

# IntroHLT: End-to-end ASR

October 24, 2023

# Speech Recognition

- Classical Methods
  - Noisy Channel

$$p_{\Theta}(\mathbf{w}|\mathbf{x}) = \frac{p_{\theta_1}(\mathbf{x}|\mathbf{w}) p_{\theta_2}(\mathbf{w})}{\sum_{\mathbf{w}} p_{\theta_1}(\mathbf{x}|\mathbf{w}) p_{\theta_2}(\mathbf{w})}$$

- ASR system has components

$$p_{\Theta}(\mathbf{w}|\mathbf{x}) = \frac{\overbrace{\sum_{\mathbf{s}} p_{\theta_1}(\mathbf{x}|\mathbf{s})}^{\mathbf{H}} \overbrace{p_{\theta_3}(\mathbf{s}|\mathbf{w})}^{\mathbf{L}} \overbrace{p_{\theta_2}(\mathbf{w})}^{\mathbf{G}}}{\sum_{\mathbf{w}, \mathbf{s}} p_{\theta_1}(\mathbf{x}|\mathbf{s}) p_{\theta_3}(\mathbf{s}|\mathbf{w}) p_{\theta_2}(\mathbf{w})}'$$

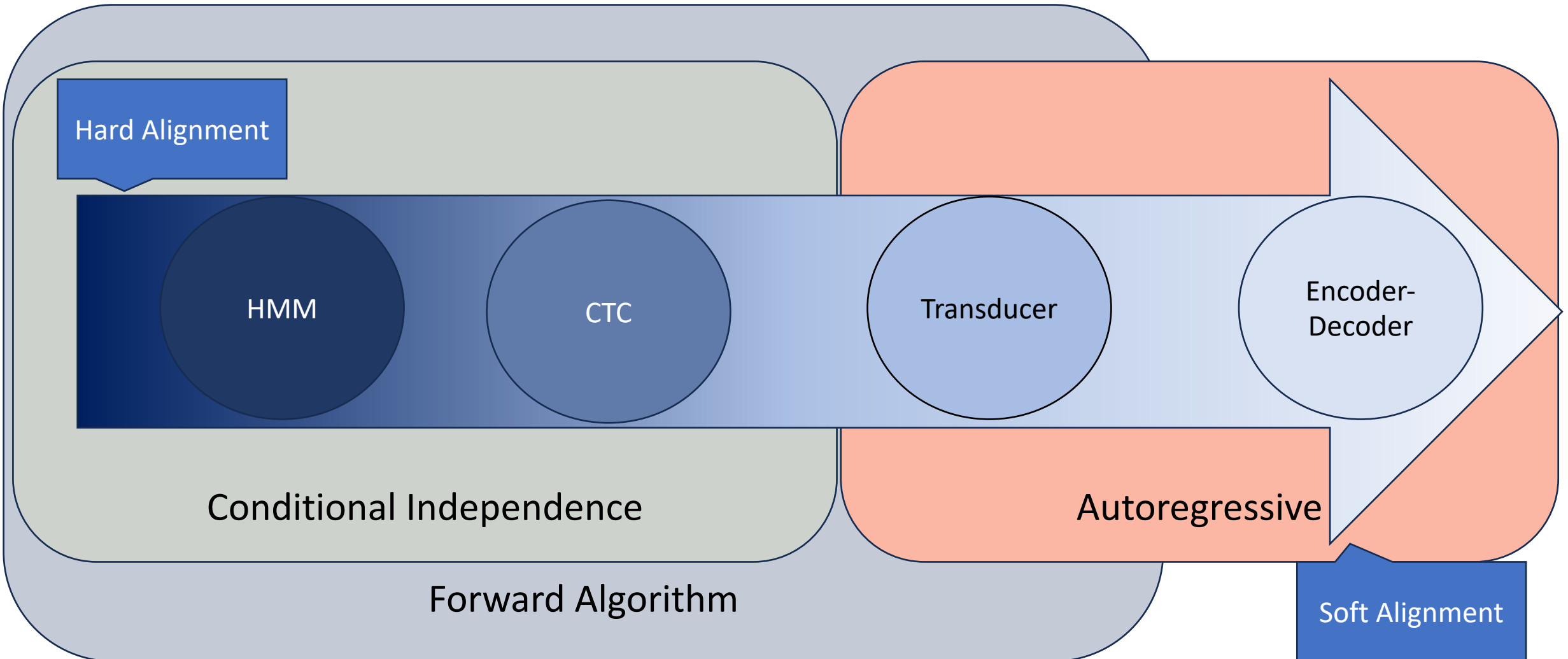
# Speech Recognition

- End-to-End Methods

$$p_{\Theta}(\mathbf{w}|\mathbf{x}) = \frac{\overbrace{\sum_{\mathbf{s}} p_{\theta_1}(\mathbf{x}|\mathbf{s})}^{\text{H}} \overbrace{p_{\theta_3}(\mathbf{s}|\mathbf{w})}^{\text{L}} \overbrace{p_{\theta_2}(\mathbf{w})}^{\text{G}}}{\sum_{\mathbf{w}, \mathbf{s}} p_{\theta_1}(\mathbf{x}|\mathbf{s}) p_{\theta_3}(\mathbf{s}|\mathbf{w}) p_{\theta_2}(\mathbf{w})}, \quad \Rightarrow \quad p_{\Theta}(\mathbf{w}|\mathbf{x}) = f_{\Theta}(\mathbf{x}, \mathbf{w})$$

- H, L, G are still present in  $f_{\Theta}(\mathbf{x}, \mathbf{w})$

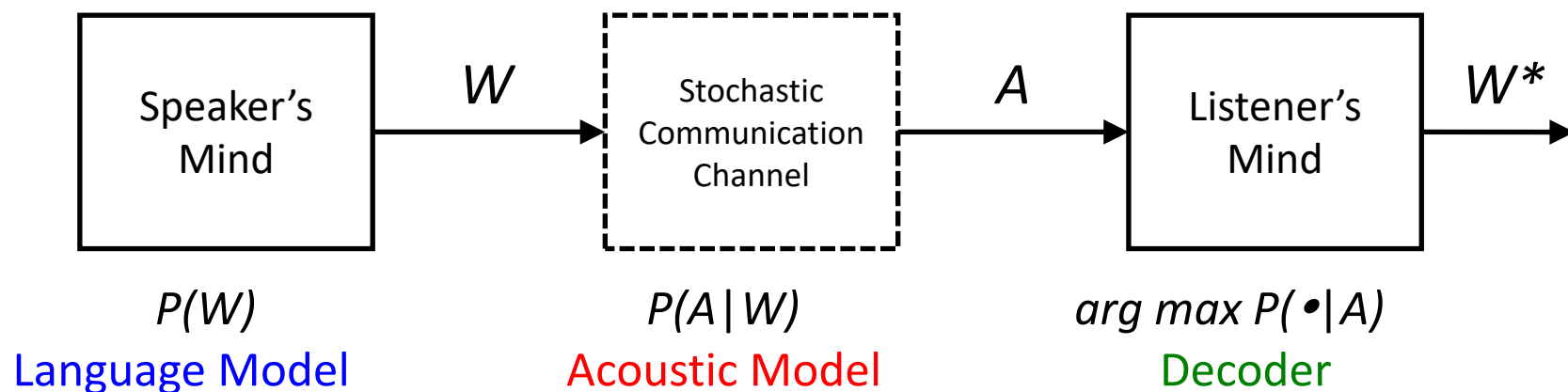
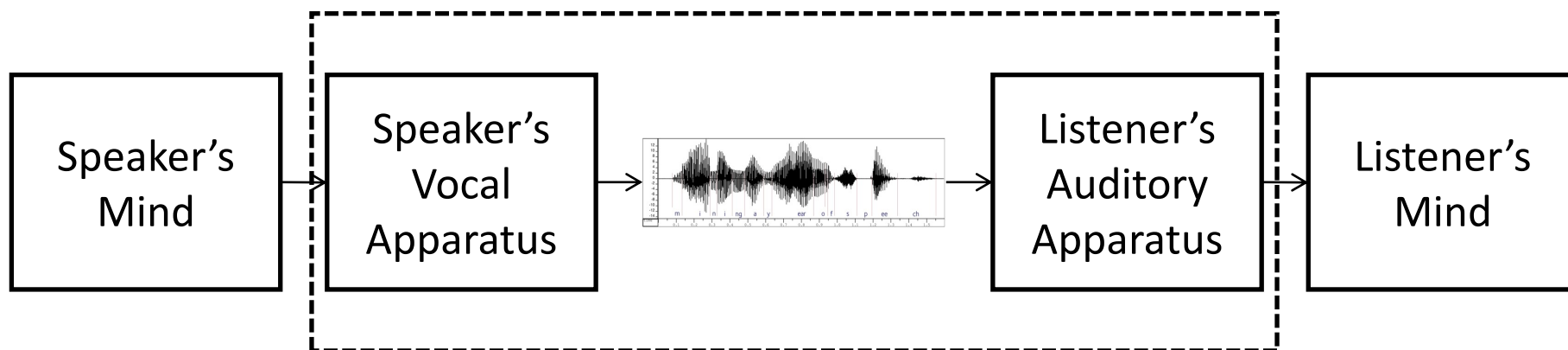
# Speech Recognition



HMM

Hidden Markov Models

# The “source-channel” model for automatic speech recognition (ASR)



# Hidden Markov models are popular as acoustic models

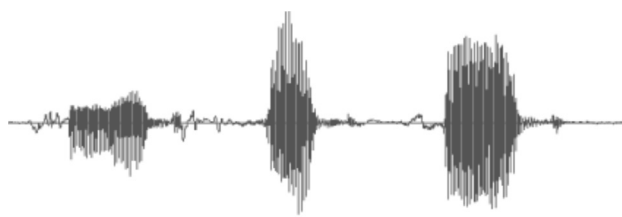
$$\begin{aligned} P(\mathbf{A} \mid \mathbf{W}) &= \sum_{\mathbf{S} \in \mathcal{S}(\mathbf{W})} P(\mathbf{A}, \mathbf{S} \mid \mathbf{W}) = \sum_{\mathbf{S} \in \mathcal{S}(\mathbf{W})} P(\mathbf{A} \mid \mathbf{S}, \mathbf{W}) P(\mathbf{S} \mid \mathbf{W}) \\ &\approx \sum_{\mathbf{S} \in \mathcal{S}(\mathbf{W})} P_E(\mathbf{A} \mid \mathbf{S}) P_T(\mathbf{S}) \\ &= \sum_{\mathbf{S} \in \mathcal{S}(\mathbf{W})} P_E(\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_T \mid s_1, s_2, \dots, s_T) P_T(s_1, s_2, \dots, s_T) \\ &= \sum_{\mathbf{S} \in \mathcal{S}(\mathbf{W})} \prod_{t=1}^T P_E(\mathbf{a}_t \mid s_t) P_T(s_t \mid s_{t-1}) \end{aligned}$$

Dynamic programming is popular for  
“decoding,” i.e. for hypothesis search

$$\begin{aligned}\widehat{\mathbf{W}} &= \arg \max_{\mathbf{W}} P(\mathbf{A} | \mathbf{W}) P(\mathbf{W}) \\ &= \arg \max_{\mathbf{W}} \sum_{\mathbf{S} \in \mathcal{S}(\mathbf{W})} P(\mathbf{A} | \mathbf{S}) P(\mathbf{S}) P(\mathbf{W}) \\ &\approx \arg \max_{\mathbf{W}} \max_{\mathbf{S} \in \mathcal{S}(\mathbf{W})} P(\mathbf{A} | \mathbf{S}) P(\mathbf{S}) P(\mathbf{W}) \\ &= \arg \max_{\mathbf{W}} \max_{\mathbf{S} \in \mathcal{S}(\mathbf{W})} \log P(\mathbf{A} | \mathbf{S}) + \log P(\mathbf{S}) + \log P(\mathbf{W}) \\ &\equiv \text{Project} \left( \text{Bestpath} \left( \text{Compose} \left( \mathbf{A}_{\log P(\mathbf{A} | \mathbf{S})} \circ \mathbf{L}_{\log P(\mathbf{S})} \circ \mathbf{G}_{\log P(\mathbf{W})} \right) \right) \right)\end{aligned}$$



# Composite HMM for "cat and hat"

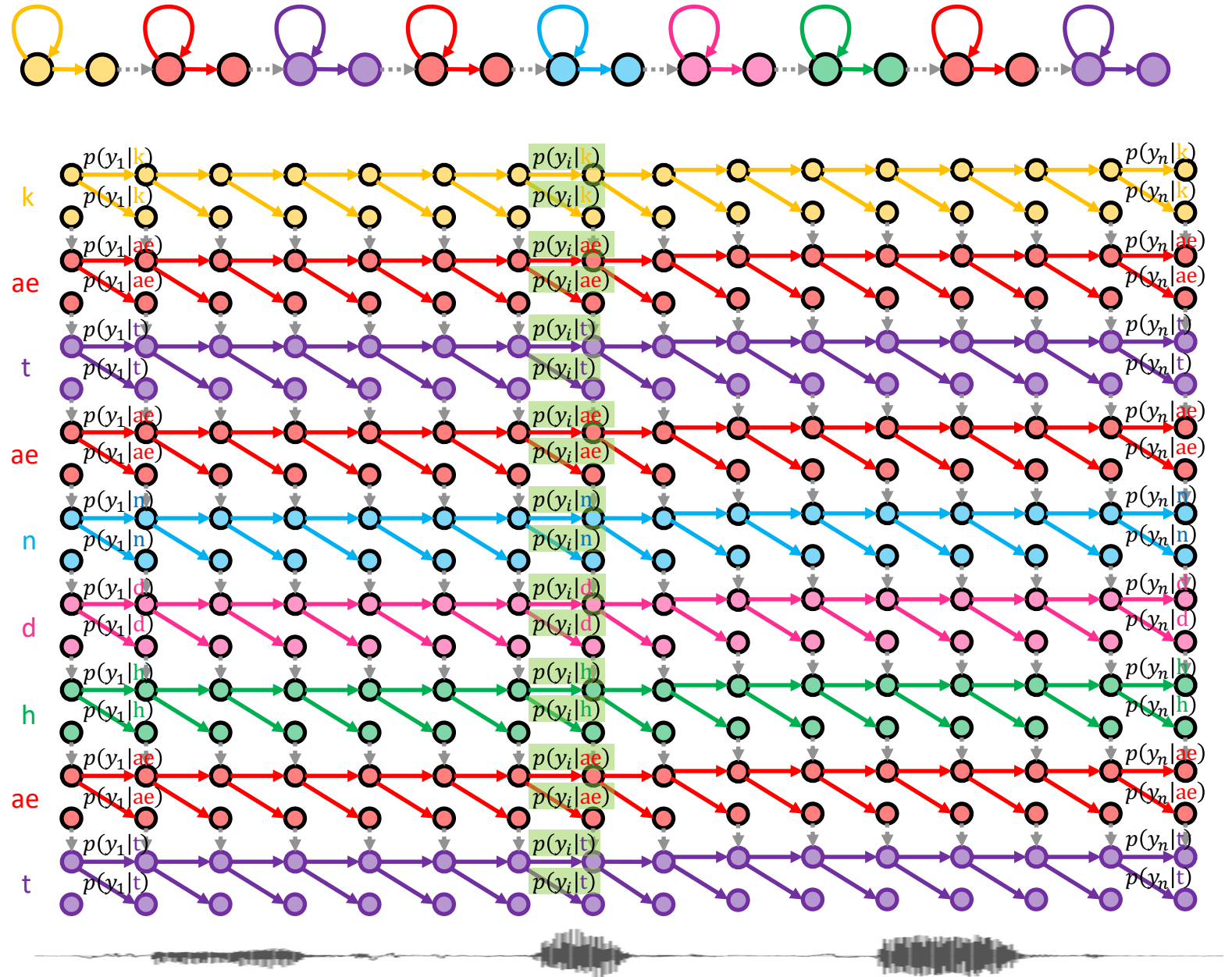
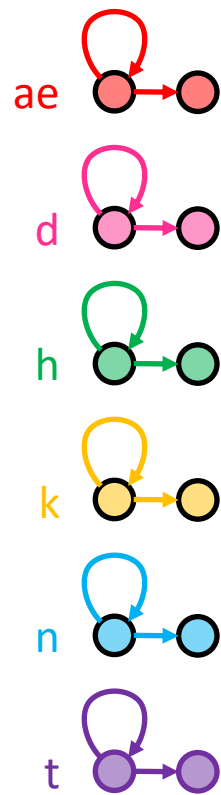


cat and hat

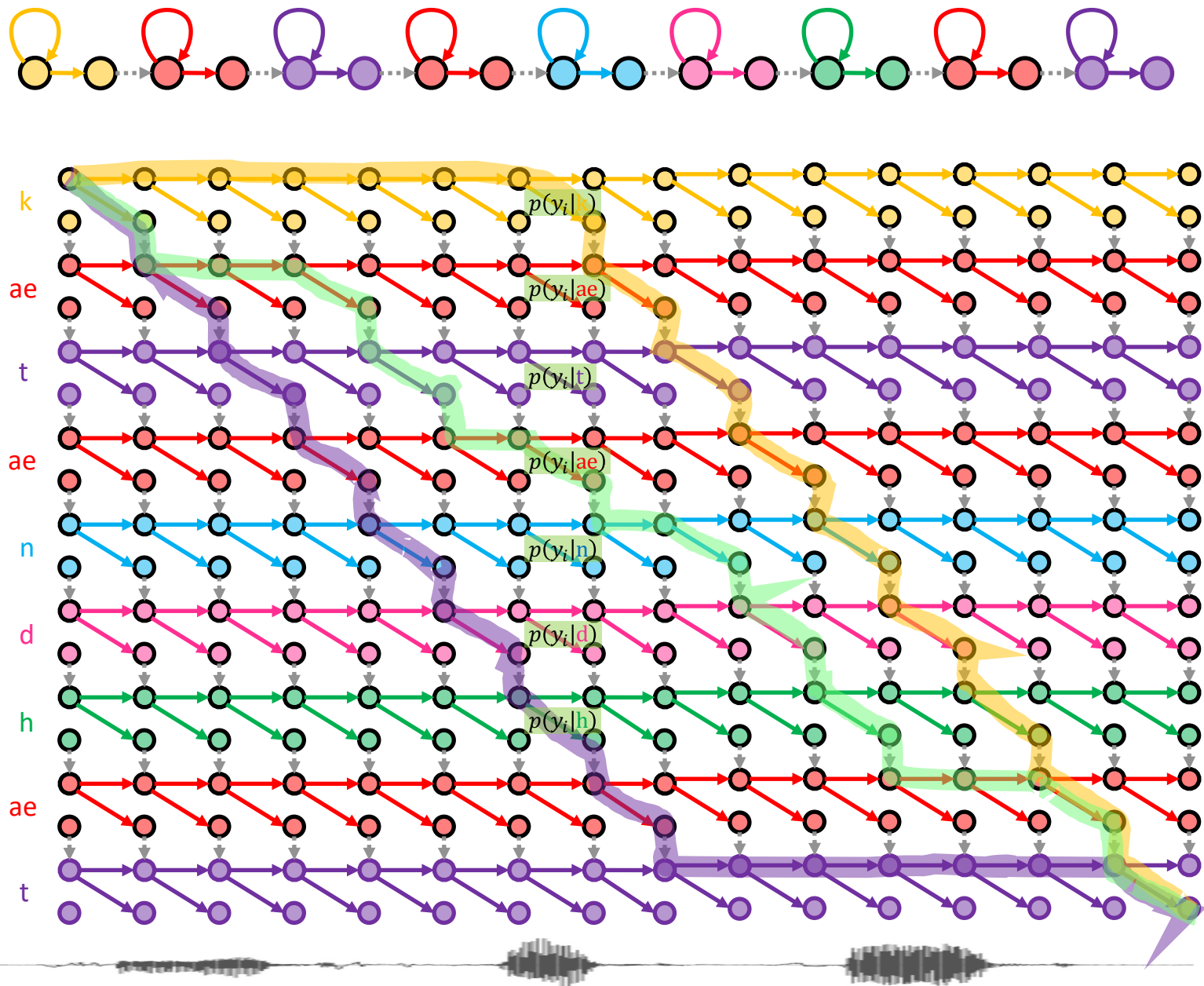
and  
cat  
hat

ae n d  
k ae t  
h ae t

Phoneme  
HMMs



# Composite HMM for “cat and hat”



## “Forward” Algorithm

$$P(\mathbf{y}|\mathbf{w}) = \sum_{\mathbf{s} \in \mathcal{S}(\mathbf{w})} P_{\vartheta}(\mathbf{y}|\mathbf{s})P_{\tau}(\mathbf{s})$$

$$= \sum_{\mathbf{s} \in \mathcal{S}(\mathbf{w})} \prod_{i=1}^n P_{\vartheta}(y_i|s_i)P_{\tau}(s_i|s_{i-1})$$

## Viterbi Algorithm

$$\hat{\mathbf{s}} = \arg \max_{\mathbf{s} \in \mathcal{S}(\mathbf{w})} P(\mathbf{s}|\mathbf{y})$$

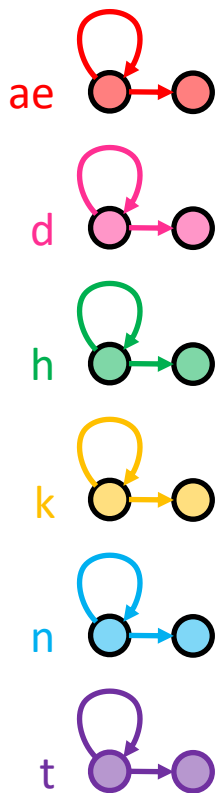
$$= \arg \max_{\mathbf{s} \in \mathcal{S}(\mathbf{w})} \frac{P(\mathbf{y}, \mathbf{s})}{P(\mathbf{y})}$$

$$= \arg \max_{\mathbf{s} \in \mathcal{S}(\mathbf{w})} \prod_{i=1}^n P_{\vartheta}(y_i|s_i)P_{\tau}(s_i|s_{i-1})$$

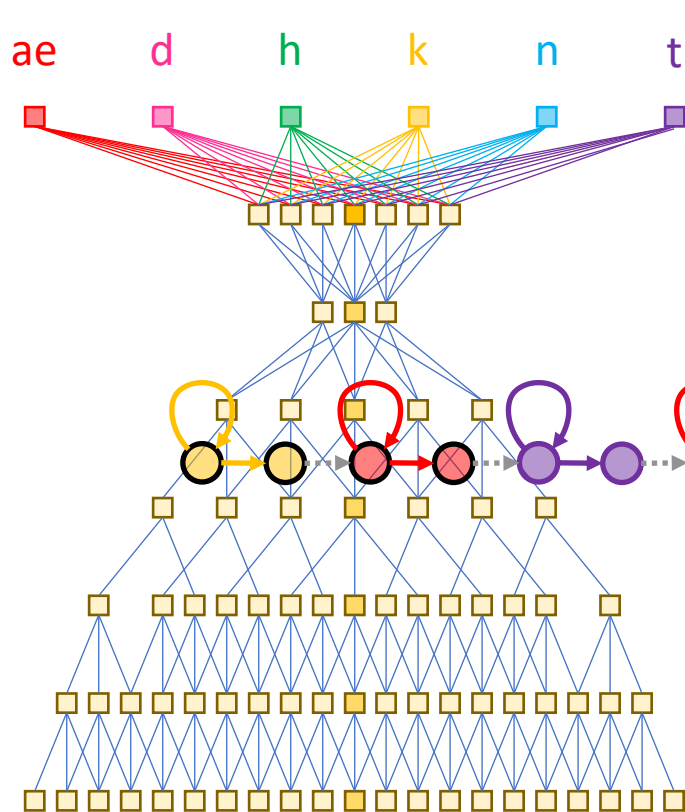
# CTC

Connectionist Temporal Classification

Phoneme  
HMMs



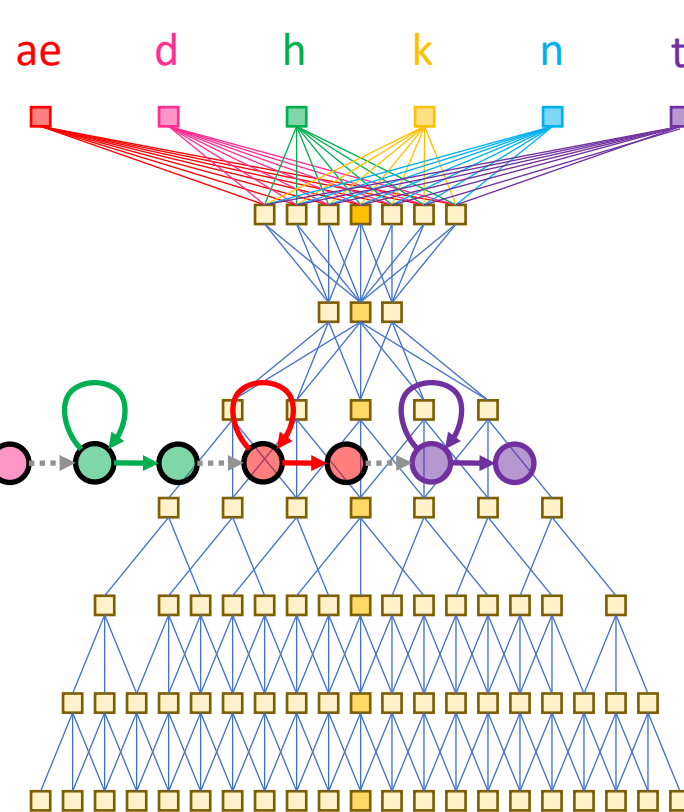
Phoneme  
Posterior Probabilities



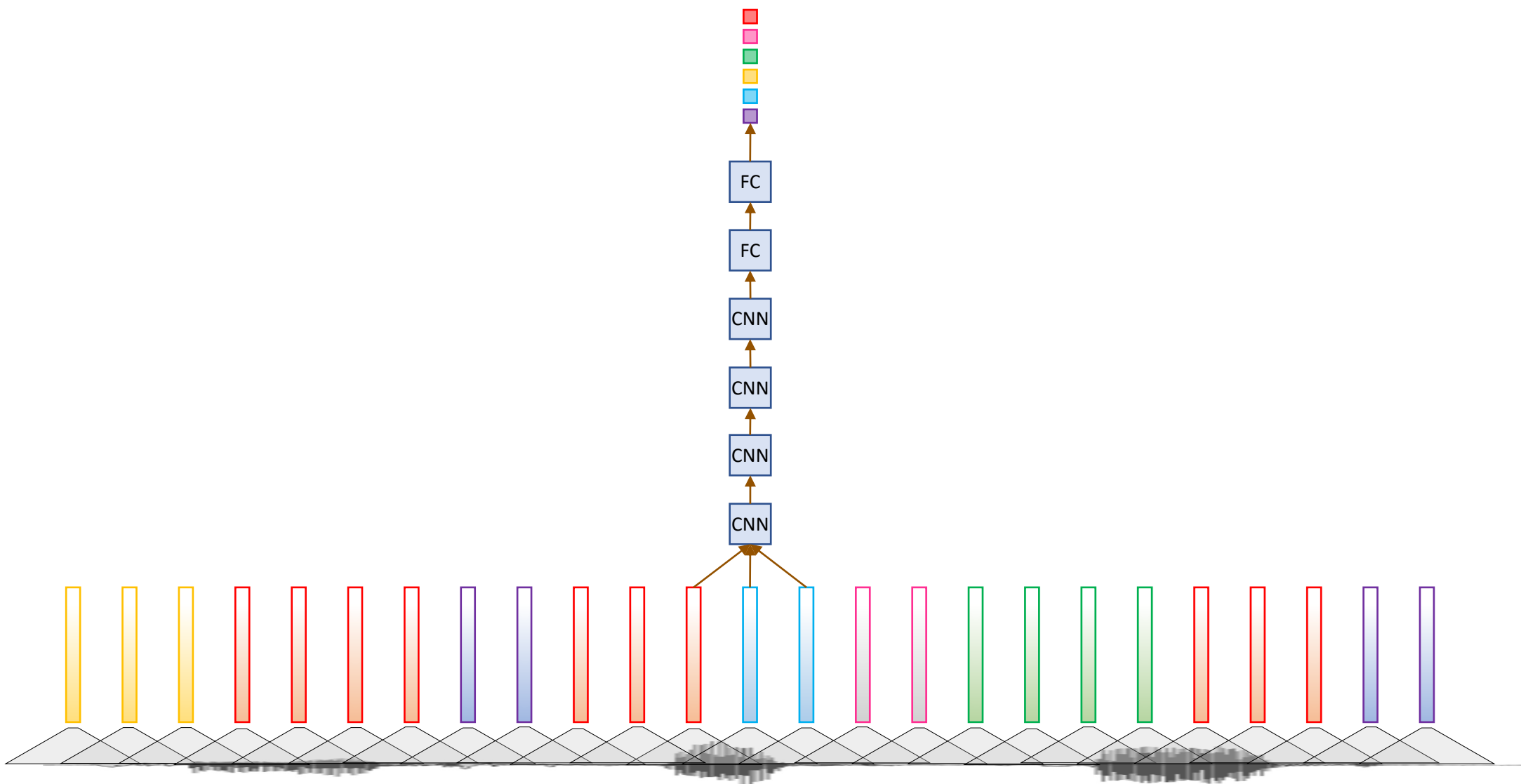
$p(\phi|y_i)$

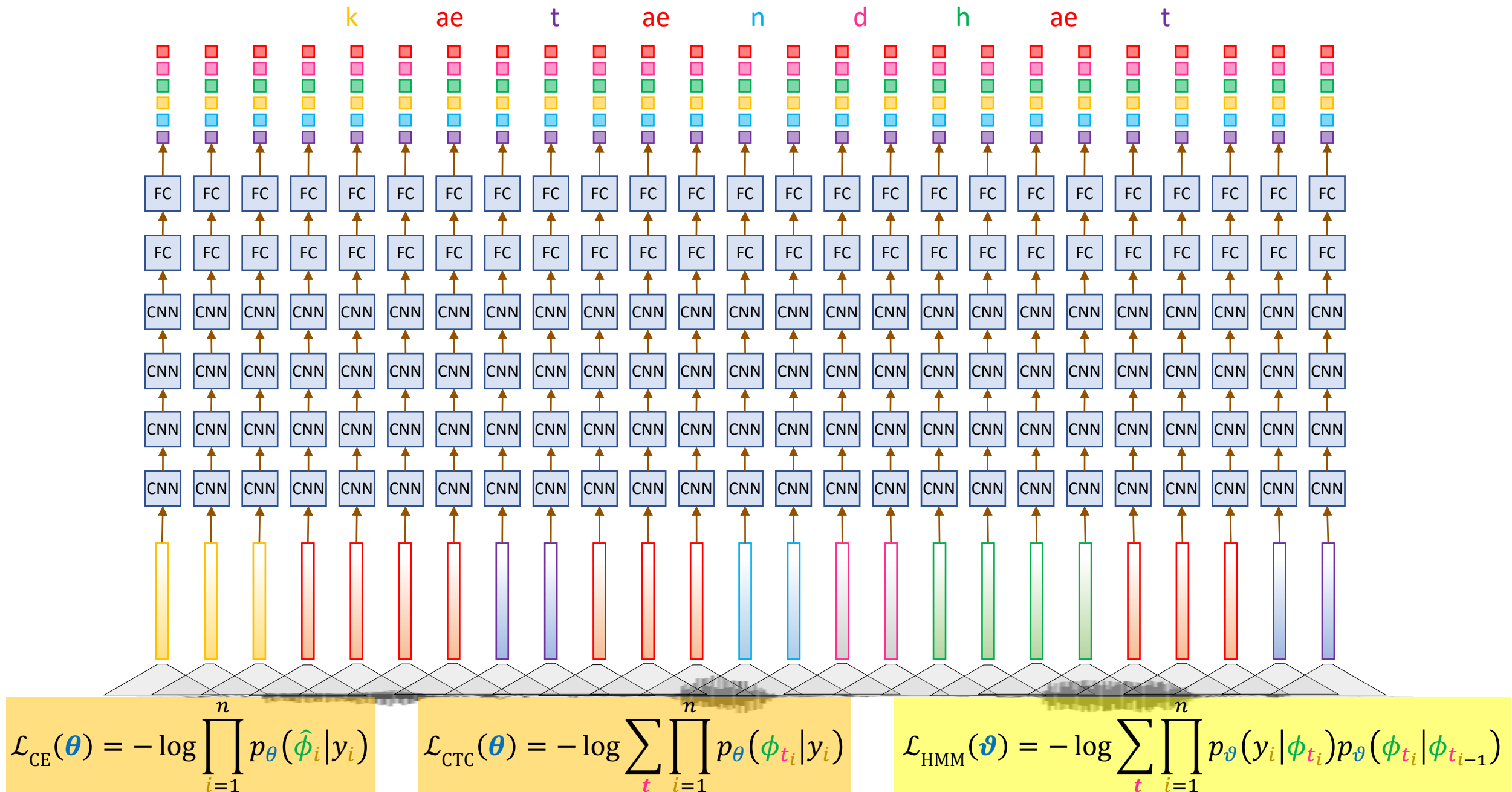
“cat and hat”

$$p(y_i|\phi) = \frac{p(\phi|y_i)p(y_i)}{p(\phi)} \propto \frac{p(\phi|y_i)}{p(\phi)}$$

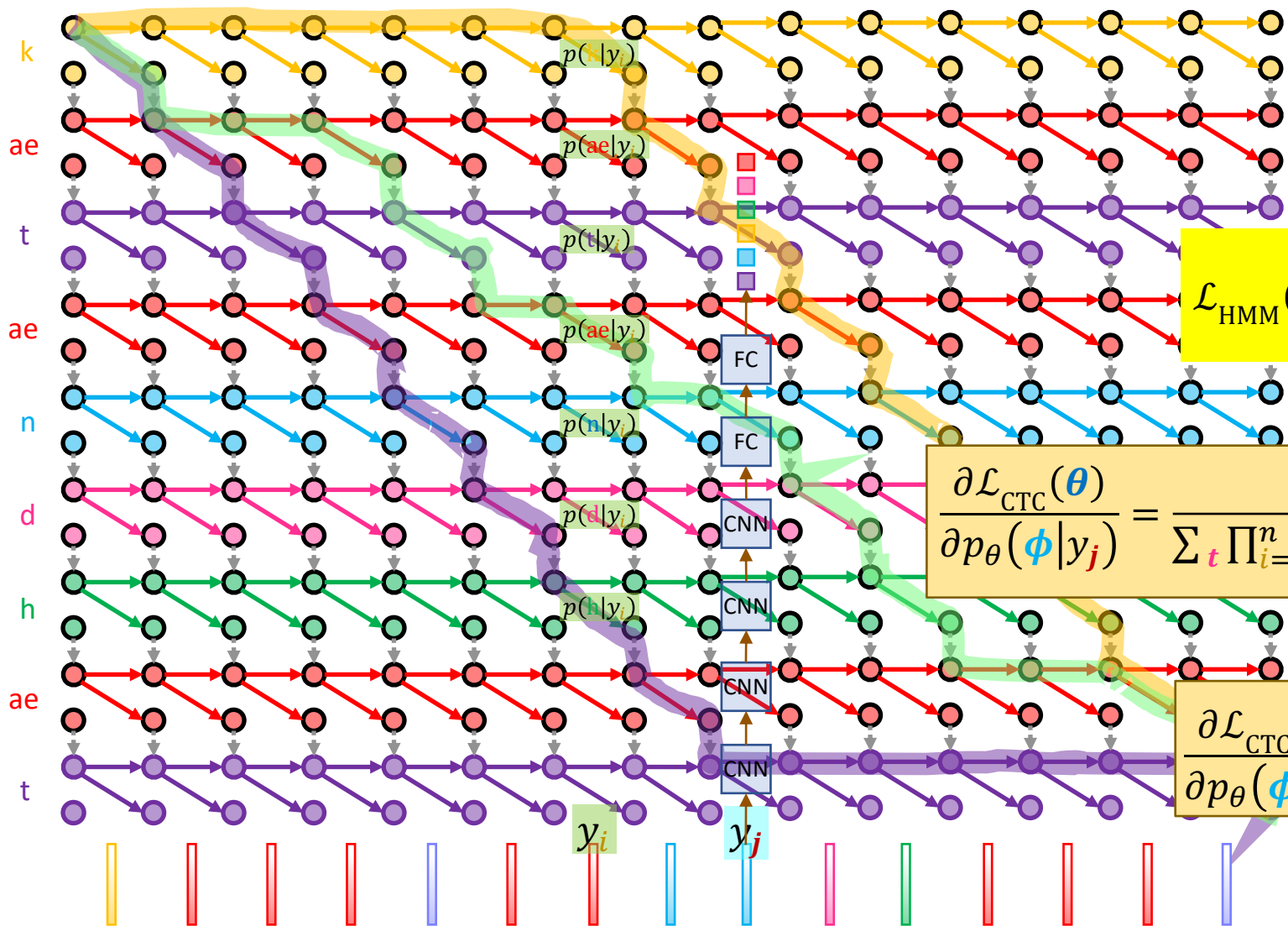
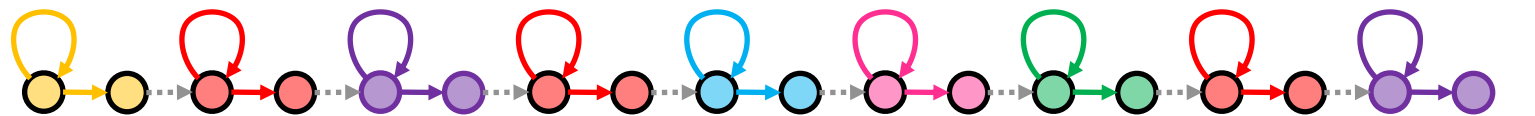


$$\mathcal{L}_{\text{CE}}(\theta) = -\log \prod_{i=1}^n p_{\theta}(\hat{\phi}_i|y_i) = -\sum_{i=1}^n \log p_{\theta}(\hat{\phi}_i|y_i)$$





Calculating the CTC loss for “cat and hat”



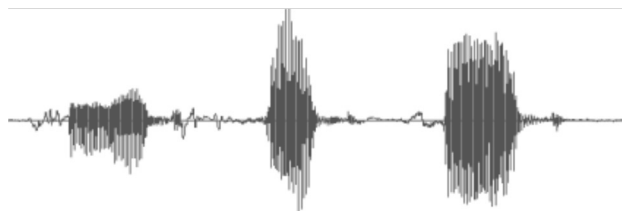
Calculating the gradient of the CTC loss

$$\mathcal{L}_{\text{CTC}}(\theta) = -\log \sum_{\mathbf{t}} \prod_{i=1}^n p_{\theta}(\phi_{t_i} | y_i)$$

$$\mathcal{L}_{\text{HMM}}(\theta) = -\log \sum_{\mathbf{t}} \prod_{i=1}^n p_{\theta}(y_i | \phi_{t_i}) p_{\theta}(\phi_{t_i} | \phi_{t_{i-1}})$$

$$\frac{\partial \mathcal{L}_{\text{CTC}}(\theta)}{\partial p_{\theta}(\phi | y_j)} = \frac{-1}{\sum_{\mathbf{t}} \prod_{i=1}^n p_{\theta}(\phi_{t_i} | y_i)} \sum_{\mathbf{t}: \phi_{t_j} = \phi} \frac{1}{p_{\theta}(\phi | y_j)} \prod_{i=1}^n p_{\theta}(\phi_{t_i} | y_i)$$

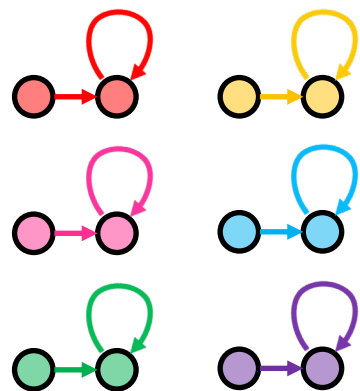
$$\frac{\partial \mathcal{L}_{\text{CTC}}(\theta)}{\partial p_{\theta}(\phi | y_j)} = -\frac{1}{p_{\theta}(\phi | y_j)} \frac{\sum_{\mathbf{t}: \phi_{t_j} = \phi} \prod_{i=1}^n p_{\theta}(\phi_{t_i} | y_i)}{\sum_{\mathbf{t}} \prod_{i=1}^n p_{\theta}(\phi_{t_i} | y_i)}$$



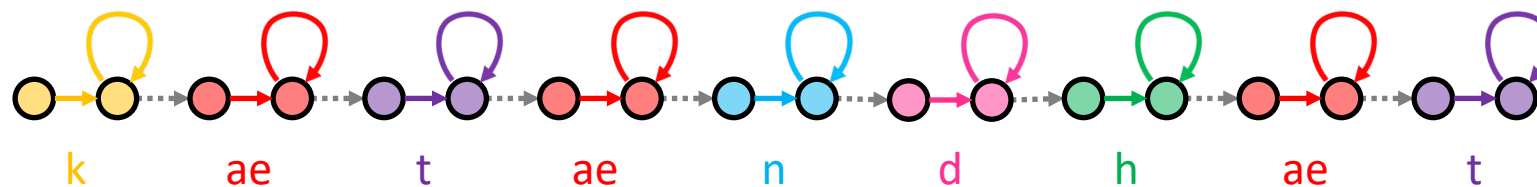
cat      and      hat

and	ae	n	d
cat	k	ae	t
hat	h	ae	t

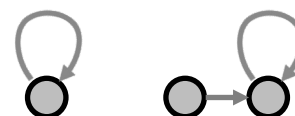
Phoneme  
HMMs



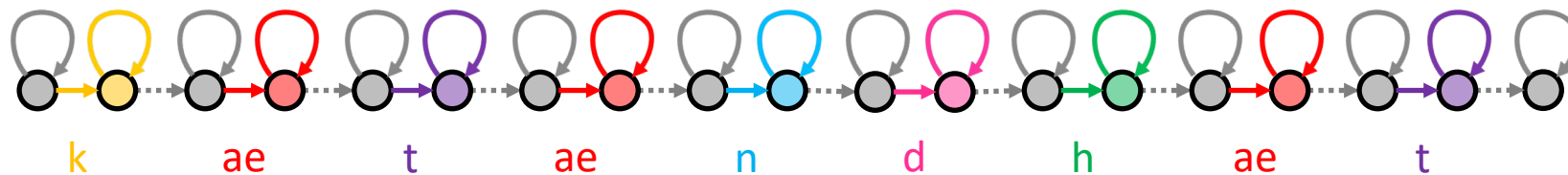
Composite HMM for “cat and hat”



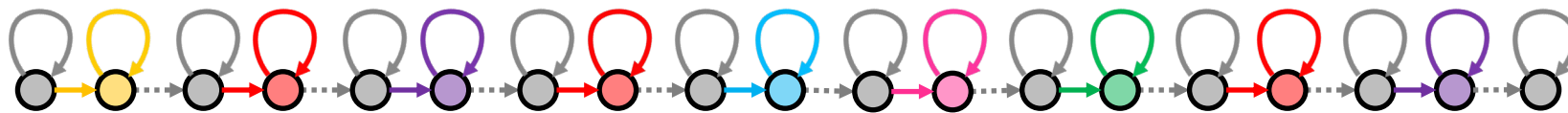
The CTC “Blank” Symbol ( $\beta$ )



FSA of permissible CTC strings for “cat and hat”





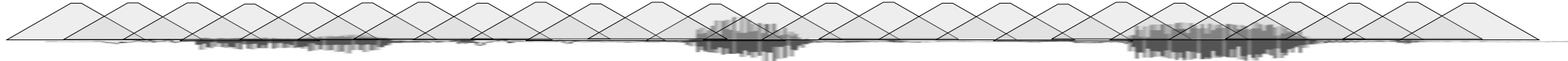
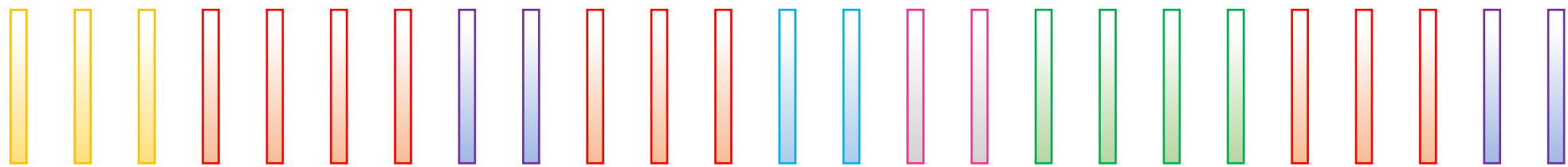
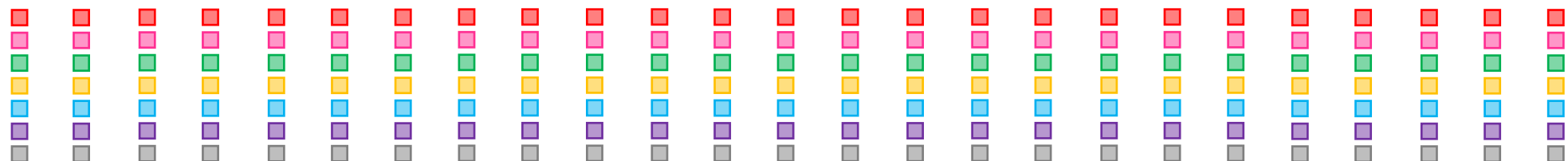


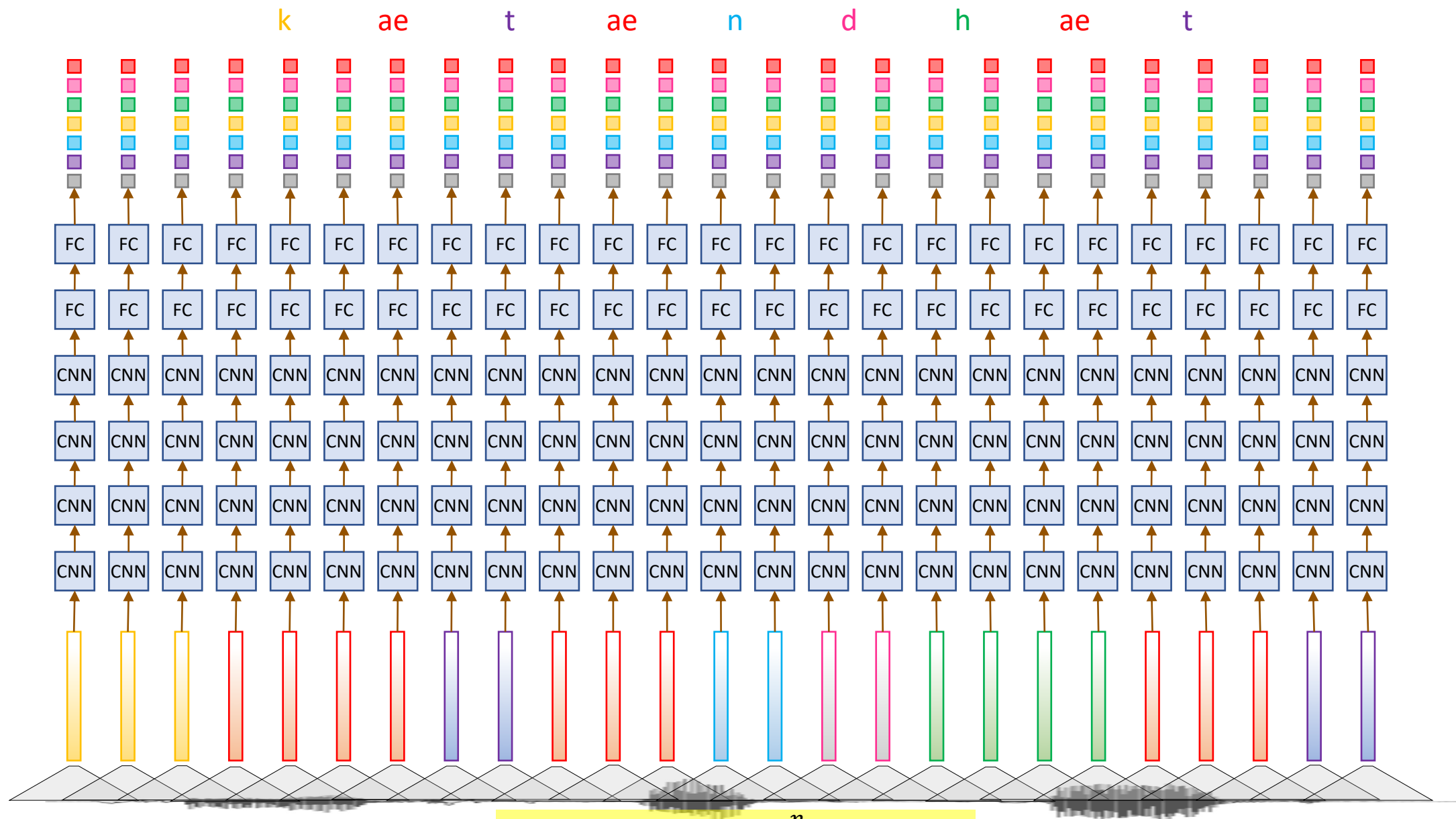
HMM State Sequences

k	k	k	k	k	k	k	k	k	k	k	k	k	k	k	k	k	ae	t	ae	n	d	h	ae	t
:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:
k	k	k	ae	ae	ae	ae	t	t	ae	ae	ae	n	n	d	d	h	h	h	h	ae	ae	ae	t	t
:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:	:
k	ae	t	ae	n	d	h	ae	ae	t	t	t	t	t	t	t	t	t	t	t	t	t	t	t	t

CTC Symbol Sequences

$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$	$\beta$
---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------	---------





$$\mathcal{L}_{\text{CTC}}(\theta) = -\log \sum_{\mathbf{t}} \prod_{i=1}^n p_{\theta}(\phi_{t_i} | y_i)$$

# Neural Speech Recognition without HMMs (aka End2End ASR)

“Purely” CTC-Based Speech Recognition Architectures

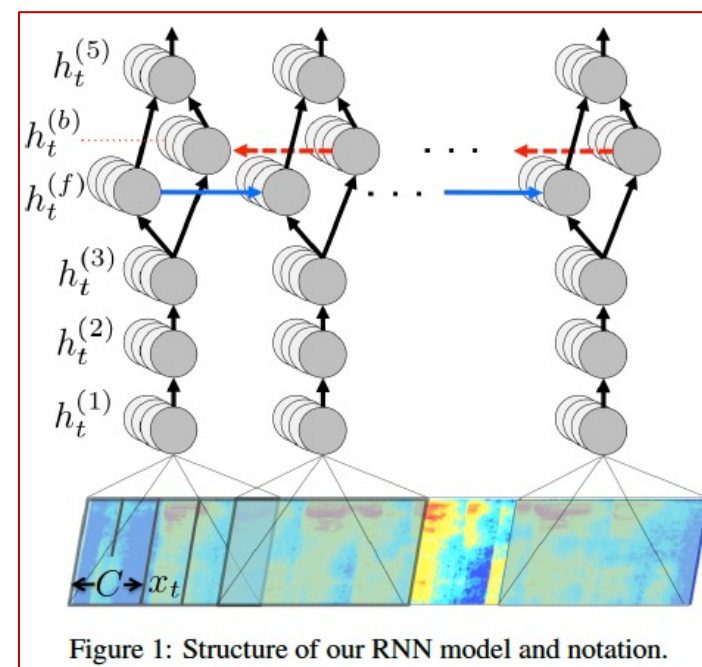
# End-to-End Speech Recognition

arXiv

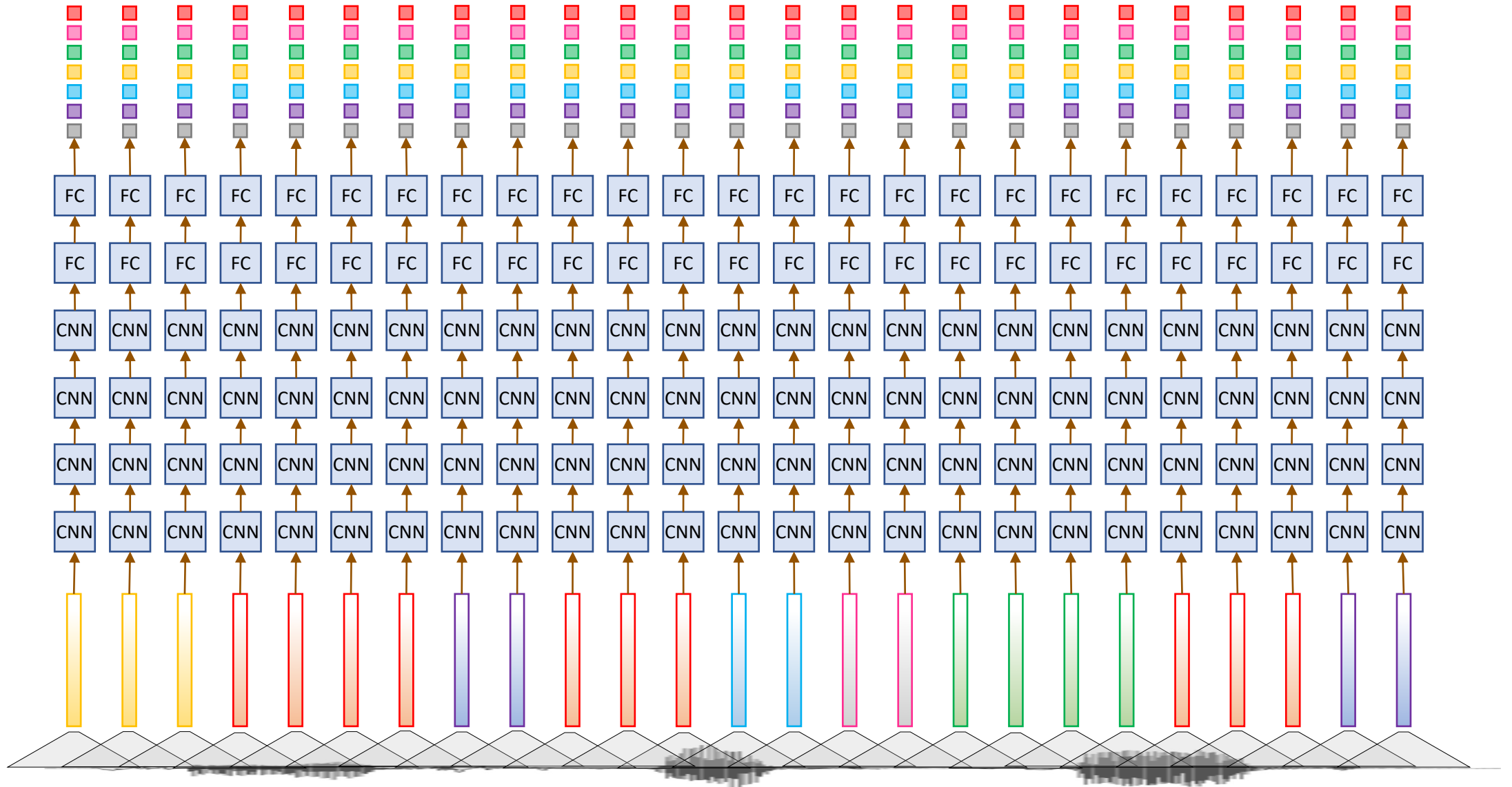
## Deep Speech: Scaling up end-to-end speech recognition

Awni Hannun\*, Carl Case, Jared Casper, Bryan Catanzaro, Greg Diamos, Erich Elsen, Ryan Prenger, Sanjeev Satheesh, Shubho Sengupta, Adam Coates, Andrew Y. Ng

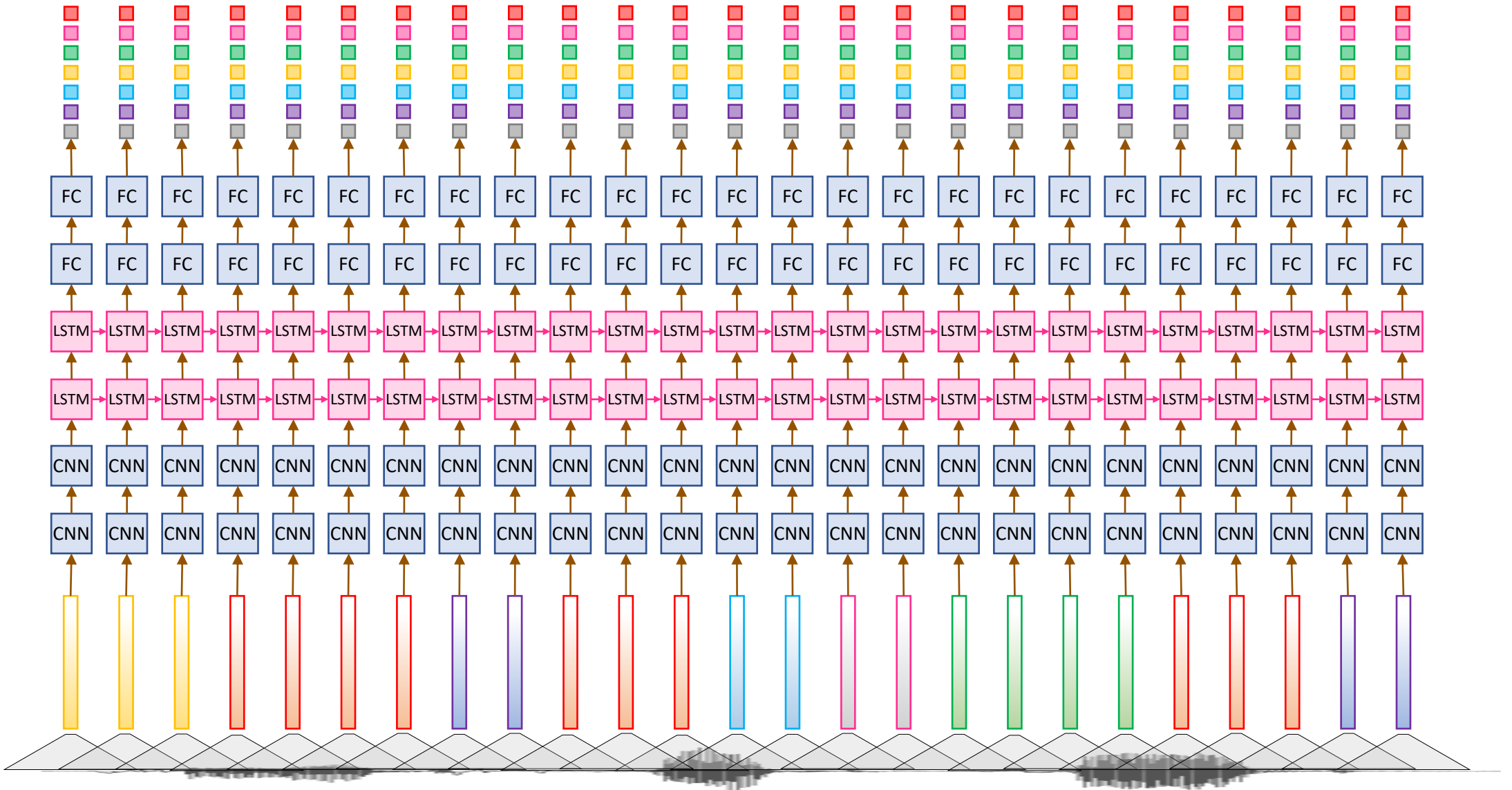
Baidu Research – Silicon Valley AI Lab



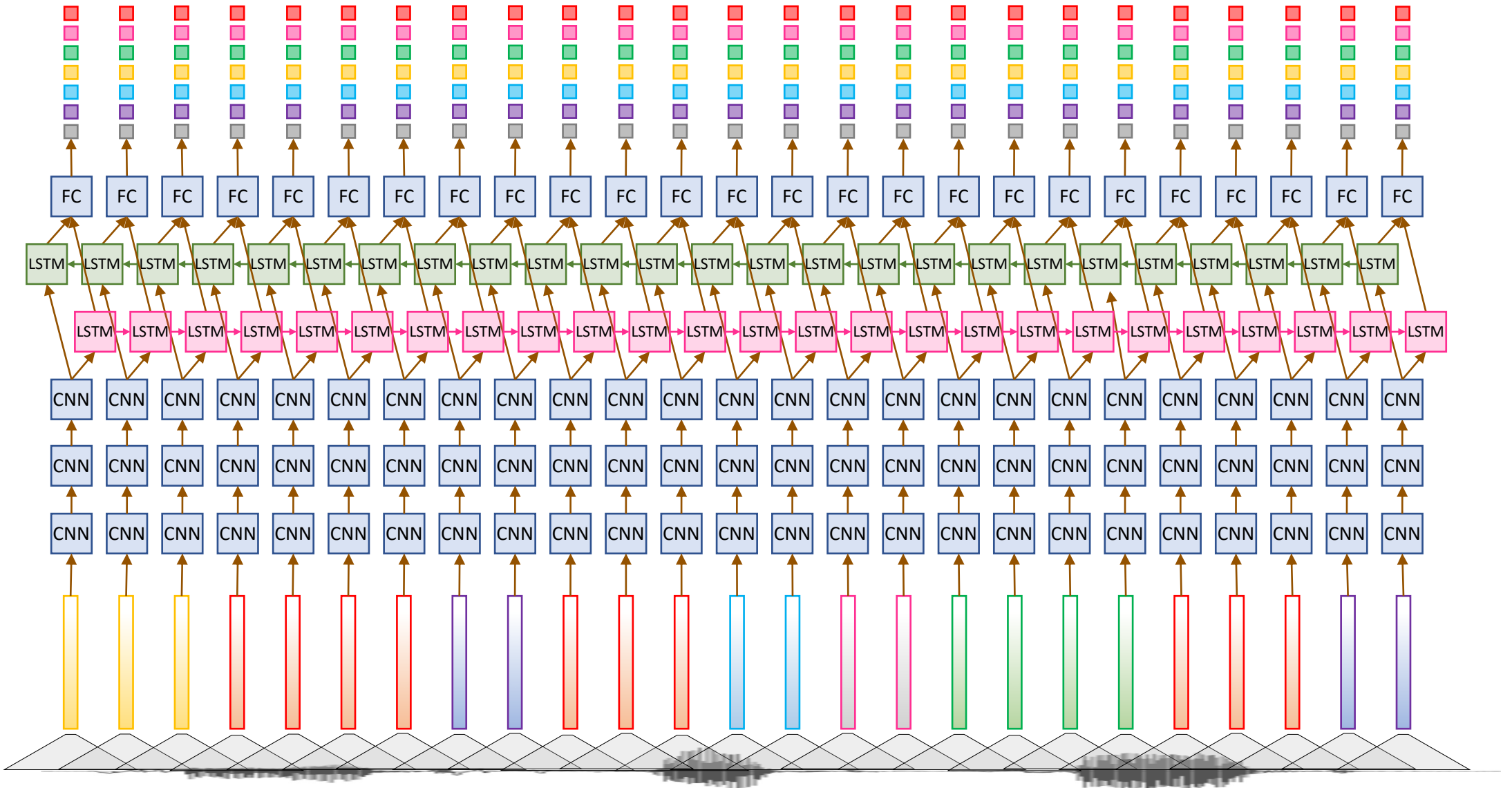
# The CNN Architecture



# The CNN+LSTM Architecture



# A Bidirectional LSTM Architecture (Deep Speech)








# Transducers

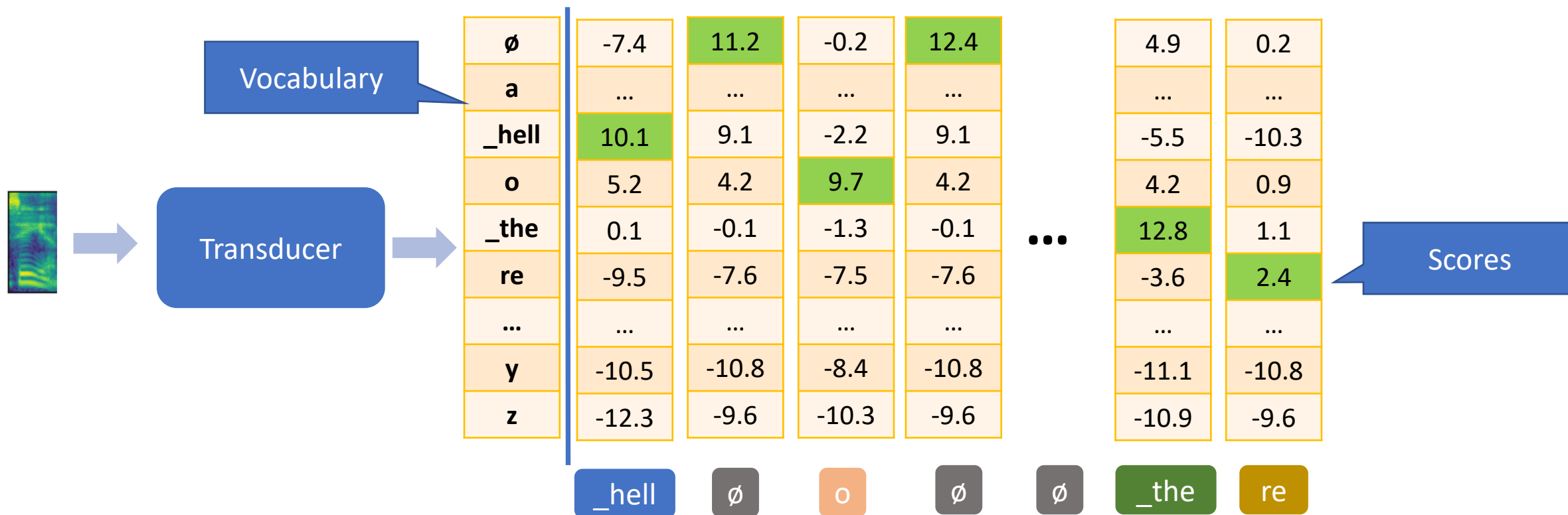
(RNN-T)



# ASR with Transducers






## The Transducer: Inference

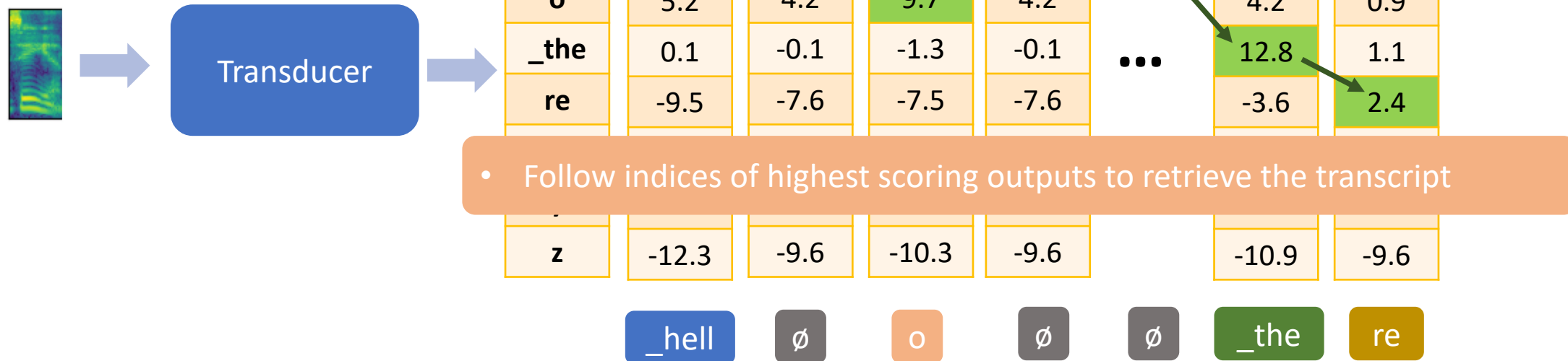
- 1 Tokenize transcripts   $\rightarrow$     
  - Transducer produces sequences of scores over a finite set of tokens, called a Vocabulary
  - Vocabulary includes a special symbol,  $\emptyset$ , which means move to the next audio sample



# ASR with Transducers

## The Transducer: Inference

- 1 Tokenize transcripts   $\rightarrow$     
  - Transducer produces sequences of scores over a finite set of tokens, called a Vocabulary
  - Vocabulary includes a special symbol,  $\emptyset$ , which means move to the next audio sample



# ASR with Transducers

## The Transducer: Training

- The posterior of **one** alignment,  $a$ , of a particular transcript  $\propto$  product of scores,  $a_i$ , along alignment path
- Many possible alignments

\_hell ø ø o ø \_the ø ø re  
\_hell ø o ø ø \_the ø ø re  
\_hell o ø ø ø \_the re ø ø  
\_hell ø o ø ø ø ø \_the re

Possible alignments of the 4 token sequence  
\_hell o \_the re  
with a 9-frame speech utterance

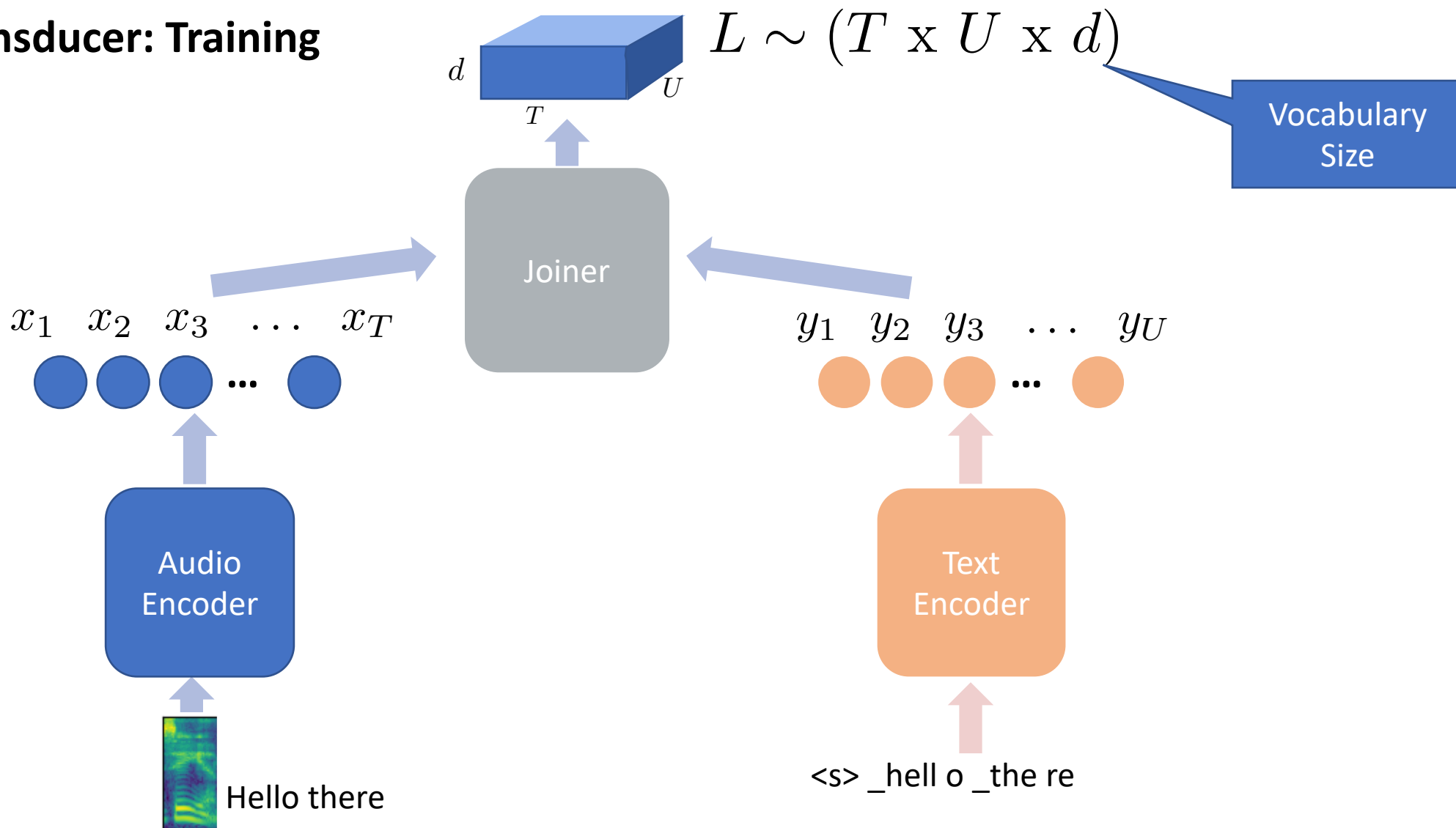
- The posterior of a particular transcript is computed by maximizing over all possible alignments
- Using normalized scores, i.e. via Softmax gives ...

$$p(y|x) = \sum_a \prod_{i=1}^{|a|} a_i(y_0^{i-1}, x)$$

No conditional independence assumption

# ASR with Transducers

## The Transducer: Training



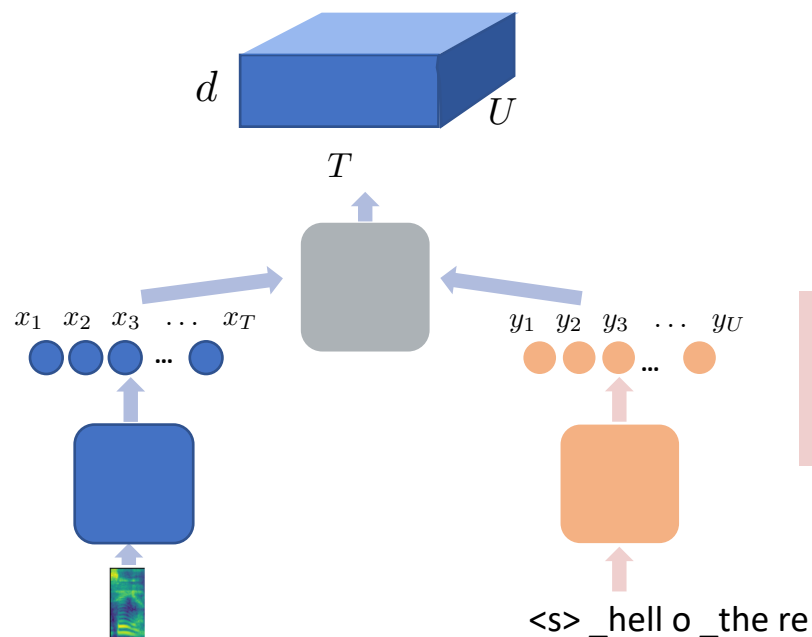
# ASR with Transducers

## The Transducer: Training Alignment

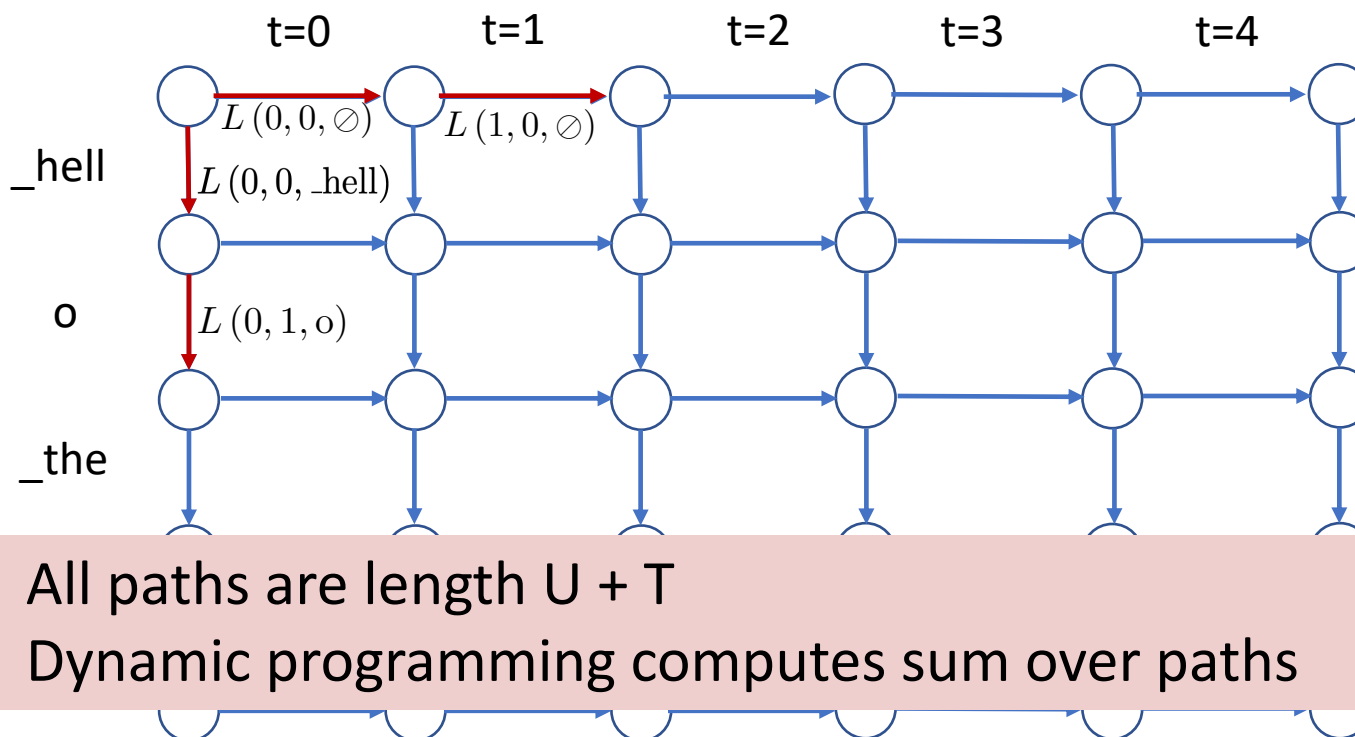
Score of aligning frame,  $t$ , to  
position,  $u$  with label,  $l$  is

$$L(t, u, l)$$

$$L \sim (T \times U \times d)$$



Align input frames with output tokens using  $\emptyset$



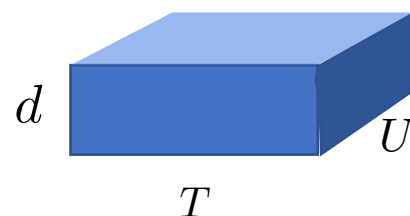
- All paths are length  $U + T$
- Dynamic programming computes sum over paths

# ASR with Transducers

---

## The Transducer: Memory Hungry

$$L \sim (T \times U \times d)$$



- When the vocabulary size is large, i.e., characters in Mandarin  $\sim 10,000$
- When the sequence length is long

Memory Explodes!

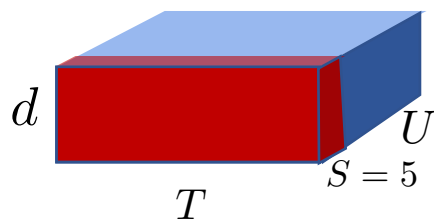
15s utterance, 30 character transcript, 10,000 types  $\rightarrow$  0.45 GB with fp16 / utterance

# ASR with Transducers

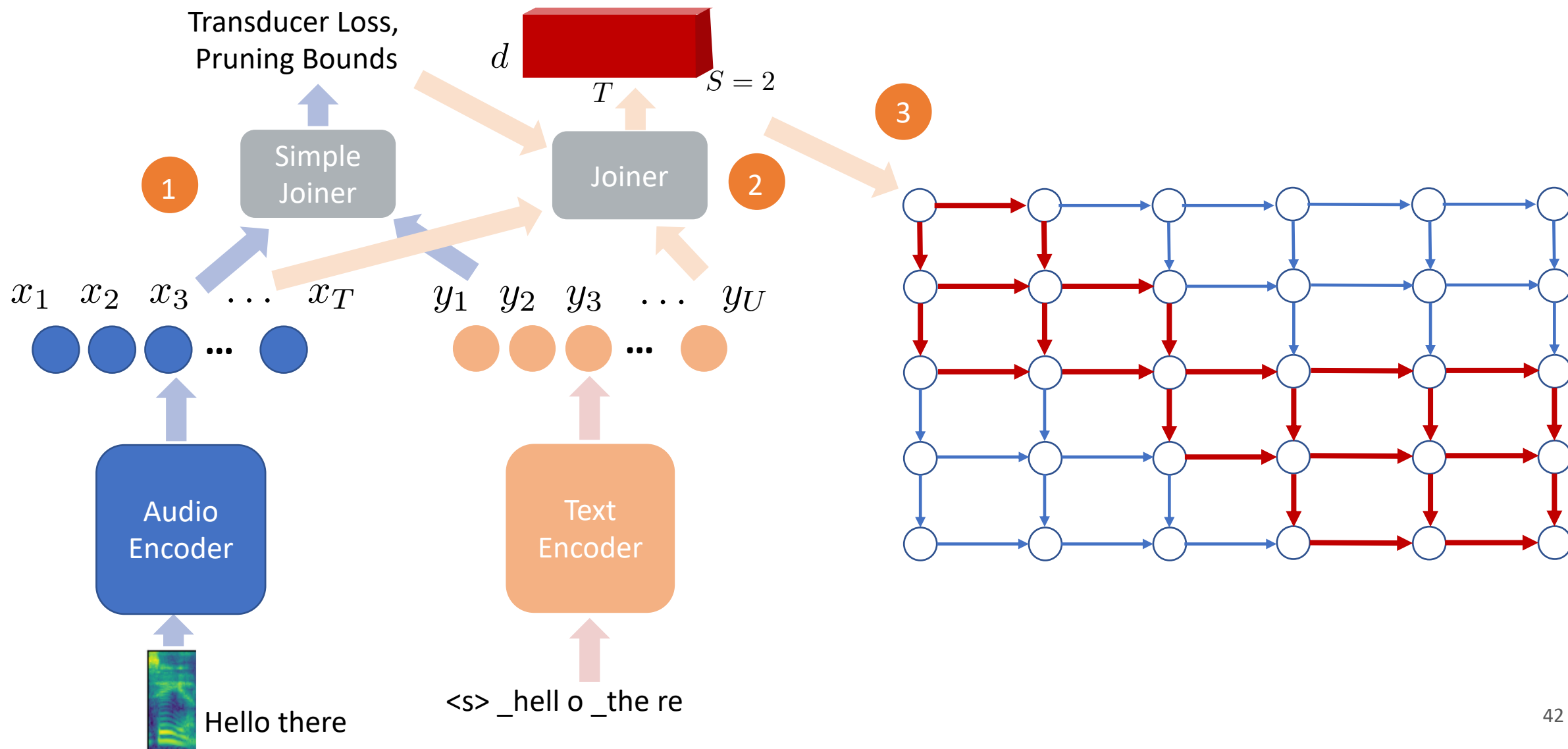
---

## The Transducer: Memory Hungry

- Solution: "Pruned RNN-T for fast, memory-efficient ASR training", Fangjun Kuang, Liyong Guo, Wei Kang, Long Lin, Mingshuang Luo, Zengwei Yao, Daniel Povey, Interspeech 2022
  - Only realize the tensor,  $L$  of scores for the top-k scoring positions at each time-step
  - Use a 2-pass method to first prune unlikely paths
  - The first pass uses a “simple” joiner that avoids realizing any large matrices
  - The new tensor of scores,  $\tilde{L}$  during the second pass has different dimensions



# ASR with Transducers

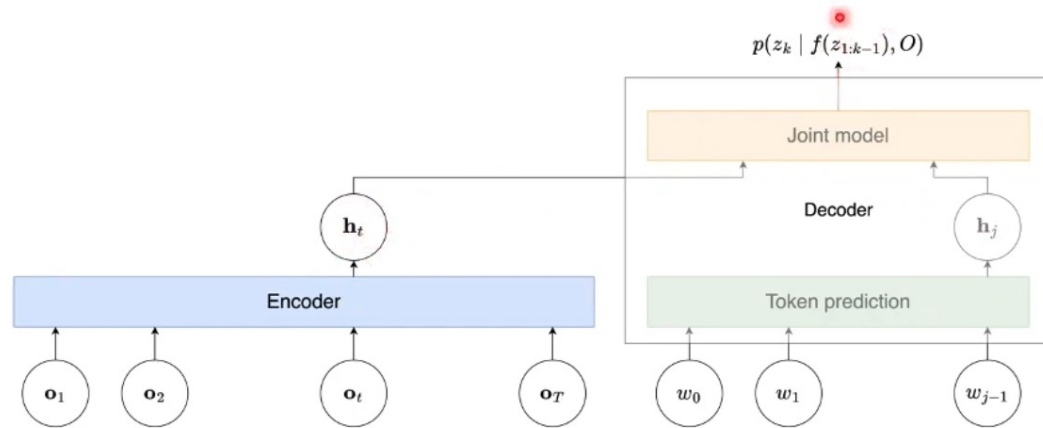




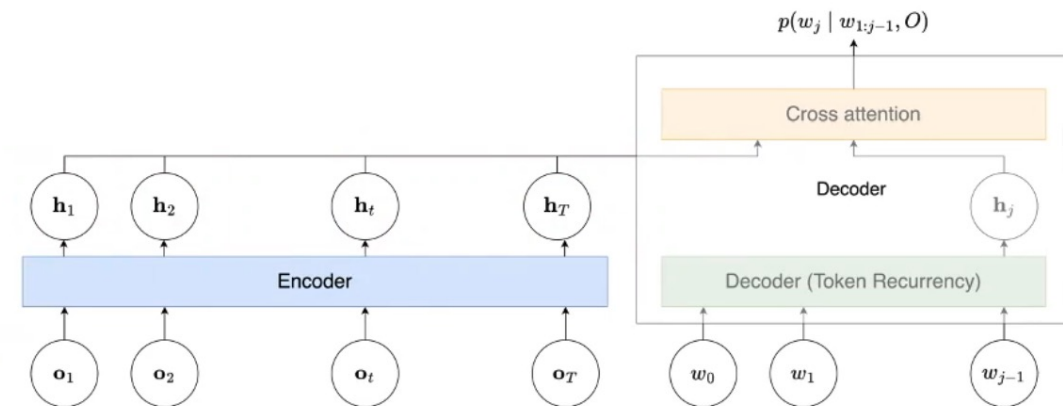
# Encoder Decoder

# Encoder-Decoder

- RNN-Transducer



- Attention



- Both do not use the explicit conditional independence assumptions
- Use of  $\mathbf{h}_t$  only versus  $\mathbf{h}_t$  for  $1, \dots, t, \dots, T \rightarrow$  Get the future information more efficiently
- The output depends on input frame  $t$  or not  $\rightarrow$  Online property

<https://www.youtube.com/watch?v=IVc46-aBnzM&list=PLfVqr2l0FG-u7chWKPQMDoT0o-l2ejxeK&index=17>

- Sequence Generation is **auto-regressive**

$$p(\mathbf{y}_0^{N-1}) = \prod_{i=0}^{N-1} p(\mathbf{y}_i | \mathbf{y}_0^{i-1})$$

- Beam-search
- Almost identical to NMT models

# Encoder-Decoder

- Whisper

## Robust Speech Recognition via Large-Scale Weak Supervision

Alec Radford<sup>\*1</sup> Jong Wook Kim<sup>\*1</sup> Tao Xu<sup>1</sup> Greg Brockman<sup>1</sup> Christine McLeavey<sup>1</sup> Ilya Sutskever<sup>1</sup>

