

Cross-Native-Language Medieval Latin Dictation

Firstname1 Surname1, Firstname2 Surname2, and Firstname3 Surname3

Affiliation1, Institute1, Address

www.website.org

{author1, author2}@institutel.org

Abstract. We present a medieval Latin charter dictation system which can be of great help for preserving language documents from the same era, since optical character recognition systems often fail to handle historic documents. Our target era and geographical regions are medieval Latin, and documents originating from the Visegrad region, for speakers with Czech, Hungarian and Polish as their native languages. Our baseline systems are separately trained grapheme-based acoustic models for all the above three languages. We introduce two pronunciation modeling techniques to outperform the separately trained models. The first one is using grapheme-to-phoneme (G2P) mapping with Latin-specific pronunciation rules applied. The second one is training a Unified Simplified Grapheme (UGS) acoustic model that can deal with cross-native-language variations. We show that our methods outperform our baseline system, reducing the WER by ...% and ...% respectively.

Keywords: pronunciation modeling, Latin, low resource speech recognition

1 Introduction

Apart from the two official pronunciations of Latin (classical and ecclesiastical), many regional pronunciations exist varying across region and era. The third most known pronunciation group is the East-Central European (ECE) one, used for medieval Latin. Although the target pronunciation is considered to be uniform in this region, it is also has to be taken into account that the acoustic base of the different native languages varies, which can lead to different speakers pronouncing the same words differently. It also has to be noted, that apart from the variations in the pronunciations, orthographic and linguistic variations are also exhibited through regions. This raises the question of how to create a dictation system which has to deal with uniform pronunciations for speakers with different native languages reading linguistically different texts. We propose a system that is suitable for medieval Latin dictation for all speakers from the ECE region. The system we develop is a unified/joint system that can deal with both the variability in the speakers' pronunciations when speaking medieval Latin, and the grammatical/lexical variabilities of the texts. Our baseline system consists of separately trained grapheme (and phoneme) based acoustic models for the different languages in the ECE region. These separately trained models work good with their respective native speakers, but perform poorly with speakers of different native languages. We apply two different pronunciation modeling techniques to develop a models that are superior to

the baseline. The first one, discussed in detail in Section 3.2, is based on the assumption that The second method we use is USGM (Unified Simplified Grapheme Modeling), where a joint/minimal/common grapheme inventory is established for all the languages participating in the joint acoustic model training. We describe this method in Section 3.3.

1.1 Related work

Similar work has been done for multi-dialectal languages such as Arabic in [5] where jointly trained acoustic models were outperformed by methods that unify dialect specific-acoustic models using knowledge distillation and multitask learning.

2 Data

2.1 Textual data

As part of our inquiry was to cover linguistic variability across the ECE region, acquiring textual data posed a few challenges. First of all, textual data are scarce for medieval Latin, and texts originating from the ECE geographical region are even more scarce. Additionally, most of the available sources mix local languages and Latin, with no metadata to separate them. For the scope of this paper, we collected monolingual texts only.

Training data A smaller amount of in-domain data (medieval charters) were collected from [3] (Monasterium), with and overall of 480k tokens. These documents are originating from the Hungarian Kingdom, from 1000 to 1524 AD. To increase the vocabulary size of the language model, we collected a relatively larger (but still small, compared to state-of-the-art language models used in speech recognition) 1.3 token corpus from [4] (LatinLibrary). This corpus consists of literary and historical texts from the post-classical era. In spite of our efforts, at the time of writing this paper, we could not gather textual data from the age and area of the Kingdoms of Bohemia and Poland.

Test data Using independent sources three-three charters were selected from the Kingdoms of Bohemia (CZ), Hungary (HU) and Poland (PL), from around 1200-1300 AD, as development and test data. The dev set was used for evaluating the language model, and the test set to test the performance of our recognizers, by having them read out loud by historians fluent in medieval Latin.

Alternate spellings One interesting feature of the acquired corpora is that they contain a significant number of spelling variants. Having spelling variants in the corpus with identical pronunciation introduces noise, and thus has a negative effect on recognition results. We obtained a unified spelling for these variants by favouring the more frequent variant in the corpus (e.g. *maiestati* to *majestati*). To detect the spelling variants we took all pairs in the pronunciation dictionary whose pronunciation were identical, and used context and expert knowledge to decide whether the pair of equivalent pronunciations

are spelling variants or homophones. Where the decision was that they are spelling variants, the less frequent one was replaced by the more frequent one. Resolving spelling variants resulted in a more consistent corpus in terms of perplexity (reducing it from 775 to 672), and reduced the OOV rate by 0.8%.

Language model The language models we built from the two corpora were estimated with the SRI Language Modeling toolkit (SRILM) [2] using modified Kneser-Ney smoothing method. After estimating the mixture parameter, linear interpolation was used.

The perplexity measures on the dev data showed that the Monasterium corpus originating from the time and era of the Hungarian Kingdom

2.2 Speech data

3 Acoustic modeling

For all the different pronunciation modeling methods, the acoustic models were trained as follows. Mel-Frequency Cepstrum + Energy features were used with Linear Discriminant Analysis (LDA) + Maximum Likelihood Linear Transformation (MLLT), with a splice context of ± 4 frames, 10 ms of frame shift. 9×40 dimensional spliced up feature vectors served as input to the feedforward, 6 hiddenlayer neural network with pnorm [16] activation function. Prior to DNN training, a Gauss Mixture Model (GMM) pre-training was performed. Clustering and Regression Tree (CART) [1] was applied to obtain acrossword context dependent shared state phone (or graph) models and their time alignment. The number of senones (and so the size of the DNN softmax output layer) was between 7.000 and 11.000 depending on the nature of the training data. The size of the hidden layers was kept constantly on 2.000. A minibatch size of 512, an initial learning rate of 0.1, and final learning rate of 0.01 was applied in 20 epochs using the KALDI toolkit [1].

3.1 Grapheme

For our baseline systems grapheme-based pronunciation models were used.

3.2 Grapheme to phoneme mapping (G2P)

Using Latin-specific pronunciation rules

3.3 Unified Simplified Grapheme Modeling

The second method we propose for cross-native-language Latin dictation is Unified Simplified Grapheme (USG) pronunciation modeling technique, which comes in play when joint acoustic models are being trained to support recognition across multiple languages.

Unified These acoustic models need a unified grapheme inventory.

Simplified

4 Experimental results

Table 1. Hungarian phoneme

Speaker	CZ	HU	PL	Avr.
CZ	19.4	6.4	28	17.9
HU	25	25.4	20.2	23.5
PL	47.4	24.6	60.5	44.2
Avr.	30.6	18.8	36.2	28.5

Table 2. USG

Speaker	CZ	HU	PL	Avr.
CZ	20.4	11.8	30.7	21
HU	20.4	14.6	25.7	20.2
PL	54.6	25.4	64.2	48.1
Avr.	31.8	17.3	40.2	29.8

Error analysis

4.1 Conclusions

In this paper, we presented two pronunciation modeling techniques for a cross-native-language medieval Latin dictation system to eliminate the efforts of digitizing medieval Latin charter data. With the objective of outperforming the separately trained grapheme-based models, we presented two approaches: an expert G2P modeling, and UGS modeling. The results showed...

Future research directions include acquiring a considerable amount of medieval speech and textual data.

References

1. Povey, D., Ghoshal, A., Boulianne, G., Burget, L., Glembek, O., Goel, N., Hannemann, M., Motlicek, P., Qian, Y., Schwarz, P., Silovsky, J., Stemmer, G., Vesely, K.: The kaldi speech recognition toolkit. In: IEEE 2011 Workshop on Automatic Speech Recognition and Understanding. IEEE Signal Processing Society (2011)
2. Stolcke, A.: Srlm – an extensible language modeling toolkit. In: IN PROCEEDINGS OF THE 7TH INTERNATIONAL CONFERENCE ON SPOKEN LANGUAGE PROCESSING (ICSLP) 2002. pp. 901–904 (2002)
3. <http://monasterium.net/mom/HU-PBFL/archive>: Monasterium.net archive
4. <http://www.thelatinlibrary.com/medieval.html>: Latin library archive
5. Waters, A., Bastani, M., Elfeky, M.G., Moreno, P., Velez, X.: Towards acoustic model unification across dialects. In: 2016 IEEE Workshop on Spoken Language Technology (2016)