

# Speech Recognition and Deep Learning

Instructor: Wei Xu  
Ohio State University

Many Slides from Adam Coates, Vinay Rao, Tony Han, Baidu Research

# Speech recognition

- Many exciting & valuable applications



Content captioning



Cars / Hands-free interfaces



Home devices

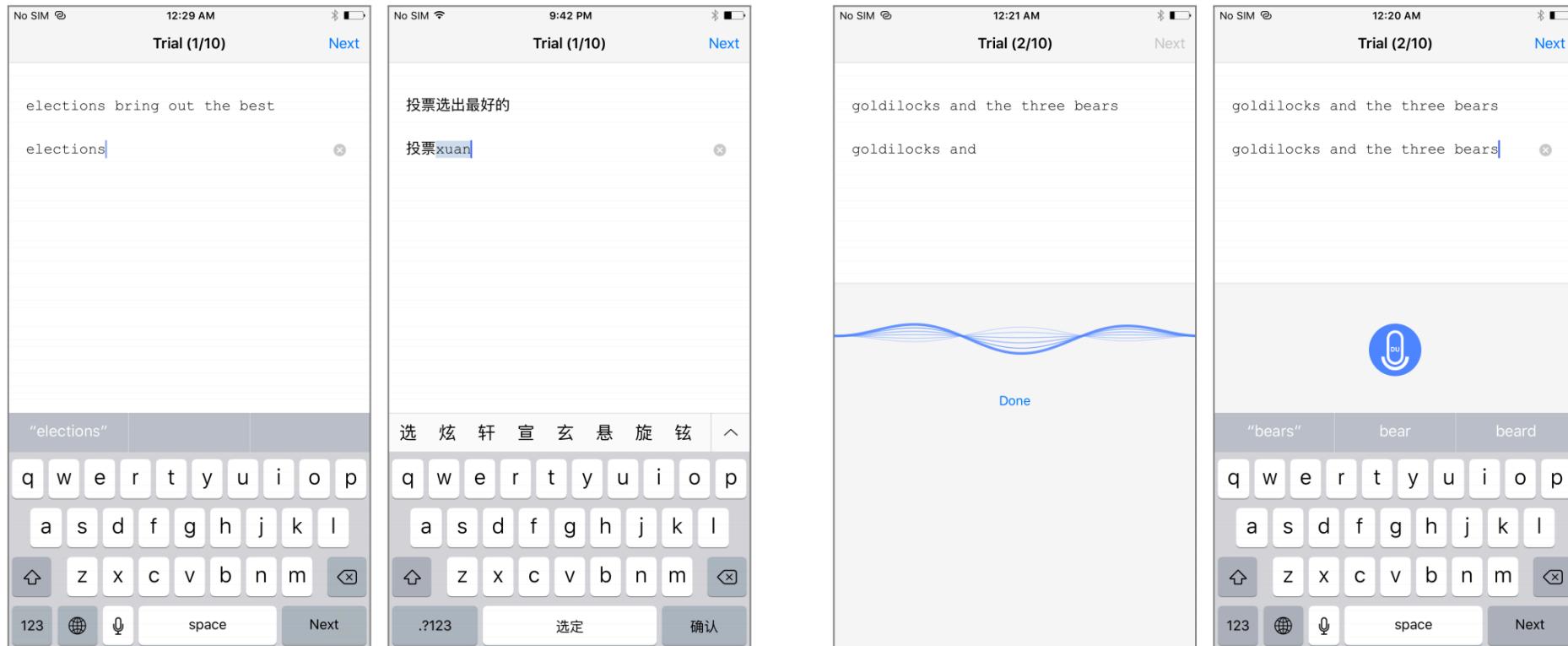


Mobile devices

# Speech recognition

- Many exciting & valuable applications

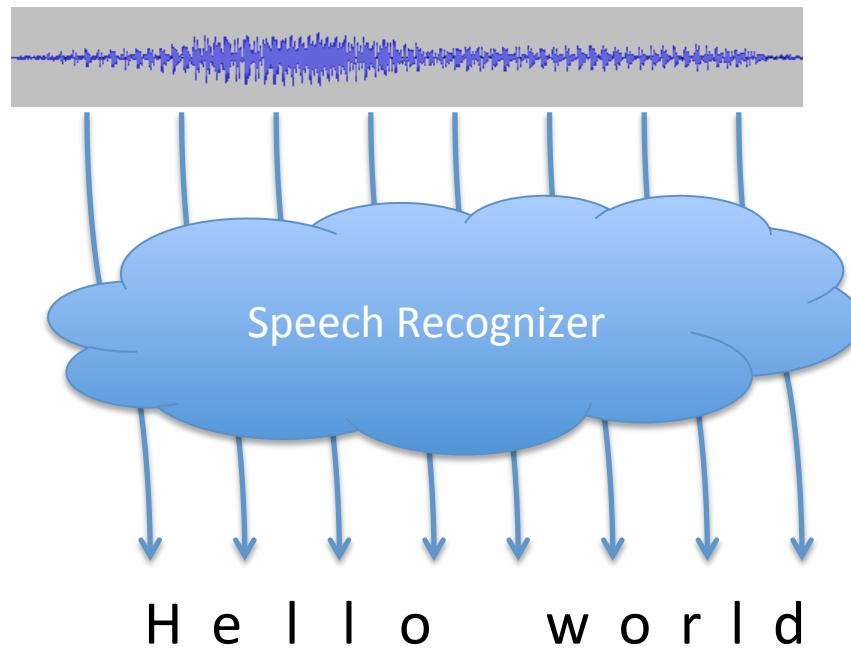
“Speech Is 3x Faster than Typing for English and Mandarin Text Entry on Mobile Devices”



[Ruan et al., 2016]

# Speech recognition

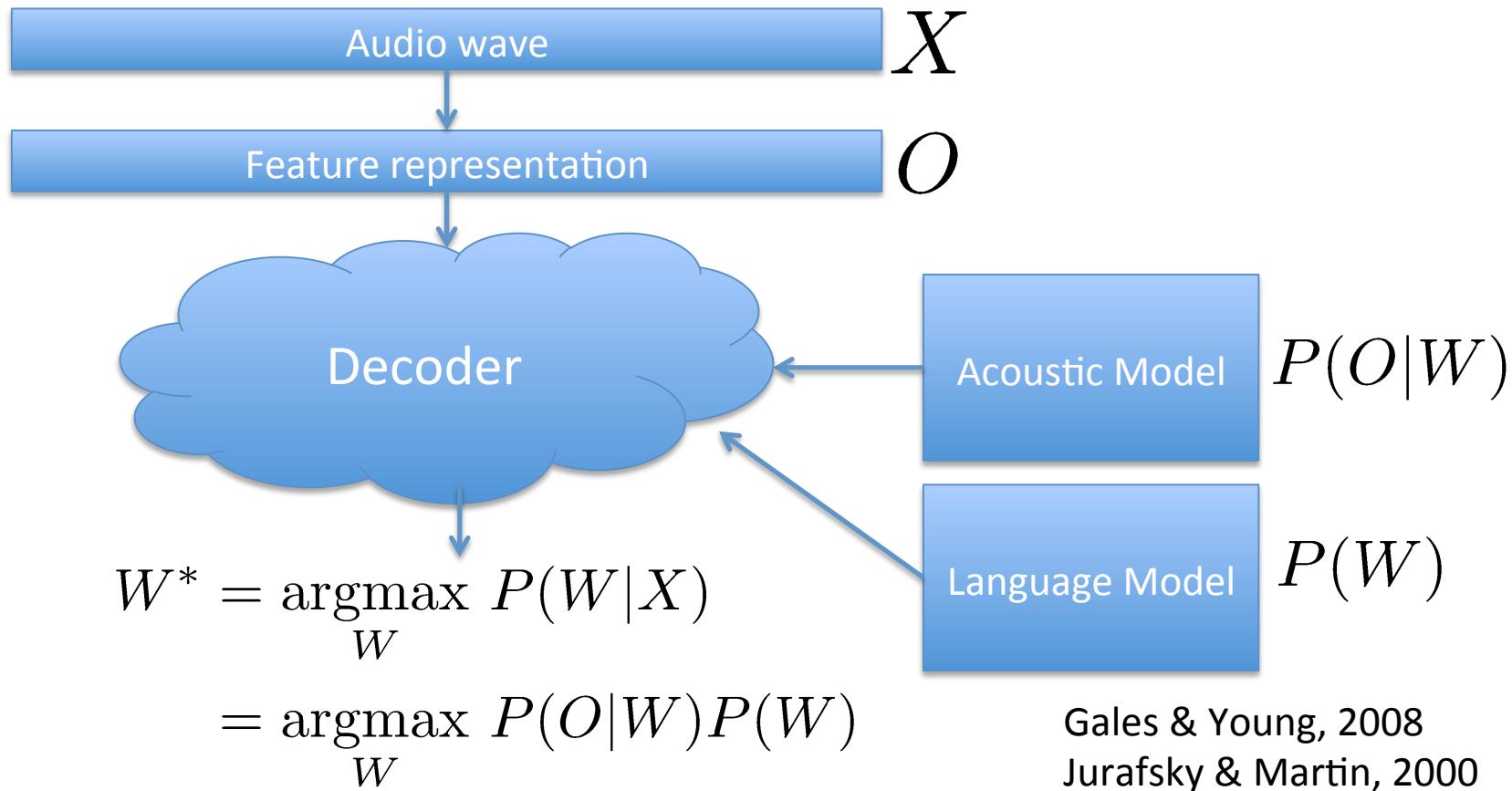
- Given speech audio, generate a transcript.



Important goal of AI: historically hard for machines, easy for people.

# Traditional ASR pipeline

- Traditional systems break problem into several key components:



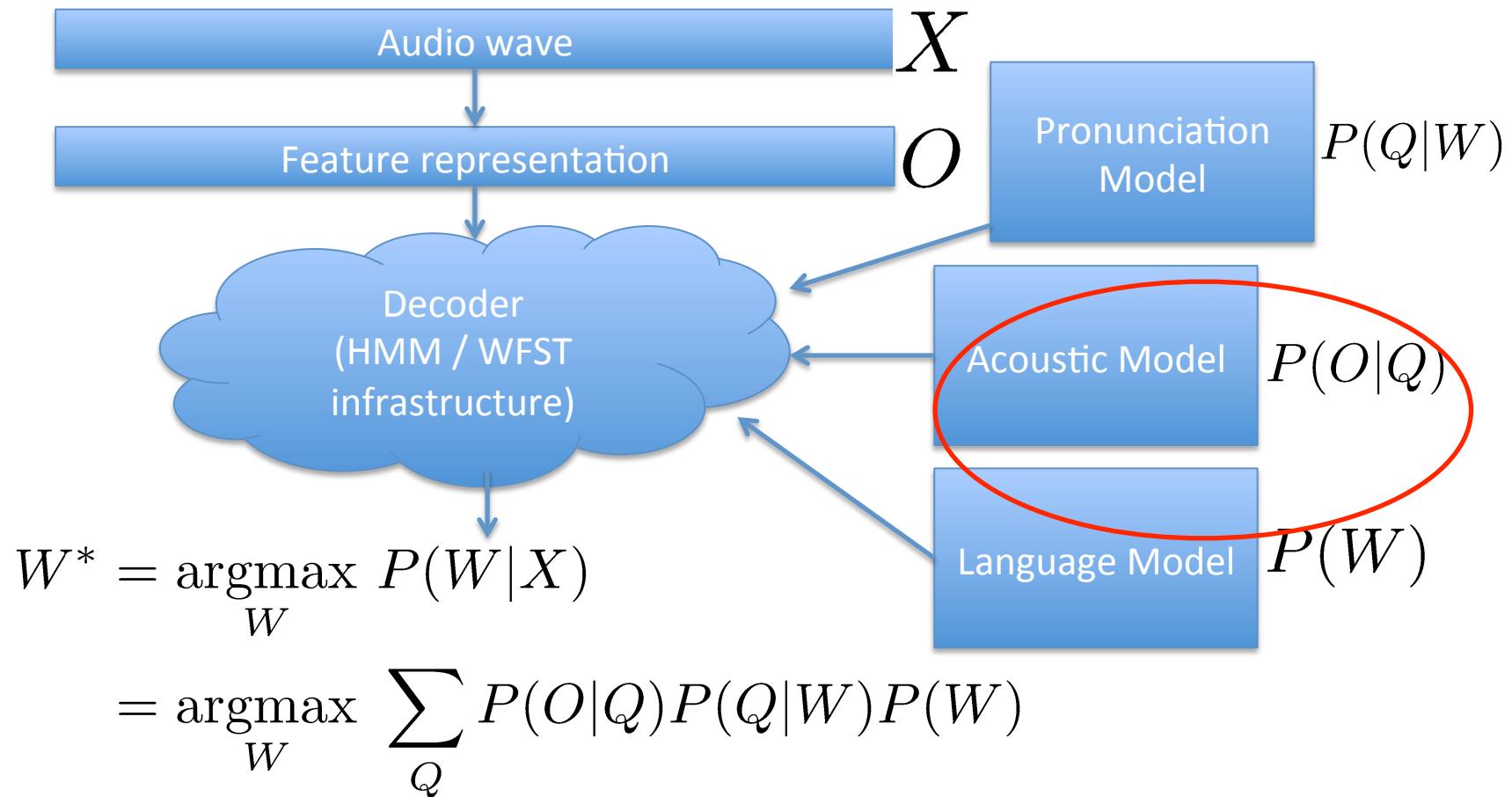
# Traditional ASR pipeline

- Usually represent words as sequence of “phonemes”:  
 $w_1 = \text{“hello”} = [\text{HH AH L OW}] = [q_1 q_2 q_3 q_4]$
- Phonemes are the perceptually distinct units of sound that distinguish words.
  - Quite approximate... but sorta standardized-ish.
  - Some labeled corpora available (e.g., TIMIT)

	Phone Label	Example		Phone Label	Example		Phone Label	Example
1	iy	beet	22	ch	choke	43	en	button
2	ih	bit	23	b	bee	44	eng	Washington

# Traditional ASR pipeline

- Traditional systems usually model phoneme sequences instead of words. This necessitates a dictionary or other model to translate.



# Traditional ASR pipeline

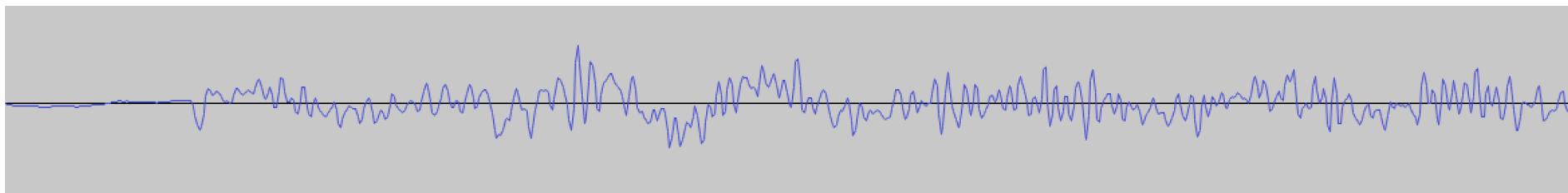
- Traditional pipeline is highly tweakable, but also hard to get working well.
- Historically, each part of system has own set of challenges.
  - E.g., choosing feature representation.

# Deep Learning in ASR

- Where to apply DL to make ASR better?
  - Good start: improve acoustic model  $P(O|Q)$

# Raw audio

- Simple 1D signal:



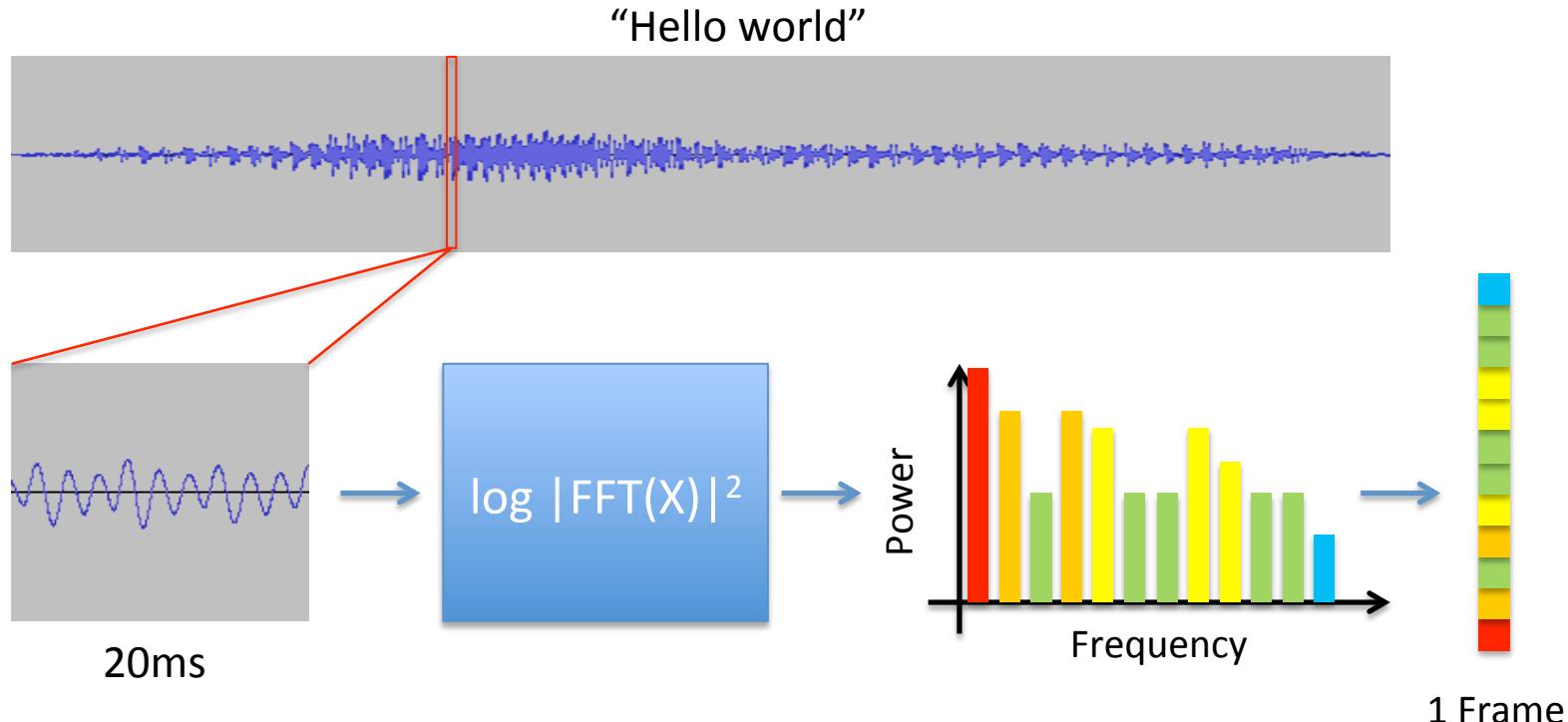
Typical sample rates for speech: 8KHz, 16KHz.

Each sample typically 8-bit or 16-bit.

- 1D vector:  $X = [x_1 x_2 \dots]$

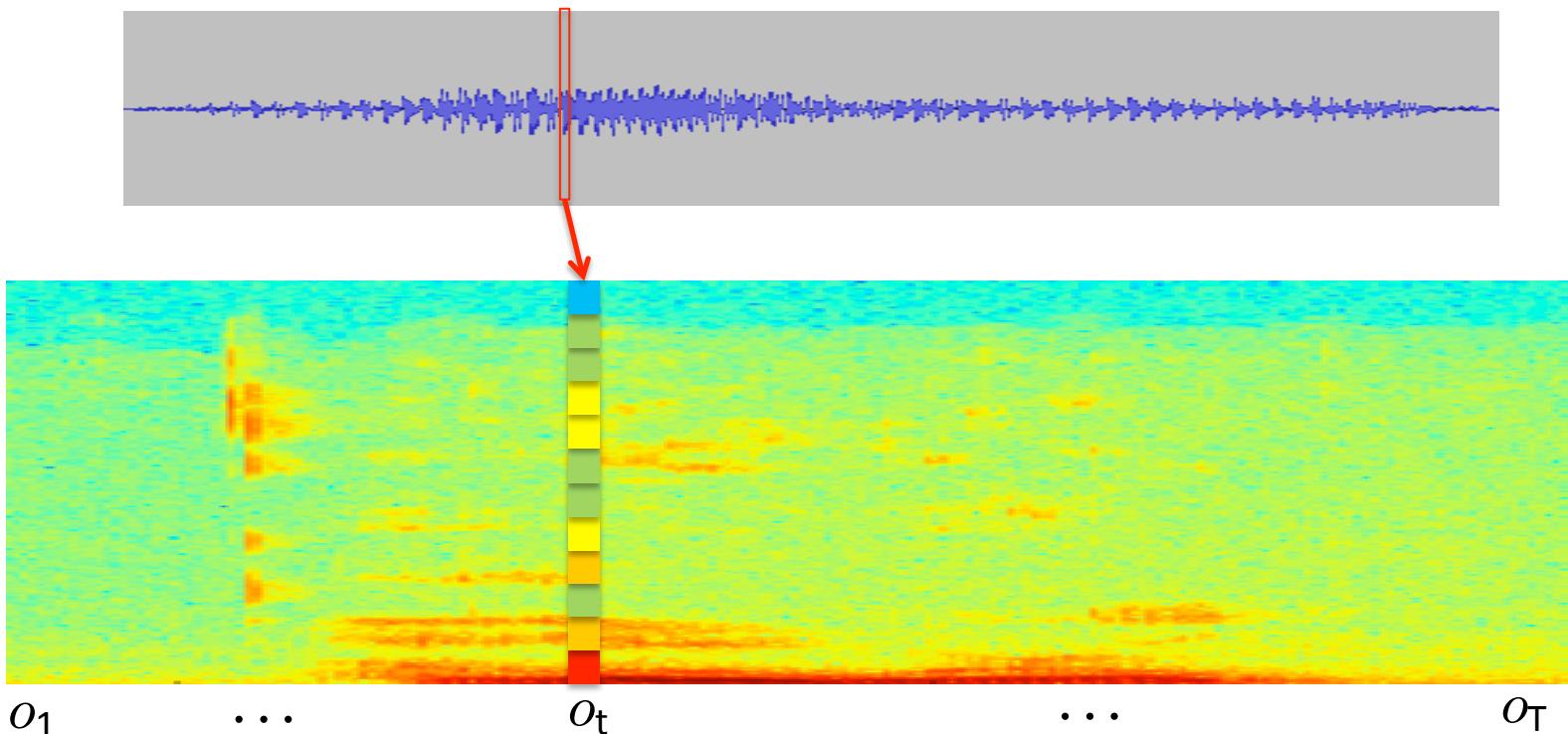
# Spectrogram

- Take a small window (e.g., 20ms) of waveform.
  - Compute FFT and take magnitude. (i.e., power)
  - Describes frequency content in local window.



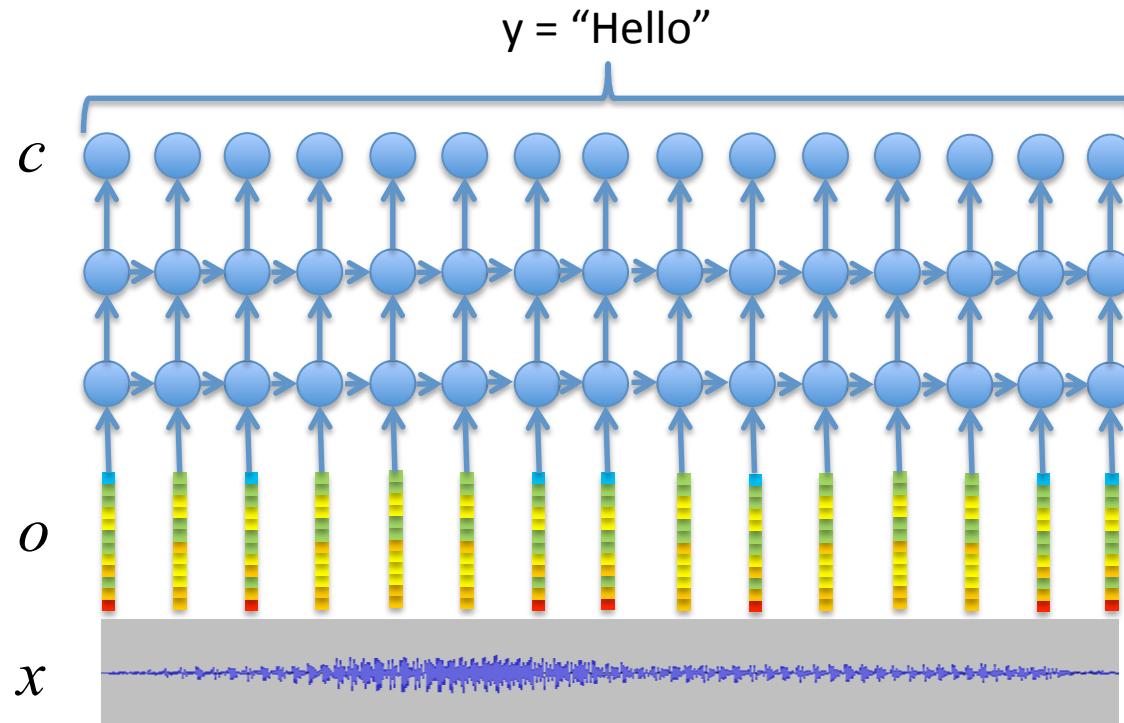
# Spectrogram

- Concatenate frames from adjacent windows to form “spectrogram”.



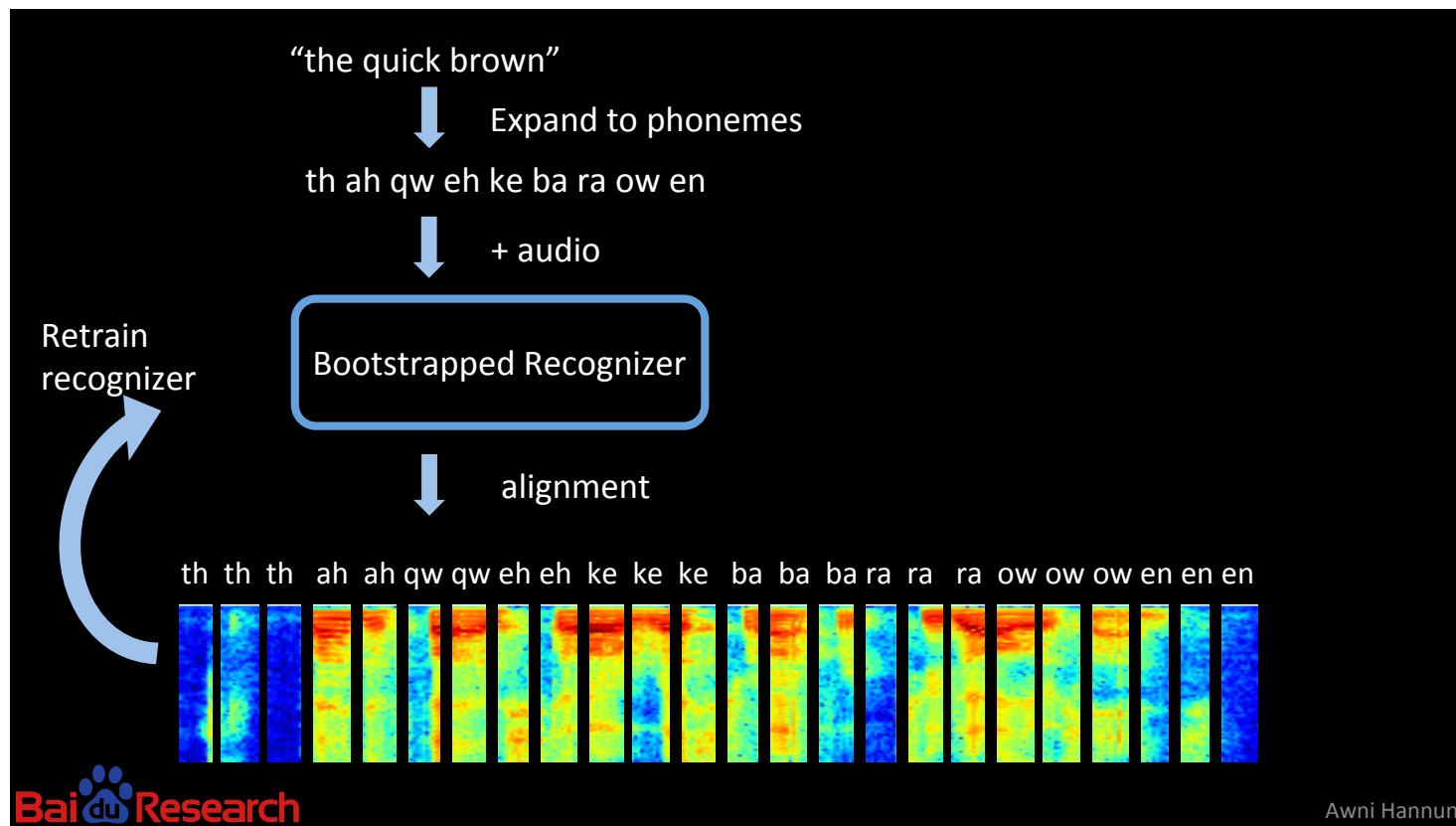
# Acoustic Model

- Goal: create a neural network (DNN/RNN) from which we can extract transcription,  $y$ .
  - Train from labeled pairs  $(x, y^*)$



# Acoustic model

- Main issue:  $\text{length}(x) \neq \text{length}(y)$ 
  - Don't know how symbols in  $y$  map to frames of audio.
  - Traditionally, try to bootstrap alignment – painful!



# DL for End-to-end Speech

- No perceptual features (MFCC). No feature transformation. No phonetic inventory. No transcription dictionary. No HMM.
- The output of the RNN are characters including space, apostrophe, (not CD phones)
- Connectionist Temporal Classification (No fixed alignment speech/character)
- Data augmentation. 5,000 hours (9600 speakers) + noise = 100,000 hours. Optimizations: data parallelism, model parallelism
- Good results in noisy conditions

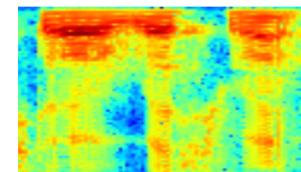
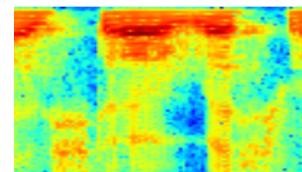
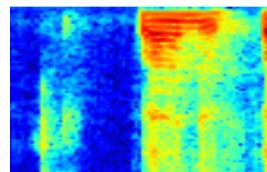
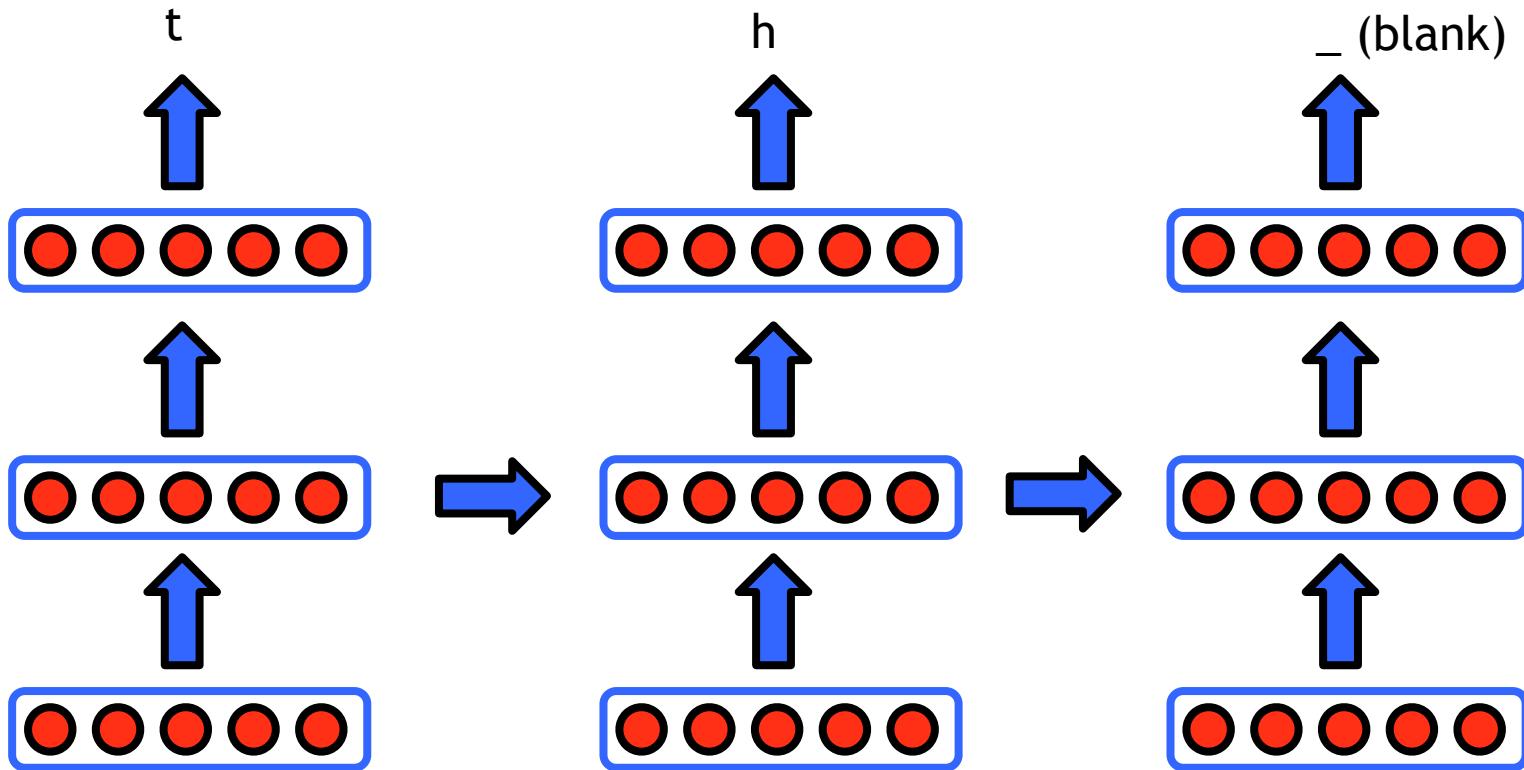
# Connectionist Temporal Classification (CTC)

[Graves et al., 2006]

- Basic idea:
  1. RNN output neurons  $c$  encode distribution over symbols. Note  $\text{length}(c) == \text{length}(x)$ .  
For phoneme-based model:  $c \in \{AA, AE, AX, \dots, ER1, \text{blank}\}$   
For grapheme-based model:  $c \in \{A, B, C, D, \dots, Z, \text{blank}, \text{space}\}$
  2. Define a mapping  $\beta(c) \rightarrow y$ .
  3. Maximize likelihood of  $y^*$  under this model.

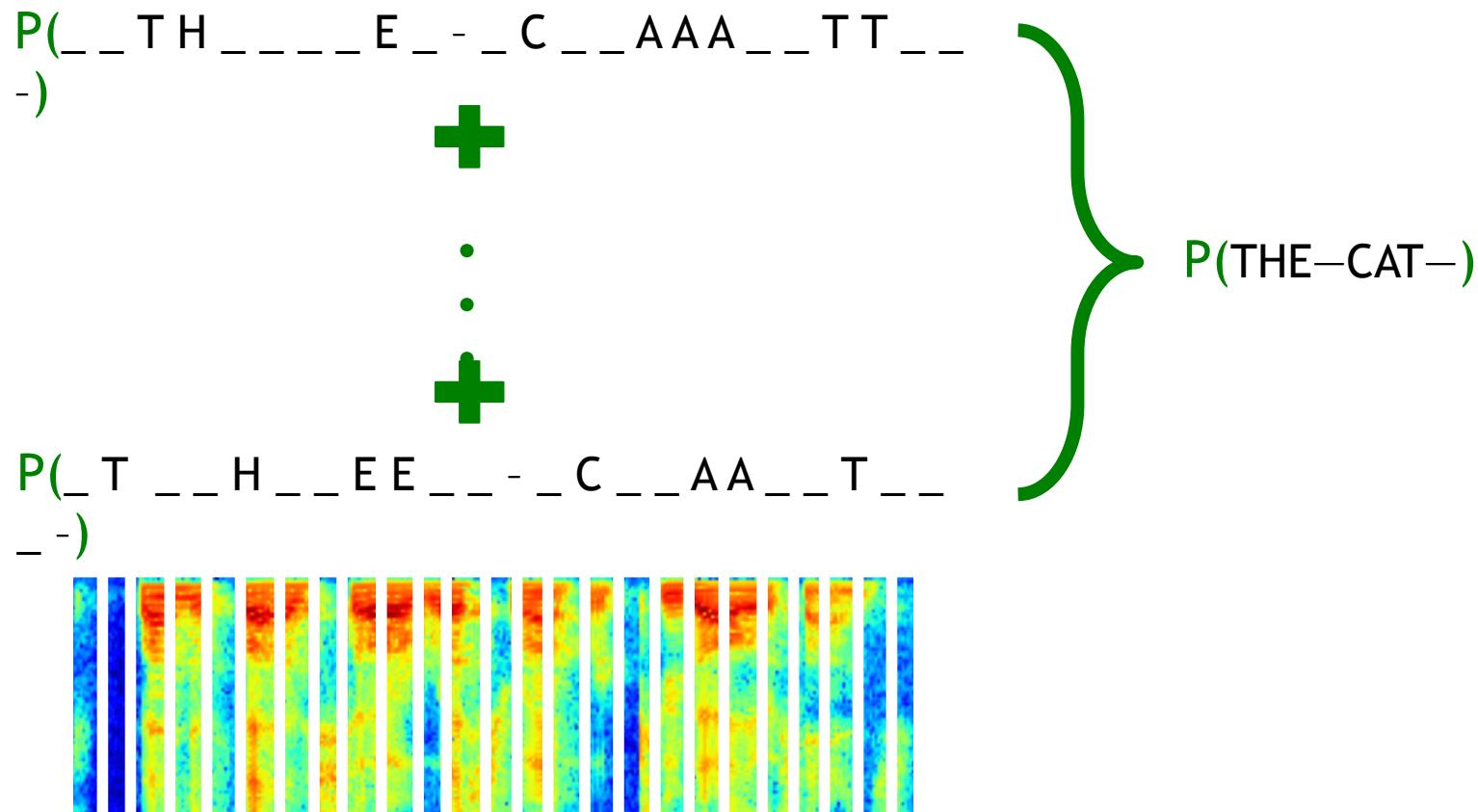
# Deep Speech - Recurrent Neural Network

Output  
alphabet,  
space,  
& blank



# Deep Speech - CTC

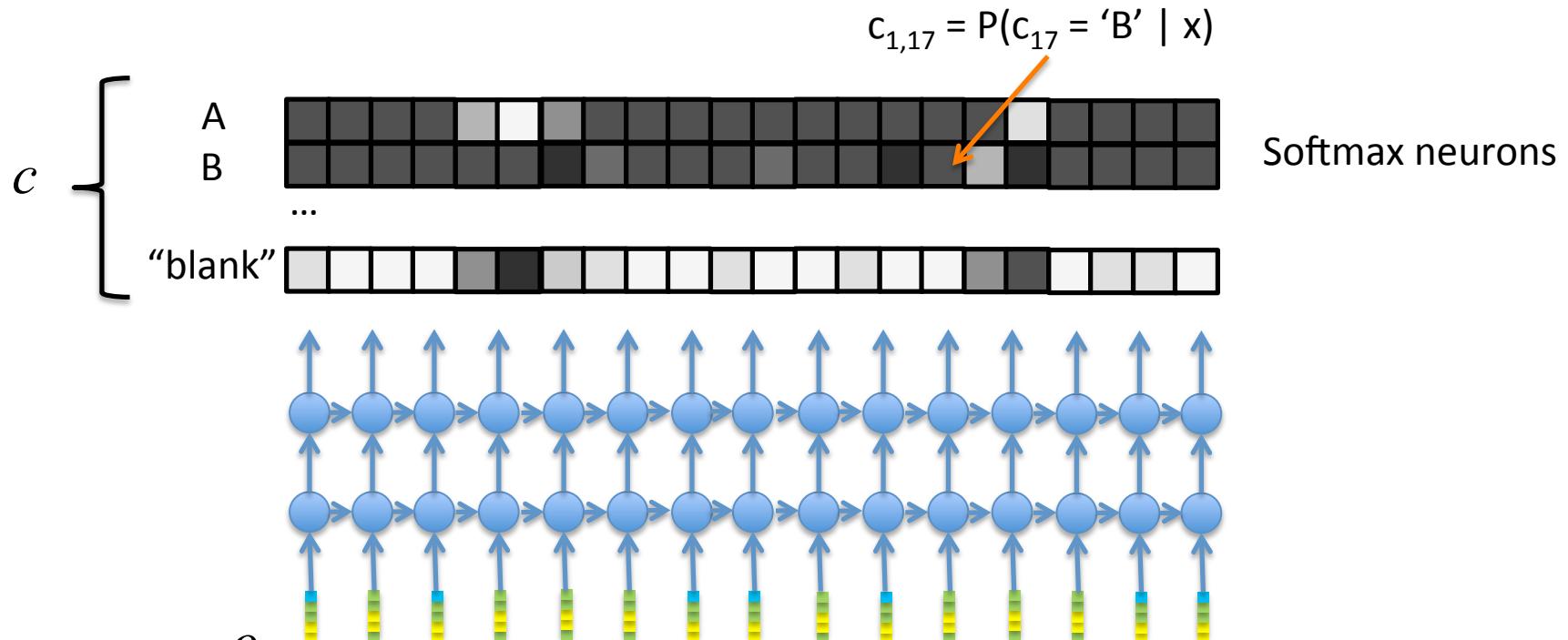
No alignment needed!



# Connectionist Temporal Classification (CTC)

1. RNN output neurons  $c$  encode distribution over symbols. Note  $\text{length}(c) == \text{length}(x)$ .

For grapheme-based model:  $c \in \{A, B, C, D, \dots, Z, \text{blank}, \text{space}\}$



# Connectionist Temporal Classification (CTC)

2. Define a mapping  $\beta(c) \rightarrow y$ .
  - Given a specific character sequence  $c$ , squeeze out duplicates + blanks to yield transcription:

$$y = \beta(c) = \beta(\text{H}\text{H}\text{H\_E\_\_LL\_LO\_\_}) = \text{"HELLO"}$$

# Connectionist Temporal Classification (CTC)

- Mapping implies a distribution over possible *transcriptions*  $y$ :

$P(c x) = \{$	0.1	HHH_E__LL_L0___	“HELLO” v.1
	0.02	HH_E__LL_L0___	“HELLO” v.2
	0.01	HHH_E__L_L_0H___	“HELL OH”
	0.01	HHH_EE_LL_L_0___	“HELLO” v.3
...		YY_E__LL_L0_W_	“YELLOW”

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

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

# Connectionist Temporal Classification (CTC)

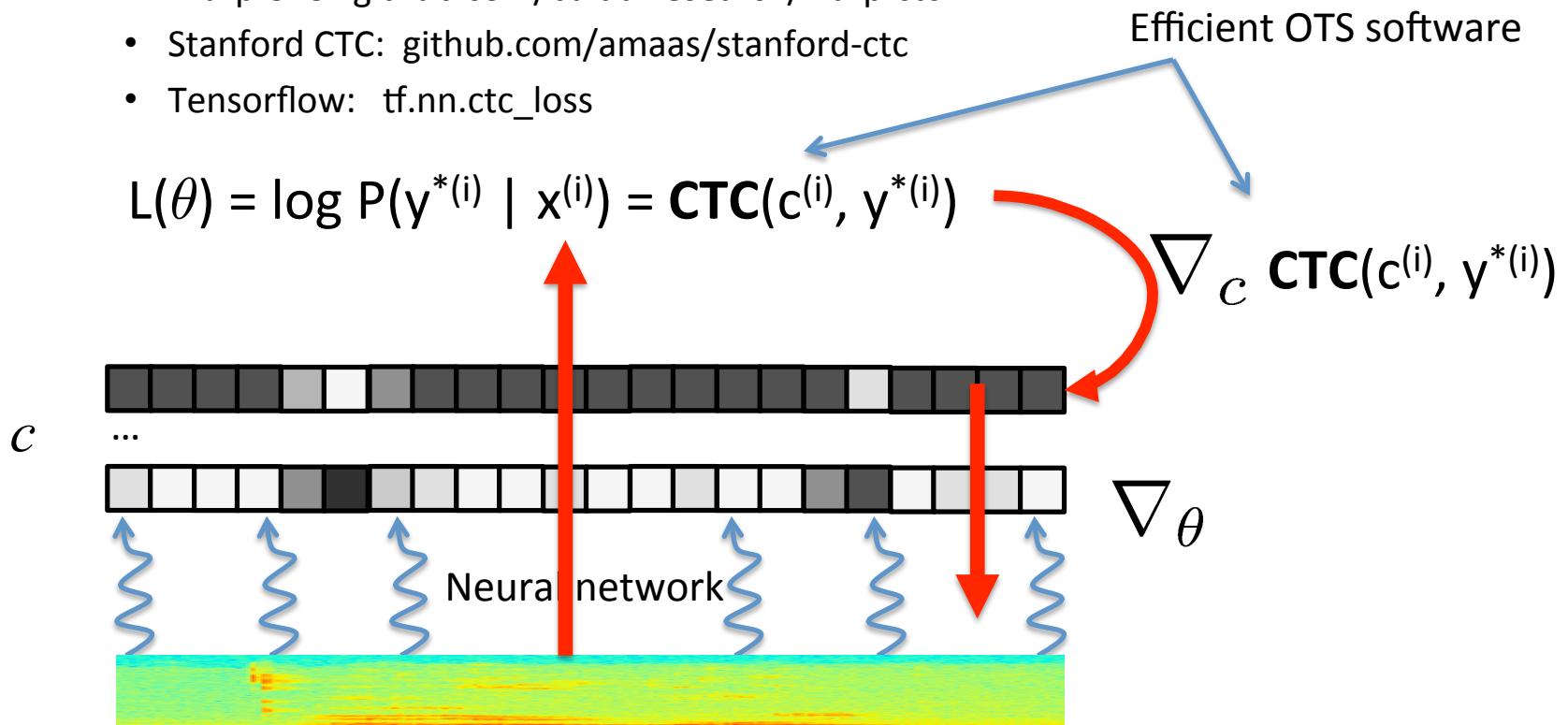
3. Update network parameters  $\theta$  to maximize likelihood of correct label  $y^*$ :

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

- [Graves et al., 2006] provides an efficient dynamic programming algorithm to compute the inner summation and its gradient.

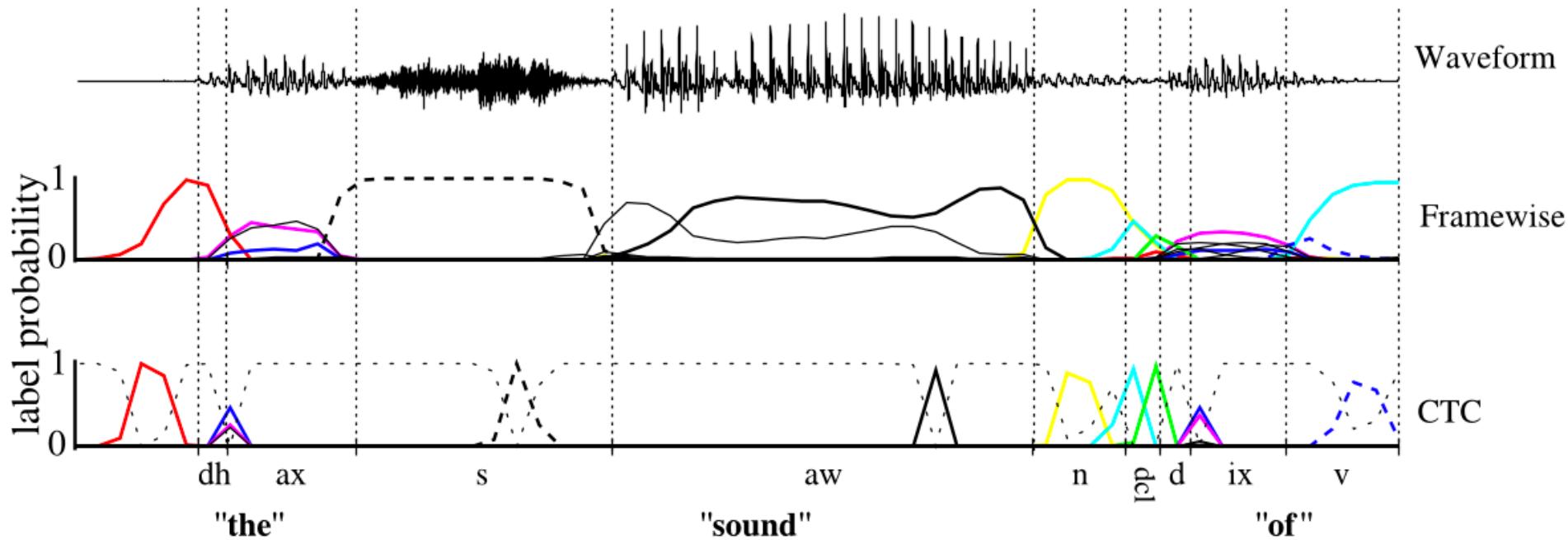
# Connectionist Temporal Classification (CTC)

- Use usual gradient descent methods to optimize.  
Tune entire network with backpropagation.
  - Given network outputs, many off-the-shelf packages to compute CTC loss (likelihood) from  $c$  and  $y^*$ , and gradient w.r.t.  $c$ .
    - Warp CTC: [github.com/baidu-research/warp-ctc](https://github.com/baidu-research/warp-ctc)
    - Stanford CTC: [github.com/amaas/stanford-ctc](https://github.com/amaas/stanford-ctc)
    - Tensorflow: `tf.nn.ctc_loss`



# Connectionist Temporal Classification (CTC)

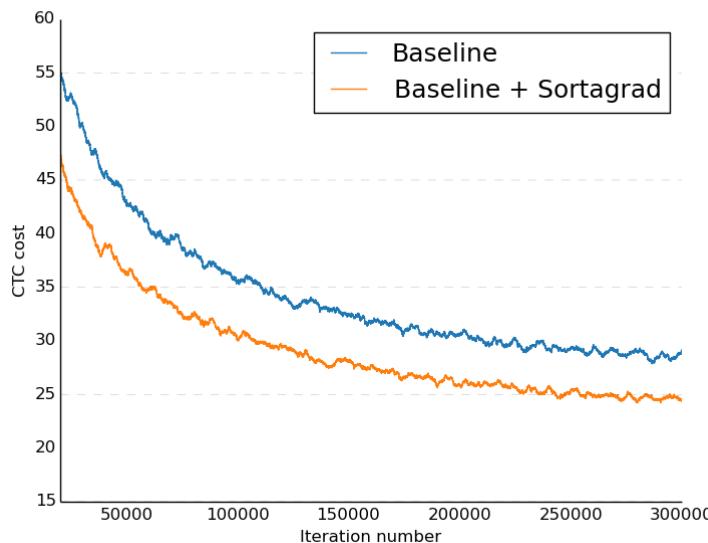
- The framewise network receives an error for misalignment
- The CTC network predicts the sequence of phonemes / characters (as spikes separated by ‘blanks’)
- No forced alignment (initial model) required for training.



# Training tricks

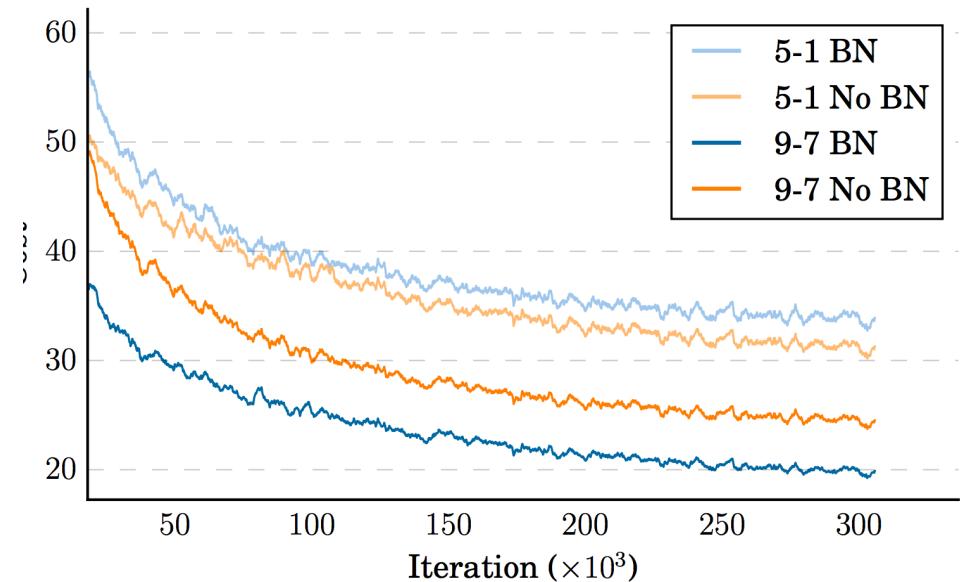
- Getting RNN to train well is tricky.

“SortaGrad”: order utterances by length during first epoch.



See [“Curriculum Learning”, Bengio et al., ICML 2009]

Batch normalization

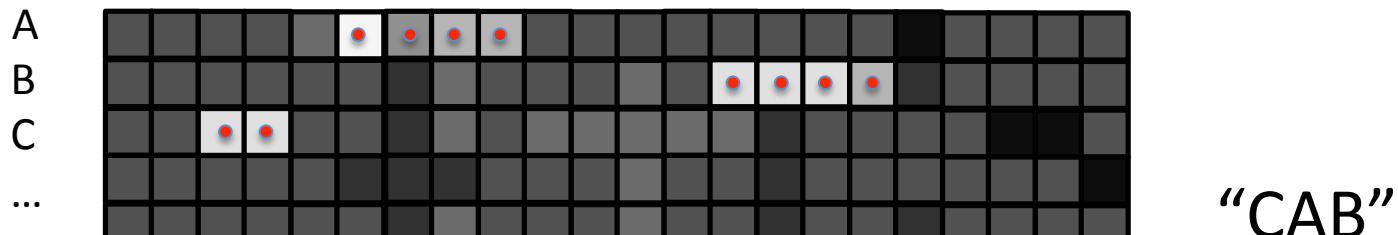


See [Ioffe & Szegedy, 2015]

# Decoding

- Network outputs  $P(c|x)$ . How do we find most likely transcription from  $P(y|x)$ ?
- Simple (approximate) solution:

$$\beta \left( \arg \max_c P(c|x) \right)$$

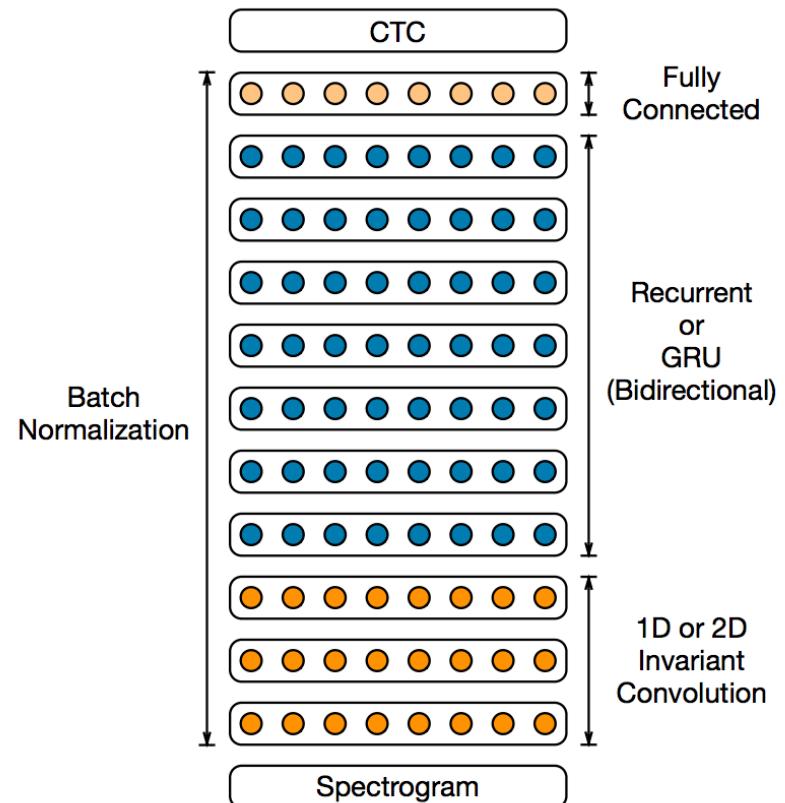


- Often terrible, but a useful diagnostic to “eyeball” models.

# Example

- RNN to predict graphemes  
(26 characters + space + blank):
  - Spectrograms as input.
  - 1 layer of convolutional filters.
  - 3 layers of Gated Recurrent Units.
    - 1000 neurons per layer.
  - 1 full-connected layer to predict  $c$ .
  - Batch normalization  
[Ioffe & Szegedy, 2015]
- CTC loss function (Warp-CTC)
- Train with SGD+Nesterov momentum.

Typical model family:



# Example

- Wall Street Journal:
  - <https://catalog.ldc.upenn.edu/LDC93S6A>
  - Reading WSJ articles.

Free alternative: LibriSpeech

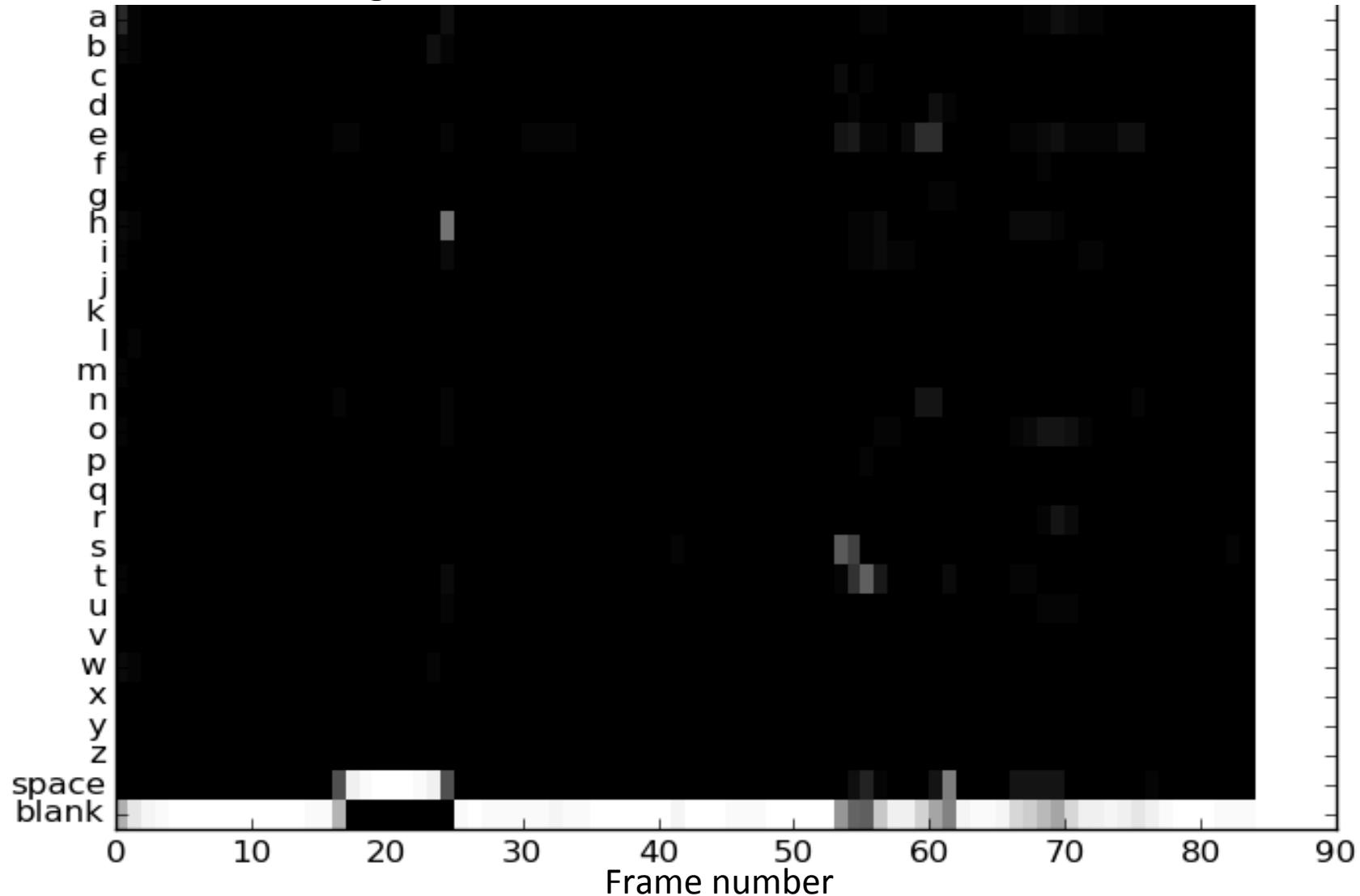
- <http://www.openslr.org/12/> [Panayotov et al., ICASSP 2015]
- Read speech from public domain audiobooks.

# Example

- What happens inside?

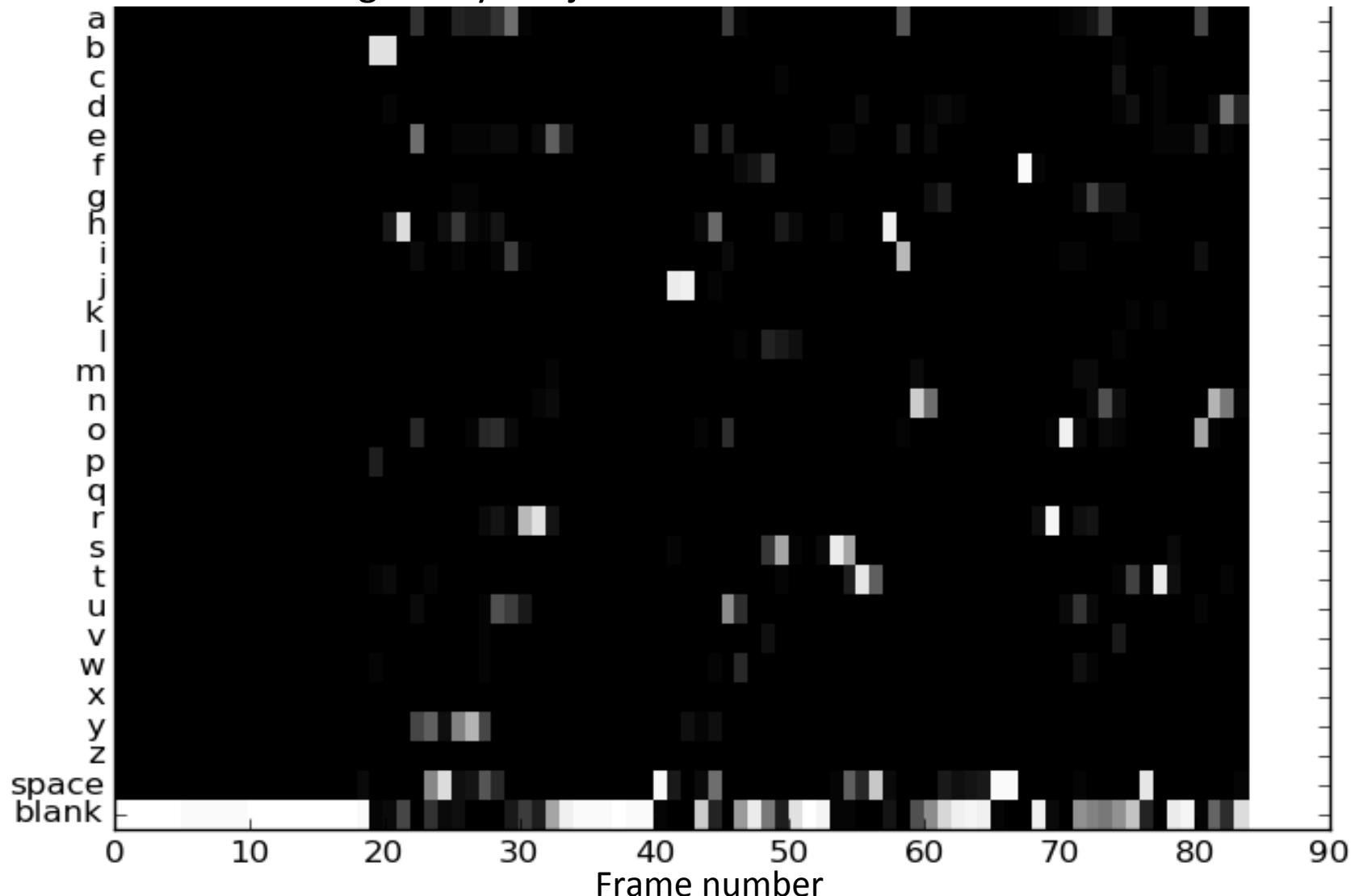
# Network outputs ( $c$ ) at Iteration 300 (Thresholding / contrast added for clarity.)

Max decoding: h



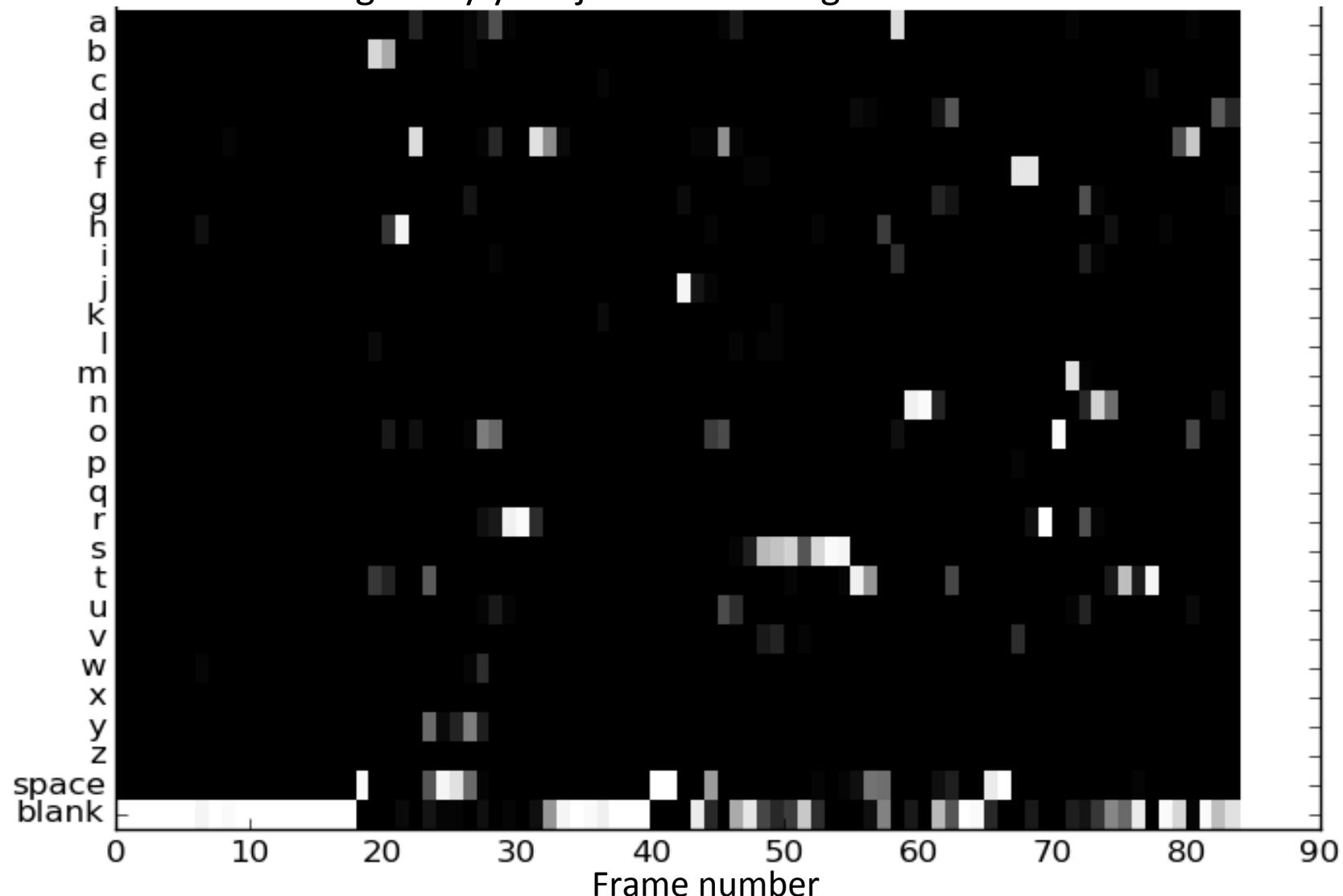
# Network outputs ( $c$ ) at Iteration 1500 (Thresholding / contrast added for clarity.)

Max decoding: bhe y uar j usst hin fro ton



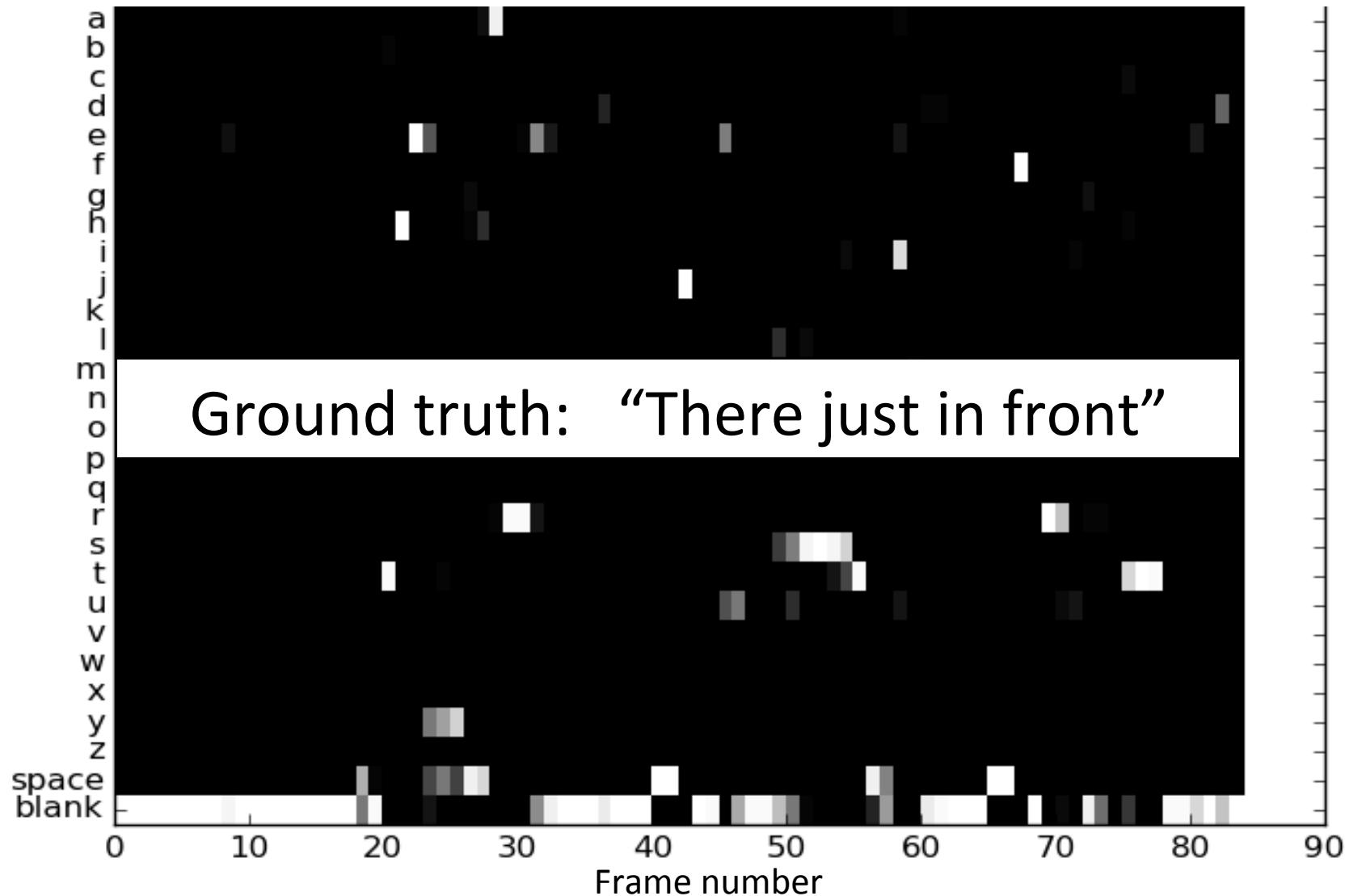
# Network outputs ( $c$ ) at Iteration 2500 (Thresholding / contrast added for clarity.)

Max decoding: bhey yore j esstand fromgnnte



# Network outputs ( $c$ ) at Iteration 5500 (Thresholding / contrast added for clarity.)

Max decoding: they ar jest in front



# Max Decoding

- Examples:

**Max Decoding**

“put pore lotttle thank and sr crits sinpt the atting to them  
having been turned ef the wal al thes years con”

**True Label**

“the poor little things cried cynthia think of them  
having been turned to the wall all these years”

**Max Decoding**

“that is true baddel gre”

**True Label**

“that is true badauderie”

# Decoding

- Network outputs  $P(c|x)$ . How do we find most likely transcription from  $P(y|x)$ ?
- No efficient solution in general. Resort to search!
  - See [Graves et al., 2006] for prefix decoding strategy.

# Language models

- Even with better decoding, CTC model tends to make spelling + linguistic errors. E.g.:

RNN output

---

what is the weather like in bostin right now  
prime miniter nerernr modi  
arther n tickets for the game

From Hannun et al., 2014.

- $P(y|x)$  modeled directly from audio.
  - But not enough audio data to learn complicated spelling and grammatical structure.
  - Only supports small vocabulary.
  - For grapheme models: “Tchaikovsky” problem.

# Language models

- Two solutions
  - Fuse acoustic model with language model:  $P(y)$
  - Incorporate linguistic data:
    - Predict phonemes + pronunciation lexicon + LM.
- Possible to train language model from massive *text corpora*.
  - Learn spelling + grammar
  - Greatly expand vocabulary
  - Elevate likely cases (“Tchaikovsky concerto”) over unlikely cases (“Try cough ski concerto”).

# Language models

- Standard approach: n-gram models
  - Simple n-gram models are common, well supported.
    - KenLM: [kheafield.com/code/kenlm/](http://kheafield.com/code/kenlm/)
  - Train easily from huge corpora.
  - Quickly update to follow trends in traffic.
  - Fast lookups inside decoding algorithms.

# Decoding with LMs

- Given a word-based LM of form  $P(w_{t+1} | w_{1:t})$ ,  
Hannun et al., 2014 optimize:

$$\arg \max_w P(w|x)P(w)^\alpha [\text{length}(w)]^\beta$$

$P(w|x) = P(y|x)$  for characters that make up  $w$ .

$\alpha$  and  $\beta$  are tunable parameters to govern weight of LM and a bonus/penalty for each word.

# Decoding with LMs

- Basic strategy: beam search to maximize

$$\arg \max_w P(w|x)P(w)^\alpha [\text{length}(w)]^\beta$$

Start with set of candidate transcript prefixes,  $A = \{\}$ .

For  $t = 1..T$ :

For each candidate in  $A$ , consider:

1. Add blank; don't change prefix; update probability using AM.
2. Add space to prefix; update probability using LM.
3. Add a character to prefix; update probability using AM.

Add new candidates with updated probabilities to  $A_{\text{new}}$ .

$A := K$  most probable prefixes in  $A_{\text{new}}$ .

# Decoding with LMs: Examples

RNN output	Decoded Transcription
what is the weather like in bostin right now prime miniter nerrenr modi arther n tickets for the game	what is the weather like in boston right now prime minister narendra modi are there any tickets for the game

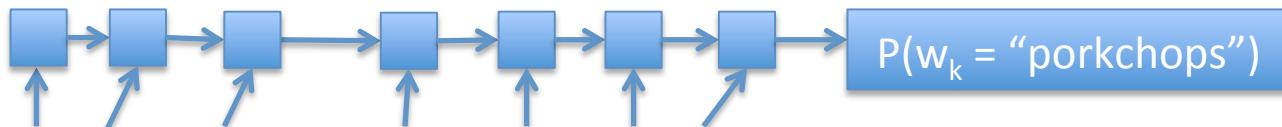
From Hannun et al., 2014.

# Rescoring

- Another place to plug in DL algorithms:  
Systems usually produce N-best list.  
Use fancier models to “rescore” this list.

# Rescoring with Neural LM

- Example: train neural language model and rescore word transcriptions.
  - Cheap to evaluate  $P(w_k | w_{k-1}, w_{k-2}, \dots, w_1)$  NLM on many sentences.
  - In practice, often combine with N-gram trained from big corpora.



1. (-25.45) I'm a connoisseur looking for wine and porkchops. -24.45
2. (-26.32) I'm a connoisseur looking for wine and port shops. -23.45
3. ...
4. ...
5. ...

# Data

- Transcribing speech data isn't cheap, but not prohibitive:
  - Roughly 50¢ to \$1 per minute.
- Typical speech benchmarks offer 100s to few 1000s of hours.
  - LibriSpeech (audiobooks)
  - LDC corpora (Fisher, Switchboard, WSJ) (\$\$)
  - VoxForge

# Types of speech data

- Application matters
  - We want to find data that matches our goals.

## Styles of speech

Read  
Conversational  
Spontaneous  
Command/control

## Issues

Disfluency / stuttering  
Noise  
Mic quality / #channels  
Far field  
Reverb / echo  
Lombard effect  
Speaker accents

## Applications

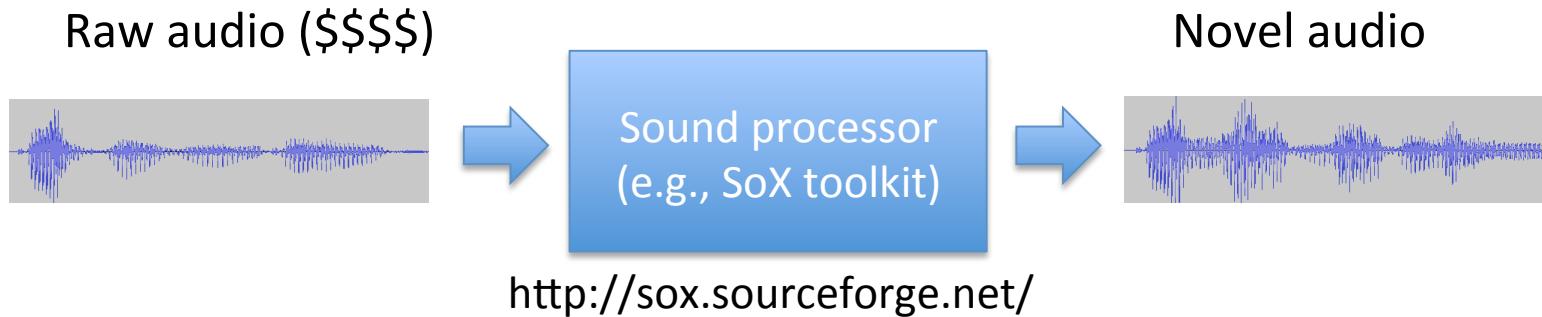
Dictation  
Meeting transcription  
Call centers  
Device control  
Mobile texting  
Home / IoT / Cars

# Read speech

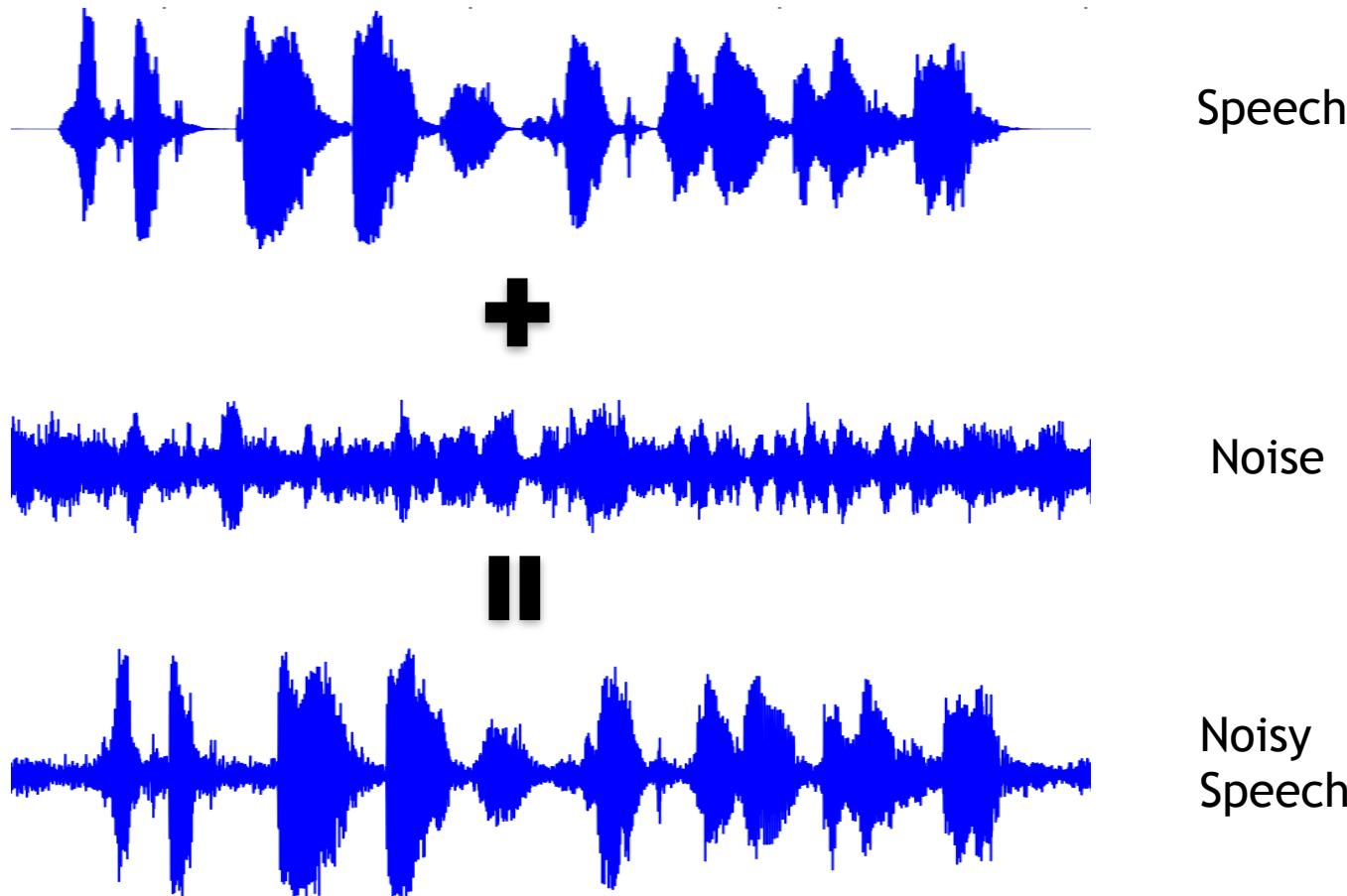
- Reading is inexpensive way to get more data.  
    < \$10/hour depending on source
- Disadvantages:
  - Misses inflection/conversational tone
  - Lombard effect
  - Speaker variety sometimes a limitation.

# Augmentation

- Many forms of distortion that model should be robust to:
  - Reverb, noise, far field effects, echo, compression artifacts, changes in tempo

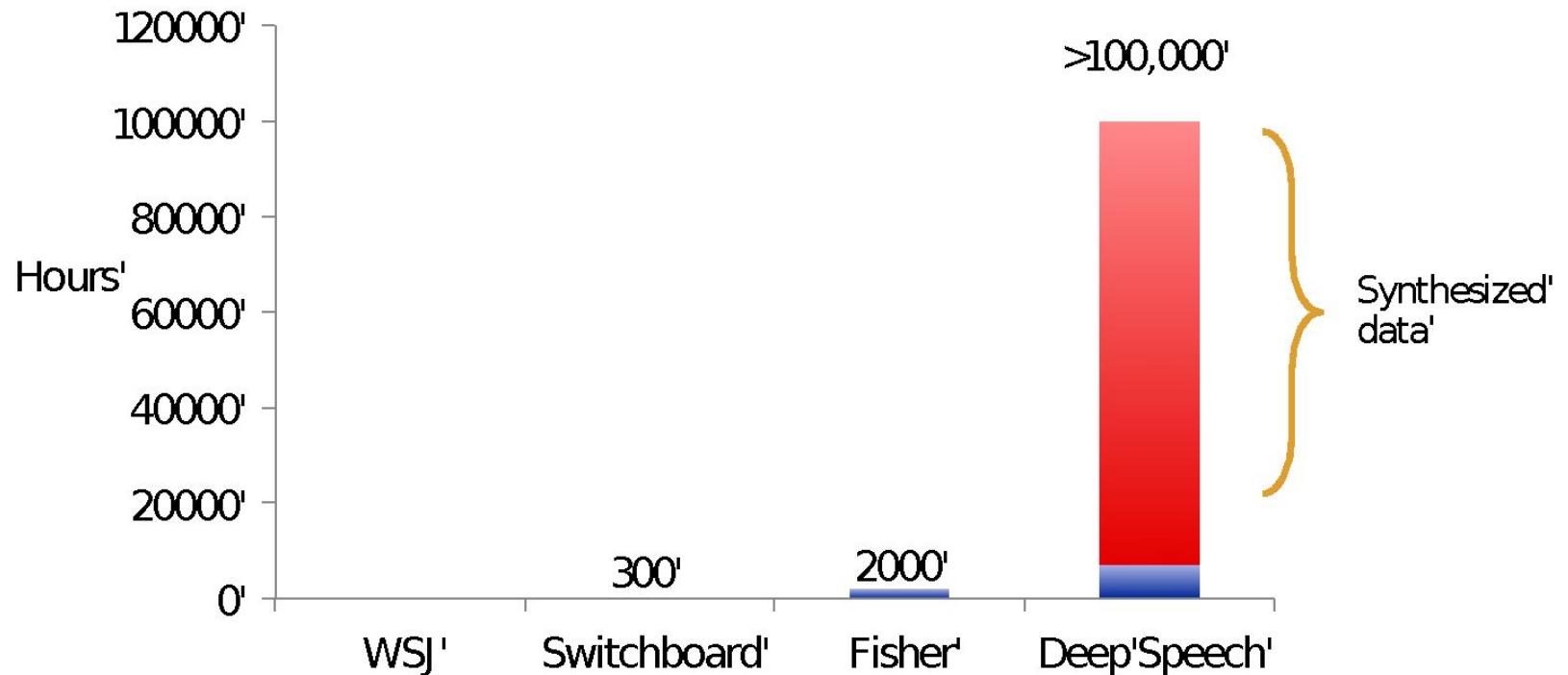


# Example: additive noise



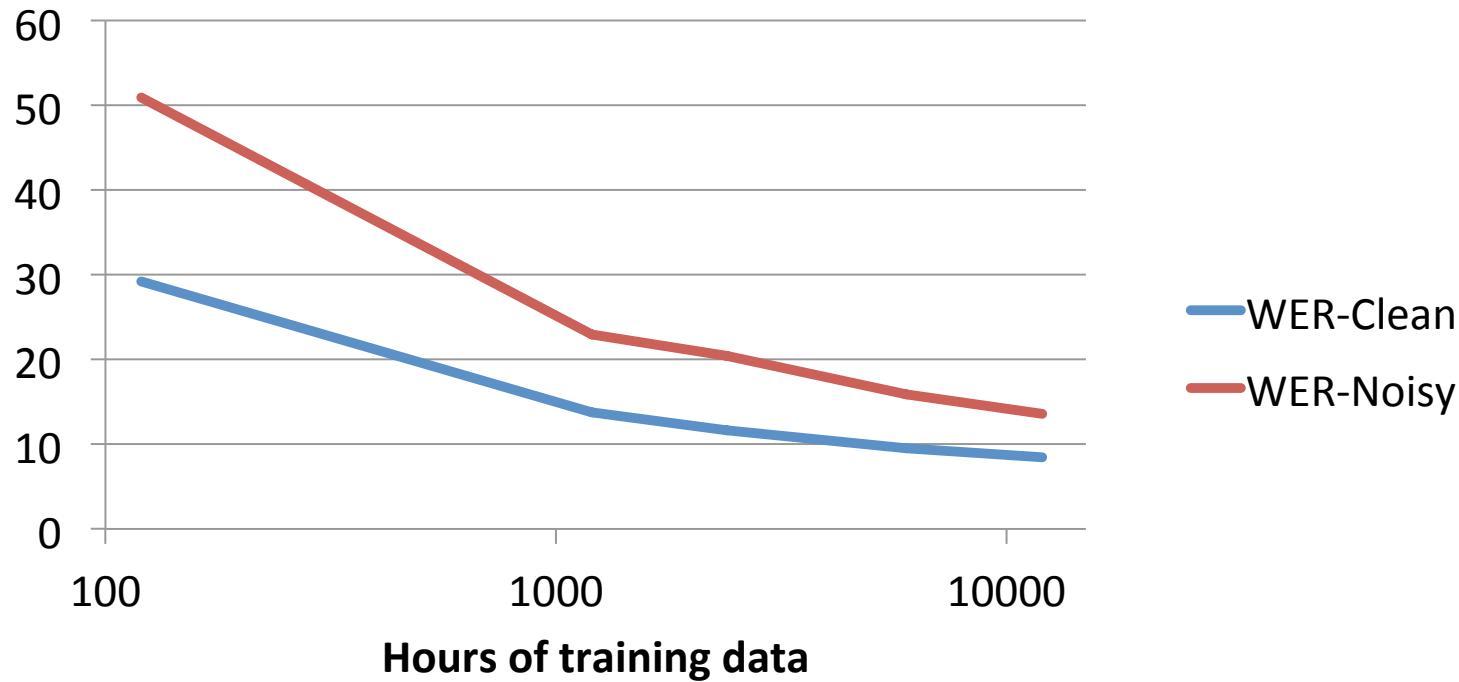
DeepSpeech 2: 10K hours of raw audio -> 100K hours of novel training data

# Data Argumentation



# Results: DeepSpeech 2

- Steady fall in error rates with new raw data.



# Computation

- How big is 1 experiment?

At least:

$$(\# \text{ connections}) \cdot (\# \text{ frames}) \cdot (\# \text{ utterances}) \cdot \\ (\# \text{ epochs}) \cdot 3 \cdot 2 \text{ FLOPs}$$

E.g.: for DS2 with 10K hours of data:

$$100\text{e}6 \cdot 100 \cdot 10\text{e}6 \cdot 20 \cdot 3 \cdot 2 = 1.2\text{e}19 \text{ FLOPs}$$

~30 days with well-optimized code on Titan X.

# Computation

CPU  
~100 GigaFLOPs



Titan X (2016) GPU  
~11 TeraFLOPS

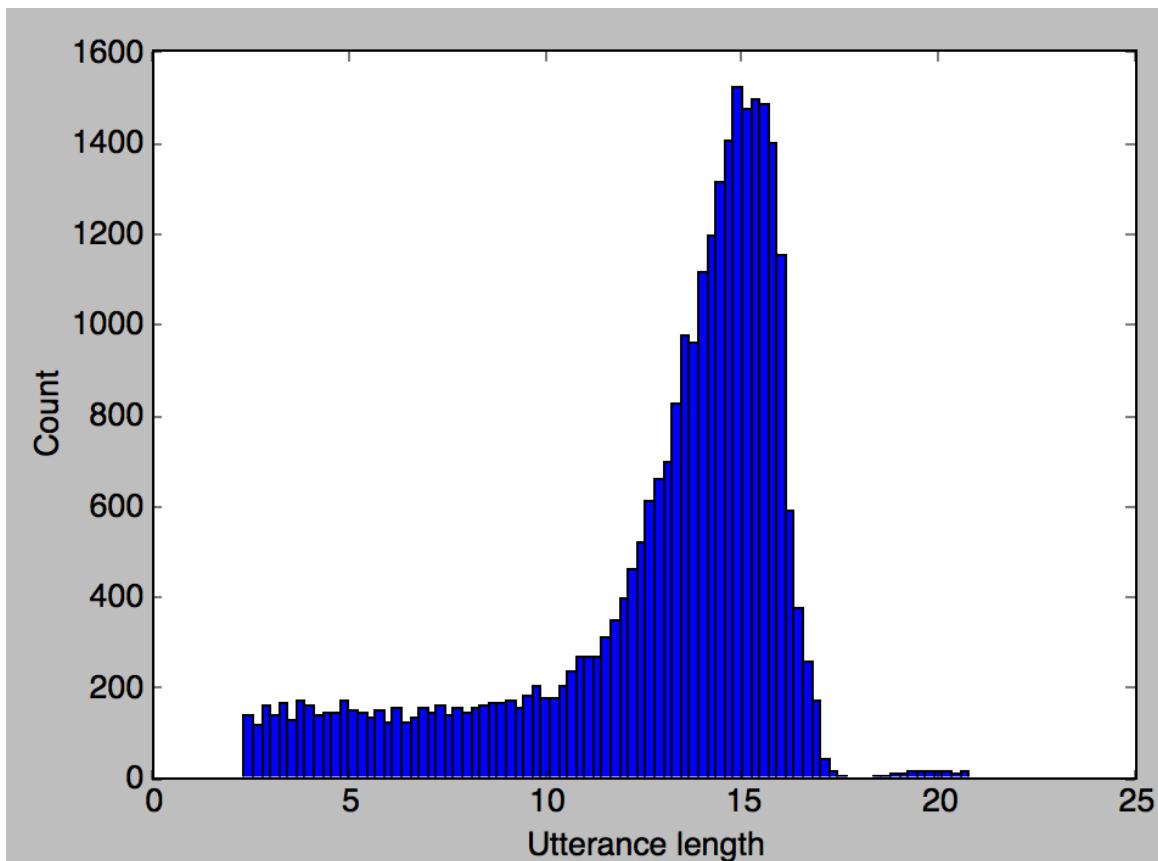


# Computation

- Easy: use more GPUs with data parallelism.
- Current servers support ~8 Titans.
  - Will get you < 1 week training time.

# Computation

- Try to keep similar-length utterances together.

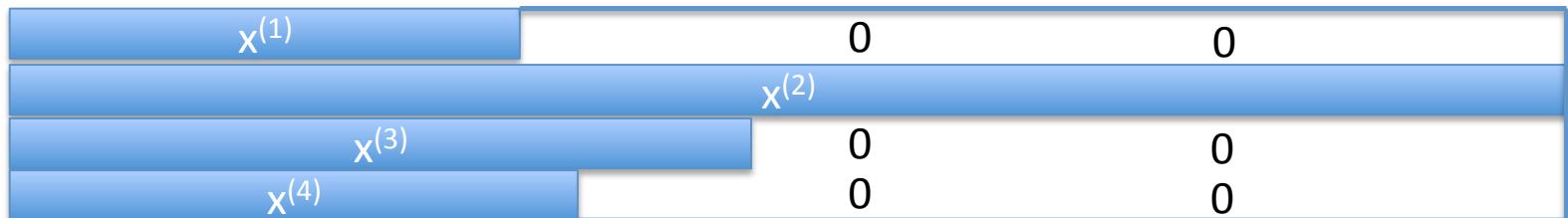


(LibriSpeech  
clean data.)

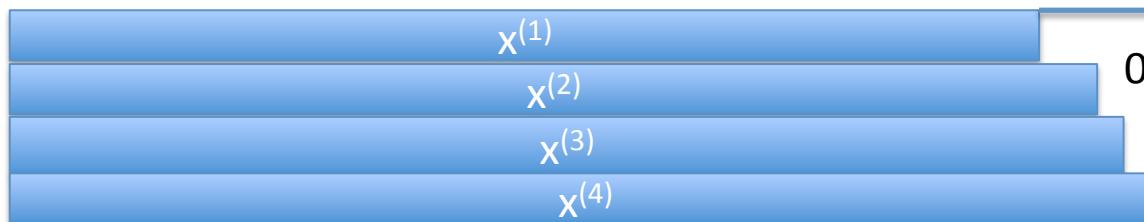
# Computation

- Try to keep similar-length utterances together.

Bad minibatch:



Good minibatch:



# Hybrid (IBM) vs. DeepSpeech 1 (Baidu)

	<b>IBM 2015</b>	<b>Baidu 2014</b>
Features	VTL-PLP, MVN, LDA, STC, fMMLR, i-Vector	80 log filter banks
Alignment	GMM-HMM 300K Gaussians	-
DNN	DNN(5x2048) + CNN(128x9x9+5x2048) ++RNN 32000 outputs	4RNN (5 x 2304) 28 outputs
DNN Training	CE + MBR Discriminative Training (ST)	CTC
HMM	32K states (DNN outputs) pentaphone acoustic context	-
Language Model	37M 4-gram + model M (class based exponential model) + 2 NNLM	4-gram (Transcripts)

# DeepSpeech 1 vs. DeepSpeech 2 (Baidu)

	<b>Deep Speech 1 (Baidu 2014)</b>	<b>DS2 (Baidu 2015)</b>
Features	80 log filter banks	?
Alignment	-	-
DNN	4RNN (5 x 2304) 28 outputs	9-layer, 7RNN, BatchNorm, Conv. Layers. (Time/Freq)
DNN Training	CTC	CTC
HMM	-	-
Language Model	4-gram	5-gram

# DeepSpeech 1 vs. DeepSpeech 2 (Baidu)

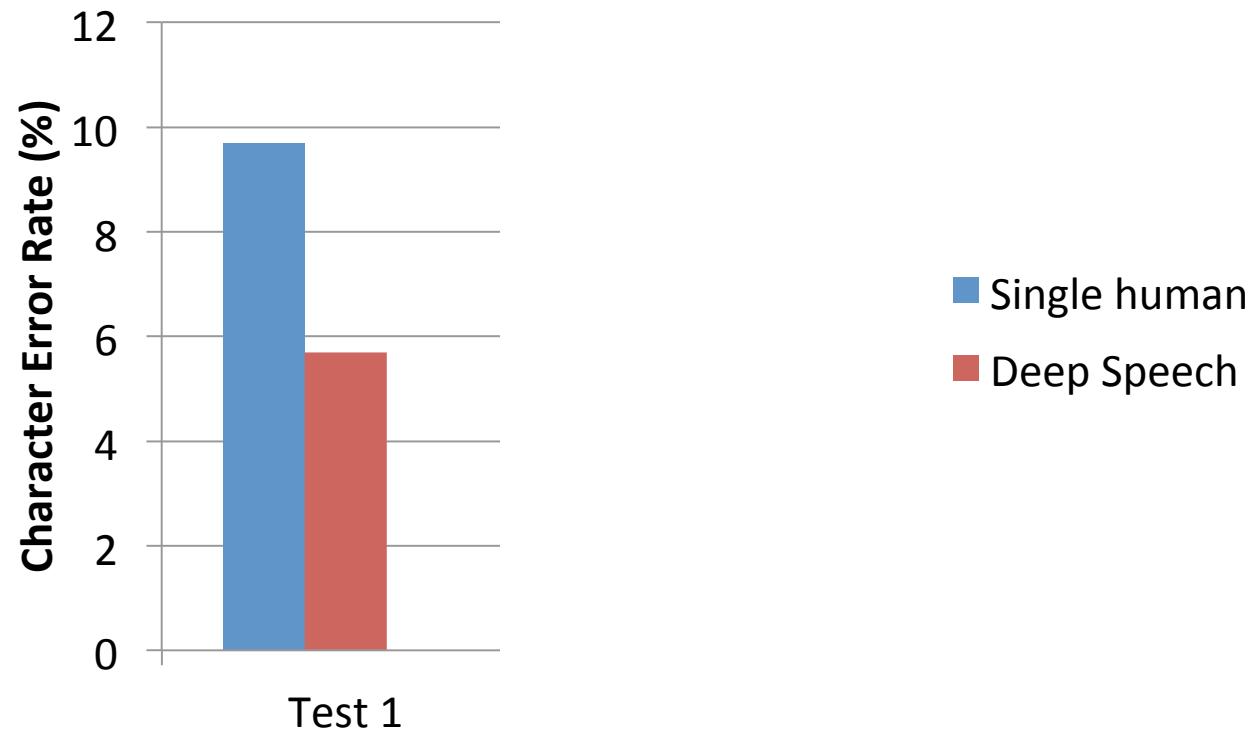
Accented Speech			
Test set	DS1	DS2	Human
VoxForge American-Canadian	15.01	7.55	4.85
VoxForge Commonwealth	28.46	13.56	8.15
VoxForge European	31.20	17.55	12.76
VoxForge Indian	45.35	22.44	22.15

**Table 14:** Comparing WER of the DS1 system to the DS2 system on accented speech.

Noisy Speech			
Test set	DS1	DS2	Human
CHiME eval clean	6.30	3.34	3.46
CHiME eval real	67.94	21.79	11.84
CHiME eval sim	80.27	45.05	31.33

# Results

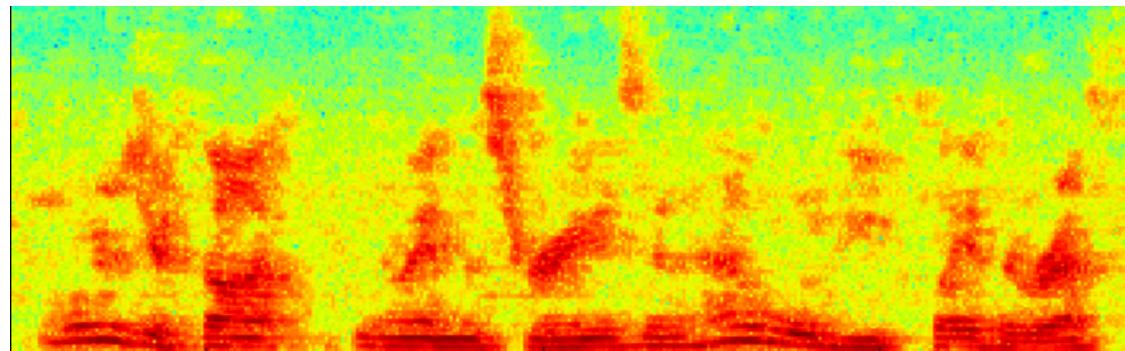
- Scaled up models in Mandarin:



# Difficulties in Mandarin Speech Recognition

## 1. Mandarin is a *tonal* language

mā má mǎ mà ·ma  
妈 麻 马 骂 吗



# Difficulties in Mandarin Speech Recognition

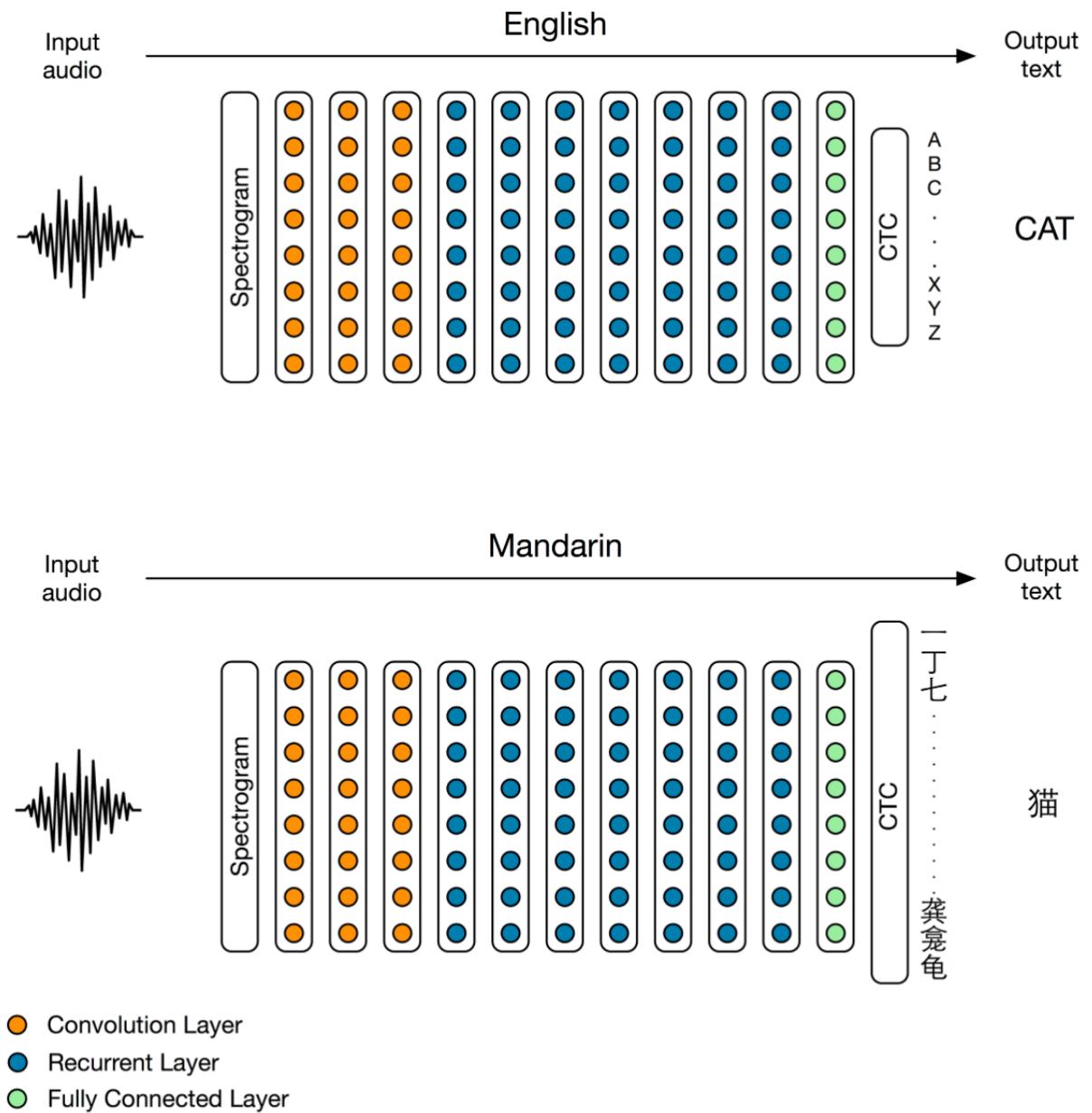
## 2. Thousands of characters! > 80K

事得真对看见加更多  
男女几各谁找子字那哪  
说着位把吧难来站每起  
被只都做已长行等再以  
所后分种将很而数天无  
吗家可件里最回万能爱  
时也还出去到他性就部  
新市与内本地这此建全  
一二三四五六六十个次元  
用之要好了年月日为名  
不在于前者会号我和你  
的人上中下大小是没有

# Difficulties in Mandarin Speech Recognition

## 3. The homophone problem is ubiquitous.





**Figure 1.** Architecture of the DS2 system used on English versus Mandarin. The only change is the larger output layer in Mandarin that accommodates the 6000 Mandarin characters, as opposed to the 29 we use in English.

# Summary

- Deep Learning makes first steps to state-of-art speech system simpler than ever.
- Performance significantly driven by data & models.
  - Focus on scaling data + compute.
  - Try more models, make more progress!
- Mature enough for production.
  - DeepSpeech model is live in Mandarin & English.

Starter code: [github.com/baidu-research/ba-dls-deepspeech](https://github.com/baidu-research/ba-dls-deepspeech)

# References

- Gales and Young. "The Application of Hidden Markov Models in Speech Recognition" Foundations and Trends in Signal Processing, 2008.
- Jurafsky and Martin. "Speech and Language Processing". Prentice Hall, 2000.
- Dzmitry Bahdanau, Jan Chorowski, Dmitriy Serdyuk, Philemon Brakel, Yoshua Bengio. "End-to-End Attention-based Large Vocabulary Speech Recognition." 2015. arxiv.org/abs/1508.04395
- Bourlard and Morgan. "CONNECTIONIST SPEECH RECOGNITION: A Hybrid Approach". Kluwer Publishing, 1994.
- William Chan, Navdeep Jaitly, Quoc V. Le, Oriol Vinyals. "Listen Attend Spell" 2015. arxiv.org/abs/1508.01211
- Jianmin Chen, Rajat Monga, Samy Bengio, Rafal Jozefowicz, "Revisiting Distributed Synchronous SGD." ICLR 2016 arXiv:1604.00981
- A Graves, S Fernández, F Gomez, J Schmidhuber. "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks." ICML, 2006.
- Dahl, Yu, Deng, Acero, "Large Vocabulary Continuous Speech Recognition with Context-Dependent DBN-HMMs". ICASSP, 2011
- Dean, J., Corrado, G., Monga, R., Chen, K., Devin, M., Mao, M., Senior, A., Tucker, P., Yang, K., Le, Q.V. and Ng, A.Y. "Large scale distributed deep networks." NIPS 2012
- Hannun, Maas, Jurafsky, Ng. "First-Pass Large Vocabulary Continuous Speech Recognition using Bi-Directional Recurrent DNNs" ArXiv: 1408.2873
- Hannun, et al. "Deep Speech: Scaling up end-to-end speech recognition". ArXiv:1412.5567
- H. Hermansky, "Perceptual linear predictive (PLP) analysis of speech", J. Acoust. Soc. Am., vol. 87, no. 4, pp. 1738-1752, Apr. 1990.
- H. Hermansky and N. Morgan, "RASTA processing of speech", IEEE Trans. on Speech and Audio Proc., vol. 2, no. 4, pp. 578-589, Oct. 1994.
- Vassil Panayotov, Guoguo Chen, Daniel Povey and Sanjeev Khudanpur, "LibriSpeech: an ASR corpus based on public domain audio books." ICASSP 2015
- H. Schwenk, "Continuous space language models", 2007.