

한국어 음소 인식을 위한 기준 프레임 추출

김 범 국* 정 현 열*
*영남대학교 전자공학과

Typical Frame Extraction for Korean Phoneme Recognition

Bum Koog Kim, Hyun Yeol Chung*

*Dept. of Electronic Eng. Yeungnam University

요 약

음소를 인식의 기본으로하는 한국어 음성인식 시스템을 구현하기 위한 기초 연구의 일환으로서 각 음소의 특징을 가장 잘 표현하는 기준프레임 추출을 위한 연구를 수행하였다.

이를 위하여 먼저 실험실과 분산비 분석을 통해서 인식에 필요한 시간 패턴의 길이를 추출한 후 이를 바탕으로 통계적 인식방법인 베이스 결정법칙을 이용하여 시단 프레임으로부터 3프레임(1프레임 10ms)씩 시합을 1프레임씩 옮기면서 인식 실험을 행하여, 각 음소별 특징이 가장 풍부한 기준프레임(3프레임)을 추출하였다. 그리고 이 기준 프레임은 중심으로 각 음소군별 인식 실험을 수행하여 그 결과를 시단을 기준으로한 경우와 비교 검토하고 한국어 전 음소별로 확장하여 인식 실험을 실시하였다.

이 실험 결과 모음의 경우 시단에서부터 5프레임, 파열음은 시단에서부터 5프레임사이, 마찰음은 3프레임에서부터 10프레임까지, 파찰음은 5프레임에서 9프레임까지, 비음과 유음의 경우 초성은 시단 프레임에서 6프레임, 종성은 종단으로부터 전 4프레임 구간이 인식율이 높게 나타나 이 부분의 특징이 인식에 가장 유효함을 알 수 있었다.

1. 서 론

음소를 인식의 기본단위로 하는 음성인식 시스템을 구현하기 위해서는 각 음소가 지속되는 시간 방향의 정보중 그 음소를 특징 지우는 변별적 특징이 가장 많이 포함된 위치(이하 기준프레임이라 함)가 어디인가를 파악할 수 있다면 음성의 합성, 인식시 매우 유용하게 이용될 수 있을 것이다.

본 연구실에서는 그 동안 한국어 음성인식 시스템을 구현하기 위한 그 기초적인 연구로서 한국어 전음소들 대상으로 한국어 모음의 분석과 인식⁽¹⁾, 파열자음의 분석, 인식^(2,3), 마찰음 및 파찰음의 분석과 인식⁽⁴⁾, 비음과 유음의 인식을 위한 특징 추출에 관한 연구⁽⁴⁾, 한국어 단음절에 포함된 음소 인식에 관한 연구⁽⁶⁾, 다차원 척도법(Multi-Dimension Scaling)을 이용한 한국어 음소 분석⁽⁷⁾등의 연구를 통하여 음소를 단위로서 분석과 인식에 관한 연구를 수행 해 왔다.

이들 연구 결과를 참고로 본 논문에서는 각 음소의 특징을 가

장 잘 표현하는 기준 프레임의 위치를 찾아 내고자 한다. 이를 위하여 먼저 실험 연구⁽¹⁻⁷⁾를 참고로 인식에 필요한 시간 패턴의 길이를 추출하고, 이를 바탕으로 시단 프레임으로부터 3프레임씩, 시점을 1프레임씩 옮겨 가면서 각 음소군별 실험을 통하여 전 음소의 기준 프레임(3프레임)을 추출한다. 이때 특징 파라미터로 21차원 cepstrum계수를 사용한다.

이 결과를 이용하여 각 음소별 기준 프레임은 중심으로 특징 파라미터를 추출하여 인식 실험을 행한 결과와 시단에서부터 프레임은 증가 시키면서 인식 실험을 행한 결과와 비교 검토한다.

2. 음소의 분석

2.1 음성 자료

단음절 자료는 현국어 억음 사전중의 단어를 출원빈도순으로 나열하여 누적빈도 90% 이내에 들어가는 단음절 501개와 음소의 수가 적은 경우에 대해서는 누적빈도가 99.9% 까지 들어가는 단음절로부터 짧은 48개를 추가한 총 549개 이다.

이 단음절 자료들 방음실에서 한국인 성인남성 3인이 랜덤하게 각각 3회씩 자연스럽게 명성한 4941개로 한다.

2.2 분석 방법

단음절 음성자료는 그림 1에 보인 바와 같이 4.5KHz LPF를 통과한 후 10KHz 12bit A/D 변환기를 통해 음성 데이터로 변환되고, 29CH BPF(Q=6, 1/6 octave, 250Hz ~ 6300Hz)를 통과시켜 분석 된다.(프레임 길이:10ms, 분석창 길이:20ms) 각 음소에는 시합에 의해 시단, 중심, 종단 프레임이라는 시간적 레이블을 부여한다.

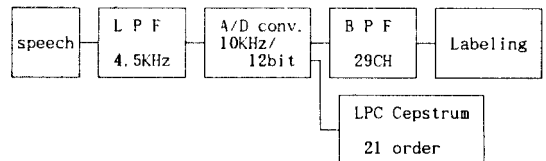


그림 1. 특징 파라미터 추출의 흐름도

3. 실험 및 고찰

3.1 기준 프레임 추출

인식 실험에 있어서는 3인 화자와 2회의 발성을 표준 패턴으로 하고 나머지 1회의 발성을 입력으로 통계적 인식 방법인 베이지 결정 법칙(Bayesian decision rule)⁽⁶⁾을 이용하여 인식 실험을 행하고 3회 반복한 평균을 인식률로 한다.

실험 실험과 분산 분석의 결과를 참고로 하여 인식에 유효한 특징추출의 위치를 조사하기 위해서는 음소군별로 각 음소의 시단으로부터 3프레임, 제 2프레임으로부터 4프레임과 같은 순서로 시단으로부터 15프레임까지 시점을 1프레임씩 옮겨 가면서 3프레임적에 대한 시간방향 특징에 대해 인식 실험을 실시하였다.

그림 3에 30ms에 대한 특징 추출 위치별 인식률의 변화를 나타내었으며, 표 3에 각 음소별 기준 프레임(3프레임)을 추출한 결과를 나타낸다.

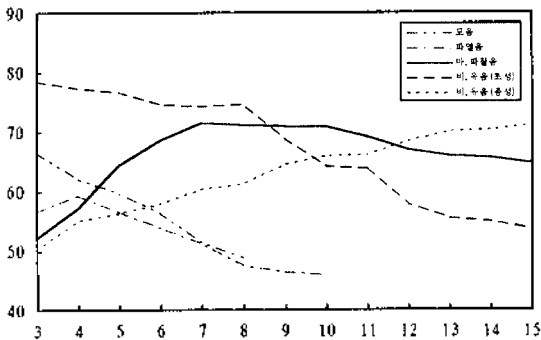


그림 3. 30ms에 대한 특징 추출 위치별 인식률의 변화.

표 3. 각 음소별 기준 프레임.

음 소	기준프레임 (3프레임)	음 소	기준프레임 (3프레임)
ㅏ (a)	1~3	ㅓ (p)	3~5
ㅑ (ä)	1~3	ㅕ (t)	1~3
ㅓ (o)	2~4	ㅋ (k)	3~5
ㅕ (u)	3~5	ㅗ (m)	4~6
ㅗ (u)	2~4	ㄴ (n)	3~5
ㅣ (i)	1~3	ㄹ (r)	1~3
ㅓ (e)	1~3	ㅛ (m*)	13~15
ㅕ (e)	2~4	ㄴ (n*)	13~15
ㅓ (p)	2~4	ㄹ (l*)	13~15
ㅓ (p')	1~3	ㅇ (g)	12~14
ㅕ (ph)	2~4	ㅜ (c)	5~7
ㅕ (t)	2~4	ㅜ (c')	7~9
ㅕ (t')	1~3	ㅜ (ch)	7~9
ㅕ (th)	2~4	ㅜ (s)	6~8
ㅕ (k)	3~5	ㅜ (s')	8~10
ㅕ (k')	2~4	ㅎ (h)	3~5
ㅕ (kh)	4~6		

그림 3과 표 3으로부터 모음의 경우는 시단 프레임에서 5프레임사이, 파열음은 시단에서 5프레임사이, 마찰음과 파찰음은 5프레임으로부터 10프레임사이, 그리고 비음과 유음의 경우 초성은 시단에서 6프레임, 종성은 중단으로부터 전 4프레임 구간이 인식에 유효한 정보가 가장 많이 포함되어 있음을 알 수 있었다.

또한, 위의 결과로부터 인식에 적당한 시간방향의 길이를 알아보기 위해서 비교적 인식에 공헌도가 높은 구간, 즉, 모음의 경우는 시단에서 8프레임, 파열음의 경우는 시단에서 10프레임, 마찰음과 파찰음의 경우는 시단에서 15프레임, 비음과 유음의 경우는 시단에서 15프레임 사이의 구간에 대해서 인식 실험을 실시하였다.

이때 기준 프레임은 포함한 구간 추출은 다음과 같이 한다.

- 1) 모음은 기준 프레임은 포함한 시단에서부터 추출하여 인식 실험을 하였다.
 - 2) 파열음은 기준 프레임 전 1프레임으로부터 프레임은 증가시키면서 추출한다.
 - 3) 마찰음, 파찰음의 경우는 기준 프레임 전 1프레임으로부터 프레임 증가시키면서 추출한다.
 - 4) 비음과 유음의 경우는 초성은 기준 프레임은 포함한 시단부터 프레임수를 증가시키면서 추출하고, 종성은 중단에서부터 전 방향으로 프레임은 증가시키면서 추출하였다.
- 이상의 결과와 시단에서 3프레임을 기준으로 2프레임씩 증가시키면서 인식 실험을 행한 결과를 비교 검토하였다. 이 결과를 각각 그림 4, 5, 6, 7에 나타내었다.

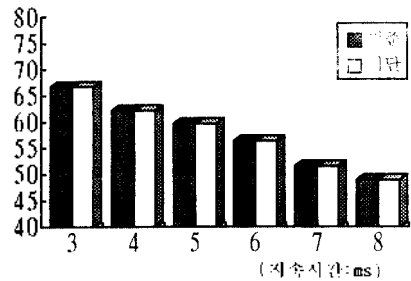


그림 4. 모음의 평균 인식률.

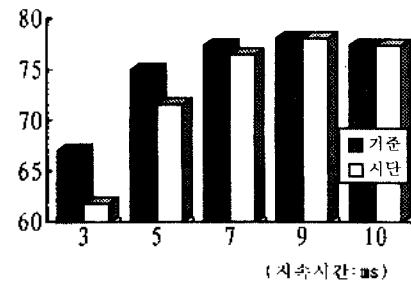


그림 5. 파열음의 평균 인식률.

한국어 음소 인식을 위한 기준 프레임 추출

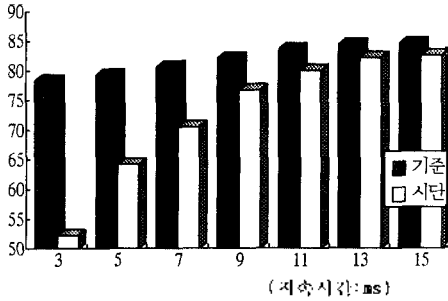


그림 6. 마찰음과 과찰음의 평균 인식률.

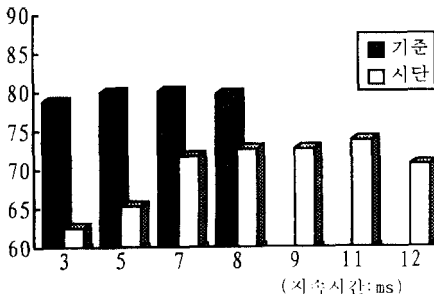


그림 7. 비음과 유음의 평균 인식률.

위의 그림으로부터 다음과 같은 사실을 알 수 있었다.

1) 모음의 경우 시단에서 5프레임, 2)파열음은 기준 프레임 전 1프레임으로부터 9프레임, 3)마찰음과 과찰음의 경우, 기음/ㅎ은 시단 프레임에서 14 프레임, 그외 음소는 기준 프레임 전 1프레임에서부터 14프레임, 4)비음과 유음의 경우는 초성은 시단 프레임에서 7프레임, 종성은 중단 프레임의 기준으로 전방향으로 7프레임에서 인식률이 높게 나타나 시단에서 프레임수를 증가시키면서 인식 실험을 행하는 경우보다 효과적인 결과를 알 수 있었다.

또한, 한국어 전 음소를 대상으로 기준 프레임의 중심으로부터 프레임수를 증가시키면서 인식 실험을 행한 결과 인식률의 향상을 기대할 수 있었다.

4. 결론

단음절에 포함된 전 음소를 대상으로 분석과 인식 실험을 통하여 각 음소별 특징을 가장 잘 포함하는 기준 프레임(3프레임)을 추출하였고, 음소군별로 보면 모음의 경우는 시단 프레임에서 5프레임사이, 파열음은 시단에서 5프레임사이, 마찰음과 과찰음은 5프레임으로부터 10프레임사이, 그리고 비음과 유음의 경우 초성은 시단에서 6프레임, 종성은 중단으로부터 전 4프레임 구간으로 나타났다.

이를 이용하여 음소군별로 기준 프레임을 포함한 인식 실험을 실시하여 다음과 같은 사실을 확인하였다.

1)모음의 경우 시단에서 5프레임, 2)파열음은 기준 프레임 전 1프레임으로부터 9프레임, 3)마찰음과 과찰음의 경우, 기음/ㅎ은 시단 프레임에서 14 프레임, 그외 음소는 기준 프레임 전 1프레임에서부터 14프레임, 4)비음과 유음의 경우는 초성은 시단 프레임에서 7프레임, 종성은 중단 프레임의 기준으로 전방향으로 7프레임에서 인식률이 높게 나타나 이 부분의 정보가 인식에 중요함을 알 수 있었다.

참고문헌

- 1) H.Y.Chung, S.Makino and K.Kido, "Analysis and recognition of Korea isolated vowels using formant frequency", J. Acoust. Jpn, 9, 5, pp. 225-232(1988).
- 2) 鄭鉉烈外, "韓國語破裂子音의 分析", 日本音響學會, 1-3-3 (Oct. 1988)
- 3) 鄭鉉烈外, "韓國語語頭破裂子音의 認識", 日本音響學會 1-2-21 (Mar. 1989)
- 4) 정석재, "마찰음 및 과찰음의 분석과 인식", 영남대학 석사학위논문, (1991)
- 5) 김병국, 정현일, "비음과 유음의 인식을 위한 특징추출에 관한 연구", 한국음향학회 영남지부 학술발표회 논문집, (1994)
- 6) 김병국, "한국어 난음절에 포함된 음소 인식에 관한 연구", 영남대학 석사학위논문(1992)
- 7) 권영욱, "MDS법을 이용한 한국어 음소 분석", 영남대학 석사학위논문(1991)
- 8) 岡田美智男, 牧野正三, 城戸健一, "スペクトルの時間變化パターンによる語頭音素の識別", 日本電子情報通信學會誌 J70A, 8, 1174-1185 (1987)
- 9) 中川聖一, "確率モデルによる音素認識", 電子情報通信學會 (1988)
- 10) 岡田美智男, "語頭音素의 統計的 分析과 認識에 關する 研究", 東北大學 博士學位 論文, 15-17 (1984)
- 11) Toshiro Haga, Shigeji Hashimoto, "回歸分析と主成分分析"(1986)