



# Hierarchical Tandem Features for ASR in Mandarin

Joel Pinto<sup>1</sup>, Mathew Magimai.-Doss<sup>1</sup>, Hervé Bourlard<sup>1,2</sup>

<sup>1</sup>Idiap Research Institute, Martigny, Switzerland

<sup>2</sup>École Polytechnique Fédérale de Lausanne (EPFL), Switzerland

{jpinto, mathew, bourlard}@idiap.ch

## Abstract

We apply multilayer perceptron (MLP) based hierarchical Tandem features to large vocabulary continuous speech recognition in Mandarin. Hierarchical Tandem features are estimated using a cascade of two MLP classifiers which are trained independently. The first classifier is trained on perceptual linear predictive coefficients with a 90 ms temporal context. The second classifier is trained using the phonetic class conditional probabilities estimated by the first MLP, but with a relatively longer temporal context of about 150 ms. Experiments on the Mandarin DARPA GALE *eval06* dataset show significant reduction (7.6% relative) in character error rates by using hierarchical Tandem features over conventional Tandem features.

**Index Terms:** Automatic speech recognition, multilayer perceptrons, Tandem features, hierarchical systems.

## 1. Introduction

Multilayer perceptron (MLP) classifier based acoustic modeling is being extensively used in state-of-the-art automatic speech recognition (ASR) systems [1][2][3][4]. The MLP classifier is typically trained using standard acoustic features such as perceptual linear predictive (PLP) coefficients with a certain temporal context. The phonetic class-conditional probabilities estimated by the MLP are transformed (e.g. logarithm and Karhunen Loeve transformation) to obtain Tandem features which are subsequently used in the standard HMM/GMM based ASR system.

The phonetic class conditional probabilities estimated by the MLP (a vector) at a particular time instant represents the instantaneous soft-decision on the underlying phoneme and carries useful information such as the probability mass assigned to the competing phonemes. A temporal context on the phonetic class conditional probabilities (a sequence of vectors) carries additional contextual information such as the transition of the estimated probabilities within a phoneme and across neighboring phonemes. To exploit this contextual information, we previously studied a hierarchical architecture to estimate the phonetic class-conditional probabilities [5]. In this hierarchical system, a second MLP classifier is trained on the phonetic class conditional probabilities (or posterior features) estimated by the MLP with a long temporal context of 150-230 ms. The estimated class conditional probabilities are used in the same way as the single MLP based estimator.

There have been works in the literature motivated towards exploiting contextual information in the posterior features using classifiers such as HMM [6] and conditional random fields [7]. Our study [5] discusses in detail the similarities and differences of the MLP based hierarchical system with other works in the literature. It also provides an in-depth analysis of the second

MLP classifier using Volterra series showing that it learns the phonetic-temporal confusion patterns in the posterior features and to a certain extent the phonotactics of the language as observed in the training data.

The effectiveness of the hierarchical system has been previously evaluated in recognition of phonemes in read speech recorded in clean conditions (TIMIT) as well as conversational speech recorded over a telephone channel (CTS) [5], where significant reduction in error rates have been observed. The reduction in word error rates in small vocabulary isolated word recognition was reported in [8]. In both these works the HMM/MLP hybrid system was used, where the phonetic class conditional probabilities estimated by the MLP were used as local scores in the states of the HMM. The objective of this work is to demonstrate the effectiveness of hierarchical Tandem features estimated using a cascade of two MLP classifiers in large vocabulary continuous speech recognition in challenging real-world scenarios.

We use the Mandarin database developed under the Global Autonomous Language Exploitation (GALE) project. On the state-of-the-art SRI-ICSI-UW ASR system [9], the hierarchical Tandem features yield an absolute (relative) reduction of 2.8% (8.7%) in character errors on broadcast conversations and a reduction of 1.1% (6.1%) in character errors on broadcast news when compared to the baseline Tandem features. The hierarchical Tandem features also outperform the standard MFCC plus pitch features. Furthermore reduction in errors is observed when hierarchical Tandem features are augmented with MFCC features.

## 2. Experimental Setup

The DARPA GALE task involves recognition of speech in audio segments acquired from various television programs broadcast in Mandarin. The broadcast segments include two types of genres, namely broadcast news (BN) and broadcast conversations (BC). The broadcast programs span a wide range of domains which include informal and colloquial language. In this section, we describe the experimental setup for the Mandarin ASR system.

### 2.1. Training and Test Data Definition

The training corpus consists of 95 hours of speech, which includes 50 hours of BN and 45 hours BC data. It is a subset of the training set of the 2008 SRI Mandarin speech-to-text system [10]. The snippet level genre classification on the training set was provided by SRI using a technique described in [11]. We use the GALE *eval06* data as the test set. The genre labels on the test set are provided by the Linguistic Data Consortium.

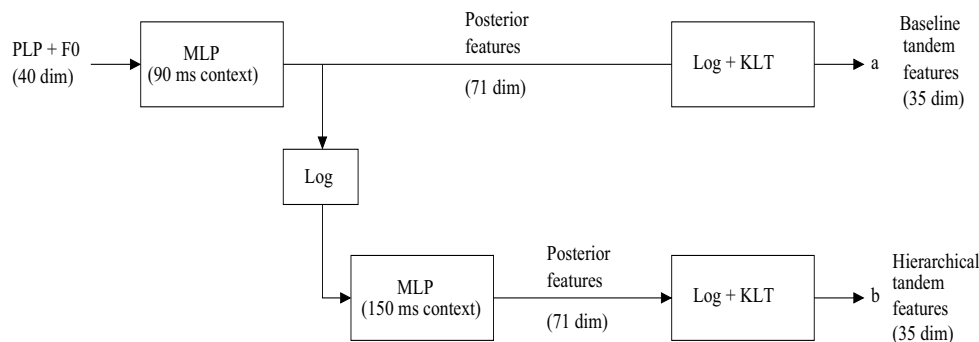


Figure 1: (a) Standard Tandem feature extraction (b) Hierarchical Tandem feature extraction.

## 2.2. Hierarchical System

Fig. 1 shows the details of baseline and hierarchical Tandem feature extraction used in the experiments. The input features to the first MLP consists of the first 13 PLP cepstral coefficients appended to their delta and delta-delta parameters. As Mandarin is a tonal language, a smoothed estimate of the log-pitch value is appended to the cepstral features [12]. The 40 dimensional combined feature vector is applied at the input of the MLP with a temporal context of 90 ms. The size of the output layer of the MLP is 71 which corresponds to the number of phonemes. Only the first 35 components of the Tandem features are retained to capture at least 95% of the total variance in the data.

In the hierarchical Tandem system, a second MLP classifier is trained on the log posterior features estimated by the first MLP with a temporal context of 150 ms. This temporal context is based on the findings from task adaptation studies reported in [8], where it was observed that the word error rates begin to saturate at a context of around 130 ms - 150 ms. The output of the second MLP is transformed in the same way as the baseline system to obtain the hierarchical Tandem features.

Three layered MLP classifiers are used with sigmoid non-linearity at the hidden layer and softmax nonlinearity at the output layer. The size of the hidden layer of the MLPs was chosen such that the total number of parameters is roughly equal to 5% of the total number of training samples.

## 2.3. Mandarin ASR System

We use the SRI-ICSI-UW Mandarin ASR system [12][9][10] developed for the DARPA GALE program. More specifically, we use the system setup described in [13] and this is briefly discussed here.

Speech-silence segmentation and automatic speaker clustering is first performed using Gaussian mixture modeling technique to derive “auto speakers”. The vocal tract length normalization factors are estimated for each auto speaker and are used in the estimation of MFCC features.

The acoustic modeling is based on the standard HMM/GMM technique. In the training phase, context independent models are first trained for each of the 71 phonemes. Context dependent models are subsequently trained and clustered down to 2000 shared Markov states, which are also known as senones. Each senone is modeled using a mixture of 32 Gaussians using phonetic decision tree based clustering. The acoustic model parameters are trained using the simple maximum likelihood criterion. Cross-word triphone

modeling and speaker adaptive training is not performed in this study.

A trigram language model, which was estimated using an assortment of text corpora totalling over a billion words was used for this study. The pronunciation dictionary consists of 60K characters, and is transcribed using 70 phonemes. A silence class was added, resulting in a total of 71 output classes.

The ASR system is evaluated in two modes - speaker independent and speaker adapted. The speaker independent system involves a single pass maximum likelihood decoding using a trigram language model. The speaker adaptive system includes an additional pass, where the features are transformed using a two class (speech/silence) constrained maximum likelihood transform, estimated for each of the auto-speakers.

## 2.4. Methodology

The HMM/GMM system is trained using three sets of features:

- mfcc-f0-42: The static feature vector consists of first 13 MFCC coefficients along with an estimate of the log pitch value (f0). The static features are appended to their first and second order temporal derivatives to obtain a 42 dimensional feature vector.
- tandem-35: The phoneme posterior probabilities estimated by the MLP classifier are transformed using logarithm and KLT, followed by dimensionality reduction to obtain a 35 dimensional feature vector. The tandem features are estimated in the conventional way using a single MLP classifier or the hierarchical approach as discussed in Fig. 1.
- mfcc-f0-tandem-77: Studies have shown that the best performance using Tandem features have been obtained when they are concatenated with the standard acoustic features [1]. To this end, we investigate an augmented feature vector formed by the concatenation of mfcc-f0-42 and tandem-35 features.

To distinguish between Tandem features estimated by standard single MLP approach and hierarchical approach, we prefix the features with qualifiers “baseline” and “hierarchical” respectively. For instance, baseline tandem-35 refers to Tandem feature estimated by the standard single MLP approach. The MLP classifiers were trained at Idiap Research Institute using the Quicknet toolkit.<sup>1</sup> Training ASR models and recognition experiments were performed at ICSI, Berkeley using the SRI Decipher system.

<sup>1</sup><http://www.icsi.berkeley.edu/Speech/qn.html>

Features	BC genre		BN genre		Both genres	
	SI (%)	SA (%)	SI (%)	SA (%)	SI (%)	SA (%)
mfcc-f0-42	33.4	31.0	20.9	19.3	27.0	25.0
baseline tandem-35	34.4	32.7	19.5	17.9	26.8	25.1
hierarchical tandem-35	<b>31.3</b>	<b>29.9</b>	<b>18.0</b>	<b>16.8</b>	<b>24.5</b>	<b>23.2</b>

Table 1: CERs obtained using mfcc-f0-42, baseline tandem-35, and hierarchical tandem-35 features. Boldface indicates the lowest CER.

Features	BC genre		BN genre		Both genres	
	SI (%)	SA (%)	SI (%)	SA (%)	SI (%)	SA (%)
baseline mfcc-f0-tandem-77	29.2	28.0	17.6	16.6	23.3	22.2
hierarchical mfcc-f0-tandem-77	<b>28.4</b>	<b>27.3</b>	<b>17.0</b>	<b>16.3</b>	<b>22.5</b>	<b>21.7</b>

Table 2: CERs obtained using baseline mfcc-f0-tandem-77 and hierarchical mfcc-f0-tandem-77 features. Boldface indicates the lowest CER.

### 3. Experimental Results

The speaker independent (SI) and speaker adapted (SA) systems are evaluated in terms of the character error rate (CER). The results are reported for the individual genres as well as on the entire test set.

Table 1 shows the character error rates on the *eval06* dataset obtained using mfcc-f0-42 features, baseline tandem-35 features, and hierarchical tandem-35 features. It can be seen that on broadcast conversations, the mfcc-f0-42 features yield a lower CER when compared to the baseline tandem-35 features. On broadcast news, the opposite trend is observed. Hence, on the entire test set, the mfcc-f0-42 and tandem-35 features yield similar performance. The hierarchical tandem-35 features yield the lowest CER on both broadcast news as well as broadcast conversations. These results clearly demonstrate the effectiveness of the MLP based hierarchical acoustic modeling in large vocabulary continuous speech recognition. The other main observations from this study are the following:

(1) The error rates on the BC genre are significantly higher when compared to the BN genre as observed in some of the previous works in the literature [11]. Recognition on the BC genre is significantly harder when compared to BN because (a) the conversational speech is spontaneous in nature and characterized by variable speaking rate, mispronunciations, false starts, hesitations, disfluencies etc. and (b) the language model, which is estimated from text is more closer to the broadcast news than conversations.

(2) The reduction in CER by using the hierarchical system is higher in the case of BC genre when compared to the BN genre. A similar trend was also observed in the recognition of phonemes [5]. On the CTS task (conversations), the hierarchical approach resulted in an absolute increase of 9.0% in the phoneme accuracy over the baseline single MLP based system. On TIMIT (read speech), the improvement in the recognition accuracy was only 3.5%. This could be because in conversational speech, there is a greater scope to correct the erroneous speech-to-phoneme mapping - caused by the artifacts in conversational speech discussed above - by exploiting the phonetic-temporal information. This aspect needs further attention.

(3) The decrease in CER obtained by using the hierarchical tandem-35 features over the baseline tandem-35 features is slightly higher in the case of speaker independent decoding when compared to the speaker adaptive decoding. It can be seen

that on the BC (BN) genre, the decrease in CER is about 3.1% (1.5%) on the speaker independent system, whereas on speaker adapted system, the decrease is about 2.8% (1.1%).

Table 2 shows the CERs obtained using the baseline mfcc-f0-tandem-77 and hierarchical mfcc-f0-tandem-77 features. The important observations from the table in conjunction with Table 1 are as follows:

(1) Significant reduction in character error rates is observed when mfcc-f0-42 features are augmented with baseline tandem-35 features (a relative decrease of 9.7% on BC and 7.3% on BN compared to the best individual stream). This shows that mfcc-f0-42 and baseline tandem-35 features bear complimentary information. Further decrease in CER is observed by using hierarchical tandem-35 features. This shows that the hierarchical tandem-35 and mfcc-f0-42 features bear complimentary information in the same way baseline tandem-35 and mfcc-f0-42 features.

(2) The improvement in performance obtained by using hierarchical Tandem features over the baseline Tandem features is reduced when these features are augmented with mfcc-f0-42 features. This suggests that the improvement in recognition accuracies obtained by feature concatenation and hierarchical processing is not exactly additive. Nonetheless, the hierarchical mfcc-f0-tandem-77 features yield the lowest error rates in both the BC and BN genres.

To summarize, in this experimental setup, the lowest CER of 21.7% on the combined test set (BN and BC) is obtained using the hierarchical mfcc-f0-tandem-77, which is an absolute (relative) decrease of 3.3% (13.2%) over the conventional mfcc-f0-42 features.

### 4. Discussion

Experimental results have shown that the hierarchical Tandem features outperform the conventional Tandem features in large vocabulary ASR. This is consistent with our previous study on recognition of phonemes. In this section, we discuss two differences in the hierarchical system for phoneme recognition and word recognition mainly based on empirical studies.

In recognition of phonemes, a temporal context of 210-230 ms on the posterior features yielded the lowest phoneme error rates [5]. In recognition of isolated words [8] as well as continuous words in this work, a temporal context of 150 ms yields the best performance. Beyond this point, the error rates begin to

saturate. This suggests that learning the phonotactics implicitly by taking a longer temporal context on the posterior features is more useful in the recognition of phonemes than in the recognition of words. This aspect needs to be investigated further.

It was observed that slightly lower error rates are obtained by training the second MLP using log posterior features rather than directly using raw posterior features. A similar observation was also made in recognition of isolated words. Taking a logarithm of the output of the MLP with a softmax output non-linearity is equivalent to taking its linear activation values, except for a constant additive factor. On recognition of phonemes, however, this did not make any difference.

Hierarchical architecture of MLPs have been found to be beneficial in task adaptation yielding lower error rates when compared to other adaptation methods such as full model re-training and partial (only hidden-to-output weights) model re-training [8] [14]. A possible extension of this work is to genre adaptation where the first MLP in the hierarchical system is trained using data from different genres but the second MLP is trained on the genre-specific data. For example, by training the first MLP using data from both BN and BC genres and training the second MLP using only BC genre, we have observed further reduction in character error rates on the BC genre [14].

## 5. Conclusions

Experimental studies on the DARPA GALE Mandarin task show that the hierarchical Tandem features yield significantly lower character error rates when compared to the standard Tandem features as well as the traditional MFCC features. This reduction in error rates is observed in both broadcast conversations and broadcast news genres. Further reduction in error rates is obtained when the hierarchical Tandem features are augmented with standard MFCC features.

## 6. Acknowledgements

This work was supported in parts by the Swiss National Science Foundation through the Indo-Swiss joint research program on keyword spotting, the Swiss National Center for Competence in Research through the Interactive Multimodal Information Management (IM2) project, and by DARPA through the GALE program. Any opinions, findings and conclusions expressed in this material are those of the authors and do not necessarily reflect the views of funding agencies. We acknowledge SRI for the permission to use the Decipher ASR system and ICSI for the computational infrastructure. We also thank Wen Wang from SRI and Suman Ravuri from ICSI for their help in setting up the experiments.

## 7. References

- [1] N. Morgan *et al.*, "Pushing the Envelope - Aside," *IEEE Signal Process. Magazine*, vol. 22, no. 5, pp. 81–88, 2005.
- [2] A. Stolcke *et al.*, "Recent Innovations in Speech-to-Text Transcription at SRI-ICSI-UW," *IEEE Trans. Audio. Speech. Language. Process.*, vol. 14, no. 5, pp. 1729–1744, 2006.
- [3] P. Fousek, L. Lamel, and J.-L. Gauvain, "Transcribing Broadcast Data using MLP Features," *Proc. of Interspeech*, pp. 1433–1436, 2008.
- [4] J. Park, F. Diehl, M. Gales, M. Tomalin, and P. Woodland, "Training and Adapting MLP Features for Arabic Speech Recognition," *Proc. of IEEE Conf. Acoust. Speech. Signal Process. (ICASSP)*, pp. 4461–4464, 2009.
- [5] J. Pinto, G. Sivaram, M. Magimai.-Doss, H. Hermansky, and H. Bourlard, "Analyzing MLP Based Hierarchical Phoneme Posterior Probability Estimator," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 2, pp. 225–241, 2011.
- [6] H. Ketabdard, J. Vepa, S. Bengio, and H. Bourlard, "Using More Informative Posterior Probabilities for Speech Recognition," *Proc. of IEEE Conf. Acoust. Speech. Signal Process. (ICASSP)*, pp. 29–32, 2006.
- [7] E. Fosler-Lussier and J. Morris, "CRANDEM Systems: Conditional Random Field Acoustic Models for Hidden Markov Models," *Proc. of IEEE Conf. Acoust. Speech. Signal Process. (ICASSP)*, pp. 4049–4052, 2008.
- [8] J. Pinto, M. Magimai.-Doss, H., and H. Bourlard, "MLP Based Hierarchical System for Task Adaptation in ASR," *Proc. of IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU)*, 2009.
- [9] M.-Y. Hwang, G. Peng, M. Ostendorf, W. Wang, A. Faria, and A. Heidel, "Building a Highly Accurate Mandarin Speech Recognizer with Language-independent Technologies and Language-Dependent Modules," *IEEE Trans. on Audio, Speech, and Language Process.*, vol. 17, no. 7, pp. 1253–1262, 2009.
- [10] X. Lei, W. Wu, W. Wang, A. Mandal, and A. Stolcke, "Development of the 2008 SRI Mandarin Speech-to-text System for Broadcast News and Conversation," in *Proc. of Interspeech*, 2009, pp. 2099–2103.
- [11] W. Wang, A. Mandal, X. Lei, A. Stolcke, and J. Zheng, "Multifactor Adaptation for Mandarin Broadcast News and Conversation Speech Recognition," in *Proc. of Interspeech*, 2009, pp. 2103–2102.
- [12] X. Lei, M. Siu, M. Ostendorf, and T. Lee, "Improved Tone Modeling for Mandarin Broadcast News Speech Recognition," *Proc. of Interspeech*, pp. 1237–1240, 2006.
- [13] F. Valente, M. Magimai.-Doss, C. Plahl, and S. Ravuri, "Hierarchical Processing of the Modulation Spectrum for GALE Mandarin LVCSR System," in *Proc. of Interspeech*, 2009, pp. 2963–2966.
- [14] J. Pinto, "Multilayer Perceptron based Hierarchical Acoustic Modeling for Automatic Speech Recognition," Ph.D. dissertation, École Polytechnique Fédérale de Lausanne (EPFL), Switzerland, July 2010.