

음성망을 이용한 한국어 연속 숫자음 인식에 관한 연구

* 이강성 이형준 변용규 김순협
광운 대학교 전자 계산기 공학과

Study on the Recognition of Spoken Korean Continuous Digits Using Phone Network

* G. S. Lee H. J. Lee Y. G. Byun S. H. Kim
Dept. of Computer Engineering, Kwang Woon Univ.

Abstract

This paper describes the implementation of recognition of speaker - dependent Korean spoken continuous digits.

The recognition system can be divided into two parts, acoustic - phonetic processor and lexical decoder. Acoustic - phonetic processor calculates the feature vectors from input speech signal and the performs frame labelling and phone labelling. Frame labelling is performed by Bayesian classification method and phone labelling is performed using labelled frame and posteriori probability. The lexical decoder accepts segments (phones) from acoustic - phonetic processor and decodes its lexical structure through phone network which is constructed from phonetic representation of ten digits.

The experiment carried out with two sets of 4continuous digits, each set is composed of 35 patterns. An evaluation of the system yielded a pattern accuracy of about 80 percent resulting from a word accuracy of about 95 percent.

제1장 서론

기계가 인간과 같은 귀를 가지는 것을 목표로 1950년대 초부터 시작된 음성 자동화의 연구는 지난 십 수년간 격리단어인식 및 연결단어 인식 분야에서 상당한 발전을 이룩해 많은 시제품까지 나오고있으며 연속음성 이해 분야에 있어서도 CMU의 Harpy 시스템을 비롯, 일본 독일 프랑스 등 각국에서 많은 연구가 진행중에 있다. 그러나 국내의 경우 연속음 인식을 위한 연구는 아직 초보적인 단계에 있으며 이를 위한 많은 연구가 시급하다고 하겠다. 연속음 인식의 적용분야중에서 숫자음은 가장 보편적인 대상어휘이므로 반드시 해결해야 할 인식 대상이다. 따라서 본 실험은 연속 숫자음을 인식 대상어휘로 선정하였다. 우선 연속음 인식 대상이 결정 되었으면 인식의 기본단위 (recognition unit)를 설정해야 하는데, 지금까지 영어를 대상으로 하는 음성인식 시스템에 대한 연구 결과에 의하면 단어를 인식의 기본 단위로 하는 것은 바람직하지 않다. [2] 숫자음은 음소의 수가 적고 어휘수가 적어 문맥에 따른 음질의 변화가 비교적 적다는 점에서 음소를 표준음성 추출의 기본단위로 설정했다. 하지만 음소나 변이음은 음의 변별적 특징 (distinctive feature)을 구별하기 위한 단위이므로 심재적으로 변화하는 음성현상과 일대일로 대응하지는 못한다. 따라서 음소사이의 과도음이나 각 음소의 변화음등의 음향현상과 대응되는 최소단위를 음성(Phone)이라 하고 본 논문에서는 이를 인식의 최소단위로 한다 [1] [3]. 즉, 다시말하면 음소단위로 추출한 표준음성을 이용하여 입력음성을 음성 (phone) 단위로 분리하고, 이를 연속 숫자음 인식의 기본 단위로 삼는 것이다.

본 논문에서의 전체적인 인식절차를 그림 1 에 보인다.

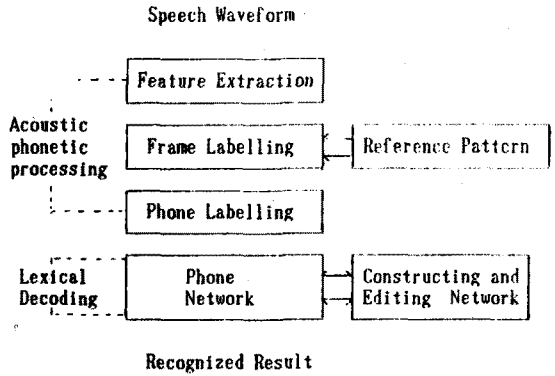


그림 1 본 논문의 구성도.

본 논문에서는 베이스 분류법을 이용하여 각 프레임에 레이블링 한것을 음성단위로 레이블링하는 방법을 제시한다. 또한 음성 레이블링에는 오류정보가 포함되기 마련인데 예기치 못한 음성의 오류는 치명적인 어휘해석의 오류를 낳게 하므로 이들 정보를 고려하는 것은 매우 중요하다. 이 문제의 해결을 위하여 음성을 노드도 하는 음성망(phone network)을 이용하여 어휘해석을 한다.

본 실험은 음성 다이얼링 시스템에 적용하도록 4연속 숫자음을 실험대상 자료로 선정하였는데 그 이유는 구내교환기의 교환번호 및 일반 전화의 번호가 국번을 따로 분리해 생각한다면 4 자리 이하라는 점에 있다. 하지만 본실험이 4연속 숫자음에만 국한된 알고리즘을 제안하는 것은 아니며, 몇자리 연속 숫자음이라도 인식할 수 있는 알고리즘을 제안 했으므로 다른 응용에도 그대로 적용 될 수 있을 것이다. 한국 전자통신연구소 (ETRI)의 음향연구실에서 작성한 자료를 평가자료로 이용한다.

제2장 표준음성의 작성

제1절 인식의 기본단위

숫자음인 경우에 고려해야할 대상음성 수가 대어휘인 경우에 비해 대폭 줄어든다. 본 실험에서 설정한 음성의 종류를 표 1에서 보인다.

설정된 표준 음성			
번호	분류	음성	예
1	모음 (vowels)	[i]	일, 이, 칠, 팔
2		[a]	삼, 사, 오
3		[o]	
4		[u]	구, 구유
5		[yu]	
6	자음 (consonants)	[m]	삼, 공, 육, 율, 팔
7		[n]	
8		[l]	
9		[s]	일, 삼, 사
10		[c ^h]	칠, 팔
11		[p ^h]	
12	[k]	구, 공	
13	목음 (silence)	[q]	

표 1 표준음성의 종류.

다음이 각 표준음성의 추출 기준이다.

- (1) 단모음 : 모음 정상부의 중심
- (2) 중모음 : spectrum의 변화시점
- (3) 묵음(silent) : 부음구간 중심
- (4) 비음 : 비음 정상부의 중심
- (5) 그외 ([s], [c^h], [p^h], [k]) : 안정부 중심

제 3 절 표준 패턴

각 음성 표준 패턴은 평균 벡터, 공분산 행렬의 역행렬과 행렬식 값이다.

제 4 절 특징 추출

LPC cepstrum 계수를 특징 벡터로 한다. 다음에 특징 벡터를 구하는 절차를 보인다.

- (1) Filter (4.5 LPF)
- (2) AD conversion (10 KHz)
- (3) pre-emphasis (1-Z⁻¹)
- (4) Hamming windowing (123 ms window for every 6.4 msec step)
- (5) LPC cepstrum analysis (order 12)

제 3 장 레이블링

제 1 절 프레임 레이블링

각 프레임의 레이블링은 베이스 식별법 (Bayesian classification method) [4] [5] 에 근거한다.

$$P(i|x) = \frac{P(x|i) P(i)}{\sum_{k=1}^N P(x|k) P(k)} \quad (3.1)$$

단, P(x|i) 는 확률함수, P(i) 는 사전(priori) 발생확률, N 는 표준 음성의 개수 x 는 입력 벡터이다. 만약 확률 분포가 Gaussian 분포를 따른다면 위의 P(x|i)를 다음식의 multivariate Gaussian 분포로서 근사할 수 있다.

$$P(x|i) = \frac{1}{(2\pi)^n |V_i|^2} \exp\left\{-\frac{1}{2}(x-\mu_i)^T V_i^{-1}(x-\mu_i)\right\} \quad (3.2)$$

이때 n은 차원 수, V는 공분산 행렬이며 μ는 평균 벡터이다. 또한 식 (3.1) 에서 p(i)는 음성 특징의 사전 확률 (prior probability) 을 나타내는데 모든 음성 특징에

대하여 같다고 가정하면 분모항은 모두 같은 값을 갖는다. 따라서 p(i|x)는 단지 p(x|i)에만 의존하게 된다. 그림 2에 베이스 분류법으로 분류한 예를 보인다. 그림의 한 개의 점은 확률 1을 17개의 gray level로 나타낸 것이며 진할수록 높은 확률을 나타낸다. 또 스펙트로그램 와도 함께 비교한다.

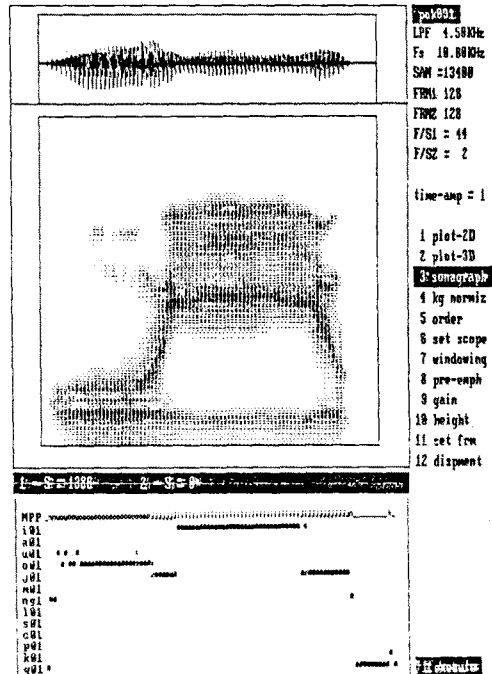


그림 2 스펙트로그램 와 가우스 분포곡선. pck031.dat (5267) 의 a) 파형 b) 스펙트로그램 c) 베이스 분류법

제 2 절 음성 레이블링 (phone labelling)

추출된 각 프레임마다 음성에 대한 확률 분포는 종종 오류 정보를 포함하고 (하거나) 불안정한 형태를 취하게 되는데 이를 일차적으로 구분하여 음성적 성질을 갖는 영역으로 나눔 필요가 있다.

2.1 대략적인 구분화 작업 (Coarse Segmentation)

이 단계에서는 프레임 레이블을 이용하여 대략적인 변이음단위의 구분화 작업을 한다. 어느 점에서 시작하여 임계값(Thn) 이상의 연속 프레임이 하나의 레이블을 갖는 구간으로 자르는 것이다. 설정된 구간내에 두개 이상의 레이블이 존재 할 때 이 구간을 안정구간이라고 하고 한개만의 레이블이 존재 할 때 안정 구간이라고 한다.

흐름도를 그림 3에 보인다.

2.2 세분 절차 (Refine procedure)

앞단에서 얻어진 세그먼트 중에서 불안정한 것 만을 취하여 그것이 몇개의 음성구간으로 다시 구분될 수 있는지의 여부를 타인한다. 불안정 구간에 존재 할 수 있는 음성구간의 수는 대체로 파라미터 Thn에 의존하는데 Thn이 클 수록 다수개의 구간이 존재 할 수 있다. 본 실험에서 설정한 것과같이 Thn이 3 일경우 최대 2개의 구간이 존재한다고 보면 충분하다. 이 구간 분할의 문제는 Brandt's GLR method[6] 의 기본 개념을 적용함으로써 해결한다.

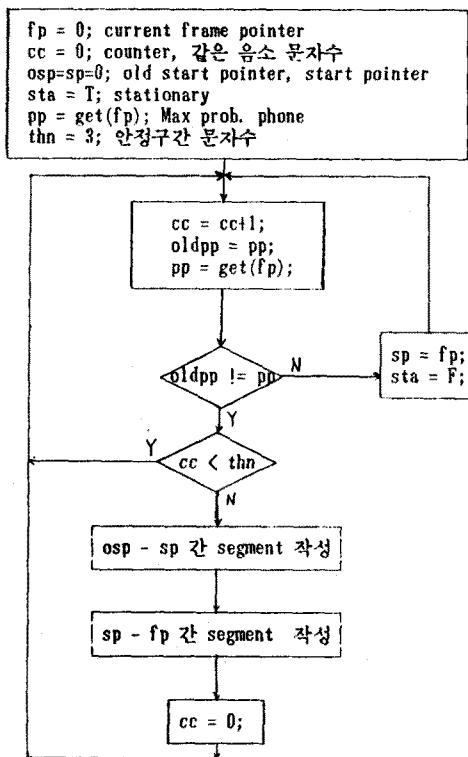


그림 3 대략적인 구분화 알고리즘.

2.3 정정 (Correction)

이 단계에서는 설정된 음성이 타당한지의 여부를 결정한다. 이 단계에서는 종종 불안정한 세그먼트(segment, phone)를 입력으로 받게되는데, 이를 보정할 필요가 있다. 먼저 [p], [s], [t], [k] 등의 자음은 음부터 끝까지 같은 음성적 성질을 갖지 않으므로 하나의 자음이 여러개의 자음 세그먼트로 분류되는데 그 패턴은 대략 일정하다. 예를 들어 /s/는 [c]+[s], [s]+[c]+[s], [c]+[s]+[p] 등의 패턴을 갖으며, /p/는 [s]+[p], /c/는 [n]+[c], [s]+[c], /k/는 [n], [p]의 패턴을 갖는다. 두성 자음구간의 세그먼트들을 위의 패턴과 시간적인 정보를 보조적으로 이용하면서 비교, 대치시킨다.

두 번째로 [yu]의 세그먼트에 대해서 살펴본다. 2 장에서 살펴본바와 같이 [yu]는 [i]에서 [u]로의 과도상태를 추출하여 만든 음성이지만 시간적인 변화 파라미터는 포함하고 있지 않고있다. 따라서 다음의 경우에 [yu]가 나타날 수 있다.

1. 음성 [yu]의 경우
2. 한음성에서 다른 음성으로가는 과도상태
3. 불명료한 발음으로 인해 안정 상태인 구간이 잘못 분류되는 경우

위의 세가지 경우를 분리하는 방법으로는 분석적인 방법과 구분적인 방법을 생각할 수 있다. 분석적인 방법은 구간에 스펙트럼 변화와 [yu]의 시간변화를 고려한 표준패턴과의 매칭거리를 파라미터로 이용한다. 구분적인 방법으로는 [yu] 세그먼트의 주변환경을 고찰함으로써 [yu] 세그먼트의 상태를 결정하는 방법이다. 전자의 방법을 [yu] 세그먼트 성격규정을 위해 사용해 보았는데 몇가지 부적당한 이유 - 불안정한 발음이나 다른 과도음 구간의 [yu] 세그먼트와 실제의 [yu] 세그먼트가 결합된 상태 또는 [yu]의 정확한 구분이 이루어지지 않은 상태에서는 위의 두가지 조건을 만족하는 절대 임계치를 설정하기 어렵다. - 로 적용할 수 없었고, 단지 구분적인 방법만을 적용하였다. 일차로 임계 프레임 갯수가 되는 세그먼트는 앞 세그먼트와 같은 세그먼트로 치환하고 [yu] 세그먼트 다음에

음성 [q], [ŋ], [i], [o]가 나오지 않는 세그먼트 [yu]는 구분에 맞게 [yu]를 대치한다. 그리고나서 음성망을 통과하는 어휘분석을 통한 인식 단계에서 [yu] 세그먼트가 '육'으로 치환될때 최소음길 길이와 비교해서 만일 '육' 음길 길이가 이것보다 작으면 실험적으로 표를 만들어 세그먼트의 성격을 판별한다. 이것이 과도 세그먼트로 판별될 경우 앞 세그먼트에 포함시킨다.

세번째로 모음의 연결(juncture)에 관해 고찰한다. 숫자음에서 문제가 되는 연결은 [i], [o]의 동일 모음 폐쇄연접(close juncture)이 일어날 경우이다. 같은 [i]이면서 '1'의 [i]와 '2'의 [i]는 조금 다르다. '2'의 [i]는 이완모음(lax vowel), '1'의 [i]는 긴장모음(tense vowel)에 속한다고 할 수 있다. 그래서 뒤 세그먼트가 [i]인가를 검사해서 [i]이라면 [i]의 긴장모음에 해당하는 부분을 파위의 dip 정보를 이용하여 추출한다.

자연스럽게 발음했을 경우 같은 이완모음 [i], [o]끼리의 연결은 에너지 상태 변화로는 구분 할 수가 없으므로 여기서는 평균 길이를 이용하여 구분한다. 이완모음 [i], [o]의 평균길이를 lmi, lmo라고 하고 세그먼트 [i], [o]의 길이를 li, lo라하면, 모음 [i], [o]의 수 ni, no는 다음과 같이 결정된다.

$$ni = li / lmi$$

$$no = lo / lmo$$

제 4 장 어휘 해석

제 1 절 음성망 (phone network)의 작성

음성신호를 적당한 단어 혹은 음절의 열로 변환하기 위하여 앞단에서 얻어진 단위의 이동 경로를 지정해주는 음성망을 작성한다. 음성망의 작성은 기본적으로 새내계를 거쳐 작성된다. 첫 번째로는 열개 숫자음의 변이음을 노드로 개별적인 브랜치를 만든다. 두 번째로는 각 단어의 같은 기호로 시작하는 노드를 묶어 트리를 형성한다. 마지막으로는 각 단어의 제일 마지막 노드로부터 시작해서 10개의 각 단어 시작 노드로 연결하는데 음운학적 규칙을 적용한다(그림 4).

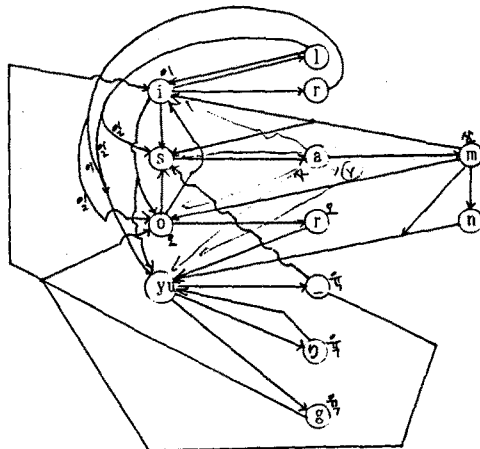


그림 4 음성망 : 음운규칙을 적용하여 단어의 끝 노드에서 모든 시작 노드와의 연결한 예.

오류정보가 포함된 세그먼트 계열을 인식하기 위해서는 모든 정보를 고려한 적당한 전이상태를 정의하지 않으면 안된다. 본실험에서 수정하여 얻은 음성망을 그림5에 보인다.

제 5 장 실험결과 및 고찰

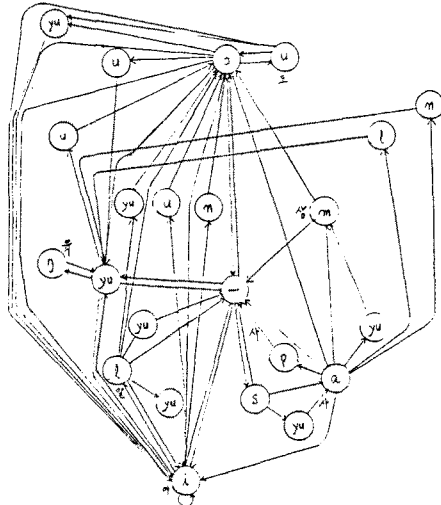


그림 5 수정된 음성망.

본 실험에서는 한국 전자 통신 연구소의 음향연구실에서 88년 1월에 작성한 4연속 숫자음 음성 대상자료로 선정하였다. 이 자료는 가능한 모든 조합을 고려해 만든 것이며 표 2 에 그 내용을 보인다. 녹음은 방음실에서 됐으며 발성 속도는 약 4.26 음절/초이다. 필터는 4.5 LPF 이며 12 Bit 10KHz 샘플링된 데이터이다.

이들중 본 실험에서는 화자 PCK의 두번 발성한 자료를 이용하였다. 첫번째 발성은 음성 표준패턴 추출용으로 사용하였으며 나머지 한개 발성을 인식 실험용으로 이용했다. 레이블링은 확대된 시간축상의 파형과 3차원 스펙트럼, 스펙트로그램, 파위를 이용했다.

실험 결과는 단어집을 기준으로 했을 경우 95%였고, 4연속 숫자음을 기준으로 했을 경우는 80%가 나왔다. 표 3 에 오인식 실험 결과를 보인다.

음성 레이블링을 거쳐서 나온 [yu] 세그먼트의 성격 규정은 구분적 방법만을 도입했는데 이 구분적 방법 뿐만 아니라 분석적 방법으로 해석이 행해져야 보다 많은 어휘 인식 시스템에 적용될 수 있을 것이다. 또한 구분화가 정확하게 이루어지지 않을 경우가 있었다는 것이 위의 분석적 방식으로 [yu] 세그먼트의 성격 규정을 하는 어려운 점이기도 했다. 음성 표준 패턴을 잡는데 있어서 추출된 벡터의 수나 추출한 부분에 따라서 베이스 분류법에 의한 사후 확률 분포가 달라질 수 있었는데 이 음성 표준 패턴을 작성하는 음성 종류의 정확한 기준이 필요할 것으로 생각된다. 프레임 레이블링의 오류나 음성 레이블링으로의 구분화 오류는 이 영향도 상당히 클 수 있다. 또한 동일모음의 연결이 일어날 경우 이의 분리에 대한 알고리즘이 미흡하다.

여기서 고려되지 않은 음성도 역시 같은 방식으로 추출 가능하며 앞으로 대어휘 연속음 인식에도 우리가 없이 적용될 수 있을 것으로 기대된다. 또한 본 실험이 앞으로의 연속음 인식분야에 유용하게 이용될 수 있으리라고 믿는다.

no.	contents	no.	contents	no.	contents
1	0287	13	1823	25	5500
2	5732	14	6378	26	6972
3	9601	15	8877	27	9861
4	4156	16	3510	28	3649
5	1199	17	3065	29	0316
6	1398	18	2934	30	7033
7	6843	19	7489	31	8194
8	0712	20	2244	32	9205
9	5267	21	4621	33	1427
10	6633	22	9176	34	2538
11	2409	23	3045	35	4750
12	7954	24	8590		

표 2 4연속 숫자 발성표.

오인식 결과	수정사항
2400 (2409)	0 → 9
8064 (8065)	0 → 5
x489 (7489)	x → 7
4x41 (4621)	x → 6, 4 → 2
3040 (3045)	0 → 5
9361 (5861)	9 → 5
4799 (4750)	99 → 50

4연속음별 인식률 : 80% 음절별 인식률 : 95%

표 3 오인식 결과.

제 6 장 결론

본 논문은 모든조합을 고려한 35종의 4연속숫자음을 대상으로 인식실험을 하였다. 음향-음성처리로서 연속음에서 추출하여 만든 음성표준 패턴을 이용, 베이스 분류법 (Bayesian classification method)을 통해 사후 확률을 구하고 프레임 레이블링(Frame Labelling)과 음성 레이블링(Phone Labelling)을 하였으며, 어휘 해석으로는 음성망을 이용해 연속음 인식을 하였다.

음성 표준패턴은 한명이 처음 발음한 연속음에서 추출하였으며 나머지 1 회 의 자료는 실험 평가용으로 사용하였다.

4연속숫자음 단위로는 80%의 인식율을 얻었으며, 음절 단위로는 95%의 인식율을 얻었다.

세그먼트의 오분류는 인식과정에서 필수적으로 따르는 것인데 오인식이 됐거나 구분적으로 오류를 정정해야 했던 이들 세그먼트의 발생은 1차적으로는 프레임 레이블링 단계에서의 오류 정보의 영향이 가장 컸으며 2차적으로는 음성 테이블 단계에서였다. 또한 세그먼트 [yu] 의 성격 규정에도 어려움이 많았는데 이를 위한 효과적인 분석적 방식이 요구된다.

본 실험은 음성 다이얼링 시스템에의 응용을 목적으로 한 것이나 Home baking 시스템 많은 숫자음 인식을 이용하는 시스템에 적용될 수 있으며 대어휘 연속음 인식 연구에도 유용하게 이용될 수 있을 것이다.

본 실험을 위해 자료와 도움을 주신 한국 전자통신 연구소 음향연구실 김경태 실장님 및 실험 여러분께 감사록 드립니다.

참고 문헌

- Jean-Paul Halton, Automatic Speech Analysis and Recognition, Reidel Publishing Company, 1982.
- 은 종관, "음성인식 기술현황," 한국 음향학회지 Vol. 7 No. 1, 1988.
- Ronald A. Cole, Perception and Production of Fluent Speech, Lawrence Erlbaum Associates, Inc., Publishers, 1980.
- 小坂 哲夫, "連續音認識のための音素及び単語の検出に 関する研究," 東北大学 修士学位論文, 1986
- Sei-ichi Nakagawa, "Speaker - Independent Phoneme Recognition in Continuous Speech by a Statistical Method and a Stochastic Dynamic Time Warping Method," Technical Report of Carnegie - Mellon University, CMU-CS-86-102, 1988.
- Regine Andre-Obrecht, "A New Statistical Approach for Automatic Segmentation of Continuous Speech Signals," IEEE, Trans. on Acoustics, Speech, Signal Processing, Vol. 36, No. 1, Jan. 1988.