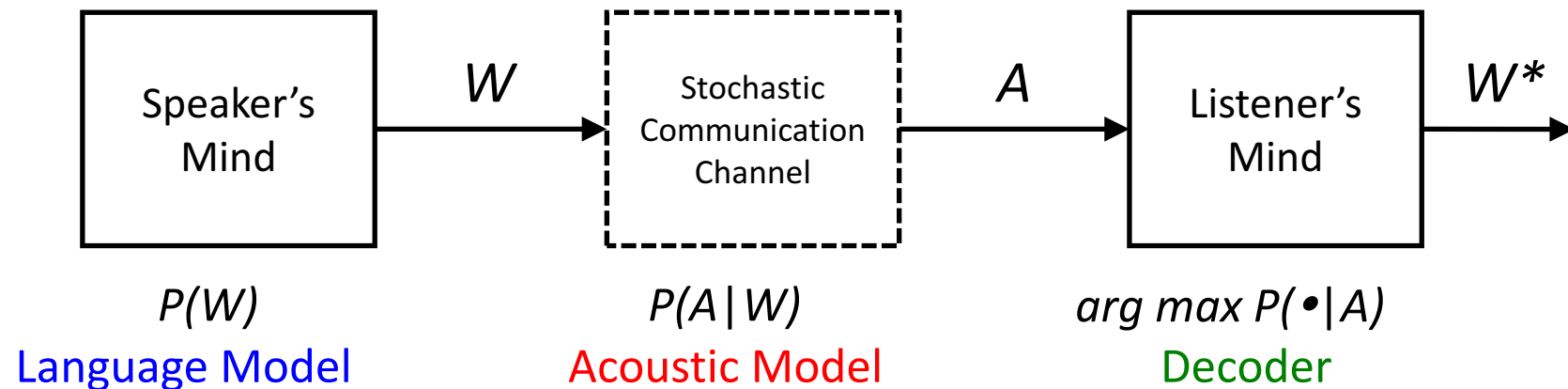
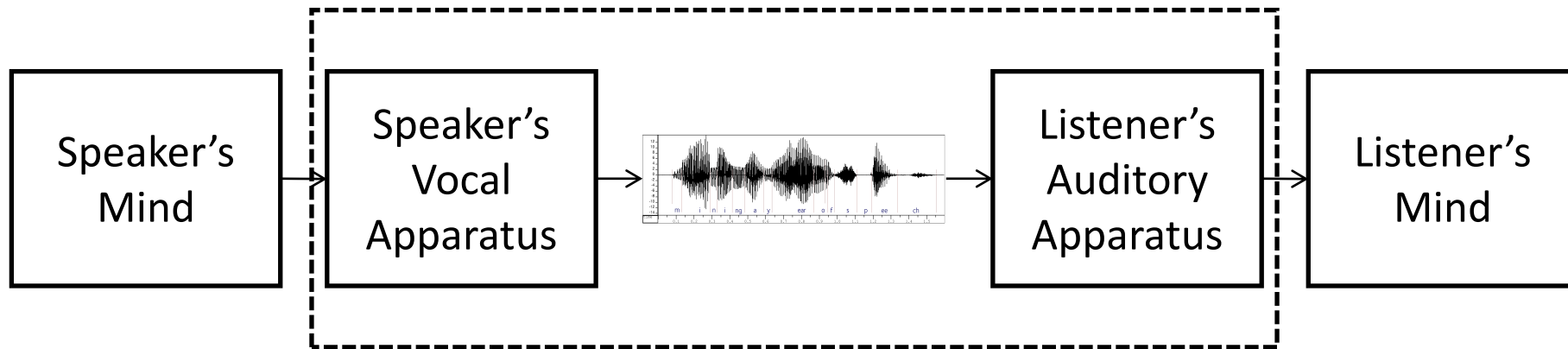


Neural Methods in Automatic Speech Recognition

IntroHLT Lecture on Nov 2, 2021

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



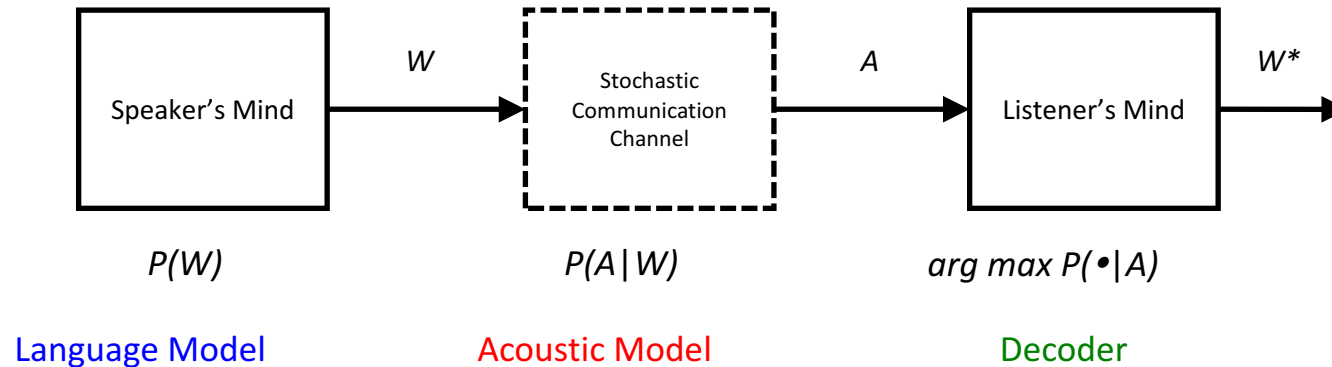
Hidden Markov models are popular as acoustic models

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

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

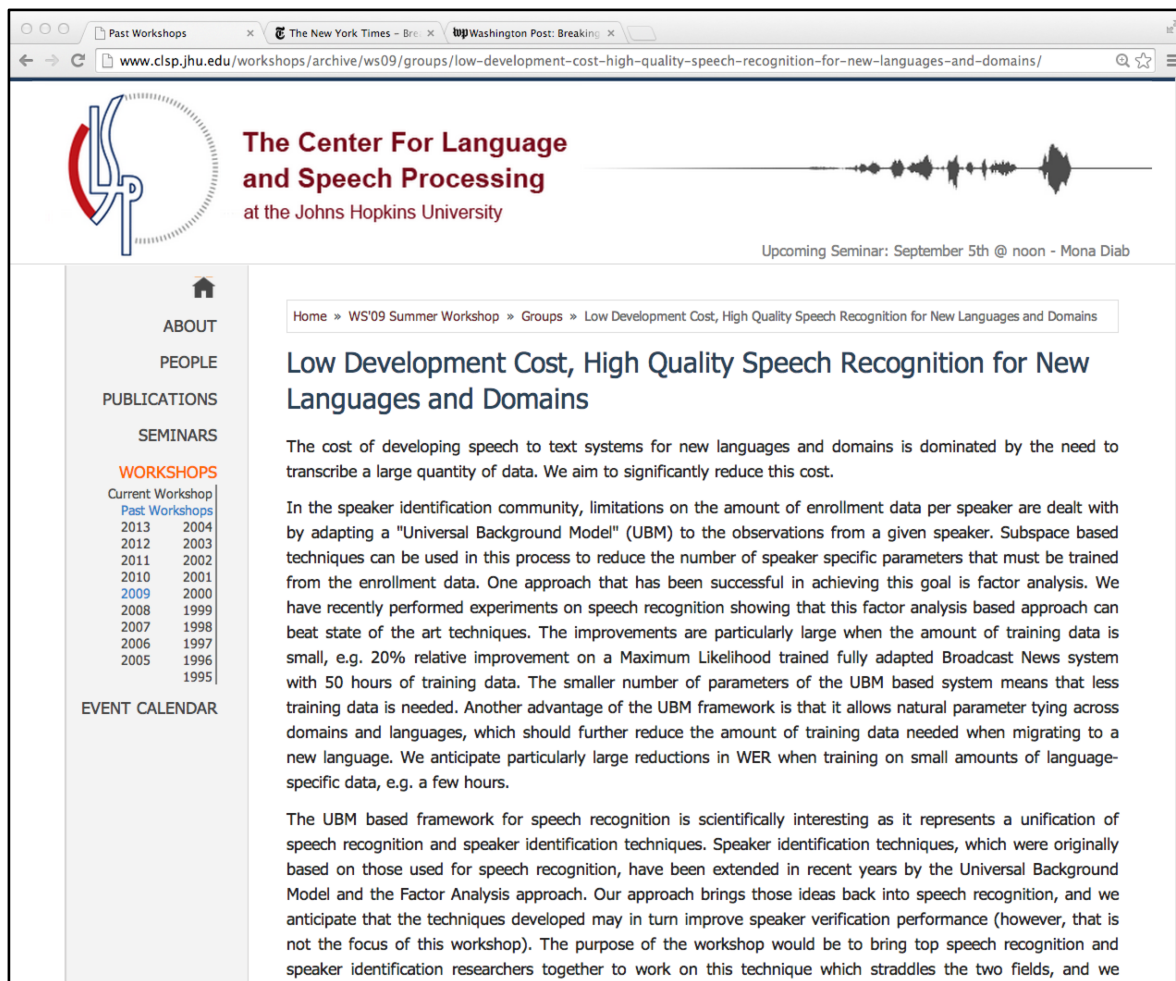
$$\begin{aligned}\widehat{\mathbf{W}} &= \arg \max_{\mathbf{W}} P(\mathbf{A} | \mathbf{W})P(\mathbf{W}) \\ &= \arg \max_{\mathbf{W}} \sum_{\mathbf{S} \in \mathcal{S}(\mathbf{W})} P(\mathbf{A} | \mathbf{S})P(\mathbf{S})P(\mathbf{W}) \\ &\approx \arg \max_{\mathbf{W}} \max_{\mathbf{S} \in \mathcal{S}(\mathbf{W})}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

The ASR Landscape in 2009



- Commercial providers had proprietary algorithms and software
- Academic software tools were mostly good only for research
 - Usually not scalable for deployment
 - Often required licensing for commercial use
- Significant barriers to entry existed for start-ups and small(er) labs
 - Algorithms were complex to understand and implement
 - Significant “black art” beyond the algorithms themselves was needed

Kaldi was born in the Summer of 2009



The screenshot shows the homepage of the Center for Language and Speech Processing (CLSP) at Johns Hopkins University. The page features a navigation menu on the left with categories like ABOUT, PEOPLE, PUBLICATIONS, SEMINARS, WORKSHOPS, and EVENT CALENDAR. The main content area is titled "Low Development Cost, High Quality Speech Recognition for New Languages and Domains" and includes a breadcrumb trail: Home » WS'09 Summer Workshop » Groups » Low Development Cost, High Quality Speech Recognition for New Languages and Domains. The text discusses the cost of developing speech-to-text systems and the use of a Universal Background Model (UBM) for speaker identification. A sidebar on the left lists "Past Workshops" from 2013 to 1995, with 2009 highlighted in blue. An audio waveform is visible in the header area.

The Center For Language and Speech Processing
at the Johns Hopkins University

Upcoming Seminar: September 5th @ noon - Mona Diab

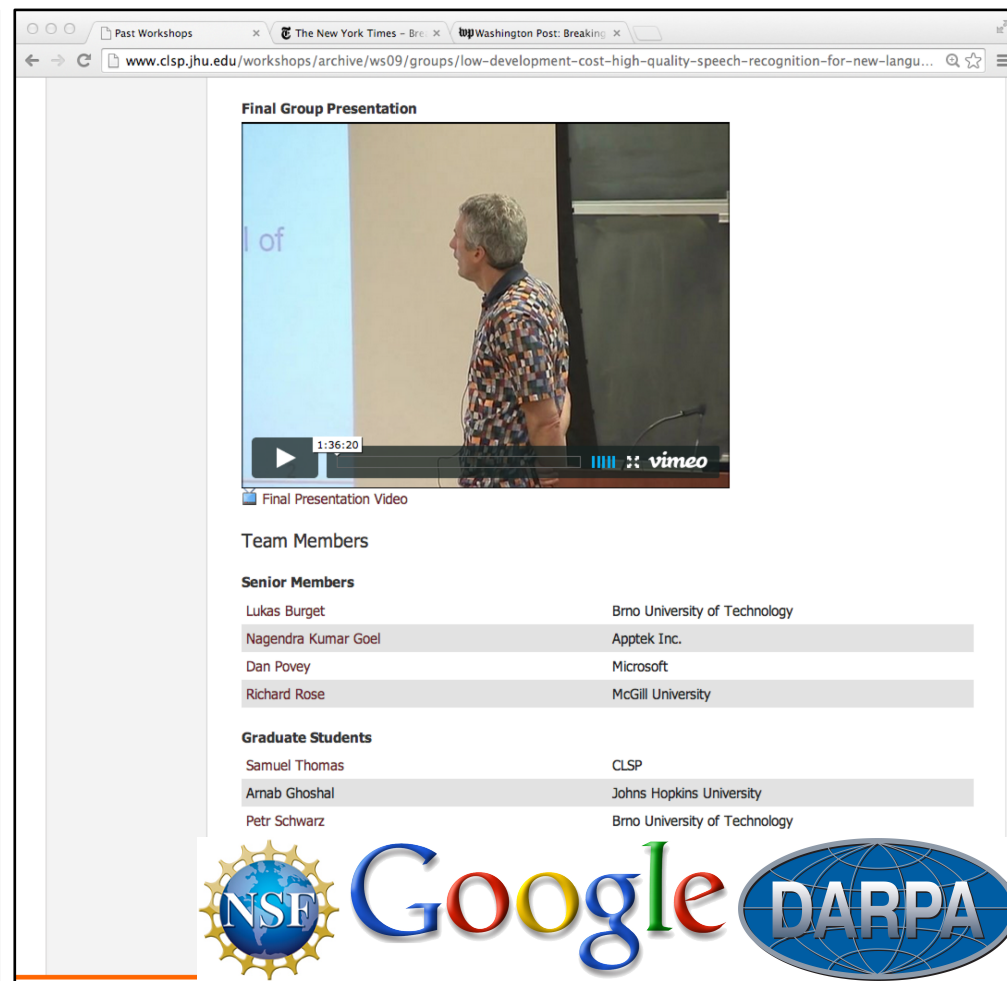
Home » WS'09 Summer Workshop » Groups » Low Development Cost, High Quality Speech Recognition for New Languages and Domains

Low Development Cost, High Quality Speech Recognition for New Languages and Domains

The cost of developing speech to text systems for new languages and domains is dominated by the need to transcribe a large quantity of data. We aim to significantly reduce this cost.

In the speaker identification community, limitations on the amount of enrollment data per speaker are dealt with by adapting a "Universal Background Model" (UBM) to the observations from a given speaker. Subspace based techniques can be used in this process to reduce the number of speaker specific parameters that must be trained from the enrollment data. One approach that has been successful in achieving this goal is factor analysis. We have recently performed experiments on speech recognition showing that this factor analysis based approach can beat state of the art techniques. The improvements are particularly large when the amount of training data is small, e.g. 20% relative improvement on a Maximum Likelihood trained fully adapted Broadcast News system with 50 hours of training data. The smaller number of parameters of the UBM based system means that less training data is needed. Another advantage of the UBM framework is that it allows natural parameter tying across domains and languages, which should further reduce the amount of training data needed when migrating to a new language. We anticipate particularly large reductions in WER when training on small amounts of language-specific data, e.g. a few hours.

The UBM based framework for speech recognition is scientifically interesting as it represents a unification of speech recognition and speaker identification techniques. Speaker identification techniques, which were originally based on those used for speech recognition, have been extended in recent years by the Universal Background Model and the Factor Analysis approach. Our approach brings those ideas back into speech recognition, and we anticipate that the techniques developed may in turn improve speaker verification performance (however, that is not the focus of this workshop). The purpose of the workshop would be to bring top speech recognition and speaker identification researchers together to work on this technique which straddles the two fields, and we



The screenshot shows a video player on the CLSP website. The video is titled "Final Group Presentation" and shows a man in a patterned shirt presenting in front of a screen. The video player includes a play button, a progress bar at 1:36:20, and a Vimeo logo. Below the video, there is a "Team Members" section with two columns of names and affiliations.

Final Group Presentation

Final Presentation Video

Team Members

Senior Members	
Lukas Burget	Brno University of Technology
Nagendra Kumar Goel	Apptek Inc.
Dan Povey	Microsoft
Richard Rose	McGill University
Graduate Students	
Samuel Thomas	CLSP
Arnab Ghoshal	Johns Hopkins University
Petr Schwarz	Brno University of Technology

NSF Google DARPA

Kaldi: Legendary Ethiopian goatherd who discovered coffee

A screenshot of the GitHub repository page for 'kaldi-asr / kaldi'. The page shows the repository name, navigation links, and a list of recent commits. The repository has 693 watchers, 8k stars, and 3.6k forks. It also has 140 issues, 72 pull requests, and 1 project. The latest commit is by danpovey, adding an option on github for kaldi10-related issues (#3796), committed 2 days ago. Other recent commits include adding an issue template and upgrading TensorFlow RNN to 2.0.

github.com

Why GitHub? Enterprise Explore Marketplace Pricing Search Sign in Sign up

kaldi-asr / kaldi Watch 693 Star 8k Fork 3.6k

Code Issues 140 Pull requests 72 Projects 1 Wiki Security Insights

This is the official location of the Kaldi project. <http://kaldi-asr.org>

kaldi c-plus-plus cuda shell speech-recognition speech-to-text speaker-verification speaker-id speech

8,815 commits 17 branches 0 packages 0 releases 298 contributors View license

Branch: master New pull request Find file Clone or download

danpovey [misc] Add option on github for kaldi10-related issues (#3796) Latest commit 1f357ce 2 days ago

.github/ISSUE_TEMPLATE [misc] Add option on github for kaldi10-related issues (#3796) 2 days ago

cmake [build,src] Upgrade TensorFlow RNN to 2.0 (#3771) 8 days ago

Kaldi today: A community of researchers cooperatively advancing ASR

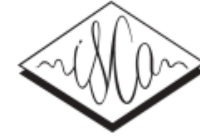
- C++ library, command line tools, several ASR recipes
 - Freely available via GitHub (Apache 2.0 license)
- Top performance in open benchmark tests
 - NIST OpenKWS 2014, IARPA ASPIRE 2015, NIST LoReHLT 2018, ..., MUCS 2021
- Widely adopted in academia and extensively used in industry
 - 300+ citations in 2014 (based on Google Scholar data)
 - 400+ citations in 2015, 600+ citations in 2016, 800+ citations in 2017, ...
 - Used (& developed further) by several US and non-US companies
- Kaldi “trunk” maintained by Dan Povey @ xiaomi and Jan Trmal @ jhu
 - Forks contain specializations by others (including other Hopkins researchers)

Staying ahead of the field: 2012-Today

- ASR technology is advancing very rapidly
 - Amazon, Apple, Baidu, Facebook, Google, Microsoft, Tencent, ...
- Kaldi leads the field with innovations, big and small, ...
 - **From SGMMs to DNNs** (2012)
 - From English to “low resource” languages (2013, IARPA BABEL)
 - Parallelization of DNN training (2014, Natural Gradient SGD)
 - From close-talking to far-field recordings (2015, IARPA ASpIRE)
 - Chain models: better, cheaper and faster (2016)
 - Backstitch: adversarial training reinterpreted (2017)
 - TDNN-F acoustic models (2018)
 - GPU acceleration of Viterbi decoding (2019)
- ... and tries to keep up with advances made by others

A paper appeared in September 2011 ...

INTERSPEECH 2011



**Conversational Speech Transcription
Using Context-Dependent Deep Neural Networks**

Frank Seide¹, Gang Li,¹ and Dong Yu²

¹Microsoft Research Asia, Beijing, P.R.C.

²Microsoft Research, Redmond, USA

{fseide, g

ICASSP 1988

Phoneme Recognition: Neural Networks vs.
Hidden Markov Models

A. Waibel

T. Hanazawa

G. Hinton *

K. Shikano

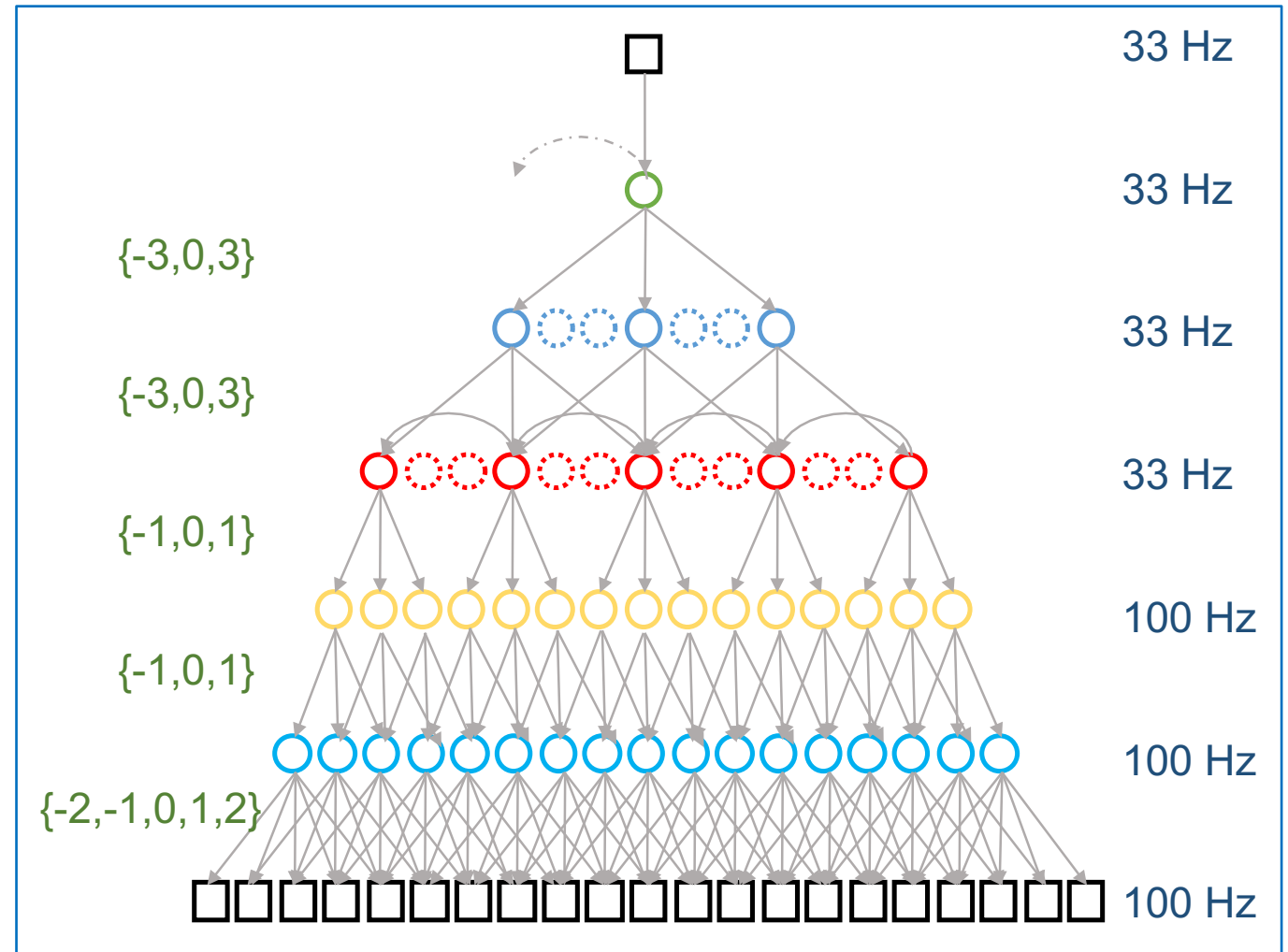
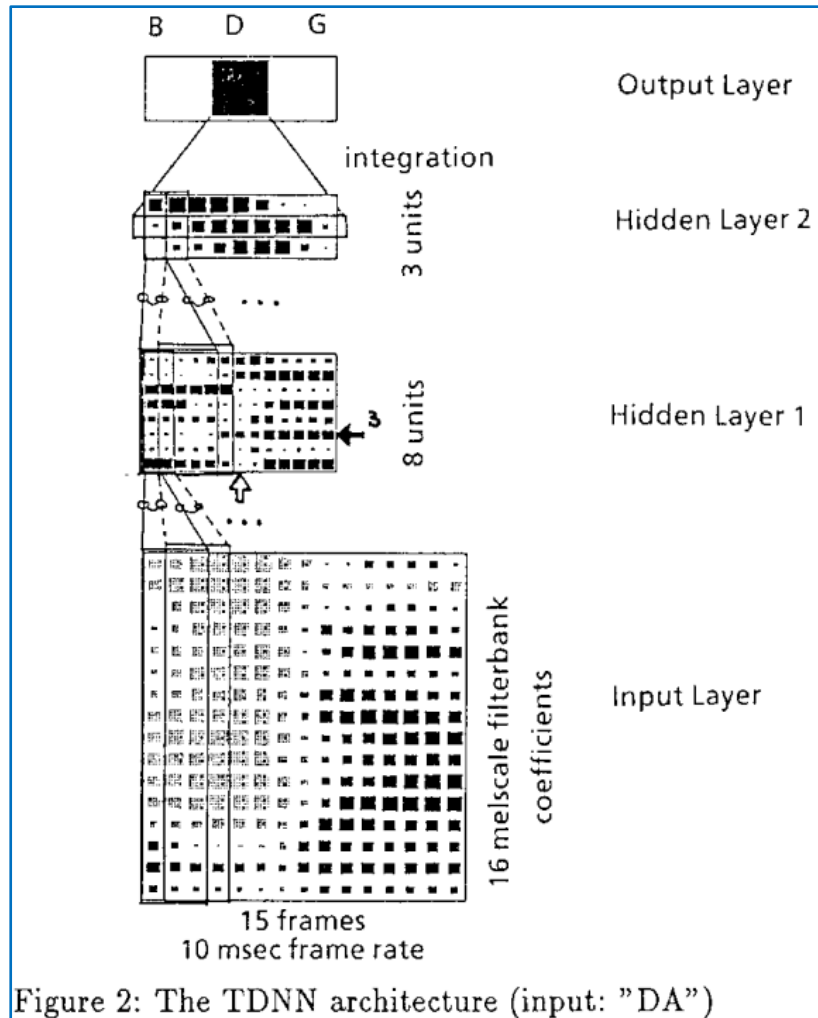
K. Lang †

ATR Interpreting Telephony Research Laboratories

*University of Toronto and Canadian Institute for Advanced Research

†Carnegie-Mellon University

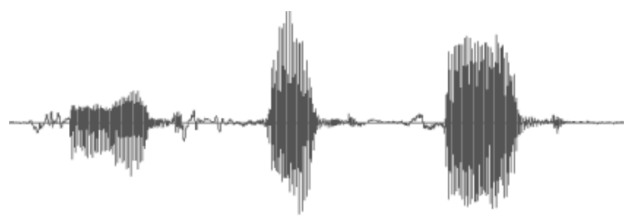
So, ~~a lot of~~ progress has been made since 1988



Acoustic Modeling with Deep Neural Networks for Hybrid ASR Systems

Repurposing Algorithms Developed for HMM-based Architectures

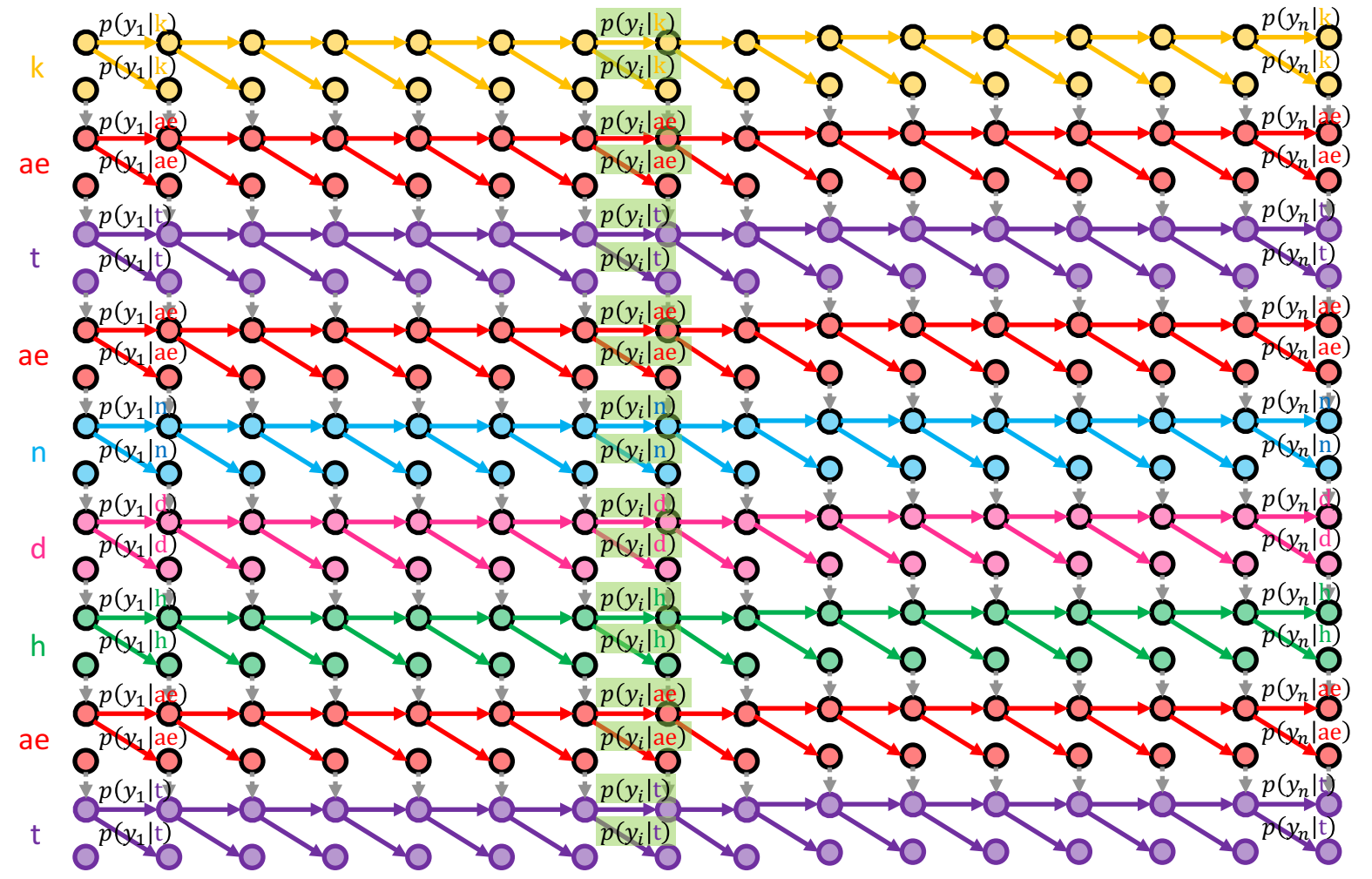
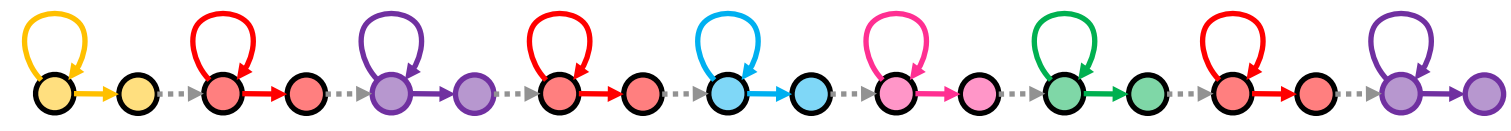
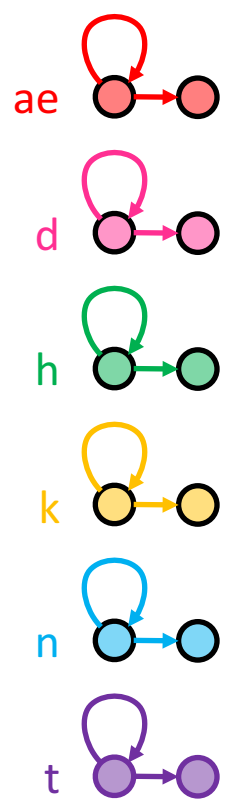
Composite HMM for "cat and hat"



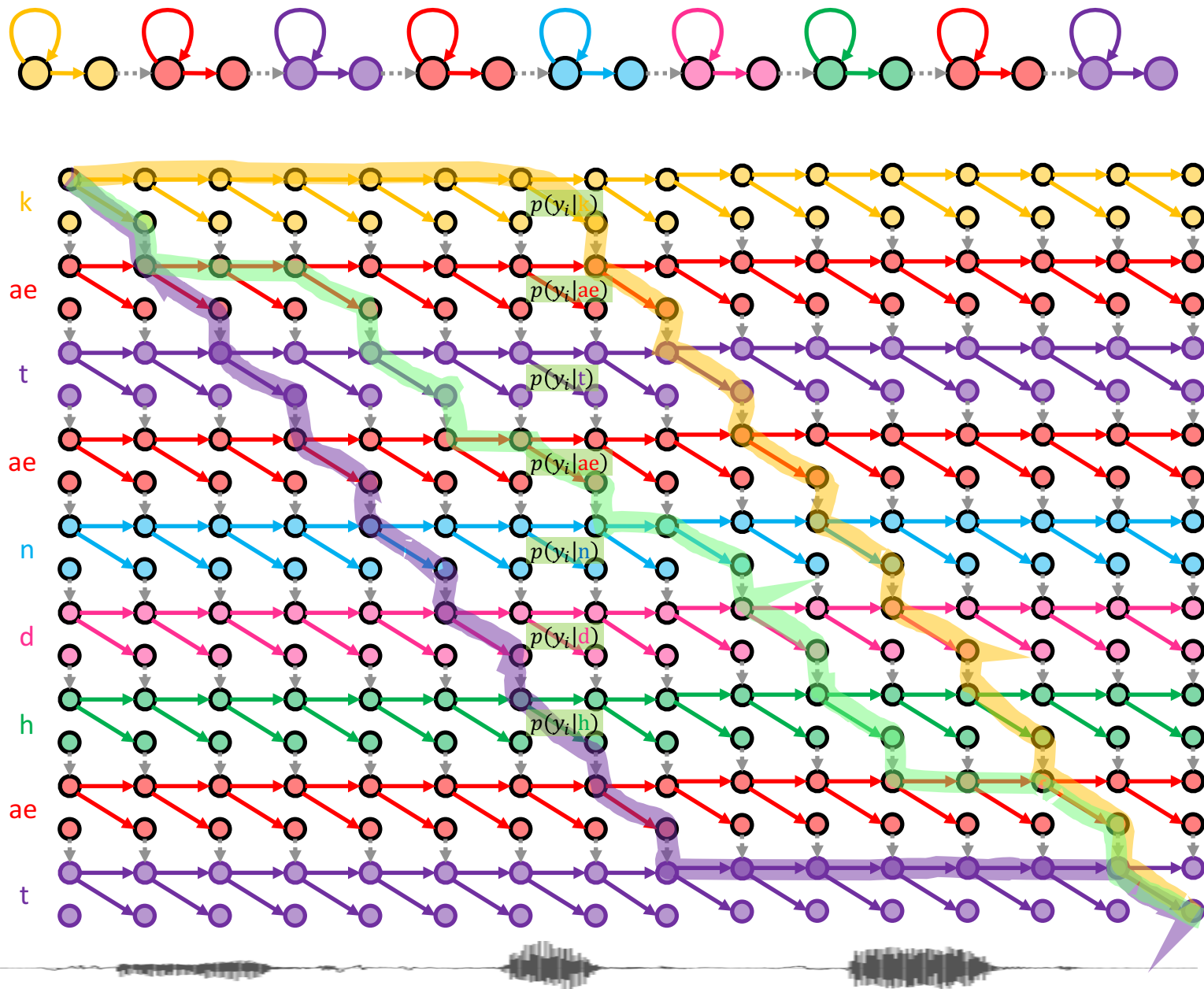
cat and hat

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

Phoneme HMMs



Composite HMM for "cat and hat"



"Forward" Algorithm

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

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

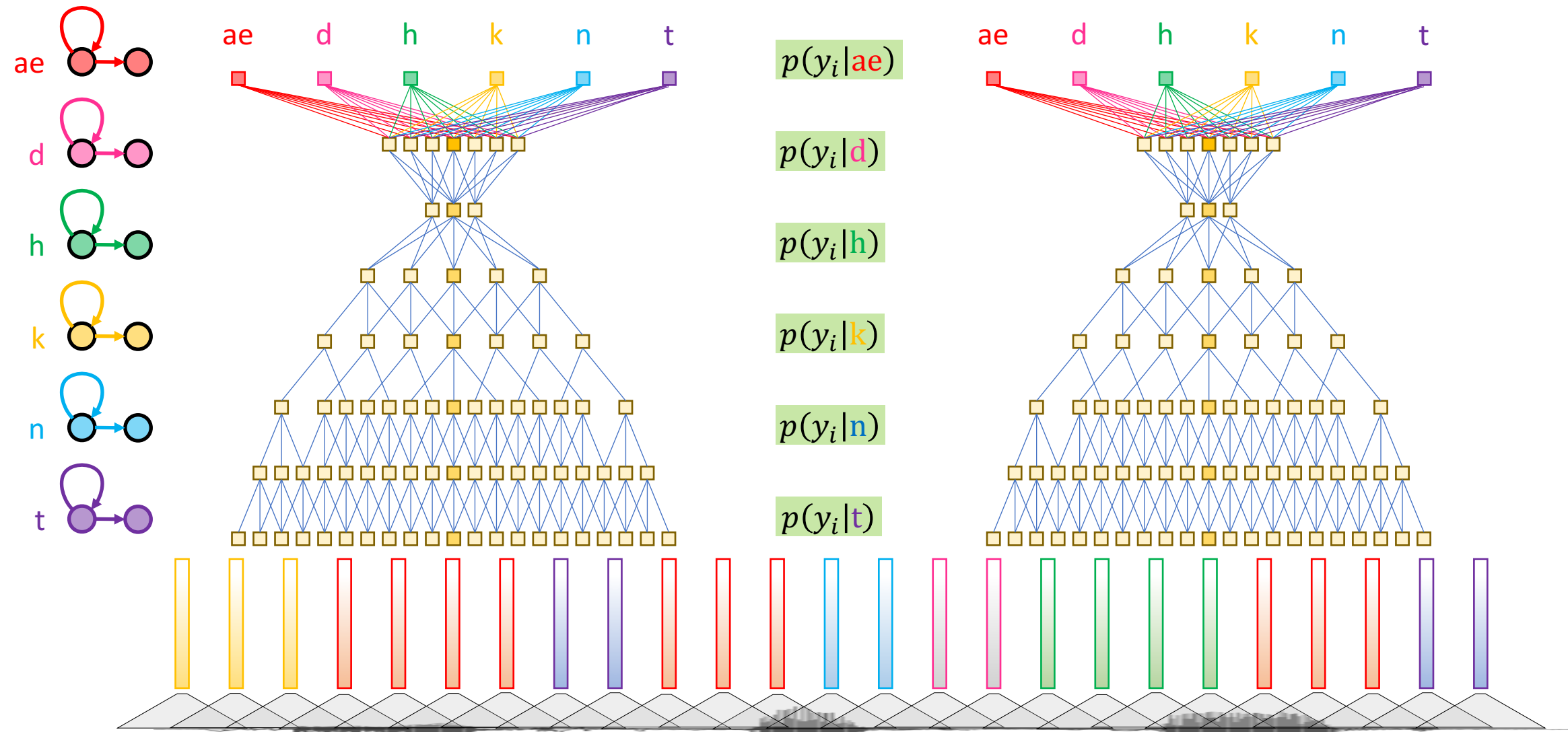
Viterbi Algorithm

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

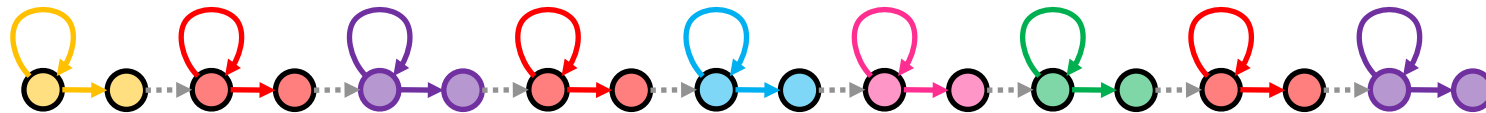
$$= \arg \max_{\mathbf{s} \in \mathcal{S}(\mathbf{w})} \frac{P(\mathbf{y}, \mathbf{s})}{P(\mathbf{y})}$$

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

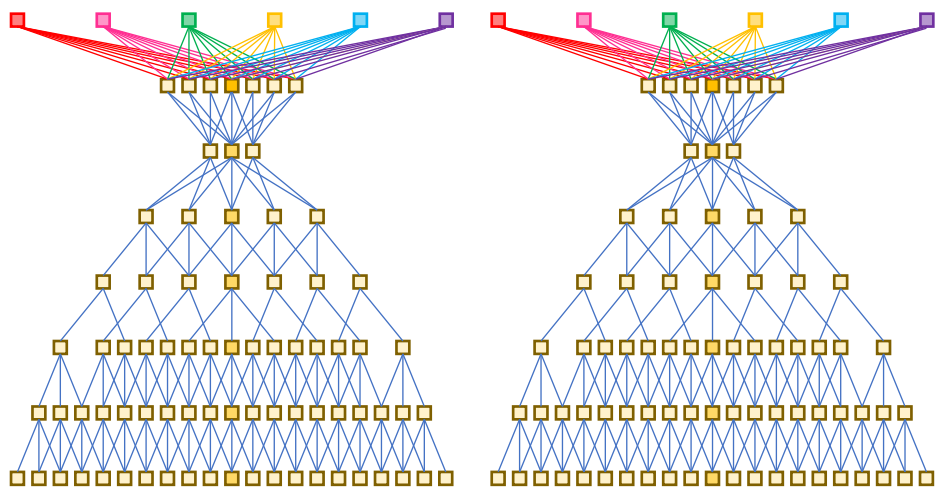
Phoneme HMMs Phoneme Posterior Probabilities Acoustic Likelihoods $p(\mathbf{h}|y_i)$ $p(y_i|\phi) = \frac{p(\phi|y_i)p(y_i)}{p(\phi)} \propto \frac{p(\phi|y_i)}{p(\phi)}$



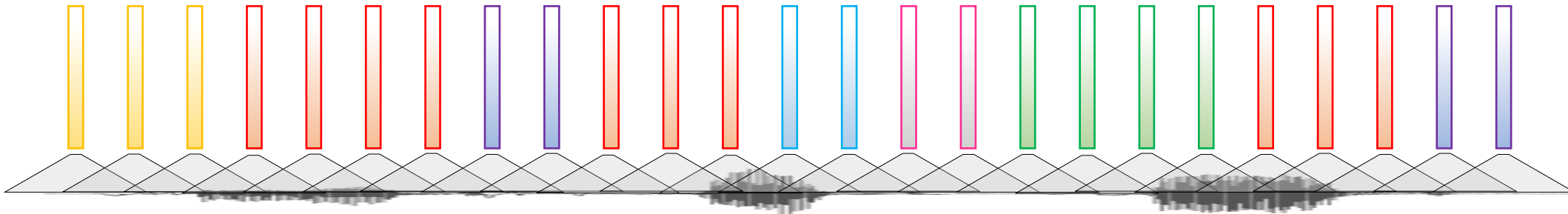
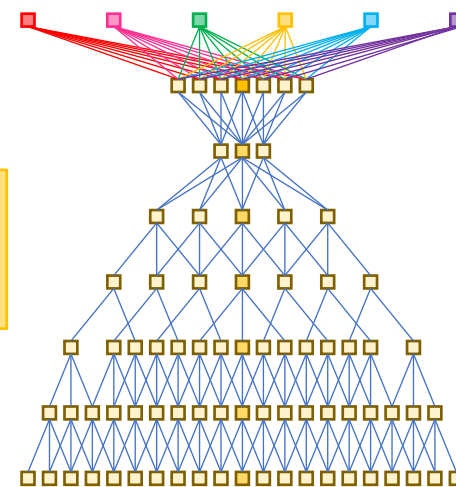
$$\mathcal{L}(\theta) = - \sum_{i=1}^n \log p_{\theta}(\hat{\phi}_i | y_i)$$



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



$$\mathcal{L}(\theta) = -\log \sum_t \prod_{i=1}^n p_{\theta}(\phi_{t_i} | y_i)$$



Language Modeling with (Recurrent) Neural Networks

Efforts to Get Further Away from GMM-HMM Architectures

Using Neural Networks to Estimate $P(w_t|h_t)$

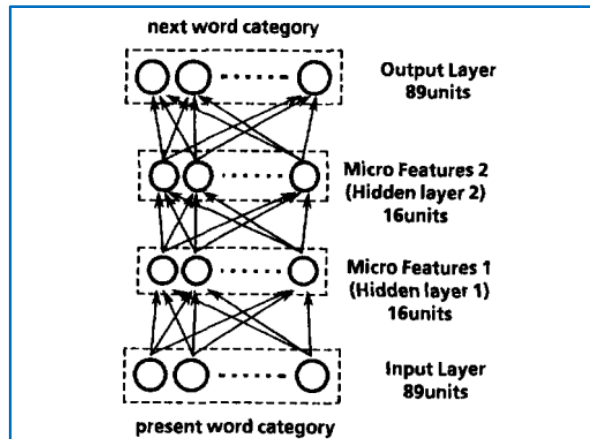
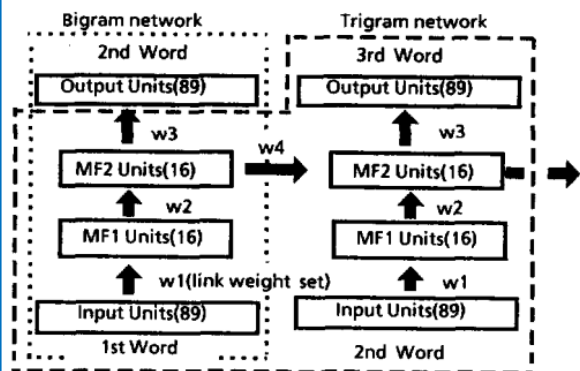


Fig.2 Basic Bigram Network for Word Category Prediction



→ every unit has a one-way connection to every unit of the next layer in this direction

Fig.3 NETgram Model 1 for Word Category Prediction

A STUDY OF ENGLISH WORD CATEGORY PREDICTION BASED ON NEURAL NETWORKS

Masami NAKAMURA, Kiyohiro SHIKANO

ATR Interpreting Telephony Research Laboratories
Seika-chou, Souraku-gun, Kyoto 619-02, JAPAN

ICASSP 1989

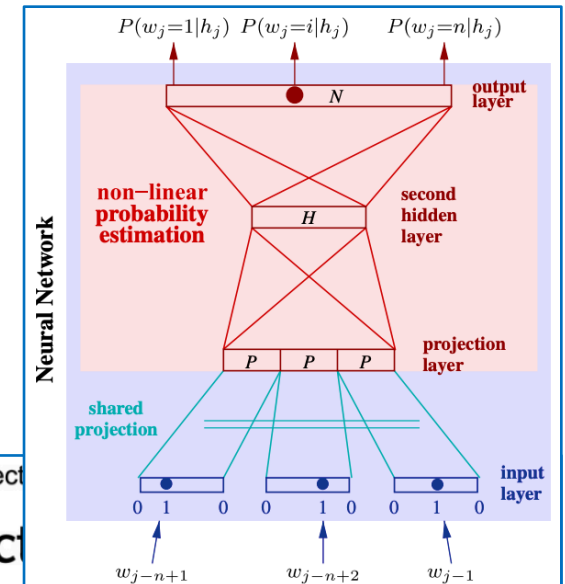


ELSEVIER

Available online at www.sciencedirect.com



Computer Speech and Language 21 (2007) 492–518



LANGUAGE

www.elsevier.com/locate/csl

Continuous space language models ☆

Holger Schwenk

Spoken Language Processing Group, LIMSI-CNRS, BP 133, 91403 Orsay cedex, France

Received 19 December 2005; received in revised form 15 September 2006; accepted 15 September 2006

Available online 9 October 2006

A paper appeared in September 2010 ...

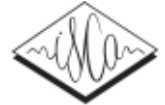
COGNITIVE SCIENCE **14**, 179–211 (1990)

Finding Structure in Time

JEFFREY L. ELMAN

University of California, San Diego

INTERSPEECH 2010



Recurrent neural network based language model

Tomáš Mikolov^{1,2}, Martin Karafiát¹, Lukáš Burget¹, Jan “Honza” Černocký¹, Sanjeev Khudanpur²

¹Speech@FIT, Brno University of Technology, Czech Republic

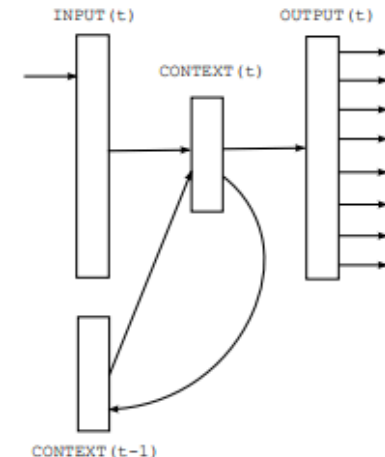
²Department of Electrical and Computer Engineering, Johns Hopkins University, USA

{imikolov,karafiat,burget,cernocky}@fit.vutbr.cz, khudanpur@jhu.edu

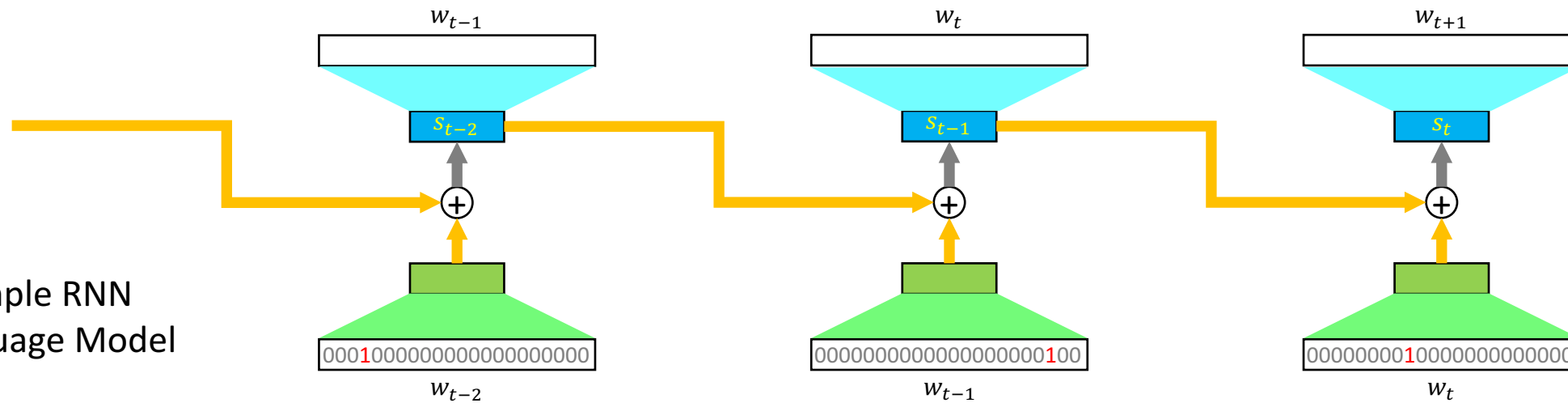
Abstract

A new recurrent neural network based language model (RNN LM) with applications to speech recognition is presented. Results indicate that it is possible to obtain around 50% reduction of perplexity by using mixture of several RNN LMs, compared to a state of the art backoff language model. Speech recognition experiments show around 18% reduction of word error rate on the Wall Street Journal task when comparing models trained on the same amount of data, and around 5% on the much harder NIST RT05 task, even when the backoff model is trained on much more data than the RNN LM. We provide ample empirical evidence to suggest that connectionist language models are superior to standard n-gram techniques, except their high computational (training) complexity.

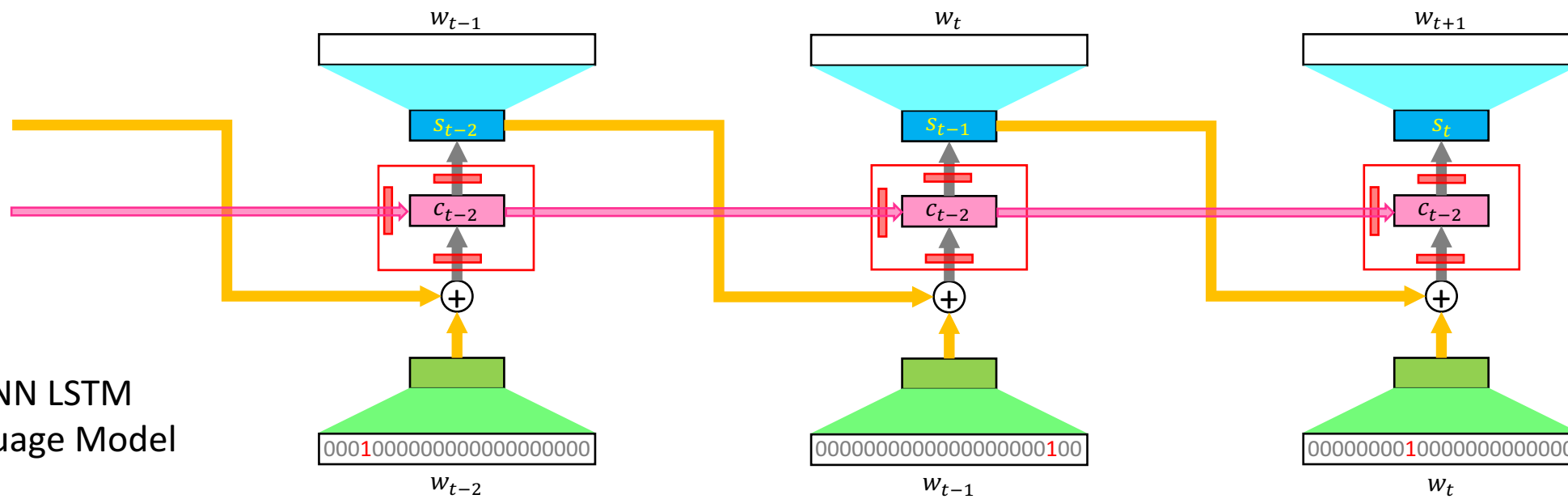
Index Terms: language modeling, recurrent neural networks, speech recognition



A Simple RNN
Language Model

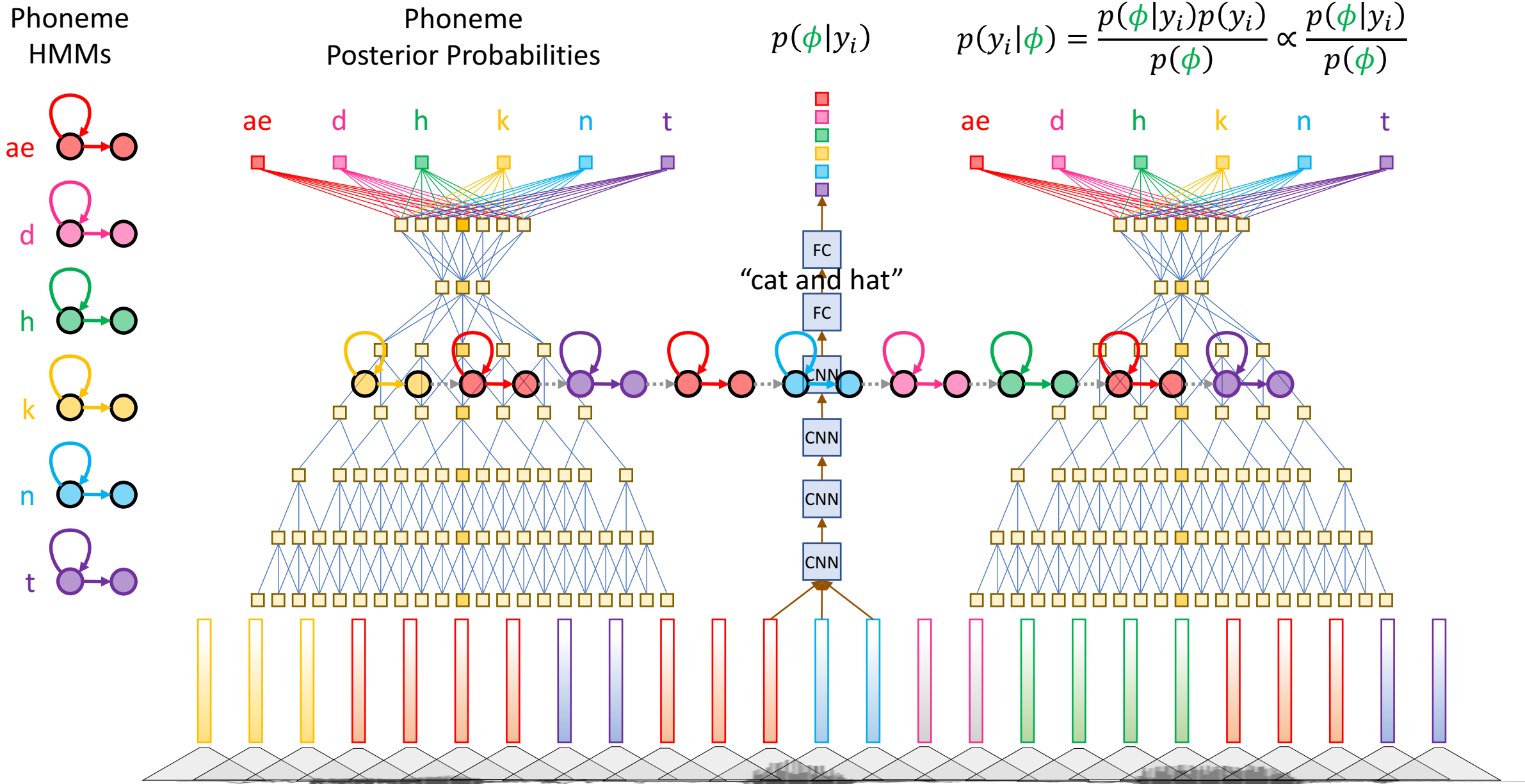


An RNN LSTM
Language Model

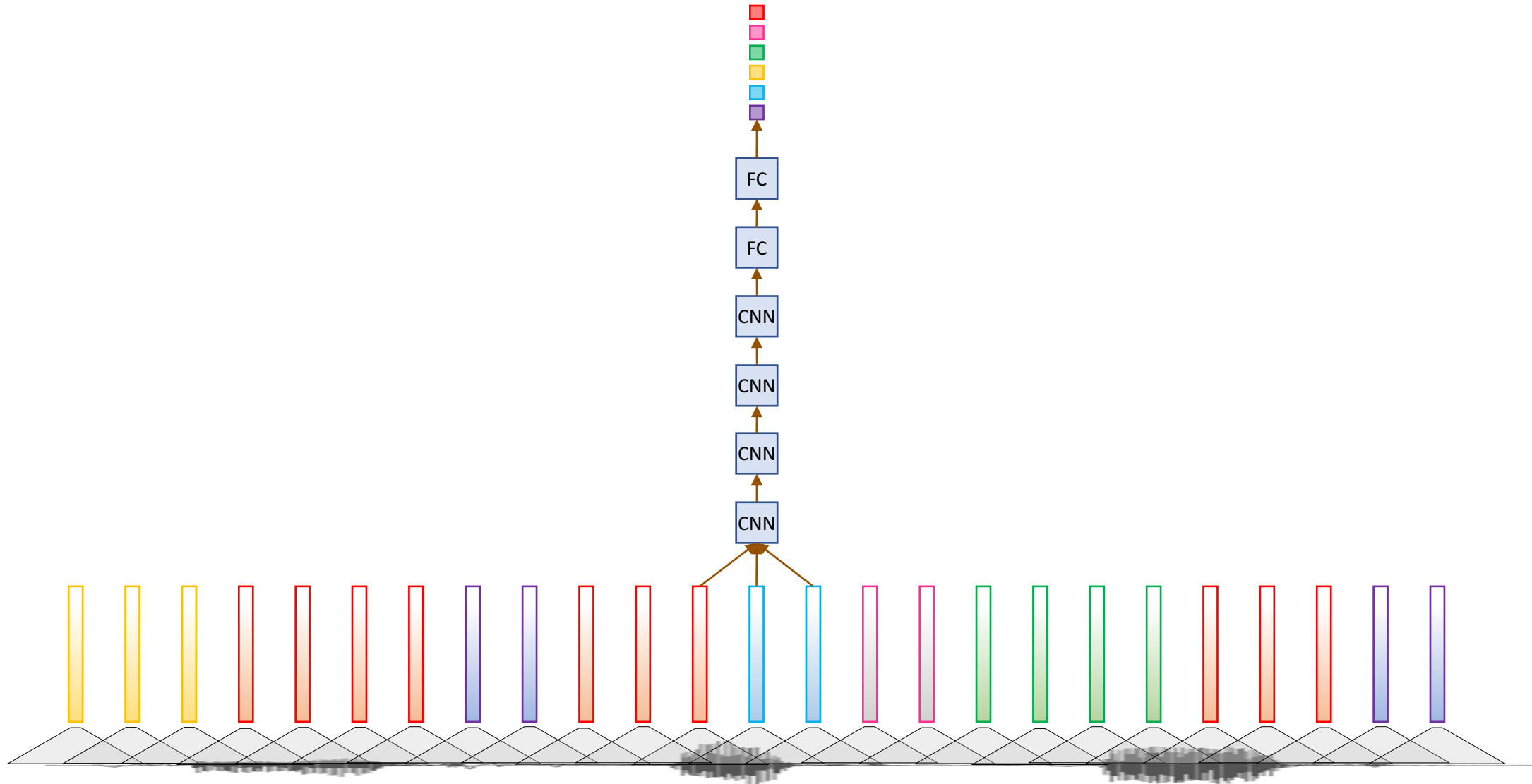


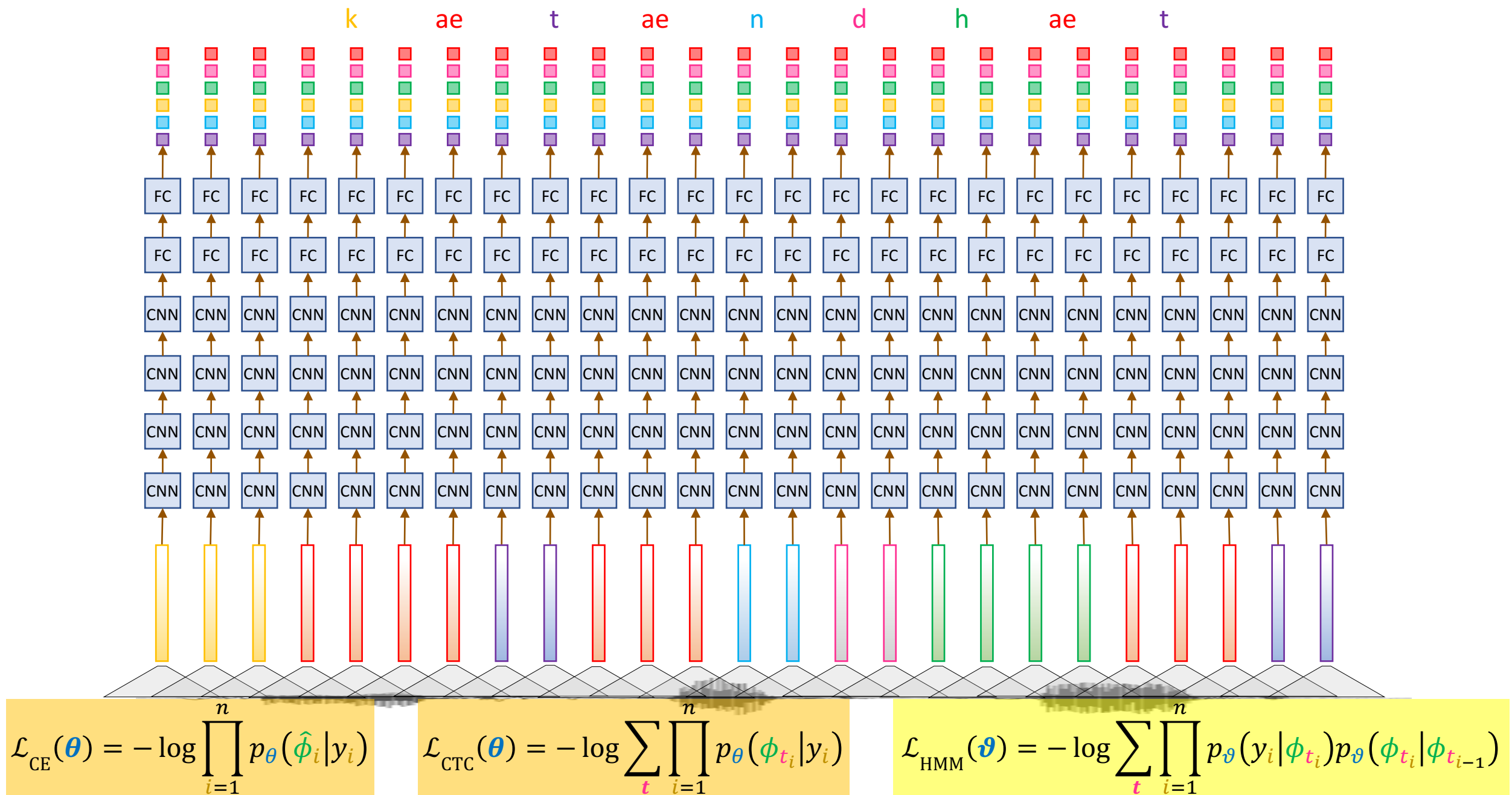
Speech Recognition without the HMM “Backend”

Efforts to Get Away from GMM-HMM Architectures



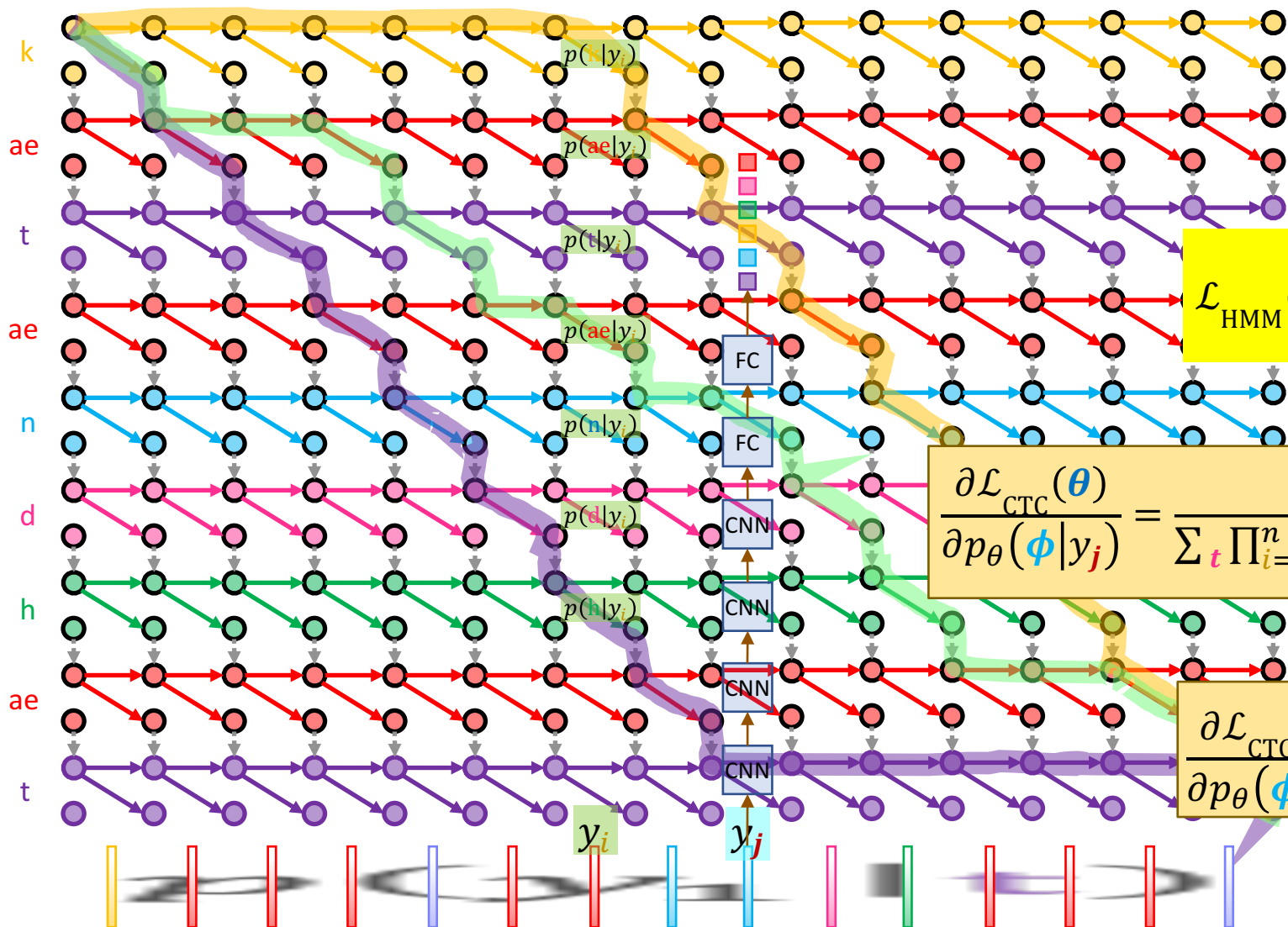
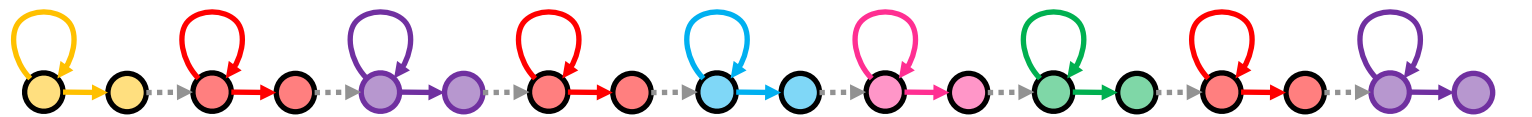
$$\mathcal{L}_{\text{CE}}(\theta) = -\log \prod_{i=1}^n p_{\theta}(\hat{\phi}_i | y_i) = -\sum_{i=1}^n \log p_{\theta}(\hat{\phi}_i | y_i)$$





Calculating the CTC loss for “cat and hat”

Calculating the gradient of the CTC loss

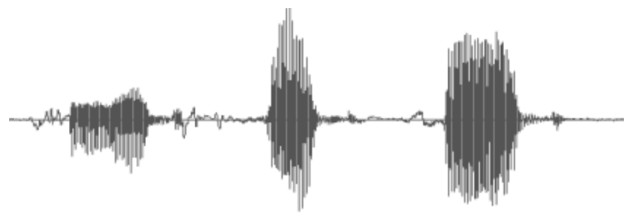


$$\mathcal{L}_{\text{CTC}}(\theta) = -\log \sum_t \prod_{i=1}^n p_{\theta}(\phi_{t_i} | y_i)$$

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

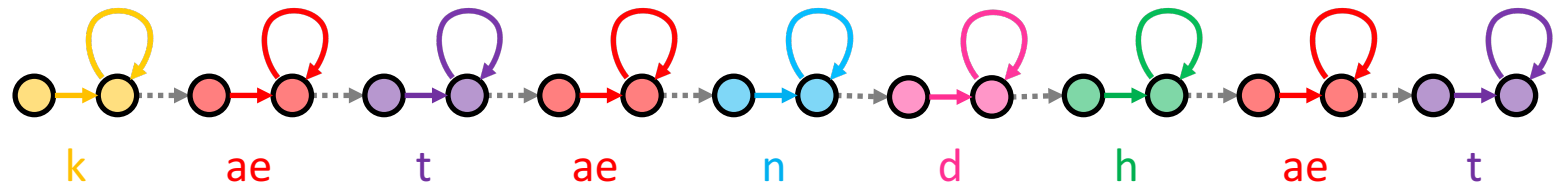
$$\frac{\partial \mathcal{L}_{\text{CTC}}(\theta)}{\partial p_{\theta}(\phi | y_j)} = \frac{-1}{\sum_t \prod_{i=1}^n p_{\theta}(\phi_{t_i} | y_i)} \sum_{t: \phi_{t_j} = \phi} \frac{1}{p_{\theta}(\phi | y_j)} \prod_{i=1}^n p_{\theta}(\phi_{t_i} | y_i)$$

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



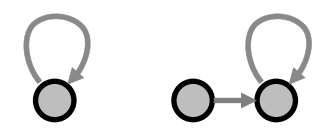
cat and hat

Composite HMM for "cat and hat"

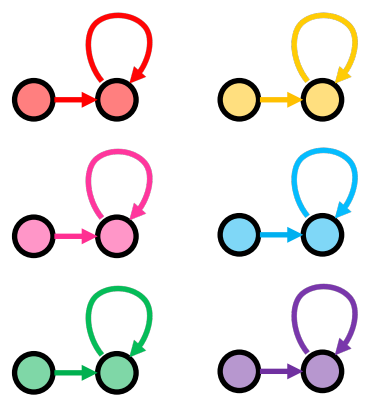


and	ae	n	d
cat	k	ae	t
hat	h	ae	t

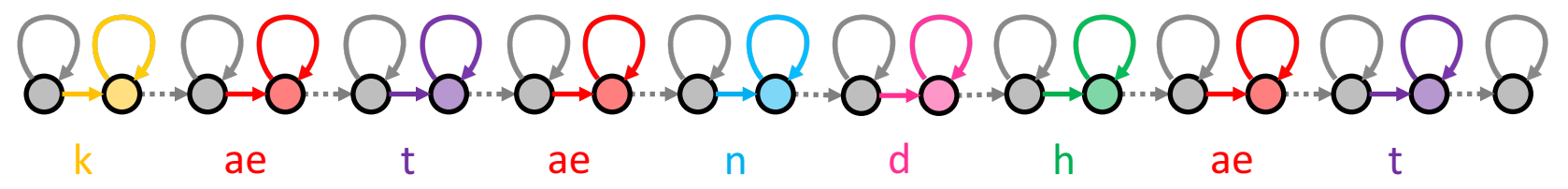
The CTC "Blank" Symbol (β)

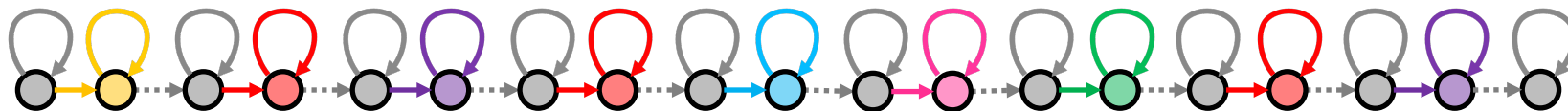


Phoneme HMMs



FSA of permissible CTC strings for "cat and hat"



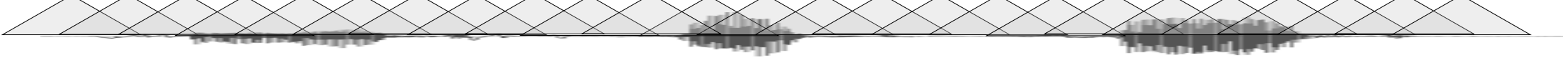
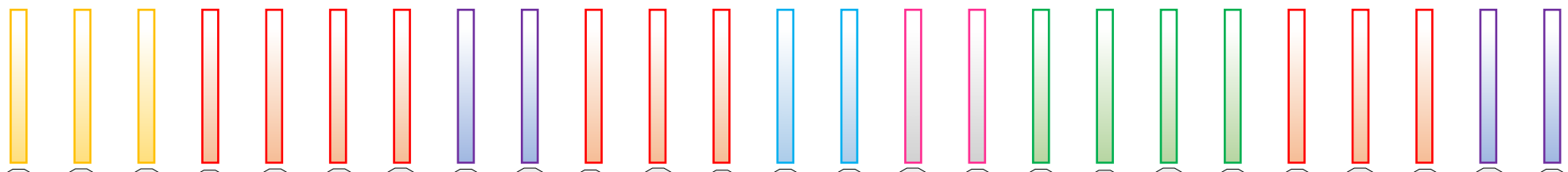
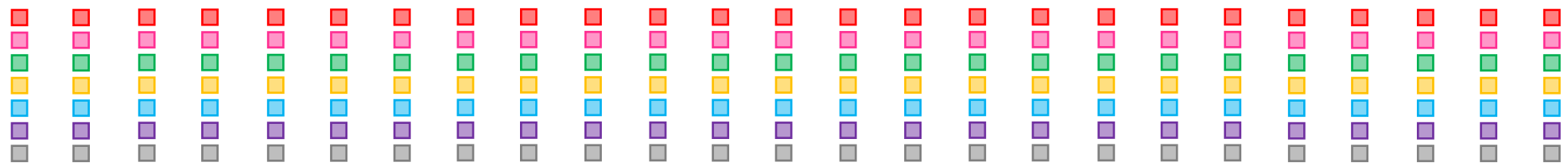


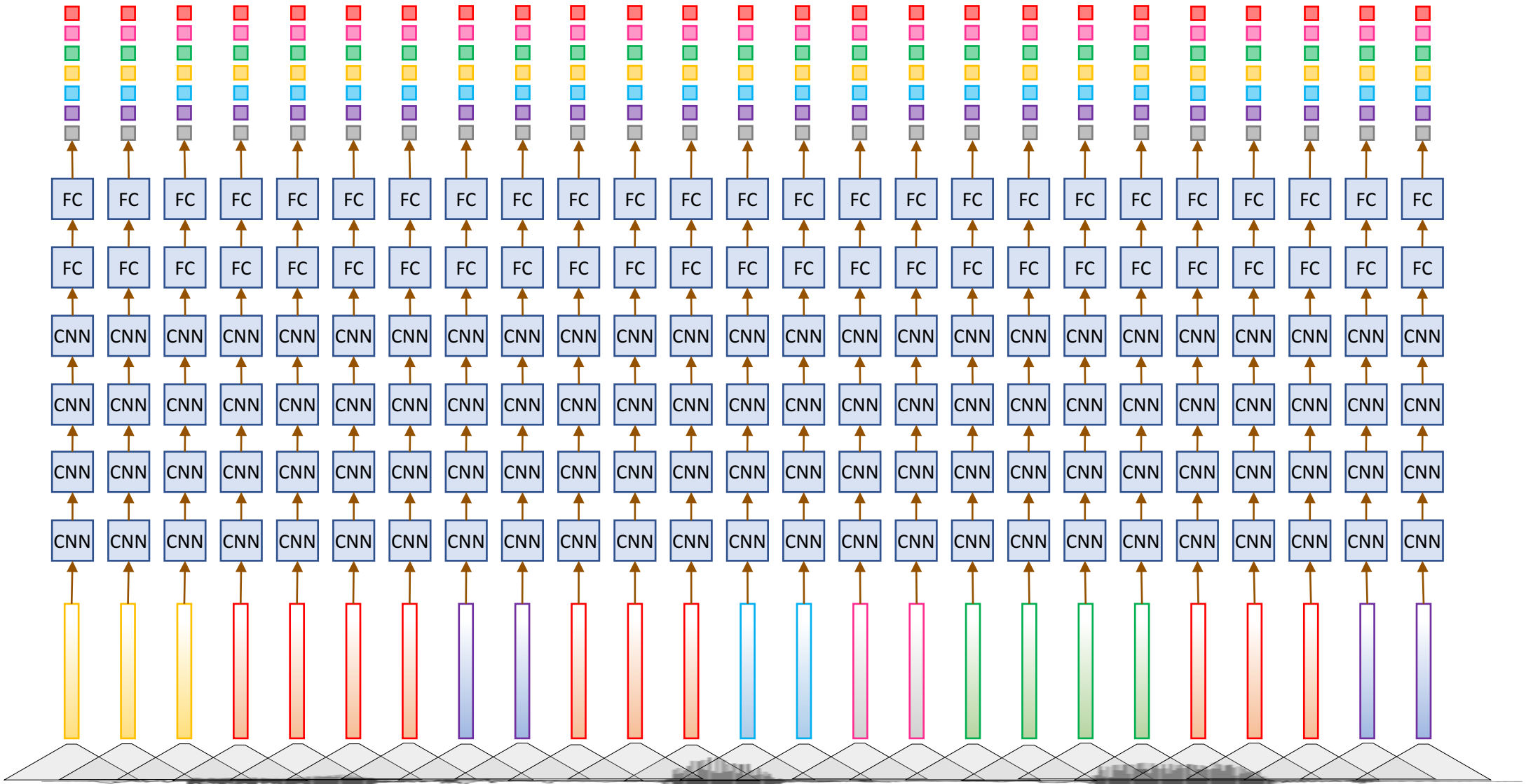
HMM State Sequences

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

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



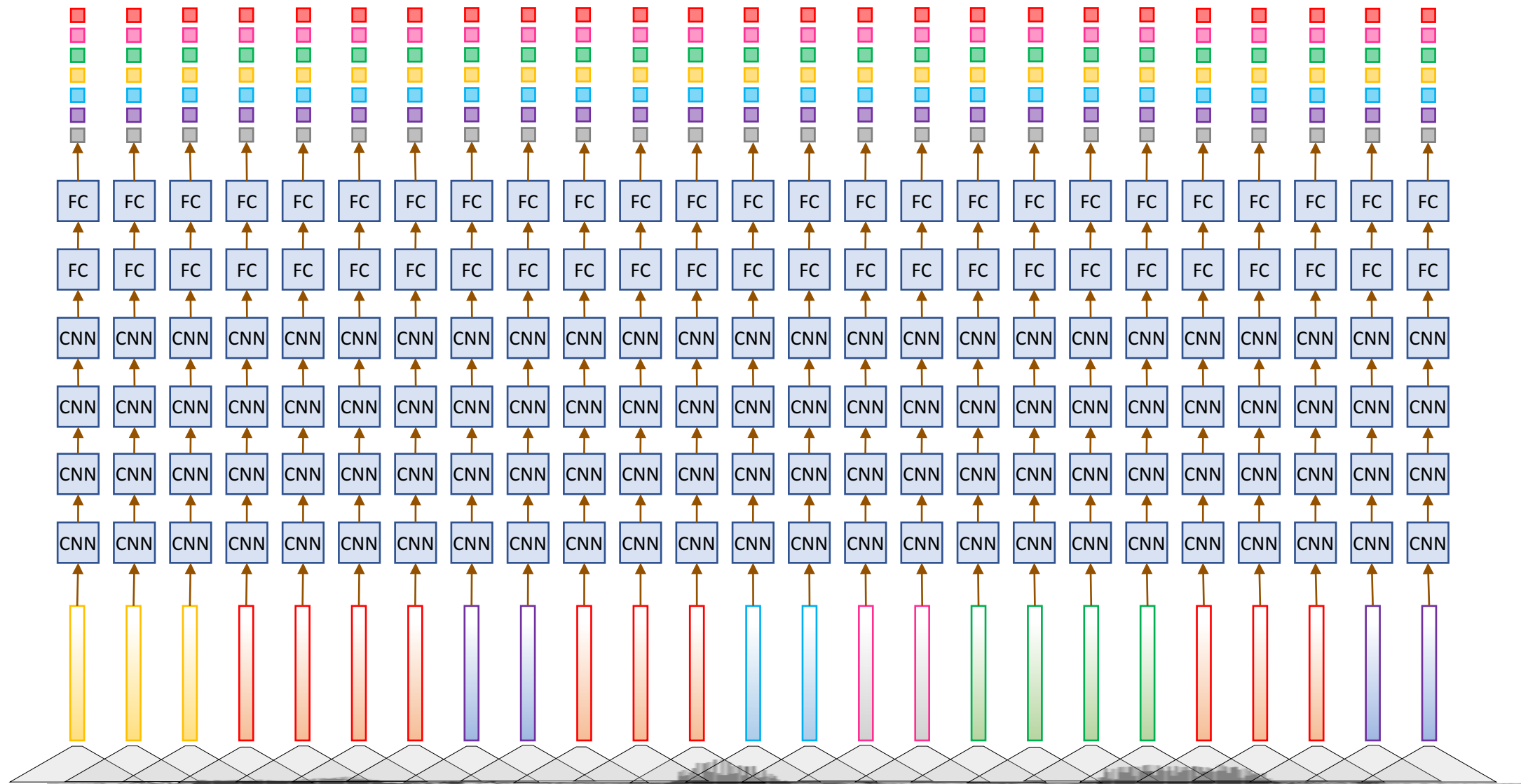


$$\mathcal{L}_{\text{CTC}}(\theta) = -\log \sum_t \prod_{i=1}^n p_{\theta}(\phi_{t_i} | y_i)$$

End-to-End Speech Recognition using Neural Networks with Attention

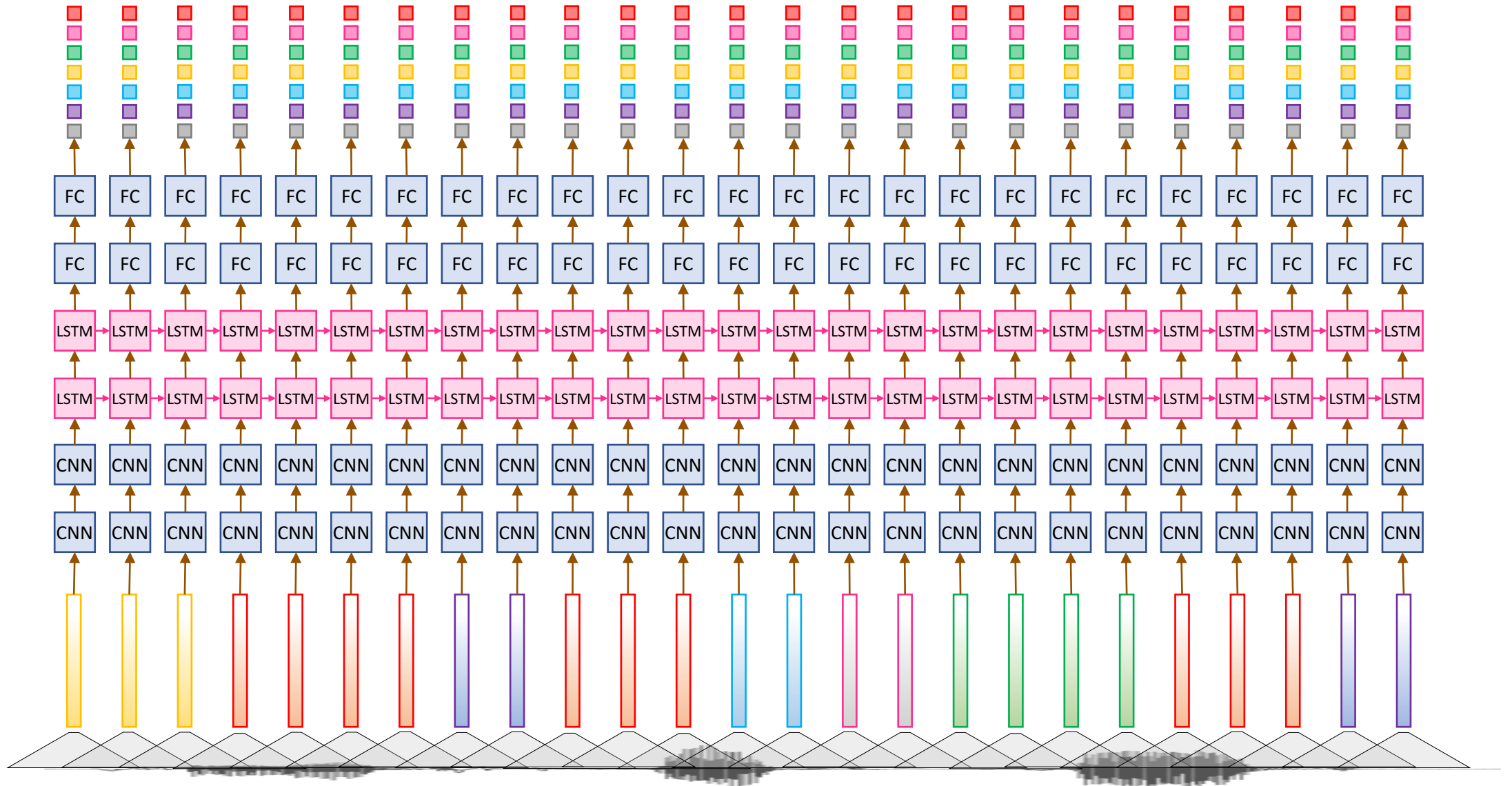
Efforts to Get Further Away from GMM-HMM Architectures

A CNN Architecture and CTC Loss

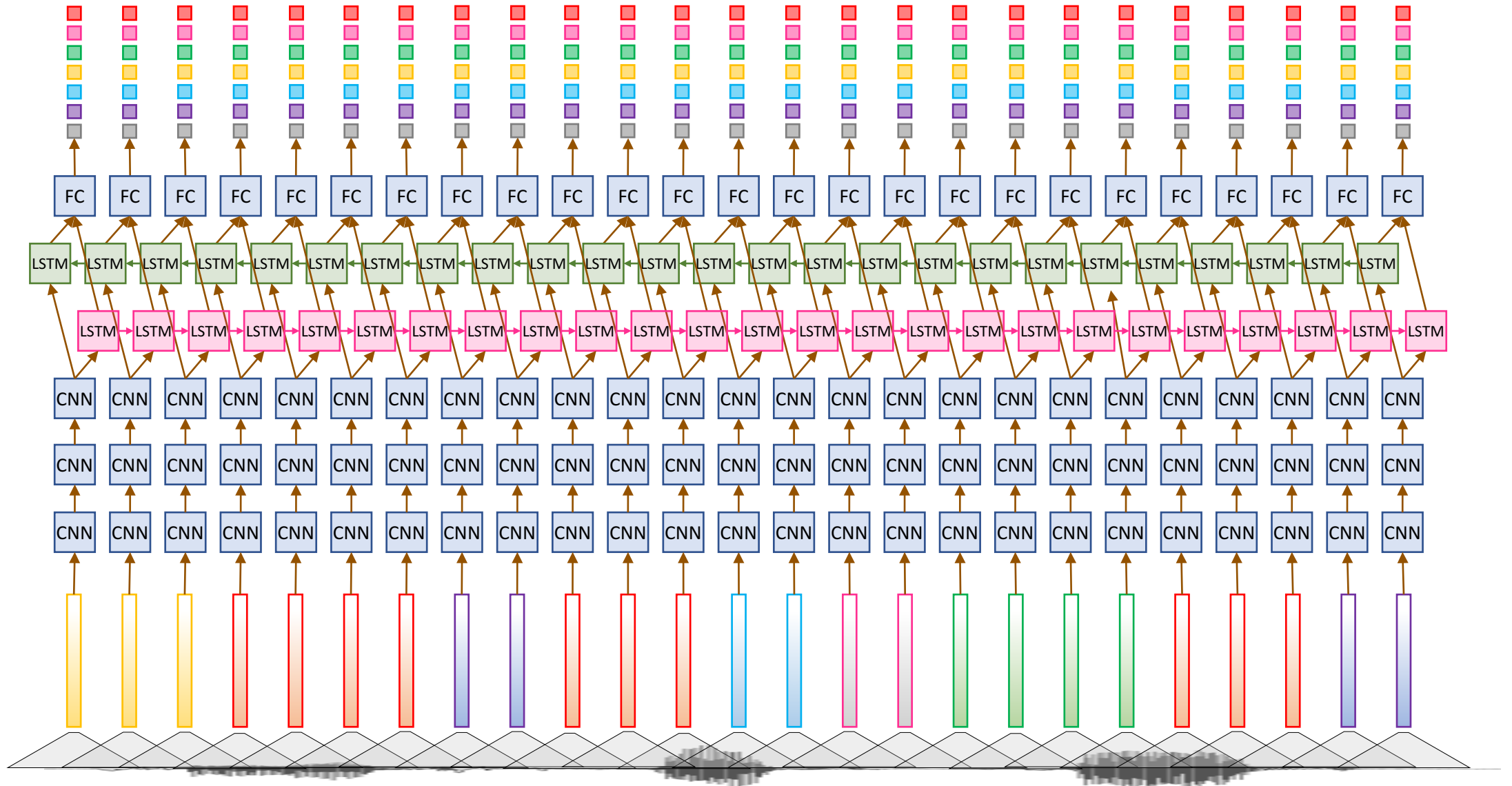


$$\mathcal{L}_{\text{CTC}}(\theta) = -\log \sum_t \prod_{i=1}^n p_{\theta}(\phi_{t_i} | y_i)$$

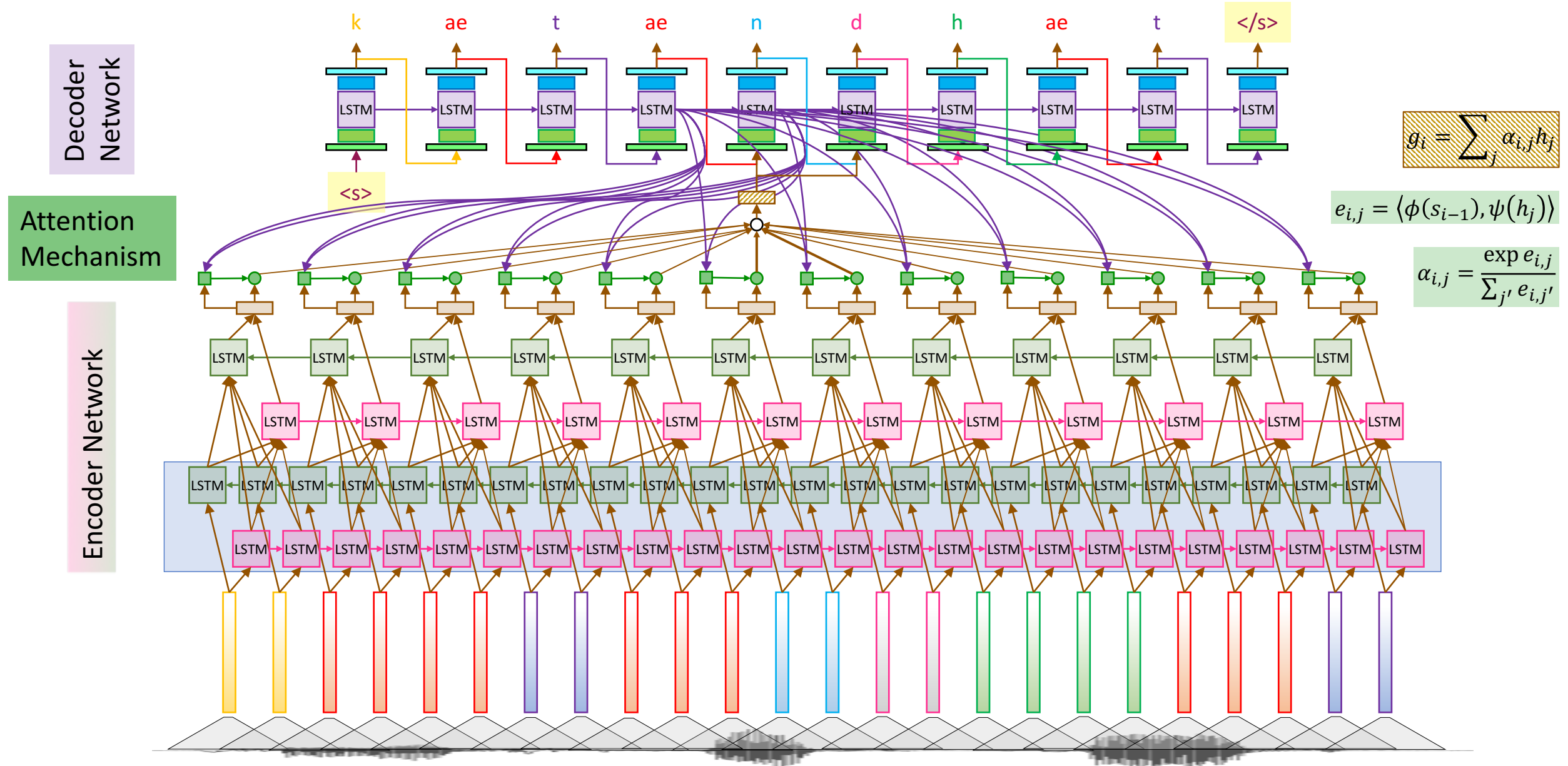
A CNN+LSTM Architecture



A Bidirectional LSTM Architecture (Deep Speech)



An Encoder-Decoder Architecture with Attention



Summary + Q&A