



TOWARDS WEAKLY SUPERVISED ACOUSTIC
SUBWORD UNIT DISCOVERY AND LEXICON
DEVELOPMENT USING HIDDEN MARKOV
MODELS

Marzieh Razavi Ramya Rasipuram
Mathew Magimai.-Doss

Idiap-RR-15-2017

APRIL 2017

Towards Weakly Supervised Acoustic Subword Unit Discovery and Lexicon Development Using Hidden Markov Models

Marzieh Razavi^{a,b,*}, Ramya Rasipuram^c, Mathew Magimai.-Doss^a

^a*Idiap Research Institute, CH-1920 Martigny, Switzerland*

^b*Ecole Polytechnique Fédérale de Lausanne (EPFL), CH-1015 Lausanne, Switzerland*

^c*Apple Inc., Cupertino, CA, USA*

Abstract

State-of-the-art automatic speech recognition and text-to-speech systems are based on subword units, typically phonemes. This necessitates a lexicon that maps each word to a sequence of subword units. Development of a phonetic lexicon for a language requires linguistic knowledge as well as human effort, which may not be always readily available, particularly for under-resourced languages. In such scenarios, an alternative approach is to use a lexicon based on units such as, graphemes or subword units automatically derived from the acoustic data. This article focuses on automatic subword unit based lexicon development using methods that are employed for development of grapheme-based systems. Specifically, we present a novel hidden Markov model (HMM) based formalism for automatic derivation of subword units and pronunciation generation using only transcribed speech data. In this approach, the subword units are derived from the clustered context-dependent units in a grapheme based system using the maximum-likelihood criterion. The subword unit based pronunciations are then generated by learning either a deterministic or a probabilistic relationship between the graphemes and the acoustic subword units (ASWUs). In this article, we first establish the proposed framework on a well resourced language by comparing it against related approaches in the literature and investigating the transferability of the derived subword units to other domains. We then show the scalability of the proposed approach on real under-resourced scenarios by conducting studies on Scottish Gaelic, a genuinely under-resourced language,

*Corresponding author

Email addresses: marzieh.razavi@idiap.ch (Marzieh Razavi),
ramya.murali@gmail.com (Ramya Rasipuram), mathew@idiap.ch (Mathew Magimai.-Doss)

and comparing the approach against state-of-the-art grapheme-based ASR approaches. Our experimental studies on English show that the derived subword units can not only lead to better ASR systems compared to graphemes, but can also be transferred across domains. The experimental studies on Scottish Gaelic show that the proposed ASWU-based lexicon development approach scales without any language specific considerations and leads to better ASR systems compared to a grapheme-based lexicon, including the case where ASR system performance is boosted through the use of acoustic models built with multilingual resources from resource-rich languages.

Keywords: automatic subword unit derivation, pronunciation generation, hidden Markov model, Kullback-Leibler divergence based hidden Markov model, under-resourced language, automatic speech recognition

1. Introduction

Speech technologies such as automatic speech recognition (ASR) systems and text-to-speech (TTS) systems typically model subword units as they are 1) more trainable compared to words and, 2) more generalizable towards unseen contexts or words. Subword modeling entails development of a pronunciation lexicon that represents each word as a sequence of subword units. Typically in the literature, the subword units are the phonemes or phones. Phonetic lexicon development requires linguistic expert knowledge about the phone set of the language and the relationship between the written form, i.e., graphemes and phonemes. Therefore, it is a time consuming and tedious task. To reduce the amount of human effort, grapheme-to-phoneme (G2P) conversion approaches have been proposed (Pagel et al., 1998; Sejnowski and Rosenberg, 1987; Taylor, 2005; Bisani and Ney, 2008). The G2P conversion approaches still require an initial phonetic lexicon in the target language to learn the relation between graphemes and phonemes through data-driven approaches. While majority languages such as English and French have well-developed phonetic lexicons, there are many other languages such as Scottish Gaelic and Vietnamese that lack proper phonetic resources.

In the absence of a phonetic lexicon, alternatively grapheme subword units based on the writing system have been explored in the literature (Kanthak and Ney, 2002a; Killer et al., 2003; Dines and Magimai.-Doss, 2007; Magimai-Doss et al., 2011; Ko and Mak, 2014; Rasipuram and Magimai.-Doss, 2015; Gales

et al., 2015). The main advantage of using graphemes as subword units is that they make development of lexicons easy. However, the success of grapheme-based ASR systems depends on the G2P relationship of the language. For languages with a regular or shallow G2P relationship such as Spanish, the performance of grapheme-based and phoneme-based ASR systems is typically comparable, whereas for languages with an irregular or deep G2P relationship such as English, the performance of a grapheme-based ASR system is relatively poor when compared to a phoneme-based system (Kanthak and Ney, 2002a; Killer et al., 2003).

Yet another way to handle lack of phonetic lexicon is to derive subword units automatically from the speech signal and build a lexicon based on that. In the literature, interest in acoustic subword unit (ASWU) based lexicon development emerged from the pronunciation variation modeling perspective, specifically with the idea of overcoming limitation of linguistically motivated subword units, i.e., phones (Lee et al., 1988; Svendsen et al., 1989; Paliwal, 1990; Lee et al., 1988; Bacchiani and Ostendorf, 1998; Holter and Svendsen, 1997). However, recently, there has been a renewed interest from the perspective of handling lexical resource constraints (Singh et al., 2000; Lee et al., 2013; Hartmann et al., 2013). A limitation of most of the existing methods for acoustic subword units based lexicon development is that they are not able to handle unseen words.

In this article, building upon the recent developments in grapheme-based ASR, we propose an approach to derive "phone-like" subword units and develop a pronunciation lexicon given limited amount of transcribed speech data. In this approach, first a set of ASWUs is derived by modeling the relationship between the graphemes and the acoustic speech signal in a hidden Markov model (HMM) framework based on two assumptions,

1. writing systems carry information regarding the spoken system. Alternately, a written text embeds information about how it should be spoken. Though this embedding can be deep or shallow depending on the language; and
2. envelope of short-term spectrum tends to carry information related to phones.

The ASWU-based pronunciation lexicon is then developed by learning the grapheme-to-ASWU (G2ASWU) relationship through the acoustic signal, and inferring pronunciations using G2ASWU conversion (analogous to G2P conversion). The G2ASWU conversion process inherently brings in the capability to

generate pronunciation for unseen words. The viability of the proposed approach has been demonstrated through preliminary studies on English (Razavi and Magimai-Doss, 2015) and Scottish Gaelic (Razavi et al., 2015), where a probabilistic G2ASWU relationship was learned and pronunciation lexicon was developed.

This article builds on the preliminary works to first extend the approach to the case where a deterministic G2ASWU relationship is learned. We then study and contrast the two G2ASWU relationship learning methods and investigate the following aspects:

1. *Domain-independency of the ASWUs*: Subword units such as phones and graphemes are by default domain-independent. This enables using a lexicon based on either of them across different domains. ASWUs are derived from a limited amount of acoustic speech signal from a domain. Furthermore, the limited data can have undesirable variabilities based on the hardware used and the conditions under which the data is collected. Therefore a question arising is whether the derived ASWUs are domain independent. Through a cross-domain study on English, we show that our approach indeed yields ASWUs that are domain independent. Furthermore, the proposed approach inherently enables transferring ASWU based lexicon developed on one domain to another.
2. *Potential of ASWUs in improving multilingual ASR*: It has been shown that both acoustic resource and lexical resource constraints can be effectively addressed by learning a probabilistic relationship between graphemes of the target languages and a multilingual phone set obtained from lexical resources of auxiliary languages using acoustic data (Rasipuram and Magimai-Doss, 2015). Success of such approaches lies on the fact that there exists a systematic relationship between linguistically motivated grapheme units and phonemes. Therefore a question that arises is: Does the ASWU-based lexicon based on the proposed approach hold the advantage over grapheme-based lexicon in such a case? Alternately, do the ASWUs exhibit similar systematic relationship to multilingual phones and can it be exploited to further improve the under-resourced language ASR? Through a study on Scottish Gaelic, a genuinely under-resourced language, we show that there exists a systematic relationship between the ASWUs and multilingual phones, which can not only be exploited to yield systems better than grapheme-based lexicons, but also to gain insight into

the derived units.

It is worth mentioning that, to the best of our knowledge, this is the first work that aims to establish these aspects in the context of ASWU-based lexicon development. Consequently, it paves the path for adopting ASWU-based lexicon development and its use for ASR technology development, especially for under-resourced languages.

The remainder of the article is organized as follows. Section 2 provides a background about the grapheme-based ASR and related approaches in the literature for subword unit derivation and pronunciation generation. Section 3 describes the proposed approach. Section 4 presents investigations on well resourced majority language English and Section 5 presents the investigations on under-resourced minority language Scottish Gaelic. Section 6 provides a brief analysis of the derived ASWUs and the generated pronunciations. Finally, Section 7 concludes the article.

2. Background

This section provides the relevant background for understanding the proposed approach for ASWU based lexicon development. Sections 2.1 and 2.2 first present a background on HMM-based ASR and grapheme-based ASR approaches, which form the basis for our proposed approach for automatic subword unit derivation and pronunciation generation. Section 2.3 then presents a survey on the existing approaches for derivation of ASWUs and lexicon development.

2.1. HMM-based ASR

In statistical automatic speech recognition, given the acoustic observation sequence $X = [\mathbf{x}_1, \dots, \mathbf{x}_t, \dots, \mathbf{x}_T]$ with T denoting the total number of frames, the goal is to find the most probable sequence of words W^* ,

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

$$= \arg \max_{W \in \mathcal{W}} p(W, X|\Theta), \quad (2)$$

where \mathcal{W} denotes the set of hypotheses and Θ denotes the set of parameters. Eqn. (2) is obtained result of applying Bayes' rule and assuming $p(X)$ to be constant w.r.t all word hypotheses. Hereafter for simplicity, we drop Θ from the equations.

HMM-based ASR approach achieves that goal by finding the most probable sequence of states Q^* representing W^* by incorporating lexical and syntactic knowledge:

$$Q^* = \arg \max_{Q \in \mathcal{Q}} p(Q, X), \quad (3)$$

$$= \arg \max_{Q \in \mathcal{Q}} \prod_{t=1}^T p(\mathbf{x}_t | q_t = l^i) \cdot P(q_t = l^i | q_{t-1} = l^j), \quad (4)$$

$$= \arg \max_{Q \in \mathcal{Q}} \sum_{t=1}^T \log(p(\mathbf{x}_t | q_t = l^i)) + \log(P(q_t = l^i | q_{t-1} = l^j)), \quad (5)$$

where \mathcal{Q} denotes all possible state sequences, q_t denotes HMM state at time frame t and $l^i \in \{l^1, \dots, l^I\}$ denotes a subword unit or lexical unit. Eqn. (4) is derived as a consequence of i.i.d and first order Markov model assumptions.

Estimation of $p(\mathbf{x}_t | q_t = l^i)$ is typically factored through latent variables or acoustic units $\{a^d\}_{d=1}^D$ as (Rasipuram and Magimai.-Doss, 2015):

$$p(\mathbf{x}_t | q_t = l^i) = \sum_{d=1}^D p(\mathbf{x}_t, a^d | q_t = l^i), \quad (6)$$

$$= \sum_{d=1}^D p(\mathbf{x}_t | a^d, q_t = l^i) \cdot P(a^d | q_t = l^i), \quad (7)$$

$$= \sum_{d=1}^D p(\mathbf{x}_t | a^d) \cdot P(a^d | q_t = l^i) \text{ (assuming } \mathbf{x}_t \perp\!\!\!\perp q_t | a^d), \quad (8)$$

$$= \mathbf{v}_t^T \mathbf{y}_i, \quad (9)$$

where $\mathbf{v}_t = [v_t^1, \dots, v_t^d, \dots, v_t^D]^T$ with $v_t^d = p(\mathbf{x}_t | a^d)$ and $\mathbf{y}_i = [y_i^1, \dots, y_i^d, \dots, y_i^D]^T$ and $y_i^d = P(a^d | q_t = l^i)$.

As presented above in Eqn. (9), estimation of $p(\mathbf{x}_t | q_t = l^i)$ can be seen as matching acoustic information \mathbf{v}_t with lexical information \mathbf{y}_i . In recent years, it has been shown that the match can also be obtained by matching posterior distributions of a^d conditioned on acoustic features and lexical information. One such approach is Kullback-Leibler divergence based HMM (KL-HMM) (Aradilla et al., 2008), where the local score is estimated as Kullback-Leibler divergence between \mathbf{y}_i and \mathbf{z}_t :

$$KL(\mathbf{y}_i, \mathbf{z}_t) = \sum_{d=1}^D y_i^d \cdot \log\left(\frac{y_i^d}{z_t^d}\right), \quad (10)$$

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

HMM-based ASR approach has been primarily built with the idea of hav-

ing a phonetic lexicon that transcribes each word as a sequence of phones. In conventional HMM-based ASR systems, lexical units $\{l^i\}_{i=1}^I$ model context-dependent phones and acoustic units $\{a^d\}_{d=1}^D$ are clustered context-dependent phone units. \mathbf{v}_t and \mathbf{z}_t are typically estimated using either Gaussian mixture models (GMMs) or artificial neural networks (ANNs); and $\{\mathbf{y}_i\}_{i=1}^I$ is a set of Kronecker delta distributions based on the one-to-one deterministic map between lexical unit l^i and acoustic unit a^d modeled by the state tying decision tree. We refer to this case where l^i and a^d are one-to-one related as *deterministic lexical modeling* framework. In (Rasipuram and Magimai.-Doss, 2015), it has been elucidated that there are HMM-based ASR approaches where the relationship between l^i and a^d is probabilistic. KL-HMM approach, probabilistic classification of HMM states (PCHMM) approach (Luo and Jelinek, 1999) and tied posterior approach (Rottland and Rigoll, 2000) are examples of *probabilistic lexical modeling* framework. In KL-HMM, \mathbf{y}_i is estimated based on \mathbf{z}_t whereas in PC-HMM and tied posterior \mathbf{y}_i is estimated based on \mathbf{v}_t . For a detailed overview on deterministic and probabilistic lexical modeling, the reader is referred to (Rasipuram and Magimai.-Doss, 2015).

2.2. Grapheme-based ASR

In the literature, the issue of lack of well developed phonetic lexicon has been addressed by using graphemes as subword units. Most of the studies in this direction have been conducted in the framework of deterministic lexical modeling, where $\{l^i\}_{i=1}^I$ model context-dependent graphemes, $\{a^d\}_{d=1}^D$ are clustered context-dependent grapheme units and \mathbf{y}_i is a decision tree learned while state tying based on either singleton question set or phonetic question set (Kanthak and Ney, 2002b; Killer et al., 2003).

In the framework of probabilistic lexical modeling, it has been shown that grapheme-based ASR systems can be built with $\{a^d\}_{d=1}^D$ based on phones of auxiliary languages or domains, and $\{l^i\}_{i=1}^I$ based on target language graphemes. More precisely, a phone class conditional probability \mathbf{z}_t estimator is trained with acoustic and lexical resources from auxiliary languages or domains, and \mathbf{y}_i , which captures a probabilistic G2P relationship, is trained on target language or domain acoustic data (Magimai.-Doss et al., 2011; Rasipuram and Magimai.-Doss, 2015). It has been shown that this approach can effectively address both acoustic resource and lexical resource constraints (Rasipuram and Magimai.-Doss, 2015; Rasipuram et al., 2013a). As a natural extension of the approach, an acoustic data-driven grapheme-to-phoneme conversion approach

has been proposed, where the G2P relationship learned in this manner through acoustics is used to infer pronunciations (Rasipuram and Magimai-Doss, 2012; Razavi et al., 2016). We dwell about the acoustic data-driven G2P conversion approach more in the paper later, as it is an integral part of the proposed ASWU based lexicon development approach.

2.3. Literature survey on ASWU derivation and pronunciation generation

The idea of using lexicons based on ASWUs instead of the linguistically motivated units has been appealing to the ASR community for three main reasons: (1) ASWUs tend to be rather data-dependent than linguistic knowledge-dependent, as they are typically obtained through optimization of an objective function using training speech data (Lee et al., 1988; Bacchiani and Ostendorf, 1998), (2) they could possibly help in handling pronunciation variations (Livescu et al., 2012), and (3) they can avoid the need for explicit phonetic knowledge (Lee et al., 2013).

Typically, the ASWU-based lexicon development process, in addition to speech signal, requires the corresponding transcription in terms of words. Alternately, the lexicon development process is weakly-supervised similar to acoustic model development in an ASR system. More recently, in the context of “zero-resourced” ASR system development, there are efforts towards developing methods that are fully unsupervised (Chung et al., 2013; Lee et al., 2015). Such methods are at very early stages and are out of the scope of this paper. In the remainder of this section, we provide a brief literature survey on weakly-supervised ASWU-based lexicon development. ASWU-based lexicon development involves two key challenges: (a) derivation of ASWUs and (b) pronunciation generation based on the derived ASWUs. The approaches proposed in the literature can be grouped into two categories based on how these two challenges are addressed. More precisely, there are approaches that decouple these two challenges and address them separately (Section 2.3.1), and there are approaches that address these two challenges in an unified manner with a common objective function (Section 2.3.2).

2.3.1. Automatic subword unit discovery followed by pronunciation generation approaches

The very first efforts approached the ASWU derivation problem as segmentation of *isolated word* speech signals into acoustic segments and clustering acoustic segments into groups each representing a subword unit (Lee et al., 1988;

Svendsen et al., 1989; Paliwal, 1990). More precisely, as shown in Figure 1, in the segmentation step, the speech utterance $X = [\mathbf{x}_1, \dots, \mathbf{x}_t, \dots, \mathbf{x}_T]$ is partitioned into I consecutive segments (with boundaries $B = \{b_1, \dots, b_i, \dots, b_I\}$) such that the frames in a segment are acoustically similar. Then in the clustering step, the acoustic segments are clustered into groups of subword units.

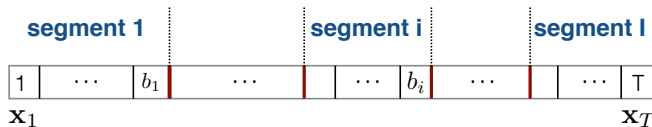


Figure 1: Segmentation of speech utterance \mathbf{x} into I segments.

In (Lee et al., 1988; Svendsen et al., 1989), the segmentation step was approached by applying dynamic programming techniques and finding the segment boundaries b_i such that the likelihood ratio distortion between the speech frames in segment i and the generalized spectral centroid of segment i (i.e., the centroid LPC vector) is minimized. The obtained acoustic segments were then clustered using the K-means algorithm in which each acoustic segment was represented by its centroid. Once a *pre-set* number of subword units was determined, a set of pronunciations for each word was found from its occurrences in the training data and were clustered to select representative pronunciations (Paliwal, 1990; Svendsen et al., 1995). The studies on isolated word recognition task on English demonstrated the potential of the approach. A limitation of these approaches is that they can generate pronunciations only for the words which are seen during training. Furthermore, these approaches need to know the word boundaries explicitly.

In (Jansen and Church, 2011), an approach was proposed in which the need for transcribed speech is limited. Specifically, given an acoustic example of each word, a spoken term discovery algorithm (Park and Glass, 2008) is exploited to search and cluster the acoustic realizations of the words from untranscribed speech. Then for each word cluster, a whole word HMM is trained in which each HMM state represents a subword unit. The number of subword units for each word is determined based on the duration of acoustic examples and the expected duration of a phone. The subword unit states are then finally clustered based on the pairwise similarities between their emission scores using a spectral clustering algorithm (Shi and Malik, 2000). The viability of the approach was limited to spoken term detection task. A limitation of the approach is that an acoustic example of each word in the dictionary is required.

Hartmann et al. (2013) proposed an approach based on the assumption that the orthography of the words and their pronunciations are related. In this approach, the subword units are obtained by clustering context-dependent (CD) grapheme models. This is achieved through a spectral based clustering approach (Ng et al., 2001), similar to (Jansen and Church, 2011). The main difference is that in this case the pairwise similarities are computed between the CD grapheme models (instead of the HMM states). The pronunciations for seen and unseen words are finally generated by employing a statistical machine translation (SMT) framework. On Wall Street Journal task, it was found that the resulting ASWU-based lexicon yields a better ASR system than the grapheme-based lexicon.

2.3.2. Joint approaches for ASWU derivation and pronunciation generation

As opposed to decoupling the ASWU derivation and pronunciation generation problems, there are also approaches which aim to jointly determine the subword units and pronunciations using a common objective function. In (Holter and Svendsen, 1997), this was done through an iterative process of acoustic model estimation and pronunciation generation. In (Bacchiani and Ostendorf, 1999, 1998), a segmentation and clustering approach was exploited for derivation of subword units, with two main differences compared to the approaches explained in Section 2.3.1: (1) in the segmentation step, pronunciation related constraints is applied such that a given word has the same number of segments across the acoustic training data, and (2) a maximum-likelihood criteria that is consistent for both segmentation and clustering is utilized. On read speech DARPA resource management task, it was shown that the proposed approach leads to improvements over the phone-based ASR system.

In (Singh et al., 2000, 2002), a maximum likelihood strategy was presented which decomposed the ASWU-based ASR system development as joint estimation of the pronunciation lexicon (including determination of ASWU set size) and acoustic model parameters. More precisely, with an initial pronunciation lexicon based on context-independent graphemes, the acoustic model parameters and the pronunciation lexicon are updated iteratively. The lexicon update step is an iterative process within itself consisting of word segmentation estimation given the acoustic model and update of the lexicon based on the segmentation. After each iteration of lexicon update and acoustic model update convergence is determined by evaluating the ASR system on cross-validation data. If not converged, the ASWU set size is increased and the process is repeated. A proof

of concept was demonstrated on DARPA Resource Management corpus.

Recently, in (Lee et al., 2013) a hierarchical Bayesian model approach was proposed to jointly learn the subword units and pronunciations. This is done by modeling two latent structures: (1) the latent phone sequence, and (2) the latent letter-to-sound (L2S) mapping rules, using an HMM-based mixture model in which each component represents a phone unit and the weights over HMMs are indicative of the L2S mappings. It was shown that the proposed approach together with the pronunciation mixture model retraining leads to improvements over the grapheme-based ASR system on a weather query task.

3. Proposed Approach

This section presents an HMM-based formulation to derive phone-like ASWUs and develop an associated pronunciation lexicon. Essentially, the formulation builds on grapheme-based ASR in deterministic lexical modeling framework as well as probabilistic lexical modeling framework. More specifically, we show that:

1. The problem of derivation of ASWUs can be cast as a problem of finding phone-like acoustic units $\{a^d\}_{d=1}^D$ given transcribed speech, i.e., the speech signal and its orthographic transcription, in the grapheme-based ASR framework. Section 3.1 dwells on this aspect.
2. Given the derived ASWUs $\{a^d\}_{d=1}^D$ and the transcribed speech, the pronunciation lexicon development problem can be cast as a problem akin to acoustic data-driven G2P conversion (Razavi et al., 2016). Section 3.2 deals with this aspect.

3.1. Automatic subword unit derivation

State clustering and tying methods in HMM-based ASR have emerged from the perspective of addressing data sparsity issue and handling unseen contexts (Young, 1992; Ljolje, 1994). However, this methodology can be adopted, as it is, to derive acoustic subword units in the framework of grapheme-based ASR. More precisely, we hypothesize and show that the clustered context-dependent grapheme units $\{a^d\}_{d=1}^D$ obtained in a context-dependent grapheme based ASR system can serve as phone-like subword units.

The reasoning behind our hypothesis is that the set of acoustic units $\{a^d\}_{d=1}^D$ is obtained by maximizing the likelihood of the training data, which is essentially

determined by estimation of $p(\mathbf{x}_t|q_t = l^i)$, as during training the sequence model for each utterance is fixed given the associated transcription and lexicon. As observed earlier in Eqn. (9), $p(\mathbf{x}_t|q_t = l^i)$ estimation involves matching of acoustic information \mathbf{v}_t with lexical information \mathbf{y}_i . We know that standard features such as cepstral features have been designed to model envelope of short-term spectrum, which carry information related to phones. In other words, standard feature such as MFCCs or PLPs for ASR primarily target modeling the spectral characteristics of vocal tract system while incorporating speech perception knowledge.

Similarly it is very well known that context-dependent graphemes capture information related to phones. This is one of the central assumptions in most of G2P conversion approaches, i.e., the relationship between context-independent graphemes and phones can be irregular but the relationship can become regular when contextual graphemes are considered. For example, as illustrated in Figure 2, in the decision tree-based G2P conversion approach (Pagel et al., 1998), given the grapheme context a decision tree is learned to map the central grapheme to a phoneme.

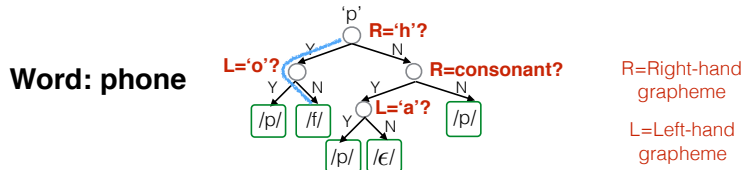


Figure 2: Example of the decision tree-based G2P conversion.

Therefore, as illustrated in Figure 3, for the likelihood of the training data to be maximized, clustered context-dependent grapheme units $\{a^d\}_{d=1}^D$ should model an information space that is common to both short-term spectrum based feature \mathbf{x}_t space and context-dependent grapheme based lexical unit l^i space, which we hypothesize it to be a phone-like subword unit space.

Our argument is further supported by an ASR study that demonstrated the interchangeability of clustered context-dependent phoneme units space and clustered context-dependent grapheme units space in the framework of probabilistic lexical modeling (Rasipuram and Magimai-Doss, 2013) as well as by earlier works on grapheme-based ASR that have explored integration of phonetic information in clustering context-dependent grapheme units and state tying (Killer et al., 2003).

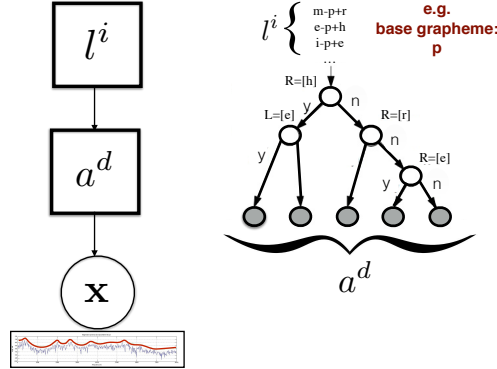


Figure 3: The clustered states a^d of a grapheme-based CD HMM/GMM system obtained through decision tree based clustering are exploited as ASWUs. As a^d should be related to both CD graphemes l^i and cepstral features \mathbf{x} , they are expected to be phone-like.

3.2. Lexicon development through grapheme-to-ASWU conversion

In order to build speech technologies with the derived ASWUs, we need a mechanism to map the orthographic transcription of words to sequence of ASWUs for both seen and unseen words. For that purpose, an approach similar to automatic G2P conversion is desirable. However, conventional G2P approaches are not directly applicable, as they necessitate a seed lexicon that maps a few word orthographies into sequence of phonemes (in our case ASWUs). More recently, it has been shown that G2P conversion can be achieved by learning the G2P relationship through acoustics using HMMs (Razavi et al., 2016). Such an approach has the inherent ability to alleviate the necessity for a seed lexicon, and thus can be exploited to develop a G2ASWU converter for lexicon development. This approach can be essentially considered as an extension of the grapheme-based ASR approach, where either a deterministic lexical model or a probabilistic lexical model $\{y_i\}_{i=1}^I$ that captures G2ASWU relationship is learned and ASWU-based pronunciations are inferred. We present below these two frameworks.

3.2.1. Deterministic lexical modeling based G2ASWU conversion

This method of lexicon development is a straightforward extension of the ASWU derivation. More precisely, in the process of ASWU derivation a deter-

ministic one-to-one map between context-dependent graphemes ($\{l^i\}_{i=1}^I$) and ASWUs ($\{a^d\}_{d=1}^D$) is learned. The pronunciations can be inferred using this information similar to the decision tree based G2P conversion approach (Pagel et al., 1998), discussed briefly earlier in Section 3.1 (Figure 2).

3.2.2. Probabilistic lexical modeling based G2ASWU conversion

Another possibility is to learn a probabilistic relationship between graphemes and ASWUs and infer pronunciations in terms of ASWUs following acoustic data-driven G2P conversion approach using KL-HMM (Rasipuram and Magimai-Doss, 2012; Razavi et al., 2016). This approach of G2ASWU conversion would involve,

1. training of an ANN-based \mathbf{z}_t estimator given the alignment of the training data in terms of $\{a^d\}_{d=1}^D$. This step is same as training a context-dependent neural network for ASR system;¹ then
2. training of a context-dependent grapheme-based KL-HMM using \mathbf{z}_t as feature observations (Magimai-Doss et al., 2011); and finally
3. inferring the pronunciations given the KL-HMM parameters $\{\mathbf{y}_i\}_{i=1}^I$ and the orthographies of the words in the lexicon. More precisely, first a sequence of ASWU posterior probability vectors is obtained from the KL-HMM given the orthography of the target word. The sequence is then decoded by an ergodic HMM in which each state represents an ASWU to infer the pronunciation.

3.3. Summary of the proposed approach

Figure 4 summarizes our approach. As illustrated, the approach consists of three phases. *Phase I* involves derivation of ASWUs. *Phase II* involves learning G2ASWU relationship given transcription and acoustic data. *Phase III* deals with lexicon development given the G2ASWU relationship and the word orthographies. *Phase II* is explicitly needed for learning probabilistic G2ASWU relationship. In the case of deterministic G2ASWU conversion, it is implicit in *Phase I*. *Phase III* can be seen as decoding a sequence of ASWU posterior probability vectors \mathbf{y}_i . It is worth mentioning that the pronunciation inference step, i.e. *Phase III*, for both deterministic and probabilistic lexical modeling

¹If the \mathbf{z}_t estimator is based on Gaussians then it would amount to going from single Gaussian to GMMs (mixture increment step) of ASR system training.

based approaches is the same. More precisely, in the case of deterministic lexical modeling based approach, the inference step is equivalent to decoding a sequence of Kronecker delta distributions resulting from the one-to-one mapping of CD graphemes (in the word orthography) to ASWU units using the decision tree (Razavi et al., 2016).

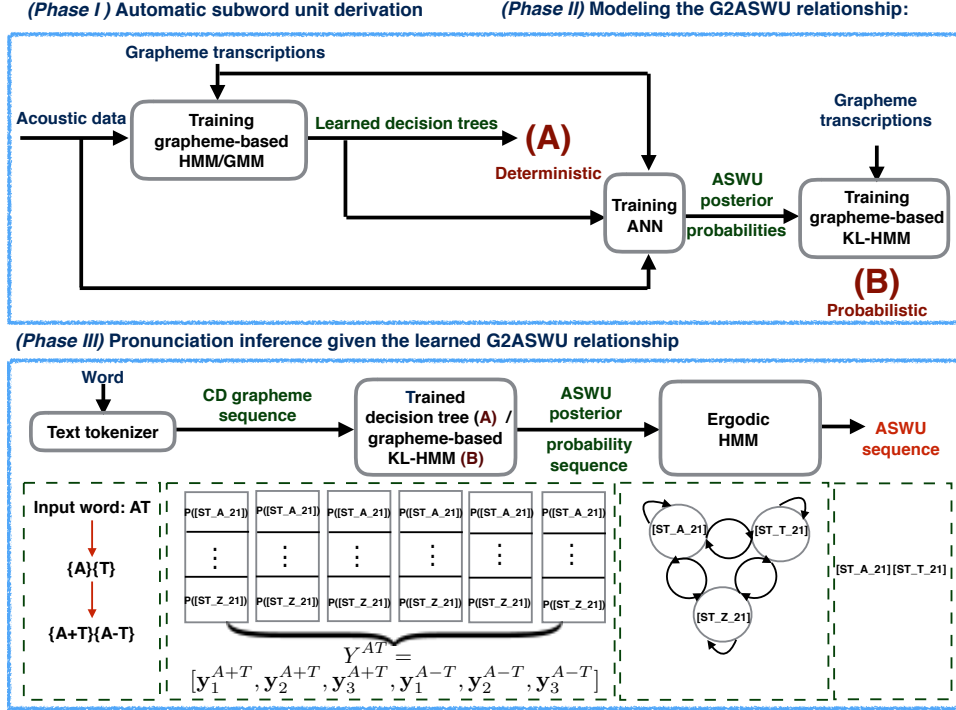


Figure 4: Block diagram of the HMM formalism for subword unit derivation and pronunciation generation. *Phase III* is shown for the case where the ASWU posterior probability vectors from KL-HMM are decoded. For the case where the ASWU posterior probability vectors are obtained from the decision trees (i.e., \mathbf{y}_i s are Kronecker delta distributions), only a single posterior probability vector per each context-dependent grapheme is generated, i.e., $Y^{AT} = [\mathbf{y}_1^{A+T}, \mathbf{y}_1^{A-T}]$

A central challenge in the proposed approach is how to determine the size of the ASWU set $\{a^d\}_{d=1}^D$. In the studies validating the proposed approach, presented in the remainder of the paper, we show that this can be achieved via cross-validation. Specifically, a range of values for acoustic units set cardinality D can be considered based on the knowledge that the ratio of number of phonemes to number of graphemes is not an extremely large value, and can be selected via cross-validation at ASR level. For instance in English, if one considers the CMU dictionary, then the ratio is $\frac{38}{26}$ or $\frac{84}{26}$ (when lexical stress is

considered). Alternately, the value of D can be chosen relative to the number of graphemes and is much smaller than the number of acoustic units considered for building context-dependent grapheme-based ASR systems, which is typically in the order of thousands.

4. In-Domain and Cross-Domain Studies on Resource-Rich Languages

In this section, we establish the proposed framework for subword unit derivation and lexicon development through experimental studies on a resource-rich language using only its word-level transcribed speech data. The rationale for studying on a well-resourced language is to enable analyzing the discovered subword units and relating them to phonetic identities. We selected English as the well-resourced language, as it is a challenging language for automatic pronunciation generation due to its irregular grapheme-to-phoneme relationship, and has been the focus of many previous works on ASWU derivation and lexicon development. Our investigations are organized as follows:

1. *Evaluation of the proposed approach through in-domain studies:* We investigate the proposed approach for derivation of ASWUs and corresponding pronunciations on two English corpora, namely Wall Street Journal (WSJ) and Resource Management (RM). We evaluate the ASWU-based lexicons through in-domain ASR studies where the performance of the ASWU-based ASR systems is compared against grapheme-based and phoneme-based ASR systems (Section 4.2).
2. *Investigating the transferability of the ASWUs through cross-domain studies:* A central challenge in ASWU based lexicon development and its adoption for wider use is ascertaining whether the ASWUs derived from limited amount of acoustic resources generalize across domains, similar to linguistically motivated subword units phonemes and graphemes. To the best of our knowledge, none of the previous works have tried to ascertain that aspect. In that sense, we go a step further to conduct cross-domain studies where the ASWUs are derived from the WSJ corpus and lexicon is developed for the RM corpus. We present three methods for development of lexicons in such a scenario, and investigate the transferability of the ASWUs by building and evaluating ASR systems using the developed lexicons (Section 4.3).

3. *Comparison to related approaches in the literature:* in Section 2.3, we discussed a few prominent approaches proposed in the literature for derivation of ASWUs and pronunciation generation. We compare the performance of the our approach with two of the related approaches in the literature studied on WSJ0 and RM corpora (Section 4.4). Indeed, one of the main reasons for selecting these two corpora is to enable comparison to these related works in the literature.

4.1. Databases

This section describes the setup on two corpora used in our experimental studies.

4.1.1. WSJ0 corpus

The WSJ corpus has been originally designed for large vocabulary speech recognition and natural language processing, and it contains a wide range of vocabulary size (Paul and Baker, 1992). The WSJ corpus (Woodland et al., 1994) has two parts - WSJ0 with 14 hours of speech and WSJ1 with 66 hours of speech. In this article, we use the WSJ0 corpus for training, which contains 7106 utterances (about 14 hours of speech) and 83 speakers. We report recognition studies on Nov92 test set, which contains 330 utterances from 8 speakers unseen during training. The training set contains 10k unique words. The recognition vocabulary size is 5k words. The language model consists of a bigram model. The grapheme lexicon was obtained from the orthography of the words and contained 27 subword units including silence. We refer to this lexicon as *Lex-WSJ-Gr-27*. The phoneme lexicon was based on UNISYN dictionary.

4.1.2. DARPA Resource Management corpus

The DARPA Resource Management (RM) task is a 1000 word continuous speech recognition task based on naval queries (Price et al., 1988). The training set consists of 3990 utterances spoken by 109 speakers amounting to approximately 3.8 hours speech data. The test set, formed by combining Feb89, Oct89, Feb91 and Sep92 test sets, contains 1200 utterances amounting to 1.1 hours of speech data. The word-pair grammer supplied with the RM corpus was used as the language model for decoding. The grapheme lexicon was obtained from the orthography of the words. In addition to the English characters, silence, symbol hyphen and symbol single quotation mark was considered as separate graphemes. Therefore, the lexicon contained 29 subword units. We refer to

this lexicon as *Lex-RM-Gr-29*. The phoneme lexicon was based on UNISYN dictionary. As mentioned earlier, the RM corpus is mainly used to investigate transferability of the ASWUs across domains. So, it is worth pointing out that 507 out of the 990 words in the RM corpus do not appear in the WSJ0 training set vocabulary.

4.2. In-domain ASR studies

In this section we first explain the setup for derivation of ASWUs and development of ASWU-based lexicons. We then present the in-domain ASR studies for evaluation of the ASWU-based lexicons.

4.2.1. ASWU derivation and lexicon development setup

The setup for subword unit derivation and lexicon development through G2ASWU conversion is as follows:

Acoustic subword unit derivation: Towards automatic discovery of subword units, cross-word single preceding and single following CD grapheme-based HMM/GMM systems were trained with 39 dimensional PLP cepstral features extracted using HTK toolkit (Young et al., 2000). Each CD grapheme was modeled with a single HMM state. The subword units were derived through likelihood-based decision tree clustering using singleton questions. Different number of ASWUs were obtained by adjusting the log-likelihood increase during decision tree based state tying. The numbers of clustered units were obtained such that they are within the range of 2 to 4 times the number of graphemes, based on the general idea explained in Section 3.3. Therefore, for the WSJ0 corpus, ASWUs of size 60, 78 and 90 were investigated, and for the RM corpus, ASWUs of size 79, 92 and 109 were studied.

Deterministic lexical modeling based G2ASWU conversion: Given the learned decision trees for each ASWU set, the pronunciation for each word was inferred by mapping each grapheme in the word orthography to an ASWU by considering its neighboring (i.e., single preceding and single following) grapheme context. We denote the lexicons in the form of *Lex-DB-Det-ASWU-M* where *DB* and *M* correspond to the database and the number of ASWUs respectively. For example, the lexicon generated on WSJ0 corpus using 78 ASWUs is denoted as *Lex-WSJ-Det-ASWU-78*.

Probabilistic lexical modeling based G2ASWU conversion: In this case, given the obtained ASWUs:

1. A five-layer multilayer Perceptron (MLP) was trained to classify the ASWUs. The input to the MLP was 39-dimensional PLP cepstral features with four preceding and four following frame context. The hyper parameters such as the number of hidden units per hidden layer were decided based on the frame accuracy on the development set. Each hidden layer had 2000 and 1000 hidden units in the WSJ0 and RM corpora respectively. The MLP was trained with output non-linearity of softmax and minimum cross-entropy error criterion using Quicknet software (Johnson et al., 2004).
2. Using the posterior probabilities of ASWUs as feature observations, a grapheme-based KL-HMM system modeling single preceding and single following grapheme context was then trained. Each CD grapheme was modeled with three HMM states. The parameters of the KL-HMM were estimated by minimizing a cost function based on the reverse KL-divergence (RKL) local score (Aradilla et al., 2008), i.e., the MLP output distribution is the reference distribution, as previous studies had shown that training KL-HMM with RKL local score enables capturing one-to-many grapheme-to-phoneme relationships (Rasipuram and Magimai.-Doss, 2013). Unseen grapheme contexts were handled by applying the KL-divergence based decision tree state tying method proposed in (Imseng et al., 2012).
3. Given the orthography of the word and the KL-HMM parameters, the pronunciations were inferred by using an ergodic HMM in which each ASWU was modeled with three left-to-right HMM states.

During pronunciation inference, some of the ASWUs with less probable G2ASWU relationships were automatically pruned or filtered out. This can be observed from Table 1, which shows the properties of the ASWU-based lexicons together with the MLPs used for the WSJ0 and RM corpora respectively. The MLPs are denoted as MLP- DB - N , with DB and N denoting the database and the size of the ASWU set respectively. Similarly, the lexicons are shown as Lex- DB -Prob-ASWU- M , with M denoting the actual number of ASWUs used in the lexicon. As an example, it can be seen that in Lex- RM -Prob-ASWU-101, from the 109 original ASWU set, only 101 remained after G2ASWU conversion.

Table 1: Summary of the ASWU-based lexicons obtained through probabilistic lexical modeling based G2ASWU conversion for WSJ0 and RM corpora.

(a) WSJ0 corpus	
Lexicon	MLP
Lex- <i>WSJ</i> -Prob-ASWU-58	MLP- <i>WSJ</i> -60
Lex- <i>WSJ</i> -Prob-ASWU-74	MLP- <i>WSJ</i> -78
Lex- <i>WSJ</i> -Prob-ASWU-88	MLP- <i>WSJ</i> -90
(b) RM corpus	
Lexicon	MLP
Lex- <i>RM</i> -Prob-ASWU-77	MLP- <i>RM</i> -79
Lex- <i>RM</i> -Prob-ASWU-90	MLP- <i>RM</i> -92
Lex- <i>RM</i> -Prob-ASWU-101	MLP- <i>RM</i> -109

4.2.2. Selection of optimal ASWU-based lexicon

Given different lexicons obtained through deterministic and probabilistic G2ASWU conversion, the optimal lexicon was determined based on the ASR accuracy on the development set. More precisely, first HMM/GMM systems using different ASWU-based lexicons were trained with 39 dimensional PLP cepstral features. Finally, the ASWU-based lexicon which led to the best performing HMM/GMM ASR system on the development set was selected.² In our experiments, in case of using the deterministic G2ASWU conversion for pronunciation generation, Lex-Det-*WSJ*-ASWU-90 and Lex-Det-*RM*-ASWU-92; and in case of using the probabilistic approach, Lex-Prob-*WSJ*-ASWU-88 and Lex-Prob-*RM*-ASWU-90 were selected as the optimal lexicons and are therefore used in the rest of the article.

4.2.3. Evaluation

To evaluate the generated ASWU-based lexicons, we compared the performance of ASWU-based ASR systems with the grapheme-based and phoneme

²It is worth mentioning that for WSJ0 and RM corpora there are no explicit development sets defined. To be more precise, in the case of RM the development set (1110 utterances) was merged with the training set (2880) to create training set of 3990 utterances in literature. So, we used the part of the data that was used for early stopping through cross validation in MLP training as the development data, and trained ASWU-based HMM/GMM systems on the remaining part of the training data. For instance, in the case of RM three HMM/GMM systems corresponding to the lexicons *Lex-RM-Prob-ASWU-77*, *Lex-RM-Prob-ASWU-90*, *Lex-RM-Prob-ASWU-101* were trained on 2880 utterances and the lexicon was selected using the 1110 utterances. We followed similar procedure for WSJ0.

based ASR systems. Toward that, we trained both context-independent and cross-word context-dependent HMM/GMM systems with 39 dimensional PLP cepstral features. Each subword unit was modeled with three HMM states. For the CI grapheme-based systems, the number of Gaussian mixtures for each HMM state was decided based on the ASR word accuracy on the cross-validation set, resulting in 256 and 128 Gaussian mixtures for WSJ0 and RM corpora respectively. In case of using ASWUs, in order to have a comparable number of parameters to the grapheme based ASR system, each HMM state was modeled with 64 and 32 Gaussian mixtures in the WSJ0 and RM corpora respectively. Similarly, for phone subword units, the number of Gaussian mixtures for each HMM state was 128 and 64 in the WSJ0 and RM corpora. In the context-dependent case, for tying the HMM states, only singleton questions were used. Each tied state was modeled by a mixture of 16 and 8 Gaussians on WSJ0 and RM corpora respectively. The number of tied states in all the systems trained on a corpus was roughly the same to ensure that possible improvements in ASR accuracy are not due to the increase in complexity.

Throughout this article, we report the ASR system performances in terms word recognition rate ($100 - \text{word error rate}$), denoted as WRR. Furthermore, for comparing the performance of different systems, we applied the statistical significant test presented in (Bisani and Ney, 2004) with the confidence level of 95%.

Table 2 presents the performance of ASR systems based on different lexicons. In the case of using CI units, the ASWU-based ASR systems perform significantly better than the grapheme-based ASR systems in both WSJ0 and RM corpora. In the case of CD units, it can be seen that for the WSJ0 corpus, the HMM/GMM system using ASWUs performs significantly better than the baseline grapheme-based ASR system. For the case of RM corpus, however, the improvements are not statistically significant. This could be due to the fact that in RM task all the words are seen during both training and evaluation. In all cases, the the ASWU based lexicon yields a system that lies between phoneme-based ASR system and grapheme-based ASR system.

When using CI subword units, it can be seen that the performance of the system using probabilistic lexical modeling based G2ASWU conversion is comparable or even better than the system using deterministic lexical modeling G2ASWU conversion, whereas when using CD subword units, this is not the case. A plausible reasoning for such a trend is that CI subword unit based systems using deterministic lexical modeling based G2ASWU conversion may

require more parameters. We tested that by building CI ASWU-based ASR systems using deterministic and probabilistic lexical modelling based pronunciations with varying number of Gaussian mixtures (from 8 to 256). We observed that the difference between the best performing CI ASR systems using deterministic and lexical modeling based G2ASWU conversion is not statistically significant³, thus indicating that the deterministic lexical modeling based G2ASWU conversion approach leads to a better ASR system compared to the probabilistic approach. A potential explanation for this difference could be that, unlike the probabilistic lexical modeling based G2ASWU conversion approach, deterministic lexical modeling based G2ASWU conversion approach avoids ASWU deletions and could therefore generate a more consistent pronunciation lexicon for English.

Table 2: HMM/GMM ASR system performances in terms of WRR using CI and CD subword units.

(a) WSJ0 corpus.			(b) RM corpus.		
Lexicon	CI	CD	Lexicon	CI	CD
Lex- <i>WSJ</i> -Gr-26	68.9	85.8	Lex- <i>RM</i> -Gr-29	84.2	94.0
Lex- <i>WSJ</i> -Det-ASWU-90	78.6	88.7	Lex- <i>RM</i> -Det-ASWU-92	89.1	94.5
Lex- <i>WSJ</i> -Prob-ASWU-88	78.7	87.3	Lex- <i>RM</i> -Prob-ASWU-90	90.7	94.2
Lex- <i>WSJ</i> -Ph-45	88.6	93.5	Lex- <i>RM</i> -Ph-45	93.5	95.9

4.3. Cross-domain ASR studies

This section presents a study that investigates the transferability of the ASWUs to a condition or domain unobserved during derivation of ASWU. As noted earlier, for ASWUs to be adopted for mainstream speech technology, this characteristic is highly desirable. Toward that we present a cross-database study where the ASWU derivation is carried out on out-of-domain (OOD) WSJ0 corpus and the lexicon is developed for target domain RM corpus. Similar to G2P conversion as elucidated in (Razavi et al., 2016), G2ASWU conversion (presented earlier in Section 3.2) can be seen as a two step process: 1) Learning

³For the WSJ0 corpus, the best performing CI ASR systems yielded WRR of 80.1 % and 79.7% ASR when using *Lex-WSJ-Det-ASWU-90* and *Lex-WSJ-Prob-ASWU-88*, respectively. For the RM corpus, the best performing CI ASR systems yielded WRR of 90.2% and 90.7% ASR word when using *Lex-RM-Det-ASWU-92* and *Lex-RM-Prob-ASWU-90*, respectively.

the relationship between the graphemes and the derived ASWUs, and 2) Inferring the ASWU sequence (pronunciation) given the word orthography and the learned G2ASWU relationship. We present three methods for cross-domain ASWU-based lexicon development based on that understanding.

Method-I: Applying standard G2P conversion approach on the seed lexicon obtained from the OOD corpus

One possible way to generate pronunciations for the in-domain RM corpus is to use the ASWU-based lexicon from the WSJ0 corpus as the seed lexicon and train a G2ASWU converter. For this purpose, we investigated one of the state-of-the-art G2P conversion approach, namely, joint multigram approach (Bisani and Ney, 2008) for G2ASWU conversion. This was done by using the Sequitur software developed at RWTH Aachen University.⁴ In our experiment, the maximum width of the grapheme used was one, and the n-gram context size was 6.⁵ As shown in Figure 5, first the G2ASWU relationship is learned on the ASWU-based lexicon for the WSJ0 corpus by training the G2ASWU converter. Then given the words in the RM corpus and the learned G2ASWU relationship, the pronunciations are inferred.⁶

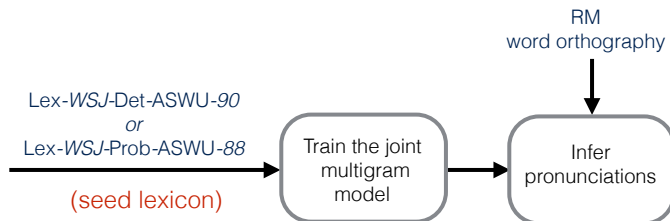


Figure 5: Diagram of joint multigram-based pronunciation generation for RM corpus using the seed lexicon trained on WSJ0 corpus (*Method-I*).

Method-II: Using the learned G2ASWU relationship on the OOD corpus for pronunciation inference on the in-domain corpus

Instead of using the ASWU-based lexicon from the WSJ0 corpus, only the learned G2ASWU relationships can be exploited for inferring pronunciations

⁴<http://www-i6.informatik.rwth-aachen.de/web/Software/g2p.html>

⁵As there are no canonical pronunciations in case of using ASWUs are available, we decided on the optimal n-gram context size based on the ASR accuracy.

⁶ The grapheme symbols such as single hyphen that appear in the RM word orthographies and have not been observed in the WSJ0 word orthographies were removed for the inference.

on the RM corpus. More precisely, we investigate use of the deterministic and probabilistic G2ASWU relationships obtained from (a) the decision trees learned on WSJ0, and (b) the KL-HMM trained on WSJ0, respectively to generate pronunciations for the RM corpus, as illustrated in Figure 6.

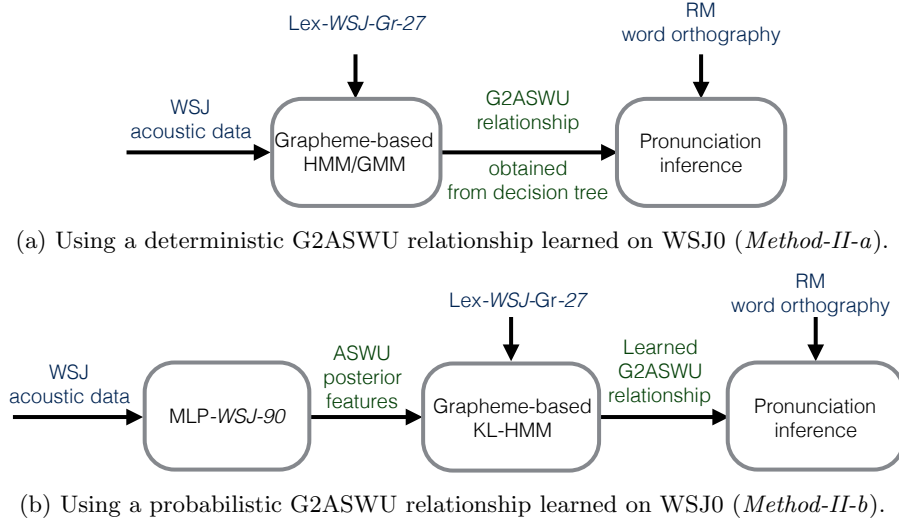


Figure 6: Illustration of pronunciation generation for RM corpus in *Method-II*.

Method-III: Learning the G2ASWU relationship on the in-domain corpus through acoustics

Instead of using the learned G2ASWU relationship on the WSJ0 corpus, we can use the trained MLP on WSJ0 corpus to estimate ASWU posterior probabilities for the RM speech data. Given the ASWU posterior probabilities as feature observations, a grapheme-based KL-HMM system can be trained on the RM corpus data. The pronunciation inference can then be done given the trained KL-HMM and the word orthographies, as shown in Figure 7.

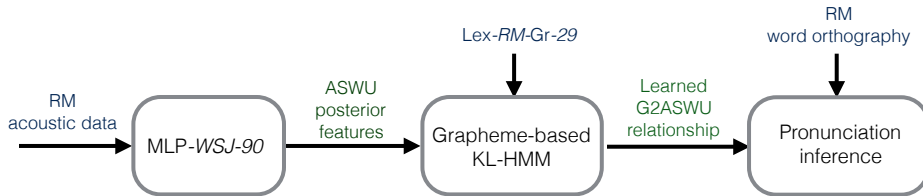


Figure 7: Illustration of pronunciation generation for RM corpus using Method III.

We generated ASWU-based lexicons for the RM corpus based on the above presented three methods. It is worth to reiterate that, in addition to acoustic differences between the two corpora, there are also differences at lexicon level, i.e. 507 out of the 990 words in the RM lexicon do not appear in WSJ0 lexicon. For each of the lexicons developed, we trained context-independent and cross-word context-dependent ASWU-based HMM/GMM system with 39 dimensional PLP cepstral features extracted using HTK toolkit. Each subword unit was modeled with three HMM states. Each CI HMM state was modeled by 32 Gaussian mixtures similar to in-domain studies in Section 4.3. Each tied HMM state was modeled by a mixture of 8 Gaussians. The HMM states were tied using singleton question set.

Table 3 presents the results in terms of WRR. It can be observed that the context-independent ASR systems, regardless of the method used for pronunciation generation, perform better than the grapheme-based CI ASR system (Table 2). The performance of the context-dependent ASR systems using the pronunciations generated through *Method-I* is inferior to the grapheme-based ASR system (Table 2). The performance of the ASR systems using *Method-II* for pronunciation generation are comparable with the ASR systems obtained through in-domain studies (Table 2). Generating pronunciations using *Method-III* also leads to a comparable system to the in-domain ASWU-based ASR systems. Comparing the performance of the systems using *Method-I* for pronunciation generation with the systems using *Method-II* and *Method-III* shows that it is better to transfer the learned G2ASWU relationship or learn the G2ASWU relationship on target domain speech. A potential reason for that is that *Method-I* relies on availability of ground truths, like availability of seed lexicon obtained through linguistic expertise in G2P conversion, which in the present scenario is not available. Overall, *Method-II* leads to the best ASR performance. It may be possible to improve *Method-III* by acoustic model adaptation techniques to adapt the MLP trained on the out-of-domain data. This is open for further research. Together these studies show that, in the proposed approach, the derived ASWUs and the G2ASWU relationship learned from one domain are transferrable to another or target domain. Alternately, the proposed approach inherently enables such transfer.

4.4. Comparison to existing approaches

In this section, we compare the present work with two existing approaches in the literature that have reported studies on the WSJ0 and RM corpora with

Table 3: ASR system performances in terms of WRR on RM corpus using different cross-domain pronunciation generation methods.

Method	G2ASWU relationship	CI	CD
<i>Method-I</i>	Deterministic	87.5	92.3
	Probabilistic	85.2	91.3
<i>Method-II</i>	Deterministic	89.0	94.4
	Probabilistic	88.8	94.0
<i>Method-III</i>	Probabilistic	89.0	94.0

the same setup as that used in our studies. More precisely, on WSJ0 corpus, Section 4.4.1 compares our approach to the spectral clustering based approach proposed in (Hartmann et al., 2013). Section 4.4.2 studies the proposed approach in comparison to the approach proposed by Bacchiani and Ostendorf in (Bacchiani and Ostendorf, 1999).

4.4.1. Comparison to Hartmann et al. (2013) approach

In essence, the proposed approach is similar to the spectral based clustering approach proposed in (Hartmann et al., 2013), as they both discover the ASWUs from the grapheme-based HMM/GMM system. However, there are two key differences between these approaches:

1. In our approach, the ASWUs are discovered through decision-tree based clustering of the HMM states, while in (Hartmann et al., 2013), the sub-word units are derived through spectral based clustering, which requires computation of similarity matrix between HMMs.
2. In our approach, the pronunciations are generated using the KL-HMM framework, while in (Hartmann et al., 2013), the pronunciations are transformed using a statistical machine translation approach.

As the experimental setup in this article on WSJ0 corpus and the work in (Hartmann et al., 2013) are the same, we provide a comparison between the baseline and the results in both works in Table 4. In (Hartmann et al., 2013) there are two grapheme baselines: one based on the standard orthography (denoted as grapheme-direct) and the other based on grapheme-to-grapheme (G2G) conversion (denoted as grapheme-transformed) employing an approach similar to machine translation. Similarly, in the ASWU based study they have two systems: one where the pronunciations are generated directly by mapping the graphemes to ASWUs based on the spectral clustering (denoted as ASWU-direct), and the other where ASWU-to-ASWU conversion is performed like G2G

case mentioned above (denoted as ASWU-transformed). We ensured that our systems have comparable number of parameters in the case of both using CI subword units and CD subword units based systems. It can be observed that the ASWU-based lexicon developed by our approach leads to a better ASR system. Furthermore, when comparing the best systems there is an absolute difference of 2.5% WRR, which indicates that the proposed approach in this article leads to a better ASR system.

Table 4: Comparison with the related work in (Hartmann et al., 2013).

Approach	Lexicon	CI	CD
	Grapheme-direct	60.1	84.2
Approach proposed in (Hartmann et al., 2013)	Grapheme-transformed	68.6	85.5
	ASWU-direct	70.7	85.6
	ASWU-transformed	76.7	86.2
Present work	Grapheme	68.9	85.8
	Lex- <i>WSJ</i> -Det-ASWU-90	78.6	88.7
	Lex- <i>WSJ</i> -Prob-ASWU-88	78.7	87.3

4.4.2. Comparison to Bacchiani and Ostendorf (1999) approach

In a broad sense, the proposed approach and the joint subword unit derivation and pronunciation generation method proposed in (Bacchiani and Ostendorf, 1999) can be considered to be similar as,

1. both approaches consist of segmentation and clustering steps, except that in our approach the segmentation and clustering is guided through graphemes during the HMM/GMM training; and
2. both approaches apply the pronunciation length constraint which ensures uniformity in the number of segments for training tokens of a word. In our approach this is automatically achieved through use of a unique grapheme sequence representation for each word.

In our studies, we have used RM corpus, which was also used in (Bacchiani and Ostendorf, 1999). However there are a few distinctions. In (Bacchiani and Ostendorf, 1999), the states of the HMMs were modeled by single Gaussian as opposed to mixture of Gaussians and the evaluation was carried out only on *Feb89* test set. So we also trained single Gaussian HMM/GMM system using the ASWU lexicon developed by our approach and evaluated on *Feb89* test set. Table 5 presents the results in the case where the two approaches are similar in

terms of number of ASWUs and clustered states. Table 6 provides a comparison between the best performance reported in (Bacchiani and Ostendorf, 1999) and the performance achieved with the lexicon based on our approach on the *Feb89* test set with 2937 clustered states. These results indicate that the ASWU lexicon developed by the proposed approach can yield ASR systems comparable to the ASWU lexicon developed by Bacchiani and Ostendorf (1999) approach, which needs additional heuristics to constrain the ASWU derivation and pronunciation generation process and necessitates all the words to be observed.

Table 5: Comparison with the related work in (Bacchiani and Ostendorf, 1999) on *Feb89* test set using single Gaussian distributions.

	# of base units	# of clustered states	WRR
Approach proposed in (Bacchiani and Ostendorf, 1999)	124	1519	86.3
Present work	92	1559	86.9

Table 6: Comparison of the best result reported in (Bacchiani and Ostendorf, 1999) on *Feb89* test set with the result using the present work on the same test set using single Gaussian distributions.

	WRR
Approach proposed in (Bacchiani and Ostendorf, 1999)	91.2
Present work	91.1

Before concluding this section, it is worth mentioning that the approach proposed in (Singh et al., 2002) was also investigated on RM corpus. Furthermore, there are also similarities w.r.t our approach, as it also exploits transcribed speech data and it uses a grapheme-based dictionary as the initial lexicon. However, the results presented in (Singh et al., 2002) can not be fairly compared against our results for the following reasons: (1) the training and test sets are different. In particular, in their studies the test set contains 1600 utterances as opposed to the standard test of 1200 utterances, and (2) their ASR system is based on semi-continuous HMMs while in the present work the ASR system is based on continuous density HMMs. Informally, it can be stated that in the present article the proposed approach has been investigated against stronger grapheme-based and phoneme-based baselines than the investigations reported (Singh et al., 2002).

5. Application to an Under-Resourced Language

In the previous section, we demonstrated the potential of the proposed framework for subword unit derivation and pronunciation generation on well-resourced language English. Most of the state-of-the-art speech recognition approaches have emerged through investigations on English. So it can be argued that our approach of deriving ASWU using grapheme-based HMM/GMM system may be well suited just for English. Furthermore, grapheme-to-phoneme relationship varies across languages. So a question arising is whether the proposed approach scalable to other languages or not.

In this section, our goal is two folds. More precisely, to show the scalability of the approach to a new language as well as its utility to under-resourced languages, specifically languages that do not have well developed phonetic resources. In that direction, we present investigations on a genuinely under-resource language, Scottish Gaelic. Unlike English, which belongs to family of Germanic languages, Scottish Gaelic belongs to family of Celtic languages. Our investigations are organized along two lines,

1. *Monolingual ASR studies*: We investigate the potential of the ASWU-based lexicons through monolingual ASR studies where we compare the performance of the ASWU-based ASR system with the alternative grapheme-based ASR system, as done in the studies on English.
2. *Multilingual ASR studies*: In (Rasipuram and Magimai.-Doss, 2015), it has been shown that performance of under-resourced ASR system can be significantly improved by (a) training a multilingual acoustic model that estimate multilingual phone posterior probabilities using resources of resource rich languages, and then (b) learning a probabilistic lexical model that captures the grapheme-to-multilingual phone relationship on the target language speech. So we also investigate if the ASWU-based lexicons hold their benefit in such a multilingual ASR system scenario as well. As a product of the study, later in Section 6, we show how phonetic identities of the derived ASWUs could be discovered.

The remainder of the section is organized as follows. Section 5.1 presents the database and experimental setup used. Section 5.2 presents the details of the ASWU-based lexicon development. Finally, Section 5.3 and 5.4 presents the monolingual ASR and multilingual ASR studies, respectively.

5.1. Database

This section first describes the characteristics of the Scottish Gaelic language. It then explains the Scottish Gaelic corpus used in our studies.

5.1.1. Scottish Gaelic language

Scottish Gaelic belongs to the class of Celtic languages. There are six Celtic languages that are still spoken. These languages are divided into two groups of Goidelic languages and Brythonic languages. Scottish Gaelic belongs to Goidelic languages along with Irish and Manx. It can be considered as a truly endangered language as it is spoken by about 60,000 people only. There are about 51 phonemes in the language (Wolters, 1997). However, the number of phonemes can change depending on the dialect. The language lacks a proper phonetic lexicon and the available transcribed speech data are also limited.

Scottish Gaelic alphabet has 18 letters, consisting of five vowels and thirteen consonants. The long vowels are represented with grave accents (À, È, Ì, Ò, Ù). There are twelve basic consonant types in Scottish Gaelic (B, C, D, F, G, I, L, M, N, P, R, S, T):

- Each consonant is either fortis or lenis (i.e., they are produced with greater or less energy). The lenited consonants are presented in the orthography with a grapheme [H] next to them.
- Each consonant is either broad (velarized) or slender (palatalized). Broad consonants are surrounded by broad vowels (A, O or U), while slender consonants are surrounded by slender vowels (E or I).

Scottish Gaelic orthography is less complicated than English. The complications partly arise due to the reason that modern orthography is based on Classical Irish orthography and the letter-to-sound rule may depend on the dialect (Wolters, 1997). The number of graphemes in Gaelic words are typically greater than the number of phones in the word due to the effect of lenited and broad/slender graphemes on the pronunciation. The grapheme-to-phoneme relationship in Scottish Gaelic can therefore be many-to-one. For example, the ratio of the number of graphemes to phonemes in the Gaelic word *SUID-HEACHADH* with pronunciation "sMj@x@G" (in the SAMPA format) is 1.7.

5.1.2. Scottish Gaelic corpus

The Scottish Gaelic corpus was collected by the University of Edinburgh in 2010 and contains recordings from broadcast news and discussion programs.⁷ In this article, the database is partitioned into training, development and test sets according to the structure provided in (Rasipuram et al., 2013b). The overview of the Scottish Gaelic corpus is given in Table 7.

Table 7: Overview of the Scottish Gaelic corpus in terms of number of utterances, hours of speech data and speakers in the train, cross-validation and test sets.

Number of	Train	Cross-validation	Test
Utterances	2389	1112	1317
Hours	3	1	1
Speakers	22	12	12

The database does not provide any phonetic lexicon. The graphemic lexicon can be simply obtained from the orthography of the words. As the corpus also contains borrowed English words, the graphemes J, K, Q, V, W, X, Y and Z are also present in the lexicon. Therefore the lexicon consists of 32 graphemes including silence as shown in Table 8. We refer to this lexicon as Lex-*SG-Gr-32*.

As the corpus does not provide a language model, we used a bigram language model trained on the sentences from the test set, as done in (Rasipuram et al., 2013b).

Table 8: Graphemes used in the Scottish Gaelic corpus.

Vowels	A, E, I, O, U, À, È, Ì, Ò, Ù
Consonants	B, C, D, F, G, H, I, L, M, N, P, R, S, T
English Graphemes	J, K, Q, V, W, X, Y

5.2. ASWU derivation and pronunciation generation setup

The setup for subword unit derivation and pronunciation generation for Scottish Gaelic is as follows:

Acoustic subword unit derivation: For automatic discovery of subword units, cross-word CD grapheme-based HMM/GMM systems were trained using 39-dimensional PLP cepstral features. Each CD grapheme was modeled with a single HMM state. Different numbers of ASWUs were obtained by adjusting

⁷<http://forum.idea.ed.ac.uk/tag/scots-gaelic>

the log-likelihood increase during decision tree clustering. The range for the number of ASWUs was decided to be similar to the range investigated in the studies on English, resulting in 85, 91 and 97 units.

Deterministic lexical modeling based G2ASWU conversion: For deterministic lexical modeling based G2ASWU conversion, the learned decision trees during ASWU derivation were exploited to map each grapheme in the word to an ASWU. We denote the lexicons generated using the deterministic lexical modeling based G2ASWU conversion as *Lex-SG-Det-ASWU-M* where *M* denotes the number of ASWUs.

Probabilistic lexical modeling based G2ASWU conversion: For probabilistic lexical modeling based G2ASWU conversion, first a five-layer MLP classifying ASWUs was trained in which each hidden layer had 1000 hidden units. Then given the ASWU posterior probabilities from the ANN as feature observations, a CD grapheme-based KL-HMM was trained. For the pronunciation inference, the ASWU posterior probabilities were decoded through the ergodic HMM in which each ASWU was modeled with three left-to-right HMM states.

Table 9 shows the properties of the ASWU-based lexicons generated using a probabilistic lexical modeling based G2ASWU conversion. Similar to the studies on English, it can be observed that some of the ASWUs are pruned out during the pronunciation generation given the probabilistic G2ASWU mapping.

Table 9: Summary of the ASWU-based lexicons obtained through probabilistic lexical modeling based G2ASWU conversion for Scottish Gaelic corpus.

Lexicon	MLP
<i>Lex-SG-Prob-ASWU-76</i>	<i>MLP-SG-85</i>
<i>Lex-SG-Prob-ASWU-82</i>	<i>MLP-SG-91</i>
<i>Lex-SG-Prob-ASWU-86</i>	<i>MLP-SG-97</i>

We selected the optimal number of ASWUs and the corresponding lexicon based on the WRR on the development set. *Lex-SG-Det-ASWU-85* and *Lex-SG-Prob-ASWU-82* yielded the best ASR systems and are therefore used in the ASR studies presented below.

5.3. Monolingual ASR system studies

As mentioned earlier, there is no well developed phonetic lexicon for Scottish Gaelic. So we evaluate the utility of the developed ASWU-based lexicon against grapheme-based lexicon by conducting monolingual ASR studies. Specifically, we compare them across two frameworks, namely, HMM/GMM framework and KL-HMM framework.

HMM/GMM framework: We trained CI and cross-word CD HMM/GMM systems with 39 dimensional PLP cepstral features extracted using HTK toolkit. Each subword unit was modeled with three HMM states. In the case of using CI subword units, the optimal number of Gaussian mixtures for the grapheme-based ASR system was 64 based on the best WRR obtained on the development set. For the ASWU-based ASR systems, the number of Gaussian mixtures was set to 16 so as to have a comparable number of parameters to the grapheme-based system. In the case of using CD subword units, for tying the HMM states singleton questions were used. Each HMM state was modeled by a mixture 8 Gaussians. The number of tied states in all the systems were roughly the same.

KL-HMM framework: This is done by using the posterior based framework of KL-HMM explained in Section 2.1 directly for speech recognition. More precisely, instead of using the KL-HMM parameters capturing a probabilistic G2ASWU relation for pronunciation inference, they are used in the KL-HMM ASR framework. In this case, we can visualize it as an approach that integrates pronunciation learning implicitly as a phase in ASR system training (Rasipuram et al., 2015). Our main motivation for performing this study was to ascertain whether doing lexicon development and ASR training as two separate stages can bring any advantage over doing direct speech recognition using grapheme-based KL-HMM system. For this purpose, we compared three KL-HMM systems, as illustrated in Figure 8, corresponding to lexicons Lex-*SG-Gr-32*, Lex-*SG-Det-ASWU-85* and Lex-*SG-Prob-ASWU-82*, respectively. All the systems use the same MLP, which is *MLP-SG-91*, as the acoustic model to estimate posterior feature observations.

Table 10 presents the HMM/GMM systems and KL-HMM systems performance in terms of WRR. It can be observed that Lex-*SG-Prob-ASWU-82* yields significantly better CI and CD systems than Lex-*SG-Gr-32* in both HMM/GMM framework and KL-HMM framework. Lex-*SG-Det-ASWU-85* yields a better system in KL-HMM framework but worse system in HMM/GMM framework

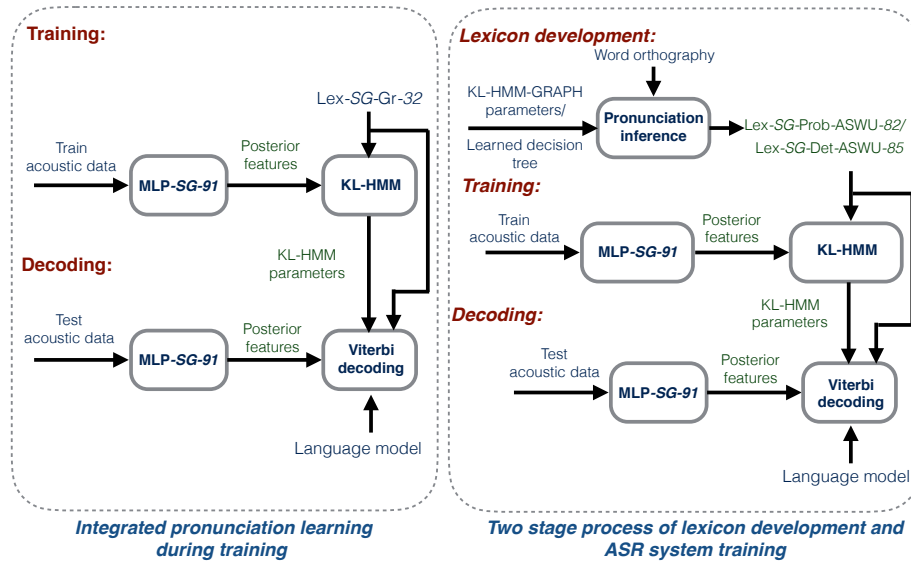


Figure 8: Illustration of KL-HMM based ASR system based on Lex-SG-Gr-32, Lex-SG-Det-ASWU-85 and Lex-SG-Prob-ASWU-82

against Lex-SG-Gr-32. A possible reason for such a trend could be that, as discussed earlier, in Scottish Gaelic the G2P relationship is many-to-one due to lenition and broad and slender consonants. So, when inferring pronunciations using the deterministic G2ASWU mappings, each grapheme in the word is invariably mapped into an ASWU. This can result in systematic erroneous pronunciations, which could lead to mismatch between acoustics and pronunciation model, as in the case of pronunciation variation. In the literature, it has been observed that KL-HMM approach is capable of handling pronunciation variation (Imseng et al., 2011; Razavi and Magimai.-Doss, 2014). As a consequence, unlike HMM/GMM framework, we observe that Lex-SG-Det-ASWU-85 yields better system than SG-Gr-32 in KL-HMM framework.

Table 10: Performance of HMM/GMM and KL-HMM systems in terms of WRR using context-independent (CI) and context-dependent (CD) subword units. For the KL-HMM systems, MLP-SG-91 is used as the acoustic model.

Lexicon	HMM-GMM		KL-HMM	
	CI	CD	CI	CD
Lex-SG-Gr-32	46.0	64.6	35.6	66.8
Lex-SG-Det-ASWU-85	54.5	63.3	52.2	69.1
Lex-SG-Prob-ASWU-82	59.6	66.4	57.5	69.5

5.4. Multilingual ASR system studies

As mentioned earlier, the under-resourced ASR system performance can be improved by using an acoustic model or ANN that classifies multilingual phones and learning a probabilistic relationship between the graphemes and multilingual phones using KL-HMM. We compared the grapheme-based lexicon and the ASWU-based lexicon in that framework by

1. first training a five-layer multilingual MLP on five auxiliary languages from SpeechDat(II) corpus namely British English, Swiss French, Swiss German, Italian and Spanish to estimate posterior probabilities of multilingual phones. The multilingual phoneset was formed by merging the phones that are shared across the aforementioned languages, leading to 117 phone units. We refer to this MLP as MLP-*MULTI*-117; and then
2. training a KL-HMM based ASR system corresponding to each lexicon Lex-*SG*-Gr-32, Lex-*SG*-Det-ASWU-85 and Lex-*SG*-Prob-ASWU-82, as illustrated in Figure 9.

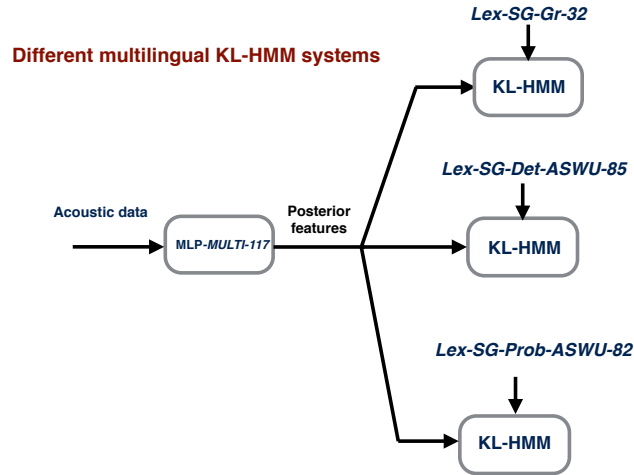


Figure 9: Illustration of KL-HMM based ASR system on Lex-*SG*-Gr-32, Lex-*SG*-Det-ASWU-85 and Lex-*SG*-Prob-ASWU-82 that exploits auxiliary multilingual resources.

Table 11 presents the performance of the different KL-HMM based systems in terms of WRR. It can be observed that the ASWU-based lexicon yields significantly better system than grapheme-based lexicon. Thus, showing that the proposed approach of ASWU-based lexicon development generalizes to multilingual resource sharing scenarios.

Table 11: Performance of KL-HMM based ASR systems exploiting auxiliary resources from resource-rich languages in terms of WRR. In these systems, MLP-*MULTI*-117 is used as the acoustic model.

Lexicon	CI	CD
Lex- <i>SG</i> -Gr-32	36.7	69.1
Lex- <i>SG</i> -Det-ASWU-85	52.1	70.7
Lex- <i>SG</i> -Prob-ASWU-82	57.7	72.6

6. Analysis

The ASR studies validated the proposed ASWU based lexicon from speech technology perspective. As explained in Section 3.1, one of our hypothesis in this article is that the ASWUs obtained from the clustered CD grapheme units are "phone-like". This section focuses on that aspect through an analysis of the derived ASWUs (Section 6.1) and the generated pronunciations (Section 6.2). It is worth mentioning that a full fledged quantitative analysis and concretely linking the derived ASWUs and lexicon to existing linguistic knowledge would need a separate investigation, and is thus out of the scope of the paper. In this section, our main goal is to provide a qualitative analysis and demonstrate how links to existing linguistic knowledge can be established to gain better understanding. We notate phones as / / and graphemes as []. Furthermore, we notate the derived ASWUs with the notation used by HTK to represent clustered CD units. For example, ASWU [ST_A_26] means a clustered CD unit with the center grapheme [A] (root node in the decision tree).

6.1. Relating the derived ASWUs to phonetic units

This section analyzes the relationship between the derived ASWUs and phonetic identities for English and Scottish Gaelic. In the case of English, the analysis uses the acoustic models of the phone-based system, while in the case of Scottish Gaelic there are no phone based lexicon. So the analysis leverages from the ASWU-to-multilingual phone relationship learned by the KL-HMM system presented in Section 5.4.

6.1.1. Studies on English

For both WSJ0 and RM corpora, we computed the KL-divergence between the Gaussian distribution modeling a mono-phone unit and the Gaussian distribution modeling an ASWU in the HMM/GMM setup. We computed the KL-divergence between single Gaussians, as this is the step at which ASWU

is derived by clustering context-dependent graphemes. The KL-divergence between the Gaussian $\mathcal{N}_0(\mu_0, \Sigma_0)$ modeling a mono-phone unit as the reference distribution and the Gaussian $\mathcal{N}_1(\mu_1, \Sigma_1)$ modeling an ASWU as the measured distribution is computed as (Duchi, 2007):

$$0.5\{\text{Tr}(\Sigma_1^{-1}\Sigma_0) + (\mu_1 - \mu_0)^T \Sigma_1^{-1}(\mu_1 - \mu_0) - K - \ln \frac{|\Sigma_0|}{|\Sigma_1|}\},$$

where μ , Σ and K are the mean vector, the covariance matrix and dimension of the vector space respectively.

Table 12 provides a few ASWUs along with the three most related phones according to the KL-divergence matrix. Furthermore, the table also provides example English words which contain the ASWUs within their pronunciations. In each example, the grapheme which has been mapped to the ASWU in the pronunciation is highlighted.

It can be observed from the table that a consistent relationship between the ASWUs and phones exists. This relationship can be clearly observed in the case of consonant graphemes (such as [L], [M], [N] and [R]). For example, the ASWUs belonging to grapheme [L] (such as [ST_L_22] and [ST_L_24] in the WSJ0 corpus) are more related to /el/ and /l/ sounds and the ASWUs belonging to grapheme [R] (such as [ST_R_25] and [ST_R_26] in the RM corpus) are more related to /r/, /axr/, and /er/ sounds. These observations here are also consistent with the empirical observations made in an earlier grapheme-based ASR study on English (Rasipuram and Magimai.-Doss, 2013), where the grapheme-to-phoneme relationship is also learned through acoustics.

6.1.2. Studies on Scottish Gaelic

As mentioned earlier, in the case of Scottish Gaelic there are no phonetic lexicon. So we analyzed the parameters or categorical distributions of the CI KL-HMM system trained with lexicon Lex-SG-Prob-ASWU-82 in the multilingual ASR studies. Table 13 provides examples of mappings between the ASWUs and multilingual phones obtained by selecting the phone with maximum probability in the categorical distribution corresponding to the ASWU. The mapped phones are shown in the SAMPA⁸ format along with the probability of the phone within the brackets. Similar to the analysis on English, we have presented example Gaelic words which contain the ASWUs within their pronunciations.

It can be observed from Table 13 that the ASWUs indeed relate to phonetic

⁸<http://www.phon.ucl.ac.uk/home/sampa/>

Table 12: Relation between example automatically derived subword units and phone units based on the KL-divergence matrix. The example pronunciations are obtained from *Lex-WSJ-Det-ASWU-90* and *Lex-RM-Prob-ASWU-90* for the WSJ0 and RM corpora respectively.

(a) WSJ0 corpus

ASWU	mapped phone	example word	ASWU	mapped phone	example word
[ST_A_26]	/eh/,/ae/,/ey/	DECELER <u>A</u> TION	[ST_L_24]	/l/,/el/,/ao/	INCL <u>I</u> NED
[ST_A_28]	/eh/,/ih/,/ae/	AHE <u>A</u> D	[ST_M_22]	/m/,/em/,/n/	CRAM <u>M</u> ING
[ST_C_21]	/z/,/s/,/zh/	DE <u>V</u> ICE	[ST_N_22]	/ng/,/en/,/n/	RAC <u>I</u> NG
[ST_C_22]	/t/,/dx/,/k/	FORTH <u>C</u> OMING	[ST_N_23]	/n/,/en/,/ng/	REMA <u>I</u> NS
[ST_D_23]	/dx/,/d/,/g/	FOUND <u>A</u> TION	[ST_O_22]	/ow/,/ao/,/aa/	QUO <u>T</u> AS
[ST_E_27]	/ih/,/eh/,/uh/	SE <u>N</u> D	[ST_R_21]	/r/,/er/,/axr/	AMER <u>I</u> CA
[ST_E_28]	/iy/,/y/,/uw/	SE <u>E</u> N	[ST_R_25]	/axr/,/r/,/uh/	ADVERTISE <u>R</u> S
[ST_F_22]	/th/,/f/,/t/	SH <u>I</u> FTE <u>D</u>	[ST_S_21]	/s/,/z/,/f/	ACCOUNT <u>S</u>
[ST_H_23]	/hh/,/dx/,/th/	H <u>A</u> D	[ST_T_21]	/t/,/th/,/dx/	AUST <u>R</u> IA
[ST_I_24]	/iy/,/ey/,/y/	INVENTOR <u>I</u> ES	[ST_U_24]	/uh/,/ax/,/ih/	ACT <u>U</u> AL
[ST_I_27]	/ih/,/uh/,/ax/	J <u>I</u> MMY	[ST_V_21]	/v/,/d/,/dh/	ACHIE <u>V</u> ED
[ST_J_21]	/dx/,/jh/,/t/	J <u>O</u> IN	[ST_W_21]	/w/,/l/,/dx/	AL <u>W</u> AYS
[ST_K_21]	/t/,/dx/,/k/	LOCK <u>E</u> D	[ST_Y_23]	/iy/,/y/,/ih/	AN <u>Y</u> BODY
[ST_L_22]	/el/,/l/,/w/	IMPOSSIB <u>L</u> E	[ST_Z_21]	/z/,/s/,/dx/	Z <u>E</u> US

(b) RM corpus

ASWU	mapped phone	example word	ASWU	mapped phone	example word
[ST_A_211]	/aa/,/aw/,/ay/	CH <u>A</u> RT	[ST_N_21]	/n/,/en/,/ng/	CAMDE <u>N</u>
[ST_A_25]	/ae/,/ey/,/ay/	TR <u>A</u> CK	[ST_O_21]	/ow/,/ao/,/ah/	LOC <u>A</u> TED
[ST_A_26]	/ey/,/eh/,/ae/	DEGR <u>A</u> DE	[ST_O_26]	/ah/,/ow/,/uh/	MON <u>D</u> AY
[ST_B_21]	/d/,/b/,/t/	B <u>A</u> D	[ST_R_25]	/er/,/axr/,/r/	SUMER <u>R</u> IZE
[ST_C_21]	/z/,/s/,/hh/	GAR <u>C</u> IA	[ST_R_26]	/r/,/axr/,/er/	TH <u>R</u> EAT
[ST_D_22]	/dx/,/em/,/d/	ADD <u>I</u> NG	[ST_S_21]	/sh/,/ch/,/s/	WAB <u>A</u> SH
[ST_E_21]	/iy/,/ey/,/uw/	SPE <u>E</u> D	[ST_S_24]	/z/,/s/,/ch/	WADSW <u>S</u> ORTH
[ST_E_25]	/axr/,/er/,/r/	SUR <u>F</u> ACE	[ST_T_21]	/t/,/th/,/dx/	WEST <u>E</u> RN
[ST_F_22]	/f/,/th/,/hh/	VANDERGRI <u>F</u> T	[ST_T_24]	/dx/,/em/,/t/	BET <u>T</u> ER
[ST_H_22]	/hh/,/dx/,/em/	H <u>A</u> D	[ST_U_21]	/ah/,/uh/,/ax/	DO <u>B</u> LE
[ST_H_24]	/dh/,/hx/,/em/	NORTH <u>E</u> RN	[ST_U_22]	/uw/,/ey/,/iy/	T <u>W</u> O
[ST_I_24]	/ih/,/eh/,/uh/	BAINBR <u>I</u> DGE	[ST_W_21]	/w/,/dx/,/em/	W <u>E</u> DNESDAY
[ST_M_21]	/m/,/n/,/ng/	BIS <u>M</u> ARK	[ST_Y_22]	/ih/,/y/,/uw/	AN <u>Y</u> BODY

units in a consistent manner. For example, the ASWU [ST_S_21] is mapped to the phone /S/ (as found in the pronunciation of the English word *SHIP*: /S/ /I/ /p/) and is used in the pronunciation of the Scottish Gaelic word *RIS* which has the slender consonant grapheme [S]. On the other hand, the ASWU [ST_S_23] is mapped to the sound /s/ (as used in the pronunciation of the English word *SKY*: /s/ /k/ /a/ /I/) and is found in the pronunciation of the Gaelic word *THUSA* which contains the broad consonant [S].⁹ Similarly the

⁹Note that in Scottish Gaelic, the broad consonant grapheme [S] is pronounced as the English sound /s/ while the slender [S] is pronounced as the English sound /S/ (web, 2016).

Table 13: Some of the ASWUs together with their mapped phones in SAMPA format and some example words.

ASWU	mapped phone	example word	ASWU	mapped phone	example word
[ST_C.21]	/x/ [0.7]	CACH	[ST_T.21]	/h/ [0.6]	THOG
[ST_C.22]	/C/ [0.7]	SMAOINICH	[ST_T.24]	/t/ [0.7]	MOTA
[ST_C.23]	/k/ [0.9]	CADAL	[ST_G.22]	/g/ [0.5]	GAD
[ST_S.21]	/S/ [0.8]	RIS	[ST_G.23]	/k/ [0.5]	LAG
[ST_S.23]	/s/ [0.8]	THUSA	[ST_R.22]	/r/ [0.4]	MAR
[ST_F.21]	/f/ [0.7]	PHÀIRT	[ST_L.21]	/l/ [0.8]	SAOIL
[ST_B.21]	/b/ [0.5]	BRIS	[ST_L.23]	/l/ [0.5]	SGEUL
[ST_B.22]	/v/ [0.4]	A-BHOS	[ST_Ò.21]	/o/ [0.3]	SPÒRS
[ST_À.21]	/a/ [0.5]	MHÀL	[ST_O.23]	/o/ [0.3]	STOC
[ST_A.212]	/@/ [0.4]	AGAD	[ST_I.23]	/I/ [0.7]	TRIC
[ST_E.21]	/@/ [0.4]	SE	[ST_I.28]	/i/ [0.2]	TRÌ
[ST_E.23]	/l/ [0.3]	WHALES			

consonant ASWUs [ST_F.21] and [ST_R.22] are related to sound units /f/ and /r/. For the vowel ASWUs such as [ST_I.28] and [ST_E.21], the ASWUs are related to the phonetic units, however with a relatively low probability. In our approach, the ASWUs are derived by clustering CD graphemes. So the low probability can be due to the reason that a CD vowel grapheme unit can get mapped to more than one phone, whereas a CD consonant grapheme can have a one-to-one relationship to a phone.

6.2. Generated pronunciations

This section provides a brief analysis on the generated pronunciations through deterministic and probabilistic G2ASWU modeling for English and Scottish Gaelic to get an understanding about the generated pronunciations along with the relation to phonetic identities inferred in the previous section.

6.2.1. English

Table 14 presents a few words selected from ASWU-based lexicons generated for WSJ0 and RM. For each word, the first pronunciation is based on deterministic G2ASWU conversion and the second pronunciation is based on probabilistic G2ASWU conversion. With the information provided in Table 12a and Table 12b, it can be observed that G2ASWU conversion approach is able to recognize different sounds of the same grapheme to provide a pronunciation

similar to what is seen in a phone-based lexicon. For example, in the case of the word *ACCENT*, the grapheme [C] first time is mapped to [ST_C_23], which in the earlier analysis was found to map to phone /k/. Whilst the second time it is mapped to [ST_C_21] in the case of deterministic G2ASWU conversion and [ST_S_25] in the case of probabilistic G2ASWU conversion, in both cases the ASWUs map to /s/. Similar trends can be observed in the example pronunciations provided for the RM corpus. For example, the grapheme [S] is mapped to [ST_S_21] when it corresponds to /sh/ (*FLASHER*) and is mapped to [ST_S_24] when it is related to the /z/ (*PRESENT*). The distinction between deterministic and probabilistic G2ASWU conversion can be very well observed through words *PHONE* and *UPHELD*. In the case of the word *PHONE*, the deterministic G2ASWU conversion maps each grapheme to an ASWU unit while probabilistic G2ASWU conversion is able to map a group of graphemes to an ASWU, i.e. *PH* to /f/ and *NE* to /n/. In the case of the word *UPHELD*, it can be observed that probabilistic G2ASWU conversion leads to deletion of an unit while deterministic G2ASWU preserves the unit. We speculate that the inferior performance of probabilistic G2ASWU conversion in the ASR studies on English is mainly due to such deletions.

6.2.2. Scottish Gaelic

Table 15 presents a few words selected from the ASWU-based pronunciations in case of using deterministic and probabilistic G2ASWU conversion. In order to help in interpreting the generated pronunciations in terms of known sound units, each ASWU in the pronunciation has been mapped to a multilingual phone with the highest probability, as explained in Section 6.1.2. Furthermore, we have provided the ‘perceived’ pronunciations for each word through informal hearing of the Gaelic words. This was done by using an online community-driven dictionary for Gaelic in which for most of the words an audio file pronouncing the word is available.¹⁰

To better understand the generated pronunciations, we first note that in Scottish Gaelic, broad consonants *MH* and *PH* are pronounced as /v/ and /f/, respectively; and the broad consonant *TH* is pronounced as /h/ (web, 2016). It can be seen that the pronunciations obtained through probabilistic lexical modeling based G2ASWU conversion can better capture the linguistic rules compared to the pronunciations obtained through a deterministic lexical mod-

¹⁰<http://www.learnghaelic.net/dictionary/index.jsp>

Table 14: Few example words together with their generated pronunciations based on a deterministic or a probabilistic lexical modeling based G2ASWU conversion on WSJ0 and RM corpora.

(a) WSJ0 corpus.

Word	Lex- <i>WSJ</i> -Det-ASWU-90					
	Lex- <i>WSJ</i> -Prob-ASWU-88					
ACCENT	[ST_A.22]	[ST_C.23]	[ST_C.21]	[ST_E.27]	[ST_N.24]	[ST_T.24]
	[ST_A.22]	[ST_C.23]	[ST_S.25]	[ST_E.27]	[ST_N.24]	[ST_T.24]
ACCORD	[ST_A.22]	[ST_C.23]	[ST_C.22]	[ST_O.21]	[ST_R.23]	[ST_D.21]
	[ST_A.22]	[ST_C.23]	[ST_C.22]	[ST_O.21]	[ST_R.23]	[ST_D.21]
ALAN	[ST_A.22]	[ST_L.24]	[ST_A.27]	[ST_N.21]		
	[ST_A.22]	[ST_L.24]	[ST_A.25]	[ST_N.21]		
ALARM	[ST_A.22]	[ST_L.24]	[ST_A.24]	[ST_R.26]	[ST_M.24]	
	[ST_A.22]	[ST_L.24]	[ST_A.24]	[ST_R.26]	[ST_M.24]	
PHONE	[ST_P.21]	[ST_H.23]	[ST_O.29]	[ST_N.24]	[ST_E.21]	
	[ST_F.22]	[ST_O.29]	[ST_N.21]			
UPHELD	[ST_U.24]	[ST_P.21]	[ST_H.23]	[ST_E.29]	[ST_L.24]	[ST_D.21]
	[ST_O.27]	[ST_P.21]	[ST_H.23]	[ST_L.24]	[ST_D.21]	

(b) RM corpus.

Word	Lex- <i>RM</i> -Det-ASWU-92					
	Lex- <i>RM</i> -Prob-ASWU-90					
CHOP	[ST_C.22]	[ST_H.22]	[ST_O.26]	[ST_P.22]		
	[ST_C.22]	[ST_H.22]	[ST_O.26]	[ST_P.22]		
CODE	[ST_C.23]	[ST_O.26]	[ST_D.22]	[ST_E.24]		
	[ST_C.23]	[ST_O.26]	[ST_D.22]			
FLASHER	[ST_F.22]	[ST_L.23]	[ST_A.21]	[ST_S.21]	[ST_H.22]	[ST_E.25]
	[ST_F.22]	[ST_L.23]	[ST_A.21]	[ST_S.21]	[ST_H.22]	[ST_E.25]
PRESENT	[ST_P.22]	[ST_R.26]	[ST_E.28]	[ST_S.24]	[ST_E.6]	[ST_N.22]
	[ST_P.22]	[ST_R.26]	[ST_E.28]	[ST_S.24]	[ST_I.27]	[ST_N.22]

eling based G2ASWU conversion. For instance, in the word *PHOS* the broad consonant *PH* is mapped to /f/ in the probabilistic lexical modeling based G2ASWU conversion, while in the deterministic approach, it is mapped to /p/ and /h/. Similarly, in the word *MHÀL*, the broad consonant *MH* corresponds to [ST_B.22] which is mapped to the /v/ in the pronunciation obtained from probabilistic G2ASWU relationship modeling, whereas it is mapped to the /v/ and /h/ sounds in the pronunciation generated through deterministic G2ASWU relationship modeling. Indeed, it can be observed that the mapped pronunciations obtained from probabilistic G2ASWU modeling corroborate well with the perceived pronunciations in several cases.

For some of the borrowed English words (e.g., *YOU* and *KATY*), on the other hand, the generated pronunciations using ASWUs seem to be influenced by Gaelic pronunciations. This could be due to a combination of factors such

as, accented English and limited number of English words in the training data.

Table 15: Example words from Scottish Gaelic together with their pronunciations obtained from *Lex-SG-Det-ASWU-91* and *Lex-SG-Prob-ASWU-82*. For each word, we have also provided the mapped pronunciation based on the sequence of multilingual phone units together with its perceived pronunciations.

Word	Lex-SG-Det-ASWU-85 Lex-SG-Prob-ASWU-82	Mapped pron.	Perceived pron.
<i>MHÀL</i>	[ST.M.21] [ST.H.27] [ST.À.21] [S.L.22] [ST.B.22] [ST.À.21] [S.L.23]	/v/ /h/ /a/ /l/ /v/ /a/ /l/	/v/ /a/ /l/
<i>THOG</i>	[ST.T.21] [ST.H.27] [ST.O.23] [ST.G.23] [ST.T.21] [ST.O.23] [ST.G.23]	/h/ /h/ /o/ /k/ /h/ /o/ /k/	/h/ /O/ /g/
<i>PHÒS</i>	[ST.P.21] [ST.H.27] [ST.Ò.21] [ST.S.23] [ST.F.21] [ST.Ò.21] [ST.S.23]	/p/ /h/ /e/ /s/ /f/ /o/ /s/	/f/ /o/ /s/
<i>VOTE</i>	[ST.V.21] [ST.O.23] [ST.T.24] [ST.E.21] [ST.B.22] [ST.O.23] [ST.T.24] [ST.E.21]	/v/ /o/ /t/ /@/ /v/ /o/ /t/ /@/	/v/ /@U/ /t/
<i>YOU</i>	[ST.Y.21] [ST.O.23] [ST.U.22] [ST.I.28] [ST.O.23]	/j/ /o/ /u/ /i/ /o/	/j/ /u:/
<i>KATY</i>	[ST.K.21] [ST.A.212] [ST.T.24] [ST.Y.21] [ST.G.23] [ST.A.212] [ST.T.24] [ST.I.28]	/k/ /@/ /t/ /j/ /k/ /@/ /t/ /i/	/k/ /eI/ /t/ /i/

7. Conclusions

This article presented a novel approach for subword unit derivation and pronunciation generation using only word level transcribed speech data. In this approach, the subword units are first derived by clustering context-dependent graphemes in an HMM-based ASR framework using maximum likelihood criteria; followed by modeling of the relationship between the graphemes and the derived units in a deterministic or probabilistic manner using acoustic data; and finally inferring pronunciations given the learned relationships and the word orthographies using an ergodic HMM. In comparison to existing approaches in the literature, a distinguishing aspect of the proposed approach is that it fits within the well-known HMM framework for ASR and speech synthesis, and is therefore fairly straight-forward to implement given the available toolkits such as HTK (Young et al., 2000) and KALDI (Povey et al., 2011). The proposed approach assumes that a correspondence between the grapheme sequence in the written form of word and the phoneme sequence in the spoken form of the word exists. For logographic languages, where the graphemes represent morphemes or words, the approach could potentially be combined with transliteration.

Our experimental studies on two languages showed that the ASWU-based lexicon can be developed in a fully data-driven manner, i.e. the set of ASWUs and the corresponding lexicon can be selected through cross validation. The ASR studies on both the languages showed that the ASWU-based lexicons consistently yield significantly better ASR systems compared to the grapheme-based lexicons. For G2ASWU conversion, we investigated two approaches, namely, decision-tree based approach and KL-HMM based acoustic G2P approach. Our experimental studies also showed that both G2ASWU approaches are equally applicable, with the acoustic G2P approach holding advantage for languages with many-to-one G2P relationship. Also, in one of the first efforts, we showed that the discovered ASWUs and the learned G2ASWU relationship can be transferred across domains in a language and the G2ASWU conversion mechanism inherently enables such transfer. Furthermore, the analysis of the learned models and the generated pronunciations showed that the derived ASWUs to a good extent are systematically related to phonetic identities. In particular, studies on Scottish Gaelic showed that the multilingual ASR approach not only helps in development of a lexicon that yields better ASR system but also enables discovery of the phonetic identities of the derived ASWUs through the use of multilingual resources. This opens potential venues for further research and development to improve phonetic and lexical resources and technologies for under-resourced languages through transfer of linguistic knowledge and data across languages.

In the proposed approach the problem of ASWU derivation was as posed as a problem of finding a latent symbol space that can be related to acoustic data and associated transcriptions (or graphemes). In this work, we used standard cepstral features that tend to carry information related to phones to find the latent symbol space. However, there are alternative features or representations that carry phone related information and could be exploited to find phone-like latent symbol space. For instance using linguistically motivated articulatory features (AFs) (Jakobson et al., 1992; Ladefoged, 1993), which may be more robust representation when compared to spectral-based features and could help in reducing the gap between ASWU-based approach and phoneme-based approach. This could be achieved without deviating from the HMM framework through the recently proposed AF-based ASR framework using KL-HMMs (Rasipuram and Magimai.-Doss, 2016), where it has been show that ASR systems can be developed by learning grapheme-to-AF relationship through acoustics. Alternately, we could cast the ASWU based lexicon development as a three step process,

where first acoustic-to-AF relationship is learned on available multilingual resources; next grapheme-to-AF relationship is learned from the target language transcribed speech and clustered to derive ASWUs using KL-HMMs; and finally G2ASWU conversion is performed, as done in the present article. Our future work will focus toward this direction on both well resourced and under-resourced languages along with development of methods to select multiple pronunciation variants.

Acknowledgment

This work was supported by Hasler foundation through the grant Flexible acoustic data driven grapheme to acoustic unit conversion (AddG2SU). All the research was conducted at the Idiap Research Institute.

References

- V. Pagel, K. Lenzo, A. Black, Letter to sound rules for accented lexicon compression, in: Proceedings of ICSLP, vol. 5, 2015–2020, 1998.
- T. Sejnowski, C. Rosenberg, Parallel networks that learn to pronounce English text, *Complex systems* 1 (1) (1987) 145–168.
- P. Taylor, Hidden Markov models for grapheme to phoneme conversion., in: Proceedings of Interspeech, 1973–1976, 2005.
- M. Bisani, H. Ney, Joint-sequence Models for Grapheme-to-phoneme Conversion, *Speech Communication* 50 (5) (2008) 434–451.
- S. Kanthak, H. Ney, Context-dependent acoustic modeling using graphemes for large vocabulary speech recognition., in: Proceedings of ICASSP, 845–848, 2002a.
- M. Killer, S. Stüker, T. Schultz, Grapheme based speech recognition, in: Proceedings of Eurospeech, 3141–3144, 2003.
- J. Dines, M. Magimai.-Doss, A study of phoneme and grapheme based context-dependent ASR systems, in: *Machine Learning for Multimodal Interaction*, Springer, 215–226, 2007.
- M. Magimai-Doss, R. Rasipuram, G. Aradilla, H. Bourlard, Grapheme-based Automatic Speech Recognition using KL-HMM, in: Proceedings of Interspeech, 2011.
- T. Ko, B. Mak, Eigentrigraphemes for under-resourced languages, *Speech Communication* 56 (2014) 132–141.
- R. Rasipuram, M. Magimai.-Doss, Acoustic and Lexical Resource Constrained ASR using Language-Independent Acoustic Model and Language-Dependent Probabilistic Lexical Model, *Speech Communication* 68 (2015) 23–40.
- M. Gales, K. Knill, A. Ragni, Unicode-based graphemic systems for limited resource languages, in: Proceedings of ICASSP, 5186–5190, 2015.
- C.-H. Lee, F. K. Soong, B.-H. Juang, A segment model based approach to speech recognition, in: Proceedings of ICASSP, 1988.

- T. Svendsen, K. Paliwal, E. Harborg, P. Husoy, An improved sub-word based speech recognizer, in: Proceedings of ICASSP, 108–111, 1989.
- K. Paliwal, Lexicon-building methods for an acoustic sub-word based speech recognizer, in: Proceedings of ICASSP, 729–732, 1990.
- M. Bacchiani, M. Ostendorf, Using automatically-derived acoustic sub-word units in large vocabulary speech recognition, in: International Conference on Spoken Language Processing, 1998.
- T. Holter, T. Svendsen, Combined optimisation of baseforms and model parameters in speech recognition based on acoustic subword units, in: Proceedings of ASRU, 199–206, 1997.
- R. Singh, B. Raj, R. Stern, Automatic generation of phone sets and lexical transcriptions, in: Proceedings of ICASSP, 1691–1694, 2000.
- C. Lee, Y. Zhang, J. R. Glass, Joint Learning of Phonetic Units and Word Pronunciations for ASR., in: Proceedings of EMNLP, 182–192, 2013.
- W. Hartmann, A. Roy, L. Lamel, J. Gauvain, Acoustic unit discovery and pronunciation generation from a grapheme-based lexicon, in: Proceedings of ASRU, 380–385, 2013.
- M. Razavi, M. Magimai-Doss, An HMM-based formalism for automatic subword unit derivation and pronunciation generation, in: Proceedings of ICASSP, 2015.
- M. Razavi, R. Rasipuram, M. Magimai-Doss, Pronunciation Lexicon Development for Under-Resourced Languages Using Automatically Derived Subword Units: A Case Study on Scottish Gaelic, in: 4th Biennial Workshop on Less-Resourced Languages, 2015.
- G. Aradilla, H. Boulard, M. Magimai-Doss, Using KL-based acoustic models in a large vocabulary recognition task., in: Proceedings of Interspeech, 928–931, 2008.
- X. Luo, F. Jelinek, Probabilistic Classification of HMM States for Large Vocabulary Continuous Speech Recognition, in: Proceedings of ICASSP, 353–356, 1999.

- J. Rottland, G. Rigoll, Tied posteriors: an approach for effective introduction of context dependency in hybrid NN/HMM LVCSR, in: Proceedings of ICASSP, 1241–1244, 2000.
- S. Kanthak, H. Ney, Context-dependent acoustic modeling using graphemes for large vocabulary speech recognition., in: Proceedings of ICASSP, 845–848, 2002b.
- M. Magimai.-Doss, R. Rasipuram, G. Aradilla, H. Bourlard, Grapheme-based Automatic Speech Recognition using KL-HMM, in: Proceedings of Interspeech, 445–448, 2011.
- R. Rasipuram, M. Razavi, M. Magimai.-Doss, Probabilistic Lexical Modeling and Unsupervised Training for Zero-Resourced ASR, in: Proceedings of ASRU, 2013a.
- R. Rasipuram, M. Magimai.-Doss, Acoustic Data-driven Grapheme-to-Phoneme Conversion using KL-HMM, in: Proceedings of ICASSP, 2012.
- M. Razavi, R. Rasipuram, M. Magimai.-Doss, Acoustic data-driven grapheme-to-phoneme conversion in the probabilistic lexical modeling framework, *Speech Communication* 80 (2016) 1–21.
- K. Livescu, E. Fosler-Lussier, F. Metze, Subword Modeling for Automatic Speech Recognition: Past, Present, and Emerging Approaches., *IEEE Signal Processing Magazine* 29 (6) (2012) 44–57.
- C.-T. Chung, C.-A. Chan, L.-S. Lee, Unsupervised discovery of linguistic structure including two-level acoustic patterns using three cascaded stages of iterative optimization, in: Proceedings of ICASSP, 8081–8085, 2013.
- C.-y. Lee, T. J. O’Donnell, J. Glass, Unsupervised lexicon discovery from acoustic input, *Transactions of the Association for Computational Linguistics* 3 (2015) 389–403.
- T. Svendsen, F. Soong, H. Purnhagen, Optimizing baseforms for HMM-based speech recognition, in: Proceedings of EUROSPEECH, 1995.
- A. Jansen, K. Church, Towards Unsupervised Training of Speaker Independent Acoustic Models, in: Proceedings of Interspeech, 1693–1692, 2011.
- A. S. Park, J. R. Glass, Unsupervised pattern discovery in speech, *IEEE Transactions on Audio, Speech, and Language Processing* 16 (1) (2008) 186–197.

- J. Shi, J. Malik, Normalized Cuts and Image Segmentation, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 22 (8) (2000) 888–905.
- A. Y. Ng, M. I. Jordan, Y. Weiss, On spectral clustering: Analysis and an algorithm, in: *Advances in Neural Information Processing Systems*, 849–856, 2001.
- M. Bacchiani, M. Ostendorf, Joint lexicon, acoustic unit inventory and model design, *Speech Communication* 29 (2) (1999) 99–114.
- R. Singh, B. Raj, R. M. Stern, Automatic generation of subword units for speech recognition systems, *IEEE Transactions on Speech and Audio Processing* 10 (2) (2002) 89–99.
- S. Young, The general use of tying in phoneme-based HMM speech recognisers, in: *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, vol. 01, 569–572, 1992.
- A. Ljolje, High accuracy phone recognition using context clustering and quasi-triphonic models, *Computer Speech & Language* 8 (2) (1994) 129–151.
- R. Rasipuram, M. Magimai-Doss, Improving Grapheme-based ASR by Probabilistic Lexical Modeling Approach, in: *Proceedings of Interspeech*, 2013.
- D. B. Paul, J. M. Baker, The Design for the Wall Street Journal-based CSR Corpus, in: *Proceedings of the Workshop on Speech and Natural Language*, 357–362, 1992.
- P. C. Woodland, J. J. Odell, V. Valtchev, S. J. Young, Large Vocabulary Continuous Speech Recognition Using HTK, in: *Proceedings ICASSP*, 125–128, 1994.
- P. Price, W. M. Fisher, J. Bernstein, D. S. Pallett, The DARPA 1000-word Resource Management Database for Continuous Speech Recognition, in: *Proceedings of ICASSP, IEEE*, 651–654, 1988.
- S. Young, D. Kershaw, J. Odell, D. Ollason, V. Valtchev, P. Woodland, *The HTK Book Version 3.0*, Cambridge University Press, 2000.
- D. Johnson, et al., ICSI Quicknet Software Package, <http://www.icsi.berkeley.edu/Speech/qn.html>, 2004.

- R. Rasipuram, M. Magimai.-Doss, Probabilistic Lexical Modeling and Grapheme-based Automatic Speech Recognition, Idiap-RR Idiap-RR-15-2013, 2013.
- D. Imseng, et al., Comparing different acoustic modeling techniques for multilingual boosting, in: Proceedings of Interspeech, 2012.
- M. Bisani, H. Ney, Bootstrap Estimates for Confidence Intervals in ASR Performance Evaluation, vol. 1, 409–412, 2004.
- M. Wolters, A Diphone-Based Text-to-Speech System for Scottish Gaelic, Master’s thesis, University of Bonn, 1997.
- R. Rasipuram, P. Bell, M. Magimai.-Doss, Grapheme and multilingual posterior features for under-resourced speech recognition: a study on Scottish Gaelic, in: Proceedings of ICASSP, 2013b.
- R. Rasipuram, M. Razavi, M. Magimai.-Doss, Integrated Pronunciation Learning for Automatic Speech Recognition Using Probabilistic Lexical Modeling, in: Proceedings of ICASSP, 5176–5180, 2015.
- D. Imseng, R. Rasipuram, M. Magimai.-Doss, Fast and flexible Kullback-Leibler divergence based acoustic modeling for non-native speech recognition, in: Proceedings of ASRU, 348–353, 2011.
- M. Razavi, M. Magimai.-Doss, On Recognition of Non-Native Speech Using Probabilistic Lexical Model, in: Proceedings of Interspeech, 2014.
- J. Duchi, Derivations for Linear Algebra and Optimization, http://www.cs.berkeley.edu/~jduchi/projects/general_notes.pdf, 2007.
- Scottish Gaelic orthography, URL https://en.wikipedia.org/wiki/Scottish_Gaelic_orthography, 2016.
- D. Povey, et al., The Kaldi Speech Recognition Toolkit, in: Proceedings of ASRU, 2011.
- R. Jakobson, G. Fant, M. Halle, Preliminaries to Speech Analysis: the Distinctive Features and their Correlates, MIT Press, 1992.
- P. Ladefoged, A Course in Phonetics, Harcourt Brace College Publishers, 1993.

R. Rasipuram, M. Magimai.-Doss, Articulatory Feature Based Continuous Speech Recognition using Probabilistic Lexical Modeling, *Computer Speech and Language* 36 (2016) 233–259.