (SID) technology uses the power of voice biometrics to recognize a speaker automatically and with high accuracy based on their voice. Its latest generation, called Deep Embeddings™, uses deep neural networks for even greater performance.'

Speaker Identification

Phonexia Speaker Identification (SID) technology uses the power of voice biometrics to recognize a speaker automatically and with high accuracy based on their voice. Its latest generation, called Deep Embeddings™, uses deep neural networks for even greater performance.

Language Identification

Automatic language identifier

Insert any text or pick a random example

Bonjour!

Speech Enhancement



[ABOUT](#)

[FAQ](#)

[DOWNLOAD](#)

**FAKIN'
THE
FUNK
?**



[BUY](#)

[FORUM](#)

[CONTACT](#)

Fakin' The Funk? is a tool that helps you to detect the true quality of your audio files in one batch.

What are other potential applications for speech processing?

Today's Topics

- Problem
- Applications
- **Speech recognition evaluation**
- Speech recognition models
- Video making tutorial

Spectrum of Tasks



Commands

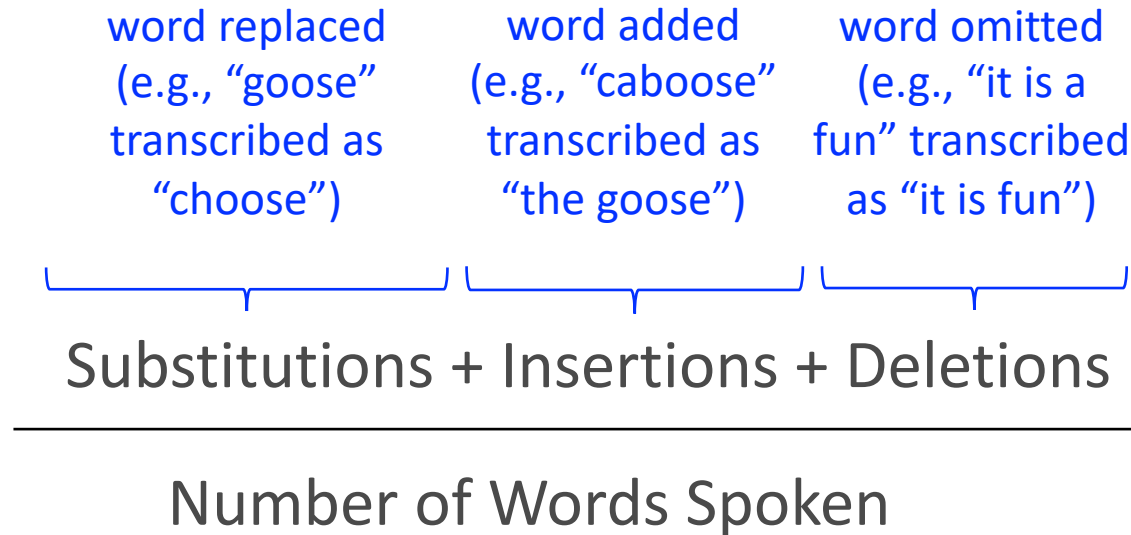
Discourse

(brief with small vocabulary)

(lengthy with large vocabulary)

Word Error Rate

- Indicates edit distance between the prediction and the target as follows:



- What indicates better performance: larger or smaller values?

Word Error Rate: Example

- Correct: The sun makes it look like uh a warm, day to go outside to adventure.
- Predicted: The son makes it to bike with a swarm to go outside to Denver today.
- Number of words spoken?
 - 15
- WER?

Substitutions + Insertions + Deletions

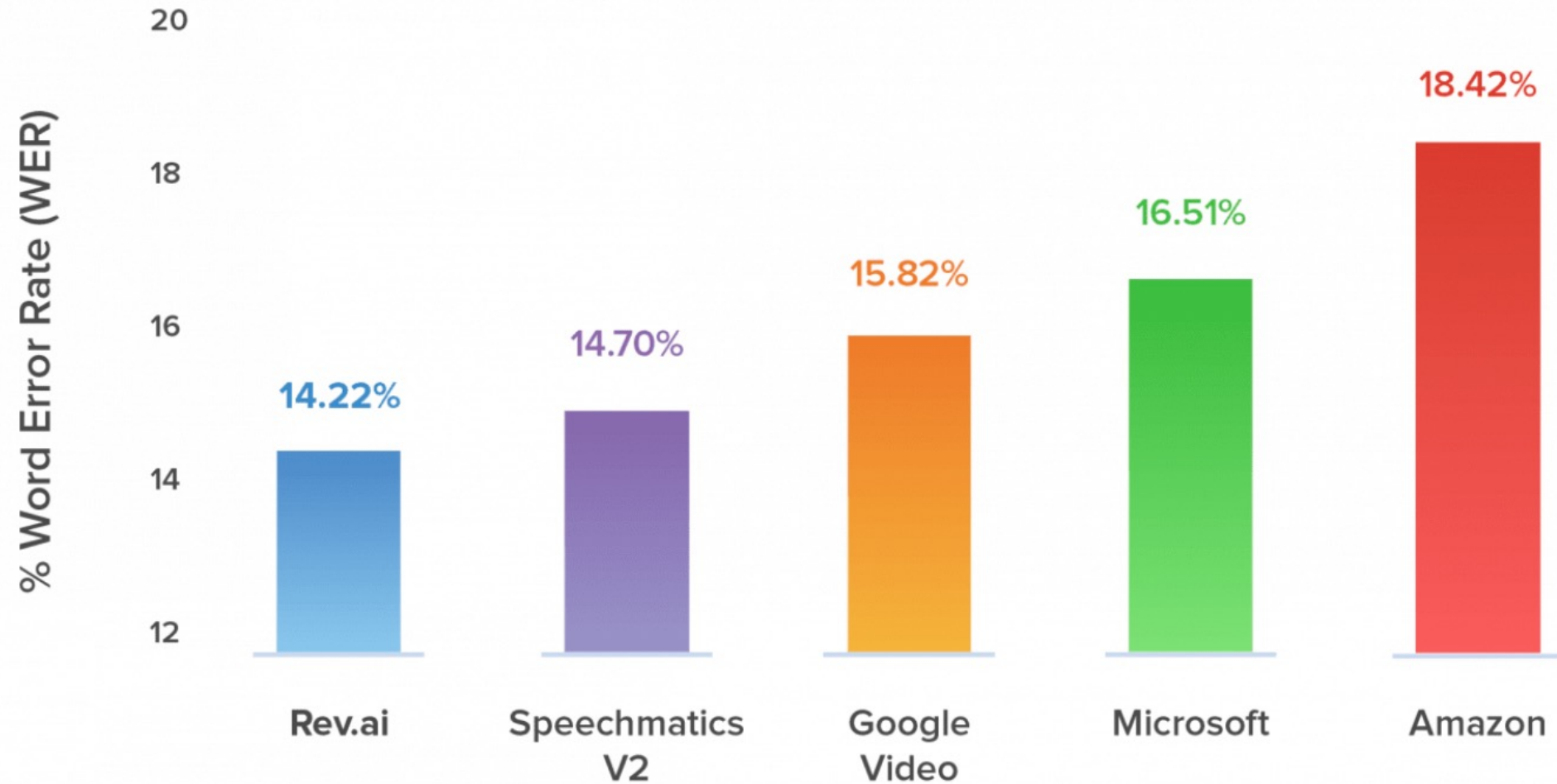
Number of Words Spoken

Word Error Rate: Example

- Correct: The sun makes it look like uh a warm, day to go outside to adventure.
- Predicted: The son makes it to bike with a swarm to go outside to Denver today.
- Number of words spoken?
 - 15
- WER?

$$\frac{6 + 1 + 1}{15} = 0.53$$

Word Error Rate: Comparison Example



Word Error Rate: What Are Its Limitations as an Evaluation Metric?

- Does not indicate why errors occur
 - Background noise (e.g., music, other talking)
 - Specialized language (i.e., words reflecting domain expertise)
 - Speaker pronunciations/accent
- Does not reflect whether transcription correctly captures:
 - Capitalization
 - Punctuation
 - Numbers
 - Paragraphs
- May indicate poor quality when humans could understand the content
- Weights all word errors equally

Today's Topics

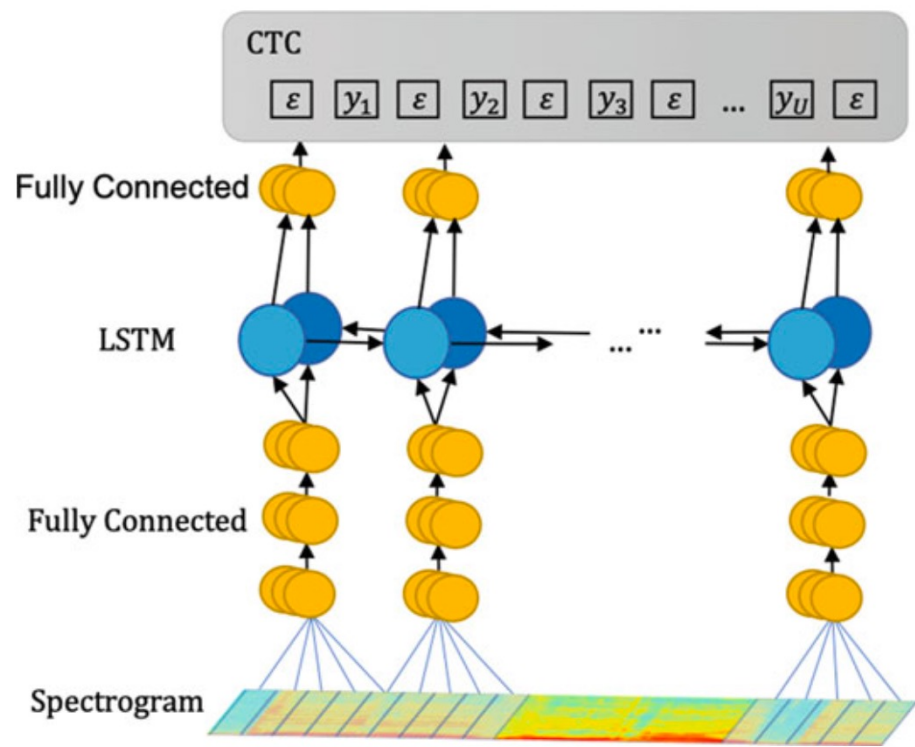
- Problem
- Applications
- Speech recognition evaluation
- **Speech recognition models**
- Video making tutorial

Popular Methods

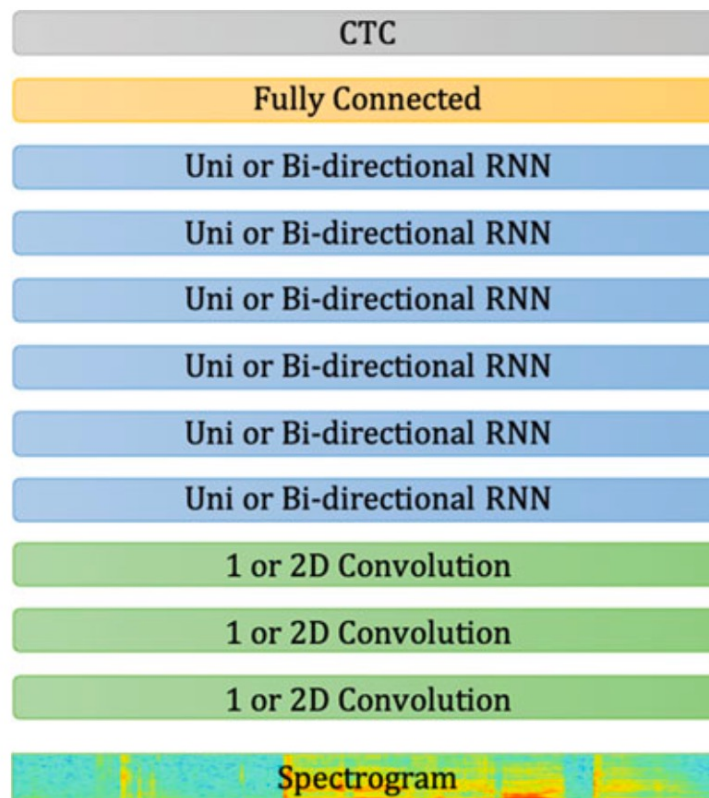
- Connectionist Temporal Classification (CTC)
- DeepSpeech
- DeepSpeech 2
- Listen, Attend, and Spell

Popular Methods

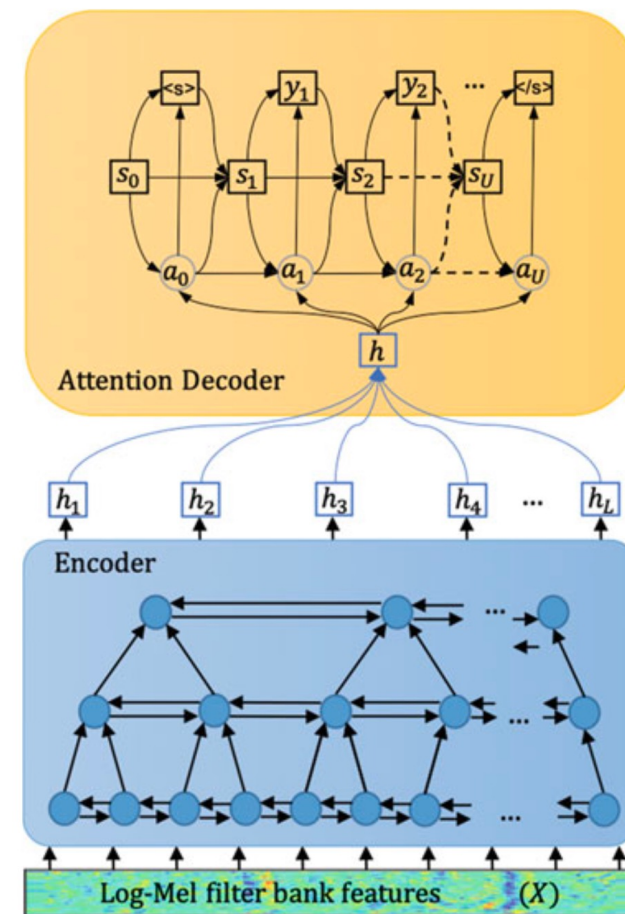
DeepSpeech



DeepSpeech2

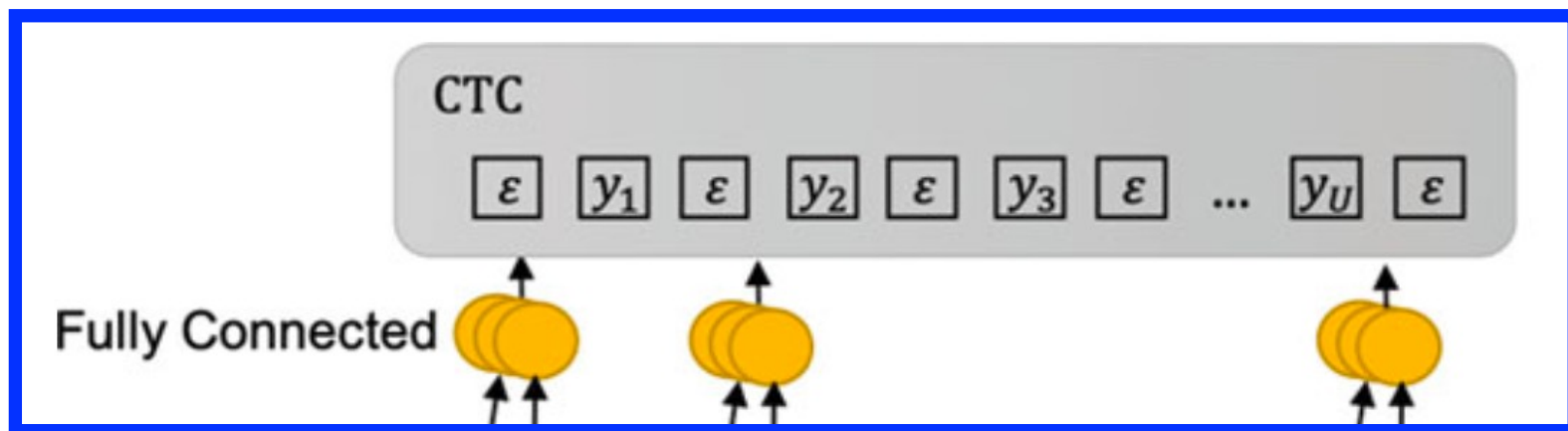


Listen, Attend, and Spell

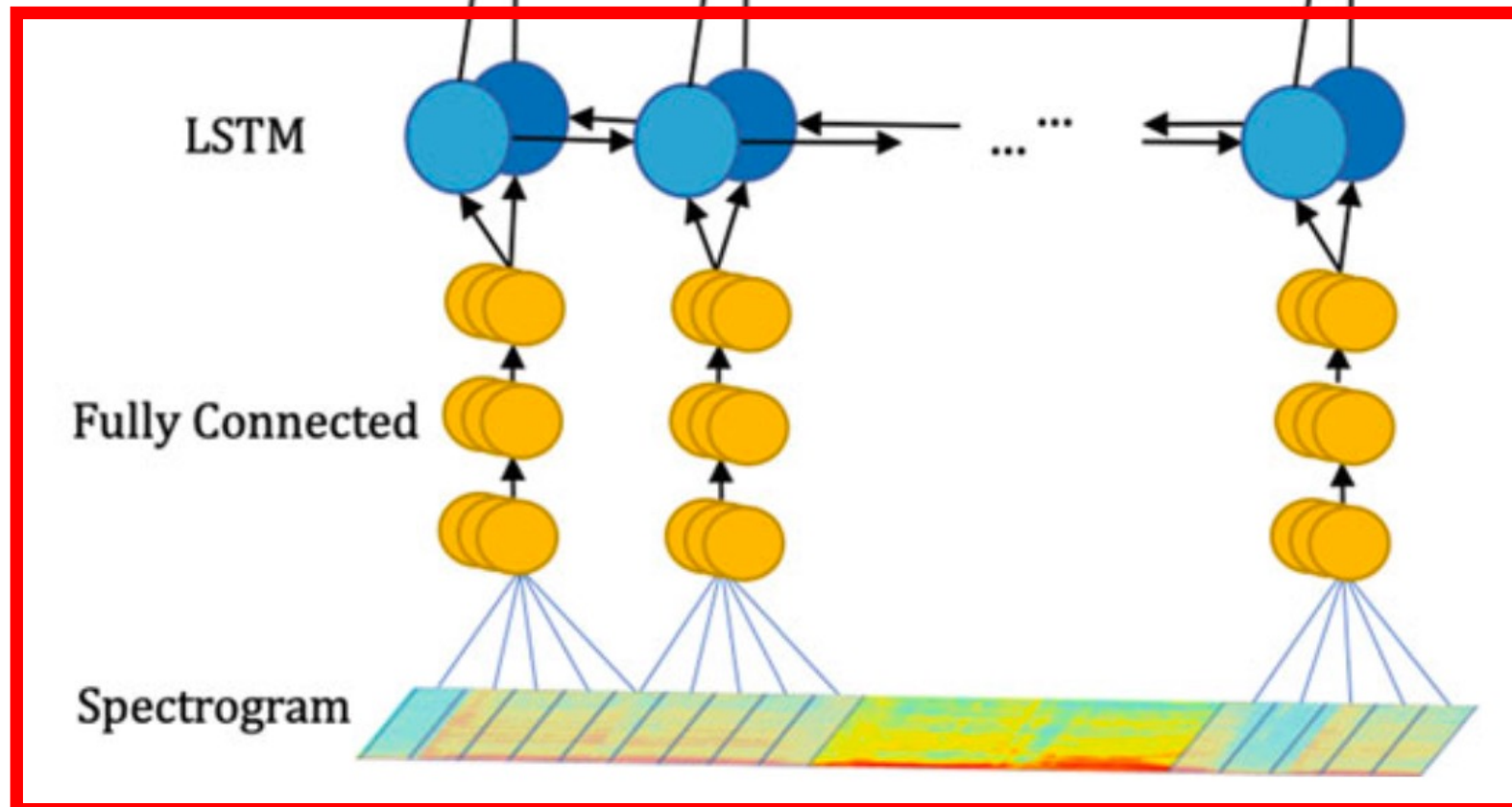


DeepSpeech

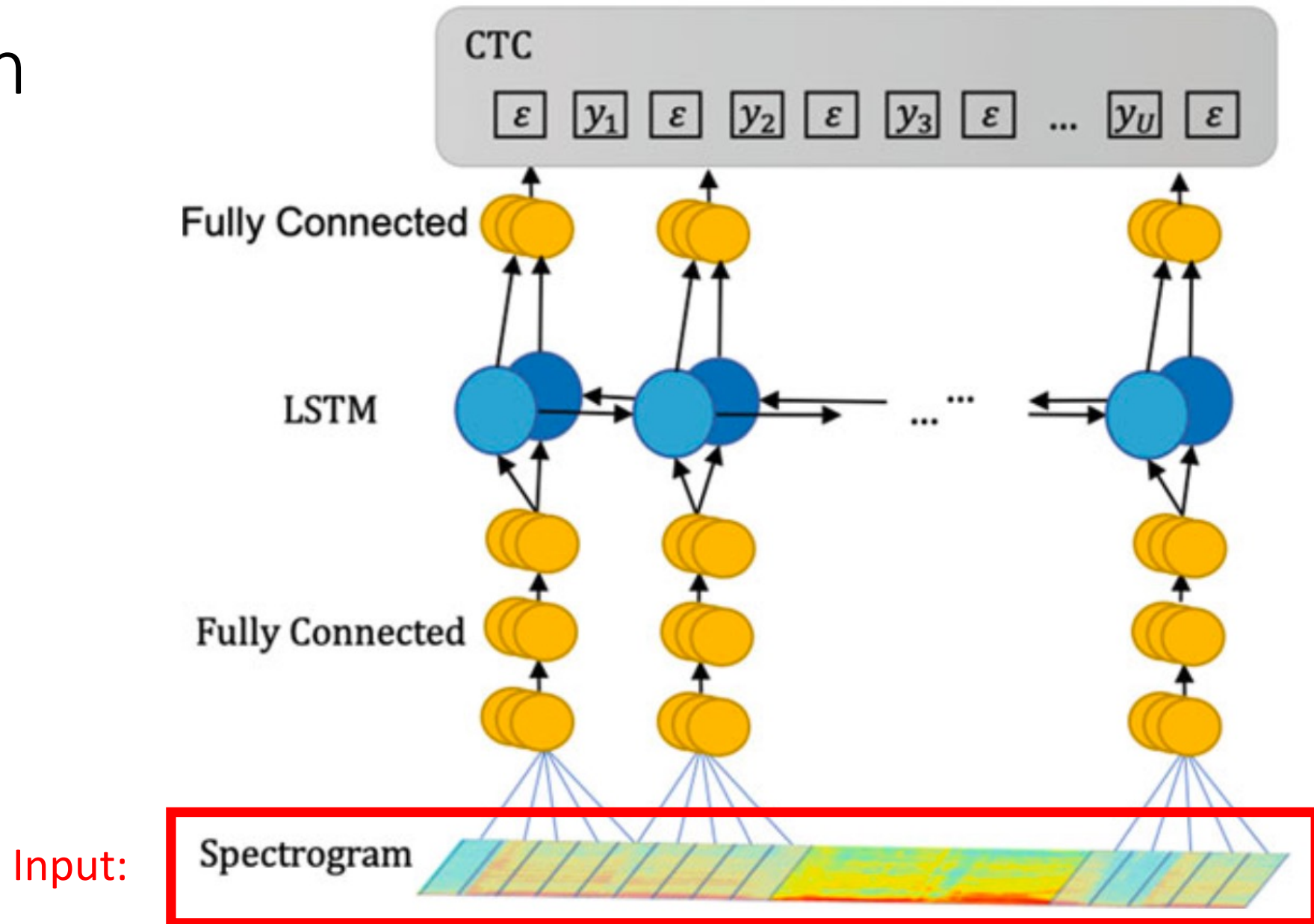
Decoder:



Encoder:

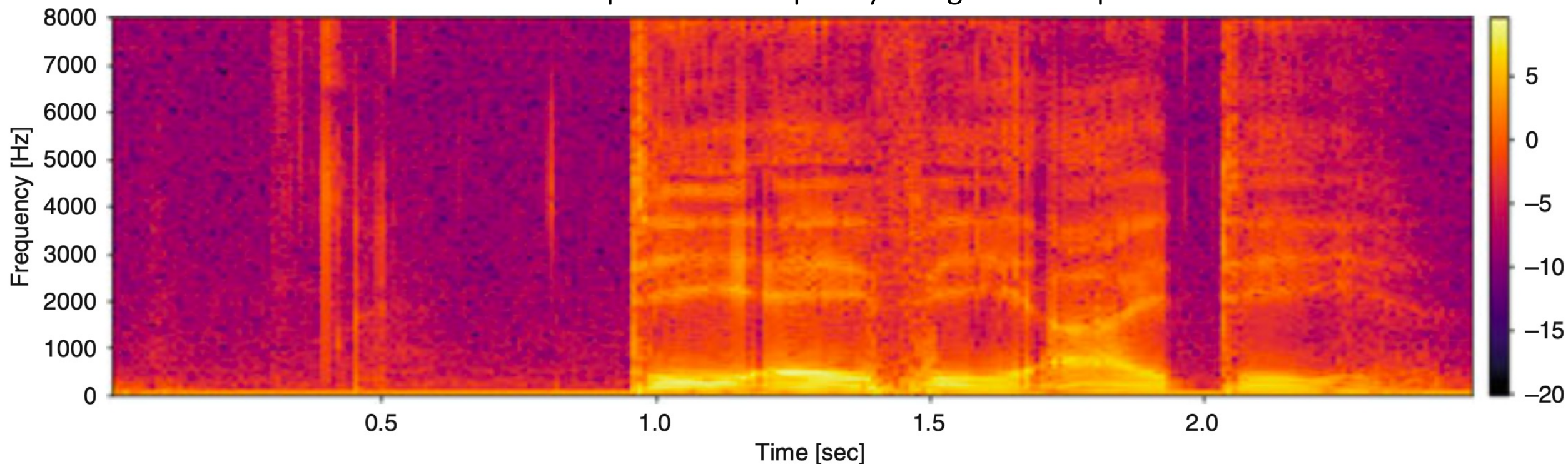


DeepSpeech



Spectrogram: Visual Representation of Audio

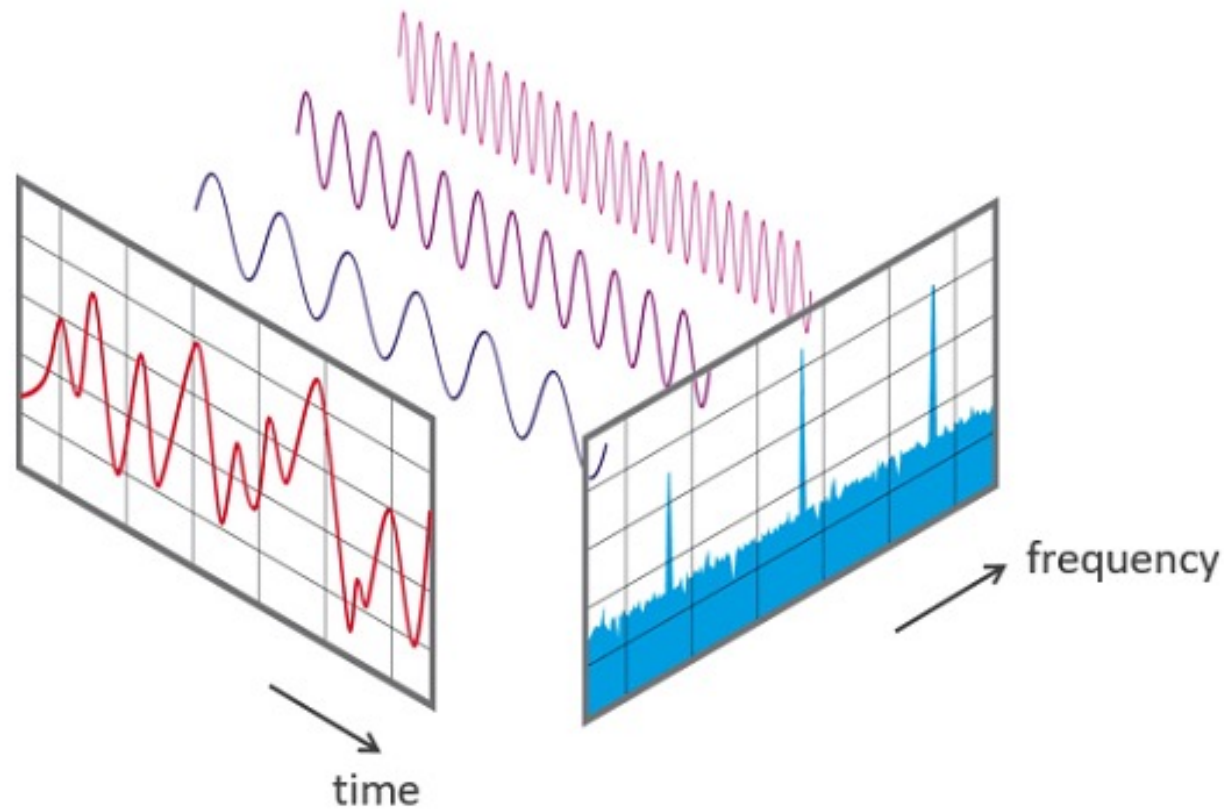
Color: amplitude of frequency at a given time point



Created by sliding a short window across the audio signal and applying a [Fourier transform](#) to each window

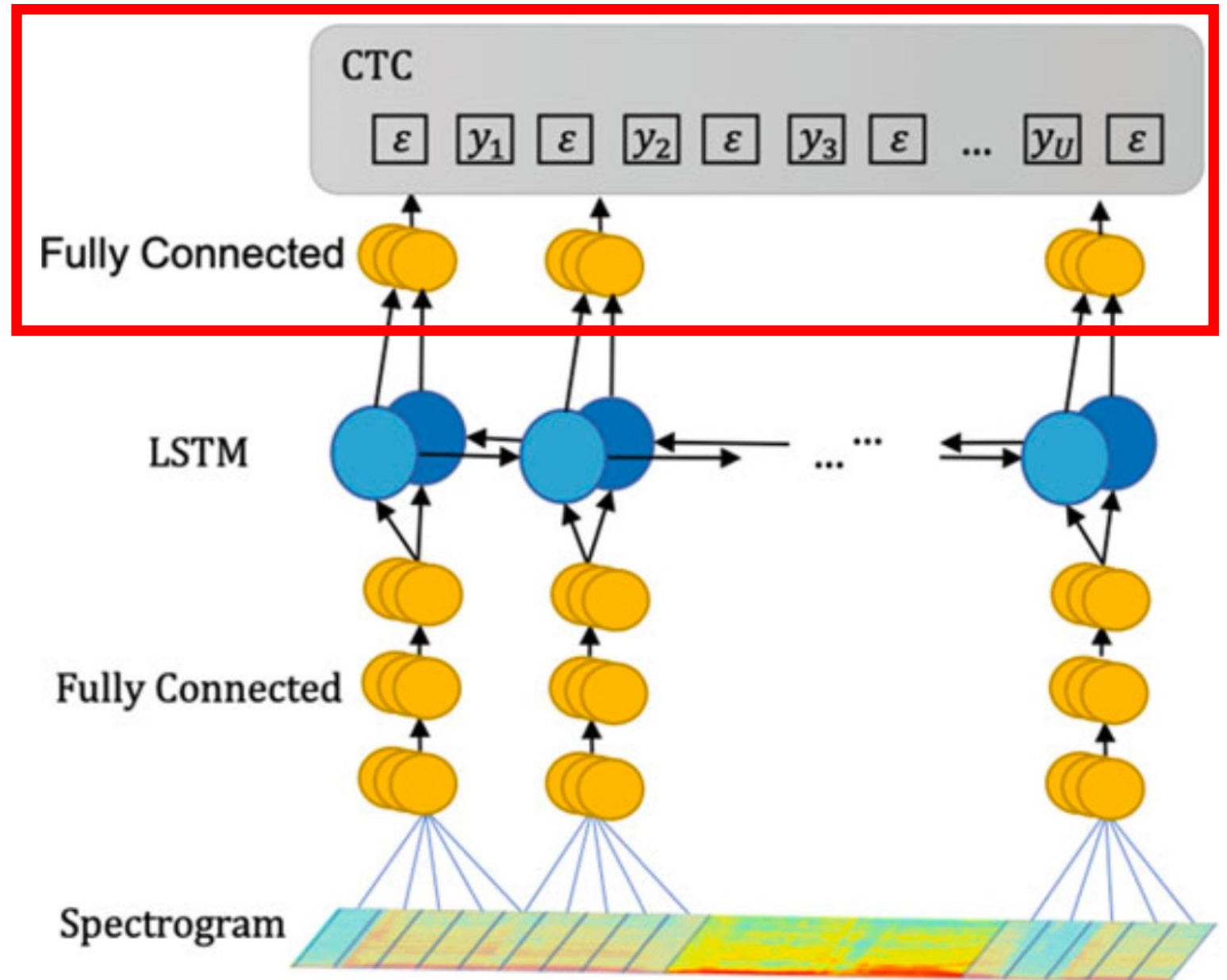
Background: Frequency Analysis of Audio Clip

Fourier transform: represents a signal as a sum of sines and cosines (*frequency-domain*):

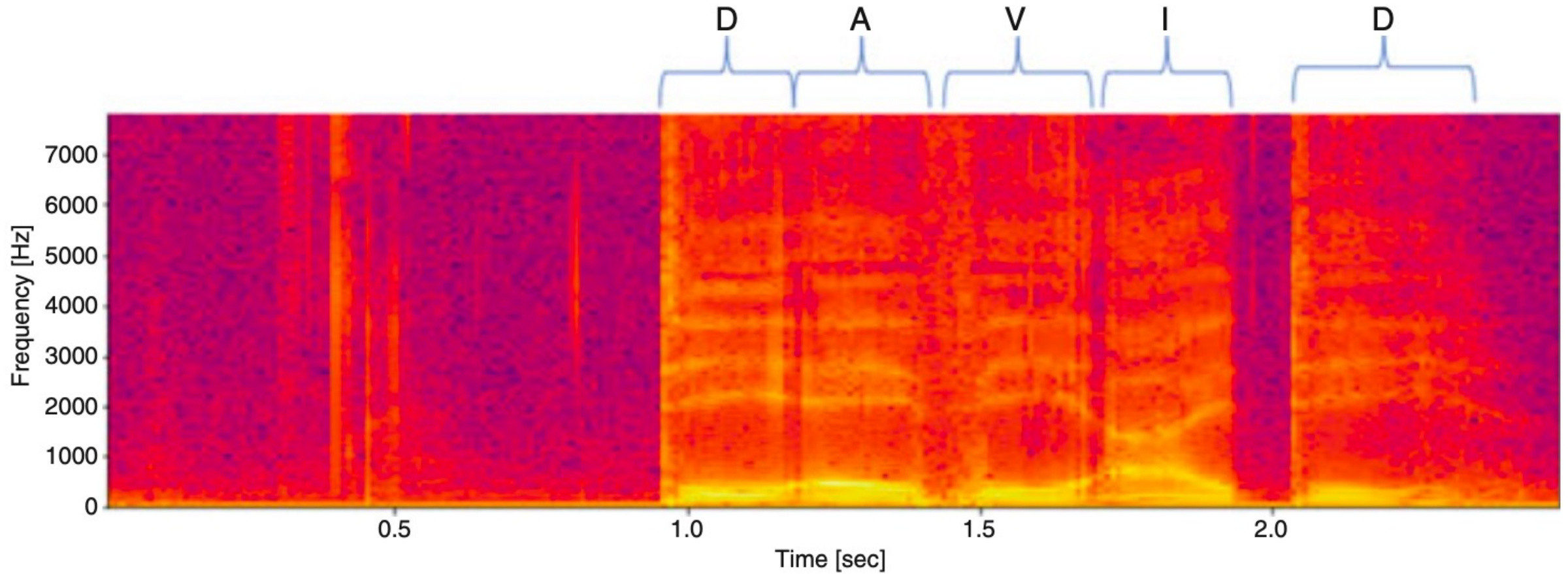


DeepSpeech

Output: character sequence predicted by a softmax layer



CTC: Input-Output Representation



CTC: Input-Output Representation

Key idea: **blank token** supports silent stretches and letter repeats (e.g., “hello” vs “helo”)

h h e ϵ ϵ l l l ϵ l l o

h e ϵ l ϵ l o

h e l l o

h e l l o

First, merge repeat characters.

Then, remove any ϵ tokens.

The remaining characters are the output.

CTC: Input-Output Representation

Key idea: **blank token** supports silent stretches and letter repeats (e.g., “hello” vs “helo”)

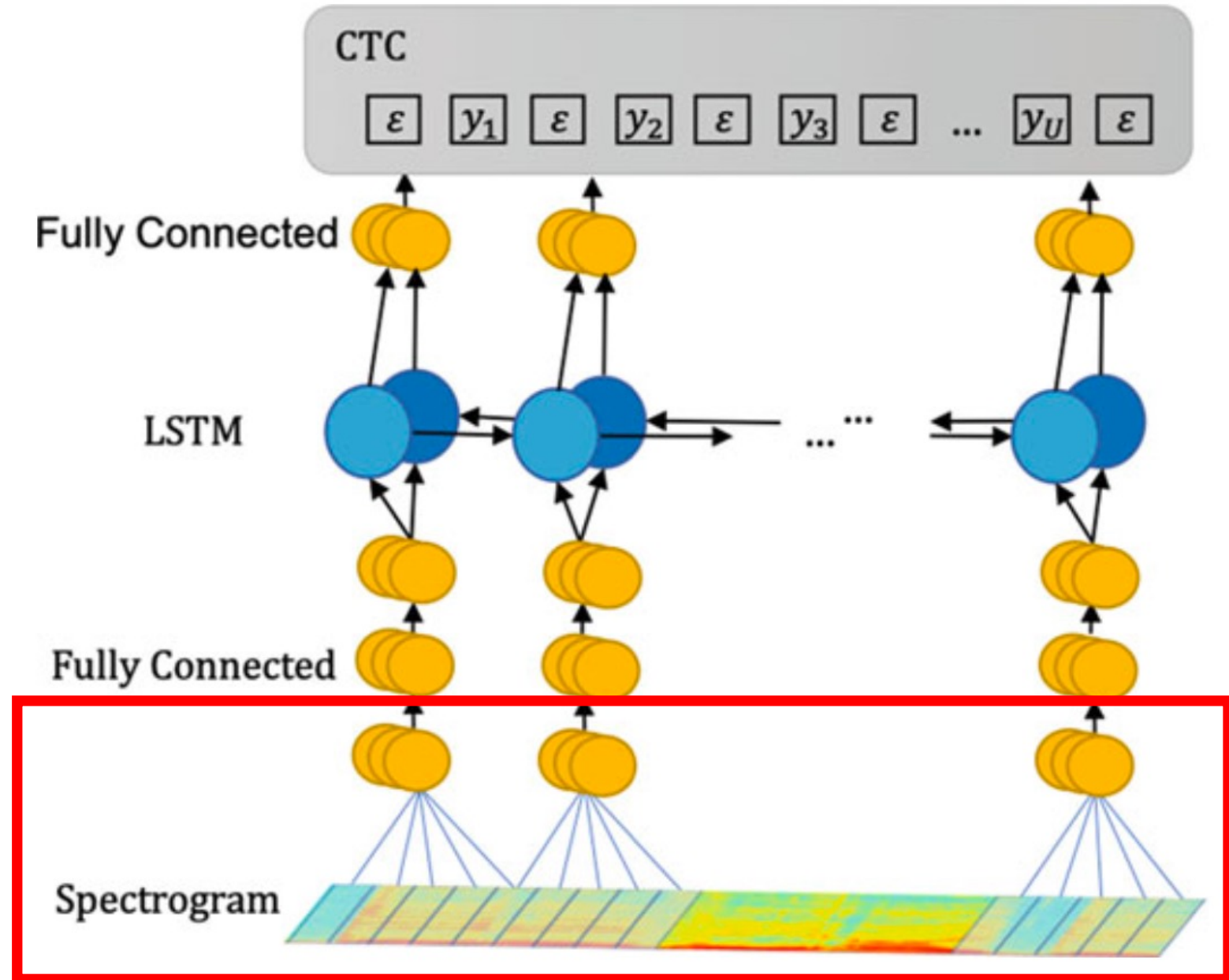
€ c c € a t

c c a a t t

c a € € € t

Supports recognizing the same word when spoken differently!

DeepSpeech



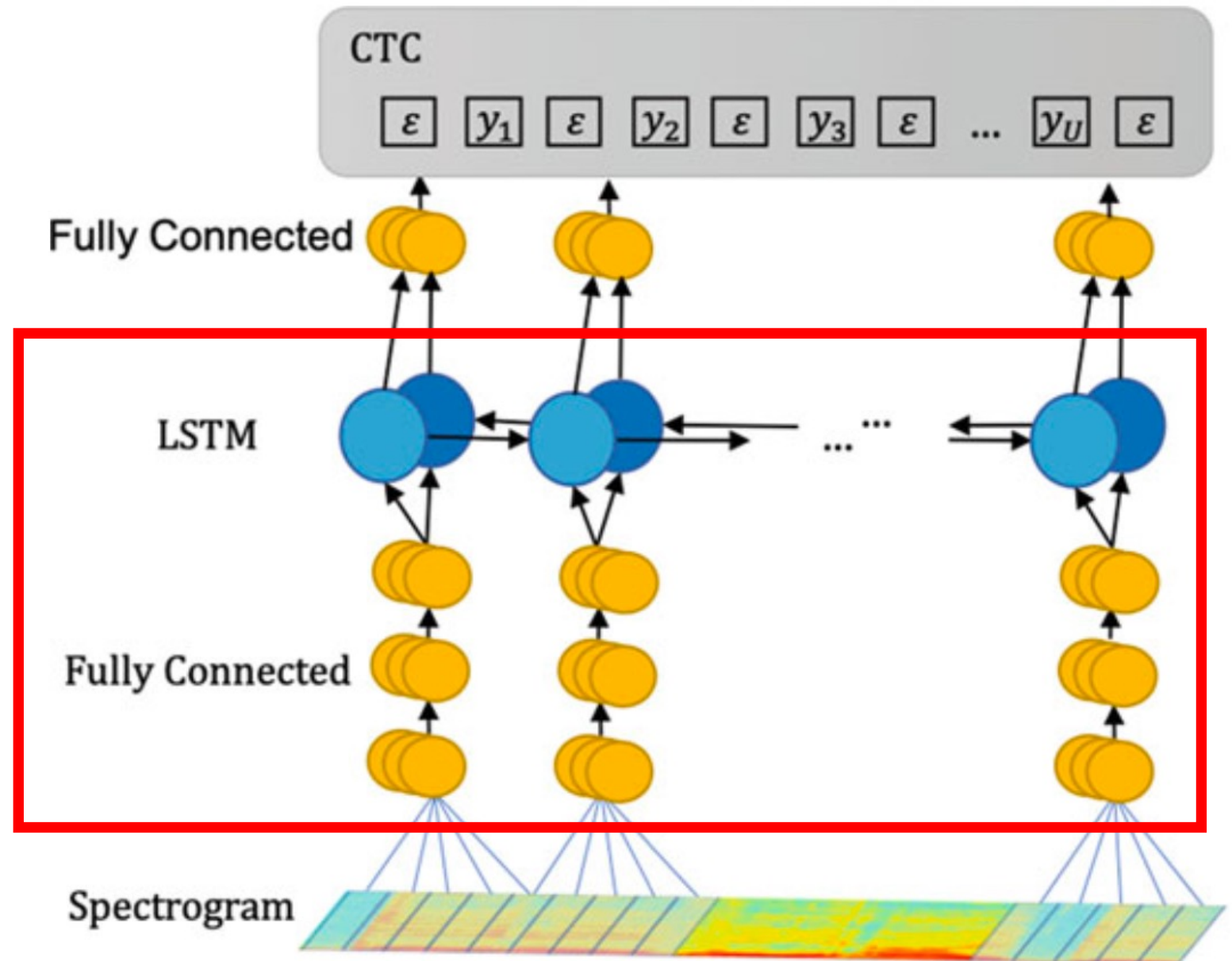
First hidden layer looks at context around input:

DeepSpeech

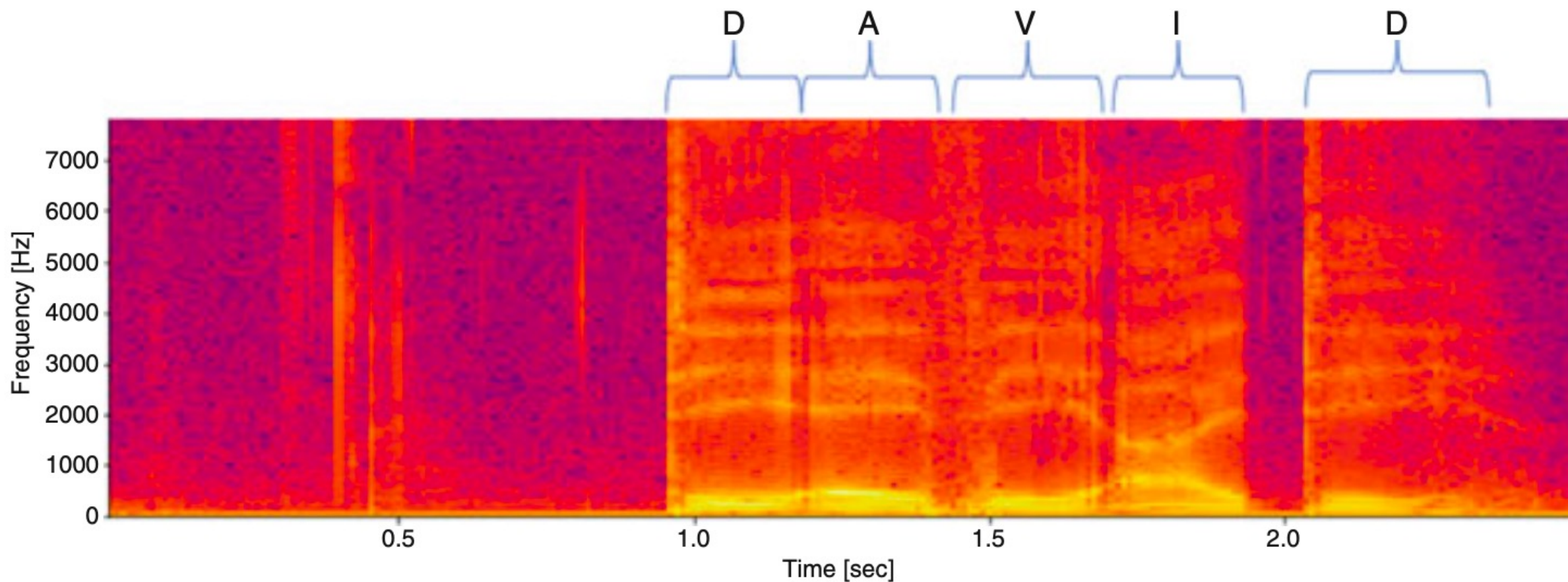
3 fully-connected layers
followed by bidirectional LSTM:

How is a bi-directional LSTM
beneficial?

How is a bi-directional LSTM
limiting?



DeepSpeech: Optimization Function (CTC)

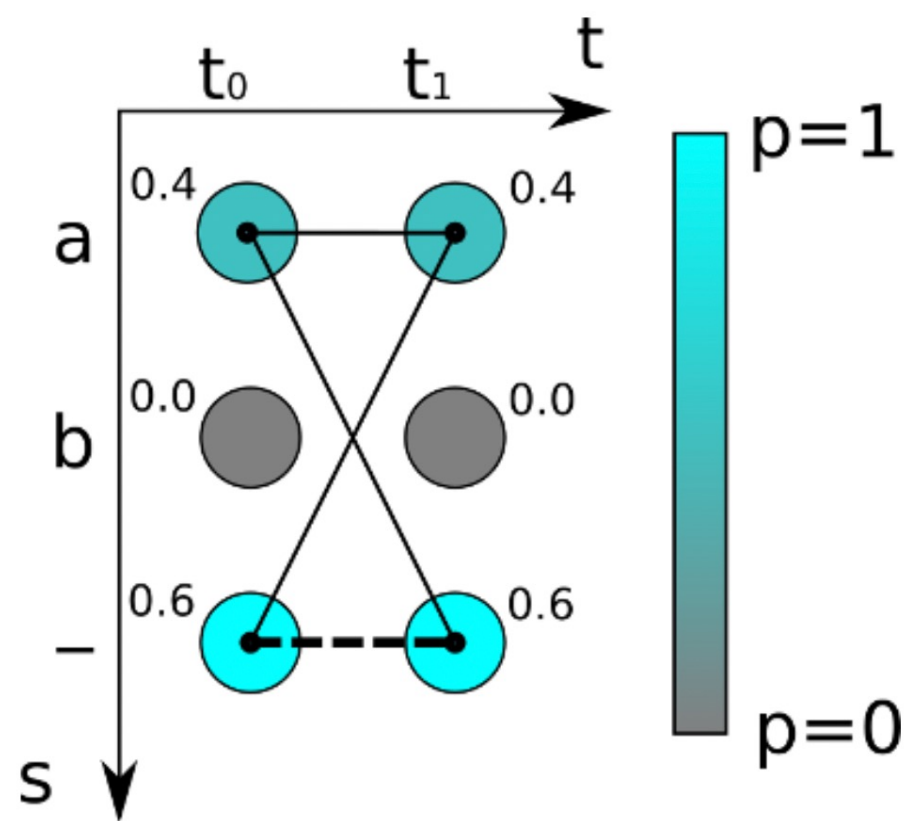


The CTC loss function enables learning output alignment without a per input label

DeepSpeech: Optimization Function (CTC)

Most plausible from all possible alignments learned;
e.g., 2 time steps with 2 potential characters and a
blank token ("-")

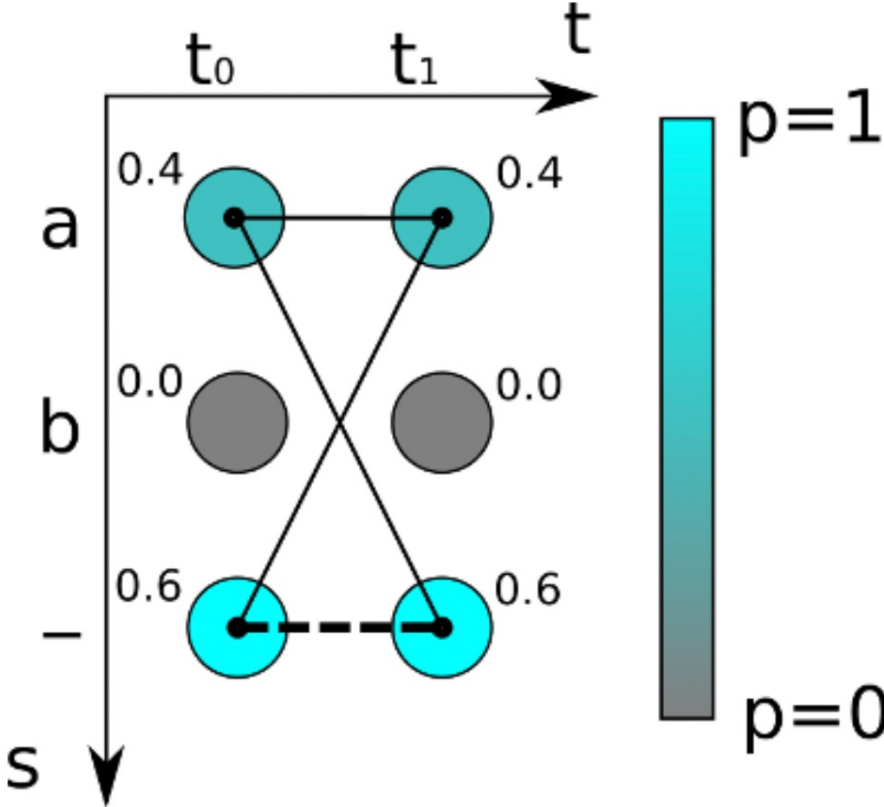
- Probability of "a" is sum of all "a" representations
 - Probability of "aa"?
 - $0.4 \times 0.4 = 0.16$
 - Probability of "a-"?
 - $0.4 \times 0.6 = 0.24$
 - Probability of "-a"?
 - $0.6 \times 0.4 = 0.24$
 - Sum: $0.16 + 0.24 + 0.24 = 0.64$



DeepSpeech: Optimization Function (CTC)

Most plausible from all possible alignments learned; e.g., 2 time steps with 2 potential characters and a blank token ("-")

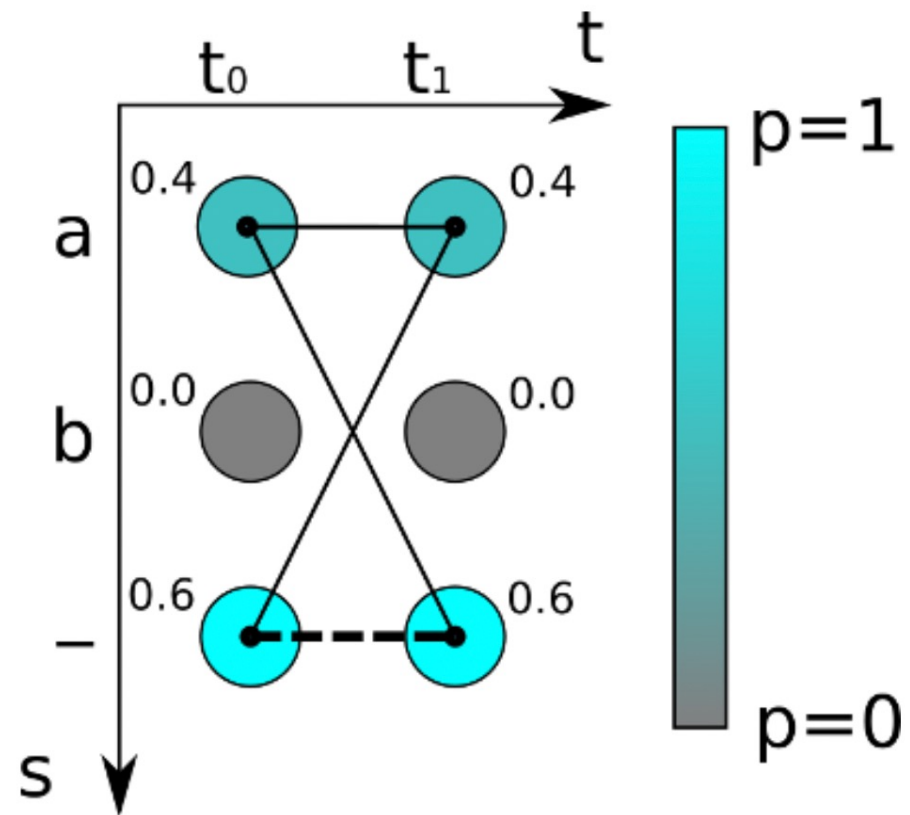
- Probability of "a": 0.64
- Probability of "" is sum of all "" representations
 - Probability of "--"?
 - $0.6 \times 0.6 = 0.36$



DeepSpeech: Optimization Function (CTC)

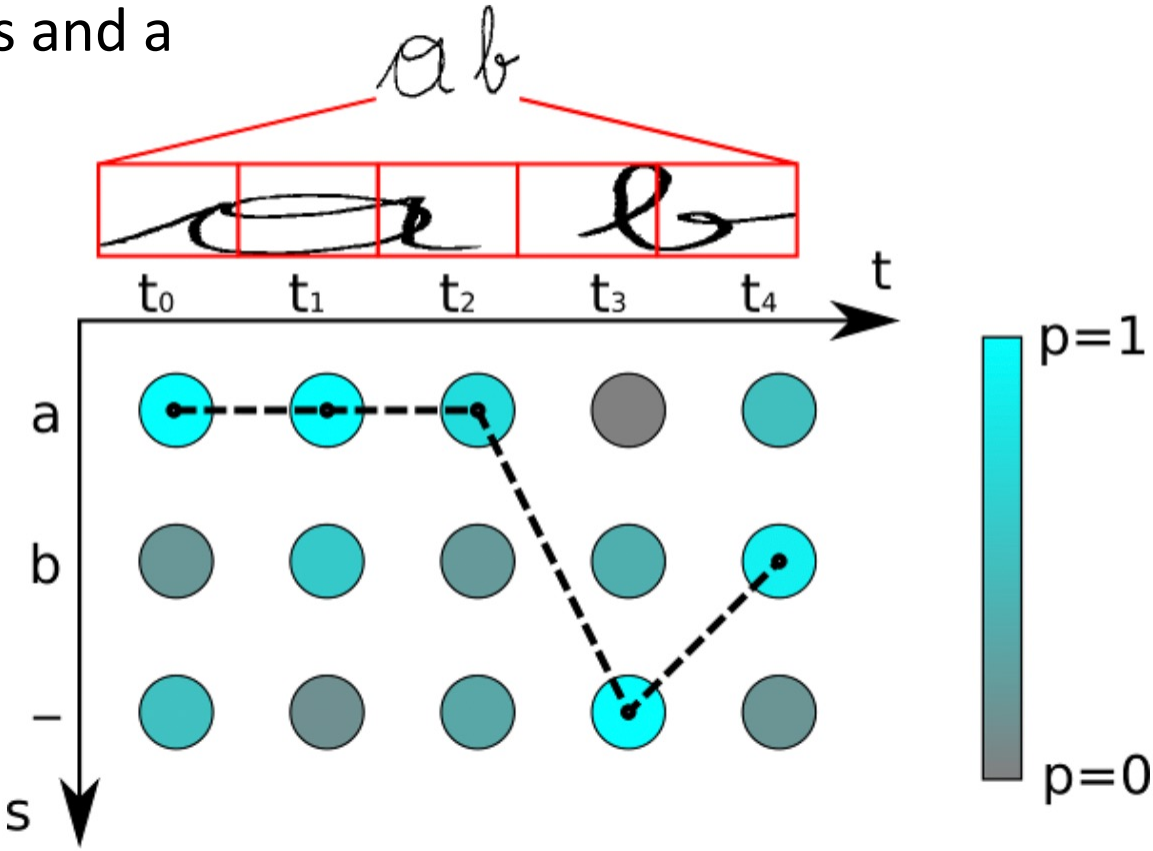
Most plausible from all possible alignments learned;
e.g., 2 time steps with 2 potential characters and a
blank token ("-")

- Probability of "a": 0.64
- Probability of "": 0.36
- And so on for all possible alignments...

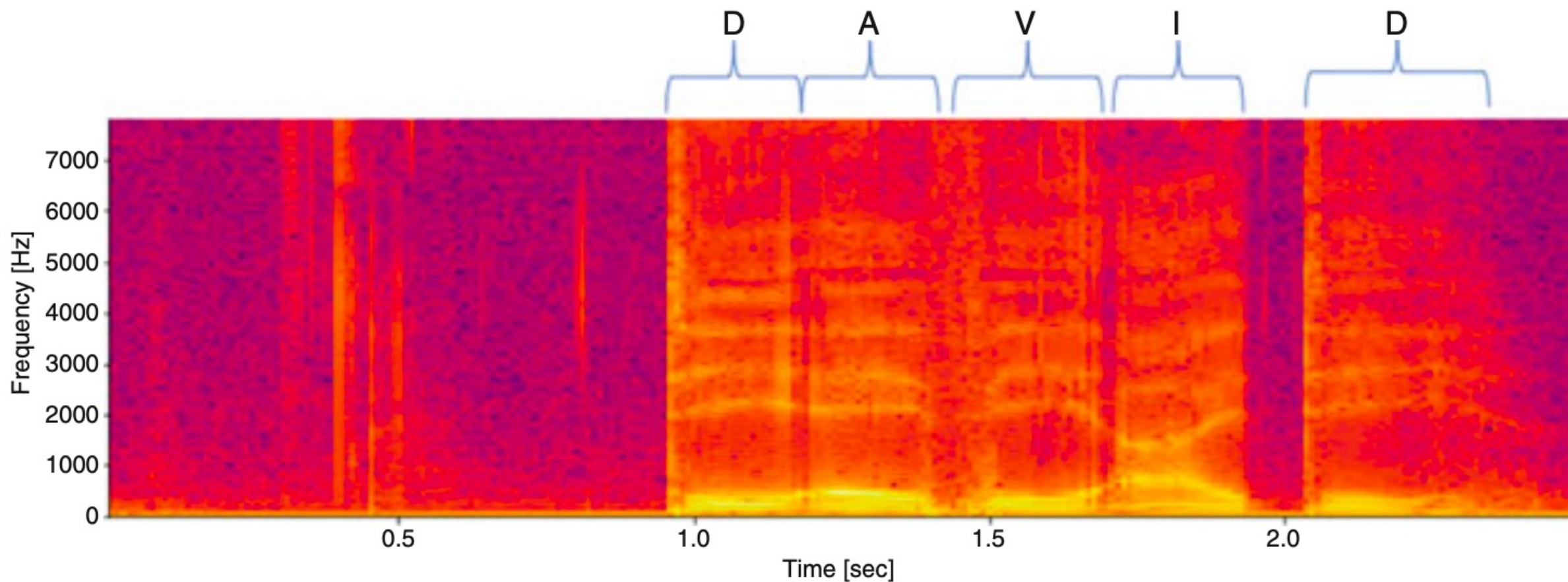


DeepSpeech: Optimization Function (CTC)

Most plausible from all possible alignments learned;
e.g., 2 time steps with 2 potential characters and a
blank token ("-") with "best path decoding"



DeepSpeech: Optimization Function (CTC)



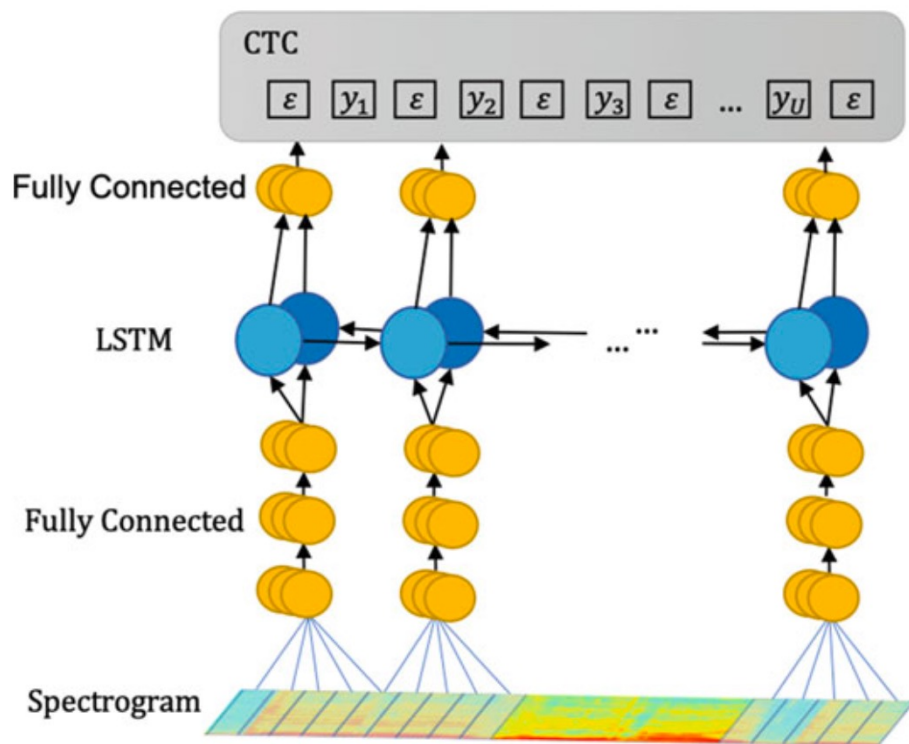
CTC uses dynamic programming to accelerate computation and is differentiable

DeepSpeech: Training (Key Ideas)

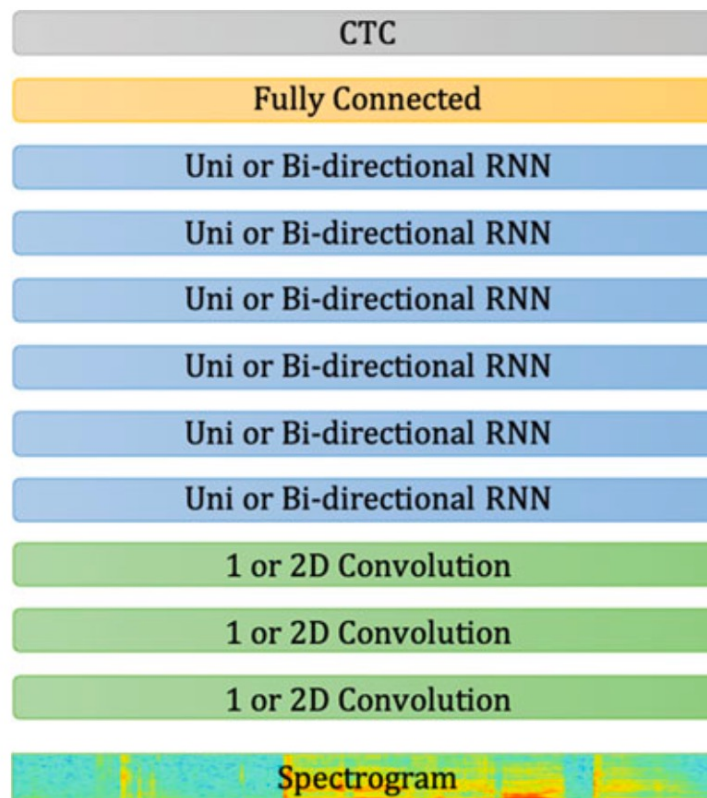
- 5000 hours from 9600 speakers
- Regularization
 - Dropout
 - Data augmentation: audio file translated 5 ms forward and backward
- Results boosted by incorporating a language model

Popular Methods

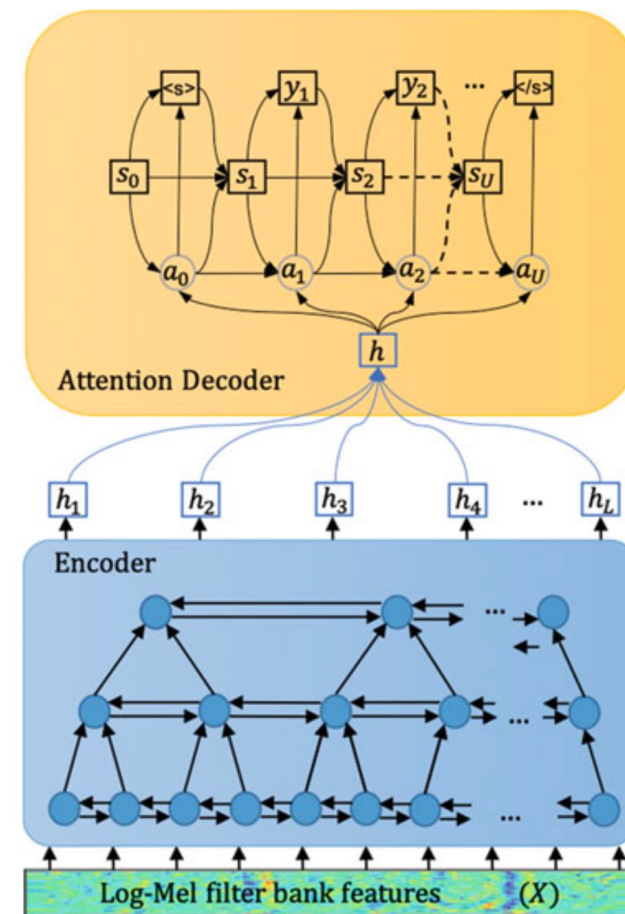
DeepSpeech



DeepSpeech2



Listen, Attend, and Spell

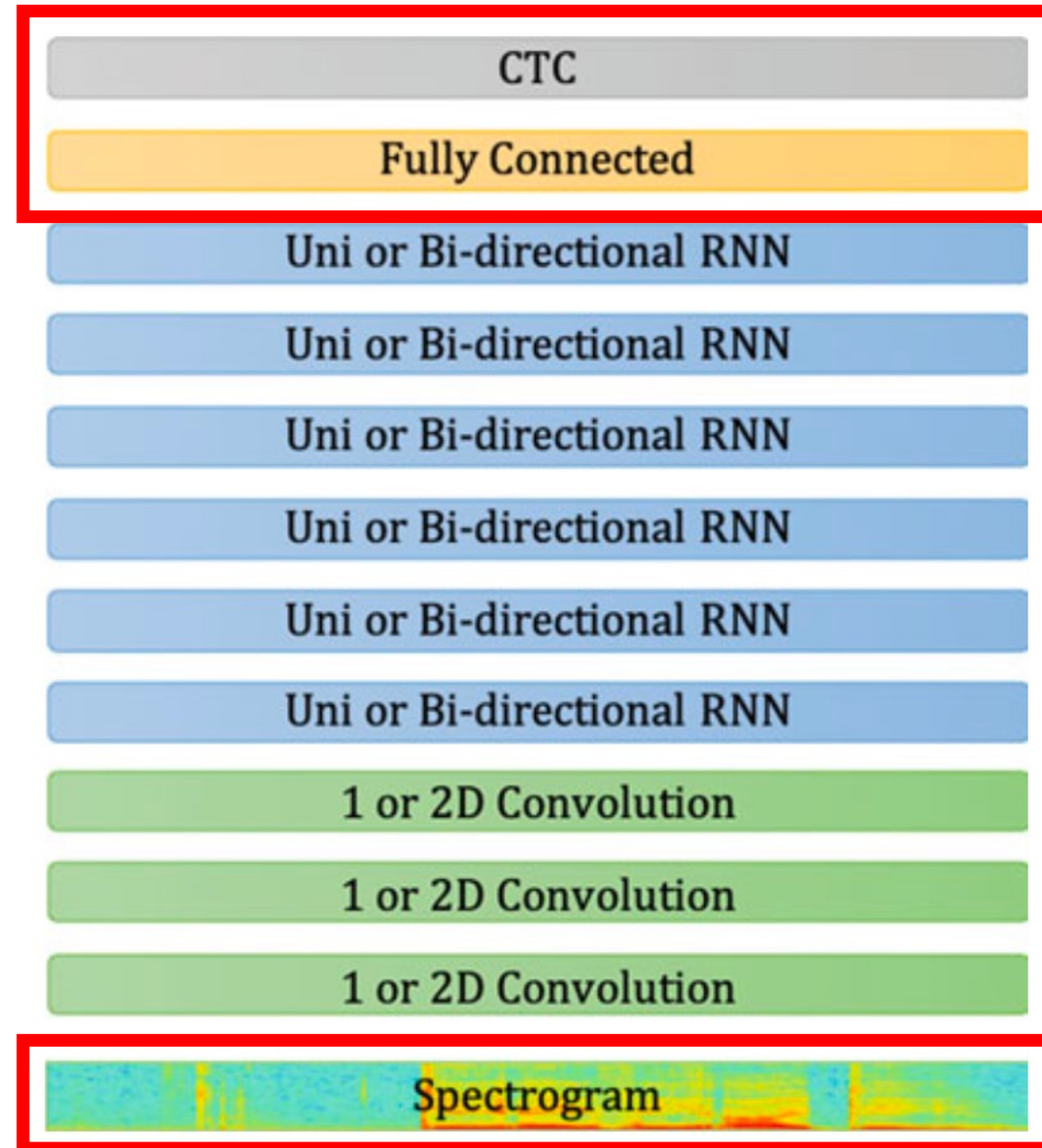


DeepSpeech2

Similar output:
(two architectures for
English and Mandarin)

Extension of DeepSpeech that
achieves a 7x speed-up and 43.4%
relative WER improvement with a
deeper architecture

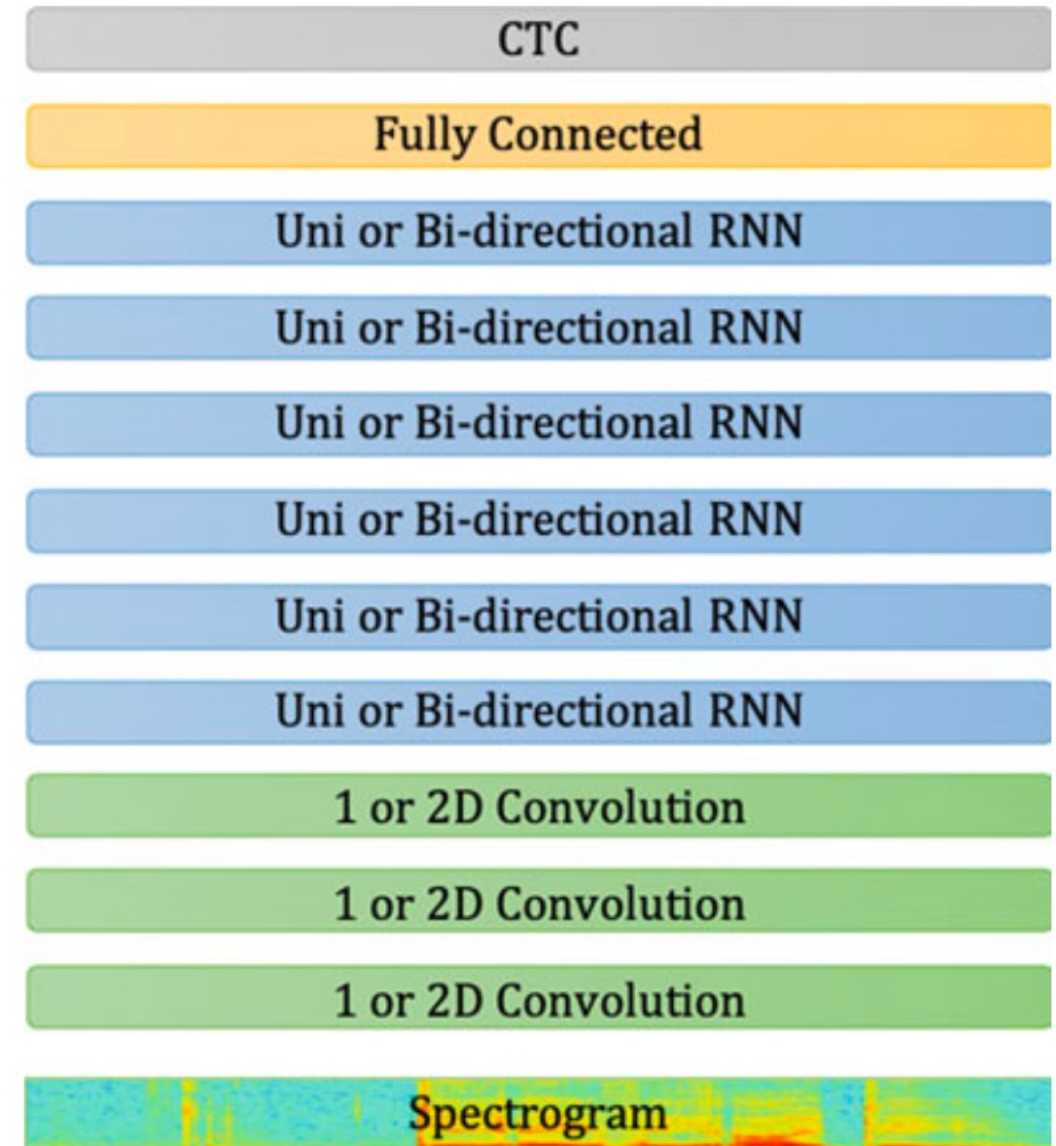
Same input:



DeepSpeech2

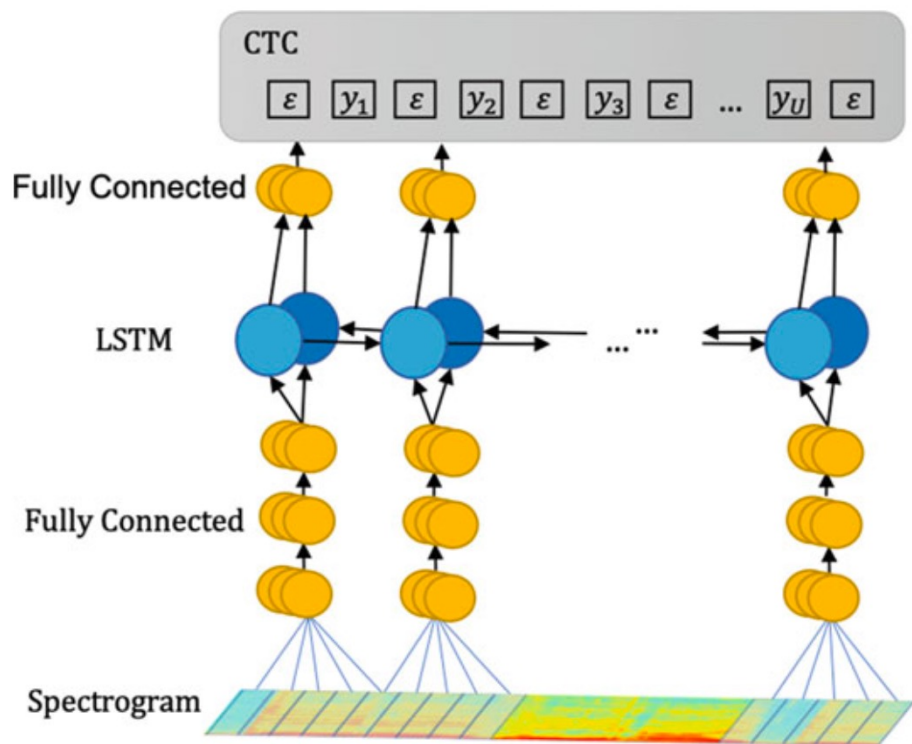
Training protocol difference from DeepSpeech:

- More training data (11,940 hours for English and 9,400 hours for Mandarin)
- Curriculum learning: trains based on length of utterances for first epoch with shorter ones first (improves WER by over 1 point)

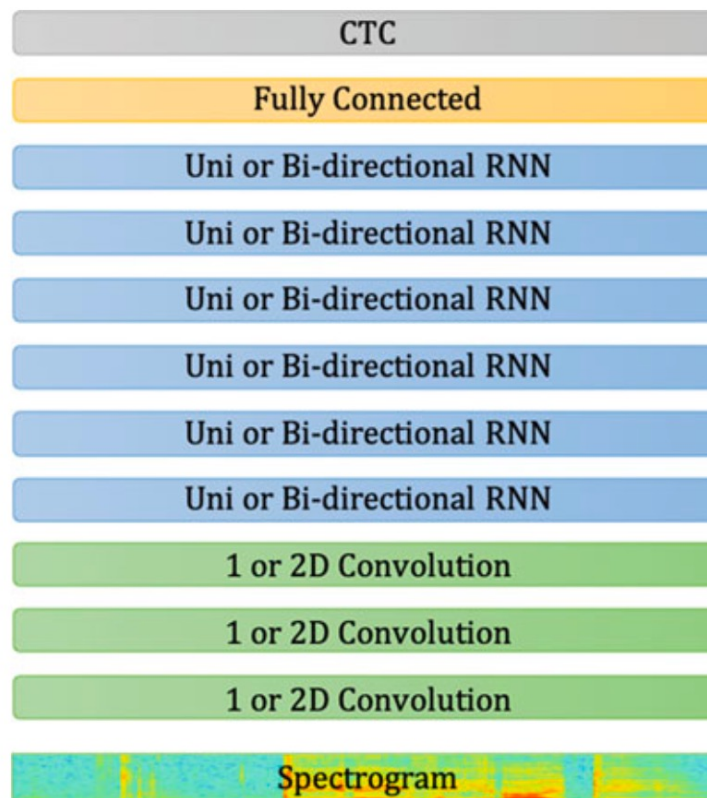


Popular Methods

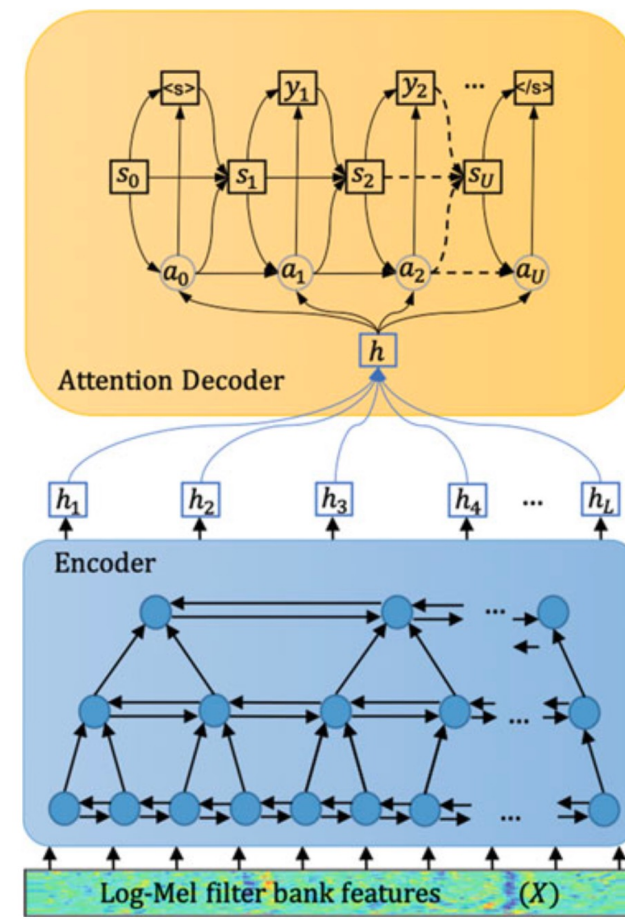
DeepSpeech



DeepSpeech2



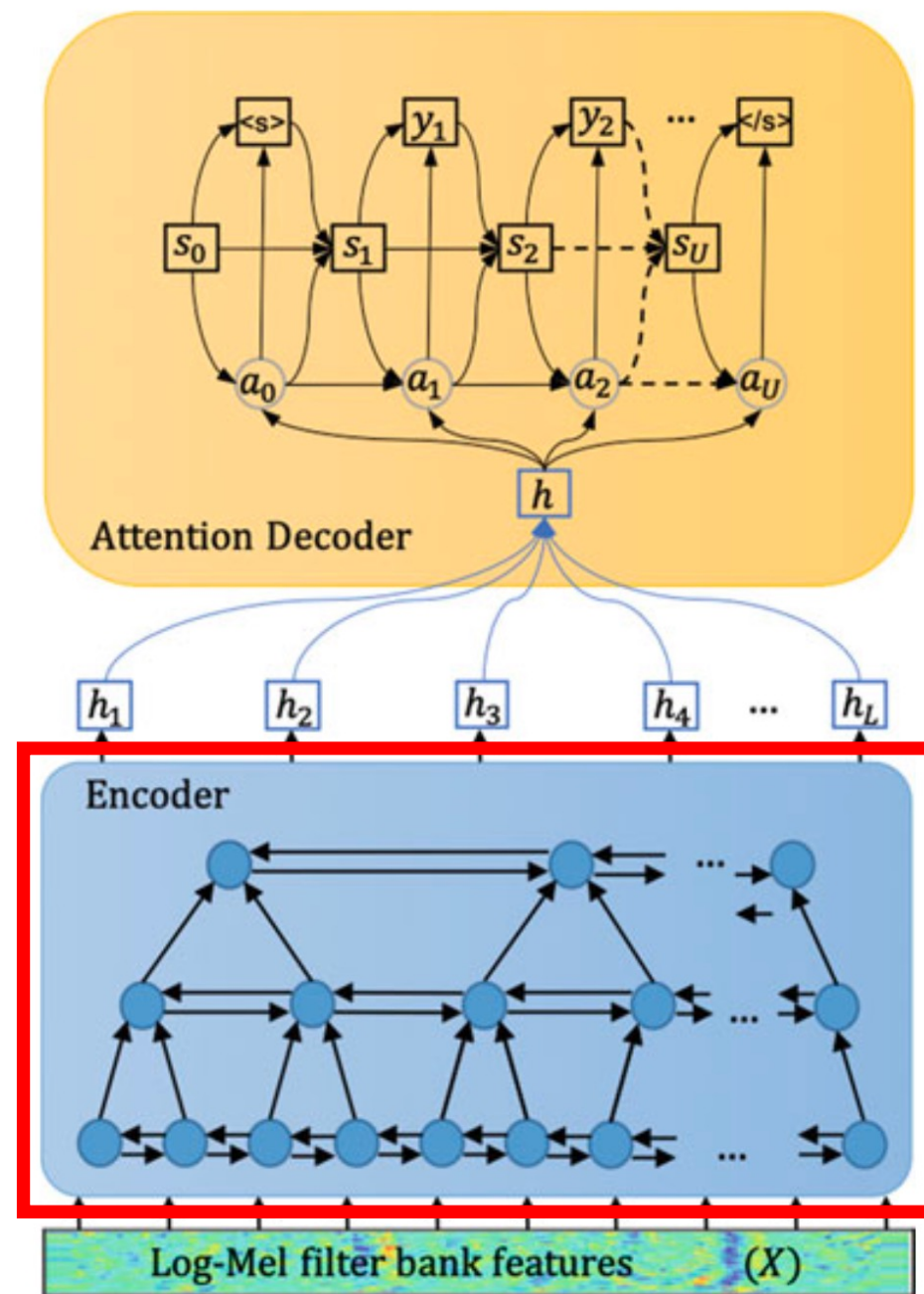
Listen, Attend, and Spell



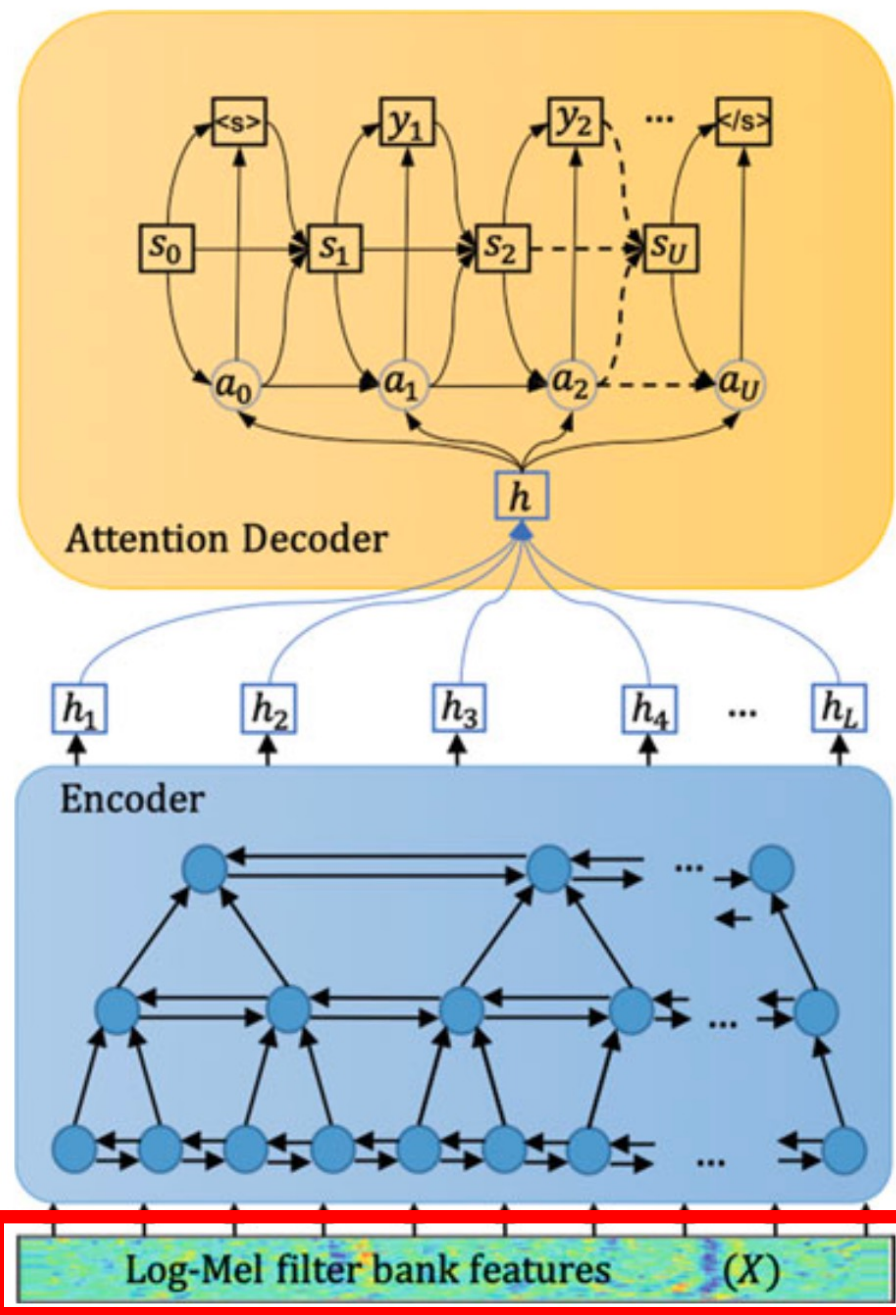
Listen, Attend, and Spell

Mimics original paper on sequence to sequence learning with attention where the decoder learns what to attend to in the encoded representation

Pyramid structure reduces number of input time steps



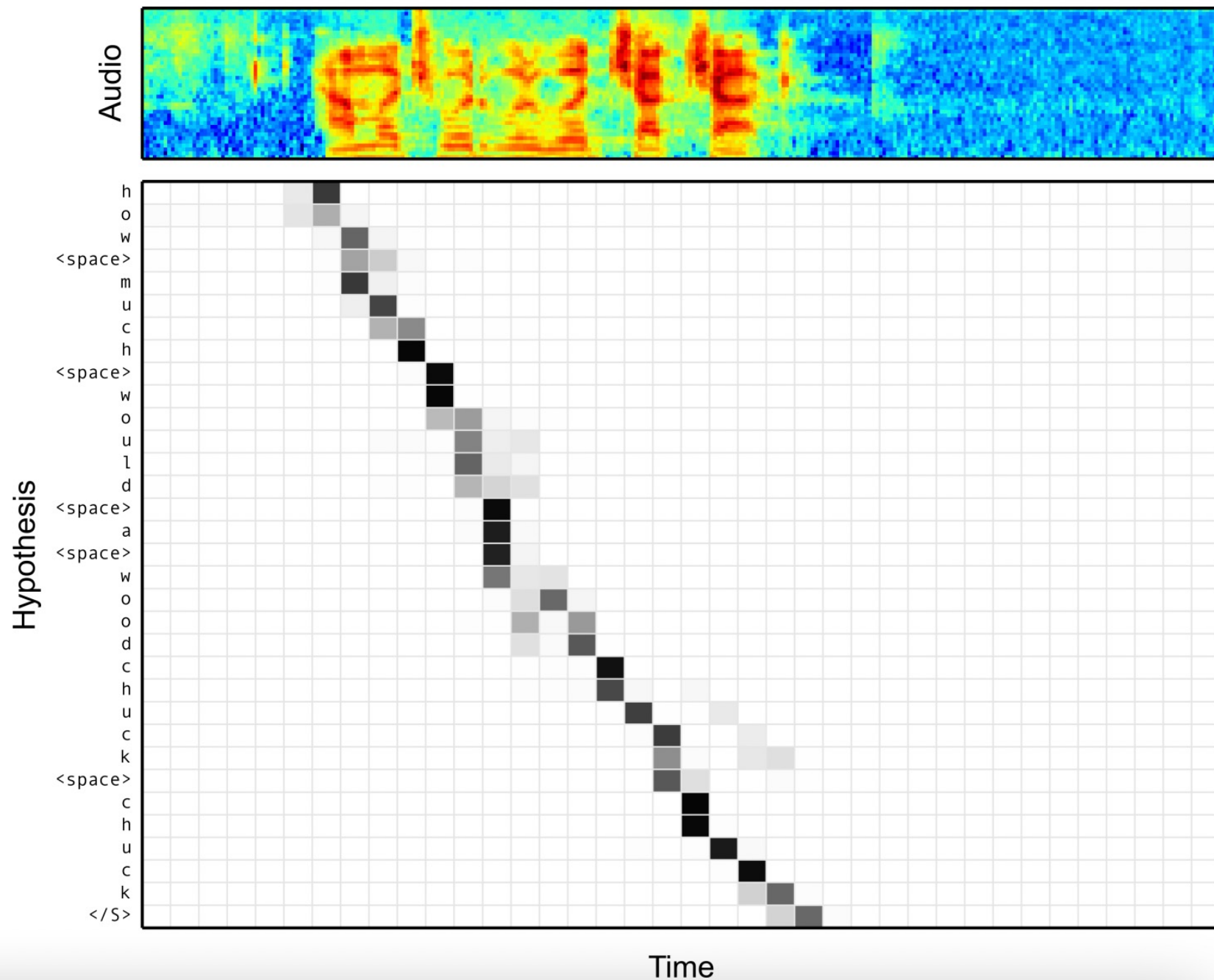
Listen, Attend, and Spell



Input: more sophisticated hand-crafted audio representation then spectrogram

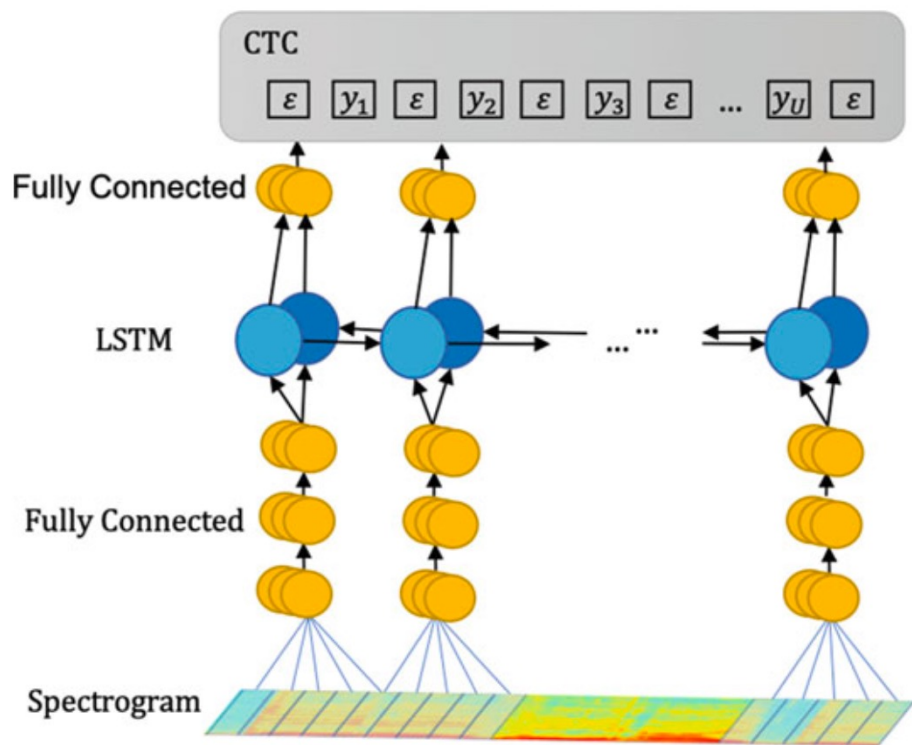
Result

Attention enables
visualizing alignment
between audio signal
and characters

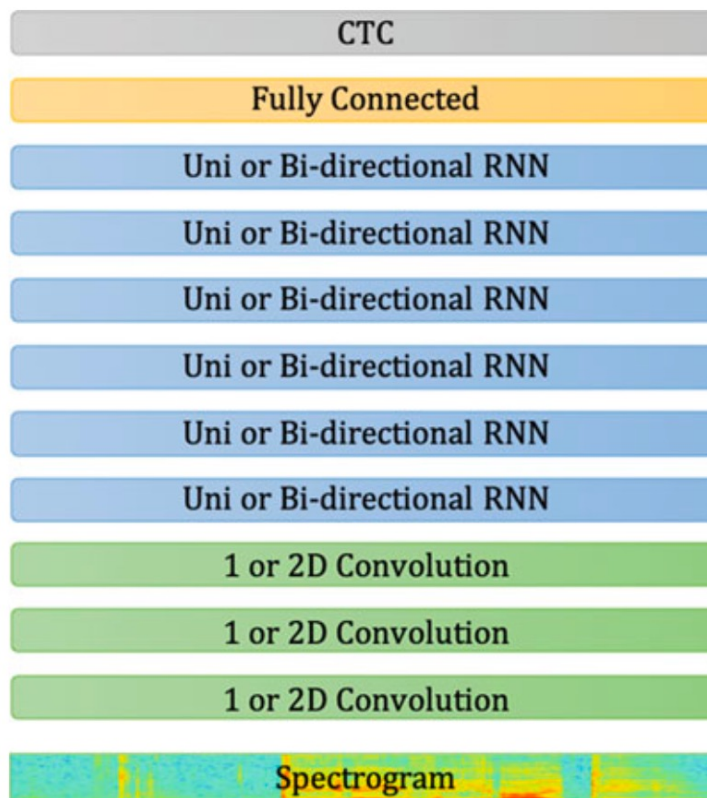


Popular Methods

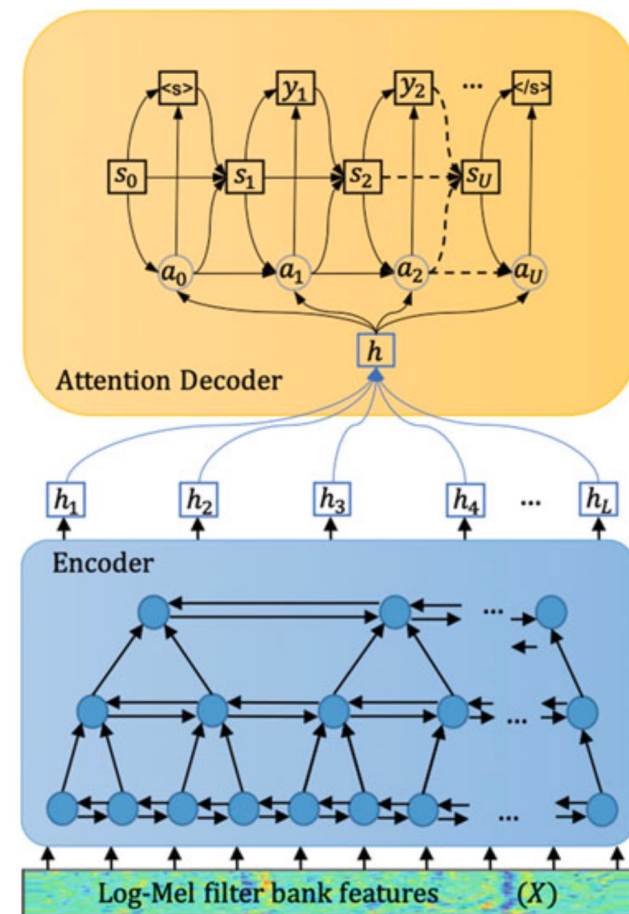
DeepSpeech



DeepSpeech2



Listen, Attend, and Spell



Today's Topics

- Problem
- Applications
- Speech recognition evaluation
- Speech recognition models
- Video making tutorial

Today's Topics

- Problem
- Applications
- Speech recognition evaluation
- Speech recognition models
- Video making tutorial

A dark gray background with a white film strip border on the left and right sides. The film strip has rectangular sprocket holes. In the center, there is a faint, glowing circular light effect. The text "The End" is written in a white, cursive script font with a slight drop shadow.

The End