## Neural Networks for Acoustic Modelling 2: Hybrid HMM/DNN systems

Peter Bell

#### Automatic Speech Recognition – ASR Lecture 11 27 February 2023

ASR Lecture 11 Neural Networks for Acoustic Modelling 2: HMM/DNN 1

### Recap: Hidden units extracting features



$$h_k = \sigma \left( \sum_{d=1}^D v_{kd} x_d + b_k 
ight) \qquad y_j = \operatorname{softmax} \left( \sum_{k=1}^K w_{jk} h_k + b_j 
ight)$$

### Simple neural network for phone classification



## Neural networks for phone recognition

- So far we have trained networks to *classify* each frame of observations
- In phone *recognition*, we need to obtain the best phone (or word) sequence
- Hybrid NN/HMM systems: in an HMM, replace the GMMs used to estimate output pdfs with the outputs of neural networks
- Train a neural network to associate a phone-state label with a frame of acoustic data (+ context)
- Can interpret the output of the network as P(phone-state | acoustic-frame)
- Use NN to obtain output probabilities in Viterbi algorithm to find most probable sequence of phones

・ 回 ト ・ ヨ ト ・ ヨ ト

#### Posterior probability estimation

- Consider a neural network trained as a classifier each output corresponds to a class.
- When applying a trained network to test data, it can be shown that the value of output corresponding to class j given an input  $x_t$ , is an estimate of the posterior probability  $P(q_t = j | x_t)$ . (This is because we have softmax outputs and use a cross-entropy loss function)
- Using Bayes Rule we can relate the posterior  $P(q_t = j | x_t)$  to the likelihood  $p(x_t | q_t = j)$  used as an output probability in an HMM:

$$P(q_t|\mathsf{x}_t) = \frac{p(\mathsf{x}_t|q_t=j)P(q_t=j)}{p(\mathsf{x}_t)}$$

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・

## Scaled likelihoods

 If we would like to use NN outputs as output probabilities in an HMM, then we would like probabilities (or densities) of the form p(x|q) – likelihoods.

We can write *scaled likelihoods* as:

$$\frac{P(q_t = j | \mathsf{x}_t)}{p(q_t = j)} = \frac{p(\mathsf{x}_t | q_t = j)}{p(\mathsf{x}_t)}$$

- Scaled likelihoods can be obtained by "dividing by the priors" – divide each network output  $P(q_t = j | x_t)$  by  $P(q_t)$ , the relative frequency of class j in the training data
- Using  $p(x_t|q_t = j)/p(x_t)$  rather than  $p(x_t|q_t = j)$  is OK since  $p(x_t)$  does not depend on the class j
- Computing the scaled likelihoods can be interpreted as factoring out the prior estimates for each phone based on the acoustic training data. The HMM can then integrate better prior estimates based on the language model and lexicon

## Hybrid NN/HMM



- NNs can naturally model *acoustic* context, but how can we model *phonetic* context?
- Early solution (Bourlard et al, 1992) separate the modelling of the primary class, *y*, and its context, *c*, with two neural networks:

$$p(y,c|x) = p(c|y,x)p(y|x)$$

or

$$p(y,c|x) = p(y|c,x)p(c|x)$$

During decoding, we need separate forward passes for each context

周下 イヨト イヨト

# Using context as input for p(y|c,x)



< 17 b

< E



ASR Lecture 11 Neural Networks for Acoustic Modelling 2: HMM/DNN 10

<ロ> <同> <同> < 同> < 同>



ASR Lecture 11 Neural Networks for Acoustic Modelling 2: HMM/DNN 10

・ロト ・回ト ・ヨト ・ヨト



ASR Lecture 11 Neural Networks for Acoustic Modelling 2: HMM/DNN 10

(4回) (1日) (日)



ASR Lecture 11 Neural Networks for Acoustic Modelling 2: HMM/DNN 10

・ロト ・回ト ・ヨト ・ヨト

Tandem scheme:

- Basic idea: use the output probabilities from the NN as input features to standard CD-HMM-GMM system
- Combines the benefits of both:
  - NNs good at modelling wide acoustic contexts, correlated input features
  - HMM-GMMs good for speaker adaptation, modelling phonetic context, sequence-training
- NN output probabilities are *Gaussianised* by taking logs and decorrelating with PCA
- Early variants used purely NN features; later variants augmented the feature vector with standard acoustic features
- Can also use "bottleneck features" (narrow, intermediate NN layers)

▲□ ▶ ▲ □ ▶ ▲ □ ▶

## Tandem scheme



ASR Lecture 11 Neural Networks for Acoustic Modelling 2: HMM/DNN 12

< ∃⇒

## Tandem scheme



ASR Lecture 11 Neural Networks for Acoustic Modelling 2: HMM/DNN 12

→ E → < E →</p>

## Tandem scheme



ASR Lecture 11 Neural Networks for Acoustic Modelling 2: HMM/DNN 12

## Monophone HMM/NN hybrid system (1993)



## Monophone HMM/NN hybrid system (1998)



- Broadcast news transcription (1998) 20.8% WER
- (best GMM-based system, 13.5%)
- Cook et al, DARPA, 1999

A B > A B
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

- Advantages of NN:
  - Can easily model correlated features
    - Correlated feature vector components (eg spectral features)
    - Input context multiple frames of data at input
  - More flexible than GMMs not made of (nearly) local components); GMMs inefficient for non-linear class boundaries

向下 イヨト イヨト

- Advantages of NN:
  - Can easily model correlated features
    - Correlated feature vector components (eg spectral features)
    - Input context multiple frames of data at input
  - More flexible than GMMs not made of (nearly) local components); GMMs inefficient for non-linear class boundaries
- Disadvantages of NNs in the 1990s:
  - Context-independent (monophone) models, weak speaker adaptation algorithms
  - NN systems less complex than GMMs (fewer parameters): RNN < 100k parameters, MLP  $\sim$  1M parameters
  - Computationally expensive more difficult to parallelise training than GMM systems

・ロト ・回 ト ・ヨト ・ヨト

#### State of the art in the year 2000

by a statistically significant margin.



tion system. The final system output is a combination of P4a, P4b, P6a and P5b.

	CU-HTK 2000	
Base model	HMM-GMM	
Acoustic context	$\Delta$ , $\Delta\Delta$ features, HLDA projection	
Phonetic context	Tied state triphones & quinphones	
Speaker adaptation	Gender-dependent models, VTLN, MLLR	
Training criterion	ML + MMI sequence training	
System architecture	6-pass system	
Other features	Multi-system combination	
Hub 2000 WER	19.3%	

< ≣

	CU-HTK 2000	
Base model	HMM-GMM	
Acoustic context	$\Delta$ , $\Delta\Delta$ features, HLDA projection	
Phonetic context	Tied state triphones & quinphones	
Speaker adaptation	Gender-dependent models, VTLN, MLLR	
Training criterion	ML + MMI sequence training	
System architecture	6-pass system	
Other features	Multi-system combination	
Hub 2000 WER	19.3%	

## No neural networks!



200

	Microsoft 2011
Base model	HMM-DNN
Acoustic context	11 frames directly modelled
Phonetic context	Tied state triphones
Speaker adaptation	None
Training criteria	Frame-level cross-entropy
System architecture	Single pass
Other features	Deep network architecture
Hub 2000 WER	16.1%

- 4 回 ト 4 三 ト 4 三 ト

## The rise of deep neural networks



- **Deeper**: Deep neural network architecture – multiple hidden layers
- Wider: Use HMM state alignment as outputs rather than hand-labelled phones – 3-state HMMs, so 3×61 states
- Used a *pretraining* scheme to improve training accuracy of models with many hidden layers
- Training many hidden layers is computationally expensive – GPUs used to provide the computational power

・ロト ・ 同 ト ・ ヨ ト ・ ヨ ト

## Acoustic features for NN acoustic models

- GMMs: filter bank features (spectral domain) not used as they are strongly correlated with each other would either require
  - full covariance matrix Gaussians
  - many diagonal covariance Gaussians
- DNNs do not require the components of the feature vector to be uncorrelated
  - Can directly use multiple frames of input context (this has been done in NN/HMM systems since 1990, and is crucial to make them work well)
  - Can potentially use feature vectors with correlated components (e.g. filter banks)
- Mel-scaled filter bank features (FBANK) found to result in greater accuracy than standard MFCCs, though higher resolution MFCCs are now used

・ 同下 ・ ヨト ・ ヨト

## TIMIT phone error rates: effect of depth and feature type



## Context-dependent units



ASR Lecture 11 Neural Networks for Acoustic Modelling 2: HMM/DNN 23

ヘロト 人間 とくほとく ほとう

#### Tied context-dependent units



ASR Lecture 11 Neural Networks for Acoustic Modelling 2: HMM/DNN 24

ヘロン 人間 とくほど くほどう

## Modelling phonetic context (3)

- In the 1990s, this was considered hard (see earlier slides)
- But in 2011, a simple solution emerged: use state-tying from a GMM system

高 ト イ ヨ ト イ ヨ ト

## Modelling phonetic context (3)

- In the 1990s, this was considered hard (see earlier slides)
- But in 2011, a simple solution emerged: use state-tying from a GMM system

# Context-Dependent Pre-Trained Deep Neural Networks for Large-Vocabulary Speech Recognition

George E. Dahl, Dong Yu, Senior Member, IEEE, Li Deng, Fellow, IEEE, and Alex Acero, Fellow, IEEE

Abstract—We propose a novel context-dependent (CD) model for large-vocabulary speech recognition (LVSR) that leverages recent advances in using deep belief networks for phone recognition. We describe a pre-trained deep neural network hidden Markov model (DNN-HMM) hybrid architecture that trains the DNN to produce a distribution over senones (tied triphone states) as its output. The deep belief network pre-training algorithm is a robust and often helpful way to initialize deep neural networks generatively that fields (CRFs) [18]–[20], hidden CRFs [21], [22], and segmental CRFs [23]). Despite these advances, the elusive goal of human level accuracy in real-world conditions requires continued, viptant research.

Recently, a major advance has been made in training densely connected, directed belief nets with many hidden layers. The resulting deep belief nets learn a hierarchy of nonlinear feature

イロト イポト イヨト イヨト

## Context-dependent hybrid HMM/DNN

- First train a context-dependent HMM/GMM system on the same data, using a phonetic decision tree to determine the HMM tied states
- Perform Viterbi alignment using the trained HMM/GMM and the training data
- Train a neural network to map the input speech features to a label representing a context-dependent tied HMM state
  - So the size of the label set is thousands (number of context-dependent tied states) rather than tens (number of context-independent phones) Each frame is labelled with the Viterbi aligned tied state
- Train the neural network using gradient descent as usual
- Use the context-dependent scaled likelihoods obtained from the neural network when decoding

・ロト ・回 ト ・ヨト ・ヨト

## Example: HMM/DNN acoustic model for Switchboard



ASR Lecture 11 Neural Networks for Acoustic Modelling 2: HMM/DNN 27

크

## Example: HMM/DNN acoustic model for Switchboard

- Alignments generated from context-dependent HMM/GMM system
- Hybrid HMM/DNN system
  - $\bullet\,$  Context-dependent 9304 output units obtained from Viterbi alignment of HMM/GMM system
  - 7 hidden layers, 2048 units per layer
  - 11 frames of acoustic context
- DNN-based system results in significant word error rate reduction compared with GMM-based system

A (1) < A (2) < A (2) </p>

## Summary

- DNN/HMM systems (hybrid systems) gave a significant improvement over GMM/HMM systems
- Compared with 1990s NN/HMM systems, DNN/HMM systems
  - model context-dependent tied states with a much wider output layer
  - are deeper more hidden layers
  - can use correlated features (e.g. FBANK) or higher resolution MFCCs
- Background reading:
  - N Morgan and H Bourlard (May 1995). "Continuous speech recognition: Introduction to the hybrid HMM/connectionist approach", *IEEE Signal Processing Mag.*, **12**(3), 24–42. http://ieeexplore.ieee.org/document/382443
  - A Mohamed et al (2012). "Understanding how deep belief networks perform acoustic modelling", Proc ICASSP-2012. http://www.cs.toronto.edu/~asamir/papers/icassp12\_ dbn.pdf