

발음열 자동 생성기를 이용한 한국어 음운 변화 현상의 통계적 분석

Statistical Analysis of Korean Phonological Variations Using a Grapheme-to-phoneme System

이 경 님*, 정 민 화*
(Kyong-Nim Lee*, Min-Hwa Chung*)

*서강대학교 컴퓨터학과 음성언어처리연구실
(접수일자: 2002년 7월 25일; 채택일자: 2002년 10월 4일)

본 논문에서는 한국어 발음열 자동 생성기를 이용하여 한국어의 음운 규칙에 대한 통계적 분석을 수행하였다. 실험에 사용한 발음열 자동 생성기는 한국어 음운 변화 현상에 대해 형태음운론에 기반한 언어학적 분석과 문교부 표준어 규정의 표준 발음법에서 유도된 필수 및 수의적 음소 변동 규칙과 변이음 규칙의 단계적 적용 모델을 사용해서 구현되었으며, 특히 연속 음성 인식을 위한 학습용 발음열과 인식용 발음사전 생성의 최적화를 목표로 하였다. 본 논문에서는 대어휘 연속음성 인식기의 음향 모델을 구축하기 위해 만들어진 삼성 PBS (Phonetically Balanced Sentence) 음성 데이터 베이스의 60,000문장에 적용된 발음열 생성기의 음소 변동 규칙들의 분포 및 그 통계를 사용해서 한국어 음운 변화 양상을 분석하였다. 적용된 빈도수를 기준으로 분석한 결과, 필수 음소 변동 규칙의 경우는 연음법칙, 경음화, 격음화, 장애음의 비음화 순으로, 수의적 음소 변동 규칙의 경우는 초성 h 탈락, 중복 자음화, 동일 조음위치 자음탈락 순으로 음운 변화가 발생하였다. 이러한 적용 규칙들의 통계적 자료를 기반으로 한국어 음운 변화 양상을 파악할 수 있었으며, 나아가 본 논문의 연구 결과는 음성 인식 시스템을 개발하는데 유용하게 사용할 수 있을 것이다.

핵심용어: 발음열 자동 생성, 한국어 음운변화 현상, 발음 사전, 음운 규칙의 통계 분석

투고분야: 음성처리 분야 (2,7)

We present a statistical analysis of Korean phonological variations using a Grapheme-to-Phoneme (GTP) system. The GTP system used for experiments generates pronunciation variants by applying rules modeling obligatory and optional phonemic changes and allophonic changes. These rules are derived from morphophonological analysis and government standard pronunciation rules. The GTP system is optimized for continuous speech recognition by generating phonetic transcriptions for training and constructing a pronunciation dictionary for recognition. In this paper, we describe Korean phonological variations by analyzing the statistics of phonemic change rule applications for the 60,000 sentences in the Samsung PBS Speech DB. Our results show that the most frequently happening obligatory phonemic variations are in the order of *liaison*, *tensification*, *aspirationalization*, and *nasalization of obstruent*, and that the most frequently happening optional phonemic variations are in the order of *initial consonant h-deletion*, *insertion of final consonant with the same place of articulation as the next consonants*, and *deletion of final consonant with the same place of articulation as the next consonant's*. These statistics can be used for improving the performance of speech recognition systems.

Keywords: Grapheme-to-phoneme system, Phonological variations, Pronunciation dictionary, Statistical analysis of Korean phonological rules

ASK subject classification: Speech signal processing (2,7)

I. 서론

음성 인식 시스템을 구성할 때 일반적으로 각 단어에 대해 가능한 발음열을 모두 고려하는 것이 인식을 향상해 도우며, 음성 합성에서도 합성음의 명료성과 자연성을 높이기 위해서 발성 상황에 따라 여러 가지 형태의 발음열 생성이 필요하다. 이러한 이유로 음성학적 발음 특징에 관한 연구를 토대로 규칙을 정의하고 음소 규칙에 의해서 입력 문장을 보다 정확한 발음표기로 변환시키는 시스템들이 개발되었으며[1,2], 발음열 자동 생성 방법에 관한 다양한 연구가 이루어지고 있다[3,4].

본 논문에서는 기존의 발음열 자동 생성 시스템[4]에서는 고려되지 않았던 표준화 규정의 일부 음운 변화 현상을 시스템에 추가 반영하고, 성능 평가 및 안정화 작업을 선행하였다. 이 생성 시스템을 활용하여 발음열 생성 과정에서 적용된 음소 변동 규칙들의 통계적 자료를 기반으로 한국어 음운 변화 현상에 대한 분석을 수행하였다. [3]을 포함하여 기존 연구들은 한글 철자에 대한 통계적 분석이 대부분이며, [7]의 경우 발음사전에 기재된 약 66만 개의 표제어에 대한 발음(음운)을 조사하여 음소와 음절들의 빈도수를 조사 분석한 통계 자료를 제시하였으나, 실제 문장에서 발생하는 형태소 및 어절 경계의 음운 변화 현상은 반영되지 않았으며 적용된 규칙에 대한 정보를 알 수 없다는 한계점이 있었다.

본 실험에 사용된 분석 대상은 트라이폰 기준으로 균형된 음소 집합을 갖도록 구축된 PBS 60,000문장으로 다양한 음운환경을 포함하며 음소열의 중복이 적고 고른 확률 분포를 갖는 문장들의 집합이다. 실험 분석은 본 논문에서 정의한 음소 변동 규칙에 따른 발생 빈도수와 음소의 경계 위치에 따른 적용 양상에 대하여 초점을 맞추었다. 적용된 음소 변동 규칙들의 통계적 자료를 기반으로 한국어 음운 변화 현상의 양상을 파악할 수 있었으며, 나아가 이러한 분석을 이용하여 음성 인식기의 성능을 향상시키기 위한 분석자료로 활용할 수 있을 것이다.

II. 한국어 음소 변동 규칙 정의

발음열 생성 시스템을 구성하는 주요 모듈은 문자열을

발음열로 변환하는 부분으로 문장, 끊어읽기 단위인 언절, 띄어쓰기 단위 어절, 그리고 단어 등 주어진 텍스트를 입력으로 받아 그에 대응하는 발음열을 생성하는 역할을 한다. 이 때 문자열에 대한 올바른 발음열을 생성하기 위해서는 해당 언어의 음운 현상에 대한 체계적이고, 정확한 분석이 필요하다. 본 시스템에서는 음성학과 음운론 연구[5,11]를 기반으로 한국어에서 발생하는 음운 변화 현상을 정리하고, 문교부에서 제정한 표준어 규정[6]의 제 2부 표준 발음법을 참고하여 한국어의 대표적인 음소 변동 규칙 중 표 1과 같이 20개의 음소 변동 규칙을 채택하여 적용하였다. 기존 생성기[4]에서는 적용 대상에서 제외된 모음 관련 규칙을 추가하였다.

표 1. 음소 변동 규칙과 표준 발음법의 대응 관계
Table 1. Correspondence between phonemic change rules and government standard pronunciation rules in[6].

| 음소 변동 규칙 | 표준 발음법 | | 음 | 항 | |
|----------|---------|------------------|-------------------|------------|--------|
| | 규칙 번호 | 규칙 명칭 | | | |
| 종성규칙 | 1 | 음절말 종화 | 117 | 9 | |
| | 2 | 지음군 단순화 | 256 | 8, 10, 11 | |
| | 3 | 격음화 | 21 | 12 | |
| | 4 | 연음법칙 | 42 | 13~15 | |
| 자음의음화 | 5 | 유음화 | 10 | 20 | |
| | 6 | 장애음의 비음화 | 34 | 18, 7장 30항 | |
| | 7 | 유음의 비음화 | 19 | 19 | |
| 자음관련 규칙 | 14 | 변지음화* | 17 | 21 | |
| | 8 | 구개음화 | 3 | 5 | |
| | 9 | 경음화 | 136 | 6 | |
| | 첨가 | 11 | L-첨가* | 30 | 7 |
| | | 13 | 중복 지음화* | 6 | 30 |
| | 탈락 | 10 | 종성 ㅇ-탈락 | 1 | 4 |
| | | 15 | 초성 ㅇ-탈락* | 5 | - |
| | | 12 | 동일 조음위치 지음탈락* | 7 | - |
| | 모음관련 규칙 | 16 | 지음 첫소리 '의' 단모음화 | 18 | 5항 다면3 |
| | | 17 | 용언의 활음형 '저, 쨌, 처' | 3 | 5항 다면1 |
| 18 | | '키' 단모음화* | 17 | 5항 다면2 | |
| 19 | | 첫 음절 외 '의' 단모음화* | 2 | 5항 다면4 | |
| 20 | | 용언 어미 '어' 이중모음화* | 1 | 5 | |

한국어의 경우 주로 자음 변화가 심하기 때문에 이에 대한 연구는 많지만 모음에 대한 연구는 체계적이지 못하다. 주로 사투리나 방언에 대한 연구가 많으며, 특히 발화 속도나 습관에 따라 변화가 다양하여 텍스트 형태와 분석 자료만을 가지고 규칙화하기 어렵다. 사투리나 방언의 경우는 텍스트에 이미 반영된 철자만을 대상으로 하고, 표준 발음법에서 제시된 변화 현상만을 그 대상으로 삼아 모음 관련 규칙을 반영하였다.

표 1에서 규칙 이름에 *가 표시된 것은 수의적 음소 변동 규칙을 나타내며, 총 13개의 필수 음소 변동 규칙과 7개의 수의적 음소 변동 규칙을 갖는다. 각 음소 변동 규칙들은 적용되는 음소 문맥 별로 다시 세부 규칙 번호가 주어지고, 이에 따라 실제 음소 문맥에 규칙이 적용된다. 음소 문맥에 따른 세부 규칙의 수는 총 816개이며, 자음과 모음을 기준으로 분류하면 자음 관련 규칙 775개와 모음 관련 규칙 41개이며, 필수 적용 규칙과 수의적 규칙으로 분류하면 필수 음소 변동 규칙 757개와 수의적 음소 변동 규칙 59개이다.

III. 형태음운론적 분석에 기반한 발음열 자동 생성

3.1. 시스템 구성

본 논문에서 사용된 발음열 자동 생성기 알고리즘은 한국어의 음운 변화 규칙을 다음과 같이 3단계로 나누어 진행된다. 해당 음소 문맥에 의해 하나의 음소가 다른 음소로 바뀌거나 탈락, 첨가되는 양상을 규칙화한 것을 음소 변동 규칙이라 정의하고, 표준 발음 생성을 위한 필수 음소 변동 규칙과 비표준 발음을 포함하여 화자의 습관 및 환경에 따라 발생 가능한 수의적 음소 변동 규칙을 단계별로 적용하였다. 마지막으로 하나의 음소가 음성 환경 말의 속도와 스타일에 따라서 여러 가지 음가를 가지는 변이음 생성 규칙을 적용하였다.

- 1 단계: 표준 발음 생성을 위한 형태소 범주에 따라 필수적 음소 변동 규칙 (표 1의 규칙 1~11, 16, 17)을 적용한다. 이 단계는 형태소의 범주에 따라 어간, 어미, 복합어, 조사, 명사·부사·관형사 (default), 다섯 개로 분리하여 수행한다. 각 형태소 범주와 적용되는 위치에 따라 5개의 필수 음소 변동 규칙들이 선택적으로 적용된다. 음소 변동 규칙으로 처리할 수 없는 발음열은 예외 발음사전을 두어 올바른 발음열이 생성되도록 하였다.

- 2 단계: 표준 발음 이외의 발생 가능한 발음열을 생성을 위한 수의적 음소 변동 규칙 (표 1의 규칙 12~15, 18~20)을 적용한다. 수의적 규칙이 적용되는 발음열은 선택적이므로 필수 규칙이 적용된 발음열은 그대로 두고 새로운 발음열을 생성하게 된다. 즉 필수 규칙이 적용된 발음열에 대해 하나의 복사본을 만든 다음 그 복사본에 대해 수의적 음소 변동 규칙을 적용한다.
- 3 단계: 해당 변이음 규칙을 적용한다. 적용된 변이음 규칙은 무파화, 탄설음화, 유성음화이다. 이 단계를 거치면 자음 28개, 모음 20개로 정의된 48개의 서강대 유사음소 단위 (PLU: Phone-like Unit)를 기준으로 발음열이 생성된다.

3.2. 연속음성 인식을 위한 발음열 자동 생성

한국어의 음운 변화는 음소의 배열과 형태소의 종류에 따라서 영향을 받는다. 같은 음소의 배열이라 하더라도 그 음소열이 '하나의 형태소 내부에 있는가', '형태소 경계에 위치하는가', 또는 '어절 경계에 위치하는가'에 따라 각기 다른 음운 변화 현상을 보여주고 있다.

고립단어 인식에서는 표제어 내부의 음운 변화 현상만을 고려하여 반영된 발음사전을 구축하면 된다. 그러나 연속 음성을 대상으로 발음사전을 구축하기 위해서는 표제어 사이의 음운 변화 현상이 반영되어야 한다. 특히 한국어 문장은 하나 이상의 형태소들이 결합된 어절들로 구성되므로 형태소를 디코딩 단위로 삼는 경우 형태소 및 어절 사이에 발생하는 변화 현상도 반영되어야 한다. 또한 같은 음소 문맥 정보를 갖더라도 발음열 생성시 현 위치와 품사 정보에 따라 변화하는 현상이 달라지기도 한다. 이러한 한국어의 특징을 잘 반영하여 발음열을 생성하려면 주어진 문장을 형태소 분석하고, 올바른 형태소열로 태깅하여 그 정보를 이용해야 한다. 그림 1은 형태음운론적인 분석을 통해 "신발을 신고, 신고하러 갔다"라는 문장이 음소 문맥과 경계 위치에 따라 해당 음소 변동 규칙이 적용되는 과정을 나타낸다.

위 문장은 끊어읽기 단위인 2개의 언절로서 총 4개의 어절과 9개의 형태소로 구성된다. 그림에서 어절 경계는 '/'으로 표현하였으며, 형태소 경계는 '+'로, 품사 태그 정보는 형태소 뒤에 '/'를 붙여 기재하였다. 위 예제에 적용된 필수 음소 변동 규칙의 세부규칙 표현은 표 2와 같다. 음소 문맥 항의 L3는 음소 변동이 일어나는 음절 경계의 앞 음절의 종성을 나타내고, R1은 뒷음절 초성을 나타낸다. 변환 코드는 해당 음소 문맥에 대한 음소의 변동 결과를 나타낸다. 적용범위는 세부규칙의 적용범위와

- 1) 규칙번호: 0.0 4.8 0.0 9.113
 적용순서: (1) (2) (3) (4)
 입력언절: 신발/ncn+을/jco // 신/pvg+고/ecs
 출력결과: [S I Y N B A A R / W W L // S I Y N / K K O W]
- 2) 규칙번호: 0.0 0.0 0.0 0.0 9.72
 적용순서: (5) (6) (7) (8) (9)
 입력언절: 신고/ncpa+하/xsv+러/ecs // 갓/pvg+다/ef
 출력결과: [S I Y N G O W / H I A A / R A X // G A A T Q / T T A A]

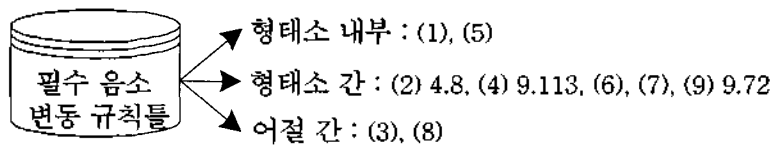


그림 1. 음소 문맥과 경계에 따른 음소 변동 규칙의 적용 예
 Fig. 1. An example of phonemic change rule applications based on phonemic contexts and boundary locations.

표 2. 그림 1에 적용된 세부 필수 음소 변동 규칙
 Table 2. Phonemic change rules used for the example in Figure 1.

| 원소 | 변동 | 변동 위치 | | 원소 | 변동 | 변동 위치 | 변동 |
|----|----|-------|----|----|-----|--------|----|
| | | 내부 | 경계 | | | | |
| ㄹ | 0 | → | ∅ | ㄹ | 8 | 110000 | 4 |
| ㅁ | ㄷ | → | ㄷ | ㅁ | 72 | 111100 | |
| ㄴ | ㄱ | → | ㄴ | ㄴ | 113 | 011000 | |

적용 양상을 나타낸다[4].

그림 1의 필수 음소 변동 규칙 적용 과정을 보면, (1)(5)에는 형태소 내부 규칙, (2)(4)(6)(7)(9)에는 형태소 경계에 적용되는 규칙, (3)(8)에는 어절간 규칙이 경계 위치에 따라 서로 다르게 적용하였다. 입력 언절 1)과 2)의 '신고'라는 단어는 서로 같은 철자이지만, '신고'가 어간과 어미의 결합인 경우 경음화 규칙 중 세부규칙 9.113에 의해 /신포/로 변환되고, 명사인 경우 변화없이 /신포/로 발화된다. 일반적으로 어절간에 일어나는 음운 변화 현상은 수의적 변이음 규칙으로 (8)에서 유성음화 규칙이 적용된 것을 볼 수 있다.

3.3. 발음열 자동 생성기 출력 형태

음성 인식 시스템은 음향 모델을 생성하기 위한 학습 부분과 실제 발성된 발화 내용이 무엇인지 알아내는 인식 부분으로 나뉘어진다. 학습 부분에서는 화자가 발성한 정확한 발음열 정보와 학습용 발음사전이 필요하며, 인식 부분에서는 발음 변화를 최대한 반영한 최적화된 인식용 발음사전이 필요하다.

학습용 발음열은 해당 입력 문자열들이 어떻게 발음되는지에 대한 PLU 리스트이다. 일반적으로 학습은 문장 단위로 수행하므로 한 문장이 끝나는 점을 기준으로 하나의 PLU 열로 출력하도록 하였다. 발음 사전의 경우는 사용 목적에 따라 학습용 발음 사전과 인식용 사전으로 나눌 수 있고, 또한 다양한 음운 변화 현상을 반영하는 정도에 따라 표준 발음사전과 수의적 발음사전으로 구분할 수 있다. 여기서 주로 학습용 발음사전으로 사용하는 경우에는 올바른 발음을 표기한 표준 발음사전을 사용하며, 인식용으로는 다양한 화자의 발화현상을 반영하기 위해 표준 발음과 함께 수의적 규칙이 반영된 다중 발음사전을 사용한다.

표 3. 발음열 자동 생성기 입출력 예
 Table 3. An example of the input and outputs of the GTP system.

| 입력 | 예제 | |
|--------------------|---------------------------------------|----------------|
| 신발/ncn+을/jco | 신/pvg+고/ecs | |
| x o 신바르 | ncn (0.0-0.0-4.8) | [0.0-0.0-0.0] |
| x o 신포바르 | ncn (0.0-0.0-4.8) | [0.0-14.1-0.0] |
| o x 을 | jco (4.8-0.0) | [0.0-0.0] |
| o o 신 | pvg (0.0-9.113) | [0.0-0.0] |
| o o 신포 | pvg (0.0-9.113) | [0.0-14.9] |
| o x 꼬 | ecs (9.113-0.0) | [14.9-0.0] |
| 출력 | S I Y N B A A R W W L S I Y N K K O W | |
| 표준/다중 (필수): | 비표준/다중 (수의적): | |
| 신발 S I Y N B A A R | 신발 S I Y N B A A R | |
| 을 W W L | 신발 S I Y M B A A R | |
| 신 S I Y N | 을 W W L | |
| 고 K K O W | 신 S I Y N | |
| | 신 S I Y N X | |
| | 고 K K O W | |

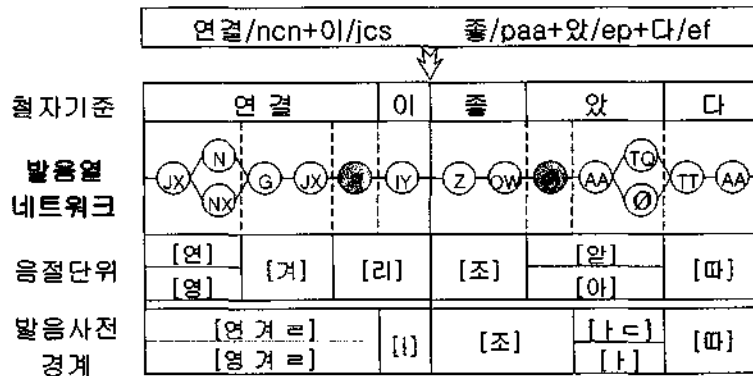


그림 2. 발음 네트워크와 발음사전 표현 예
Fig. 2. Pronunciation networks and pronunciation dictionary for “연결이 좋았다”.

형태소, 어절, 언절 또는 문장 등의 다양한 형태의 입력에 대해 발음열을 생성하도록 구성하였다. 표 3은 “신발을 신고”라는 입력 언절에 대해 다양한 출력한 결과로서 이를 이용하여 학습용 발음열과 발음사전을 구성한다. 중간 출력 결과는 적용된 세부 규칙 번호와 출력 형태를 한글로 볼 수 있도록 구성하였으며 선택 사항이다. 앞뒤에 형태소의 결합 여부를 o/x로 표기하여 경계 위치에 대한 정보를 나타내었으며, 표제어 별로 형태소 태그 정보를 포함하였다. ()안에는 적용된 필수 음소 변동 규칙의 세부 규칙번호를, []안에는 수의적 음소 변동 규칙의 세부 규칙 번호를 표기하였다.

일반적으로 학습용 발음열은 변화된 순서대로 나열하면 되지만, 발음사전의 경우에는 출력 음소열이 어느 표제어에 해당하는지 경계 구분이 중요하다. 그림 2와 같이 연음규칙과 중성 ‘ㅎ’ 탈락 규칙이 적용되어 뒷부분 형태소의 결과값으로 음소열이 출력되는 경우가 발생한다. 이러한 대표적인 규칙으로는 연음 규칙, 중성 ‘ㅎ’ 탈락, 구개음화 규칙 등이 있으며, 일부 격음화 규칙의 경우도 앞 단어의 마지막 음소열이 탈락되면서 뒤 음소열에 영향을 미치는 규칙 등이 있다.

그림 2의 예제는 /ㄹ/이 뒷음절 초성으로 이동되고 앞 음절 종성 /ㅎ/이 탈락되어 제로격으로 출력된 경우로,

변화된 음소를 뒤 형태소의 초성으로 표기하지 않고 앞 형태소의 가장 뒷부분에 단독으로 출력 결과를 표기하도록 하였다.

표 4의 예제는 앞 음절 종성 /ㄷ/과 /ㄱ/의 음가가 뒷음절 초성 /ㅇ/와 /ㅎ/에 영향을 받아 /ㅊ/과 /ㅋ/으로 변화하여 뒷음절 음소열로 실현된 것으로 표제어 경계 정보에 따라 출력 음소열은 앞부분에 표기한다.

IV. 실험 결과 및 분석

4.1. 발음열 자동 생성기 성능 평가

한국어에서 빈번히 발생하는 음운 변화 현상을 테스트 하기 위한 대상으로 표준어 규정[6]의 제 2부 표준 발음법에 명시된 예제를 사용하였다. ‘붙임’과 ‘다만’ 항목을 포함하여 총 364개의 언절을 사용하였으며, 이 중 수의적으로 적용되는 예제가 16개였다.

위의 표 5의 결과는 어절을 한 단어로 입력하여 생성하는 경우와 형태소 태깅을 통해 품사 정보와 경계 정보를 이용하여 발음열을 생성한 결과를 보여준다. 이 중 접미사 정보까지 사용해야 올바른 결과를 얻는 경우가 9개 있었으며, 예외사전과 수의적 발음으로 생성되는 예제

표 4. 발음사전 생성을 위한 출력 예
Table 4. Output phoneme sequences for pronunciation dictionary.

| 입력 예제 | 출력 음소열 | | 적용 규칙 (세부규칙) |
|---------------|--------|-------------------|--------------|
| | 발음열 | 발음사전 | |
| 꿀/ncn+0i/jcs | [크차] | 꿀 [크차] 이 [리] | 구개음화 (8.2) |
| 부탁/nopa+해/xsv | [부타카] | 부탁 [부타카] 해 [해] | 격음화 (3.1) |

표 5. 표준 발음법 예제를 통한 발음열 자동 생성기 평가
Table 5. Evaluation of the GTP system using examples in the government standard pronunciation rules.

| 입력형태 | 올바르게 생성된 예제 수 (세부규칙) | | | 미생성 |
|-----------|----------------------|----|---------|-----|
| | 필수 | 수의 | 예외 발음사전 | |
| 어절 단위 | 191 | 5 | 0 | 168 |
| 형태소 분석 결과 | 232 | 21 | 112 | 11 |

중 중복되는 경우를 포함하여 총 353개의 올바른 발음열을 생성할 수 있었다. 예외사전 참조가 많은 이유는 일반적인 예제를 대표적으로 들고 주로 예외적인 사항을 항목에 포함하고 있기 때문으로 실제 문장에 적용되는 경우 사전 참조 비율은 낮다. 여기서 사용된 예외 발음사전은 낭독체 연속 음성 인식의 사용 목적으로 구축되었으며, 사전 크기는 약 10 K이다. 미 생성된 경우는 주로 어절 사이에 일어나는 소리의 첨가가 제대로 반영되지 않았기 때문으로 주로 어절 사이에서 발생하는 경음화와 'ㄴ', 'ㄹ' 첨가 현상이 발화자의 습관이나 끊어 읽기에 따라 다르게 적용되는 예제들로 언어학적으로 규칙화하기 어려운 부분이었다.

수치적인 평가 외에 실제 음성 인식 실험에서도 형태소 분석을 수행하여 적용한 학습용 발음열과 발음사전을 사용한 경우 단어 인식이 상대적으로 12% 정도 더 좋은 결과를 얻었다[4].

4.2. 실험 데이터 베이스

실험 결과의 합당성을 뒷받침하기 위해서는 본 실험에 사용한 데이터 베이스의 검증 및 분석이 필요하다. 본 논문에서는 발생 가능한 모든 음운 현상을 포함하며, 가능한 다양한 트라이폰 모델을 포함하도록 설계된 삼성 PBS (Phone Balanced Sentence) 음성 데이터 베이스의 문장을 실험에 사용하였다. 구축된 문장에 대한 형태소 분석

표 6. 실험에 사용된 삼성 PBS 음성 데이터 베이스 문장의 특징
Table 6. Characteristics of the sentences in the Samsung PBS Speech DB used for the experiment.

| 문장 | 60,000 | 60,000 | 9.2 어절 |
|--------|-----------|---------|---------|
| 어절 | 551,820 | 170,419 | 2.1 형태소 |
| 형태소 | 1,160,597 | 44,303 | 1.9 음절 |
| 음절 | 2,230,845 | 167,949 | - |
| 형태소 경계 | 608,777 | - | - |

표 7. 트라이폰 발생 빈도수 및 범위
Table 7. Frequencies and coverages of triphones.

| 형태소 수 | 17K | 1.16M | 7.29M |
|--------|---------------------|--------------------|--------------------|
| 트라이폰 수 | 발생횟수>1 (31.25%) | 6,840 (17.196%) | 17,196 (78.54%) |
| | 발생횟수>5 (27.56%) | 4,589 (12.318%) | 12,318 (73.97%) |
| | 발생횟수>10 (25.13%) | 3,712 (25.13%) | 10,524 (71.23%) |

결과는 표 6과 같다. 문장 분석은 형태소 분석 결과에 품사 태그가 부착된 형태를 기준으로 하였다.

표 7은 KKWON PBS 문장[1]과 이 논문에서 사용된 삼성 PBS 60,000문장, 그리고 연속음성 인식기의 언어모델 생성을 위해 수집 가공한 신문 및 방송 뉴스 대상의 7M 형태소 텍스트 코퍼스로부터 발음열 생성기를 사용해서 얻은 트라이폰 수를 보여준다. 각각 1회, 5회 및 10회 이상 나온 트라이폰에 대한 분석 결과이다. 텍스트 코퍼스 크기는 형태소를 기준으로 각각 17 K, 1.16 M, 7.29 M 형태소이며, 표 7의 괄호 안의 백분율은 7 M 형태소를 기준으로 보았을 때 해당 데이터 베이스가 포함하는 트라이폰 비율을 나타낸다.

6배 이상의 텍스트 크기를 갖는 7M 형태소를 기준으로 실험에 사용된 삼성 PBS 60,000 문장은 약 79%의 트라이폰을 포함하고 있으며, 가능한 트라이폰을 균형적으로 포함하도록 설계되었기 때문에 일반 텍스트에서 발생하는 현상보다 신뢰성있는 결과를 보여준다.

4.3. 형태소 범주에 따른 규칙 적용

한국어는 형태소의 범주에 따라 서로 다른 음소열로 발음열이 실현된다. 예를 들어 '감기'라는 어절이 명사인 경우 /K AA M G IY/ (감기)로 발음되고 어간과 어미의 결합인 경우 /K AA M KK IY/ (감기)로 발음된다. 두 번째 음절의 초성 /K/가 형태소 경계 정보와 범주에 따라 /KK/ 혹은 /G/로 달리 발음된다. 여기서는 필수 음소 변동 규칙이 형태소의 범주에 따라 어간, 어미, 조사, 명사·부사·관형사 (default), 복합어로 분리하여 수행된 결과를 분석하였다. 표 8은 규칙 적용 범위에 따라 분류된 음소 변동 규칙 오토마타를 참조하여 얻은 결과로 명사 프로세스의 경우 입력 형태소 중 34.4%가 변동 규칙이 적용되어 다른 음소열로 변화하였다.

표 8. 형태소 범주 별로 적용된 음소 변동 규칙 분석
Table 8. Analysis of phonemic change rule applications with respect to morpheme categories.

| 명사 | 593,666 | 79,940 | 124,120 | 34.4% |
|-----|-----------|---------|---------|-------|
| 어간 | 119,501 | 14,289 | 26,222 | 33.9% |
| 어미 | 210,741 | 32,348 | 1,871 | 16.2% |
| 조사 | 236,649 | 14,513 | 1,692 | 6.5% |
| 복합어 | 40 | 39 | 5 | 110% |
| 합계 | 1,160,597 | 141,129 | 153,910 | 25.4% |

4.4. 적용된 음소 변동 규칙의 통계적 분석

그림 3은 13개 범주의 필수 음소 변동 규칙이 적용된 결과를 분석한 것으로 가로축은 적용된 음소 변동 규칙이며, 세로축은 형태소 내부와 형태소 경계에서 적용된 규칙의 발생 횟수이다. 분석 결과 가장 많이 적용된 규칙은 ‘연음법칙’이며 ‘경음화’, ‘격음화’, ‘장애음의 비음화’ 순이다.

발생 빈도수를 기준으로 가장 많이 적용된 필수 음소 변동 세부 규칙은 형태소 내부의 경우 규칙 번호 4.4 (연음규칙; $L+O \rightarrow \phi+L$)이며, 형태소 경계의 경우 9.72 (경음화; $\#+C \rightarrow C+\#$)이다. 규칙 4.4는 예를 들면 ‘은영/ncn \rightarrow 우녕/’ 과 같이 받침 ‘ㄴ’이 연음이 되어 다음 음절의 초성으로 이동하는 현상이다. 9.72는 형태소 경계의 종성 ‘ㅍ’이 초성 ‘ㄷ’을 만나 경음화규칙이 적용된 것으로 ‘있/paa+다/ef \rightarrow 인/+/따/’의 경우를 들 수 있다.

그림 4는 7개 범주의 수의적 음소 변동 규칙이 적용된 결과를 분석한 것으로 필수 음소 변동 규칙보다 발생 빈도수가 비교적 적다. 수의적 음소 변동은 형태소 경계 정보에 따라 발화 현상이 달라지지는 않으나, 경계에 따라 발음사전에 기재되는 음소열이 변화하므로 분류하여 분석하였다. 다만 모음화 규칙 18, 19, 20은 음절의 종성

변화 규칙으로 형태소 경계에서는 발생하지 않는다.

수의적 음소 변동 규칙인 ‘동일 조음위치 자음 탈락’과 ‘중복 자음화’는 발화 속도에 따른 발생 환경이 서로 상반되는 음운 변화 현상이다. ‘있/paa+다/ef’의 표준 발음은 /인따/이지만 빠르게 말하는 경우 자음이 탈락하여 /이따/로 발화되며, ‘부터/jxc’의 경우 천천히 또박또박 발화하는 경우 중복 자음 ‘ㄷ’이 종성에 추가되어 /분터/로 발화하거나 ‘아파트’가 /압파트/로 발화되기도 한다. 이와 같이 화자가 빠르게 발화하는 경우 ‘중복 자음화’가 잘 발생하지 않으며, 천천히 또박또박 발화하는 경우에는 ‘동일 조음위치 자음 탈락’에 의한 현상이 잘 나타나지 않는다[5]. 이러한 현상은 개발하는 음성 인식 시스템의 종류나 환경에 따라 조절해야 하는 요소로 필요에 따라 사용할 수 있도록 출력 선택 기능을 추가하였다.

분석 결과 형태소 내부에서는 관형격 조사 ‘의’가 ‘에’로 변화하는 규칙이 가장 많이 발생하였다. 실제 낭독체 문장을 대상으로 녹음할 때에는 임의로 ‘에’로 발성하도록 유도하였다. 그 다음으로는 규칙번호 15.1(초성 ϕ -탈락; $\phi+\# \rightarrow \phi+\phi$)으로 ‘화해/ncn/화애/’나 ‘공개/ncpa+해/xsv \rightarrow /공개애/’와 같이 약한 유성음 ‘ㅎ’이 탈락되는 규칙이 형태소 내부나 경계에서 빈번하게 적용되었다.

다음 그림 5는 적용된 필수 음소 변동 세부 규칙을 고빈

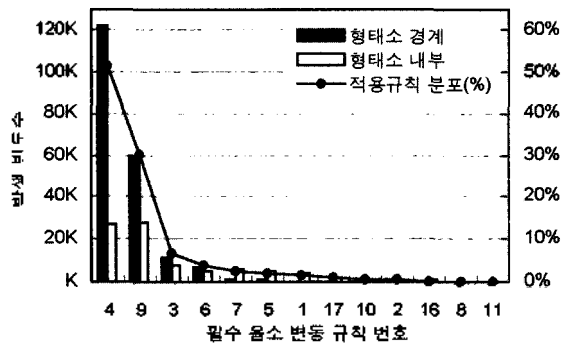


그림 3. 필수 음소 변동 규칙 발생 횟수 및 분포
Fig. 3. Occurrence frequency and distribution of the obligatory phonemic rules.

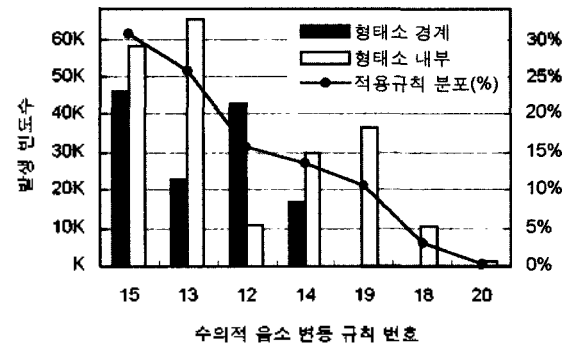


그림 4. 수의적 음소 변동 규칙 발생 횟수 및 분포
Fig. 4. Distributions of optional phonemic rule applications.

표 9. 적용 횟수 기준 상위 5개의 필수 음소 변동 규칙
Table 9. Top 5 obligatory phonemic change rules.

| 순위 | 형태소 내부 | | 형태소 경계 | |
|----|--------|-----------------------------|--------|-----------------------------|
| | 규칙번호 | 적용 규칙 | 규칙번호 | 적용 규칙 |
| 1 | 4.4 | L+O \rightarrow ϕ +L | 9.72 | #+C \rightarrow C+# |
| 2 | 4.8 | L+O \rightarrow ϕ +L | 4.4 | L+O \rightarrow ϕ +L |
| 3 | 9.4 | ㄱ+ㄱ \rightarrow ㄱ+ㄱ | 4.1 | ㄱ+O \rightarrow ϕ +ㄱ |
| 4 | 9.1 | ㄱ+ㄱ \rightarrow ㄱ+ㄱ | 4.8 | L+O \rightarrow ϕ +L |
| 5 | 5.1 | L+ㄹ \rightarrow ㄹ+ㄹ | 4.17 | O+O \rightarrow ϕ +O |

표 10. 적용 횟수 기준 상위 5개의 수의적 음소 변동 규칙
Table 10. Top 5 optional phonemic change rules.

| 순위 | 형태소 내부 | | 형태소 경계 | |
|----|--------|---|--------|---|
| | 규칙번호 | 적용 규칙 | 규칙번호 | 적용 규칙 |
| 1 | 19.2 | 의 (조사) \rightarrow 에 | 12.3 | C+# \rightarrow ϕ +# |
| 2 | 15.1 | ϕ +ㅎ \rightarrow ϕ + ϕ | 15.1 | ϕ +ㅎ \rightarrow ϕ + ϕ |
| 3 | 13.4 | ϕ +E \rightarrow C+E | 15.5 | O+ㅎ \rightarrow O+ ϕ |
| 4 | 15.5 | O+ㅎ \rightarrow O+ ϕ | 13.2 | ϕ +ㄱ \rightarrow ㄱ+ㄱ |
| 5 | 13.7 | ϕ +ㅍ \rightarrow C+ㅍ | 15.2 | L+ㅎ \rightarrow ϕ +L |

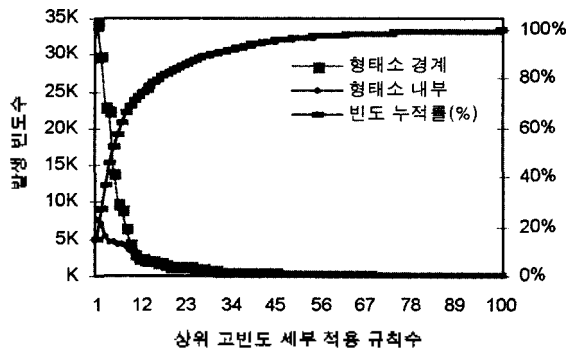


그림 5. 적용된 필수 음소 변동 규칙의 분포도
Fig. 5. Distributions of obligatory phonemic rule applications.

도 순으로 정렬한 상위 100개의 규칙 분포도이다. 필수 음소 변동 세부 규칙 757개 중 192가지의 규칙이 삼성 PBS 60,000문장에서 적용되었으며, 총 289,169회의 변동 규칙이 발생하였다. 형태소 경계와 내부에서 모두 포함하여 1000번 이상 발생한 규칙은 상위 36번째 규칙까지이며, 100번 이상은 상위 82번째까지이다. 이 중 평균 상위 100개의 규칙으로 약 99.67%의 적용률을 보였다.

음성인식을 수행할 때 발음변화를 모델링하는 방법으로는 발음사전을 사용하는 것이 대표적이다. 표제어 내부에서 일어나는 음운 변화 현상은 발음사전에 등록하여 해결할 수 있으나, 경계 부분에서 발생하는 변화 현상을 반영하기 위해 발음사전에 가능한 모든 발음을 등록하는 경우에는 표제어 수가 증가함에 따라 인식 속도와 인식률에 나쁜 영향을 미치게 된다. 이를 해결하기 위한 방안으로 이 논문에서 소개된 분석 결과를 활용하여 빈번히 발생하는 음운 변화 현상만을 발음사전에 추가하여 활용할 수 있을 것이다.

또 다른 접근 방법으로는 인식 단위(표제어)의 경계 부분에서 일어날 수 있는 음운 변화 현상을 인식 네트워크에 적용시키는 방법으로 인식 전에 미리 앞뒤에 가능한 음소 문맥을 적용시켜서 인식 네트워크를 만드는 방법이 있다. 인식 단위의 경계 부분에서 일어날 수 있는 모든 가능한 음소 문맥을 인식 전에 미리 인식 네트워크에 적용하는 방법으로 앞 표제어의 종성과 뒤의 초성의 쌍으로 나타낼 수 있는 모든 쌍에서 음운 변화 현상이 일어나는 것이 아니라 일정한 규칙에 따라 특정한 쌍에서만 일어나게 된다. 특히 형태소 내부와 형태소 경계에서 발생하는 현상이 다를 뿐만 아니라 음소 문맥에 따라 발생 가능한 네트워크만을 확장하는 것이 효율적이므로 이 논문에서 소개된 자료를 활용하여 인식기의 성능을 향상시킬 수 있다. [2]의 연구에서는 이러한 분석 자료를 이용하여 트

리 구조의 인식 네트워크의 공유 효율을 높이고 이로 인해 네트워크의 크기를 줄일 수 있도록 인식 중에 음소 문맥을 이용해 인식 네트워크에 음운 변화 현상을 적용시키는 방법을 제안하였다.

V. 결론

정확한 발음열을 생성하기 위해 한국어가 가지는 언어학적 지식과 문교부 제정 표준어 규정을 기반으로 음운 변화 규칙을 분석하고, 이를 통해 정의된 음소 변동 규칙과 변이음 규칙을 다단계로 적용하여 가능한 모든 발음열을 생성하였다. 정의된 음소 변동 규칙들이 실제 적용되는 현상을 분석하기 위하여 휴대폰 기반의 PBS 60,000 문장에 발음열 자동 생성기를 적용하여 나온 결과를 통계적으로 분석하였다. 실험은 음소변동을 모델링한 분류에 따른 빈도수와 음소의 경계 위치에 따른 적용양상에 대하여 초점을 맞추었다. 적용된 음소 변동 규칙들의 통계적 자료를 기반으로 한국어 음운 변화 현상 양상을 파악할 수 있었으며, 나아가 이러한 분석을 이용하여 음성 인식기의 성능을 향상시키기 위한 자료로 활용할 수 있을 것이다.

일반적으로는 가능한 모든 음운 변화 현상을 분석하여 모델링하는 것이 정확한 음운변이를 반영할 수 있으나, 혼잡도 증가와 변별력 감소 문제 및 인식 네트워크 확장 시 가능한 음소 문맥을 적용하는 경우 적용 규칙수가 필요 이상으로 많아지기 때문에 본 논문에서 통계적으로 분석된 음운 변화 현상을 사용함으로써 시스템 개발에 유용하게 사용할 수 있을 것이다.

감사의 글

본 실험에 사용한 삼성종합기술원의 PBS 음성 데이터 베이스 사용허가에 감사드립니다.

참고 문헌

1. 이경남, 전재훈, 정민화, "한국어 연속음성 인식을 위한 발음열 자동 생성," 한국음향학회지, 20 (2), 35-43, 2001.
2. H. Elovitz, R. Johnson, A. Mchugh, and J. Shore, "Letter-to-sound rules for automatic translation of english text to phonetics," *IEEE Trans. Acoust., Speech, Signal Processing*,

- ASSP-24 (6), 446-459, 1976.
- 3 B. Kim, W. Lee, G. Lee, and J. Lee, "Unlimited vocabulary grapheme-to-phoneme conversion for Korean TTS," *Proc. of ACL-COLING 98*, 675-679, 1998.
 - 4 H. Strik and C. Cucchiari, "Modeling pronunciation variation for ASR: Overview and comparison of methods," *Proc. of the ESCA workshop 'Modelling pronunciation variation for automatic speech recognition'*, 137-144, 1998.
 - 5 김홍규, 강범모, *한글 사용빈도의 분석*, 고려대학교 민족문화연구원, 1997.
 - 6 한국방송공사, *표준 한국어 발음 대사전*, 1993.
 - 7 이기문, 김진우, 이상역, *국어음운론*, 학연사, 2000.
 - 8 J. Clark and C. Yallop, *An Introduction to Phonetics and Phonology*, Oxford, 1995.
 - 9 표준어 규정, 문교부 고시, 88 (2), 1988.
 - 10 김봉완, 김종진, 김선태, 김태환, 김영일, 이용주, "공동 이용을 위한 단어음성DB의 구축 및 PBS 설계에 관한 검토," 제13회 음성통신 및 신호처리워크샵 논문집, 13 (1), 256- 261, 1996.
 - 11 김한준, 음소 문맥과 음운 변화 현상을 이용한 한국어 연속 음성 인식, 서강대학교 컴퓨터학과 석사학위 논문, 2001.

저자 약력

● 이 경 님 (Kyong-Nim Lee)



1996년 2월: 명지대학교 컴퓨터공학과 졸업 (공학사)
 1998년 8월: 서강대학교 대학원 컴퓨터학과 졸업 (공학석사)
 1998년 9월 ~ 현재: 서강대학교 대학원 컴퓨터학과 박사과정
 * 주관심분야: 연속음성인식, 어휘/음운모델링, 발음사전

● 정 민 화 (Min-Hwa Chung)

한국음향학회지 제20권 제2호 참조