



INSTITUT FÜR LOGIK, KOMPLEXITÄT
UND DEDUKTIONSSYSTEME
UNIVERSITÄT KARLSRUHE

AM FASANENGARTEN 5
D-76128 KARLSRUHE

Datengetriebene Bestimmung von
Vokabulareinheiten für koreanische
Spracherkennung auf großen Wortschätzen

Diplomarbeit von
DANIEL KIECZA

Betreuer:
Prof. Dr. Alex Waibel
Tanja Schultz



angefertigt am
Computer Science Department
Interactive Systems Laboratories
Carnegie Mellon University
Pittsburgh, PA 15213, U.S.A.

daniel@kiecza.de

October 1999

Hiermit erkläre ich, daß ich die vorliegende Arbeit selbständig erstellt und keine anderen als die angegebenen Quellen und Hilfsmittel verwendet habe.

Pittsburgh, den 21.10.1999

Daniel Kieczka

Zusammenfassung

Koreanisch zählt zu den agglutinierenden Sprachen. Verwendet man die aus dem Agglutinationsprozeß entstehenden Einheiten – genannt Eojeols – als Wörterbucheinträge für ein Spracherkennungssystem, wächst die Größe des Vokabulars etwa linear mit der Größe des vorliegenden Textmaterials. Außerdem erreicht die out-of-vocabulary (OOV) Rate Größenordnungen, die ein leistungsfähiges Spracherkennungssystem unmöglich machen. Es werden folglich geeignete sub-Eojeol-Einheiten als Wörterbucheinträge benötigt.

Bisherige Arbeiten verwenden hierzu aufwendige Expertensysteme, die die Eojeols in ihre Morphemkomponenten zerlegen. Diese Morphemeinheiten werden dann als Wörterbucheinträge verwendet.

In der vorliegenden Arbeit wird ein neues Verfahren zur Findung geeigneter Einheiten vorgestellt, das ausschließlich *datengetrieben* arbeitet. Ausgangspunkt ist Textmaterial, in dem sämtliche Eojeols in ihre Silbenkomponenten zerlegt sind. Dann werden wiederholt bestimmte Paare von Einheiten zusammengefügt. Die Wahl der zu verbindenden Paare wird dabei so gefällt, daß die akustische Verwechselbarkeit der Einheiten reduziert wird.

Die Erkennungsleistung der resultierenden Systeme wird präsentiert und mit der Erkennungsleistung entsprechender morphembasierter Systeme verglichen. Das beste, auf dem datengetriebenen Ansatz basierende System erreicht eine Eojeol-Fehlerrate von 24.6%. Das entspricht einer Silben-Fehlerrate von 14.5% und einer Phonem-Fehlerrate von 9.9%. Dies vergleicht sich mit einer Eojeol-Fehlerrate von 24.0%, einer Silben-Fehlerrate von 13.0% und einer Phonem-Fehlerrate von 9.4% für das beste morphembasierte System. Die beiden Ansätze bringen vergleichbare Leistung, und man kann folglich auf die aufwendige Konstruktion von Morphemzerlegungssystemen verzichten.

Abstract

Korean is an “agglutinating” language, i.e. verbs and nouns consist of a stem and various appended suffixes which have grammatical function. Using the units that result from the agglutination process, called *eojeols*, as dictionary entries for a speech recognition system makes the vocabulary size grow linearly in the task size. Furthermore, an extremely high out-of-vocabulary (OOV) rate has to be dealt with which makes the development of a high performance speech recognition system almost impossible. A solution to this problem is to work with sub-*eojeol* units.

Previous work has used complex expert systems that split *eojeol* units into their morpheme components. These morpheme components are then used as dictionary units.

This work presents a new *data-driven* approach to determine appropriate dictionary units. The approach starts with a text corpus where the *eojeols* are split up into their syllable components. Then it repeatedly merges certain pairs of units. The choice of the unit pair to merge is done so as to reduce acoustic confusability of the units.

The recognition performance of the resulting systems is presented and is compared to the performance of morpheme based recognition systems. The best data-driven system we present has an *eojeol* error rate of 24.6%. This corresponds to a syllable error rate of 14.5% and a phone error rate of 9.9%. The best morpheme based system has an *eojeol* error rate of 24.0% corresponding to a syllable error rate of 13.0% and a phone error rate of 9.4%. Thus, the two approaches show comparable performance yet we did not need to rely on expert knowledge to generate an *appropriate* set of dictionary units.

Acknowledgements

This work would not have been possible without the help from a number of people.

I would like to thank Professor Alex Waibel for the opportunity to conduct my research at the Interactive Systems Laboratory at Carnegie Mellon University in Pittsburgh.

A big thank you to the supervisor of my thesis, Tanja Schultz, for her advice and guidance during this project. Although our physical separation (me being in Pittsburgh and Tanja being in Karlsruhe) did not make the task of supervision easy we managed to cooperate and communicate effectively and had many very interesting discussions on the telephone and by email.

I also want to mention the many fruitful and inspiring discussions I had with Michael Finke. There, he often suggested solutions when I still tried to understand the problem. These discussions helped me to develop a deeper understanding of the problems involved in this project.

Thank you to Detlef Koll. Sharing his insight was of great pleasure to me and his ideas and suggestions were always very useful.

Then I would like to thank Iain Matthews and Victoria Maclaren for proof-reading and *de-germishifying* this thesis and for helping me to improve my English.

Thank you to Oh-Wook Kwon from the ETRI labs in Seoul, Korea for providing the morpheme segmentation of our text data and also for several interesting discussions via email.

Furthermore, thank you to Monika Woszczyna and Gang-Seong Lee for interesting discussions on aspects related to this work.

Thanks to our assistants at the Interactive Systems Labs in Pittsburgh and Karlsruhe, Debbie Clement, Nadine Reaves and Silke Dannenmaier who always helped to solve administrative issues quickly.

As every computer scientist will understand, I especially want to thank our system administrators Frank Dreilich, Weiyi Yang and Eric Carraux for supporting me in my constant battle with *the hardware*.

Thanks to Sang-Hun Shin, Keal-Chun Cho and Kyung-Kyu Lee for their enthusiasm during collection and validation of the database.

Last but not least, I want to thank my parents who supported me in every possible way during my whole education, especially during the time I spent in Pittsburgh.

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 1 |
| 1.1 | Motivation | 1 |
| 1.2 | Problem Definition | 1 |
| 1.3 | Contribution of this Work | 3 |
| 1.4 | Organization of the Thesis | 3 |
| 2 | The Korean Language | 4 |
| 2.1 | Historical Remarks | 4 |
| 2.2 | 한글 (Han-geul) – The Korean Writing System | 6 |
| 2.3 | Syllable Construction | 9 |
| 2.4 | Phonetics and Phonology | 11 |
| 2.5 | Romanisation | 15 |
| 3 | Speech Recognition | 17 |
| 3.1 | Motivation | 17 |
| 3.2 | Overview | 18 |
| 3.3 | Front-End | 19 |
| 3.4 | Language Modeling | 21 |
| 3.5 | Acoustic Modeling | 22 |
| 3.6 | Decoding | 24 |
| 3.7 | Speaker Adaptation | 25 |
| 3.8 | The Janus Speech Recognition Toolkit | 25 |
| 4 | System Description | 27 |
| 4.1 | Speaker Database | 27 |
| 4.2 | Language Model Data | 28 |
| 4.3 | Pronunciation Generation | 32 |
| 4.4 | HMM Recognizer Structure | 33 |

| | | |
|----------|--|-----------|
| 4.4.1 | Phone Set | 33 |
| 4.4.2 | Speech Preprocessing | 34 |
| 4.4.3 | Context Dependent Phone Modeling | 35 |
| 4.4.4 | Model Initialization | 35 |
| 5 | Morpheme Based Recognition | 36 |
| 5.1 | Motivation | 36 |
| 5.2 | Determination of Units | 37 |
| 5.3 | Speech Recognition Systems | 37 |
| 6 | Data-Driven Unit Determination | 43 |
| 6.1 | Motivation | 43 |
| 6.2 | Determination of Units | 44 |
| 6.2.1 | Preprocessing | 44 |
| 6.2.2 | Unit Merging | 45 |
| 6.3 | Speech Recognition Systems | 47 |
| 7 | Experiments | 54 |
| 7.1 | Introduction | 54 |
| 7.1.1 | Recognition Accuracy | 54 |
| 7.1.2 | Lattice Rescoring | 55 |
| 7.2 | Context Modeling | 57 |
| 7.3 | Results And Discussion | 57 |
| 7.3.1 | Baseline | 57 |
| 7.3.2 | Pronunciation Variant LM | 60 |
| 7.3.3 | More Data | 61 |
| 7.3.4 | Corrected Speech Database | 62 |
| 7.3.5 | Discussion | 63 |
| 8 | Conclusions | 65 |
| 8.1 | Conclusions | 65 |
| 8.2 | Future Work | 66 |
| A | Transcription Systems for 한글 | 67 |
| B | Text Corpus Mappings | 69 |

| | |
|---|-----------|
| <i>CONTENTS</i> | III |
| C Janus Toolkit | 71 |
| C.1 Phoneme Models | 71 |
| C.2 Phoneme Context Question Sets | 73 |
| D Coverage | 74 |

List of Figures

| | | |
|-----|--|----|
| 3.1 | Example of a left-to-right hidden Markov model. | 24 |
| 4.1 | Self coverage of language model corpora. | 30 |
| 4.2 | Cross coverage of <i>Test</i> with different language model corpora. | 31 |
| 4.3 | Vocabulary growth in corpus <i>Chosun+Train</i> | 32 |
| 5.1 | Self coverage of <i>Morph</i> -based language model corpora. | 39 |
| 5.2 | Self coverage of <i>MorphTag</i> -based language model corpora. | 40 |
| 5.3 | <i>Morph</i> -based cross coverage of <i>Test</i> with different language model corpora. | 40 |
| 5.4 | <i>MorphTag</i> -based cross coverage of <i>Test</i> with different language model corpora. | 41 |
| 5.5 | Morpheme based vocabulary growth in corpus <i>Chosun+Train</i> | 42 |
| 6.1 | Vocabulary growth for the four merge based systems in corpus <i>PartChosun+Train</i> | 49 |
| 6.2 | <i>MergeIntraMax</i> -based self coverage of <i>PartChosun+Train</i> and cross coverage of <i>Test</i> and <i>Test-Utts</i> | 50 |
| 6.3 | <i>MergeIntraAll</i> -based self coverage of <i>PartChosun+Train</i> and cross coverage of <i>Test</i> and <i>Test-Utts</i> | 50 |
| 6.4 | <i>MergeInterMax</i> -based self coverage of <i>PartChosun+Train</i> and cross coverage of <i>Test</i> and <i>Test-Utts</i> | 51 |
| 6.5 | <i>MergeInterAll</i> -based self coverage of <i>PartChosun+Train</i> and cross coverage of <i>Test</i> and <i>Test-Utts</i> | 51 |
| 6.6 | Vocabulary growth for <i>MergeIntraMaxWhole</i> based system in corpus <i>Chosun+Train</i> | 53 |
| 6.7 | <i>MergeIntraMaxWhole</i> -based self coverage of <i>Chosun+Train</i> and cross coverage of <i>Test</i> and <i>Test-Utts</i> | 53 |

| | | |
|-----|---|----|
| D.1 | Eojeol-based cross coverage of <i>Test-Utts</i> with different language model corpora. | 74 |
| D.2 | <i>Morph</i> -based cross coverage of <i>Test-Utts</i> with different language model corpora. | 75 |
| D.3 | <i>MorphTag</i> -based cross coverage of <i>Test-Utts</i> with different language model corpora. | 75 |
| D.4 | <i>MergeIntraMax</i> -based self coverage of <i>Train</i> and cross coverage of <i>Test</i> and <i>Test-Utts</i> | 76 |
| D.5 | <i>MergeIntraAll</i> -based self coverage of <i>Train</i> and cross coverage of <i>Test</i> and <i>Test-Utts</i> | 76 |
| D.6 | <i>MergeInterMax</i> -based self coverage of <i>Train</i> and cross coverage of <i>Test</i> and <i>Test-Utts</i> | 77 |
| D.7 | <i>MergeInterAll</i> -based self coverage of <i>Train</i> and cross coverage of <i>Test</i> and <i>Test-Utts</i> | 77 |
| D.8 | <i>MergeIntraMaxWhole</i> -based self coverage of <i>Train</i> and cross coverage of <i>Test</i> and <i>Test-Utts</i> | 78 |

List of Tables

| | | |
|-----|---|----|
| 2.1 | The three basic vowels. | 6 |
| 2.2 | Primary vowels. | 7 |
| 2.3 | Secondary vowels. | 7 |
| 2.4 | Compound vowels. | 8 |
| 2.5 | The five basic consonants and their extensions. | 9 |
| 2.6 | Original and actual form of the consonant symbols and their romanisation. | 10 |
| 4.1 | Summary of acoustic database. | 27 |
| 4.2 | Mapping table for acronyms. | 29 |
| 4.3 | Summary of language model corpora. | 29 |
| 4.4 | Summary of OOV rates and, in parantheses, OOV words. . . . | 31 |
| 5.1 | Summary of morpheme based language model corpora. | 38 |
| 5.2 | Summary of OOV rates and, in parantheses, OOV words. . . . | 41 |
| 6.1 | Number of character pairs that were merged. | 47 |
| 6.2 | Summary of merge based language model corpora. | 48 |
| 6.3 | Summary of OOV rates and, in parantheses, OOV words. . . . | 48 |
| 6.4 | Summary of characteristics of <i>MergeIntraMaxWhole</i> corpora, and summary of OOV rates and, in parantheses, OOV words. . . . | 52 |
| 7.1 | Recognition task complexities on <i>Chosun+Train</i> | 58 |
| 7.2 | Summary of recognition error rates, %. | 59 |
| 7.3 | Recognition task complexities on <i>Chosun+Train</i> , pronunciation variants not mapped on baseform in language model. . . . | 60 |
| 7.4 | Summary of recognition error rates using a variant-LM, %. . . | 61 |
| 7.5 | Recognition task characteristics. | 62 |
| 7.6 | Summary of recognition error rates on <i>Chosun+Train</i> , %. . . . | 62 |

| | | |
|-----|--|----|
| 7.7 | Summary of recognition error rates on <i>Chosun+Train</i> with corrected speech database, %. | 63 |
| 7.8 | Summary of the performance of the recognition systems. | 64 |
| A.1 | Transcription systems for the 한글 consonants. | 67 |
| A.2 | Transcription systems for the 한글 vowels. | 68 |
| B.1 | Some examples of acronym mappings. | 69 |
| B.2 | Mapping table for units. | 70 |
| C.1 | Korean vowel models. | 71 |
| C.2 | Korean consonant models. | 72 |
| C.3 | Phone sets used for decision tree based context dependent phone modeling. | 73 |

Chapter 1

Introduction

1.1 Motivation

Speech is the primary mode of communication among humans. In order to make the interaction between humans and computers more user-friendly, speech must be considered an essential input component.

Over the last two decades, intensive research in the field of automatic speech recognition and a significant increase in available and affordable computing power have led to practical solutions to this problem. Current systems achieve a recognition accuracy between 90% and 100% under certain restricted conditions. Examples of these conditions are a limited vocabulary size in a well defined task domain, a good-quality microphone and a non-noisy environment. One focus of current speech research is to eliminate these limitations.

A great deal of effort is spent developing high performance speech recognition systems for large vocabulary recognition tasks. Each new language that is explored for large vocabulary speech recognition has its own characteristics which can impose difficulties to the development of a high-performance recognition system. This thesis describes our efforts to develop a large vocabulary speech recognition system for the Korean language.

1.2 Problem Definition

For most Western languages, e.g. English, the *word* is an appropriate choice of vocabulary unit. But for Korean, the choice is much more difficult. Korean

is an “agglutinating” language. The structure of verbs and nouns is subject to the following rules: A verb consists of a stem and various appended morphemes that have certain grammatical functions, i.e. indicate tempus, modality or social relationship between the communication partners. The syntactic role of a noun is indicated by appending case suffixes to its stem. We call the resulting “meta”-units *eojeols*¹.

Using *eojeols* as vocabulary units² makes the size of the vocabulary grow linearly in the task size. Furthermore, this leads to an extremely high out-of-vocabulary (OOV) rate. This makes the development of a high performance large vocabulary continuous speech recognition (LVCSR) system almost impossible.

Eojeols are built from only about 3600 different syllables³. Each such syllable consists of one to four phonemes. Using these syllables as vocabulary units provides a small vocabulary and the OOV rate is 0%. Unfortunately, due to their shortness, two problems arise when this approach is used:

- acoustic confusability of syllable units is very high,
- a standard 3-gram language model has very limited scope.

To overcome the difficulties of using either *eojeols* or syllables we have to find a more “appropriate” set of vocabulary units that lie inbetween these two extremes.

Recent work on Korean LVCSR makes use of expert systems which are commonly referred to as morpheme tagging systems. These split *eojeol* units into their morpheme components. As presented in [33], the resulting morpheme units can indeed be used as an appropriate set of vocabulary units for a Korean speech recognition system. Unfortunately, it requires a great deal of effort to develop such a morpheme analyzing system as it involves a lot of *apriori* expert knowledge about the morphological structure of the Korean language.

¹Pronunciation: [əjəl].

²We use the same set of units in the vocabulary, the dictionary and the language model. These notions are used synonymously throughout this work.

³Synonymously called *characters*.

1.3 Contribution of this Work

This work presents a new *data-driven* approach to determine appropriate vocabulary units for Korean LVCSR. The morphemic structure of the Korean language is ignored for the unit determination. Instead, the problem is approached from an acoustically motivated side. Roughly speaking, the data-driven procedure works as follows: In a first pass, each eojeol is split into its character components. Then, repeatedly, syllable pairs are merged in order to reduce acoustic confusability of the phone transitions between units. A vocabulary size of 64k units is used as termination criterion.

The recognition performance of the resulting systems is presented and is compared to the performance of morpheme based recognition systems. This demonstrates whether appropriate vocabulary units for Korean LVCSR may be determined using data-driven methods, without extensive use of *a priori* expert knowledge.

1.4 Organization of the Thesis

This thesis is organized as follows: chapter 2 describes the structural elements of the Korean language which are relevant to this work. In chapter 3 we give a short overview of state-of-the-art speech recognition. The intention is not to give an extensive introduction but to present key concepts and important notations that are needed as a foundation for this thesis. In chapter 4 we describe the speech and text data we used. In addition, we explain how automatic dictionary generation is done and present the baseline structure of our speech recognition system. Chapters 5 and 6 describe the morpheme based approach and our new data-driven approach, respectively. In chapter 7 we present and discuss the recognition performance results of the systems. Finally, chapter 8 draws conclusions and suggests ideas for future work on Korean speech recognition.

Chapter 2

The Korean Language

2.1 Historical Remarks

A common hypothesis is that the Korean language belongs to the Altaic family [18]. Other languages in this family include Manchu, Mongolian and Turkish. But other theories consider Korean an isolated language. It is not clear how long the Korean language has been spoken. Despite geographic proximity of the countries, Korean is quite unlike Chinese and Japanese.

The Chinese writing system and the Chinese culture came to Korea about two thousand years ago. The Chinese influence was very strong and the existing cultural tradition was widely repressed. It is unknown what writing systems existed in Korea before the Chinese era.

The Chinese system was used in Korea until the 15th century. Due to its complexity and difficulty to learn, most people remained illiterate. People from the upper class were able to use the “foreign” Chinese system because of their education, but there is evidence that even they found it difficult. This was due to the significant differences of the two languages, concerning their phoneme inventories as well as their sentence structure.

The Chinese language consists of innumerable single characters which are pronounced monosyllabically and have a morphemic function. They can be pronounced using different tones in order to distinguish homophonic characters. The language is said to be *tonal*. In contrast, Korean is a *syllabic*

language, i.e. words are generated by combination of syllables.

After unsuccessful attempts to solve this conflict, King Se-Jong (reign 1419 – 1450) initiated the development of a totally new writing system for the Korean language in the 15th century. In 1420, a Royal institute was established to develop a Korean alphabet system. Thirteen years later these scientists presented their results to King Se-Jong. In 1446 – after this language system had been tested for three years – the King introduced it in a publication known as *훈민정음* (Hun-min Jeong-eum)¹. A book was published with the same name, containing the background on this system, the reasons for its creation, information on its usage and so forth. King Se-Jong wrote in the book: “Our language is unique, different from that of China therefore we needed an alphabet of our own. Only scholars or people from the upper class get the opportunity to read or write. I felt it is my responsibility to enlighten people as a King so I have made a new 28 letters and I wish everyone can learn this new alphabet easily and use these new letters comfortably in daily life”.

The *훈민정음* alphabet was built on a purely phonetic basis and consisted of 28 letters. It has evolved into the modern Korean alphabet, *한글* (han-geul), which has 24 letters. Although Chinese characters are still used today in conjunction with the Korean alphabet, especially in newspapers, high-level communication in Korea is possible without using them.

That is why the Korean language often offers two expressions for a notion, a Chinese one and a pure Korean one. The Chinese characters that are used in Korean evolved over time, especially concerning their pronunciation. Their writing form remained mostly unchanged since the Han dynasty (206 BC – 220 AD).

For over 400 years the new alphabet has been ill-treated by nobilities who still claimed that the Chinese language was the only option. But then the new system became gradually popular among people of literature.

During the 36 years of Japanese invasion in the early 20th century, it was forbidden to read or write han-geul or to speak Korean. During that time the name *한글* was used for the first time to refer to the Korean system. Af-

¹The right sound to teach people.

ter independence of Japan, 한글 became the official alphabet for the Korean country.

Currently, Korean is natively spoken by about 67 million speakers, 63 million of which live inside Korea. The Korean language is commonly divided into six main dialects. The Seoul dialect is usually referred to as “standard Korean”.

2.2 한글 (Han-geul) – The Korean Writing System

The complete Korean alphabet consists of 40 letters. There are 10 basic vowels and 14 basic consonants. In addition, there are 11 compound vowels which are combined among the basic vowels, and 5 double consonants.

| Vowel | Symbolizes | Actual Form | Phoneme | Romanisation |
|-------|------------|-----------------------|---------|--------------|
| · | Heaven | Not used in isolation | – | – |
| — | Earth | — | /i/ | eu |
| ㅣ | Man | ㅣ | /i/ | i |

Table 2.1: The three basic vowels.

This section describes the definition of the letters of the Korean alphabet. The following tables show the original form of the letters, their actual form, the respective IPA [1] phoneme symbol for their pronunciation and the respective roman transliteration symbol that will be used throughout this thesis. Romanisation of the Korean language will be discussed in section 2.5.

The symbols of the three most important cosmic elements heaven “·”, earth “—” and man “ㅣ” are used as a foundation for the vowel system. With these three symbols as building blocks a total of ten basic vowels are created. Two of these are the *earth* and *man* symbols themselves, pronounced as /i/ and /i/, respectively. These three basic symbols are shown in table 2.1.

| Old Form | Actual Form | Phoneme | Romanisation |
|----------|-------------|---------|--------------|
| ㅏ | ㅑ | /o/ | o |
| ㅓ | ㅕ | /u/ | u |
| ㅗ | ㅛ | /a/ | a |
| ㅛ | ㅜ | /ə/ | eo |

Table 2.2: Primary vowels.

Four primary vowels are then created by building combinations of *earth* and *man* with *heaven*. *Earth* is combined with *heaven* above it to form the letter ㅏ, which is pronounced as /o/. Situating *heaven* below *earth* creates the vowel ㅓ, pronounced as /u/. Combining *man* with *heaven* to its left or right results in the letters ㅗ and ㅛ which are pronounced as /a/ and /ə/, respectively. The primary vowels are summarized in table 2.2.

| Old Form | Actual Form | Phoneme | Romanisation |
|----------|-------------|---------|--------------|
| ㅛ | ㅝ | /io/ | yo |
| ㅜ | ㅠ | /iu/ | yu |
| ㅗ | ㅜ | /ia/ | ya |
| ㅛ | ㅟ | /iə/ | yeo |

Table 2.3: Secondary vowels.

Four more vowels – the secondary vowels – are built by combining *earth* and *man* with two *heaven* symbols. Following the same combination concept for the primary vowels, this results in the letters ㅝ, ㅠ, ㅜ and ㅟ with two *heaven* symbols to the top, bottom, right and left, respectively. The pronunciations for these vowels are /io/, /iu/, /ia/ and /iə/. See table 2.3 for a summary.

These ten basic vowels are combined in pairs to form a total of 11 compound vowels. A summary of the original and actual form of all these letters along

with their IPA phoneme and roman transliteration symbol is given in table 2.4.

| Old Form | Actual Form | Phoneme | Romanisation |
|----------|-------------|---------|--------------|
| ㅐ | ㅑ | /ɛ/ | ae |
| ㅒ | ㅓ | /iɛ/ | yae |
| ㅕ | ㅖ | /e/ | e |
| ㅗ | ㅛ | /ie/ | ye |
| ㅜ | ㅠ | /oa/ | wa |
| ㅡ | ㅟ | /oɛ/ | wae |
| ㅝ | ㅞ | /uə/ | weo |
| ㅝ | ㅞ | /ø/ | oe |
| ㅟ | ㅠ | /y/ | wi |
| ㅞ | ㅠ | /ue/ | we |
| ㅟ | ㅠ | /ii/ | yi |

Table 2.4: Compound vowels.

Five basic symbols were created as a foundation for building consonant letters. These five symbols were shaped after organs in the human articulation system which produce the sounds. Namely these five basic symbols are ㄱ, ㄴ, ㅁ, ㅅ and ㅇ. The tongue is pressed against the **molar teeth** to produce the phoneme /k/ associated with ㄱ. The phoneme /n/ is represented by ㄴ which shows the shape of the **tongue** while this sound is produced. A front view of the **lips** looks like ㅁ while producing an /m/. **Teeth** and tongue are used to produce an /s/ sound. The symbol ㅅ represents this “tooth” sound. And finally, the circle shaped **pharynx** (throat) is used to produce the phoneme /ŋ/. The pharynx is symbolized by ㅇ.

Several more consonants are built by increasing the level of articulation among these five basic ones. This increasing articulation is represented by adding further strokes to a letter.

As stated before, the 한글 letters were developed on a pure phonetic basis.

While the pronunciation of vowel letters is fixed and an IPA sound symbol can be associated easily, we can not provide one simple sound symbol for consonants. This is because the acoustic realisation of a letter is highly dependent on its context. We call this a “variant realisation” or an “allophone”. So, to introduce the consonant letters here we associate them with their roman transliteration symbol only. For a detailed discussion of the phonetic and allophonic characteristics of the consonant letters see section 4.4.1.

The definition form of all 19 consonant symbols is displayed in table 2.5. These symbols are grouped according to the five basic consonant sets defined above. Table 2.6 lists the original symbols again along with their actual form and their roman transliteration symbol.

| Producing Organ | Basic Consonants | Extensions | Exceptions |
|------------------|------------------|------------|------------|
| Molar Tooth | ㄱ | ㅋ ㆁ | |
| Tongue | ㄴ | ㄷ ㄹ | ㄹ |
| Lips | ㅁ | ㅂ ㅃ ㅆ | |
| Tooth | ㅇ | ㆁ ㆅ ㅆ ㅈ | |
| Pharynx (Throat) | ㅇ | ㅎ | |

Table 2.5: The five basic consonants and their extensions.

2.3 Syllable Construction

The Korean writing system is letter based which makes it fairly easy to learn. The letters are not written sequentially like in the roman writing systems but arranged in syllable complexes, also called *characters*.

The notion of a syllable in Korean is different than in English. A Korean consonant by itself can not form a syllable. This is because a consonant, if not followed by a vowel, can not be released. However, a vowel by itself can form a syllable. In English, a syllable is a sound or a group of sounds accompanied by one of four stresses, whereas in Korean it is a sound or group of sounds which takes up a certain relative space of time like metronome beats.

| Old Form | Actual Form | Romanisation | Old Form | Actual Form | Romanisation |
|----------|-------------|--------------|----------|-------------|--------------|
| ㄱ | ㄱ | k | ㅁ | ㅁ | m |
| ㅋ | ㅋ | kh | ㅂ | ㅂ | p |
| ㆁ | ㆁ | kk | ㅅ | ㅅ | ph |
| ㄴ | ㄴ | n | ㅆ | ㅆ | pp |
| ㄷ | ㄷ | t | ㅈ | ㅈ | s |
| ㅌ | ㅌ | th | ㅊ | ㅊ | c |
| ㄸ | ㄸ | tt | ㅌ | ㅌ | ch |
| ㄹ | ㄹ | r/l | ㅍ | ㅍ | ss |
| ㅇ | ㅇ | ng | ㅑ | ㅑ | cc |
| ㅎ | ㅎ | h | | | |

Table 2.6: Original and actual form of the consonant symbols and their romanisation.

A syllable is built according to one of the three structural types CV, CVC and CVCC, where C stands for consonant and V stands for vowel. These letters are arranged in an imaginary square according to the following rules:

- The nine vowels that have the stroke standing in a vertical position, namely ㅏ, ㅑ, ㅓ, ㅕ, ㅗ, ㅛ, ㅜ, ㅠ, ㅡ and ㅣ, are situated to the right side of the syllable's initial consonant, e.g. 나, 더, 비.
- The five horizontal vowels ㅜ, ㅠ, ㅡ, ㅝ and ㅞ are situated below the syllable's initial consonant, e.g. 소, 주, 그.
- The seven diphthongs ㅑ, ㅓ, ㅕ, ㅗ, ㅛ, ㅜ, ㅠ and ㅡ have the syllable's initial consonant on the top of their left sides, e.g. 되, 꿩, 휘.
- The final consonant (CVC) or consonant pair (CVCC) is put at the bottom edge of the character.
 - With a horizontal vowel it is centered horizontally, e.g. 물, 습, 손, 뚝.
 - With a vertical vowel it is located slightly to the right of the horizontal center position, e.g. 각, 강, 섬, 법, 값, 넓.

“set A minus set B”.

• The phonemes ㅂ, ㄷ and ㄱ are lenis plosives that are articulated with little tension and without aspiration. They are bilabial, apico-alveolar and dorso-velar, respectively. For each of these phonemes, three allophones exist:

1. Voiceless: [p], [t], [k]
 - (a) As word initial sound, e.g. ㅂ, 다 or 그.
 - (b) Inside a word after **C\S**, e.g. ㅂ in 학부.
2. Voiced: [b], [d], [g]
 - (a) Inside a word between vowels, e.g. ㅂ in 나ㅂ.
 - (b) Inside a word after **S**, e.g. ㅂ in 공부.
3. Unreleased: [p̚], [t̚], [k̚]
 - (a) Word final position, where the word is followed by an articulation break, e.g. ㅂ in 입.
 - (b) Inside a word before **C\{ㄱ, ㄴ, ㄷ, ㅎ}**, e.g. ㄱ in 학부.

• The phoneme ㅈ is a lamino-alveolar affricate articulated with little tension and without aspiration. The variants are:

1. Voiceless: [t͡ɕ]
 - (a) As word initial sound, e.g. 잘.
 - (b) Inside a word after **C\S**, e.g. 극장.
2. Voiced: [d͡ʒ]
 - (a) Inside a word between vowels, e.g. 가족.
 - (b) Inside a word after **S**, e.g. 남자.
3. Unreleased: [t͡ɕ̚]
 - (a) Word final position, where the word is followed by an articulation break, e.g. 빛.
 - (b) Inside a word before **C\{ㄱ, ㄴ, ㅎ, ㅅ, ㅆ}**, e.g. 멋다.

• The phonemes ㅍ , ㅌ , ㅋ and ㆁ are strongly aspirated sounds corresponding to p , t , k and s , respectively. Two positional variants exist:

1. Voiceless: $[\text{p}^h]$, $[\text{t}^h]$, $[\text{k}^h]$, $[\text{ç}^h]$
 - (a) As word initial sound, e.g. 팔 .
 - (b) Inside a word after **C** or between vowels, e.g. 남포 , 앞에 .
2. Unreleased: $[\text{p}^ʔ]$, $[\text{t}^ʔ]$, $[\text{k}^ʔ]$, $[\text{t}^ʔ]$
 - (a) Word final position, where the word is followed by an articulation break, e.g. 밭 .
 - (b) Inside a word before $\text{C} \setminus \{\text{ㄱ}, \text{ㄴ}, \text{ㄷ}, \text{ㄹ}, \text{ㅅ}, \text{ㅆ}\}$, e.g. 앞방 .

• The phoneme ㅅ is a voiceless apico-alveolar fricative. Before ㅣ and ㅑ it is slightly palatalised. The variants are:

1. Voiceless: $[\text{s}]$
 - (a) As word initial sound, e.g. 소 .
 - (b) Inside a word between vowels, e.g. 무슨 .
 - (c) Inside a word after **C**, e.g. 항상 .
2. Unreleased: $[\text{t}^ʔ]$
 - (a) Word final position, where word is followed by an articulation break, e.g. 빛 .
 - (b) Inside a word before $\text{C} \setminus \{\text{ㄱ}, \text{ㄴ}, \text{ㅎ}, \text{ㄹ}, \text{ㅅ}, \text{ㅆ}\}$, e.g. 햇것 .

• The phonemes ㅃ , ㅆ , ㄱ , ㅋ and ㅅ are voiceless unaspirated fortis sounds produced with a partially constricted glottis and additional subglottal pressure. They are bilabial, apico-alveolar and dorso-velar stops, lamino-alveolar affricate and apico-alveolar fricative, respectively. Two variants occur:

1. Glottalized: $[\text{p}^ʔ]$, $[\text{t}^ʔ]$, $[\text{k}^ʔ]$, $[\text{ç}^ʔ]$, $[\text{s}^ʔ]$
 - (a) As word initial sound, e.g. 뿌리 .
 - (b) Inside a word after **C** or between vowels, e.g. 맏딸 , 글씨 .

2. Unreleased: [k̚], [t̚]

- (a) Syllable final position, followed by an articulation break, only ㄱ and ㄷ appear in this position, e.g. 밖, 있고.

• The phoneme ㅁ is a voiced bilabial nasal. It is pronounced as [m] in all positions.

• The phoneme ㄴ is a voiced apico-alveolar nasal. It is pronounced as [n]. In connection with ㄷ it is subject to assimilation effects. See further below.

• The phoneme ㅇ occurs only in syllable final position and is pronounced as [ŋ].

• The phoneme ㄹ is a liquida phoneme which has the variants [ɾ], a voiced apico-alveolar flap, and [l], a voiced apico-alveolar lateral. These variants appear in the following situations:

1. Flap: [ɾ],

- (a) As word initial sound, e.g. 라디오.
 (b) Inside a word between vowels, e.g. 나라.
 (c) Inside a word before ㅎ, e.g. 말하다.

2. Lateral: [l],

- (a) Word final position, e.g. 달.
 (b) Inside a word before C \ {ㅎ}, e.g. 살다.
 (c) Inside a word after ㄹ, e.g. 물로.

• The phoneme ㅎ is a voiceless glottal fricative. Generally it is pronounced as [h]. In front of ㅏ, ㅑ, ㅓ, ㅕ, ㅗ, ㅛ, ㅜ, ㅠ, ㅡ and ㅟ it is slightly palatalised. In front of ㅓ, ㅕ, ㅗ, ㅛ, ㅜ, ㅠ and ㅡ it sounds like a bilabial [f]. The positional variants are:

1. Word initial position, according to above rules, e.g. 하나.

2. Inside a word, syllable initial position:

- (a) After {ㄱ, ㅋ, ㆁ}, according to above rules.
 - (b) After ㅍ\{ㅍ}, aspiration of preceding consonant, e.g. 잊히다.
 - (c) After ㅂ or ㅃ, ㅍ weakens or even mutes, e.g. 마흔.
3. Inside a word, syllable final position:
- (a) In front of ㄱ, ㅋ and ㆁ, aspiration of that consonant, e.g. 좋고.
 - (b) In front of ㅍ and ㅃ, ㅍ mutes; ㅍ is glottalized, ㅃ is lengthened, e.g. 좋소, 좋니.
 - (c) In front of ㅂ, ㅍ weakens or even mutes, e.g. 좋은.

These allophonic rules specify a consonant's pronunciation. In the case of two neighbouring consonants, so called assimilation effects additionally take place which further alter the pronunciation of the involved consonants. These effects can be divided into three groups: 1) neighbouring of two sonorants, 2) neighbouring of an obstruent and a sonorant and 3) neighbouring of two obstruents. These rules are rather complicated. Instead of covering them extensively, here are some examples. A detailed discussion of the assimilation effects can be found in [18, 53].

Two neighbouring sonorants: If ㄱ is followed by ㄴ then this ㄴ will change its pronunciation to [n].

A sonorant neighbours an obstruent: If ㄴ is preceded by ㅁ then this ㅁ will change its pronunciation to [m].

Two neighbouring obstruents: If ㅍ is followed by ㅃ then this ㅃ will change its pronunciation to [s].

2.5 Romanisation

The definition and development of the Unicode character coding system since 1991² made it possible to represent the letters and characters of most languages on a computer system. But up to now, most systems have relied on

²The Unicode Standard Consortium was founded in 1991. Its goal is to develop “a character coding system designed to support the interchange, processing, and display of the written texts of the diverse languages of the modern world.”. Visit <http://www.unicode.org> for further information.

the 7 or 8-bit ASCII standard which only provides codes for the characters of Western languages. To represent non-Western writing systems, for example the Korean one, it is necessary to define a mapping of its characters onto (combinations of) ASCII characters. This is called a *transcription system* or a *romanisation*.

Many systems exist for the Korean language, but none is internationally accepted as standard. The *McCune-Reischauer* transcription system is considered the traditional one and despite its inconsistencies and use of diacritic marks, it is the most widely used.³ Further transcription systems include the ones proposed by the *North Korean* and *South Korean* governments. The *Yale* transcription system has been developed with simplicity and consistency. At present, this system is used by most technical papers in linguistics.

The unix tool *hcode* [72] can convert between several code representations for 한글 characters. We use this tool to convert all our 한글 text data into the *hcode*-specific transcription system. This system is different from the above mentioned ones, yet it is very similar to the one proposed by the South Korean government.

Tables A.1 and A.2 in the appendix show a summary of these transcription systems.

³Libraries have adopted the *McCune-Reischauer* system as their standard. Switching to a new transcription system would be very expensive as for instance all the library references where alphabetical ordering is used would have to be rewritten or updated.

Chapter 3

Speech Recognition

3.1 Motivation

Speech input is more natural than keyboard input as it is less intrusive to the user and he is not bound to a hardware device and does not have to keep his focus on the screen constantly. Therefore, to make the interaction of humans and computers more user-friendly we have to consider speech an essential input component. However, there are tasks, like drawing, where other input modalities are superior to speech.

Applications for speech recognition systems are widespread. For the average user it would be much more convenient to enter text to a *word processing* program via talking than via typing. However, it would, for example, not be feasible to let several employees use speech input for their word processing purposes in a big shared office. Speech recognition systems can be used as part of a *speech-to-speech translation* system which would allow people to communicate using different languages [31]. Large amounts of speech data (broadcast news, interviews, meetings, talk shows etc.) are stored on audio tape and speech recognition could be employed to create a *transcription* of these sources [55,66]. This transcription can then be fed into a database and be indexed for *information retrieval* purposes [25]. *Car* drivers can keep both hands on the steering wheel while they access speech-recognition enabled control instruments like navigation system, radio and cellular phone [59].

Intensive research efforts are still underway to further improve speech recog-

tion technology. Areas of focus include large vocabulary speech recognition [46,47,50–52], recognition of spontaneous speech [12] and robust speech recognition in noisy environments [59].

3.2 Overview

Current large vocabulary speech recognition systems are based on the principles of statistical pattern recognition. A front-end acoustic processor (3.1) converts an unknown speech signal S into a sequence of feature vectors $X = x_1, x_2, \dots, x_m$. The speech recognition system has to find the most probable sequence of words $W = w_1, w_2, \dots, w_n$ given the parameterised acoustic signal X , i.e. it has to find the word sequence \bar{W} that maximizes $P(W | X)$. Using Bayes' rule [10] (3.2) and the fact that $P(X)$ is a constant term in respect to the maximization (3.3), the desired probability can be decomposed as follows:

$$\begin{aligned} \bar{W} &= \operatorname{argmax}_W P(W | S) \\ &= \operatorname{argmax}_W P(W | X) \end{aligned} \tag{3.1}$$

$$= \operatorname{argmax}_W \frac{P(W) \cdot P(X | W)}{P(X)} \tag{3.2}$$

$$= \operatorname{argmax}_W (P(W) \cdot P(X | W)) \tag{3.3}$$

The term $P(W)$ represents the *a priori* probability of observing the word sequence W , independent of the observed speech signal. For instance, the word sequence “how are you” is *a priori* much more likely than the sequence “are how you”. A *language model* (LM) is used to capture this information.

The second term $P(X | W)$ represents the probability of observing the signal X given the word sequence W . This value is determined by an *acoustic model* (AM).

The practical determination of the most likely word sequence \bar{W} requires the solution of a number of difficult problems. The process for finding \bar{W} is called *decoding* and the design of efficient decoders is crucial to the realisation of practical LVCSR systems.

In the description above we assumed that a sentence can be decomposed into word units $W = w_1, w_2, \dots, w_n$. While this is straightforward for languages

like English where the notion of a *word*, as a syntactical spacing unit, is very suitable, the choice of the w_i 's is not at all clear in many other languages. Some languages do not use spacing at all and therefore do not have the notion of a *word*. Other languages, like Korean, do have syntactical spacing units, but using these units as w_i 's is not suitable for a LVCSR system. This thesis focusses on the determination of appropriate w_i 's for a Korean large vocabulary speech recognition system.

In the next four sections we briefly describe each of the components of a speech recognition system. In 3.7 the technique of speaker adaptation is described. It is used to improve recognition performance of speaker independent speech recognition systems. We end the chapter with a short introduction to the Janus Speech Recognition Toolkit in section 3.8.

3.3 Front-End

Computers can only work with discrete data. So, to make computer based speech recognition possible a digital representation of the continuous speech signal is needed. A microphone transforms sound into a continuous electrical signal, then an analog-to-digital converter transfers the continuous electrical signal into discrete time slices and discrete amplitude values. Commonly, this yields 16,000 16-bit samples per second.

It is not usual to use this signal representation directly for the recognition process. First, the amount of data is rather large. Second, the signal still contains a lot of unwanted information like background noise, speaker properties or microphone channel. Thus, several preprocessing steps are used to reduce the amount of data and to extract only the relevant speech information from the signal. The first step in preprocessing is to transform the signal into a spectral representation. Usually, preprocessing systems use a sliding window with a length between 5 ms and 20 ms to extract "frames" of samples from the speech waveform. These frames are commonly extracted every 10 to 20 ms, i.e. the frames usually overlap. Then the discrete Fourier transform [6] is used to transform a frame of time samples into a spectral representation.

Knowledge about the sensitivity of the human ear can be applied to reduce the, commonly over 100, resulting spectral coefficients to about 16 values.

This mapping results in features such as Mel-scale [45].

In order to capture signal changes the feature vector is usually appended with information about neighbour frames. Generally, first and second order differences between the successor and predecessor feature vectors are used. Further features such as signal energy can be added to the vector.

To summarize, the preprocessing is used to obtain a feature vector of around 40 coefficients for about every 10 milliseconds of speech.

Commonly, additional techniques are applied to the resulting feature vector. Among these are:

Mean Normalization The additive stationary parts which are introduced by the channel noise are removed by subtracting the mean of all observation vectors.

Linear Discriminant Analysis (LDA) This is a very efficient technique to reduce the dimension of the feature vector without loss of relevant information. A linear transformation matrix is created based on the classes of the feature vectors. The goal is to build a matrix, with which the feature can be projected into a subspace while keeping or increasing the separability of the classes. The usefulness of linear discriminant analysis in a speech recognition front-end has been widely shown [2, 57, 68].

Vocal Tract Length Normalization (VTLN) One major source of interspeaker variability is the variation in vocal tract shape. Different speakers have different vocal tract lengths. Different vocal tract lengths imply different pitch and formant frequencies. In order to normalize for the length of the vocal tract a maximum likelihood linear or piecewise-linear warping in the frequency axis of the speech signal is performed for each speaker. See [11, 54, 58, 69–71] for more information.

3.4 Language Modeling

The language model captures the probability of a sequence of words $W = w_1, w_2, \dots, w_n$ and is given by:

$$\begin{aligned} P(W) &= P(w_1) \cdot P(w_2 | w_1) \cdot P(w_3 | w_1, w_2) \cdots P(w_n | w_1, \dots, w_{n-1}) \\ &= \prod_{i=1}^n P(w_i | w_1, \dots, w_{i-1}) \end{aligned}$$

A word's likelihood is calculated based on its *word history* w_1, \dots, w_{i-1} . For a vocabulary of size L there are L^{i-1} different histories, and so to specify the probability completely, L^i values would have to be estimated. This is an extremely large number for practical values of L . As a consequence, the histories must be considered to belong to only a manageable number of different equivalence classes. Let $\pi(\cdot)$ denote a mapping of histories into some number of equivalence classes. Then the probability $P(W)$ may be approximated by

$$P(W) \approx \prod_{i=1}^n P(w_i | \pi(w_1, \dots, w_{i-1}))$$

In practice, different equivalence class definitions can be used. The most widely employed approach is the so called *3-gram* model. Here histories are considered equivalent if they end in the same two words. Thus

$$P(W) \approx \prod_{i=1}^n P(w_i | w_{i-2}, w_{i-1})$$

The estimation of the basic probabilities $P(w_i | w_{i-2}, w_{i-1})$ is not trivial. For example, complicated smoothing techniques are used to avoid zero probabilities for 3-grams that do not occur in the text corpus.

For a detailed discussion of LM probability estimation and further LM related issues see [24, 28].

3.5 Acoustic Modeling

Acoustic probabilities can be calculated by representing each word unit with a statistical model. Using a separate model for each word has proven unfeasible for large vocabulary speech recognition systems. It is necessary to obtain several samples of a word from different speakers to train reasonable speaker-independent models for each word, and there are simply too many words to be trained. Also the extension of the recognition vocabulary is very complicated as the training process must be repeated for each new word.

This problem is solved by creating acoustic models for *subword* units. Common choices for subword units are syllables or phones. In fact, phones are most frequently used. As there is only a small and fixed set of phones per language (e.g. about 50 for English) these models can be trained well with a reasonable amount of training data. Adding new words to the recognition vocabulary is then as simple as defining their pronunciation in terms of phones.

Unfortunately, the acoustic realisation of a phone highly depends on the neighbouring phones. This can be addressed by using context-dependent phone models. The first such proposed models took into account the direct predecessor and successor phone and were called triphones [3,34]. Models that consider the two predecessors and the two successors are called quintphones. A general context-dependent phone model is called polyphone. However, the acoustic features of a phone unit are not constant but change within its boundaries. Consequently, many state-of-the-art recognition systems split each phone into a number of states. These are subsequently called *subphonetic* units. Context-dependent subphonetic units are called *subpolyphone models* [19,21]. For our systems, we use subquintphone models, meaning we split each phone into three states and model each state with a context width of two.

The key issue in building context-dependent models is to maintain a balance between the desired model complexity and the number of parameters which can be robustly estimated from the available training data [64]. The solution here is to use sufficiently complex models, more than can be trained robustly, and then to cluster similar ones. This can be done bottom-up or top-down. The bottom-up approach is easier to implement, as it works data-

dependently. The top-down method involves linguistic expert knowledge but it has proven much more powerful in providing suitable models for contexts that do not appear in the training material [35].

The main algorithmic instrument to implement the top-down approach is a phonetic decision tree [5, 8, 38]. At the root of the tree is the set of all polyphones corresponding to a phone. Each node has a binary “question” regarding their left and right contexts. These questions are created using expert knowledge and are designed to capture classes of contextual effects. In general, they are of the form “Does the {next, second-next, previous, before-previous, ...} phone belong to phonetic class x ”, for example “Is the previous phone a consonant?”.

Constructing such a tree is a sequential optimization process which recursively partitions the set of states based on a goodness-of-split criterion. The tree leaves contain the acoustic models. To find the corresponding leaf for a specific polyphone model, the tree is traversed by answering the questions attached to each node, until a leaf node is reached. All models that fall in the same leaf are then represented by the same acoustic model.

Most state-of-the-art systems use hidden Markov models (HMMs) as a statistical representation of the subword units [31, 37, 60]. An HMM is a set of discrete states connected by transitions. Each state can produce an observed feature vector with a certain *output probability*. Each HMM transition from any state i to state j has a static *transition probability*. Mainly left-to-right HMM models are used for speech recognition because speech only “goes forward” and hence no back transitions are needed. Commonly, a phone model is represented by a three state left-to-right HMM as three frames roughly correspond to the average length of a phone unit. Each state has a self-loop, a transition to the next state and sometimes even transitions that skip one or several states. See figure 3.1 for a simple left-to-right hidden Markov model. An HMM state represents a small subspace of the overall feature space. The shape of this subspace is sufficiently complex that it is commonly characterized by a multivariate Gaussian distribution.

Word models can be built by retrieving their subword model representation from a pronunciation dictionary and connecting the HMMs of those subword units.

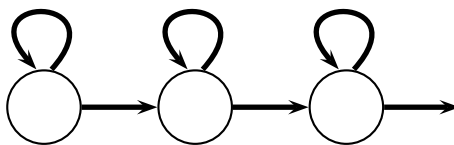


Figure 3.1: Example of a left-to-right hidden Markov model.

There are several reasons why HMMs are the most widely used approach to calculating $P(X | W)$. Speech production can be modelled as a stochastic process. A phone is pronounced differently by different speakers; even the same speaker will pronounce a phone differently at different times. The statistical nature of HMMs along with their ability to model temporal processes make them very suitable for this task. Furthermore, HMMs have been researched for a long time and efficient algorithms exist for training and evaluation of these models.

The theory and the practical use of HMMs in speech recognition is covered very well in the literature. The first theoretic work on HMMs was done by Baum [4] in 1972. A detailed introduction to the theory of HMMs can be found in [20, 42]. A general overview of the speech recognition problem may be found in [43, 44, 61]. They also have very extensive bibliography sections. A good introduction into statistical speech recognition is given in [23]. A very good overview of large vocabulary speech recognition systems is given in [65], and [56] contains a broad spectrum of selected speech recognition related papers.

3.6 Decoding

Finding the most likely sequence of words \bar{W} is a complicated search problem. As with all search problems, there are two main approaches: *depth-first* and *breadth-first*. In depth-first designs, the most promising hypothesis is

pursued until the end of the utterance is reached. Depth-first decoders are commonly called *stack-decoders* or *A*-decoders* [22, 26, 40, 41]. In breadth-first designs, all hypotheses are pursued in parallel. Breadth first decoders are often referred to as *Viterbi decoders* [62]. Sophisticated pruning techniques are employed to reduce the very complex search space. A process called *beam search* is typically used for this purpose [17].

3.7 Speaker Adaptation

Speaker adaptation is a technique used to adjust the acoustic models of a speaker independent recognition system to a specific speaker in order to increase recognition performance. Its general usefulness has been widely proven [36]. Adaptation is divided into two classes: supervised and unsupervised. In the case of supervised adaptation the real textual transcription of an utterance is known. In the case of unsupervised adaptation the transcription is unknown.

We used speaker adaptation for two purposes: Firstly, supervised adaptation is employed to improve the “labeling”¹ of the training utterances. This technique is also referred to as *label-boosting*. Secondly, we use speaker adaptation in the recognition stage to improve performance. As the actual transcription of a test utterance is not known beforehand, the following strategy is used: in a first recognition pass the best hypothesis is calculated. This hypothesis is assumed to be the actual transcription of the test utterance. Using this transcription, supervised speaker adaptation can be performed. Finally, a second recognition pass is done with the speaker adapted acoustic models.

3.8 The Janus Speech Recognition Toolkit

The Janus Recognition Toolkit (JRTk) [12, 31, 32, 62, 63, 67, 68] is a speech recognition system for research and development. It has been developed at the Interactive Systems Labs jointly at Carnegie Mellon University and the University of Karlsruhe. This toolkit is embedded into a Tcl/Tk interpreter

¹For time efficiency reasons, generally the mapping of speech frames to acoustic models of a training utterance is not done dynamically during training but a fixed mapping is calculated before the training stage. This process is usually called “labeling”.

which allows a user to create a complex recognition system easily while still retaining control over every single variable of the recognition backend. Main benefits from this scripting language interface are a maximum flexibility and a very good extensibility.

The JRTk allows for a large variety of recognition system architectures. It can handle semi-continuous to fully-continuous Gaussian mixture observation models. The Gaussians can be modeled with radial, diagonal or full covariance. Neural nets can also be used to do the acoustic modeling [14,15].

A common architecture of a JRTk recognition system uses three sub-polyphone models per phone. These sub-polyphones are clustered to around 3,000 sub-allophones [13]. To model each of them a multivariate Gaussian with 16 to 48 diagonal covariance components is used.

JRTk uses a multi-pass search strategy for the decoding of an utterance. The two first passes are Viterbi beam search based and produce a word hypothesis graph, also referred to as word lattice. The third pass, which is called lattice rescoring, can make use of higher order language models to extract hypotheses from that word lattice. See [13,61,62] for further details.

Chapter 4

System Description

4.1 Speaker Database

We use the Korean portion of the *GlobalPhone* database [46–52] for the development and evaluation of our systems. This section consists of 20 hours of speech data spoken by 100 native Korean speakers. Each speaker read several articles from a Korean national newspaper. The articles were chosen from the areas: national politics, international politics and economy. The speech data was recorded in stereo at a sampling rate of 48kHz using a close-talking microphone connected to a DAT-recorder. After the sound data was transferred from the DAT-recorder to a hard disc, it was downsampled to 16kHz, 16bit.

| | Training | Test | Test subset |
|-----------------------------|----------|---------------|--------------|
| Speakers | 80 | 10 | 10 (same) |
| Utterances | 6,350 | 798 | 84 |
| Vocabulary (eojjeols) | 41,876 | 7,338 | 923 |
| OOV rate (OOV words) | – | 40.07% (4535) | 41.43% (440) |
| Total utterances | 7,148 | | |
| Total vocabulary (eojjeols) | 45,983 | | |

Table 4.1: Summary of acoustic database.

Eighty of the speakers were used for training the acoustic models. They spoke a total of 6,350 utterances with a vocabulary size of 41,876 eojjeols.

Ten speakers were chosen as test speakers. The remaining ten speakers are kept as a further cross-validation set. The test speakers spoke a total of 798 utterances with a vocabulary size of 7,338 eojeols and an OOV rate of 40.07%.

A subset of 84 uniformly selected utterances from the ten test speakers was used to carry out our experiments. The vocabulary size of this test set is 923 at an OOV rate of 41.43%. See table 4.1 for an overview of the database.

4.2 Language Model Data

To overcome the sparse data problem in language model generation, we collected a large corpus of text data from the internet. We retrieved the online articles of the Korean newspaper *Chosunilbo* [72] from October 1995 to August 1998. A text preprocessing script cleaned the text data by removing all HTML-related code. Numbers were mapped onto their textual transcription. Acronyms were replaced by mapping each letter onto a transcription of its pronunciation. See table 4.2 for the mapping table. Table B.1 in the appendix shows some examples of acronyms along with their transcription. We mapped acronyms that are pronounced as a word rather than a sequence of letters (e.g. FIFA, OPEC, RAM, UEFA) onto a transcription of that pronunciation. Table B.2 in the appendix summarizes mappings that translate units like *mm*, *GB*, *%* or *kbps* into an appropriate textual description. The text processing script finally dropped all sentences which still contained non-한글 characters (such as Chinese) as our speech recognition system is based on a pure 한글 database.

In the following description, we will refer to the transcription text data of the training speakers as *Train*. The transcription text data of the test speaker transcription will be called *Test* and the subset thereof which is actually used for system evaluation (84 utterances) will be called *Test-Utts*. Refer to table 4.1 for an overview.

The text corpus that was retrieved from the internet plus the corpus *Train* will be referred to as *Chosun+Train*. This corpus has a total size of 14,770,769 eojeols and consists of 1,494,509 different eojeols. In terms of characters the total corpus size is 43,332,100 and the vocabulary size is 3,583. To ensure a time efficient evaluation of our unit determination process and the resulting

| Roman letter | 한글 and transcription | Roman letter | 한글 and transcription |
|--------------|----------------------|--------------|----------------------|
| A | 에이 (e-i) | N | 엔 (en) |
| B | 비 (pi) | O | 오 (o) |
| C | 씨 (ssi) | P | 피 (phi) |
| D | 디 (ti) | Q | 큐 (khyu) |
| E | 이 (i) | R | 알 (al) |
| F | 에프 (e-pheu) | S | 에스 (e-seu) |
| G | 지 (ci) | T | 티 (thi) |
| H | 에이치 (e-i-chi) | U | 유 (yu) |
| I | 아이 (a-i) | V | 브이 (peu-i) |
| J | 제이 (ce-i) | W | 더블유 (teo-peul-yu) |
| K | 케이 (khe-i) | X | 엑스 (ek-seu) |
| L | 엘 (el) | Y | 와이 (wa-i) |
| M | 엠 (em) | Z | 제트 (ce-theu) |

Table 4.2: Mapping table for acronyms.

systems we decided to use only about 15% of the large corpus plus the corpus *Train*. This data is referred to as *PartChosun+Train*. The large corpus is only used for some selected experiments. *PartChosun+Train* has a total size of 2,354,072 eojeols. It consists of 417,648 different eojeols. In terms of characters the corpus size is 6,854,294 and the vocabulary size is 3,002. Table 4.3 summarizes the language model corpora information.

| | Chosun+ Train | PartChosun+ Train | Train |
|---------------------------|------------------|----------------------|---------|
| Number of eojeols | 14,770,769 | 2,354,072 | 92,378 |
| Eojeol vocabulary size | 1,494,509 | 417,648 | 41,876 |
| Number of characters | 43,332,100 | 6,854,294 | 303,203 |
| Character vocabulary size | 3,583 | 3,002 | 1,963 |

Table 4.3: Summary of language model corpora.

The *self-coverage* of a text corpus is a one-dimensional function $f(x)$ where x runs from 0 to the size of the vocabulary (*vocabSize*). $f(x)$ lies between 0 and 1 and indicates what fraction of the words of a text is covered when those x words which occur the most often in that corpus are known. Figure 4.1 shows

the eojeol-based self-coverage for the corpora *Train*, *PartChosun+Train* and *Chosun+Train*.

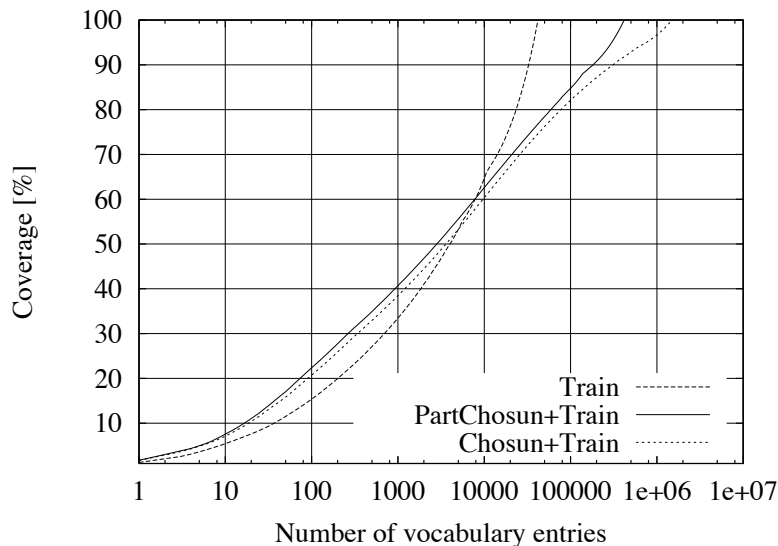


Figure 4.1: Self coverage of language model corpora.

The *cross-coverage* indicates the percentage of *Test* text words covered when the x words are known which occur the most often in the language model corpus. Again, $f(x)$ lies between 0 and 1 but it is not necessarily $f(vocabSize) = 1$. This is because OOV words in the test text cannot be covered as they do not appear in the language model corpus. Figure 4.2 shows the eojeol-based cross-coverage of *Test* using the corpora *Train*, *PartChosun+Train* and *Chosun+Train*, respectively.

Figure D.1 shows the eojeol-based cross-coverage of *Test-Utts* using the same corpora as above. Figures 4.2 and D.1 illustrate the problems using eojeol-based Korean speech recognition. The complete vocabulary of *Train* (41,876 words) covers around 60% (59%) of *Test* (*Test-Utts*) which corresponds to a OOV rate as high as 40.07% (41.43%). *PartChosun+Train* covers up to 80.21% (79.38%) of *Test* (*Test-Utts*). In terms of OOV this is a rate of 19.79% (20.62%). When the complete vocabulary of over 1.2 million eojeols of *Chosun+Train* is used 87.67% (88.80%) of *Test* (*Test-Utts*) are covered. This corresponds to an OOV rate of 12.33% (11.20%). Even this OOV rate is

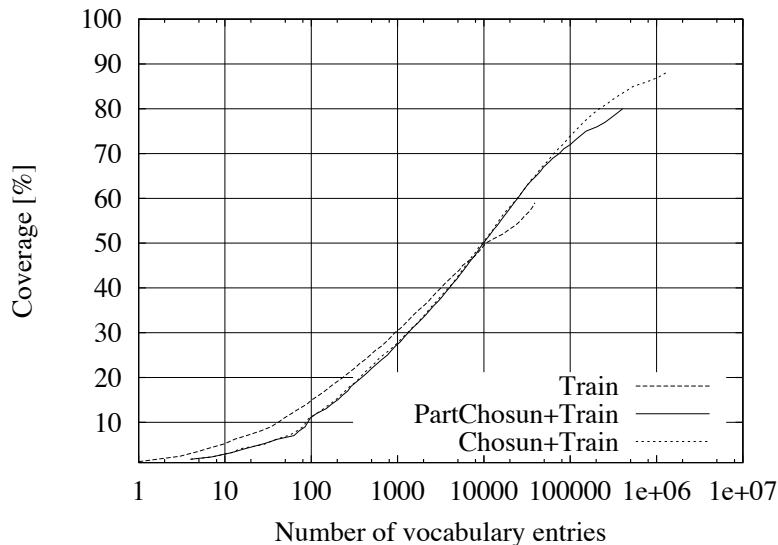


Figure 4.2: Cross coverage of *Test* with different language model corpora.

still unacceptably high for an LVCSR system. The Janus speech recognition toolkit allows for a maximum recognition vocabulary of 64k words. In this case, neither of the mentioned language model corpora cover more than 70% of *Test* or *Test-Utts*, an OOV rate of 30% and above. See table 4.4 for a summary of the OOV information.

| | <i>Test</i> | <i>Test-Utts</i> |
|-------------------------|---------------|------------------|
| <i>Train</i> | 40.07% (4535) | 41.43% (440) |
| <i>PartChosun+Train</i> | 19.79% (2239) | 20.62% (219) |
| <i>Chosun+Train</i> | 11.20% (1268) | 12.33% (131) |

Table 4.4: Summary of OOV rates and, in parentheses, OOV words.

Assuming that a speech recognition toolkit did not have any limit for the recognition vocabulary, then gathering as much text data as possible, as the above figures suggest, could be a means of lowering the OOV rate. Unfortunately, this would not work as the vocabulary grows almost linearly in the text data size as figure 4.3 shows. Thus, satisfactory eojeol-based text

coverage can not be reached in practice.

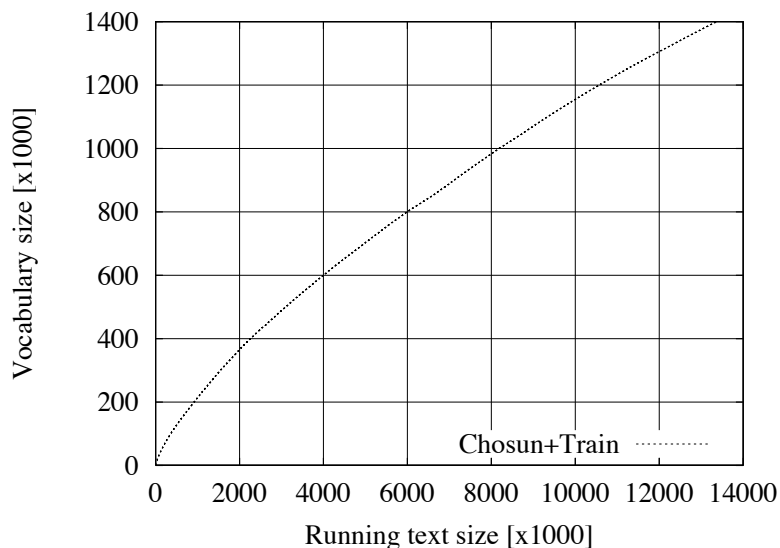


Figure 4.3: Vocabulary growth in corpus *Chosun+Train*.

4.3 Pronunciation Generation

An essential component of an HMM based speech recognizer is the pronunciation dictionary. This dictionary defines a sequence of HMM phone models for each vocabulary unit. For many languages, such a dictionary can only be built by manually editing, which is a strenuous task. Fortunately, this is not the case for Korean. A complete collection of the rules for the allophonic variants of consonant phonemes, as described in 2.4, and the phonological rules, like assimilation, reinforcement and weakening for neighbouring consonants, also in 2.4, makes it possible to generate the pronunciation dictionary automatically.

The basic strategy for automatic dictionary generation is to take a corpus of text data and then generate the pronunciation for each word found by applying the above set of rules.

As an example, we consider the word “cheon-ku-paek-o-sip-nyeon”¹. Applying the set of rules returns the phone sequence [čəngubəgosimmniən]. As we can see, the allophonic rules are processed correctly, for example the letter *k* of the syllable *paek* is mapped on the voiced sound [g] because of its intervocalic position. The neighbouring letters *p* and *n* between the characters *sip* and *nyeon* show an example for the correct treatment of the phonological rules. The two letters are mapped onto the phones [m] and [n], respectively.

Handling phonological changes inside a vocabulary unit is straightforward – simply apply the defined set of rules. However, phonological changes can also occur at unit boundaries. To handle these cases, we extract the last character of the preceding unit and the first character of the succeeding unit and connect them respectively to the beginning and end of the current unit. Now the set of rules can be easily applied, phonological changes happen within the newly created meta-unit. After the corresponding sequence of phones is created, the phones that belong to the two added characters are removed. As a result, we obtain the pronunciation of the current unit in the given context.

Consider the sentence “u-ri cip ro-cheon-ne mal-i-ci”. To create a pronunciation for “ro-cheon-ne” in this context, we first build the meta unit “cip-ro-cheon-ne-mal” by connecting the two neighbouring characters. Applying the allophonic and phonological rules results in the phone sequence [jimnočənnemal]. Removing those phones that belong to the neighbouring syllables which here are [jim] and [mal] leaves us [nočənnə] as the pronunciation. Of course this procedure might return different pronunciations for a specific unit depending on the context. These are handled as pronunciation variants in the recognizer’s dictionary.

4.4 HMM Recognizer Structure

4.4.1 Phone Set

Based on the results presented in section 2.4 a total of 41 phones, 10 vowels, 8 diphthongs and 23 consonants were defined. A list of these phone models can be found in the appendix. The vowel models are listed in table C.1, the consonants are presented in table C.2. The two tables show for each model

¹Korean for: *year 1950*.

its IPA representation, the Janus phone model name and an example English word which contains (an approximation of) the specific phone.

We chose not to represent the diphthongs [oɛ], [iɛ] and [ue] as separate acoustic models but rather to break them up into the two respective monophthong models. These three diphthongs are not well enough represented in the acoustic training material to ensure a reliable model parameter estimation. For the same reason we joined the consonants [p] and [p^h], [t] and [t^h] and also [k] and [k^h].

In addition to the phone models we have one silence model (SIL) and one acoustic model that represents human non-speech articulatory noises (+hGH).

The phone models and the noise model are represented as a three-state, left-to-right HMM (see figure 3.1). For the silence model one four-state left-to-right HMM is used. The output probability of each of these HMMs is modeled with a mixture of 16 diagonal 24-dimensional Gaussians.

4.4.2 Speech Preprocessing

The general problem of turning speech data into a form that can be processed by an automatic speech recognition system was described in section 3.3. Based on that section we describe here the final speech feature vector and which preprocessing steps were used.

A window of size 20ms was shifted over the discretized speech data with an offset of 10ms. For each window 13 Mel-frequency cepstral coefficients were calculated. Mean subtraction was applied to remove the stationary characteristics of the recording channel. Then a composite 43 dimensional feature vector was generated from the 13 Mel-coefficients, their first and second order derivatives and zero crossing plus logarithmic signal energy and its first and second order derivative. The final 24 dimensional feature vector is computed by an LDA transformation of this 43 dimensional feature vector. Vocal tract length normalization is applied to minimize speaker differences. See section 3.3 for references to literature covering these techniques.

4.4.3 Context Dependent Phone Modeling

All context-dependent systems presented in this work consist of 3000 sub-quintphone models. Crossword models were used to capture contextual effects between words. For algorithmic reasons, maximum context width across words is one. The decision tree used for the context dependent models was generated using a set of 63 phone sets. These sets are listed in table C.3 in the appendix. Further information on context dependent modeling can be found in section 3.5.

4.4.4 Model Initialization

An initial context-independent Korean recognition system was trained using the labels generated by a multi-lingual speech recognition system. Among the languages of this system are German, English, Japanese and Spanish. The Korean phone models were initialized with heuristically chosen close equivalents in the multi-lingual system.

Chapter 5

Morpheme Based Recognition

5.1 Motivation

As described in chapter 1 an eojeol unit consists of a combination of morpheme components. For large vocabulary speech recognition, this structure is a severe problem. Consider as an example the noun 학교 which means *school*. It can occur in a text corpus in many different ways. In the subject form the word would appear as 학교가, in the genitive form it would appear as 학교의, in the direct object form it would appear as 학교를 and so forth. Each time a grammatical suffix is added to the noun's stem to indicate its case. In fact, suffix appending to nouns is not only used to indicate their case, but also to indicate the number or prepositional function. For verbs, among the functions of a suffix are indication of tempus, modality and social relationship between the communication partners. It is easy to see that one or more forms of a noun may appear in the text corpus that is used to build the recognition system. Also, when presented with new text to recognize, we can expect to encounter a previously unseen form of that noun. This is also the case with verbs.

A straightforward approach to Korean speech recognition is to use morpheme units instead of eojeols as the base vocabulary. The OOV rate of morpheme based systems is below 5% which makes these units much more suitable. As described in [29] the morpheme units can be used effectively as dictionary units for Korean LVCSR.

5.2 Determination of Units

It is very time consuming and unfeasible to edit very large text corpora by hand in order to split up each eojeol into its morpheme components. Automatic morpheme analyzing systems must be built which break up an eojeol into its morpheme components. These systems are also referred to as part-of-speech (POS) tagging systems and generally use a knowledge based approach. Each resulting morpheme unit receives a part-of-speech tag which indicates its grammatical function.

Unfortunately, the development of a morpheme analyzer involves a lot of human expert knowledge about the morphological structure of the Korean language. And as is always the case with rule-based approaches, these systems include errors and omissions. But generally, these errors are tolerable as they are consistent throughout the underlying text data.

The morpheme segmentation was provided by Oh-Wook Kwon from the ETRI lab in Seoul, Korea.

5.3 Speech Recognition Systems

We evaluated two morpheme unit based approaches. The first one – called *Morph* – uses the pure morpheme components that are generated by the morpheme analyzing system as vocabulary units. The second system – called *MorphTag* – makes use of the POS tag which is associated with each morpheme unit. In this case, the combination of each morpheme with its POS-Tag, *morpheme+<POS-tag>* is used as vocabulary units. The idea is to give the recognition system a means, via the language model, to determine which POS-tags can follow each other, and thus increase the score for hypotheses that consist of valid morpheme sequences.

As the system *Morph* does not use POS information but only the plain morpheme units it can be compared directly to the units generated by our data-driven approach. It allows us to compare whether the morphemes are a better choice than an automatically generated set of units. The system *MorphTag* will demonstrate whether the use of the additional POS information helps to increase the performance of the morpheme based recognition system.

For the baseline comparison of the morpheme approach and the data-driven approach, we used the corpus *PartChosun+Train* instead of *Chosun+Train* for efficiency reasons. In further selected experiments we also evaluated the morpheme systems on the complete corpus *Chosun+Train*. These systems will be referred to as *MorphWhole* and *MorphTagWhole*.

In the rest of this section we will describe the characteristics of the morpheme based systems. Performance analysis and discussion of the two unit determination approaches is presented in section 7.3.

| | <i>Chosun+Train</i> | <i>PartChosun+Train</i> | <i>Train</i> | <i>Test</i> |
|---------------------------------|---------------------|-------------------------|--------------|-------------|
| Number of morphemes | 26,781,814 | 4,249,723 | 185,608 | 21,973 |
| <i>Morph</i> vocabulary size | 349,264 | 116,529 | 17,408 | 4,981 |
| <i>MorphTag</i> vocabulary size | 411,697 | 134,081 | 20,038 | 5,621 |

Table 5.1: Summary of morpheme based language model corpora.

The corpus *Train* has a size of 185,608 morphemes, its *Morph*-based vocabulary size is 17,408 units, for *MorphTag* it is 20,038. The *Test* corpus contains 21,973 morphemes, 4,981 different ones and 5,621 *MorphTag* units. The total size of *Chosun+Train* is 26,781,814 morphemes and it consists of 349,264 morphemes or 411,697 *MorphTag* units. *PartChosun+Train* contains 4,249,723 morphemes at a morpheme vocabulary size of 116,529 and a *MorphTag* vocabulary of 134,081 units. Of course, the vocabulary size of *MorphTag* is bigger than *Morph*. This is because a *Morph* unit can have more than one morphemic function and so can be associated with several POS tags, resulting in more than one *MorphTag* unit. Table 5.1 summarizes this information.

The *Morph*-based self-coverage for the corpora *Train*, *PartChosun+Train* and *Chosun+Train* is displayed in figure 5.1. Figure 5.2 shows the self-coverage for the corpora *Train*, *PartChosun+Train* and *Chosun+Train* in the *MorphTag*-based setup. The *Morph* and *MorphTag*-based cross-coverage of *Test* using the corpora *Train*, *PartChosun+Train* and *Chosun+Train*, respectively is shown in figure 5.3 and 5.4, respectively.