

COM4511 Speech Technology: Advanced Acoustic Models

Anton Ragni

February 24, 2020



- ▶ Given a parameterised audio sequence, infer the underlying latent representation
 - ▶ **parameterised audio**: sequences of feature vectors $\mathbf{O}_{1:T} = \mathbf{o}_1, \dots, \mathbf{o}_T$
 - ▶ **latent representation**: sequences of words $\mathbf{w}_{1:L} = w_1, \dots, w_L$

- ▶ Options for inference
 - ▶ maximum-a-posteriori

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \left\{ P(\mathbf{w} | \mathbf{O}_{1:T}) \right\}$$

- ▶ yields most probable sequence of words (sentence-level)
- ▶ minimum Bayes' risk

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \left\{ \sum_{\mathbf{w}'} P(\mathbf{w}' | \mathbf{O}_{1:T}) \mathcal{L}(\mathbf{w}, \mathbf{w}') \right\}$$

- ▶ yields sequence of words with the smallest expected loss (word or character level)
- ▶ Need to know:
 - ▶ how to model the posterior and perform inference (search)

Speech recognition centric view – later lectures will cover speech synthesis

- ▶ **Generative approach** models posterior **indirectly**

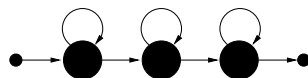
$$P(\mathbf{w}_{1:L} | \mathbf{O}_{1:T}) = \frac{1}{p(\mathbf{O}_{1:T})} p(\mathbf{O}_{1:T} | \mathbf{w}_{1:L}) P(\mathbf{w}_{1:L})$$

- ▶ standalone acoustic model $p(\mathbf{O}_{1:T} | \mathbf{w}_{1:L})$
- ▶ **Discriminative approach** models posterior **directly**
 - ▶ acoustic model is integrated into posterior — integrated or end-to-end approaches
 - ▶ **BUT** so is the language model
- ▶ Both approaches model the posterior probability
 - ▶ what is the difference? which one is better?

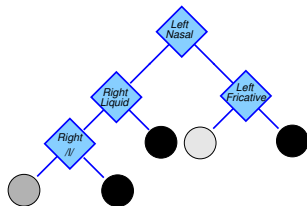
State of The Art in Generative Modelling — All Variants of HMM



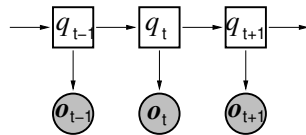
~ 1980	GMM-HMM	L. R. Rabiner, "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition", Proc IEEE, 1989
~ 2000	Tandem (FFNN)	H. Hermansky <i>et al</i> , "Tandem Connectionist Feature Extraction for Conventional HMM systems", Proc ICASSP, 2000
~ 2010	Hybrid (FFNN)	G. Hinton <i>et al</i> , "Deep Neural Networks for Acoustic Modeling in Speech Recognition", IEEE Sig Proc Mag, 2012
~ 2013	Hybrid (RNN)	H. Sak <i>et al</i> , "Long Short-Term Memory Recurrent Neural Network Architectures for Large Scale Acoustic Modeling", Proc Interspeech, 2014
~ 2015	Hybrid (CNN and RNN)	T. Sainath <i>et al</i> , "Convolutional, Long Short-Term Memory, fully connected Deep Neural Networks", Proc ICASSP, 2015 V. Peddinti <i>et al</i> , "Low latency acoustic modeling using temporal convolution and LSTMs", IEEE Sig Proc Let, 2018
~ 2018	Hybrid (CNN)	D. Povey, "Semi-Orthogonal Low-Rank Matrix Factorization for Deep Neural Networks", Proc Interspeech, 2018



(a) Standard ASR topology



(b) Phonetic decision trees



(c) Dynamic Bayesian Network

- ▶ Observations conditionally independent of other observations given state
- ▶ States conditionally independent of other states given past state

$$p(\mathbf{O}_{1:T} | \mathbf{w}_{1:L}) = \sum_{\mathbf{q}_{1:T} \in \mathbf{Q}_{1:T}^{(\mathbf{w}_{1:L})}} \prod_{t=1}^T p(\mathbf{o}_t | q_t) P(q_t | q_{t-1})$$

- ▶ transition probabilities, $P(q_t | q_{t-1})$, probability mass functions (discrete distributions)
- ▶ output distributions, $p(\mathbf{o}_t | q_t)$, probability density functions (continuous distributions)

- ▶ Form of probability density function

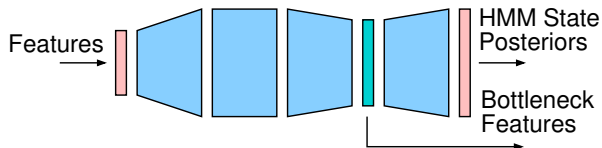
$$p(\mathbf{o}|q) = \sum_{m=1}^M c_m \mathcal{N}(\mathbf{o}; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$$

- ▶ mixture "weights" $\mathbf{c}_{1:M}$ such that $c_m \geq 0$ for all m and $\sum_{m=1}^M c_m = 1$
- ▶ how many parameters are in the model?
- ▶ Simple form of Gaussian distribution enables easy adjustment
 - ▶ speakers, noise and environment, e.g. maximum likelihood linear regression

$$\boldsymbol{\mu}_m^{(s)} = \mathbf{A}^{(s)} \boldsymbol{\mu}_m + \mathbf{b}^{(s)}, \quad \boldsymbol{\Sigma}^{(s)} = \mathbf{H}^{(s)\top} \boldsymbol{\Sigma}_m \mathbf{H}^{(s)}$$

- ▶ decision boundary, e.g., maximum mutual information

$$\mathcal{L}(\mathcal{D}; \boldsymbol{\theta}) = -\frac{1}{|\mathcal{D}|} \sum_{(\mathbf{w}, \mathbf{o}) \in \mathcal{D}} \log(P(\mathbf{w}|\mathbf{o}))$$



- Feed-forward neural network

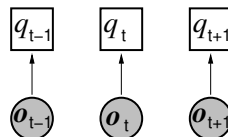
$$\mathbf{h}^{(l)} = \phi^{(l)}(\mathbf{A}^{(l)}\mathbf{h}^{(l-1)} + \mathbf{b}^{(l)})$$

- initialise $\mathbf{h}^{(0)} = \mathbf{o}$, terminate $P(q = i | \mathbf{o}) = h_i^{(L)}$, $\phi^{(L)}$ softmax, $\phi^{(l < L)}$ sigmoid/ReLU
- Augment hand-crafted features with features **learnt** by the neural network
 - options available what data to use (matched, mismatched, multi-lingual)
- **BUT** learnt features may be correlated
 - if Gaussians with diagonal covariance matrices used need to decorrelate

$$\hat{\mathbf{o}} = \begin{bmatrix} \mathbf{I} & \mathbf{I} \\ \mathbf{I} & \mathbf{S} \end{bmatrix} \begin{bmatrix} \mathbf{o} \\ \mathbf{h}^{(l)} \end{bmatrix}$$

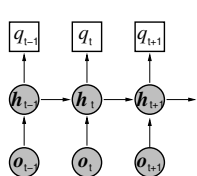


(a) Standard ASR architecture

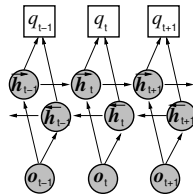


(b) DNN DBN

- ▶ Use feed-forward neural network to predict HMM state posteriors
 - ▶ efficient use of parameters (c.f. decision trees)
 - ▶ enables complex forms of classifiers (c.f. lecture on neural networks)
 - ▶ less sensitive to correlated features (c.f. MFCC)
- ▶ Possible to integrate into generative and discriminative sequence models
 - ▶ **generative**: HMM (discussed later)
 - ▶ **discriminative**: CTC (discussed later)
- ▶ Other types of neural networks are currently more popular
 - ▶ recurrent neural network, convolutional neural network and their combinations



(a) (Unidirectional) RNN



(b) Bidirectional RNN

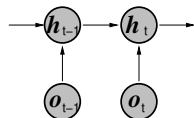
- Introduce dependence on all past observations using learnt history representation

$$P(\mathbf{q}_{1:T} | \mathbf{O}_{1:T}) = \prod_{t=1}^T P(q_t | \mathbf{q}_{1:t-1}, \mathbf{O}_{1:T}) \approx \prod_{t=1}^T P(q_t | \mathbf{O}_{1:t}) \approx \prod_{t=1}^T P(q_t | \mathbf{h}_t) \quad (1)$$

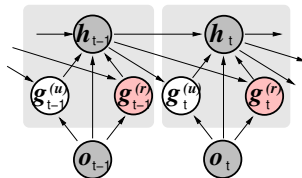
- assumes future observations are of no use
- Use another recurrent cell to condition on future observations (bi-directional RNN)

$$P(\mathbf{q}_{1:T} | \mathbf{O}_{1:T}) \approx \prod_{t=1}^T P(q_t | \vec{\mathbf{h}}_t, \overleftarrow{\mathbf{h}}_t)$$

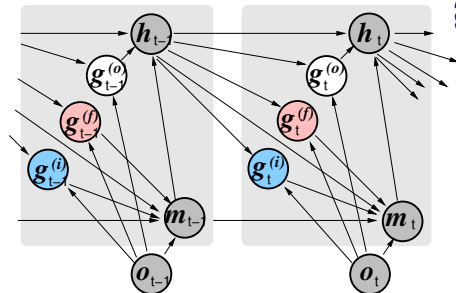
- BUT states remain conditionally independent



(a) Simple Recurrent Unit



(b) Gated Recurrent Unit



(c) Long Short-Term Memory

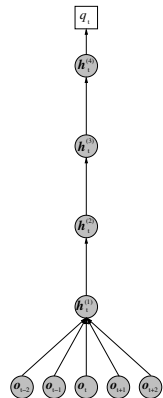
- Simple recurrent unit is known to have unstable behaviour
 - use "gates" to pass or block

$$\mathbf{g}_t = \sigma(\mathbf{A}^{(g)} \mathbf{o}_t + \mathbf{C}^{(g)} \mathbf{h}_{t-1} + \mathbf{b}^{(g)}), \quad \text{where } \sigma \text{ is a sigmoid}$$

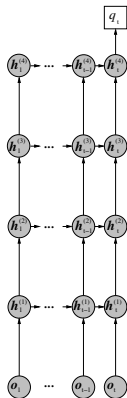
- Multiple "smart" recurrent units can be devised, e.g. GRU

$$\mathbf{h}_t = \mathbf{g}_t^{(u)} \odot \mathbf{h}_{t-1} + (\mathbf{1} - \mathbf{g}_t^{(u)}) \odot \phi(\mathbf{A} \mathbf{o}_t + \mathbf{C}(\mathbf{g}_t^{(r)} \odot \mathbf{h}_{t-1}) + \mathbf{b})$$

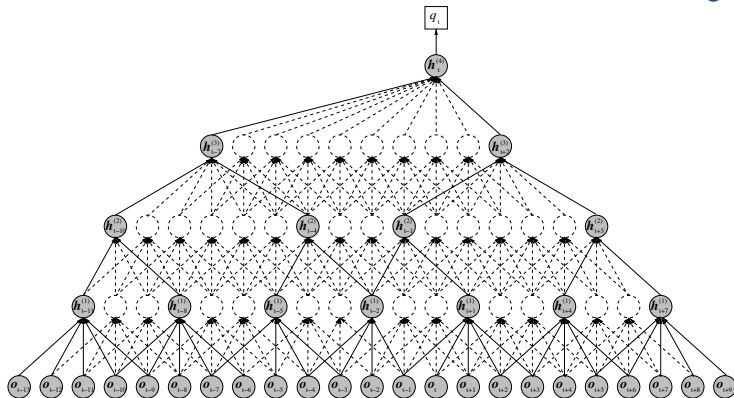
Time-Delay Neural Networks Models



(a) Feed-Forward

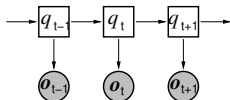


(b) Recurrent

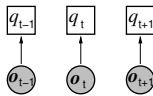


(c) Time Delay

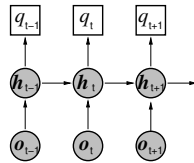
- ▶ Long-term dependency modelling using standard networks is challenging
 - ▶ feed-forward models are highly inefficient
 - ▶ recurrent models are hard to train (full history, instability)
- ▶ Use a **sub-sampled pyramidal** structure to increase scope of possible dependencies
 - ▶ sharing parameters in each layer is akin to a filter (**convolutional network!**)



(a) HMM DBN



(b) DNN DBN



(c) RNN Pseudo-DBN

- **Feed-forward** types inherit standard conditional independence assumptions

$$p(\mathbf{O}_{1:T} | \mathbf{q}_{1:T}) = \prod_{t=1}^T p(\mathbf{o}_t | \mathbf{O}_{1:t-1}, \mathbf{q}_{1:T}) \approx \prod_{t=1}^T p(\mathbf{o}_t | q_t) = \prod_{t=1}^T \overbrace{\frac{P(q_t | \mathbf{o}_t) p(\mathbf{o}_t)}{P(q_t)}}^{\text{FFNN}}$$

- **Recurrent** types relax conditional independence assumptions of observations

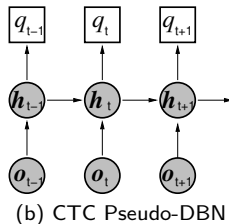
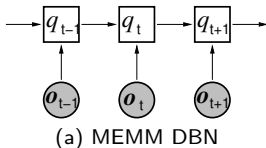
$$p(\mathbf{O}_{1:T} | \mathbf{q}_{1:T}) = \prod_{t=1}^T p(\mathbf{o}_t | \mathbf{O}_{1:t-1}, \mathbf{q}_{1:T}) \approx \prod_{t=1}^T p(\mathbf{o}_t | \mathbf{O}_{1:t-1}, q_t) = \prod_{t=1}^T \overbrace{\frac{P(q_t | \mathbf{O}_{1:t}) p(\mathbf{o}_t | \mathbf{O}_{1:t-1})}{P(q_t | \mathbf{O}_{1:t-1})}}^{\text{RNN}}$$

- **BUT** not the assumptions about states

State of The Art in Discriminative Modelling



~ 2010	SCRF/CAug	G. Zweig <i>et al</i> , "A Segmental CRF Approach to Large Vocabulary Continuous Speech Recognition", ASRU, 2009 A. Ragni <i>et al</i> , "Derivative kernels for noise-robust ASR", ASRU, 2011
~ 2012	CTC	A. Graves <i>et al</i> , "Towards End-to-End Speech Recognition with Recurrent Neural Networks", ICML, 2014
~ 2014	Encoder-Decoder	I. Sutskever, "Sequence to Sequence Learning with Neural Networks", NIPS, 2014
~ 2016	Attention	W. Chan, "Listen, Attend and Spell", ICASSP, 2016
~ 2018	RNN-Transducer	Y. He <i>et al</i> , "Streaming end-to-end speech recognition for mobile devices", ICASSP, 2019
~ 2018	Transformer	A. Zayer <i>et al</i> , "A Comparison of Transformer and LSTM Encoder Decoder Models for ASR", ASRU, 2019



- ▶ Direct approach to using neural network for character sequence prediction

$$P(\mathbf{w}_{1:L} | \mathbf{O}_{1:T}) \approx \sum_{\mathbf{q}_{1:T} \in Q^{(\mathbf{w}_{1:L})}} P(\mathbf{q}_{1:T} | \mathbf{O}_{1:T}) \approx \sum_{\mathbf{q}_{1:T} \in Q^{(\mathbf{w}_{1:L})}} \prod_{t=1}^T P(q_t | \mathbf{h}_t)$$

- ▶ relaxes assumptions about observations but not about latent variables
- ▶ possible to use MEMM (HMM) style forward-backward algorithm
- ▶ Need to decide on the form of latent variable model (LVM)
 - ▶ only need to be able to generate latent variable sequences — options available?

Example: MEMM/HMM Latent Variable Model



- ▶ Latent variable model "rules"
 1. each symbol must appear at least once
 2. symbol order must be preserved
 3. must have mechanism to track symbol changes
- ▶ Given character sequence w_1, w_2, w_2, w_3 , possible latent variable sequences are

Length (T)	Latent Variable Sequences $\mathbf{q}_{1:T}$
4	$w_1, w_2, w_2, w_3; \dots$
5	$w_1, w_2, w_2, w_2, w_3; \dots$
6	$w_1, w_2, w_2, w_2, w_2, w_3; \dots$
7	$w_1, w_2, w_2, w_2, w_2, w_2, w_3; \dots$
8	$w_1, w_2, w_2, w_2, w_2, w_2, w_2, w_3; \dots$
9	$w_1, w_1, w_2, w_2, w_2, w_2, w_2, w_2, w_3; \dots$
\vdots	\dots

- ▶ explicit start/end states enable to know symbol boundaries
- ▶ write a regular expression to express all possible latent variable sequences

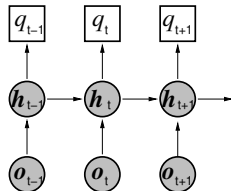
Example: CTC "Latent Variable Model"



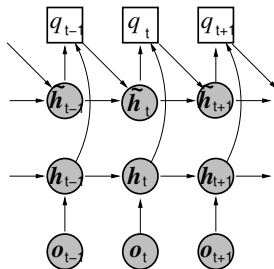
- ▶ Two observations can be made about latent variable handling in MEMM/HMM
 - ▶ symbol change is ambiguous only for two identical symbols
 - ▶ forcing each q_t to take one of $\mathbf{w}_{1:L}$ values may be problematic
- ▶ Introduce a new symbol ϵ to play the role of a delimiter and an uncertain symbol
 - ▶ encourage each q_t to take values in $\mathbf{w}_{1:L}$ only if highly certain
 - ▶ BUT has to artificially enforce ϵ between two identical symbols
- ▶ Given character sequence w_1, w_2, w_2, w_3 , possible latent variable sequences are

Length (T)	Latent Variable Sequences $\mathbf{q}_{1:T}$
4	—
5	$w_1, w_2, \epsilon, w_2, w_3; \dots$
6	$w_1, \epsilon, w_2, \epsilon, w_2, w_3; \dots$
7	$w_1, \epsilon, \epsilon, w_2, \epsilon, w_2, w_3; \dots$
8	$w_1, w_1, \epsilon, \epsilon, w_2, \epsilon, w_2, w_3; \dots$
9	$w_1, w_1, \epsilon, \epsilon, w_2, \epsilon, w_2, w_3, w_3; \dots$

- ▶ no need to have an explicit latent variable model
- ▶ write a regular expression to express all possible latent variable sequences



(a) CTC

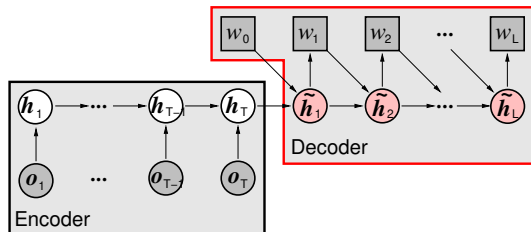


(b) RNN-T

- Break conditional independence assumption among latent variables

$$P(\mathbf{w}_{1:L} | \mathbf{O}_{1:T}) \approx \sum_{\mathbf{q}_{1:T} \in \mathbf{Q}(\mathbf{w}_{1:L})} \prod_{t=1}^T P(q_t | \mathbf{q}_{1:t-1}, \mathbf{O}_{1:t}) \approx \sum_{\mathbf{q}_{1:T} \in \mathbf{Q}(\mathbf{w}_{1:L})} \prod_{t=1}^T P(q_t | \tilde{\mathbf{h}}_t, \mathbf{h}_t)$$

- reported to yield competitive performance given large amounts of data



Predict next word using recursion

$$\tilde{h}_l = \phi(\tilde{h}_{l-1}, w_{l-1})$$

Options for setting \tilde{h}_0 :

- ▶ use output of another recursion

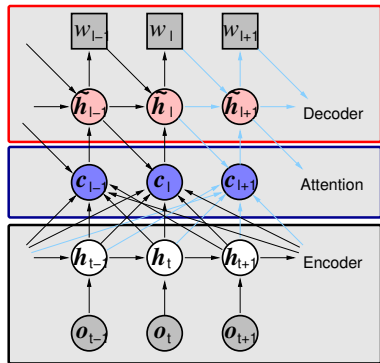
$$h_T = \phi(h_{T-1}, o_T)$$

- ▶ name other common options

- ▶ Use recurrent units to handle variable length observation and word sequences

$$P(\mathbf{w}_{1:L} | \mathbf{O}_{1:T}) = \prod_{l=1}^L P(w_l | \mathbf{w}_{1:l-1}, \mathbf{O}_{1:T}) \approx \prod_{l=1}^L P(w_l | \mathbf{w}_{1:l-1}, h_T) \approx \prod_{l=1}^L P(w_l | \tilde{h}_l)$$

- ▶ **encoder**: maps observation sequence $\mathbf{O}_{1:T}$ to fixed length representation h_T
- ▶ **decoder**: generates word sequence $\mathbf{w}_{1:L}$ given h_T
- ▶ **BUT** hard to ensure relevant information propagated/used for prediction



Position-dependent history

$$\mathbf{c}_l = \sum_{t=1}^T \alpha_{l,t} \mathbf{h}_t, \quad \text{where} \quad \alpha_{l,t} = \frac{\exp(z_{l,t})}{\sum_{t=1}^T \exp(z_{l,t})}$$

Unnormalised attention weights

$$z_{l,t} = \mathbf{d}^T \phi(\mathbf{A}\tilde{\mathbf{h}}_{l-1} + \mathbf{C}\mathbf{h}_t)$$

► "IR": query $\tilde{\mathbf{h}}_{l-1}$, key \mathbf{h}_t , values $\mathbf{h}_{1:T}$

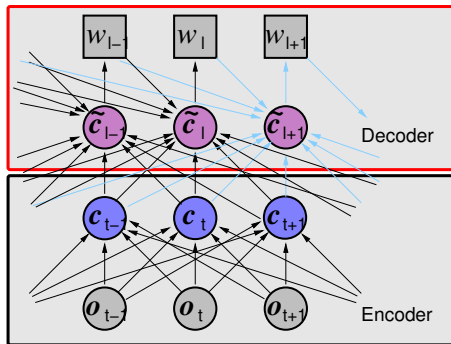
(other forms possible)

► Improve encoder-decoder by position-dependent input representation

$$P(\mathbf{w}_{1:L} | \mathbf{O}_{1:T}) = \prod_{l=1}^L P(w_l | \mathbf{w}_{1:l-1}, \mathbf{O}_{1:T}) \approx \prod_{l=1}^L P(w_l | \mathbf{w}_{1:l-1}, \mathbf{c}_l) \approx \prod_{l=1}^L P(w_l | \tilde{\mathbf{h}}_l)$$

► reported to yield competitive performance given large amounts of data

Transformer (Not a movie!)



- Replace all recurrent units with attention
 - multi-head attention: multiple parallel mechanisms to increase modelling power
 - positional encoding: add positional offset to encode temporal information

$$P(\mathbf{w}_{1:L} | \mathbf{O}_{1:T}) = \prod_{l=1}^L P(w_l | \mathbf{w}_{1:l-1}, \mathbf{O}_{1:T}) \approx \prod_{l=1}^L P(w_l | \mathbf{w}_{1:l-1}, \mathbf{c}_{1:T}) \approx \prod_{l=1}^L P(w_l | \tilde{\mathbf{c}}_l)$$

- Reported to yield competitive performance given large amounts of data

- ▶ Generative approaches provide a natural framework to integrate language models

$$P(\mathbf{w}_{1:L} | \mathbf{O}_{1:T}) \propto p(\mathbf{O}_{1:T} | \mathbf{w}_{1:L}) P(\mathbf{w}_{1:L})$$

- ▶ enables to plug-in any suitable language model
- ▶ Less obvious how to integrate language model into discriminative models
 - ▶ heuristic probability combination

$$s(\mathbf{w}_{1:L}, \mathbf{O}_{1:T}) = P(\mathbf{w}_{1:L} | \mathbf{O}_{1:T}) P(\mathbf{w}_{1:L})$$

- ▶ fusion between decoder and language model recurrent units

$$\tilde{\mathbf{h}}_l = \phi(\tilde{\mathbf{h}}_{l-1}, \tilde{\mathbf{h}}_{l-1}^{(LM)}, w_{l-1})$$

- ▶ multi-task training, pre-training and many more!

- ▶ The past two lectures have
 - ▶ introduced hidden Markov models (HMM)
 - ▶ discussed how HMMs can be used as an acoustic model
- ▶ This lecture has discussed more complex forms of acoustic models
 - ▶ generative models (Tandem, Hybrids)
 - ▶ discriminative models (CTC, encoder-decoder, transformer)
- ▶ Next lectures will look at language models