

새로운 스펙트럼 완만화에 의한 합성 음질 개선

(Improvement of Synthetic Speech Quality using a New Spectral Smoothing Technique)

장 효 증 [†] 최 형 일 ^{**}
(Hyo-Jong Jang) (Hyung-Il Choi)

요 약 본 논문에서는 단위음소로 다이폰을 사용하여 음성을 합성하는 방법에 관하여 기술한다. 음성 합성은 기본적으로 단위음소들의 연결을 통하여 이루어지는데, 이때 발생하는 가장 큰 문제점은 두 단위음소 사이의 연결부분에서 불연속이 발생하는 것이다. 이 문제를 해결하기 위하여 본 논문에서는 포먼트 궤적뿐 아니라 스펙트럼의 분포특성과 인간의 청각적인 특성을 반영하여 스펙트럼을 완만화하는 방법을 제안한다. 즉, 제안하는 방법은 단위음소의 연결 구간에서 인간의 청각신경 특성을 고려하여 완만화의 양과 범위를 결정한다. 다음, 두 다이폰 경계의 스펙트럼 분포를 시간에 따라 가중치를 다르게 주어 스펙트럼 완만화를 수행한다. 이 방법은 불연속을 제거하며 완만화로 인하여 발생할 수 있는 음성의 왜곡을 최소화한다. 제안하는 방법의 성능을 평가하기 위하여 ETRI 음성 DB 샘플과 개인별로 자체 녹음한 총 20여개의 문장에서 추출한 약 500여 개의 다이폰에 대하여 실험을 수행하였다.

키워드 : 포먼트 궤적, 스펙트럼 완만화, 청각적인 특성, 음성합성

Abstract This paper describes a speech synthesis technique using a diphone as an unit phoneme. Speech synthesis is basically accomplished by concatenating unit phonemes, and it's major problem is discontinuity at the connection part between unit phonemes. To solve this problem, this paper proposes a new spectral smoothing technique which reflects not only formant trajectories but also distribution characteristics of spectrum and human's acoustic characteristics. That is, the proposed technique decides the quantity and extent of smoothing by considering human's acoustic characteristics at the connection part of unit phonemes, and then performs spectral smoothing using weights calculated along a time axis at the border of two diphones. The proposed technique reduces the discontinuity and minimizes the distortion which is caused by spectral smoothing. For the purpose of performance evaluation, we tested on five hundred diphones which are extracted from twenty sentences using ETRI Voice DB samples and individually self-recorded samples.

Key words : formant trajectory, spectral smoothing, acoustic characteristic, speech synthesis

1. 서 론

정보 통신의 발달로 인간은 컴퓨터를 이용하여 다양한 방법으로 정보를 교환하고 있다. 보다 신속하고 정확한 정보 교환을 위하여 인간과 컴퓨터 사이의 의사 교환 또한 중요한 문제로 대두되고 있다. 이러한 관점에서 음성은 인간의 가장 자연스러운 의사 전달 수단이므로 음성을 통한 컴퓨터와의 정보교환 기술은 매우 중요하

다. 음성을 통한 정보교환에는 인간의 음성을 문자정보로 변환하는 음성인식과 문자정보를 음성으로 출력하는 음성합성에 의한 방법이 있다. 본 논문에서는 음성합성에 관하여 기술한다.

음성합성시스템에 의하여 생성된 합성음이 자연스럽지 못한 경우는 청자로 하여금 이질감을 느끼게 하고, 부정확한 경우는 정확한 의사전달을 방해하게 된다. 따라서, 음성합성시스템의 궁극적인 목적은 정확한 정보를 전달하기 위하여 보다 자연스럽고 정확한 합성음을 생성하는데 있다. 현존하는 음성합성시스템 대부분은 기본적으로 음성데이터베이스로부터 단위음소를 추출하고 이를 연결하여 합성음을 얻는 방법을 취한다. 이 접근 방법의 대표적인 문제점은 단위음소의 연결 부분에서 불연속이 발생하는 것이다. 그 원인은 데이터베이스에

· 본 연구는 숭실대학교 교내연구비 지원으로 이루어졌음

[†] 학생회원 : 숭실대학교 컴퓨터학과
ozjh@vision.soongsil.ac.kr

^{**} 종신회원 : 숭실대학교 미디어학부 교수
hic@computing.soongsil.ac.kr

논문접수 : 2003년 4월 1일

심사완료 : 2003년 8월 2일

있는 단위음소들이 모든 경우에 대한 문맥적인 차이나 변화들을 나타낼 수 없기 때문이다.

단위음소의 연결 부분에서 발생하는 불연속을 해결하기 위한 기존의 주요 연구에서는 트라이폰과 같은 큰 합성 단위를 사용하는 방법[1], 스펙트럼 불연속이 최소화된 연결 위치를 찾아 이를 연결하여 불연속을 최소화하는 방법[2], 파형 또는 스펙트럼 완만화 등을 통해 스펙트럼의 불연속을 제거하는 방법[3] 등이 있다. 이러한 방법들의 문제점을 살펴보면 다음과 같다. 첫 번째 방법은 단위 음소의 연결 합성 시에 단위 음소로 트라이폰을 사용함으로써 모노폰이나 다이폰을 사용할 경우보다 불연속을 줄일 수 있으며 큰 단위를 사용함으로써 문맥에 민감하기 때문에 합성의 질을 높일 수 있다는 장점이 있다. 하지만 트라이폰을 사용해도 불연속이 완전히 없어지는 것이 아니라 빈도가 줄어들 뿐이다. 또한 큰 합성 단위를 쓰기 때문에 필요한 자료양이 많아 데이터베이스의 크기를 증가시키는 단점이 있다. 두 번째 방법은 스펙트럼 불연속이 발생하는 연결부분으로부터 좌우로 불연속 정도를 탐색하여 그 정도가 가장 낮은 부분을 서로 연결함으로써 불연속을 제거하는 방법이다. 이 방법은 연결 부분의 포먼트 궤적이 수평이 아니라 가정에 근거한다. 이 가정은 스펙트럼의 불연속이 두드러지게 나타나는 모습에서는 포먼트의 궤적이 수평으로 나타나기 때문에 연결 위치의 변화만으로는 불연속을 제거할 수 없는 단점이 있다. 세 번째 방법은 완만화를 수행할 때 연결부분의 불연속 정도에 상관없이 고정적으로 완만화 필터를 적용하여 완만화를 수행한다. 따라서 불연속 정도가 변함에 따라 고정적으로 수행되는 완만화 필터의 적용 결과를 신뢰하기가 어렵다. 이 외에도 기존의 방법 중에는 단위음소의 연결 부분에서 나타나는 스펙트럼의 차이를 수정하여 불연속을 제거하는 여러 시도가 있다[4,5]. 그 중에 몇몇 방법은 스펙트럼 완

만화가 오히려 합성의 질을 떨어뜨리는 결과를 보여주기도 한다. 예를 들면, 연결 부분에서 갑자기 나타나거나 혹은 사라지는 스펙트럼의 피크와 같은 스펙트럼의 왜곡이 생기는 경우이다[4]. 이러한 단점을 해결하기 위해 데이터베이스로부터 이상적인 퓨전 유닛을 추출하여 이를 스펙트럼 스무딩에 이용하는 방법이 있다[5]. 즉, 이 방법에서 사용되는 퓨전 유닛을 불연속 부분에 대체하여 합성합성함으로써 불연속을 제거하는 방법이다. 그러나 이 방법 또한 이상적인 퓨전 유닛을 추출하는 것이 어려운 문제점이 있다. 본 논문에서는 포먼트 궤적 뿐 아니라 스펙트럼의 분포특성과 인간의 청각적인 특성을 반영하여 스펙트럼을 완만화하는 방법을 제안한다.

제안하는 방법은 단위음소의 연결 구간에서 인간의 청각신경 특성을 고려하여 완만화의 양과 범위를 결정한다. 다음, 두 다이폰 경계의 스펙트럼 분포를 시간에 따라 가중치를 다르게 주어 스펙트럼 완만화를 수행한다. 이 방법은 불연속을 효과적으로 제거하며 완만화로 인하여 발생할 수 있는 음성의 왜곡을 최소화한다.

본 논문의 구성은 다음과 같다. 제2절에서는 인간의 청각신경특성을 반영한 완만화의 양과 범위 결정에 대해서 기술하고, 제3절에서는 결정된 완만화의 양과 범위를 사용하여 포먼트 궤적과 포먼트 주변의 스펙트럼 분포특성을 반영하여 완만화를 수행하는 방법에 관하여 설명한다. 제4절에서는 실험결과를 보이며, 마지막으로 제5절에서는 결론 및 향후연구에 관하여 논술한다.

2. 인간의 청각신경 특성을 반영한 완만화의 구간길이와 주파수 범위 결정

이 절에서는 인간의 청각신경 특성을 반영하여 스펙트럼 완만화를 위한 구간길이와 주파수 범위를 불연속의 정도에 적용적으로 결정하는 방법에 대하여 기술한다. 그림 1에서는 구간길이가 Δt 와 주파수 범위 Δf 의 예

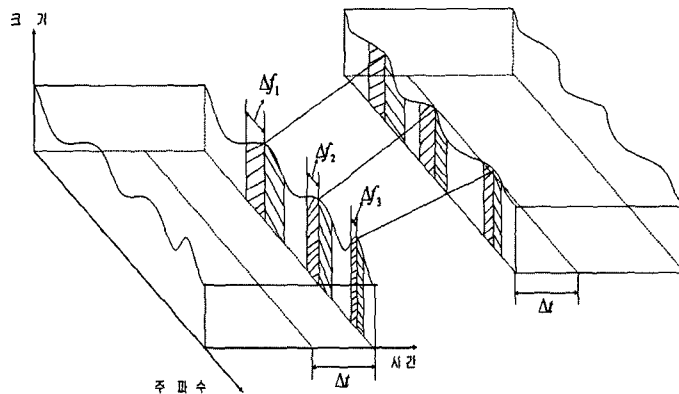


그림 1 완만화를 위한 구간길이와 주파수 범위

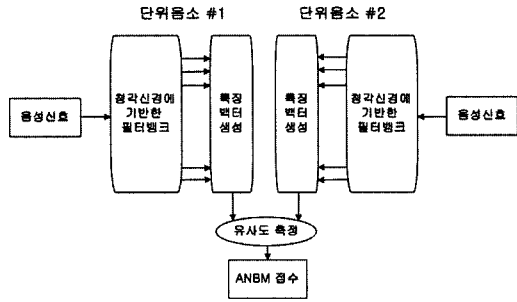


그림 2 ANBM 점수를 산출하는 과정

를 보인다.

인간의 귀에 있는 달팽이관의 해부학적 구조에 의하여 인간은 고주파의 변화보다는 저주파의 변화에 더 민감하다는 사실을 알 수가 있다[6]. 따라서 이러한 특성을 스펙트럼 완만화에 적용하기 위해, 본 논문에서는 인간의 청각신경 특성에 기반한 ANBM(Auditory Nerve Basilar Membrane) 점수[7]라는 척도를 이용하고자 한다. 즉, ANBM 점수는 인간의 귀 내에 있는 청각신경 세포(Basilar Membrane)의 분포가 모든 대역에 대하여 일정하지 않은 것을 반영하여 단위음소들의 연결부분에서 발생하는 불연속의 정도를 나타내는 척도이다 그림 2는 ANBM 점수를 산출하는 과정을 나타낸 것이다.

ANBM 점수를 산출하는 방법은 다음과 같다. 먼저, 두 개의 음성신호 각각을 청각신경에 기반한 필터뱅크를 사용하여 대역을 나눈 다음, 각 대역별로 최대강도를 가지는 두 주파수의 차를 합하여 얻어진다. 이 과정을 수식으로 표현하면 다음의 식과 같다.

$$d(x, y) = \sum_{k=1}^N |x_k - y_k| \quad (1)$$

여기서, x와 y는 두 개의 단위음소에 대한 음성신호이고, x_k 와 y_k 는 k 대역에서 최대강도를 가지는 두개의 주파수이다. 식 (1)의 의미는 주어진 음성신호에서 모든 채널에 대한 최대 크기를 가지는 주파수의 차를 모두 누적시킨 것이므로, ANBM 점수가 낮을 경우는 불연속에 대한 청각인지도가 낮은 상태, 높은 경우는 불연속에 대한 청각인지도가 높은 상태를 의미한다. 필터뱅크는 참고문헌 [6]에서 제시된 방법을 사용한다. 이 방법에 의하여 나누어진 주파수 대역은 표 1과 같다.

본 논문에서는 위에서 구한 ANBM 점수를 스펙트럼 완만화의 구간길이와 주파수 범위를 결정하기 위하여 아래와 같이 3단계로 분류한다.

- ANBM 점수 High : ANBM 값 = 4, 완만화 수행 구간 30ms
- ANBM 점수 Middle : ANBM 값 = 2, 완만화 수행 구간 20ms

표 1 청각신경에 기반한 필터뱅크

대역번호	중심주파수(Hz)	대역폭(Hz)
1	50	0-100
2	150	100-200
3	250	200-300
4	350	300-400
5	450	400-510
6	570	510-630
7	700	630-770
8	840	770-920
9	1000	920-1080
10	1170	1080-1270
11	1370	1270-1480
12	1600	1480-1720
13	1850	1720-2000
14	1250	2000-2320
15	2500	2320-2700
16	2900	2700-3150
17	3400	3150-3700

• ANBM 점수 Low : ANBM 값 = 1, 완만화 수행 구간 10ms

상기와 같이 등급화된 ANBM값을 사용하여 완만화가 수행되는 주파수의 범위는 식 (2)를 통하여 산출된다. 여기서, W는 포먼트 f_0 가 존재하는 बैं크의 대역폭이다. 즉, Δf 값은 포먼트 f_0 가 존재하는 बैं크의 대역폭과 등급화된 ANBM값에 따라 적용적으로 결정된다. 실험에서 ANBM점수를 고려하지 않고 실험할 경우에도 ANBM값은 1로 하여 완만화를 수행한다.

$$\Delta f = W \cdot ANBM \text{ 값} \quad (2)$$

$$\int_{f_0}^{f_0 + \Delta f} f(x) dx = C \quad (3)$$

식 (2)에서 결정된 Δf 는 식 (3)에서 스펙트럼 분포를 이용한 완만화에 사용되는 단위 넓이 C값을 결정하는데 사용된다.

3. 비선형적 가중치에 의한 스펙트럼 스무딩의 완만화

이 절에서는 2절에서 계산된 ANBM 점수를 이용하여 적용적으로 결정된 대역에 대하여 포먼트 궤적과 그 주위의 분포를 고려하여 스펙트럼 완만화를 수행하는 방법에 대하여 기술한다. 스펙트럼 완만화는 보간 포인트를 결정하는 단계와 보간값을 산출하는 단계로 구성된다. 먼저, 포먼트를 추출하는 방법에 관하여 설명한 다음, 보간 포인트를 결정하는 방법과 보간값을 산출하는 방법 순으로 기술한다.

포먼트 주파수는 특정한 음성이 발생될 때 성도의 공명주파수를 의미하는데, 이를 단순히 포먼트라고도 한

다. 이 공명주파수를 스펙트럼 관점에서 보면 봉우리 형태로 나타나게 되고, 저주파에서부터 나타나는 순서대로 제1포만트, 제2포만트 등으로 표기한다. 포만트의 중심 주파수를 산출하는 방법에는 퓨리에 변환이나 필터뱅크의 출력, 또는 선형예측 등을 이용하여 스펙트럼 영역에서 봉우리를 찾는 봉우리 선택(peak-picking) 방법이 사용된다[8].

보간 포인트는 봉우리 선택 방법을 사용하여 포만트를 찾은 다음 그 주변의 스펙트럼 분포를 고려하여 결정된다. 보간을 수행하기 위해서는 두 음성신호의 연결 부분에서 좌측과 우측에 있는 포만트들의 대응관계를 형성하여야 하는데, 순서가 서로 일치하도록 하여야 한다. 즉, 좌측 1번 포만트와 우측 1번 포만트가 대응하게 되며, 2, 3, 4번 포만트에도 같은 방법으로 대응관계를 형성한다. 보간을 위한 포만트들의 대응관계가 정의된 다음에는 보간을 수행한 주파수 범위를 정의하는 보간 포인트를 산출한다. 그림 3을 참조하여 이 과정을 설명하면 다음과 같다.

대응하는 두개의 포만트를 f_0 와 g_0 라고 할때, 포만트를 중심으로 상하로 스펙트럼의 주파수 강도값들을 적분한 결과가 단위 넓이 C가 되는 곳이 보간 포인트이다. 보간 포인트의 수는 포만트를 중심으로 상하 각각 m개인데, 본 논문에서는 3 즉, 7개의 보간 포인트들을 사용하였다. 여기서, 단위 넓이 C는 2절에서 언급한 것처럼 불연속정도에 따른 인간의 청각신경 특성에 따라 적용적으로 결정된다. 각 보간 포인트는 식 4와 같이 그림 3에서 빗금 친 부분의 넓이가 $A=A'$ 와 $B=B'$ 가 되도록 결정된다. 식 (4)에서 $f(x)$ 와 $g(x)$ 는 두 음성신호의 주파수 강도 분포를 나타낸다. i 는 포인트의 인덱스이다.

$$\int_{f_0}^{f_1} f(x)dx = \int_{g_0}^{g_1} g(x)dx = \left[\frac{i+1}{2} \right] C \quad (4)$$

보간 포인트들이 결정된 다음에는 대응하는 포인트들에 의하여 형성되는 보간구간 Δt 에서 스펙트럼 강도와 주파수 위치를 산출한다. 보간을 수행할 때는 원만화의 대상 구간인 Δt 을 n개의 시간 위치로 나누어 시간축을 따라서 양쪽 스펙트럼의 주파수 강도 분포에 대한 가중치를 적용하여 식 5와 식 6과 같이 새로운 주파수 강도와 위치를 산출한다. 식 5와 식 6에서는 왜곡을 최소화하기 비선형적으로 보간을 수행한다. 식 5의 결과인 M 은 비선형을 보간을 통하여 얻어진 새로운 주파수 강도이고, 식 6은 주파수의 새로운 위치이다.

$$M_n^i(k) = \frac{f_j^m \cdot NL \cdot (n-k) + g_j^m \cdot NL \cdot k}{n} \quad (0 \leq k \leq n, k \in Z) \quad (5)$$

$$F_n^i(k) = \frac{f_j \cdot NL \cdot (n-k) + g_j \cdot NL \cdot k}{n} \quad (0 \leq k \leq n, k \in Z) \quad (6)$$

식 5와 6에서 j 는 포인트의 인덱스이고 f_j^m 은 j 번째 포인트에서의 주파수의 강도이고 f_j 는 그 포인트가 위치하는 주파수이다. NL은 비선형 함수인데, 본 논문에서는 식 7과 같이 B-스플라인(spline)을 사용한다[9]. 식 7에서 l 의 범위는 구간길이에 따라 정해지며 $2l = \Delta t$ 가 되도록 정한다.

$$f(x) = \begin{cases} \frac{1}{2}|x|^3 - |x|^2 + \frac{2}{3} & 0 \leq |x| < 1 \\ -\frac{1}{6}|x|^3 + |x|^2 - 2|x| + \frac{4}{3} & 1 \leq |x| < 2l \\ 0 & 2l \leq |x| \end{cases} \quad (7)$$

그림 4와 5에서는 식 5와 6 그리고 7을 사용하여 스펙트럼 강도와 주파수 위치의 보간하는 과정을 가시화

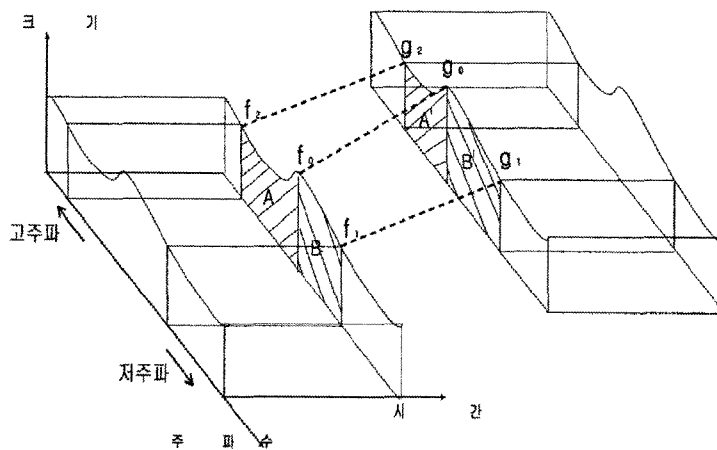


그림 3 스펙트럼 분포를 고려한 완만화

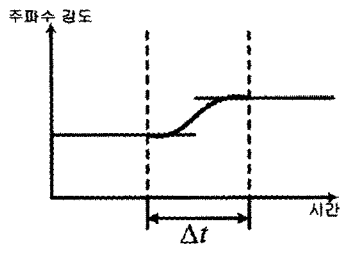


그림 4 주파수 강도의 보간

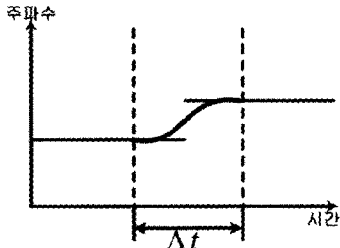


그림 5 주파수의 보간

한 것이다.

4. 실험 및 결과

본 논문에서는 실험을 위해서는 ETRI 음성 DB 샘플과 직접 녹음하여 제작한 총 20여개의 문장에서 추출한 약 500여 개의 다이폰을 사용하였다. 각 음성 샘플은 8KHz로 샘플링한 16bit 모노 음성 샘플이고 단위음소의 선택은 개인별로 추출한 다이폰을 사용하였다. 이렇게 추출된 다이폰을 사용하여 ‘ㄱ’, ‘ㄴ’, ‘ㄷ’ 세 부분에 대하여 각 5개를 합성하여 실험하였다.

그림 6에서 음성합성의 예를 스펙트로그램(spectrogram)으로 보여준다. 이때, 사용된 단위음소는 ‘목음 |’와 ‘| 귀’이다. (a)는 아무런 처리 없이 단위음소들을 단순히 연결한 경우 이고, (b)는 기존의 포맷만을 사용하여 완만화를 수행한 경우 이며, (c)는 제안된 방법으로 완만화를 수행한 경우이다. 스펙트로그램 상에서는 기존 방법과 제안된 방법 모두는 아무 처리를 하지 않은 경우 보다 자연스러운 스펙트럼의 연결을 보여주고 있다. 그러나, (b)에서 작은 두 개의 타원으로 표시된 양쪽의 각 부분의 밀도가 약간 다르게 나타남을 볼 수 있다. 이는 기존의 방법이 스펙트럼의 포맷에 초점을 맞추어 완만화를 수행하기 때문이다. (c)는 포맷 제적뿐 아니라 그 주변의 분포도 완만화 되었음을 볼 수가 있다.

성능평가는 주관적인 관점과 객관적인 관점에서 수행하였다. 주관적인 관점에서의 성능평가를 위한 척도 MOS(Mean Opinion Score)을 사용하였고, 5명을 대상으로 합성한 문장들을 들려주어 합성된 음성이 듣기에 자연스러운 정도를 평가하였다[4]. 객관적인 관점에서의

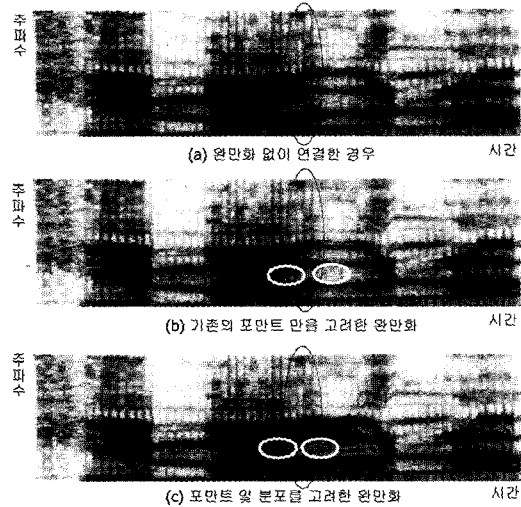


그림 6 기존 완만화 방법과 제안된 방법의 스펙트로그램 비교

성능평가는 단위음소의 연결 부분에서 식 8의 KL-거리를 사용하여 불연속 정도를 측정함에 의하여 이루어진다[10].

$$KL(f, g) = \int f(x) \log \left(\frac{f(x)}{g(x)} \right) dx \quad (8)$$

식 (8)에서 f(x)와 g(x)는 각각 두 단위음소의 경계에 있는 스펙트럼의 주파수 강도의 확률분포함수이다. 이 두 함수의 관계 f(x)=g(x) 이면 두 경계의 유사도가 매우 높음을 의미한다. 즉, 실제적인 KL-거리는 연결 부분에서 주파수 강도의 밴드별 평균과 표준편차 사이의 차이를 누적함에 의하여 산출된다.

그림 7에서는 기존의 스펙트럼 완만화 방법과 제안된 방법의 객관적인 관점에서의 성능의 비교를 보여준다. 이때 사용된 음소는 다이폰 기반의 음성합성에서 왜곡이 일어나기 가장 쉬운 음소인 ‘ㄱ’, ‘ㄴ’, ‘ㄷ’ 이다[11]. 그림 7에 의하면 본 논문에서 제안된 방법이 보다 더 우수함을 알 수 있다.

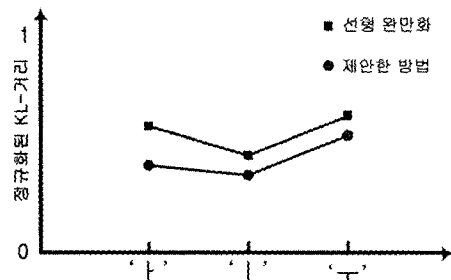


그림 7 ‘ㄱ’ ‘ㄴ’ ‘ㄷ’의 음성합성 결과에 대한 KL-거리 비교

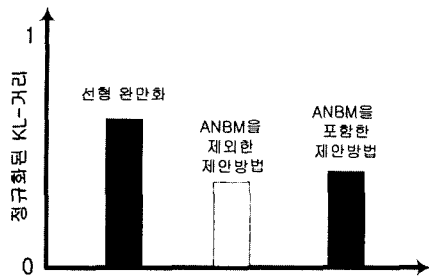


그림 8 3가지 방법에 대한 KL-거리 비교

그림 8은 같은 척도를 사용하여 기존의 방법과 비교하는 동시에 제안된 방법에서 ANBM을 적용한 방법과 적용되지 않은 방법을 비교하였다. 여기서 주목할 사실은 제안된 방법에 ANBM 점수를 적용한 경우와 적용하지 않은 경우를 비교했을 때, 오히려 ANBM 점수를 적용하지 않은 방법이 불연속을 더 줄이는 것으로 나타난다.

KL-거리가 적다할지라고 음성합성을 통하여 산출된 결과 음성신호를 청취할 때 왜곡으로 인해 자연스럽지 못한 경우가 많다. 따라서, 합성된 결과 음성을 사람이 실제 청취한 후 주관적으로 성능을 평가하는 방법이 요구된다. 이를 위하여 본 논문에서는 5명을 대상으로 MOS를 측정하였다. 그 결과는 표 2와 같다.

표 2 MOS 테스트

알고리즘	MOS
자연 음성	4.54
가공하지 않은 연결	3.21
기존의 선형적인 완만화	3.36
ANBM을 사용하지 않은 제안방법	3.61
ANBM을 사용한 제안방법	3.82

그림 8과 표 2에 의하면, ANBM 점수를 적용하지 않은 방법이 객관적인 불연속 정도를 줄일 수는 있었으나 실제 사람이 듣는 합성음성의 질에서는 ANBM 점수를 적용한 방법이 더 나은 결과를 보여줌을 알 수 있다.

5. 결론 및 향후 연구 과제

본 논문에서는 음성합성을 할 때 연결부위에서 생겨나는 스펙트럼의 불연속을 효과적으로 제거하는 방법을 제안하고 실험을 통한 성능평가를 수행하였다. 제안하는 방법은 단위음소의 연결 구간에서 인간의 청각신경 특성을 고려하여 완만화의 양과 범위를 결정한 다음, 두 다이폰 경계의 스펙트럼 분포를 시간에 따라 가중치를 다르게 주어 스펙트럼 완만화를 수행한다. 성능평가는 객관적인 관점과 주관적인 관점에서 수행하였다. 그 결

과 본 논문에서 제안한 방법이 기존의 주요 방법 보다 우수함을 알 수 있었다.

본 논문에서 제안한 완만화 방법은 음성 데이터베이스의 용량이 클수록 더욱 자연스러운 음성을 생성할 수 있는 한계를 가지고 있다. 따라서, 향후 연구로는 이러한 한계를 극복하고 소용량의 음성 데이터베이스를 사용하여 자연스러운 음성 합성을 수행하는 방법에 대한 연구라 고려된다.

참고 문헌

- [1] R.E. Donovan, P.C. Woodland, A hidden Markov model based trainable speech synthesizer, Computer Speech and Language, pp1-19, 1999.
- [2] Conkie, A.D., Isard S., Optimal coupling of diphones Progress in Speech Synthesis, Springer, New York, Chapter 23, pp293-304, 1997.
- [3] Kleijn W.B., Haagen J., Waveform interpolation for coding and synthesis, Speech Coding and Synthesis, Chapter 5, pp175-207, 1995.
- [4] David T. Chappell, John H.L. Hansen, A Comparison of Spectral Smoothing methods for segment concatenation based speech synthesis, Speech Communication 36, pp343-374, 2002.
- [5] Wouters, J., Macon, M.W., Control of Spectral Dynamics in Concatenative Speech Synthesis, Speech and Audio Processing, IEEE Transactions on, Vol 9, No. 1, pp30-38, Jan 2001.
- [6] Hossein Najafzadeh-Azghandi, Perceptual Coding of Narrowband Signals, Ph.D Thesis, Department of Electrical & Computer Engineering, McGill University, Montreal, Canada, April 2000.
- [7] John H. L. Hansen and David T.Chappell, An Auditory-Based Distortion Measure with Application to Concatenative Speech Synthesis, Speech and Audio Processing, IEEE Transactions on, Vol 6, No. 5, pp489-495, Sep 1998.
- [8] L. R. Rabiner, R. W. Schafer, Digital Processing of Speech Signals, Prentice-hall, 1978.
- [9] H. S. Hou and H. C. Andrews, Cubic Splines for Image Interpolation and Digital Filtering, IEEE Trans. Acoustics, Speech, and Signal Processing, ASSP-26,6, December 1978, 508-517.
- [10] Esther Klabbbers, Raymond Veldhuis, Reducing Audible Spectral Discontinuities, IEEE Transactions on Speech and Audio Processing, Vol 9, No. 1, Jan 2001.
- [11] H. van den Heuvel, B.Cranen, T.Rietveld, Speaker variability in the coarticulation of /a,i,u/, Speech Communication 18, pp113-130, 1996.



장 효 중

2001년 2월 숭실대학교 컴퓨터학부 졸업 (공학사). 2003년 2월 숭실대학교 대학원 컴퓨터학과 졸업(공학석사). 2003년 3월~현재 숭실대학교 대학원 컴퓨터학과 박사과정. 관심분야는 컴퓨터비전, 음성 처리, 영상처리, 패턴인식, 3D모델링 등



최 형 일

1979년 2월 연세대학교 전자공학과 졸업 (공학사). 1982년 6월 미시간대학교 전산 공학과 졸업(공학석사). 1987년 6월 미시간대학교 전산공학과 졸업(공학박사) 1987년 9월~현재 숭실대학교 미디어학 부 교수. 관심분야는 컴퓨터비전, 패턴인식, 퍼지이론, 비디오검색, 인터페이스 에이전트 등.