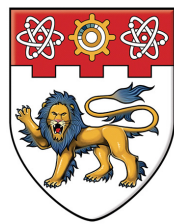


very well written
with depth of analysis
and coherence!
well done!

DOMAIN ADAPTATION OF LANGUAGE MODEL FOR SPEECH RECOGNITION



NANYANG
TECHNOLOGICAL
UNIVERSITY

A Confirmation Report
Submitted to the School of Computer Science and Engineering
of the Nanyang Technological University

by

Yerbolat Khassanov

for the Confirmation for Admission
to the Degree of Doctor of Philosophy

January 7, 2017

Abstract

Acknowledgments

I would like to express my sincere thanks and appreciation to my supervisor Dr. Chng Eng Siong for his invaluable guidance, support and suggestions. His knowledge, suggestions, and discussions help me to become a capable researcher. His encouragement also helped me to overcome the difficulties encountered in my research.

I also want to thank my colleagues in Rolls-Royce@NTU Corporate lab for their generous help. I want to thank Chong Tze Yuang for his generous help to write my first paper and prepare presentation slides. I also want to thank Benjamin Bigot for introducing me to the speech recognition systems.

I am very grateful to the members of our RT1.1 team. It is a pleasure to collaborate with my team mates, Kyaw Zin Tun and San Linn.

Last but not least, I want to thank my family in Kazakhstan, for their constant love and encouragement.

Contents

Abstract	i
Acknowledgments	ii
List of Figures	v
List of Tables	vi
List of Abbreviations	vii
1 Introduction	1
1.1 Motivation	1
1.2 Contributions	6
1.3 Report Organization	7
2 Introduction to Language Model Adaptation for ASR	8
2.1 Background	9
2.1.1 Automatic Speech Recognition	9
2.1.2 Statistical Language Models	10
2.1.3 Domain Mismatch Problem	14
2.2 General LM Adaptation Framework	15
2.2.1 Supervised vs. Unsupervised	15
2.2.2 Cross-domain vs. Within-domain	16
2.2.3 Re-decoding vs. N-best and Lattice Re-scoring	16
2.3 Review of Unsupervised LM Domain Adaptation Techniques	17
2.3.1 Cache-based	18
2.3.2 Topic-mixture	19
2.3.3 Query-based	21
2.4 Summary	22

3	Review of Data Selection	23
3.1	Overview	23
3.1.1	Data availability.	24
3.1.2	Application scenarios.	24
3.1.3	Domain adaptation by data selection.	25
3.2	Data Selection Techniques	25
3.3	Applications	30
3.4	Summary	32
4	LM Adaptation by Data Selection for ASR	34
4.1	Proposed Framework	35
4.1.1	Overview	35
4.1.2	Data Selection	36
4.2	Experiment and Discussion	37
4.2.1	Data	37
4.2.2	The ASR System	38
4.2.3	Experiment Setup and Results	39
4.3	Summary	43
5	Conclusions and Future Work	45
5.1	Contributions	45
5.2	Future Directions	47
5.2.1	Extracting Richer Linguistic Information	47
5.2.2	Domain Tracking	48
	Publication	50
	References	51

List of Figures

2.1	Architecture of automatic speech recognition system.	10
2.2	General LM adaptation framework.	15
2.3	Architecture of cache-based adaption techniques for ASR.	19
2.4	Architecture of topic-mixture based adaption techniques for ASR.	20
2.5	Architecture of query-based adaption techniques for ASR.	21
3.1	Data selection framework.	24
4.1	Proposed LM adaptation framework based on data selection.	35
4.2	WER results obtained by the proposed LM adaptation framework.	40
4.3	Perplexity results of target domain LMs computed on reference data.	40
4.4	WER results for 2-gram feature.	42
4.5	WER results for BOW feature.	42

List of Tables

4.1	TED-LIUM corpus characteristics.	37
4.2	TED-LIUM corpus test set details.	38

List of Abbreviation

AM	Acoustic Model
ASR	Automatic Speech Recognition
BOW	Bag-of-words
CE	Cross Entropy
CED	Cross Entropy Difference
DNN	Deep Neural Networks
fMLLR	Feature Space Maximum Likelihood Linear Regression
GMM	Gaussian Mixture Model
HMM	Hidden Markov Model
IDF	Inverse Document Frequency
KN	Kneser-Ney
LDA	Latent Dirichlet Allocation
LM	Language Model
LSA	Latent Semantic Analysis
LVCSR	Large Vocabulary Continuous Speech Recognition
MFCC	Mel-Filterbank Cepstral Coefficient
ML	Maximum Likelihood
MLLT	Maximum Likelihood Linear Transform
MT	Machine Translation
NER	Named-entity Recognition
NLP	Natural Language Processing
POS	Part-of-speech
PPL	Perplexity
RBM	Restricted Boltzmann Machine
RNN	Recurrent Neural Network
SLM	Statistical Language Model
SMT	Statistical Machine Translation
sMBR	State-level Minimum Bayes risk
TF	Term Frequency
TM	Translation Model
WER	Word Error Rate
WFST	Weighted Finite State Transducers
WWW	World Wide Web

Chapter 1

Introduction

1.1 Motivation

①

A brief history of speech recognition systems. Designing a machine that can mimic complex human behaviors such as understanding spoken language and responding accordingly has been envisioned long before advancement of computers. The major step in fulfilling this vision is to develop automatic speech recognition (ASR) systems which have attracted a substantial amount of effort over the last few decades [1].

What is ASR how related to human needs.

②

Given the complexity of human language, the speech recognition technology evolved gradually. The first speech recognition systems focused on simple tasks such as recognizing numbers. For example, in 1952, Bell Laboratories designed *Audrey* [2] which is a first known and documented speech recognizer. *Audrey* could recognize ten digits spoken isolatedly by a single speaker with an accuracy of 97-99%. In 1962, IBM demonstrated *Shoebox*¹, a system which could recognize sixteen words, including ten digits and six arithmetic operations. Over the next decade, speech recognition technology advanced progressively from a simple machine that can recognize a few words to a sophisticated system that can recognize speech with a large vocabulary. Notably, in 1971, DARPA initiated Speech Understanding Research program which was responsible for Carnegie Mellon's *Harpy* [3] system. *Harpy* could recognize speech using a vocabulary of 1,011 words, approximately the vocabulary of an average three-year-old.

discuss, talk about the History of ASR. discussion in chronological order.

③

In these large vocabulary systems, however, the complexity of the task had considerably increased, particularly the confusion attributed to homophones. For example, the

What are the problems.

¹www-03.ibm.com/ibm/history/exhibits/specialprod1/specialprod1.7.html

Focus on one that relate to the thesis.

words ‘buy’, ‘bye’ and ‘by’ comprise same phoneme sequence ‘B AY’ (based on ARPA-bet² phoneme set). Distinguishing such words was an infeasible task for the early speech recognition systems that mainly relied on acoustic information. Thus, the recognition capability of large vocabulary systems was limited.

4 Introduction of language models in ASR. The use of only acoustic information proved to be insufficient to achieve human-like performance. Hence, other sources of knowledge were required. ~~Therefore,~~ ¹ in 1975, Jelinek et al. [4] proposed to incorporate a grammar structure of the natural language into the speech recognizer. The grammar structure was encoded into a language model (LM) based on statistical principles. The function of the statistical LM was to encapsulate syntactic, semantic, and pragmatic properties of the language considered. In the speech recognition system, the encapsulated knowledge was used to constrain search in a decoder by limiting the number of possible words to follow at any one point. The consequence was faster search and higher recognition accuracy. Since then, ~~the~~ statistical LMs have become an indispensable part of large vocabulary speech recognition systems. We will provide a thorough ~~explanation~~ ^{discussion} of the state-of-the-art statistical LMs in chapter 2.

5 The domain mismatch problem. The statistical LMs ~~retain~~ ^{encapsulate} knowledge in the form of probability distribution of linguistic units (e.g. words, sentences) learned from textual training data [5]. It is desirable for ~~this~~ ^{the} training data to possess characteristics ^{similar to} input utterances submitted to the ASR system. For example, covering similar topics, speaking styles or both. Otherwise, the distribution learned by LM ~~might~~ ^{will} mismatch with the ~~target~~ ^{the target} domain distribution of input utterances. ~~As a result,~~ ^{the} the ASR output will be corrupted [5]. For example, a LM trained on industry domain data, but applied to input utterances from the math domain, ~~might~~ ^{may} confuse the ASR to recognize ‘COFACTOR IS’ as ‘COW FACTORIES’ (the ~~Hamming~~ ^{Hamming} distance between phoneme sequence of these phrases is 1). Therefore, for reliable performance of ASR ~~system~~ ^{the} the distribution learned by LM should ^{match} fit the target domain.

In ASR systems, however, maintaining a LM that fits the distribution of input test data is a challenging task. Specifically, in the cases where input utterances cover several

²<http://www.speech.cs.cmu.edu/cgi-bin/cmudict>

previous example illustrates problem when only acoustic data is exploited. Hence LM is important.

stat. description LM.

What is the domain mismatch problem? the text should be similar to the target domain.

5A

Lecorvé et al. [7] used only incorrectly recognized parts of the ASR output to perform adaptation. Surprisingly, they obtained more than 10% relative perplexity improvement. They concluded that some ~~misrecognized words are still in domain words aiding to capture appropriate domain information and others are harmless.~~ ^{may} Thus, despite the presence of errors, the ASR output contains valuable information which can be effectively utilized by adaptation techniques.

A brief review of existing adaptation techniques. Since the introduction of statistical LMs, several LM domain adaptation techniques have been proposed to alleviate the effect of distribution mismatch [8]. In practice, LMs can be adapted in two different stages of recognition process: online and offline. In online adaptation, the LM is adapted during the decoding process. However, the decoding process itself is a highly complex mechanism involving intensive computations which make the online LM adaptation impractical. In offline adaptation, on the other hand, the generic LM is first applied to produce initial ASR output (word lattice). Then, ^{the} ~~produced~~ ASR output is utilized to ^{predict} ~~generate~~ target domain information, in ^{either} supervised or unsupervised manner, ^{then} which is employed to adapt the generic LM offline. Lastly, the adapted LM ^{can be} is applied to re-decode the input utterances, or to re-score the word lattice (or N-best list). ^{either}

9A Given the complexity of decoding process, re-decoding the input utterances is a tedious task. Hence, only a few LM types, for which fast decoding algorithms are available, are eligible for this task, ^{e.g.} such as backoff n-gram models [10, 11]. The backoff n-gram model is a predominant choice for decoding in the state-of-the-art ASR systems due to its effectiveness and simplicity [12] (generic LM is a backoff n-gram model). Whereas, ^{for} more complex ^{LMS} models, such as neural network based [13], are usually employed to re-score the word lattices [9]. While the complex models are expected to have greater prediction power, the efficacy of re-scoring is constrained by the quality of generated word lattice which contains only subset of all possible hypotheses. For example, an inadequate LM used during the decoding stage might ^{have already} discard correct hypotheses, as a result, a deficient word lattice will be produced [14]. Hence, in this work, we will ^{first} focus on adapting backoff n-gram models which can be effectively applied to re-decode the input utterances.

10 The three popular ~~backoff~~ ^{By ?} n-gram model adaptation techniques applied to ASR systems are cache-based, topic-mixture and query-based. These techniques employ domain-specific information to tune distribution of generic LM so that it better matches the target domain. For example, the cache-based techniques [6, 15–18] are based on the hypothesis that a word used in a recent past is more likely to be used again. Hence, the probabilities of recognized words are increased within the generic LM. In the topic-mixture techniques [6, 18–21], the ~~generic~~ ^{adapted} LM is ~~dismantled~~ ^{formed by interpolation} into several sub-domain (or sub-topic) LMs ~~interpolated together~~. Here, the domain of the final interpolated LM ~~can be controlled~~ ^{is formed} by tuning interpolation weights of the sub-topic LMs. ~~Hence, the ASR output is used to find the ‘closest’ sub-topic LMs to increase their weights.~~ ^{The strategy is to estimate the weights from} The query-based techniques [7, 22–24], on the other hand, use ASR output to generate queries which are submitted to external sources, such as world wide web (WWW), to retrieve ‘similar’ data. The retrieved data is then used to update parameters of the generic LM. For example, by training new ‘pseudo’ in-domain LM from the retrieved data and interpolating it with the generic LM. These LM adaptation techniques have been shown effective to improve recognition performance of ASR systems. The complete review of these techniques will be given in chapter 2.

11 **The proposed adaptation approach based on data selection.** The existing adaptation techniques typically adjust the distribution learned by generic LM to match the target domain distribution. This adjustment is performed by directly changing parameters of the generic LM, for instance, by increasing or decreasing probabilities of individual words (or n-grams). ~~Changing~~ ^{Although} parameters of LM ~~might help to achieve desired distribution,~~ ^{the way} ~~however,~~ ^{the} the adapted LM ~~most probably won’t~~ ^{will not} represent a distribution corresponding to the natural text produced by human. ~~Consequently, the encapsulated knowledge might be corrupted.~~ Thus, in this work, rather than directly updating parameters of the generic LM, we will examine other adaptation methods that preserves the ‘natural’ distribution of linguistic units.

11A In particular, we propose to manipulate the training data used to build generic LM. As was mentioned previously, the training data consist of text assembled from various domains. Hence, we will employ data selection techniques [25] to select a subset of

training data ‘similar’ to the ASR output (broad overview of data selection techniques will be exposed in chapter 3). *The strategy is to discard* As a result, out-of-domain sentences *while retaining* ~~will be discarded~~, leaving *only* in-domain sentences. *set of formal* The in-domain sentences are then used to train new LM which is expected to better converge *to* with the target domain distribution. In addition, the new LM *is* ~~will represent~~ *an* adapted version of *the* generic LM, since it was build from the same, but pruned data. More importantly, the adapted LM produced in this way will encapsulate *the* appropriate linguistic knowledge which complies with the regularities of natural language. *as complete sentences or paragraphs are retained during data selection.* To evaluate the effectiveness of proposed approach we conducted several experiments on TED-LIUM speech corpus which will be described in chapter 4.

1.2 Contributions

In this thesis, we proposed unsupervised LM adaptation framework to address domain mismatch problem inherent in generic ASR systems. The proposed framework is based on data selection technique which customizes generic background corpus to produce domain-specific LM. The novelties of the proposed framework are listed below:

1) Existing LM adaptation techniques aim to tune parameters of the generic model to diverge its focus towards target domain. Different from them, the proposed approach employs ASR output and data selection techniques to perform adaptation at the data level. This work shows that a LM adapted in this way possesses a strong discriminative ability that results in substantial WER reduction.

2) Although the generic background corpus is sufficiently large and contains data from various domains, several adaptation techniques (e.g. query-based) still require in-domain data retrieved from external sources such as WWW. Unlike these approaches, our method efficiently utilizes available background corpus by intelligently selecting in-domain sentences. Hence, the proposed method doesn’t rely on any external source which might be unavailable for some tasks involving private corporation or military domains.

Experiments performed on TED-LIUM speech corpus show that proposed adaptation framework can produce domain-specific LM that achieves up to 10% relative WER reduction. When we adapted LM to a more specific domain the WER reduction up to 12% was observed. Moreover, we compared our approach against standard adaptation

method based on linear interpolation which directly updates parameters of a LM, and observed better WER.

The work on unsupervised LM adaptation by data selection was accepted by ACIIDS conference [26].

1.3 Report Organization

The report is organized as follows:

In Chapter 2, we provide background information on ASR systems, statistical LMs, and domain mismatch problem. We describe general LM adaptation framework, followed by a review of popular LM adaptation techniques applied to ASR systems.

In Chapter 3, we provide an overview of the current state-of-the-art data selection techniques including linguistic features used to represent data and similarity metrics. We briefly review other natural language processing (NLP) applications where data selection has been employed.

In Chapter 4, we propose data selection based unsupervised LM adaptation framework for ASR systems. We explain experiment setup and data. Lastly, obtained results are discussed.

Chapter 5 concludes the report and lists future research directions.