

유/무성음 결정에 따른 가변적인 시간축 변환

“손 단 영”, 김 원 구*, 윤 대 희*, 차 일 환*
*연세대학교 전자공학과, **군산대학교 전기공학과

Variable Time-Scale Modification with Voiced/Unvoiced Decision

Dan-Young Son*, Weon-Goo Kim**, Dae-Hee Youn*, Il-Whan Cha*
Dept. of Electronics Eng. Yonsei Univ.
Dept. of Electrical Eng. Kunsan National Univ.

Abstract

In this paper, a variable time-scale modification using SOLA(Synchronized OverLap Add) is proposed, which takes into consideration the different time-scaled characteristics of voiced and unvoiced speech. The conventional method performs time-scale modification at a uniform rate for all speech.

For this purpose, voiced and unvoiced speech duration at various talking speeds were statistically analyzed. A clipping autocorrelation function was applied to each analysis frame to determine voiced and unvoiced speech to obtain respective variation rates. The results were used to perform variable time-scale modification.

To evaluate performance, a MOS test was conducted to compare the proposed voiced/unvoiced variable time-scale modification and the uniform SOLA method. Results indicate that the proposed method produces sentence quality superior to that of the conventional method.

I. 서론

음성 변환(voice transformation)[1-3]은 음성신호를 표현할 수 있는 몇개의 특징 변수(parameter)를 분석 및 변화시킴으로써 음성을 인위적으로 합성해내는 것을 말한다. 특히, 시간축 변환(time-scale modification)은 음성을 특징지우는 여러 변수들 중에서 소리의 길이를 변화시켜주는 것으로써, 단기간 Fourier 해석에 의한 방법[4], 정현파 모델에 의한 방법[5,6], SOLA(Synchronized OverLap and Add) 방법[1], PSOLA(Pitch Synchronized OverLap and Add) 방법[7] 등과 같이 여러 방법으로 연구되어 왔다. 음성 신호의 시간축 변환은 청각 장애, 언어 장애가 있는 사람을 위한 시스템과 언어 학습을 위한 시스템[9] 등에 이용될 수 있고, 미트웬을 줄이기 위한 음성 부호화 시스템[8]에 응용될 수 있다.

여러가지 시간축 변환 방법 중 SOLA 방법은 음질 저하를 줄이는데 효과적이며 적은 계산량으로 우수한 성능을 갖는다[1]. 그러나, SOLA 방법은 효과적으로 고음질의 변환신호를 합성하지만, 변환 비가 클 경우에는 사람이 느리게 혹은 빨리 말할 때와는 다른 부자연스러운 소리로 합성하게 된다. 느리게 말하는 경우 사람의 발음 특성상 유성음은 길게 늘어나지만 무성음은 비교적 적은 비율로 늘어나는데, 기존의 시간축 변환에서는 SOLA 방법을 이러한 현상에 대한 고려없이 일률적으로 적용하여 음성신호를 늘여주거나 줄여주기 때문이다.

본 연구에서는 유성음 구간과 무성음 구간에 다른 길이의 비로 SOLA 방법을 적용하여 음성 신호를 압축하거나 신장하는 가

변적인 시간축 변환을 제안하였다. 유/무성음의 결정은 적은 계산량으로 좋은 결과를 갖는 클리핑 자기 상관 함수 방법(Modified Autocorrelation with Clipping Method)[3]을 이용하였다. 유/무성음 각각의 압축 또는 신장 비율을 판단하기 위해 특정 문장을 사람이 속도를 변화시키거나 말음하도록 하여 통계적 특성을 조사하였다. 얻어진 결과를 이용하여 유/무성음에 따라 가변적으로 SOLA 방법을 적용해 시간축 변환을 함으로써 합성 신호를 얻었다. MOS 테스트를 통해 구현된 알고리즘의 성능을 평가하였다.

II장에서는 SOLA 방법의 이론을 살펴보고, 유/무성음 판단 방법에 대해 설명한 후, SOLA를 가변적으로 적용하기 위한 시간축 변환 시스템을 제안하였다. III장에서는 말음속도 변화에 따른 유/무성음 길이 변화의 통계적 특성을 조사하여 실제로 SOLA 방법에 가변적으로 적용한 다음, MOS 테스트로 성능향상을 확인해 보았다. 또한, 위에 자연스럽게 들리는 유/무성음 변화 비율을 조사하였다. 그리고, 마지막으로 IV장에서 결론을 맺었다.

II. 유/무성음 결정에 따른 가변적인 시간축 변환

음성신호의 시간축 변환은 음성신호의 기본 주파수와 성도 모델 스케프링을 보존하여 원래의 신호특성은 그대로 유지하면서, 말음속도만 변화시키는 것이다[1,4,7]. 변환된 음성은 원래의 음성보다 빠르게 또는 느리게 발음되는 것처럼 들리게 된다. 기존에 제안된 시간축 변환 방법[1,4-7] 중에서 본 연구에서는 적은 계산량으로 우수한 성능을 나타내는 SOLA 방법을 이용하였다 [1].

SOLA 방법은 LSEE-MSTFTM(Least Squared Error Estimation-Modified Short Time Fourier Transform Magnitude) 알고리즘에 기본을 두고있는 방법인데, LSEE-MSTFTM 알고리즘은 반복적인 계산을 통해서 변환된 신호의 Fourier 변환 크기와 원 신호의 Fourier 변환 크기의 차이를 최소화시키 시간적으로 변형된 고음질의 음성신호를 만드는 최소 차등 오차 추정 방법이다[2]. 그림 1에서처럼 임펄스 신호 $x(n)$ 이 비례상수 $\alpha = S_1/S_2$ 에 의해 속도가 바뀐 신호 $y(m)$ 으로 변환될 경우를 예로 들어 두 방법을 비교해 보면 다음과 같다(여기서 S_1 는 분석구간의 이동값이고 S_2 는 합성구간의 이동값으로, $\alpha > 1$ 보다 크다는 것은 말음속도가 늦어지는 것이고 1보다 작다는 것은 빨라지는 것이다).

LSEE-MSTFTM 알고리즘은 중첩가산(OverLap and Add) 과정을 포함하고 있어 고음질의 신호를 합성하기도 하지만, 원하는 신호

유/무성음 결정에 따른 가변적인 시간축 변환

로 완벽히 수렴하기 위해서 보통 100의 정도의 많은 반복 계산을 필요로한다. 그러나, SOLA 방법은 5회 이하의 적은 반복계산으로도 수렴을 보장할 수 있다. 그림 1(b)에서 볼 수 있듯이 LSEE-MSTFTM 알고리즘은 분석단에서 S_n 만큼씩 분석창을 이동하여 분석구간들을 얻고 이들을 S_n 간격으로 일괄적으로 재배치하여 중첩가산하는 과정으로 신호를 합성하게 되므로 초기 계산 결과로는 그림 1(b)의 (iv)와 같은 예상 밖의 신호를 얻게 되지만, SOLA 방법은 중첩가산하기 전에 수렴이 보장되는 초기 추정치를 연속되는 두 신호구간 사이의 동기가 일치하는 점으로 구하여 합성구간을 이 위치로 재배치한 뒤 더해주기 때문이다. 신호구간 사이의 동기가 일치하는 점은 상호 상관 함수가 최대가 되는 k 로 찾아 준다. 그림 1(c)는 적당한 k 에 의해 연속되는 합성구간의 이동이 조정되어 더해지는 과정을 보여 주며, 이때 합성신호는 중첩가산된 S_n 개 샘플이 차례로 정규화되어 출력된다. 원신호와 일치하는 평균된 결과 파형이 그림 1(c)의 (iv)와 같이 얻어진다[1,2].

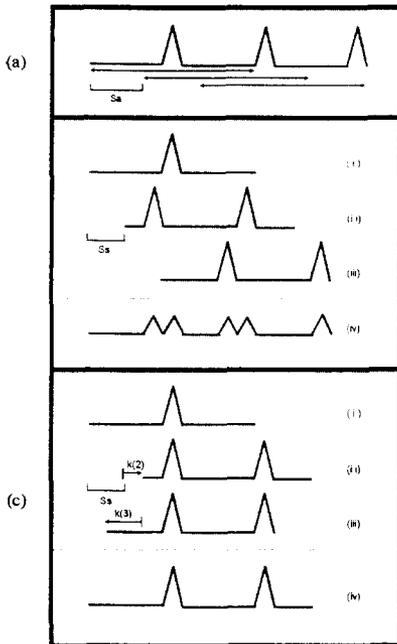


그림 1. 시간축 변환을 위한 LSEE 방법 및 SOLA 방법
(a) 원 신호 (b) LSEE 방법 (c) SOLA 방법

기존의 시간축 변환은 SOLA 방법을 음성 신호를 일괄적으로 압축하거나 신장시키는 방법으로 적용해왔다. 이 경우, 고음질의 변환신호가 합성되기는 하지만, 큰 비율로 신호를 늘이거나 줄일 때에는 사람이 듣기에 부자연스러운 소리를 합성하게 된다. 이것은 사람이 느리게거나 빠르게 발음했을 경우에는 사람의 발음기관 특성상 음성음의 길이는 많이 변화하지만 무성음의 길이 변화량은 상대적으로 적은 현상을 고려하지 않고, 유/무성음 구별없이 일괄적으로 시간축 변환을 적용하였기 때문이다. 그러므로, 본 연구에서는 유/무성음을 결정한 후에, 가변적으로 시간축 변환을 해주는 시스템을 제안하였다. 그림 2는 제안된 시스템의 블럭도이다. 신호의 분석은 신호의 특성이 일정하게 유지되는 구간(L)을 75% 중첩가산(이동값 S_n)으로 수행하였다.

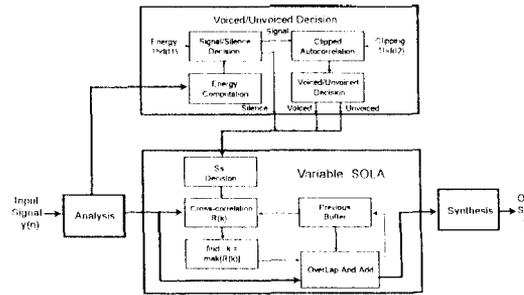


그림 2. 유/무성음 결정에 따른 가변적인 시간축 변환 시스템의 블럭도

유/무성음의 결정은 적은 계산으로도 효과적인 결과를 얻을 수 있는 클리핑 자기 상관 함수 방법을 이용하여 신호의 매 분석구간에서 이루어지게 되는데, 클리핑 상관 함수는 음성신호의 클리핑 문턱치를 찾아, 이를 기준으로 신호를 클리핑하고 그 결과 신호를 이용하여 자기 상관 함수를 계산하는 것이다[3]. 분석구간은 300샘플(L)로 한다. 그림 2의 블럭도에 나타난 것처럼, 유/무성음을 결정하기 전에 분석구간의 에너지를 구하여 에너지 레벨이 에너지 문턱치보다 작으면 현재 구간을 묵음구간으로 간주하고 유/무성음을 구별하지 않게 된다. 에너지 문턱치는 배경잡음 레벨을 기준으로하여 얻을 수 있다[3]. 클리핑 문턱치는 처음의 100샘플과 뒤의 100샘플 구간에서 각각 가장 큰 절대값을 갖는 값을 찾아, 두 값을 비교한 후, 작은 값의 64%로 잡아 준다. 클리핑 신호는 신호가 클리핑 문턱치보다 크면 +1을, 이것의 음수값보다 작으면 -1을, 그외의 경우는 0을 주어서 얻게 된다. 그림 3은, 클리핑 신호로 구한 자기 상관 함수를 정규화(normalization)시키고, 최대값을 찾아 준다. 구해진 최대값이 유/무성음 결정 문턱치인 0.3을 초과하면 현재 분석구간을 음성음 구간으로 선정하고, 반대의 경우는 무성음으로 결정한다[3].

유/무성음이 결정되면 실험을 통해 통계적으로 얻은 유/무성음의 길이 변화율에 따라 그림 2의 블럭도와 같이 SOLA를 가변적으로 적용한 시간축 변환을 할 수 있다. 사람이 라는 비율로 빠르게 혹은 느리게 발음한다고 가정할 때 음성음의 길이 변화율은 α 로, 무성음의 길이 변화율을 ω 로 표시하기로 한다(이렇게 방음한 경우에는 α 는 α 보다 큰 값이, ω 는 작은 값이 될 것이고, 빠르게 발음했을 경우는 반대가 될 것이다). 가변적인 SOLA의 적용이 일관적인 것까 다른 점은 후자의 경우 합성시의 이동값 S_n 가 고정되는데 반하여 전자는 음성음, 무성음인가 혹은 묵음인가에 따라 변화한다는 것이다. 즉, α 라는 비율로 발음속도가 변환될 때 일정한 적용시에는 S_n 가 $\alpha \cdot S_n$ 라는 한가지 값은 갖지만, 가변적으로 적용할 때는 유/무성음에 따라 $\alpha \cdot S_n$ 와 $\omega \cdot S_n$ 라는 각각 다른 값을 갖게 된다는 것이다. 그리고, 가변적인 적용의 경우 현재구간에서 결정된 S_n 만큼 신호를 이동시켜 전구간의 신호와 상호 상관이 최대가 되는 위치 k 를 찾아 동기를 맞춰 준다. 그리고, 동기가 일치하는 점으로 신호가 재배치되면 현재구간의 신호를 전 구간의 신호에 중첩가산해서, S_n 개 샘플은 정규화하여 출력하고, 나머지 신호는 S_n 만큼 이동시켜 저장한다. 이러한 과정을 음성신호 전 구간에 걸쳐 수행하면 원하는 합성신호를 얻을 수 있다. 합성된 신호는 가변적인 S_n 만큼씩 정규화되어 출력됨으로써 결과적으로 유/무성음의 길이가 다른 비율로 변환된다. 그러나, 유/무성음의 길이 변화율을 다르게 해주기 때문에, 원신호

에서 유/무성음이 차지하는 퍼센트가 같다면 합성 신호의 길이가 일률적으로 변환해 줄 때와 같게 되겠지만 다른 경우에는 일률적인 변환과는 다르게 될 것이다. 예를 들어, 일률적인 변환을 할 때의 합성 신호의 길이가 (1)의 L_s 일 때, 유성음이 무성음에 비해 많이 포함되어 있다면 가변적인 변환의 전체 길이 L_{ms} 는 L_s 보다 긴 (2)가 될 것이다.

$$L_s = \alpha * S_a * T \quad (1)$$

$$L_{ms} = \alpha * S_a * T_s + \alpha_v * S_a * T_v + \alpha_w * S_a * T_w \quad (2)$$

여기서, T 는 전체 프레임 수, T_s 는 목음의 프레임 수, T_v 는 유성음의 프레임 수이고, T_w 는 무성음의 프레임 수이다. 목음의 변화비율은 전체 변화비율과 같은 α 를 적용한다.

원신호에서 유/무성음이 차지하는 퍼센트가 다르면 가변적인 변환시의 전체 길이 L_{ms} 가 일률적인 변환시의 전체 길이 L_s 와 다르게 되는데, L_{ms} 를 L_s 와 일치시키기 위해서는 비례상수 $m = L_s/L_{ms}$ 를 L_{ms} 에 곱해 주면 된다. 그렇다면, (2)는 다음과 같이 변형될 것이다.

$$L_s = m * L_{ms} = \alpha' * S_a * T_s + \alpha'_v * S_a * T_v + \alpha'_w * S_a * T_w \quad (3)$$

여기서, $\alpha' = m * \alpha$, $\alpha'_v = m * \alpha_v$ 이고, $\alpha'_w = m * \alpha_w$ 이다.

그러므로, 변환 비율을 α 로 하여 가변적인 시간축 변환을 수행할 때 유/무성음의 길이 변화 비율은 통계적으로 얻어진 α , α_v 가 아닌, 조정된 α' , α'_v 를, 목음의 경우는 α' 를 이용하면 전체 길이가 α 만큼 변환된 원하는 신호가 얻어질 것이다.

III. 실험 및 결과 고찰

본 실험에서는 유/무성음 결정에 따른 가변적인 시간축 변환 알고리즘을 제안하여 구현하고 그 동작을 확인하고자 컴퓨터로 모의 실험을 시행하였다.

각 분석구간에서 유/무성음을 결정하기 위해, 클리핑 자기상관 함수 알고리즘을 이용하였다. 그리고, 실제로 사람이 느리게 혹은 빠르게 발음했을 때, 유/무성음의 길이가 각각 어떻게 변화하는가를 통계적으로 조사해 보았다. 얻어진 통계적 특성 결과를 이용해 가변적으로 시간축 변환을 하였으며, 이 결과와 SOLA를 이용하여 적용했을 때의 결과를 비교하여 MOS (Mean Opinion Score) 테스트를 행하였다. 그리고, 유성음의 변화율을 고정시키고, 무성음의 변화율을 바꿔가면서 귀로 자연스럽게 들리는 유/무성음 변화 비율을 조사하였다.

모의 실험에 사용된 음성 데이터는 25세 남성화자가 비교적 조용한 연구실 환경에서 녹음하였다. 실험에 사용된 문장들은 아래와 같다. 유/무성음의 길이가 다른 비율로 변환하는 것을 확인하기 위하여, 문장 1), 2)는 유/무성음이 골고루 포함된 것을 보통속도, 2배 빠르게, 2배 느리게 각각 5번씩 발음하게 하였다. 그리고, 실제적인 적용을 위하여 유/무성음 결정 알고리즘을 이용하여 유/무성음 길이 변화율이 특정 속도에서 각각 어떤 변화율을 갖는가를 조사하였는데, 음성 데이터는 한국에서 사용 빈도가 높은 음소[10]를 골고루 포함하고있는 문장 3)~7)을 0.7, 1.0, 1.3, 1.5, 1.8배 속도로 녹음하여 사용하였다. MOS 테스트를 위하여 통계적 조사에서 사용되지 않은 유/무성음이 골고루 포함된 문장 8), 9)를 보통속도로 녹음하여 참가하였다.

- 1) 차의 범도는 호프러린 마음가짐을 바로 잡아 준다.
- 2) 새터는 교통이 편리한 커다란 대학촌이다.
- 3) 유성은 인간과 기계 사이에 가장 효과적인 정보 전달 수단의 하나이다.

- 4) 새터는 문화적으로 본다면 서울의 중심부이다.
 - 5) 현대는 문화 예술의 시대이다.
 - 6) 차의 범도는 사육의 호프러린 마음가짐을 바로 잡아 준다.
 - 7) 오늘은 어제의 열매이며 내일의 씨앗이다.
 - 8) 분수처럼 흘러지는 푸른 종소리.
 - 9) 파란 이파리 사이로 함초롬한 꽃방울이 피어난다.
- 음성 데이터는 4.5KHz 차단 주파수(cut-off frequency)를 갖는 저역 통과 필터에 통과 시킨 후 16비트 PCM으로 A/D 변환하여 실험에 사용하였다. A/D 변환시의 샘플링 주파수는 10KHz로 하였다. 음성신호의 분석은 300샘플(L) 단위로 분석구간을 잡아, 75% 중첩가산(분석사의 이동값 $S_s = 75$)으로 수행하였다.

1. 유/무성음 결정

본 실험에서는 유/무성음 결정을 위해, 선처리 작업으로 배경 잡음의 평균을 구한 후 이 값에 3dB를 더한 값을 그림 2의 목음과 신호를 구별하는 문턱치(T)로 잡아주었다. 각 구간의 에너지 계산 전에는 $1-0.95z^{-1}$ 의 선처리(pre-emphasis) 과정을 거쳐 고주파 부분은 강조하였다. 유/무성음의 결정은 우선 현재의 분석구간 에너지를 구하여 에너지 문턱치를 넘으면 신호로, 그렇지 않으면 목음으로 결정해 주고, 신호라고 판단되면 클리핑 자기 상관을 계산하여 유/무성음을 결정해 주었다. 클리핑 문턱치(T_2)는 분석구간의 신호에서 앞의 100샘플 구간과 뒤의 100샘플 구간 각각에 자기 절대값의 최대치를 찾아, 두 값 중 작은 값의 64%로 잡아주었다. 유/무성음 결정 문턱치는 정규화된 자기 상관의 30%값인 0.3으로 하였다. 얻어진 결과의 변화율을 줄이기 위해 5-포인트 미디언 필터를 통과시켜 주었다. 최종적으로 얻어진 결과에 따라 유성음에는 2, 무성음, 목음에는 각각 1, 0 값을 주었다.

위의 유/무성음 결정 알고리즘은 문장 3)~7)에 적용되었다. 그림 3는 결정 결과 중 문장 6)에 해당하는 것으로, 위쪽이 문장의 파형이고, 아래쪽이 유/무성음의 결정 결과 파형이다. 0으로 나타나는 구간이 목음구간이고, 1 값을 갖는 구간이 무성음, 2 값을 갖는 구간이 유성음 구간이다. 한 두 구간에서 오차가 발생하였지만 비교적 좋은 결정 결과를 보여 본 실험에서 제안한 가변적인 시간축 변환을 수행하기에 적합함을 알 수 있다.

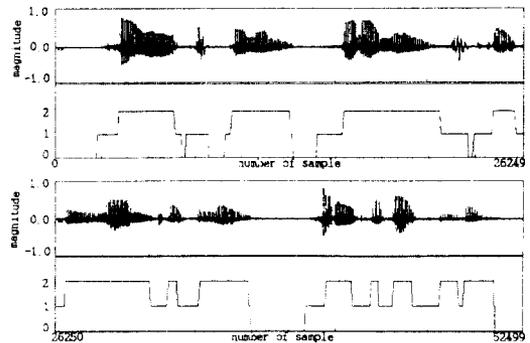


그림 3. 문장 6에 대한 유/무성음 결정 결과 (유성음 : 2, 무성음 : 1, 목음 : 0)

2. 음성의 통계적 특성

사람이 속도를 다르게하여 발음하면 무성음과 유성음이 다른 비율로 늘어나게 되는 현상을 고려하여, 본 실험에서는 실제로 사람이 느리게 혹은 빠르게 발음했을 때, 유/무성음이 각각 어떻게 변화하는가를 통계적으로 조사해 보았다. 우선, 사람의 발음

유/무성을 결정에 따른 가변적인 시간축 변환

속도가 변할 때 유/무성의 변화 비율이 다름을 확인하기 위해 음소를 직접 손으로 검출하여 통계해 보았다. 그리고, 실제적인 적용을 위하여 유/무성을 결정 알고리즘을 이용하여 유/무성을 각각의 변화하는 비율이 특정 속도에서 어떤 변화율을 갖는지 조사하였다.

2.1. 음소 구분을 통한 유/무성을 통계적 특성

문장 1), 2)를 보통 속도, 2배 빠르게, 2배 느리게 각각 5번씩 반복하게 하여 얻어진 음성 데이터에서 각 음소를 손으로 직접 검출하는 방법으로 유성음과 무성음의 길이 변화율을 통계적으로 조사하여 보았다[10].

음소마다 늘거나 줄어드는 비율이 다르지만 크게 유성음과 무성음으로 분류하여 각각에 대한 평균을 보면, 2배 이상 느리게 발음했을 때 유성음은 3배 이상까지도 늘어나고 무성음은 1.5배 정도까지만 늘어났다. 빠르게 발음했을 경우에는 유성음은 0.7배 정도로 줄었으나 무성음은 0.8배 정도까지만 줄었다.

2.2. 유/무성을 결정 알고리즘을 이용한 통계적 특성

본 실험에서는 유/무성음 결정 알고리즘을 사용해 기계적으로 유/무성음을 결정한 후 이 결과의 통계적 특성을 조사해 보았다. 사용된 문장은 3)~7)이다.

유/무성음의 늘어난 비율은 1.30, 1.52, 1.83배 속도 변화시 각각 1.12/1.37, 1.30/1.61, 1.41/2.01 등으로 유성음이 무성음보다 큰 비율로 변화함을 확인할 수 있었다. 그리고, 이 결과를 도시한 그림 4의 그래프에서는 유성음과 무성음의 변화율을 보이는 직선이 비교적 선형적인 특성을 나타내며 특히 유성음은 무성음보다 큰 기울기를 가지고 변화함을 알 수 있었다.

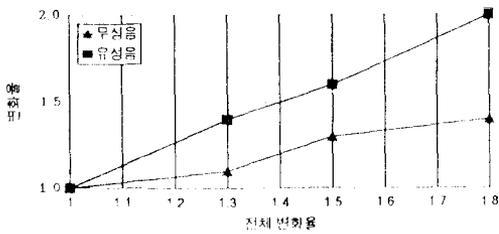


그림 4. 속도에 따른 유/무성음 변화율 그래프

3. 유/무성을 결정 결과를 적용한 시간축 변환

앞 절의 통계적 조사에서 얻어진 결과와 유/무성음 결정 결과인 SOLA 방법에 적용해 가변적인 시간축 변환을 해보았다. 사용된 문장은 앞 절에 인용된 문장들 중 3), 6)과 테스트를 위해 첨가한 문장 8), 9) 등 4 문장이다.

전체 과정이 0.7, 1.3, 1.5, 1.8배가 되도록 시간축 변환을 하기 위해, 각각의 비율에 대한 유/무성음의 길이 변화 비율은 앞 절의 실험에서 얻어진 0.7/0.9, 1.4/1.1, 1.6/1.3, 2.0/1.4를 사용하였다. 그림 5 (b) : 문장 6)의 가변 합성 결과이다.

4. MOS 테스트

3절에서 얻은 결과들 가지고 SOLA를 일괄적으로 적용했을 경우와 가변적으로 적용했을 경우를 MOS 테스트로 비교하는 실험을 하였다.

1) 실험 대상 : 음성 신호 처리 연구과 실험에 대한 경험이 많은 대학원생 15명

(연령은 24세에서 31세 사이, 성별은 남성 14명에 여성 1명)

2) 선호도 조사 : 0.7, 1.3, 1.5, 1.8 배 속도에 기존의 SOLA 방법과 제안된 방법으로 얻은 합성 결과 비교

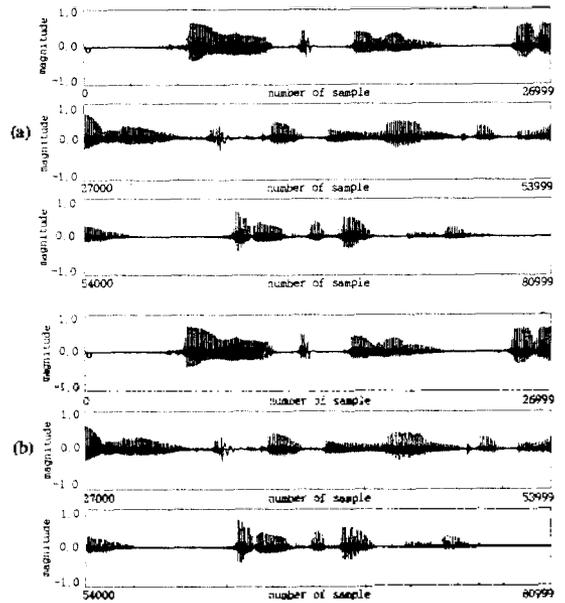


그림 5. SOLA 및 가변 시간축 변환 (문장 6에 대한)

(a) SOLA에 의한 시간축 변환 (1.5 배 속도)

(b) 제안된 가변 시간축 변환 (유성음 1.6 배, 무성음 1.3 배)

기존의 MOS 테스트가 음질의 좋고 나쁨을 평가하는 것인데 본 실험은 유/무성음 결정 결과에 가지고 가변적으로 시간축 변환을 해준 것이 얼마나 변환된 결과보다 자연스럽게 들리는지를 확인하려는 것이 목적이었으므로, 음질이 좋고 나쁨이 아닌 자연스러움은 기준으로 좋고 나쁨을 선명하게 하겠다[11]. 자연스러움의 기준을 위해 테스트 전에 사람이 보통 속도로 자연스럽게 발음한 것을 들려주고 테스트를 시작하였다. 테스트 결과, 자연스럽게 들리는 것에는 1의 값을 그렇지 않은 것에는 0의 값을 주어 주어진 다음, 이러한 결과로부터 선호도 (제안된 방법에 대한)를 피싱의 값으로 표현하여 표 1에 나타내었다. 문장에 따라 좋고 나쁨의 선호도에 차이가 있는데, 이것은 유성음인가 무성음인가에 따라지만이 아니고 음소별로도 속도 변화 대상이 다르기 때문이다. 전체 결과를 보면 1.5배 전정서에 제안된 방법의 성능이 85%로 가장 좋게 나온을 알 수 있고, 1.3 배 전정서 78.3%, 1.8배 전정서 68.3%로 두 경우 모두 비교적 좋게 나오고, 0.7배 압축시에도 66.7%로 좋은 결과가 나왔다. 그러므로, 위의 같은 MOS 테스트 결과로 사람이 듣기에도 유/무성음 결정 결과에 따라 차별적으로 시간축 변환을 해주는 것이 더 자연스러움은 알 수 있다.

표 1. MOS 테스트 결과 (통계적 특성에 따라 변화시킨 결과)

무성음 변화율	문장 3	문장 6	문장 8	문장 9	평균
0.7	66.7	86.7	53.3	60	66.7
1.3	93.3	80	73.3	66.7	78.3
1.5	86.7	100	66.7	86.7	85
1.8	73.3	80	53.3	66.7	68.3

위의 결과는 유/무성음을 구별하여, 사람이 속도를 달리하여 발음하였을 때와 같이 무성음은 유성음보다 작은 비율로 변화시

킬 때 사림이 들기에 자연스럽다는 것을 보여주고 있다. 다음으로, 유성음이 일정한 비율로 변화할 때 귀는 무성음이 얼마의 비율로 변화하면 자연스럽다고 느끼는가를 정취실험을 통해 조사해 보았다. 앞 결과 같은 문장, 기준, 그리고, 대상으로 유성음을 전부 2배로 늘려 주고, 무성음을 변화시키지 않았을 때, 앞 절에서 얻은 통계적 비율인 1.4배로 가변적으로 변환시켰을 때, 그리고, 마지막으로 2배로 일괄적으로 변환시켰을 때의 새가시름 비교하여 MOS 테스트를 수행하였다. 각각의 방법에 대한 전체적인 선호도 결과를 표 2에 보였다. 무성음을 변화시키지 않았을 때가 70%로 가장 좋고, 다음은 1.4배로 변화시켰을 때 30%였으며, 일괄적인 변환은 0%로 전혀 선호도를 보이지 않았다. 이것으로 무성음의 변화는 저울수복 자연스럽게 들림을 알 수 있다.

표 2. MOS 테스트 결과 (무성음의 경우만 변화시킨 경우)

유성음 변화율	1.0 배	1.4 배	2.0 배
선호도 (%)	70.0	30.0	0.0

V. 결 론

본 연구에서는 시간축 변환을 할 때, 사람의 발음 특성상 무성음이 유성음보다 크게 들어난다는 현상을 고려하여, 기존의 시간축 변환이 음성신호를 일괄적으로 압축 또는 신장하는 것에 반하여, 유/무성음을 결정한 후 이 결과에 따라 SOLA 방법을 가변적으로 적용하여 시간축 변환을 수행하는 방법을 제안하였다.

사람이 발음할 때 속도가 변화하면 실제로 유/무성음의 변화 비율이 다를 수 통계적으로 조사하여 보았다. 컴퓨터 자기 상관 함수 방법을 사용하여 유/무성음을 결정하고, 통계적으로 얻어진 변화 비율로써 가변 시간축 변환을 수행하였다. 제안된 방법의 성능을 평가하기 위하여, SOLA 방법을 가변적으로 적용한 시간축 변환 결과와 일괄적으로 적용한 변환 결과를 MOS 테스트로 비교하였고, 높은 자지도를 얻어 제안된 방법의 성능이 우수함을 확인하였다. 그리고, 무성음의 변화는 저울수복 자연스러움을 알 수 있었다.

앞으로의 연구 과제로 음소들 각각의 속도에 따른 변화 비율을 조사하여 각 음소별로 세분화된 변환 기준을 마련하고, 이것을 시간축 변환에 이용한다면, 제안된 방법 이상의 성능 향상을 기대할 수 있을 것이다.

참고 문헌

[1] Salim Roucos and Alexander M. Witgud, "High Quality Time-Scale Modification for Speech," in Proc. ICASSP, pp.493-496, Apr. 1986.

[2] D. W. Griffin and J. S. Lim, "Signal Estimation from Modified Short-Time Fourier Transform," *IEEE Trans Acoust., Speech, Signal Processing*, vol. ASSP-32., No. 2, pp.236-243, Apr. 1984.

[3] Il-Hyun Nam, *Voice Personality transformation*, Ph. D. Thesis, Rensselaer Polytechnic Institute, Jan. 1991.

[4] Michael R. Portnoff, "Time-Scale Modification of Speech Based on Short-Time Fourier Analysis," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-29, No. 3, pp. 374-390, Jun. 1981.

[5] Thomas F. Quatieri and Robert J. McAulay, "Shape Invariant Time-scale and Pitch Modification of Speech," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 40, No. 3, pp. 497-510, Mar. 1992.

[6] Thomas F. Quatieri and R. J. McAulay, "Speech Transformations Based on a Sinusoidal Representation," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-34, No. 6, pp. 1449-1464, Dec. 1986.

[7] E. Moullines and F. Charpentier, "Pitch-Synchronous Waveform Processing Technique for Text-to-speech Synthesis using Diphones," *Speech Communication*, vol. 9 (5/6), pp. 453-467, 1990.

[8] John Makhoul, "Time-scale Modification in Medium Low Rate Speech Coding," in Proc. ICASSP, pp.33.7.1 ~ 33.7.4, 1986.

[9] Tohbu TAKAGI and Eiichi MIYASAKA, "A Speech Prosody conversion with a High Quality Speech Analysis-Synthesis Method", *EUROSPEECH*, vol. 2, pp. 991-994, Sep. 1993.

[10] 한국 방송공사, "표준 한국어 발음 대사전," 이문각, 1993.

[11] D. Sinha and A. H. Tewfik, "Low Bit Rate Transparent Audio Compression using Adaptive Wavelets," *IEEE Trans Acoust., Speech, Signal Processing*, vol. ASSP-41, No.12, pp. 3463-3477, Dec. 1993.

* 본 논문은 1993년도 한국과학기술단의 연구비 지원으로 이루어진 것임