



Available online at www.sciencedirect.com

ScienceDirect

Computer Speech and Language xxx (2015) xxx–xxx

COMPUTER
SPEECH AND
LANGUAGE

www.elsevier.com/locate/csl

Articulatory feature based continuous speech recognition using probabilistic lexical modeling[☆]

Ramy Rasipuram^{a,b,*}, Mathew Magimai.-Doss^a

^a *Idiap Research Institute, Martigny, Switzerland*

^b *Ecole Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Switzerland*

Received 3 July 2014; received in revised form 12 February 2015; accepted 23 April 2015

Abstract

Phonological studies suggest that the typical subword units such as phones or phonemes used in automatic speech recognition systems can be decomposed into a set of features based on the articulators used to produce the sound. Most of the current approaches to integrate articulatory feature (AF) representations into an automatic speech recognition (ASR) system are based on a deterministic knowledge-based phoneme-to-AF relationship. In this paper, we propose a novel two stage approach in the framework of probabilistic lexical modeling to integrate AF representations into an ASR system. In the first stage, the relationship between acoustic feature observations and various AFs is modeled. In the second stage, a probabilistic relationship between subword units and AFs is learned using transcribed speech data. Our studies on a continuous speech recognition task show that the proposed approach effectively integrates AFs into an ASR system. Furthermore, the studies show that either phonemes or graphemes can be used as subword units. Analysis of the probabilistic relationship captured by the parameters has shown that the approach is capable of adapting the knowledge-based phoneme-to-AF representations using speech data; and allows different AFs to evolve asynchronously.

© 2015 Elsevier Ltd. All rights reserved.

Keywords: Automatic speech recognition; Articulatory features; Probabilistic lexical modeling; Kullback–Leibler divergence based hidden Markov model; Phoneme subword units; Grapheme subword units

1. Introduction

Articulatory features describe the properties of speech production, i.e., each sound unit of a language, a phone or a phoneme, can be decomposed into a set of features based on the articulators used to produce it. The use of articulatory feature (AF) representations in an automatic speech recognition (ASR) system is motivated by their abilities such as:

- Better pronunciation modeling: AFs are hypothesized to capture acoustic variation at a finer level than the phoneme-based representation (Deng et al., 1997; Richardson et al., 2003; Livescu et al., 2008).

[☆] This paper has been recommended for acceptance by Karen Livescu.

* Corresponding author at: Idiap Research Institute, Martigny, Switzerland. Tel.: +41 277217711; fax: +41 277217712.
E-mail addresses: ramya.rasipuram@idiap.ch (R. Rasipuram), mathew@idiap.ch (M. Magimai.-Doss).

- Robustness to noise: Different AFs may have variable noise sensitivity. The “divide and conquer” approach provides a framework to exploit the variable noise sensitivity of AFs (Kirchhoff et al., 2002).
- Multilingual and crosslingual portability: AFs can provide better sharing capabilities than phonemes across languages (Stüker et al., 2003; Lal and King, 2013; Siniscalchi et al., 2012).

To incorporate the articulatory knowledge in an ASR system, the following main concerns have to be addressed:

1. AF representations: There exist different types of articulatory representations of speech, e.g., binary features, multivalued features, and government phonological features. AFs defined by Chomsky and Halle (1968) are binary valued features, for example +voice and –voice, +sonorant and –sonorant. However, according to Ladefoged (1993), it is more natural to allow AFs to take multiple values. In government phonological feature system, speech sounds are destructed into a set of primes and can be represented by fusing them structurally (Harris, 1994). In the paper, we have used the multivalued AFs, because it has been argued that they better represent non-binary parameters such as height of vowels and place of articulation (Ladefoged, 1993).
2. Estimation of AFs from acoustic speech signal: In the literature, many approaches have been explored to extract AFs from the acoustic speech signal. For example, techniques based on acoustic-to-articulatory feature codebooks (Hogden et al., 1996; Suzuki et al., 1998), artificial neural networks (Livescu et al., 2008; Kirchhoff et al., 2002; Chang, 2002; Rasipuram and Magimai-Doss, 2011), support vector machines (Juneja and Espy-Wilson, 2004; Scharenborg et al., 2007), Gaussian mixture models (Metze and Waibel, 2002; Stüker et al., 2003), hidden Markov models (Hiroya and Honda, 2004), conditional random fields (Prabhavalkar et al., 2011), nearest neighbour (Næss et al., 2011), dynamic Bayesian networks (Frankel and King, 2005; Frankel et al., 2007) are used.
3. Integration: Integrating AFs into the conventional hidden Markov model (HMM) based ASR framework is a challenging task mainly because of the multiple AF estimators. The dynamic Bayesian network (DBN) based approaches for AF integration preserve the articulatory representation in DBN state space (Livescu and Glass, 2004; Livescu et al., 2008; King et al., 2007). These approaches have shown promising results in lexical access¹ experiments. Posterior probabilities of AFs can be transformed for use as features in tandem speech recognition systems (Cetin et al., 2007, 2007; Lal and King, 2013). Posterior probabilities of AFs are also used to enhance phoneme-based acoustic models (Kirchhoff et al., 2002; Siniscalchi et al., 2012). These approaches however lose other benefits of articulatory representation such as finer granularity and asynchronous evolution.

In this paper, we propose an approach in the framework of probabilistic lexical modeling to integrate multivalued AFs. In a probabilistic lexical model based ASR system, the relationship between subword units in the lexicon and acoustic feature observations is factored into two models using latent variables: An acoustic model which models the relationship between acoustic feature observations and the latent variables; and a lexical model which models a probabilistic relationship between the subword units in the lexicon and the latent variables. In this paper, we show that by choosing the latent variables as multiple multivalued AFs, the approach effectively integrates AFs into the HMM-based ASR framework. The lexical model parameters in the proposed approach capture a probabilistic relationship between subword units and AFs learned through transcribed speech data.

The potential of the proposed approach for AF integration is demonstrated on a continuous speech recognition task through experiments and comparisons with the tandem approach. In the proposed framework we explore the use of domain-independent data for acoustic model training; and phonemes and graphemes as subword units. Furthermore, through the analysis of the lexical model parameters we show that the approach adapts the knowledge-based phoneme-to-AF or grapheme-to-AF relationship and allows different AFs to evolve asynchronously.

The rest of the paper is organized as follows: Section 2 gives an overview of the HMM-based ASR and the framework of probabilistic lexical modeling. Section 3 presents the literature review of approaches that integrate multivalued AFs for ASR in the light of the background information given in Section 2. In Section 4, the approach for AF integration is presented and the contributions of the present paper with respect to prior work are elaborated. Sections 5 and 6 present the experimental setup and results, respectively. Section 7 presents an analysis on the subword-unit-to-AF relationship captured by the lexical model parameters. Finally, in Section 8 we provide a discussion and conclusion.

¹ The task of lexical access involves predicting a word given its phonetic or broad phonetic transcription (Huttenlocher and Zue, 1984).

2. Background

The goal of ASR is to find the most likely word sequence $W^* = [w_1, \dots, w_m, \dots, w_M]$ given the acoustic observation sequence $X = [x_1, \dots, x_t, \dots, x_T]$ where M is the number of words in the utterance and T represents the number of frames in the speech signal. The most likely word sequence W^* given the acoustic observation sequence is obtained as follows:

$$W^* = \arg \max_{W \in \mathcal{W}} P(W|X) \quad (1)$$

$$W^* = \arg \max_{W \in \mathcal{W}} p(X|W)P(W) \quad (2)$$

where \mathcal{W} denotes the set of all possible word sequences. The first term on the right hand side of Eq. (2) denotes the acoustic likelihood and the second term denotes the language model probability.

In general, speech recognition systems model words as a sequence of subword units, which are further modeled as a sequence of HMM states. The sequence of subword units for a word is given by its pronunciation model as specified in the pronunciation lexicon. The acoustic likelihood in an HMM-based ASR system is computed as follows:

$$p(X|W, \Theta_A) = \sum_{Q \in \mathcal{Q}} p(X|Q, W, \Theta_A)P(Q|W, \Theta_A) \quad (3)$$

$$p(X|W, \Theta_A) = \sum_{Q \in \mathcal{Q}} p(X|Q, \Theta_A)P(Q|W, \Theta_A) \quad (4)$$

$$p(X|W, \Theta_A) \approx \max_{Q \in \mathcal{Q}} p(X|Q, \Theta_A)P(Q|W, \Theta_A) \quad (5)$$

$$p(X|W, \Theta_A) \approx \max_{Q \in \mathcal{Q}} \left[\prod_{t=1}^T p(x_t|q_t, \Theta_A)P(q_t|q_{t-1}, \Theta_A) \right] \quad (6)$$

The set Θ_A includes the parameters of the acoustic likelihood estimator. In Eq. (3), the acoustic likelihood is obtained by summing over all possible state sequences \mathcal{Q} where each $Q = [q_1, \dots, q_t, \dots, q_T]$ denotes a sequence of HMM states corresponding to a word sequence hypothesis. Eq. (4) assumes that the acoustic likelihood is independent of words given the state sequence. In Eq. (5), a Viterbi approximation is employed where the sum over all possible state sequences is replaced with the most probable state sequence. Eq. (6) arises from the two HMM assumptions, i.e., acoustic feature observations are conditionally independent of each other and the HMM state at time t depends only on the HMM state at time $t - 1$.

In subword unit based ASR systems, HMM states represent lexical units, i.e., $q_t \in \mathcal{L} = \{l^1, \dots, l^i \dots l^I\}$ and I is the number of lexical units. If context-independent phonemes are used as subword units then the number of lexical units $I = M \times K$ where K is the number of context-independent subword units in the lexicon and M is the number of HMM states for each context-independent phoneme. If context-dependent phonemes are used as subword units then the number of lexical units $I = M \cdot K^{c_r + c_l + 1}$ where c_l is the preceding context length, c_r is the following context length. Typically, each context-independent or context-dependent phoneme is modeled with three HMM states, i.e., $M = 3$. In HMM-based ASR systems, the relationship between lexical units and acoustic features is not always modeled directly. As we will see in the remainder of the section, the relationship is typically modeled through intermediate acoustic units (for example, clustered context-dependent subword states).

2.1. Framework of probabilistic lexical modeling

In the framework of probabilistic lexical modeling (Rasipuram and Magimai-Doss, 2015), relationship between the acoustic feature observation x_t and the lexical unit l^i is factored through a *latent* variable a^d as follows:

$$p(x_t|q_t = l^i, \Theta_A) = \sum_{d=1}^D p(x_t, a^d|q_t = l^i, \Theta_A) \quad (7)$$

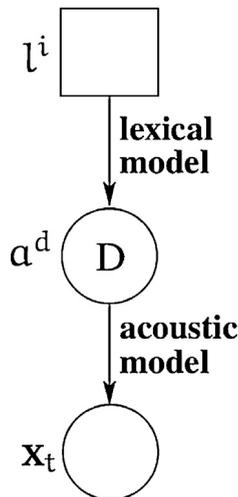


Fig. 1. The graphical model representation of a system incorporating probabilistic lexical modeling. The D in the middle of the node indicates that there are D acoustic units.

$$p(\mathbf{x}_t | q_t = l^i, \Theta_A) = \sum_{d=1}^D p(\mathbf{x}_t | a^d, q_t = l^i, \theta_a, \theta_l) \cdot P(a^d | q_t = l^i, \theta_l) \tag{8}$$

$$p(\mathbf{x}_t | q_t = l^i, \Theta_A) = \sum_{d=1}^D \underbrace{p(\mathbf{x}_t | a^d, \theta_a)}_{\text{acoustic model}} \cdot \underbrace{P(a^d | q_t = l^i, \theta_l)}_{\text{lexical model}} \tag{9}$$

The parameters of the acoustic likelihood estimator Θ_A encompass the *acoustic model* (θ_a), the *pronunciation lexicon* (θ_{pr}) and the *lexical model* (θ_l) parameters, therefore, $\Theta_A = \{\theta_a, \theta_{pr}, \theta_l\}$. The relationship in Eq. (9) is as a result of the assumption that given a^d , $p(\mathbf{x}_t | a^d, q_t = l^i, \theta_a, \theta_l)$ is independent of $q_t = l^i$. In Eq. (9), $p(\mathbf{x}_t | a^d, \theta_a)$ is the acoustic unit likelihood and $P(a^d | q_t = l^i, \theta_l)$ is the probability of the latent variable given the lexical unit. We refer to $p(\mathbf{x}_t | a^d, \theta_a)$ as the acoustic model, $P(a^d | q_t = l^i, \theta_l)$ as the lexical model, the latent variable a^d as the acoustic unit, the set of acoustic units $\mathcal{A} = \{a^1, \dots, a^d, \dots, a^D\}$ and D as the number of acoustic units.

Fig. 1 shows the Bayesian network of an ASR system that uses the factorization of Eq. (9). The lexical unit is given deterministically by the current word and its subword units. The lexical unit is mapped to all the acoustic units and the acoustic feature observation is conditioned on all the acoustic units.

2.2. Lexical and acoustic units

In the case of context-independent ASR systems, the lexical unit set \mathcal{L} and the acoustic unit set \mathcal{A} are knowledge driven and defined based on the subword units in the pronunciation lexicon. The number of lexical units or acoustic units $I = D = M \times K$, typically, $M = 3$.

In the case of context-dependent ASR systems, the number of lexical units $I = M \cdot K^{c_r+c_l+1}$. Generally, not all context-dependent subword units will appear sufficiently often in the training data. Hence a sharing approach is used to enable multiple lexical units to share an acoustic model. This is done using a decision-tree based state clustering and tying technique that uses a pronunciation lexicon, linguistic knowledge and acoustic data to prepare a phonetic question set (Young et al., 1994). The number of acoustic units D varies depending on hyper parameters such as the state occupancy count and the log-likelihood threshold that are used during decision-tree based state clustering. However, the number of acoustic units D is well below the number of lexical units I . The resulting acoustic units are typically referred as clustered context-dependent states or tied-HMM states.

Other possibilities for the choice of the acoustic units are fenones (Bahl et al., 1988), senones (Hwang and Huang, 1992), automatically derived units from the acoustic data (Holter and Svendsen, 1997), etc. In this paper, we show that HMM-based ASR systems can be built using multivalued AFs as the acoustic units.

2.3. Acoustic model

The acoustic model which models the relationship between the acoustic feature observation \mathbf{x}_t and the acoustic unit a^d can be trained using either generative approaches such as Gaussian mixture models (GMMs) or discriminative approaches such as artificial neural networks (ANN).

GMMs (Rabiner, 1989): In the GMM approach, the acoustic model score $p(\mathbf{x}_t|a^d, \theta_a)$ is estimated given a mixture of C^d Gaussians that model an acoustic unit a^d , i.e.,

$$p(\mathbf{x}_t|a^d, \theta_a) = \sum_{c=1}^{C^d} w_c^d \mathcal{N}(\mathbf{x}_t, \mu_c^d, \Sigma_c^d) \quad (10)$$

where w_c^d , μ_c^d and Σ_c^d are the weight, means and covariances of the mixture component c of the acoustic unit a^d . In the literature, there are various variants of the GMM approach such as semi-continuous GMMs (Huang and Jack, 1989), subspace GMMs (Povey et al., 2011).

ANN (Morgan and Bourlard, 1995): The ANN computes the probability of acoustic units given the acoustic feature observations $P(a^d|\mathbf{x}_t, \theta_a)$ which is then converted to scaled-likelihood, i.e.,

$$p_{sl}(\mathbf{x}_t|a^d, \theta_a) = \frac{p(\mathbf{x}_t|a^d, \theta_a)}{p(\mathbf{x}_t)} = \frac{P(a^d|\mathbf{x}_t, \theta_a)}{P(a^d)} \quad (11)$$

2.4. Lexical model

The lexical model in an ASR system can be deterministic or probabilistic.

2.4.1. Deterministic lexical model based ASR approaches

When the lexical model is deterministic, each lexical unit l^i is deterministically mapped to an acoustic unit a^j ($l^i \mapsto a^j$), i.e.,

$$P(a^d|q_t = l^i, \theta_l) = \begin{cases} 1, & \text{if } d = j \\ 0, & \text{otherwise} \end{cases} \quad (12)$$

As a result of the deterministic mapping, the only term contributing to the summation in Eq. (9) is the acoustic unit that is mapped to a lexical unit.

In standard HMM-based ASR approaches such as HMM/GMM (Rabiner, 1989) and hybrid HMM/ANN (Morgan and Bourlard, 1995), the lexical model is deterministic. In the case of context-independent ASR systems, it is a *knowledge-based look-up table* that maps each lexical unit to an acoustic unit. In the case of context-dependent ASR systems, typically the lexical model is obtained through clustering and state tying using *decision trees*. The decision trees map each context-dependent subword unit to a tied HMM state (or an acoustic unit). Fig. 2 illustrates various steps in an HMM-based ASR system where the lexical model is deterministic. In this illustration, each context-dependent subword unit is composed of one lexical unit. However, normally each context-dependent subword unit is represented with three lexical units.

2.4.2. Probabilistic lexical model based ASR approaches

The two conditions, namely, $0 \leq P(a^d|l^i, \theta_l) \leq 1$ and $\sum_{d=1}^D P(a^d|l^i, \theta_l) = 1$ in Eq. (9) characterize an ASR approach where each lexical unit is probabilistically related to all acoustic units. In (Rasipuram and Magimai-Doss, 2015), we showed that there are approaches such as probabilistic classification of HMM states (Luo and Jelinek, 1999), tied posteriors (Rottland and Rigoll, 2000) and Kullback-Leibler divergence based hidden Markov model (KL-HMM) (Aradilla et al., 2008) where the relationship between lexical units and acoustic units is probabilistic. The experimental studies in our previous work indicated that the KL-HMM approach performs better than that of the tied posterior

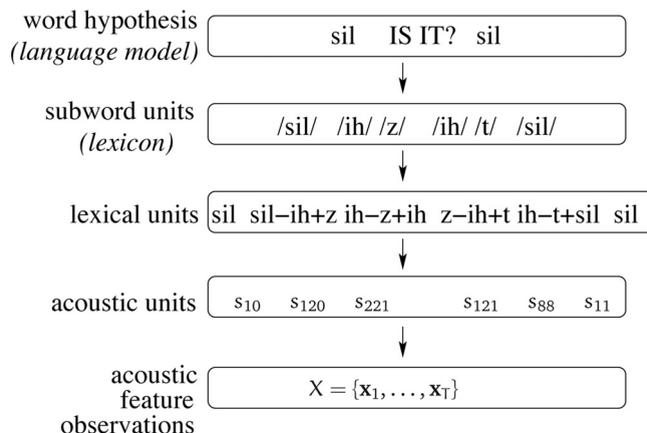


Fig. 2. Sequence of steps taken in a context-dependent HMM-based ASR system where the lexical model is deterministic.

approach on various ASR tasks (Rasipuram and Magimai-Doss, 2015). Therefore, in this paper, we use the KL-HMM approach for probabilistic lexical modeling. The KL-HMM approach is briefly explained in Section 4.1.

2.5. Advantages of probabilistic lexical modeling

Pronunciation variability modeling: Standard HMM-based ASR systems like HMM/GMM and hybrid HMM/ANN deterministically model the relationship between lexical units and acoustic feature observations. As a result of the deterministic relationship, these systems rely on a well developed phoneme lexicon to handle the variability in the acoustic training data. However, when the pronunciations in the lexicon do not reflect the underlying speech data then such a model may poorly represent the training data. For example, this can happen in the case of non-native speakers (when pronunciations normally reflect native speakers) or in the case of spontaneous and conversational speech (when spoken words are pronounced differently from lexicon pronunciations) or in the case of a grapheme lexicon (where pronunciations are based on the orthography of the word). To account for such a variation, typically, phoneme-based ASR systems add pronunciation variants to the lexicon (Strik and Cucchiaroni, 1999). However, the manual addition of pronunciation variants may require explicit human knowledge. In the context of pronunciation variability modeling, it has been shown that the limitation of the standard HMM/GMM system imposed by deterministic mapping can be handled by modeling a probabilistic relationship between lexical and acoustic units (Hain and Woodland, 1999; Saraclar et al., 2000; Hain, 2005; Rasipuram and Magimai-Doss, 2013, 2015).

Resource optimization: In probabilistic lexical model based approaches such as KL-HMM and tied posteriors, the acoustic model and the lexical model are trained one after another and can be trained on an independent set of resources. For example, the acoustic model can be trained on resources from resource-rich languages and domains whereas the lexical model can be trained on a relatively small amount of target language data (Imseng et al., 2012; Rasipuram and Magimai-Doss, 2015).

Flexibility in the choice of acoustic and lexical units: In probabilistic lexical model based ASR approaches, it is not necessary that the subword unit set used for defining the acoustic units should be the same as the subword unit set used for defining the lexical units. The lexical model can capture the relationship between the distinct subword unit sets through acoustics. This flexibility has been exploited to build ASR systems where the acoustic unit set is based on phonemes and the lexical unit set is based on graphemes (Magimai-Doss et al., 2011; Rasipuram and Magimai-Doss, 2013). Furthermore, lexical and acoustic units can model different contextual units. For instance, lexical units can be based on context-dependent subword units while the acoustic units can be based on context-independent subword units (Rottland and Rigoll, 2000; Magimai-Doss et al., 2011; Imseng et al., 2012).

In Section 4.2, we propose a novel approach to integrate articulatory features into HMM-based ASR in the framework of probabilistic lexical modeling that can exploit all the above advantages.

3. Literature survey

There has been a sustained interest in incorporating speech production knowledge into an ASR system for reasons already stated in Section 1. AFs have been incorporated at various levels of an ASR system. Here we provide a brief overview of ASR systems that used multivalued AFs according to the background of the previous section.

3.1. *Lexical units are phonemes and acoustic units are AFs*

In these works, the acoustic units are based on AFs or both AFs and phonemes (Metze and Waibel, 2002; Stüker et al., 2003; Juneja and Espy-Wilson, 2004; Livescu et al., 2008). Each acoustic unit is modeled with a GMM or with discriminative classifiers like ANNs or support vector machines. The lexical model is deterministic, i.e., each phoneme-based lexical unit is deterministically mapped to its AF attributes. The scores from different AF-based acoustic models are combined to arrive at the local emission score $p(x_t|q_t = l^i)$. On continuous speech recognition and cross-lingual adaptation tasks, the use of AF-based acoustic models in combination with phoneme-based acoustic models resulted in a relative reduction in word error rate (WER) of about 5–10% compared to the use of phoneme-based acoustic models alone (Metze and Waibel, 2002; Stüker et al., 2003).

3.2. *Lexical units are AFs and acoustic units are AFs*

These are the systems analogous to standard HMM-based ASR systems where both lexical and acoustic units are either based on context-independent or context-dependent subword units. However, unlike standard ASR systems, the subword units are AFs determined from the AF-based pronunciation lexicon (Deng et al., 1997; Richardson et al., 2003; Kirchhoff, 1996; Wester et al., 2004; Livescu et al., 2008). The AF-based pronunciation lexicon transcribes each word in terms of the positions of the articulators. Each AF is associated with its own hidden state variable. The multiple hidden state variables can follow an independent path to a certain extent and can allow certain amount of asynchrony. In the initial works, hidden state variables of various AFs were required to re-synchronize at phoneme level (Deng et al., 1997; Richardson et al., 2003) or syllable level (Kirchhoff, 1996). In more recent works, the flexible DBN framework allows synchronization to happen at word level or even across word boundaries (Livescu and Glass, 2004; Livescu et al., 2008). When the lexical units are based on context-dependent AFs, the acoustic units are typically clustered context-dependent subword states obtained using decision tree-based state tying methods. These systems have obtained improvements in lexical access experiments (Livescu and Glass, 2004; Livescu et al., 2008; Jyothi et al., 2011).

An approach was proposed by Jyothi et al. (2013) to convert a DBN-based pronunciation model into an equivalent set of factored WFSTs. The utility of this approach was demonstrated using a phoneme-based pronunciation model on isolated word and continuous word speech recognition tasks; and using an AF-based pronunciation model on lexical access tasks. Along the similar lines, an approach was outlined to convert an AF-based DBN pronunciation model into equivalent WFSTs for ASR (Jyothi, 2013, Chapter 8).

3.3. *Lexical units are phonemes and acoustic units are phonemes*

In these systems, similar to standard HMM-based ASR systems, both lexical and acoustic units are based on phonemes. However, AF representations are used as auxiliary information to enhance the performance of the acoustic model (Kirchhoff et al., 2002; Siniscalchi et al., 2012). For example, the acoustic model can be seen as two stage classifier. In the first stage, a set of AF-based ANNs model the relationship between acoustic features and AFs. In the second stage, a phoneme-based ANN models the relationship between all AFs and phonemes. The resulting phoneme-based ANN is used as an acoustic model in hybrid HMM/ANN systems. These systems have achieved a relative reduction in WER of about 5–6% on noise robust ASR tasks and cross-lingual ASR tasks compared to the systems where acoustic-to-phoneme information is directly modeled (Kirchhoff et al., 2002; Siniscalchi et al., 2012).

Alternatively, AF-based neural networks have also been used in tandem speech recognition systems (Cetin et al., 2007, 2007; Livescu et al., 2008; Lal and King, 2013). In the tandem approach, the posterior probabilities of AFs and/or phonemes replace conventional cepstral features in HMM-based ASR systems (Hermansky et al., 2000). In order to model the output of an ANN that is typically non-Gaussian with GMMs, posterior probabilities of the acoustic units are Gaussianized using the log function and then decorrelated using the Karhunen–Loeve transform (KLT).

For tandem systems, the use of AF-based ANNs trained on language-independent data was investigated and compared with the use of AF-based ANNs trained on language-dependent data (Cetin et al., 2007). Cetin et al. (2007) observed that the AF-based ANNs trained on language-independent data reduced the performance (i.e., the WER increased by about 2% relative) of the tandem system compared to the phoneme-based ANNs trained on the same language-independent data. Work by Lal and King (2013) compared the use of AF-based ANNs trained on data from multiple languages (also including the target-language) with AF-based ANNs trained on data from the target-language for tandem systems. It was observed that irrespective of whether the ANNs are trained on data from multiple languages or on the target-language, the AF-based ANNs performed better (relative improvement in WER of 1–9%) than the phoneme-based ANNs. However, the AF-based (or phoneme-based) ANNs trained on data from the target-language performed better than the AF-based (or phoneme-based) ANNs trained on data from multiple languages (Lal and King, 2013).

The difference in conclusion by Cetin et al. (2007) and Lal and King (2013) could be because of the differences in the number of languages used to train the MLPs and their relationship with the target-language. Cetin et al. (2007) used AF-based ANNs trained on English to generate tandem features for Mandarin ASR task; and English and Mandarin belong to different language families. Whereas Lal and King (2013) used ANNs trained on data from multiple languages that also included the target-language.

To summarize, most of the approaches to integrate AFs into an ASR system are based on the deterministic knowledge-based phoneme-to-AF relationship. The approaches summarized in Section 3.1 use the knowledge-based phoneme-to-AF relationship to define the deterministic lexical model parameters. The approaches described in Section 3.2 allow asynchronous evolution of various AFs using an AF-based pronunciation lexicon. However, the AF-based pronunciation lexicon is prepared using the knowledge-based phoneme-to-AF relationship.

In the next section, we present an ASR approach that integrates a model for lexical access using AFs into the HMM-based ASR framework. The approach adapts the knowledge-based phoneme-to-AF relationship using transcribed speech data and incorporates a probabilistic phoneme-to-AF relationship in the model parameters.

4. Proposed approach

In probabilistic lexical model based ASR systems, a lexical unit is probabilistically related to all acoustic units (Section 2.4.2). For each lexical unit l^i , let \mathbf{y}_i be the D -dimensional probability vector or the categorical variable that captures a probabilistic relationship between the lexical unit l^i and D acoustic units, i.e.,

$$\mathbf{y}_i = [y_i^1, \dots, y_i^d, \dots, y_i^D]^T \quad \text{where} \quad y_i^d = P(a^d | l^i) \quad (13)$$

Therefore, the lexical model parameter set $\theta_l = \{\mathbf{y}_i\}_{i=1}^I$. In this paper we use the KL-HMM approach to estimate the parameters of the lexical model (Aradilla et al., 2007, 2008; Aradilla, 2008).

4.1. KL-HMM approach

The KL-HMM approach for lexical model parameter estimation is summarized below:

- The approach assumes that the acoustic unit set \mathcal{A} is known and a trained acoustic model is available. It has been shown that the acoustic units can be modeled using an ANN (Aradilla et al., 2007, 2008; Rasipuram and Magimai-Doss, 2015) or using GMMs (Rasipuram and Magimai-Doss, 2013). In this paper, we use an ANN as the acoustic model.

- Given the acoustic model, the probabilities of acoustic units or the acoustic unit posterior probability vectors for the training data are estimated. At time t , the acoustic unit posterior probability vector \mathbf{z}_t is a D dimensional probability vector.

$$\mathbf{z}_t = [z_t^1, \dots, z_t^d, \dots, z_t^D]^T = [P(a^1|\mathbf{x}_t), \dots, P(a^d|\mathbf{x}_t), \dots, P(a^D|\mathbf{x}_t)]^T \quad (14)$$

- The acoustic unit probability vector sequences are used along with the pronunciation lexicon and word level transcriptions to train the parameters of the probabilistic lexical model. More precisely, the acoustic unit probability vector sequences are used as feature observations to train an HMM where the states represent the lexical units. Each state l^i is parameterized by a categorical variable \mathbf{y}_i that captures a probabilistic relationship between a lexical unit l^i and D acoustic units.
- As both feature observations and state distributions are probability vectors, the local score at each HMM state can be computed as the KL-divergence between the feature observation \mathbf{z}_t and the categorical variable \mathbf{y}_i . KL-divergence being an asymmetric measure, there are the following three possible ways to estimate the KL-divergence:
 1. KL-divergence (S_{KL}): In this case, the state distribution \mathbf{y}_i is the reference distribution.

$$S_{KL}(\mathbf{y}_i, \mathbf{z}_t) = \sum_{d=1}^D y_i^d \log \left(\frac{y_i^d}{z_t^d} \right) \quad (15)$$

2. Reverse KL-divergence (S_{RKL}): In this case, the acoustic unit probability vector \mathbf{z}_t is the reference distribution.

$$S_{RKL}(\mathbf{y}_i, \mathbf{z}_t) = \sum_{d=1}^D z_t^d \log \left(\frac{z_t^d}{y_i^d} \right) \quad (16)$$

3. Symmetric KL-divergence (S_{SKL}): The local score S_{SKL} is the average of the local scores S_{KL} and S_{RKL} .

$$S_{SKL}(\mathbf{y}_i, \mathbf{z}_t) = \frac{1}{2} \cdot [S_{KL}(\mathbf{y}_i, \mathbf{z}_t) + S_{RKL}(\mathbf{z}_t, \mathbf{y}_i)] \quad (17)$$

- The parameters $\{\mathbf{y}_i\}_{i=1}^I$ are trained using the Viterbi training algorithm optimizing a function based on one of the KL-divergence based local scores.
- Decoding is performed using the standard Viterbi decoder where the log-likelihood based score is replaced with the KL-divergence based local score.

The details of the parameter estimation are elaborated in [Appendix A](#). The details on the role of different local scores in estimating the lexical model parameters can be found in the work by [Rasipuram and Magimai-Doss \(2013\)](#) whereas the role of local scores during decoding can be found in the work by [Rasipuram and Magimai-Doss \(2015\)](#).

4.2. Integrating AFs using KL-HMM

In the proposed approach, the relationship between lexical units and acoustic features is factored into two parts through the use of AFs as latent variables or acoustic units:

1. *The acoustic model* where the relationship between AFs and acoustic features is modeled.
2. *The lexical model* where a probabilistic relationship between AFs and lexical units is modeled.

The proposed approach exploits the advantage of probabilistic lexical modeling that the subword unit set used for defining the acoustic units need not be the same as the subword unit set used for defining the lexical units (Section 2.5). The lexical units can be based on context-independent or context-dependent subword units. The proposed approach for AF integration can be summarized in the following steps:

- The acoustic unit set consists of AFs such as manner and place of articulation. Therefore, the acoustic unit set can be seen as a superset of the individual AF sets, i.e.,

$$\mathcal{A} = \{\{\mathcal{A}_1\}, \dots, \{\mathcal{A}_F\}\} \quad (18)$$

where $\{\mathcal{A}_1\}, \dots, \{\mathcal{A}_F\}$ denote the individual AF sets and F , the total number of AFs. For example, the set $\{\mathcal{A}_1\}$ may consist of all the classes specifying the manner of articulation such as vowel, stop, fricative, and so on; the set $\{\mathcal{A}_2\}$ may consist of all the classes specifying the place of articulation such as alveolar, back, dental, dorsal, front and so on. The total number of acoustic units

$$D = D_1 + \dots + D_F \tag{19}$$

where D_1, \dots, D_F denote the cardinality of the individual AFs.

- Each AF is associated with an acoustic model, in our case, with an ANN.
- Given the AF-based acoustic models, posterior probabilities of AFs are estimated. The posterior probability estimates of various AFs are concatenated to produce a D dimensional acoustic unit probability vector \mathbf{z}_t , i.e.,

$$\mathbf{z}_t = [\mathbf{z}_{t,1}, \dots, \mathbf{z}_{t,F}]^T \quad \text{where} \tag{20}$$

$$\mathbf{z}_{t,1} = [z_{t,1}^1, \dots, z_{t,1}^{D_1}] \text{ similarly } \mathbf{z}_{t,F} = [z_{t,F}^1, \dots, z_{t,F}^{D_F}]$$

- The lexical model parameters $\theta_l = \{\{\mathbf{y}_{i,f}\}_{f=1}^F\}_{i=1}^l$ where $\mathbf{y}_{i,f}$ captures a probabilistic relationship between the lexical unit l^i and the AF classes $\{\mathcal{A}_f\}$. That is,

$$\mathbf{y}_{i,1} = [y_{i,1}^1, \dots, y_{i,1}^{D_1}], \sum_{d=1}^{D_1} y_{i,1}^d = 1$$

$$\mathbf{y}_{i,F} = [y_{i,F}^1, \dots, y_{i,F}^{D_F}], \sum_{d=1}^{D_F} y_{i,F}^d = 1$$

- The parameters of the lexical model are trained using the KL-HMM approach. For a lexical unit l^i , the state distribution \mathbf{y}_i can be seen as a stack of F categorical variables, i.e.,

$$\mathbf{y}_i = [\mathbf{y}_{i,1}, \dots, \mathbf{y}_{i,F}]^T \tag{21}$$

- The local score at each HMM state is the KL-divergence between the feature observation and the state distribution. If the local score is S_{RKL} , then the KL-divergence is computed as:

$$S_{\text{RKL}}(\mathbf{y}_i, \mathbf{z}_t) = \sum_{d=1}^{D_1} z_{t,1}^d \log \left(\frac{z_{t,1}^d}{y_{i,1}^d} \right) + \dots + \sum_{d=1}^{D_F} z_{t,F}^d \log \left(\frac{z_{t,F}^d}{y_{i,F}^d} \right) \tag{22}$$

Fig. 3 illustrates the proposed AF-based ASR approach. The ANNs are trained to classify various AFs. The posterior probabilities of the individual AFs are stacked and used as feature observations to train an HMM. As mentioned in Section 2.2, each context-independent or context-dependent subword unit is modeled with three-HMM states.

4.3. Previous findings

The proposed approach for AF integration was studied for phoneme recognition using the TIMIT corpus (Rasipuram and Magimai-Doss, 2011, 2011, 2011). In our previous work the notion of probabilistic lexical modeling was not introduced and the studies were presented from the perspective of acoustic modeling. However, as shown in our recent work, KL-HMM is a lexical modeling approach (Rasipuram and Magimai-Doss, 2013, 2015). Given the formulation of Section 4.2, the findings from the previous studies are refined and re-summarized below:

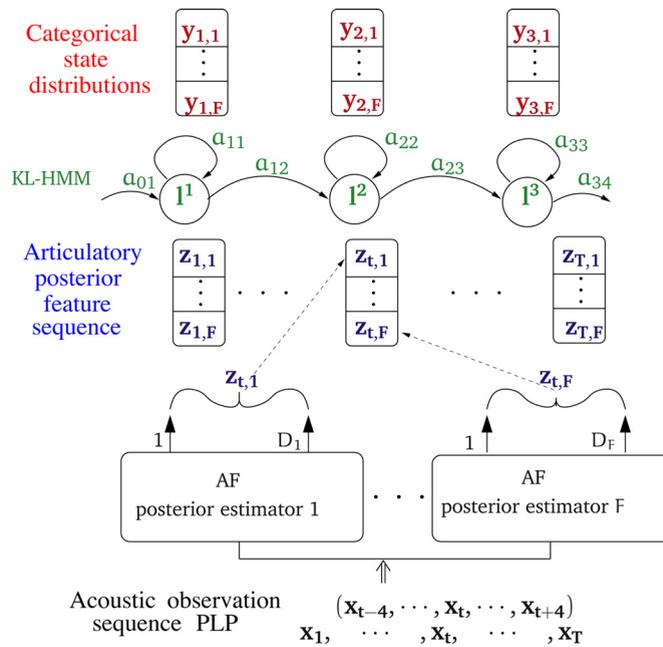


Fig. 3. Illustration of the KL-HMM approach for AF integration.

- It was demonstrated that using AF-based acoustic models in the KL-HMM approach results in better system than using the same acoustic models in the hybrid HMM/ANN approach (Rasipuram and Magimai-Doss, 2011). The study illustrated that it is beneficial if the lexical model is probabilistic.
- Performance of the KL-HMM system using both phoneme-based and AF-based acoustic models was better than the KL-HMM system using only the phoneme-based acoustic model. The study indicates that the KL-HMM approach has the potential to reduce the error rates by incorporating articulatory knowledge into an ASR system.
- The phoneme recognition performance of a system can be improved by improving the AF acoustic models. It was observed that the AF classification accuracy can be improved by modeling the inter-feature dependencies using multistage MLP classifiers and/or multitask learning. In doing so, the performance gap between systems using phonemes as acoustic units and AFs as acoustic units was greatly reduced.

4.4. Other contributions of the paper

The main contribution of the paper is the formulation of AF-based speech recognition in the framework of probabilistic lexical modeling. In addition to that and the previous findings summarized in Section 4.3, the contributions of the paper are as follows:

- The proposed approach for AF integration is evaluated on a continuous speech recognition task (see Section 5 for the experimental setup). In the evaluation, acoustic models are ANNs estimating various multivalued AFs. The lexical units are based on context-dependent phonemes and the lexical model is probabilistic. The proposed approach is compared with the tandem approach for AF integration (Section 6.1). The tandem systems also use the same ANNs as the KL-HMM systems, but as feature extractors.
- Exploiting the resource optimization advantage of probabilistic lexical modeling (see Section 2.5), we study the case where AF-based acoustic models are trained on domain-independent resources and the lexical model parameters are trained on domain-dependent resources (see Section 6.2). Furthermore, the use of multistage MLP classifiers is studied for a continuous speech recognition task (Section 6.3).
- In the framework of probabilistic lexical modeling it is possible to build grapheme-based ASR systems where the lexical units are based on graphemes and the acoustic units are based on phonemes (Magimai-Doss et al., 2011; Rasipuram and Magimai-Doss, 2013, 2015). These grapheme-based ASR systems, where the lexical model

parameters capture a probabilistic graphemes-to-phoneme relationship, performed better than the grapheme-based ASR approaches where the lexical model is deterministic. Motivated by this, in this paper we hypothesize that it is possible to build grapheme-based ASR systems where the lexical units are based on graphemes and the acoustic units are AFs. In this case, the lexical model parameters capture a probabilistic grapheme-to-AF relationship. The resulting grapheme-based ASR approach, in addition to exploiting the advantages of AFs, can also address resource constrained speech recognition with limited or no pronunciation lexicon.

- The lexical model parameters capture a probabilistic relationship between phonemes and AFs. We analyze the parameters of the lexical model to understand if the captured phoneme-to-AF relationship associates well with the knowledge-based phoneme-to-AF relationship (see Section 7).

5. Experimental studies

In this paper, we evaluate the proposed approach for AF integration on a speaker-independent continuous speech recognition task using the DARPA Resource Management (RM) corpus.

5.1. Database

The RM corpus consists of read queries on the status of Naval resources (Price et al., 1988). The task is artificial in aspects such as speech type, range of vocabulary and grammatical constraint. The training set consists of 2880 utterances and the development set 1110 utterances. The training and development sets together consist of 3990 utterances spoken by 109 speakers corresponding to approximately 3.8 h of speech data.

There are four test sets provided by DARPA, namely, feb89, oct89, feb91, and sep92. Each test set contains 300 utterances spoken by 10 speakers. The test set used in this work is obtained by combining the four test sets and thus contains 1200 utterances amounting to 1.1 h in total. The test set is completely covered by a word pair grammar included in the task specification.

The phoneme-lexicon consists of 991 words. In this paper, we followed the RM setup used by Dines and Magimai-Doss (2007) for the pronunciation lexicon to have fairer comparison with the previous work. The UNISYN² lexicon with general American accent was used. There are 45 context-independent phonemes including silence. For ASR experiments, the 45 phonemes were converted to a set of 42 phonemes by merging some allophones (/em/, /en/ and /el/) to their broader phoneme class (/m/, /n/ and /l/). About 35 words in the phoneme lexicon have more than one pronunciation.

The grapheme-lexicon for the RM task is transcribed using 79 graphemes. The first grapheme and the last grapheme of a word are treated as separate graphemes. Therefore, the grapheme set consists of 26 English graphemes ({[A],[B],...[Z]}), 26 English graphemes occurring at the begin of word ({[b_A],[b_B],...,[b_Z]}), 26 English graphemes occurring at the end of word ({[e_A],[e_B],...,[e_Z]}) and silence.

5.2. MLPs

KL-HMM systems use ANNs, more specifically, multilayer perceptrons (MLPs) as the acoustic models whereas tandem systems use the same MLPs as feature extractors. The input to all the MLPs is 39-dimensional perceptual linear prediction (PLP) features with a nine frame temporal context (i.e., four frame preceding and four frame following context). We use three-layer (one input, one hidden and one output layer) MLPs that are trained with the frame level cross entropy error criteria using the Quicknet software.³ The number of hidden units of the MLPs are selected based on the optimal frame accuracy on the development set.

The target labels for the MLPs with phonemes as output units were obtained from the HMM/GMM system. The target labels for the MLPs with AFs as output units are obtained from the phoneme-to-AF map given in Table B.1. The AFs consist of manner, place, height, and vowel. The phoneme-to-AF map is adapted (to distinguish all phonemes of the RM task) from the mapping defined by Hosom (2009). The place class is expanded by adding features like

² <http://www.cstr.ed.ac.uk/projects/unisyn/>.

³ <http://www.icsi.berkeley.edu/Speech/qn.html>.

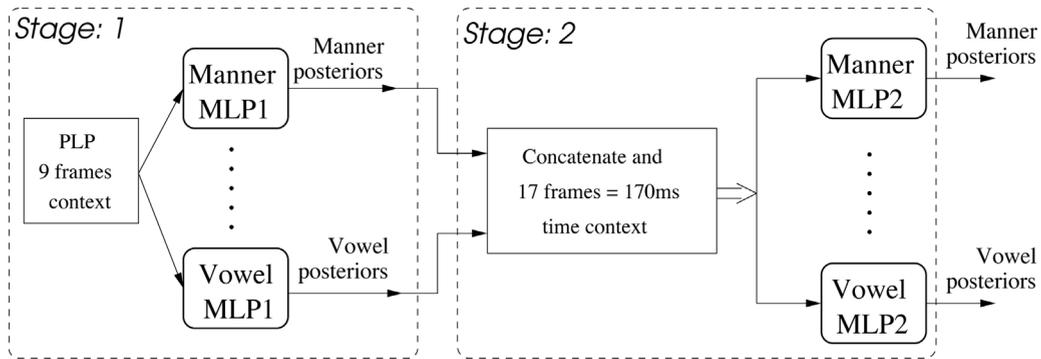


Fig. 4. Multistage AF-based MLP classifiers.

mid-front and mid-back. The height class is expanded by adding features like mid, mid-low, mid-high. Also, the vowel AF is added.

5.2.1. Domain-dependent MLPs

We use the following MLPs trained on the RM corpus:

1. *MLP-RM-PH*: An *off-the-shelf* MLP (Dines and Magimai-Doss, 2007) trained on the RM corpus to classify 45 context-independent phonemes.
2. *MLP-RM-AF*: Set of MLPs trained on the RM corpus to classify AFs.
3. *MLP-RM-PH+AF*: The phoneme and the AF MLPs together are referred as *MLP-RM-PH+AF*.

5.2.2. Domain-independent MLPs

Exploiting the resource optimization advantage of probabilistic lexical modeling (Section 2.5) we study the case where the acoustic model is trained on domain-independent data whereas the lexical model is trained on target-domain data. This also allows us to study the data-invariance aspect of AFs. Similarly, a tandem system can use the MLP trained on domain-independent data for feature extraction. In this paper, we use the Wall Street Journal (WSJ) corpus (Paul and Baker, 1992; Woodland et al., 1994) as domain-independent data to train the MLPs whereas the RM corpus is used as domain-dependent data for which we are interested to build an ASR system. The WSJ corpus consists of two parts – WSJ0 with 14 hours of speech (7193 utterances from 84 speakers) and WSJ1 with 66 hours of speech (29,322 utterances from 200 speakers). In this paper, we use only the WSJ1 corpus as the domain-independent data. Among 1000 words present in the RM task, WSJ corpus consists of only 568 words. The phoneme-lexicon for the WSJ corpus was also obtained from the UNISYN lexicon. Therefore, the phoneme sets and the AF sets for the RM and WSJ corpora are identical. In this paper, we use the following MLPs trained on the WSJ corpus:

1. *MLP-WSJ-PH*: An *off-the-shelf* MLP (Aradilla et al., 2008) trained on the WSJ corpus to classify 45 context-independent phonemes.
2. *MLP-WSJ-AF*: Set of MLPs trained on the WSJ corpus to classify AFs.
3. *MLP-WSJ-PH+AF*: The phoneme and articulatory MLPs together are referred as *MLP-WSJ-PH+AF*.

5.2.3. Multistage MLPs

In this paper, we use the multistage MLP classifier illustrated in Fig. 4. In the first stage, a set of parallel MLPs are used to estimate posteriors of the AFs. Each MLP receives PLP features as input and is trained to classify a specific AF. In the second stage, to model the temporal contextual information and inter-feature dependencies of AFs, a new set of MLPs are trained using articulatory posteriors estimated by the first stage of MLPs (along with other AFs) and a longer temporal context (eight frame preceding and eight frame following context) as input. In this paper, we use the following multistage MLPs:

Table 1

Overview of different systems. In the table *ph*, *af* and *af+ph* denote that acoustic units are phonemes, AFs, and both AFs and phonemes, respectively. cCD denotes clustered context-dependent subword states. In the tandem approach, the ANNs trained to classify the acoustic units (*ph*, *af* or *af+ph*) are used to extract features for the HMM/GMM system. This is indicated through (ANN) notation. *Det* and *Prob* denote that the lexical model is deterministic and probabilistic, respectively.

System	Acoustic units	Acoustic model	Lexical units	Lexical model
KL-HMM	{ph} or {af} or {af+ph}	ANN	CD	Prob
Tandem	cCD	(ANN)+ GMM	CD	Det
HMM/GMM	cCD	GMM	CD	Det

1. *MULTI-WSJRM-PH*: The phoneme-based multistage MLP, where the first stage MLP is trained on the WSJ corpus to classify context-independent phonemes and the second stage MLP is trained on the RM corpus to classify context-independent phonemes.
2. *MULTI-WSJRM-AF*: Set of AF-based multistage MLPs as shown in Fig. 4, where the first set of AF-based MLPs are trained on the WSJ corpus and the second set of AF-based MLPs are trained on the RM corpus.
3. *MULTI-WSJRM-PH+AF*: The phoneme and articulatory multistage MLPs together are referred as *MULTI-WSJRM-PH+AF*.

The input, output and hidden layer sizes of all the MLPs used in the paper along with the total number of parameters and frame accuracies are given in Table B.2 of Appendix B.

5.3. Systems

We compare the KL-HMM, tandem and HMM/GMM systems trained on both the training and development sets (3990 utterances) of the RM corpus. All the systems model crossword context-dependent subword units (either phonemes or graphemes) and each subword unit is modeled as a 3-state HMM. Table 1 summarizes the different systems and their capabilities.

KL-HMM systems: The acoustic units in the KL-HMM system can be context-independent phonemes, or AFs or both context-independent phonemes and AFs. KL-HMM systems use an MLP as the acoustic model. The lexical units are based on context-dependent phonemes and the lexical model is probabilistic. The lexical model is trained using the local score (S_{SKL} or S_{RKL} or S_{KL}) that results in minimum KL-divergence on the training data compared to other local scores. It was observed that for phoneme-based KL-HMM systems, the local score that resulted in minimum KL-divergence was S_{SKL} whereas for grapheme-based KL-HMM systems it was either S_{SKL} or S_{RKL} .

Tandem systems: The tandem systems use the MLPs as feature extractors. The MLPs used for feature extraction are the same MLPs that are used as the acoustic model in the KL-HMM systems. The output of the MLPs is Gaussianized with log transformation followed by KLT. The dimensionality of the features is reduced by retaining the feature components that contribute to 99% of the variance. The resulting features are used to train an HMM/GMM system where the acoustic units are the clustered context-dependent subword states and the lexical units are based on the context-dependent subword units. The lexical and acoustic units are deterministically related. Each acoustic unit is modeled with an eight component Gaussian mixture model as it resulted in an optimal ASR performance.

HMM/GMM systems: The 39-dimensional PLP features used to train the MLP are also used to train the HMM/GMM systems where the acoustic units are clustered context-dependent subword states and the lexical units are based on context-dependent subword units. The lexical and acoustic units are deterministically related. Each acoustic unit is modeled with an eight mixture Gaussian. The state tying resulted in 1611 tied states for the phoneme-based system and 1536 tied states for the grapheme-based system.

The details of various systems, such as the number of lexical units, acoustic units, dimensionality of tandem features, tied states in tandem systems, total number of parameters in different systems are given in Table B.3.

Table 2

Performance in terms of word accuracy on the test set of the RM corpus for the crossword context-dependent phoneme-based and the grapheme-based KL-HMM and tandem systems. The MLPs *MLP-RM-PH*, *MLP-RM-AF* and *MLP-RM-PH+AF* are trained on the RM corpus. The MLPs *MLP-WSJ-PH*, *MLP-WSJ-AF* and *MLP-WSJ-PH+AF* are trained on the domain-independent WSJ corpus. The multistage MLPs *MULTI-WSJRM-PH*, *MULTI-WSJRM-AF* and *MULTI-WSJRM-PH+AF* are trained on both RM and WSJ corpora. Performance on the test set of the RM corpus for the standard crossword context-dependent phoneme and grapheme subword based HMM/GMM systems is 95.9% and 94.8% word accuracy, respectively.

MLPs		Phoneme-based systems		Grapheme-based systems	
		KL-HMM	Tandem	KL-HMM	Tandem
Domain-dependent	<i>MLP-RM-PH</i>	95.6	95.7	95.1	94.6
	<i>MLP-RM-AF</i>	94.9	95.3	94.4	94.1
	<i>MLP-RM-PH+AF</i>	96.2	95.3	95.8	94.3
Domain-independent	<i>MLP-WSJ-PH</i>	95.9	95.8	95.5	94.5
	<i>MLP-WSJ-AF</i>	95.5	95.4	94.9	94.8
	<i>MLP-WSJ-PH+AF</i>	96.4	94.9	96.0	94.5
Multistage	<i>MULTI-WSJRM-PH</i>	96.1	96.0	95.8	95.6
	<i>MULTI-WSJRM-AF</i>	95.4	95.2	95.8	95.2
	<i>MULTI-WSJRM-PH+AF</i>	96.6	95.1	96.4	95.0

6. Results

The performance of the phoneme-based and grapheme-based KL-HMM and tandem systems using the various MLPs in terms of word accuracy on the test set of the RM corpus is given in Table 2. Performance on the test set of the RM corpus for the standard crossword context-dependent phoneme-based HMM/GMM system is 95.9% word accuracy and the crossword context-dependent grapheme-based HMM/GMM system is 94.8% word accuracy. The performances of the baseline phoneme-based HMM/GMM systems of this paper and the phoneme-based HMM/GMM systems in the literature reported by Hain and Woodland (1999) and Povey et al. (2011) on the RM corpus are the same despite the difference in the phoneme-lexicon used. In this work, we used a phoneme-lexicon based on the UNISYN lexicon whereas Hain and Woodland (1999) used a phoneme-lexicon based on the CMUDict.

In the remainder of the section, the results of the phoneme and grapheme-based KL-HMM and tandem systems are analysed according to the various MLPs used (as discussed in Section 5.2).

6.1. Domain-dependent MLPs

For both phoneme and grapheme subword units, the results indicate that:

- The KL-HMM systems using the phoneme-based MLP (*MLP-RM-PH*) significantly ⁴ outperform the KL-HMM systems using the AF-based MLPs (*MLP-RM-AF*). The KL-HMM systems using both phoneme-based and AF-based MLPs (*MLP-RM-PH+AF*) significantly outperform the KL-HMM systems using either phoneme-based or AF-based MLPs.
- The difference in performance between the KL-HMM and tandem systems is significant when both phoneme-based and AF-based MLPs (*MLP-RM-PH+AF*) are used. However, the difference in performance between the KL-HMM and tandem systems employing either the phoneme-based (*MLP-RM-PH*) or the AF-based (*MLP-RM-AF*) MLPs is not significant.
- The difference in performance between the KL-HMM and tandem systems using the phoneme-based MLP (*MLP-RM-PH*), and the baseline HMM/GMM system is not statistically significant.

In this section, the AF-based MLPs are trained on a limited amount of domain-dependent data. In the next section, we verify that AF estimation using MLPs is data-invariant and can benefit from a larger domain-independent data set.

⁴ In the paper, “significant difference” implies that the difference in performance between the compared systems is statistically significant with confidence greater than 95.0%. Statistical significance tests for the systems are performed using the approach proposed by Bisani and Ney (2004).

6.2. Domain-independent MLPs

In this case, the KL-HMM and tandem systems use the MLPs trained on the domain-independent WSJ corpus. The results show that:

- The phoneme-based KL-HMM system using the AF-based MLPs (*MLP-WSJ-AF*) trained on a large set of domain-independent data (the WSJ corpus) significantly outperforms the phoneme-based KL-HMM system using the AF-based MLPs (*MLP-RM-AF*) trained on a small set of domain-dependent data (the RM corpus). For all the other phoneme-based KL-HMM systems and all the grapheme-based KL-HMM systems, there is no significant difference in performance between the systems using domain-dependent MLPs and the systems using domain-independent MLPs.
- The grapheme-based tandem system using domain-dependent AF-based MLPs (*MLP-WSJ-AF*) performs significantly better than the grapheme-based tandem system using domain-independent AF-based MLPs (*MLP-RM-AF*). However, for all the other grapheme-based and all the phoneme-based tandem systems, there is no significant difference in performance between the systems using domain-dependent MLPs and the systems using domain-independent MLPs.
- Unlike the domain-dependent MLPs case, the difference in performance between the phoneme-based KL-HMM systems using the phoneme-based MLP (*MLP-WSJ-PH*) and the AF-based MLPs (*MLP-WSJ-AF*) is not statistically significant.
- Similar to the domain-dependent MLPs case, the phoneme-based and grapheme-based KL-HMM systems using both the phoneme-based and AF-based MLPs (*MLP-WSJ-PH+AF*) significantly outperform all other KL-HMM systems (employing other domain-independent MLPs) and all the tandem systems.

6.3. Multistage MLPs

In this case, the KL-HMM and tandem systems use the multistage MLPs trained on both WSJ and RM corpora. The results show that:

- The grapheme-based KL-HMM systems using either the AF-based multistage MLPs or both AF and phoneme-based multistage MLPs (i.e., *MULTI-WSJRM-AF* and *MULTI-WSJRM-PH+AF*) significantly outperform the grapheme-based KL-HMM systems using the respective single-stage MLPs trained on the WSJ corpus (*MLP-WSJ-AF* and *MLP-WSJ-PH+AF*). However, for all the phoneme-based KL-HMM systems and all other grapheme-based KL-HMM systems, there is no significant difference in performance between the systems using the multistage MLPs and the systems using domain-independent MLPs.
- Unlike the domain-dependent and domain-independent MLPs case, the difference in performance between the grapheme-based KL-HMM system using the phoneme-based multistage MLP (*MULTI-WSJRM-PH*) and the AF-based multistage MLPs (*MULTI-WSJRM-AF*) is not statistically significant.
- The performance gap between grapheme-based and phoneme-based KL-HMM systems is greatly reduced when multistage MLPs are used.

6.4. Summary of the experimental results

To summarize, the following conclusions can be drawn from the experimental study:

1. The proposed approach for AF integration resulted in an ASR system that performs similar to the tandem approach for AF integration. Though both approaches perform similarly, there are two main advantages of the proposed approach.

Firstly, in the proposed approach, the articulatory representations are kept intact in the model parameters in the form of probabilistic phoneme-to-AF or grapheme-to-AF relationship learned from the transcribed speech data. Furthermore, as we will see in the next section, the approach also adapts the knowledge-based phoneme-to-AF and grapheme-to-AF relationship on the transcribed speech data, and allows different AFs to evolve asynchronously.

However, the tandem approach tends to lose the two primary benefits of articulatory representation, i.e., finer granularity and asynchronous evolution.

Secondly, the KL-HMM system achieves similar performance as the tandem system but uses fewer parameters than the tandem approach (about 30–40% relative for the single-stage MLPs trained on the RM corpus, and about 10% relative for the single-stage MLPs trained on the WSJ corpus and multistage MLPs as observed in Table B.3).

2. The performance of the grapheme- or phoneme-based KL-HMM systems using both AFs and phonemes as acoustic units is always better (about 10–12% relative reduction in WER) than the KL-HMM system using either of them as acoustic units. We speculate the following reasons:
 - (a) When both AFs and phonemes are used as acoustic units, not only the acoustic model is improved but also the lexical model, as both probabilistic phoneme-to-phoneme and phoneme-to-AF (or grapheme-to-phoneme and grapheme-to-AF) relationships are modeled.
 - (b) The AF-based and phoneme-based MLPs are trained independently. However, the probabilistic phoneme-to-phoneme and phoneme-to-AF (or grapheme-to-phoneme and grapheme-to-AF) relationships are learned together during lexical model training. Therefore, the approach can learn the inter-feature dependencies among various AFs and phonemes or graphemes.
 - (c) The information captured by the lexical model with phonemes as acoustic units and with AFs as acoustic units is complementary.

The proposed approach significantly outperforms the tandem approach if both AFs and phonemes are used as acoustic units. The performance of the grapheme- or phoneme-based tandem systems using both AFs and phonemes as acoustic units is always worse than the tandem system using either of them as acoustic units. This could be because the statistical characteristics of phoneme and articulatory feature probabilities are different. It may be necessary to model phonemes and AFs separately such as in the factored observation model (Cetin et al., 2007).

3. The results indicate that AFs estimation using MLPs is data and domain invariant. Further, the use of AF-based MLPs trained on a larger domain-independent data set helps both phoneme-based and grapheme-based KL-HMM systems.
4. In our previous phoneme recognition studies it was observed that the multistage articulatory MLPs improved the phoneme recognition accuracy (Rasipuram and Magimai.-Doss, 2011). In this paper, multistage AF-based MLPs did not improve the performance of the phoneme-based KL-HMM systems but improved the performance of the grapheme-based KL-HMM systems. We conjecture the following two reasons for the observed trends.

Firstly, the multistage AF-based MLP was motivated from the work by Pinto et al. (2011). In that work it was shown that the second MLP in the multistage MLP classifier learns the phonetic temporal patterns (i.e., the phonetic confusions at the output of the first MLP) and the phonotactics of the language observed in the training data. In our case, the second set of MLPs in the multistage AF-based MLP classifier could model phonotactic constraints at the articulatory feature level.

- When the lexical units are based on context-dependent phonemes, the lexical model incorporates phonotactic constraints. The results indicate that it may be redundant to model the phonotactic constraints twice, once at the acoustic model level and again at the lexical model level.
- In the case of context-dependent grapheme-based KL-HMM systems, the lexical model is modeling graphemic constraints and the acoustic model is modeling phonotactic constraints, which could be complementary to each other especially given the fact that the grapheme-to-phoneme relationship in English is irregular.

Secondly, though the frame accuracy and phoneme recognition accuracy are considered as important factors for word recognition, they may not be the only indicators of word level performance (Greenberg et al., 2000). The relationship between frame accuracy and word accuracy depends on more than one factor: the pronunciation lexicon, the acoustic model, the lexical model and the language model components of an ASR system. In other words, lexical constraints and syntactic constraints incorporated while decoding can handle the shortcomings of the acoustic model or may render some of the gains obtained through the acoustic model redundant.

7. Analysis

In the proposed approach, with phonemes or graphemes as lexical units and AFs as acoustic units, the parameters of the lexical model capture a probabilistic phoneme-to-AF or grapheme-to-AF relationships. In this section, we analyze

Table 3

The phoneme-to-manner of articulation and phoneme-to-place of articulation relationship captured by the lexical model parameters of context-independent phonemes.

Synchronous Examples				
Phoneme	AF	state1	state2	state3
aa	Manner	vow		
	Place	bck		
d	Manner	vst		
	Place	alv		
l	Manner	app		
	Place	lat		
m	Manner	nas		
	Place	lab		
Asynchronous Examples				
Phoneme	AF	state1	state2	state3
aw	Manner	vow		
	Place	midf	bck	midb
ay	Manner	vow		
	Place	bck	midf	
ow	Manner	vow		
	Place	bck		midb
oy	Manner	vow		
	Place	bck	lat	midf
ah	Manner	vow		
	Place	bck	mid	
uh	Manner	vow		
	Place	bck	midb	
b	Manner	sil	vst	vow
	Place	lab		
g	Manner	vst		vow
	Place	dor		
ch	Manner	stp		frc
	Place	alv	fnt	
jh	Manner	vst		vow
	Place	alv	fnt	
er	Manner	app		vow
	Place	ret		
ng	Manner	nas		
	Place	dor		alv

the parameters of the lexical model at subword level (Section 7.1) and word level (Section 7.2) to understand the following:

- Is the phoneme-to-AF or grapheme-to-AF relationship captured by the lexical model parameters close to the knowledge-based phoneme-to-AF or grapheme-to-AF relationship?
- Does the model allow different AFs to evolve asynchronously?

The target labels for the MLPs with AFs as output units are obtained from the knowledge-based phoneme-to-AF map. Hence, during training, the MLPs do not account for the AFs changing asynchronously. However, as shown by King and Taylor (2000), it is typical that at the output of the MLP, different AFs change at different times (especially at the phoneme boundaries) and exhibit asynchronous behaviour. Since the KL-HMM system models the output of a set of AF-based MLPs, we hypothesize that the lexical model parameters can capture asynchronous AF behaviour especially at the HMM states modeling subword unit boundaries (i.e., the first and third HMM states of a subword unit).

Table 4
Examples of the context-dependent phonemes that exhibit synchronous and asynchronous manner and place of articulation at the HMM state level.

Synchronous Examples					Asynchronous Examples				
Phoneme	AF	st1	st2	st3	Phoneme	AF	st1	st2	st3
aa-n+t	Manner	nas			w-ey+r	Manner	vow		
	Place	alv				Place	bck	midf	ret
ah-z+sh	Manner	vow	frc		sil-jh+ae	Manner	sil	vst	vow
	Place	midf	frt			Place	sil	fnt	

Table 5
For various KL-HMM systems, the percentage of context-dependent phonemes where the changes in manner and place of articulations between the three HMM states are synchronous and asynchronous.

MLP	Synchronous	Asynchronous
MLP-RM-AF	59.49%	40.50%
MLP-WSJ-AF	61.56%	38.40%
MULTI-WSJRM-AF	67.50%	32.49%

7.1. Subword level analysis

Table 3 shows the manner and place of articulation for context-independent phonemes in three HMM states captured by the lexical model parameters. The analysis is presented only on manner and place of articulation for the sake of simplicity. However, similar trends are observed even for height of articulation. The denoted AF values correspond to the dimension with maximum probability captured by the lexical model parameters of context-independent subword units. The first and second parts of the table presents examples of context-independent phonemes where the manner and place of articulation between three HMM states are synchronous and asynchronous, respectively. It can be observed that the phoneme-to-AF relationship of Table 3 relates well with the knowledge-based relationship given in Table B.1.

Table 3 indicates that for diphthongs such as /aw/, /ay/, /ow/, and /oy/ and for vowels such as /ah/ and /uh/ the captured place of articulation changed between the HMM states whereas the captured manner of articulation is the same, i.e., “vowel”. For phonemes that are voiced-stops, i.e., /b/, and /g/, the captured place of articulation is the same across the three HMM states whereas the captured manner of articulation changed between the HMM states. More specifically, the initial states of /b/ and /g/ captured a “voiced-stop” and the third state captured a “vowel”. For phonemes /ch/ and /jh/, the captured manner of articulation changed at the second HMM state whereas the captured place of articulation changed at the first HMM state.

We have computed the percentage of context-dependent phonemes where the lexical model parameters exhibited asynchrony between manner and place of articulation at the HMM state level. A context-dependent phoneme model is said to be synchronous if manner and place of articulation change at the same state transition or remain the same across three HMM states. A context-dependent phoneme model is said to be asynchronous if manner and place of articulation change at different HMM states. For example, in Table 4, the context-dependent phonemes /aa-n+t/ and /ah-z+sh/ exhibit synchronous changes in manner and place of articulation whereas the phonemes /w-ey+r/ and /sil-jh+ae/ exhibit asynchronous changes at the HMM state level.

In Table 5, the first column indicates the set of MLPs used to train the KL-HMM system, and the second and the third columns indicate the percentage of context-dependent phonemes where the changes in manner and place of articulations at the HMM state level are synchronous and asynchronous, respectively. It is important to note that the classification of context-dependent phoneme models in terms of synchronous and asynchronous does not take into account the errors in the phoneme-to-AF map captured by the lexical model parameters. For example, in Table 7, the place and height of articulation for the grapheme model [W] are asynchronous. However, the captured place of articulation in the first state of [W] as “lateral” could be considered as an error. Therefore, the percentage of context-dependent phoneme models exhibiting synchronous or asynchronous behaviour are only an indicative of the asynchronous nature of the context-dependent phoneme models.

Table 5 indicates that asynchronous articulatory movements among manner and place of articulation are relatively lower when multistage AF-based MLP classifiers are used. We argue the following two reasons for this. Firstly, as

Table 6

The phoneme-to-AF relationship captured by the lexical model parameters of context-dependent phonemes for the word “BELOW”. Pronunciation of the word BELOW in the phoneme lexicon is /b/ /ih/ /l/ /ow/. AF values for manner, place and height of articulation are given.

Phoneme	/b/			/ih/			/l/			/ow/		
	st1	st2	st3	st1	st2	st3	st1	st2	st3	st1	st2	st3
Manner	sil	vst	vow			app			vow			
Place	sil	lab		midf		lat			bck		midb	
Height	max		high			vhi			mid	vhi		

Table 7

The grapheme-to-AF relationship captured by the lexical model parameters of context-dependent graphemes for the word “BELOW”. Pronunciation of the word BELOW in the grapheme lexicon is [B] [E] [L] [O] [W]. AF values for manner, place and height of articulation are given.

Grapheme	[B]			[E]			[L]			[O]			[W]		
	st1	st2	st3	st1	st2	st3	st1	st2	st3	st1	st2	st3	st1	st2	st3
Manner	vst		vow			app			vow						
Place	sil	lab	midf		lat			bck		<i>lat</i>	midb	sil			
Height	max		high			vhi			mid	<i>midl</i>	vhi		sil		

discussed in Section 6.4, the second stage of AF-based MLPs in the multistage MLP classifiers model the phonotactics of the language, therefore various AFs may be more synchronous. Secondly, the frame accuracies of the multistage AF-based MLP classifiers are better than the frame accuracies of other MLPs. As a result the number of context-dependent phonemes exhibiting asynchronous changes because of the errors in the captured relationship could be relatively less.

The analysis of the parameters indicated that the model is able to capture asynchronous AF configurations at the subword unit level. In the next section we will see that asynchronous articulatory configurations are more meaningful at the word level as the context-dependent subword models also capture information of the neighbouring phonemes.

7.2. Word level analysis

The phoneme-to-AF and grapheme-to-AF relationships captured by the lexical model parameters of phoneme-based and grapheme-based KL-HMM systems for the word “BELOW” are given in Tables 6 and 7, respectively. The tables indicate the manner, place and height of articulation.

It can be observed that the phoneme-to-AF relationship of Table 6 relates well with the knowledge-based relationship given in Table B.1. Table 7 shows that even if subword units are graphemes, articulatory patterns similar to the system using phoneme subword units are captured. The two differences between grapheme-to-AF and phoneme-to-AF are indicated in red italic font in Table 7. The number of subword units in the pronunciation of the word “BELOW” are five in the case of graphemes and four in the case of phonemes. It can be observed from the table that this irregularity in the grapheme pronunciation has been accounted for, as the sequence of graphemes [O] and [W] together capture the information of phoneme /ow/.

Furthermore, the tables also indicate that various AFs evolve asynchronously. For example, in Table 6, the captured manner of articulation in the second state of /b/ is “voiced-stop” and in the third state of /b/ is “vowel”, whereas the place of articulation in both second and third states of /b/ is “labial”.

It was observed by Magimai-Doss et al. (2011) and Rasipuram (2014)[Chapter 4] that the lexical model parameters of the context-dependent graphemes tend to model the transition information to the next context-dependent grapheme in the sequence. Similar observations can also be made from Tables 6 and 7, but at the finer articulatory feature level. For example, in Table 6, the captured manner of articulation in the third state of /b/ is “vowel” which corresponds to the next phoneme in the sequence, i.e., /ih/. Similarly in Table 7, the captured manner, place and height of articulation in the third state of grapheme [B] correspond to the next grapheme in the sequence, i.e., [E]. The analysis shows that the lexical model parameters of the context-dependent phonemes and graphemes are capable of capturing some information about preceding and following articulatory configurations.

8. Discussion and conclusion

In this paper, we proposed an approach to integrate articulatory feature representations into HMM-based ASR in the framework of probabilistic lexical modeling. The proposed approach involves two stages: acoustic model and lexical model. The acoustic model is a posterior probability estimator that models the relationship between acoustic feature observations and AFs. The lexical model, models a probabilistic relationship between the lexical units and the AFs. The approach has the following potential advantages:

- **Lexical access:** As opposed to knowledge-based approaches, the parameters of the lexical model in the proposed approach are learned using transcribed speech data by training an HMM whose states represent lexical units and the parameters of the state capture a probabilistic relationship between a lexical unit and AFs. Consequently, the approach integrates a model for lexical access using AFs into the HMM-based ASR framework.
- **Asynchrony of AFs:** As observed in Section 7, the model also allows different AFs to evolve asynchronously. Thus, overcoming some of the limitations of knowledge-based approaches.
- **Combination of various AFs:** A challenge often faced in using articulatory features for ASR is the combination of evidences from different AFs. In that regard, the proposed approach can be seen as a multi-channel approach where each AF serves as a separate channel and various AFs are combined at the local score computation level. Also, as seen in this paper, the multi-channel approach can be trivially extended to combine other relevant information such as the phoneme information.

Our investigations on a continuous speech recognition task have shown that the proposed approach effectively integrates AFs into the HMM-based ASR framework; improves ASR performance if combined with phoneme-based acoustic models; exploits domain-independent resources; and offers flexibility to use either phonemes or graphemes as subword units.

The probabilistic grapheme-to-AF relationship captured in the lexical model parameters of the KL-HMM system with acoustic units as AFs and lexical units based on context-dependent graphemes can be exploited to generate an AF-based pronunciation lexica using the acoustic data-driven grapheme-to-phoneme conversion approach proposed by [Rasipuram and Magimai.-Doss \(2012\)](#). The AF-based pronunciation lexica can be used in DBN-based approaches for AF integration that require such lexica ([Livescu and Glass, 2004](#); [Livescu et al., 2008](#)) or in AF-based text-to-speech synthesis systems.

In this paper, we focussed mainly on the integration of AFs into an ASR system. More precisely, we focussed on the lexical model aspect of the proposed approach. The three-layer or multistage MLPs classifying context-independent AFs (or phonemes) were used as acoustic models. The approach can be potentially improved by improving the acoustic model along the following directions:

1. **Context-dependent AFs:** The output of MLPs or the acoustic units could be context-dependent AFs that take into account the neighbouring articulatory context.
2. **Deep architectures for AF estimation:** More recently, ANNs with deep architectures have gained lot of attention ([Dahl et al., 2012](#); [Hinton et al., 2012](#)). In similar a vein, the articulatory feature model can be based on deep ANN architectures ([Siniscalchi et al., 2012](#)).

In this paper, we have shown that the AF-based acoustic models can be trained on domain-independent data whereas the lexical model can be trained on domain-dependent data. In our recent work we found that in the framework of probabilistic lexical modeling, the acoustic model can be trained on language-independent resources and the lexical model on a relatively small amount of language-dependent data ([Rasipuram and Magimai.-Doss, 2015](#)). AFs are considered to be more language-independent and effective for cross-linguistic adaptation ([Lal and King, 2013](#); [Siniscalchi et al., 2012](#)). Therefore, we hypothesize that the use of articulatory feature based language-independent acoustic model

in the proposed approach can offer potential advantages in building ASR systems for under-resourced and minority languages.⁵ Our future work will focus on extending the proposed approach along this direction.

Acknowledgments

This work was partly supported by the Swiss NSF through the grants “Flexible Grapheme-Based Automatic Speech Recognition (FlexASR, grant numbers 124985 and 146229)” and partly by the Commission for Technology and Innovation (CTI) on “Automatic scoring and adaptive pedagogy for oral language learning (ScoreL2: CTI project 15990.2 PFES-ES)”.

Appendix A. Parameter estimation in the KL-HMM approach

Given a trained ANN and a training set of N utterances $\{X(n), W(n)\}_{n=1}^N$, the set of acoustic unit probability vectors $\{Z(n), W(n)\}_{n=1}^N$ are estimated. For each training utterance n , $X(n)$ represents the sequence of cepstral features, $W(n)$ represents the sequence of underlying words, $Z(n)$ represents a sequence of acoustic unit probability vectors and $T(n)$ represents the number of cepstral features or the number of acoustic unit probability vectors.

The KL-HMM system is parameterized by $\Theta_{kull} = \{\{y_i\}_{i=1}^I, \{a_{ij}\}_{i,j=1}^I\}$. The lexical model parameters $\{y_i\}_{i=1}^I$ are initialized uniformly, i.e., initially $y_i^d = \frac{1}{D} \forall i, d$. The training data $\{Z(n), W(n)\}_{n=1}^N$ and the current parameter set Θ_{kull} , are used to estimate the new set of parameters $\hat{\Theta}_{kull}$ by the Viterbi algorithm. In the case of the local score S_{RKL} the cost function minimized is,

$$\hat{\Theta}_{kull} = \arg \min_{\Theta_{kull}} \left[\sum_{n=1}^N \min_{Q \in \mathcal{Q}} \sum_{t=1}^{T(n)} [S_{RKL}(y_{q_t}, z_t(n)) - \log a_{q_{t-1}q_t}] \right] \quad (23)$$

where $\mathcal{Q} = \{q_1, \dots, q_t, \dots, q_{T(n)}\}$, $q_t \in \{1, \dots, I\}$ and \mathcal{Q} denotes set of all possible HMM state sequences.

The training process involves iteration over the segmentation and the optimization steps until convergence. Given the current set of parameters, the segmentation step yields an optimal state sequence for each training utterance using the Viterbi algorithm. Given optimal state sequences and acoustic unit posterior vectors belonging to the states, the optimization step then estimates new set of model parameters by minimizing Eq. (23) subject to the constraint that $\sum_{d=1}^D y_i^d = 1$.

The optimal state distribution for the local score S_{RKL} (Eq. (16)), is the arithmetic mean of the training acoustic unit probability vectors assigned to the state, i.e.,

$$y_i^d = \frac{1}{M(i)} \sum_{z_t(n) \in Z(i)} z_t^d(n) \quad \forall d \quad (24)$$

where $Z(i)$ denotes the set of acoustic unit probability vectors assigned to state i and $M(i)$ is the cardinality of $Z(i)$. More details about the parameter estimation for the KL-HMM approach can be found in the thesis by Aradilla (2008).

Appendix B. Details of the experimental study

B.1. Phoneme-to-articulatory feature map

See Table B.1.

⁵ A language that lacks one or more resources required to build an ASR system is referred to as under-resourced language (Besacier et al., 2014). Minority languages are languages spoken by a minority of population. A minority language need not be under-resourced and an under-resourced language may or may not be a minority language.

Table B.1
Knowledge-based phoneme-to-articulatory feature map used in this paper.

Phoneme	Manner	Place	Height	Vowel
sil	sil	sil	sil	sil
aa	vowel	back	low	aa
ae	vowel	mid-front	low	ae
ah	vowel	mid	mid	ah
ao	vowel	back	mid-low	ao
aw1	vowel	mid-front	low	aw1
aw2	vowel	mid-back	high	aw2
ax	vowel	mid	mid	ax
axr	approximant	retroflex	mid	consonant
ay1	vowel	back	low	ay1
ay2	vowel	mid-front	high	ay2
b	voiced-stop	labial	max	consonant
ch	stop	front	max	consonant
d	voiced-stop	alveolar	max	consonant
dh	voiced-fricative	dental	max	consonant
eh	vowel	mid-front	mid	eh
el	approximant	lateral	very-high	consonant
em	nasal	labial	max	consonant
en	nasal	alveolar	max	consonant
er	vowel	mid	mid	er
ey1	vowel	front	mid-high	ey1
ey2	vowel	mid-front	high	ey2
f	fricative	labial	max	consonant
g	voiced-stop	dorsal	max	consonant
hh	aspirated	unknown	max	consonant
ih	vowel	mid-front	high	ih
iy	vowel	front	very-high	iy
jh	voiced-stop	front	max	consonant
k	stop	dorsal	max	consonant
l	approximant	lateral	very-high	consonant
m	nasal	labial	max	consonant
n	nasal	alveolar	max	consonant
ng	nasal	dorsal	max	consonant
ow1	vowel	back	mid	ow1
ow2	vowel	mid-back	high	ow2
oy1	vowel	back	mid-low	oy1
oy2	vowel	mid-front	high	oy2
p	stop	labial	max	consonant
r	approximant	retroflex	mid-low	consonant
s	fricative	alveolar	max	consonant
sh	fricative	front	max	consonant
t	stop	alveolar	max	consonant
th	fricative	dental	max	consonant
uh	vowel	mid-back	high	uh
uw	vowel	back	very-high	uw
v	voiced-fricative	labial	max	consonant
w	approximant	back	very-high	consonant
y	approximant	front	very-high	consonant
z	voiced-fricative	alveolar	max	consonant
zh	voiced-fricative	front	max	consonant

B.2. Details of the MLPs

The number of hidden units of the MLPs are selected based on the optimal frame accuracy on the development set of the RM or WSJ corpora. Hence, the total number of parameters are different for various MLPs. For the AF-based MLPs trained on the RM corpus, optimal frame accuracy was obtained when the total number of MLP parameters

Table B.2

Overview of the various MLPs used in the paper. The number of input, output and hidden units of MLPs are denoted as i , h and o . The total number of parameters of an MLP is computed as $h * (i + o + 1)$. The frame accuracies of the MLPs *MLP-RM-PH*, *MLP-RM-AF*, *MULTI-WSJRM-PH* and *MULTI-WSJRM-AF* are computed on the development set of the RM corpus. The frame accuracies of the MLPs *MLP-WSJ-PH*, *MLP-WSJ-AF* are computed on the development set of the WSJ corpus.

MLP	Acoustic units	Input size (i)	Hidden size (h)	Output size (o)	MLP size	Frame acc.
<i>MLP-RM-PH</i>	phonemes	$39 * 9 = 351$	1260	45	0.5M	73.77
<i>MLP-WSJ-PH</i>	phonemes	$39 * 9 = 351$	3652	45	1.5M	69.34
<i>MULTI-WSJRM-PH</i>	phonemes	$17 * 45 = 765$	730	45	2.1M	80.16
	manner	$39 * 9 = 351$	682	10	0.2M	82.48
<i>MLP-RM-AF</i>	place	$39 * 9 = 351$	676	13	0.2M	76.20
	height	$39 * 9 = 351$	685	8	0.2M	79.26
	vowel	$39 * 9 = 351$	658	23	0.2M	79.38
<i>MLP-WSJ-AF</i>	manner	$39 * 9 = 351$	4005	10	1.4M	89.18
	place	$39 * 9 = 351$	3972	13	1.4M	86.83
	height	$39 * 9 = 351$	4027	8	1.4M	88.10
	vowel	$39 * 9 = 351$	3866	23	1.4M	88.29
<i>MULTI-WSJRM-AF</i>	manner	$17 * 54 = 918$	1063	10	6.6M	86.47
	place	$17 * 54 = 918$	1059	13	6.6M	82.42
	height	$17 * 54 = 918$	1065	8	6.6M	84.25
	vowel	$17 * 54 = 918$	1048	23	6.6M	84.36

were about 15% of the number of training frames. For the phoneme-based and AF-based MLPs trained on the WSJ corpus, it was observed that the optimal frame accuracy on the development set was obtained when the total number of parameters were about 5% of the number of training frames (28M). Therefore, the total number of parameters of the MLPs trained on the WSJ corpus is about $28M * 0.05 \approx 1.5M$. The total number of parameters in the multistage MLPs consist of the parameters of the first stage MLP(s) and the second stage MLP (Table B.2).

B.3. Details of various systems

See Table B.3.

Table B.3

The details of the various KL-HMM systems, namely the total number of MLP parameters (N_{θ_a}), lexical units (I), acoustic units (D), lexical model parameters (N_{θ_l}), KL-HMM system parameters ($N_{\theta_a} + N_{\theta_l}$). The details of the various tandem systems, namely the number tied states (ts), Gaussian components in each tied state (ts), dimensionality of the tandem features that contribute to 99% of the variance (v), tandem system parameters excluding the MLP parameters ($N_{\theta_{tan}}$) and tandem system parameters with the MLP parameters ($N_{\theta_a} + N_{\theta_{tan}}$). The total number of parameters in the cross-word context-dependent phoneme-based and grapheme-based HMM/GMM systems are about 1.0M.

MLP	N_{θ_a}	KL-HMM system				Tandem system			
		#lexical Units (I)	#acoustic Units (D)	N_{θ_l} $I * D$	Total $N_{\theta_a} + N_{\theta_l}$	#tied states (ts), #Gaussians (g)	dim (v)	Tandem ($N_{\theta_{tan}}$) $ts * (g * (2v + 1))$	Total $N_{\theta_a} + N_{\theta_{tan}}$
<i>MLP-RM-PH</i>	0.5M	10,988	45	0.5M	1.0M	1726, 8	39	1.1M	1.6M
<i>MLP-RM-AF</i>	0.8M	8127	54	0.4M	1.2M	1737, 8	42	1.2M	2.0M
<i>MLP-RM-PH+AF</i>	1.3M	11678	99	1.1M	2.4M	1723, 8	75	2.1M	3.4M
<i>MLP-WSJ-PH</i>	1.5M	11,723	45	0.5M	2.0M	2367, 8	37	1.4M	2.9M
<i>MLP-WSJ-AF</i>	5.6M	8059	54	0.4M	6.0M	1751, 8	42	1.2M	6.8M
<i>MLP-WSJ-PH+AF</i>	7.1M	9753	99	1.0M	8.1M	1718, 8	75	2.1M	9.2M
<i>MULTI-WSJRM-PH</i>	2.1M	10,857	45	0.5M	2.6M	2546, 8	36	1.5M	3.6M
<i>MULTI-WSJRM-AF</i>	9.6M	7753	54	0.4M	10.0M	1807, 8	40	1.2M	10.8M
<i>MULTI-WSJRM-PH+AF</i>	11.7M	9975	99	1.0M	12.7M	1784, 8	73	2.1M	13.8M

References

- Aradilla, G., 2008. *Acoustic Models for Posterior Features in Speech Recognition*. EPFL, Switzerland (Ph.D. thesis).
- Aradilla, G., Boulard, H., Magimai-Doss, M., 2008. Using KL-based acoustic models in a large vocabulary recognition task. In: *Proceedings of Interspeech*, pp. 928–931.
- Aradilla, G., Vepa, J., Boulard, H., 2007. An acoustic model based on Kullback–Leibler divergence for posterior features. In: *Proceedings of ICASSP*, pp. 657–660.
- Bahl, L., Brown, P., de Souza, P., Picheny, M., 1988. Acoustic Markov models used in the Tangora speech recognition system. In: *Proceedings of ICASSP*, pp. 497–500 vol.1.
- Besacier, L., Barnard, E., Karpov, A., Schultz, T., 2014. Automatic speech recognition for under-resourced languages: a survey. *Speech Commun.* 56, 85–100.
- Bisani, M., Ney, H., 2004. Bootstrap estimates for confidence intervals in ASR performance evaluation. In: *Proceedings of ICASSP*, pp. I-409–12.
- Cetin, O., Kantor, A., King, S., Bartels, C., Magimai-Doss, M., Frankel, J., Livescu, K., 2007. An articulatory feature-based tandem approach and factored observation modeling. In: *Proceedings of ICASSP*, pp. 645–648.
- Cetin, O., Magimai-Doss, Livescu, K., Kantor, A., King, S., Bartels, C., Frankel, J., 2007. Monolingual and crosslingual comparison of tandem features derived from articulatory and phone MLPs. In: *Proceedings of ASRU*, pp. 36–41.
- Chang, S., 2002. *A Syllable, Articulatory-Feature, and Stress-Accent Model of Speech Recognition*. University of California, Berkeley (Ph.D. thesis).
- Chomsky, N., Halle, M., 1968. *The Sound Pattern of English*. MIT Press.
- Dahl, G., Yu, D., Deng, L., Acero, A., 2012. Context-dependent pre-trained deep neural networks for large-vocabulary speech recognition. *IEEE Trans. Audio Speech Lang. Process.* 20, 30–42.
- Deng, L., Ramsay, G., Sun, D., 1997. Production models as a structural basis for automatic speech recognition. *Speech Commun.* 22, 93–111.
- Dines, J., Magimai-Doss, M., 2007. A study of phoneme and grapheme based context-dependent ASR systems. In: *Proceedings of Machine Learning for Multimodal Interaction (MLMI)*, pp. 215–226.
- Frankel, J., King, S., 2005. A hybrid ANN/DBN approach to articulatory feature recognition. In: *Proceedings of Interspeech*, pp. 3045–3048.
- Frankel, J., Wester, M., King, S., 2007. Articulatory feature recognition using dynamic Bayesian networks. In: *Computer Speech & Language*, pp. 620–640.
- Greenberg, S., Chang, S., Hollenback, J., 2000. An introduction to the diagnostic evaluation of Switchboard-corpus automatic speech recognition systems. In: *Proceedings of the NIST Speech Transcription Workshop*.
- Hain, T., 2005. Implicit modelling of pronunciation variation in automatic speech recognition. *Speech Commun.* 46, 171–188.
- Hain, T., Woodland, P.C., 1999. Dynamic HMM selection for continuous speech recognition. In: *Proceedings of EUROSPEECH*, pp. 1327–1330.
- Harris, J., 1994. *English Sound Structure*. Blackwell.
- Hermansky, H., Ellis, D., Sharma, S., 2000. Tandem connectionist feature extraction for conventional HMM systems. In: *Proceedings of ICASSP*, pp. 1635–1638.
- Hinton, G., Deng, L., Yu, D., Dahl, G., Mohamed, A.R., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath, T., Kingsbury, B., 2012. Deep neural networks for acoustic modeling in speech recognition: the shared views of four research groups. *IEEE Signal Process. Mag.* 29, 82–97.
- Hiroya, S., Honda, M., 2004. Estimation of articulatory movements from speech acoustics using an HMM-based speech production model. *IEEE Trans. Speech Audio Process.* 12, 175–185.
- Hogden, J., Lofqvist, A., Gracco, V., Zlokarnik, I., Rubin, P., Saltzman, E., 1996. Accurate recovery of articulator positions from acoustics: new conclusions based on human data. *J. Acoust. Soc. Am.* 100, 1819–1834.
- Holter, T., Svendsen, T., 1997. Incorporating linguistic knowledge and automatic baseform generation in acoustic subword unit based speech recognition. In: *Proceedings of EUROSPEECH*.
- Hosom, J., 2009. Speaker-independent phoneme alignment using transition-dependent states. *Speech Commun.* 51, 352–368.
- Huang, X., Jack, M., 1989. Semi-continuous hidden Markov models for speech signal. *Comput. Speech Lang.* 3, 239–251.
- Huttenlocher, D., Zue, V., 1984. A model of lexical access from partial phonetic information. In: *Proceedings of ICASSP*, pp. 391–394.
- Hwang, M.Y., Huang, X., 1992. Subphonetic modeling with Markov states – Senone. In: *Proceedings of ICASSP*, pp. 33–36 vol. 1.
- Imseing, D., Dines, J., Motlicek, P., Garner, P., Boulard, H., 2012. Comparing different acoustic modeling techniques for multilingual boosting. In: *Proceedings of Interspeech*.
- Juneja, A., Espy-Wilson, C., 2004. Significance of invariant acoustic cues in a probabilistic framework for landmark-based speech recognition. In: *From Sound to Sense: Fifty+ Years of Discoveries in Speech Communication*. MIT Press, Cambridge, MA.
- Jyothi, P., 2013. *Discriminative and Articulatory Feature-Based Pronunciation Models for Conversational Speech Recognition*. The Ohio State University, Ohio (Ph.D. thesis).
- Jyothi, P., Fosler-Lussier, E., Livescu, K., 2013. Discriminative training of WFST factors with application to pronunciation modeling. In: *Proceedings of Interspeech*, pp. 1961–1965.
- Jyothi, P., Livescu, K., Fosler-Lussier, E., 2011. Lexical access experiments with context-dependent articulatory feature-based models. In: *Proceedings of ICASSP*, pp. 4900–4903.
- King, S., Frankel, J., Livescu, K., McDermott, E., Richmond, K., Wester, M., 2007. Speech production knowledge in automatic speech recognition. *J. Acoust. Soc. Am.* 121, 723–742.
- King, S., Taylor, P., 2000. Detection of phonological features in continuous speech using neural networks. *Comput. Speech Lang.* 14, 333–353.
- Kirchhoff, K., 1996. Syllable-level desynchronisation of phonetic features for speech recognition. In: *Proceedings of the International Conference on Spoken Language Processing*.

- 26 *R. Rasipuram, M. Magimai.-Doss / Computer Speech and Language xxx (2015) xxx–xxx*
- Kirchhoff, K., Fink, G.A., Sagerer, G., 2002. Combining acoustic and articulatory feature information for robust speech recognition. *Speech Commun.* 37, 303–319.
- Ladefoged, P., 1993. *A Course in Phonetics*. Harcourt Brace College Publishers.
- Lal, P., King, S., 2013. Cross-lingual automatic speech recognition using tandem features. *IEEE Trans. Audio Speech Lang. Process.* 21, 2506–2515.
- Livescu, K., Cetin, O., Hasegawa-Johnson, M., King, S., Bartels, C., Borges, N., Kantor, A., Lal, P., Yung, L., Bezman, A., Dawson-Haggerty, S., Woods, B., 2008. *Articulatory Feature-Based Methods for Acoustic and Audio-Visual Speech Recognition: 2006 JHU Summer Workshop Final Report*. Technical Report. John Hopkins University, Center for Language and Speech Processing.
- Livescu, K., Glass, J.R., 2004. Feature-based pronunciation modeling with trainable asynchrony probabilities. In: *Proceedings of the International Conference on Spoken Language Processing*.
- Luo, X., Jelinek, F., 1999. Probabilistic classification of HMM states for large vocabulary continuous speech recognition. In: *Proceedings of ICASSP*, pp. 353–356.
- Magimai-Doss, M., Rasipuram, R., Aradilla, G., Boulard, H., 2011. Grapheme-based automatic speech recognition using KL-HMM. In: *Proceedings of Interspeech*, pp. 2693–2696.
- Metze, F., Waibel, A., 2002. A flexible stream architecture for ASR using articulatory features. In: *Proceedings of the International Conference on Spoken Language Processing*.
- Morgan, N., Boulard, H., 1995. Continuous speech recognition: An introduction to the hybrid HMM/connectionist approach. *IEEE Signal Process. Mag.*, 25–42.
- Næss, A.B., Livescu, K., Prabhavalkar, R., 2011. Articulatory feature classification using nearest neighbors. In: *Proceedings of Interspeech*, pp. 2301–2304.
- Paul, D., Baker, J., 1992. The design for the wall street journal-based CSR corpus. In: *DARPA Speech and Language Workshop*. Morgan Kaufmann Publishers.
- Pinto, J., Sivaram, G.S.V.S., Magimai-Doss, M., Hermansky, H., Boulard, H., 2011. Analysis of MLP based hierarchical phoneme posterior probability estimator. *IEEE Trans. Audio Speech Lang. Process.* 19, 225–241.
- Povey, D., Burget, L., Agarwal, M., Akyazi, P., Kai, F., Ghoshal, A., Glembek, O., Goel, N., Karafiát, M., Rastrow, A., Rose, R., Schwarz, P., Thomas, S., 2011. The subspace Gaussian mixture model – a structured model for speech recognition. *Comput. Speech Lang.*
- Povey, D., Ghoshal, A., Boulianne, G., Burget, L., Glembek, O., Goel, N., Hannemann, M., Motlicek, P., Qian, Y., Schwarz, P., Silovsky, J., Stemmer, G., Vesely, K., 2011. The Kaldi speech recognition toolkit. In: *Proceedings of ASRU*.
- Prabhavalkar, R., Fosler-Lussier, E., Livescu, K., 2011. A factored conditional random field model for articulatory feature forced transcription. In: *Proceedings of ASRU*, pp. 77–82.
- Price, P.J., Fisher, W., Bernstein, J., 1988. The DARPA 1000-word resource management database for continuous speech recognition. In: *Proceedings of ICASSP*, pp. 651–654.
- Rabiner, L., 1989. A tutorial on hidden Markov models and selected applications in speech recognition. In: *Proceedings of IEEE77*, pp. 257–286.
- Rasipuram, R., 2014. *Grapheme-based Automatic Speech Recognition using Probabilistic Lexical Modeling*. EPFL, Switzerland (Ph.D. thesis).
- Rasipuram, R., Magimai-Doss, M., 2011. Improving articulatory feature and phoneme recognition using multitask learning. In: *Artificial Neural Networks and Machine Learning – ICANN 2011*.
- Rasipuram, R., Magimai-Doss, M., 2011. Integrating articulatory features using Kullback–Leibler divergence based acoustic model for phoneme recognition. In: *Proceedings of ICASSP*.
- Rasipuram, R., Magimai-Doss, M., 2011. Multitask Learning to Improve Articulatory Feature Estimation and Phoneme Recognition. *Idiap Research Report, Idiap-RR-21-2011*.
- Rasipuram, R., Magimai-Doss, M., 2012. Acoustic data-driven grapheme-to-phoneme conversion using KL-HMM. In: *Proceedings of ICASSP*.
- Rasipuram, R., Magimai-Doss, M., 2013. Improving grapheme-based ASR by probabilistic lexical modeling approach. In: *Proceedings of Interspeech*.
- Rasipuram, R., Magimai-Doss, M., 2013. Probabilistic Lexical Modeling and Grapheme-based Automatic Speech Recognition. *Idiap Research Report*. http://publications.idiap.ch/downloads/reports/2013/Rasipuram_Idiap-RR-15-2013.pdf
- Rasipuram, R., Magimai-Doss, M., 2015. Acoustic and lexical resource constrained ASR using language-independent acoustic model and language-dependent probabilistic lexical model. *Speech Commun.* 68, 23–40, URL: <http://www.sciencedirect.com/science/article/pii/S0167639314000995>
- Richardson, M., Bilmes, J., Diorio, C., 2003. Hidden-articulatory Markov models for speech recognition. *Speech Commun.* 41, 713–716.
- Rottland, J., Rigoll, G., 2000. Tied posteriors: an approach for effective introduction of context dependency in hybrid NN/HMM LVCSR. In: *Proceedings of ICASSP*, pp. 1241–1244.
- Saraclar, M., Nock, H., Khudanpur, S., 2000. Pronunciation modeling by sharing Gaussian densities across phonetic models. In: *Comput. Speech Lang.* pp. 137–160.
- Scharenborg, O., Wan, V., Moore, R.K., 2007. Towards capturing fine phonetic variation in speech using articulatory features. *Speech Commun.* 49, 811–826.
- Siniscalchi, S., Lyu, D.C., Svendsen, T., Lee, C.H., 2012. Experiments on cross-language attribute detection and phone recognition with minimal target-specific training data. *IEEE Trans. Audio Speech Lang. Process.* 20, 875–887.
- Strik, H., Cucchiari, C., 1999. Modeling pronunciation variation for ASR: a survey of the literature. *Speech Commun.* 29, 225–246.
- Stüker, S., Metze, F., Schultz, T., Waibel, A., 2003. Integrating multilingual articulatory features into speech recognition. In: *Proceedings of Eurospeech*, pp. 1033–1036.
- Suzuki, S., Okadome, T., Honda, M., 1998. Determination of articulatory positions from speech acoustics by applying dynamic articulatory constraints. In: *Proceedings of the International Conference on Spoken Language Processing*.

- Wester, M., Frankel, J., King, S., 2004. [Asynchronous articulatory feature recognition using dynamic Bayesian networks](#). In: [Proceedings of IEICI Beyond HMM Workshop](#).
- Woodland, P.C., Odell, J.J., Valtchev, V., Young, S.J., 1994. [Large vocabulary continuous speech recognition using HTK](#). In: [Proceedings of ICASSP](#), pp. 125–128.
- Young, S.J., Odell, J.J., Woodland, P.C., 1994. [Tree-based state tying for high accuracy acoustic modelling](#). In: [Proceedings of the Workshop on Human Language Technology \(HLT\)](#), pp. 307–312.