

## 필기체 혼합 문서 인식에 관한 연구

正會員 沈 東 圭\* 正會員 金 仁 權\* 正會員 咸 永 國\* 正會員 朴 來 弘\*  
 正會員 李 昌 範\*\* 正會員 金 庠 仲\*\* 正會員 尹 炳 楠\*\*

## A Study on the Recognition of Handwritten Mixed Documents

Dong-Gyu Sim\*, In Kwon Kim\*, Young Kug Ham\*, Rae-Hong Park\*,  
 Chang Bum Lee\*\*, Sang Joong Kim\*\*, Byeong Nam Yoon\*\* *Regular Members*

### 要 約

본 논문에서는 그래픽을 포함한 필기체 한글과 영숫자로 구성된 혼합문서 인식시스템을 제안하였다. 전처리 과정에서 제안한 국부적응 이진화 알고리즘으로 이진화를 수행하며, 연결요소와 체인코드를 이용하여 그래픽 영역을 분리하고 한글의 문자유형, 크기 그리고 수직모음의 부분적인 인식을 이용하여 개별문자를 분리한다. 인식단계에서는 DP 정합 비용함수값에 따른 branch and bound 알고리즘을 이용하여 한글 문자를 인식하며, 또한 몇개의 안정한 특징값을 이용하여 영숫자를 인식하였다. 또한 인식단계에서의 정보와 단어사전의 정합을 통하여 인식기의 오류를 정정하였다. 컴퓨터 모의실험을 통하여 제안한 시스템이 그래픽을 포함한 필기체 한글과 영숫자를 효과적으로 인식함을 보였다.

### ABSTRACT

This paper proposes an effective recognition system which recognizes the mixed document consisting of handwritten Korean/alphanumeric texts and graphic images. In the preprocessing step, an input image is binarized by the proposed thresholding scheme, then graphic and character regions are separated by using connected components and chain codes. Separated Korean characters are merged based on partial recognition and their character types and sizes. In the character recognition step, we use the branch and bound algorithm based on DP matching costs to recognize Korean characters. Also we recognize alphanumeric characters using several robust features. Finally we use a dictionary and information of a recognition step to correct wrong recognition results. Computer simulation with several test documents shows that the proposed algorithm recognizes effectively handwritten mixed texts.

\*西江大學校 電子工學科  
 Dept. of Electronic Eng., Sogang Univ.

\*\*韓國電子通信研究所  
 Electronics and Telecommunications Research Institute

論文番號 : 9430  
 接受日字 : 1994年 1月 27日

## I. 서론

컴퓨터의 정보처리 능력은 정보화 사회에 있어 중요한 역할을 하고 있다. 이것은 컴퓨터의 빠른 연산 처리 능력과 거대한 정보의 저장, 검색을 효과적으로 처리할 수 있는 능력 때문이다. 이러한 장점으로 인하여 앞으로도 컴퓨터의 사용은 급속히 증가할 것이며 그 응용분야 또한 다양해질 것이다. 그러나 컴퓨터의 정보 표현 방식이 인간의 정보 표현 방법과 달리 정보의 교환에 있어서는 많은 단점을 가지고 있다. 이 때문에 컴퓨터에 입력할 자료의 자동화에 대한 요구가 급증하고 있으며, 자료 입력에 해당하는 시각능력과 청각능력을 컴퓨터로 구현하고자 하는 노력이 꾸준히 진행되어 왔다. 이 중 시각능력의 관점에서 문자인식, 악보인식, 도형이나 기하학적인 물체의 인식, 설계도인식 등이 컴퓨터와 인간의 인터페이스 기술로 발전하였다. 이 중에서 문서는 기존에 많이 사용되어 온 정보 표현수단으로, 문서와 컴퓨터 사이의 효과적인 정보 교환이 요구된다. 특히 기존의 많은 필기체 문서의 처리를 위한 입력 시스템이 필요할 것으로 생각된다.

이러한 문서의 자동 입력과 인식을 위해서는 무엇보다도 먼저 기본적으로 문자인식이 이루어져야 하는데 문자인식은 패턴인식 응용의 한 분야로 오래전부터 연구되어 왔다. 지금까지 한글과 영문자 인식에 대해서는 많은 연구가 이루어져왔고, 규격화된 문서나 제한적인 인쇄체 한글, 영문 및 그래픽 등으로 이루어진 혼합문서에 대한 연구가 이루어지기 시작했으나 필기체를 포함한 한글, 영문 및 그래픽 등으로 구성된 혼합문서 인식에 대한 연구는 미미한 상태이다.<sup>(1)</sup>

한글문서는 일반적으로 도표와 함께 한글, 영숫자 그리고 기호 등이 포함된 것으로 정의할 수 있다. 이것을 인식하기 위해서는 문자와 그림영역을 분리하고, 분리된 문자영역으로부터 각 문자를 추출하여 인식하여야 한다. 본 연구에서는 필기체 한글문서중 제한된 필기체 한글, 영숫자 그리고 그림으로 구성된 혼합문서의 인식 시스템을 개발하였다. 혼합문서가 입력되면 먼저 전처리 단계에서 국부적 이진화 과정을 수행하고 연결화소를 이용하여 문자부분과 그래픽부분을 분리하였다. 다음 연결화소의 유형과 수직모음의 인식을 통해 문자부분에서 개별문자를 추출하여 인식단으로 넘겨준다. 인식단에서는 개별문자

에 대해 연결화소의 개수와 끝점의 수를 이용하여, 한글과 영문의 구분이 가능한 경우에는 분리하고, 분리가 모호한 경우에는 투영값과 흑백화소 변화수, 거리 등의 특징에 의하여 영문이나 숫자로 분리되면 영숫자 인식단에서 인식하며, 그렇지 않은 경우 한글 문자인식단에서 인식하였다. 한글 문자인식은 DP (Dynamic Programming) 정합에 의한 비용함수값에 바탕을 두어 branch and bound 알고리즘에 의하여 입력패턴의 획을 탐색하여 기본자소를 추출하였다. 다음 기본자소로부터 발생할 수 있는 나머지 자소를 구조해석법을 이용하여 인식하고, 자소간의 접촉 가능성을 고려함으로써 접촉으로 인식이 어려운 문자를 옳게 인식하였다. 또한 한 문자 자체만으로 인식이 불가능한 경우에도 단어사전과의 정합으로 인식을 하였고, 이로 인해 인식률이 향상되었다.

본 논문에서는 II 장에서 문서의 전처리 단계에 대해 서술하고, III 장에서는 제안된 필기체 혼합문서 인식과 후처리과정에 대하여 서술하였으며, IV 장에서는 실험결과 및 분석에 대해 서술하였다. 마지막 V 장에서 결론을 맺었다.

## II. 전처리

제한한 시스템은 크게 전처리 단계와 문자인식 단계로 나누어진다. 전처리 단계에서는 입력영상에서 개별문자를 추출하는 단계까지를 말한다. 전처리 과정에서는 256 그레이 레벨 입력문서를 이진화한 후 비문자부분을 분리하고 개별문자를 추출한다.

### 1. 이진화

문서인식 시스템에서 영상 이진화는 전처리에서 매우 중요한 과정이다. 그러나 일반적인 문자인식 시스템들은 스캐너 등의 입력장치가 제공하는 이진영상을 사용하기 때문에 응용예에 따라서 이진영상이 불만족스러운 경우가 있다. 인식 시스템의 경우 입력장치에 의해 입력되는 영상의 상태에 따라 전체 인식 시스템의 성능이 달라진다. 즉 입력상태에 따른 문자들의 검침, 끊어짐 그리고 파손 등으로 인하여 문서 영상의 영역화나 인식과정이 영향을 받는다. 따라서 이진화 작업은 그레이 레벨 영상으로 입력받은 문서 영상에서 문자에 대한 정보를 정확하게 추출할 수 있어야 한다.

본 논문에서는 밝기값의 골을 찾기 위하여 물의 흐름에 기반을 둔 방법을 사용하였다. 즉 밝기값의 표

면을 지형으로 간주하여 먼저 여기에 비를 뿌린다. 일정량의 비가 쏟아진 후, 낮은 곳에서는 물웅덩이가 형성된다. 물의 양이 많은 웅덩이가 즉, 골이 깊은 지형은 문자영역으로 처리하며 그외의 지형은 배경영역으로 처리한다. 지형의 특성에 따라 얇은 물웅덩이가 형성된 곳을 제거하기 위하여 비가 멈춘 뒤에는 어느 정도 햇빛을 가하여 물을 증발시킨다. 본 논문에서는 이러한 자연적 현상을 문서의 이진화에 적용하였으며 물체와 배경만으로 이루어진 영상에 대해 효과적으로 물체와 배경을 분리할 수 있다.

그림 1은 맑기값 영상을 1차원 단면으로 표현한 것이다. 그림에서 깊은 부분은 빗물에 의해 고인 물을 나타내며 얇은 부분은 햇빛에 의해 증발된 물을 나타낸다. 이러한 방법을 사용하면 맑기값의 형태에 대한 특성을 이용하여 지역적으로 영상을 이진화할 수 있다. 여기서 점선은 일정한 문턱값을 이용하여 진역적으로 이진화시킨 경우를 나타낸 것으로 영상의 국부적인 특성을 고려하지 않아 실제 문자부분에 해당하는 골부분을 찾아내지 못하게 된다. 반면 골의 특징을 이용하여 낮은 파인 부분을 검출하는 방법은 전체적인 영상의 맑기값 분포가 다르더라도 효과적으로 문자 부분을 검출할 수 있다.

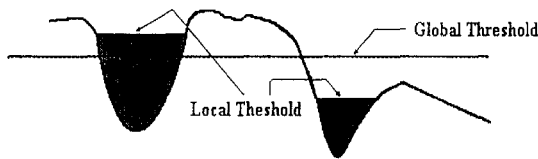


그림 1. 제안한 국부적용 이진화 방법  
Fig. 1. Proposed local thresholding scheme.

## 2. 비문자 부분의 분리

비문자 영역으로 분류된 영역의 경우 의미있는 문자이외의 나머지 영역을 분리하며 비문자 영역이 제거된 문자만의 영상에 대하여 개별문자 추출이 이루어진다. 추출된 개별문자는 인식을 위하여 인식과정으로 넘겨진다. 비문자 영역으로 판단하기 위하여 글꼴들의 크기나 형태 그리고 문서의 구조 등을 이용한다. 본 연구에서는 간단히 크기와 형태만을 이용하였으며 영역의 크기가 큰 경우나 그 형태가 세로나 가로로 긴 영역을 비문자로 추출하였다.

문서 영상을 영역화하기 위하여 연결화소에 기반

을 둔 방법을 사용하였다. 연결화소로 이루어지는 각 영역들은 하나의 영역으로 분리된다. 본 논문에서는 연결화소를 구하는 시간을 줄이기 위하여 전체 연결화소를 검출하지 않고 연결화소로 이루어지는 연결요소의 윤곽선만을 추출하였다. 물체의 윤곽선을 검출하기 위하여 8방향의 체인코드를 사용하였으며, 추출된 체인 코드의 길이와 영역에 대한 최적의 사각형의 크기, 모양 등을 이용하여 비문자 영역과 문자 영역으로 분리하였다. 비문자 영역으로 분리된 연결요소들은 영역내에 문자가 포함되는 경우를 고려하여 외곽선이 추출된 부분의 연결화소만을 비문자로 분리하였다. 따라서 비문자영역내에 문자가 포함된 경우 비문자 영역과의 겹침이 없는 문자들은 문자영역으로 남게 된다. 비문자로 판단하기 위해 사용한 크기값은 분류된 각각의 영역에 대한 크기값을 이용하여 계산하였다. 보통의 문서에서는 문자영역이 비문자 영역보다 많이 존재하며 각 문자들도 여러 자소로 이루어지기 때문에 영역화 과정에서 분류되는 영역의 크기는 대부분 문자의 크기값을 갖는다. 따라서 발생빈도가 큰 구간의 평균값은 그 문서의 문자 크기를 나타낸다. 이값을 이용하면 문자와 다른 크기값을 갖는 비문자 영역을 판단할 수 있다.

## 3. 개별문자 분리

개별문자의 분리는 비문자 영역으로 추출된 영상에 대하여 이진의 영역화 과정에서 추출한 체인 코드를 이용하여 분리한다. 각각의 영역으로 분리된 연결화소들은 하나의 연결화소가 하나의 문자를 이루는 경우도 있으나 대부분 여러 개가 하나의 의미있는 개별문자를 이룬다. 따라서 인식과정에 필요한 개별문자들을 추출하기 위해서는 분리된 영역들을 하나의 의미있는 개별문자로 군집화하는 과정이 필요하다. 영역화된 연결화소들을 하나의 개별문자로 추출하기 위하여 각 연결화소들간의 위치 관계와 형태 등을 이용하였다.

자소조합에 의한 한글 형태를 살펴보면 아래위로의 위치관계를 갖는 연결화소는 하나의 개별문자로 판단할 수 있다. 따라서 먼저 이러한 위치관계를 갖는 각각의 연결화소들이 수직으로 반이상 겹치는 경우 하나의 영역으로 군집화하였다. 이 경우 아래위로 각 연결화소가 겹치는 정도와 떨어져 있는 거리에 따라서 결합 여부를 판단하게 되므로 이웃하는 다른 개별문자에 의해 투영상에서 겹치는 경우에도 각각의 개별문자들을 분리할 수 있다. 다음 위아래의 위치관

계에 의해서는 하나의 영역으로 군집화되지 않는 수직 모음에 의한 경우를 고려하기 위하여 세로 형태의 연결화소에 대하여 앞뒤 관계를 이용하여 군집화하였다. 이 경우에 있어서 문제가 되는 것은 단순히 앞뒤의 위치관계만을 고려할 경우 영어의 'd', 'h', 'l', 'p', 'i', 그리고 숫자 '1'과 같은 글자들에 의해 하나의 영역으로 분리되어야 하는 영문자가 서로 군집화되어 하나의 개별문자로 분리되는 것이다. 따라서 이러한 오류를 해결하기 위하여 단순한 위치관계외에 한글의 수직모음에 대한 인식과정을 함께 사용하였다. 인식의 전처리 과정에서 인식과정을 필요로 하는 상호 교차성을 줄이기 위하여 영역화 과정에서 세로 형태로 분리된 연결화소만을 필요한 인식대상으로 하였다. 또한 인식의 정확도보다 한글의 수직모음류인지 영문인지를 구별하는 정도에 관심을 두었다. 한글의 수직모음을 인식하기 위하여 먼저 세선화 과정을 통하여 연결화소에 대한 끝점과 분기점, 그리고 교차점의 개수를 얻었다.

그림 2는 혼합문서 인식에 대한 본 연구의 전처리 과정을 보여주고 있다. 전처리 단계에서 그림과 문자부분을 분리하고, 또한 개별문자를 추출하여 인식단으로 넘겨주면 인식단에서는 여러가지 특징을 이용

하여 전처리 단계에서 넘어 온 입력문자를 인식한다. 출력단에서는 전처리 단계에서 분리한 그래픽부분과 인식한 문자부분을 결합하여 출력한다.

### III. 제안한 필기체 혼합문서인식 및 후처리

본 연구에서는 기본적으로 획분석에 의한 구조해석법을 이용하였다. 기존의 인쇄체 문자 인식법에선 특징값에 의한 방법이 주로 사용되어 왔으나, 필기체의 경우 문자의 왜곡에 의하여 특징추출 자체에 문제가 있어 사용하기 어렵다. 따라서 본 연구에서는 구조 해석법에 의하여 입력문자를 분석하고 이로부터 특징을 추출하는 방법을 사용하였다. 이러한 방법은 구조적방법과 통계적방법을 효과적으로 결합한 것이다. 또한 구조해석법의 적용에 있어서도 branch and bound 알고리즘을 사용하고, 여기에 비용함수로서 DP정합법을 사용하여 최적의 자소추출을 행할 수 있다. 이러한 자소추출 과정을 기친 후 자소간의 위치관계를 이용하여 문자의 유형에 맞는지 확인한다. 이러한 방법으로 개별문자를 인식한 후 후처리 과정으로서 단어사전을 이용한 정합을 통하여 최종문자를

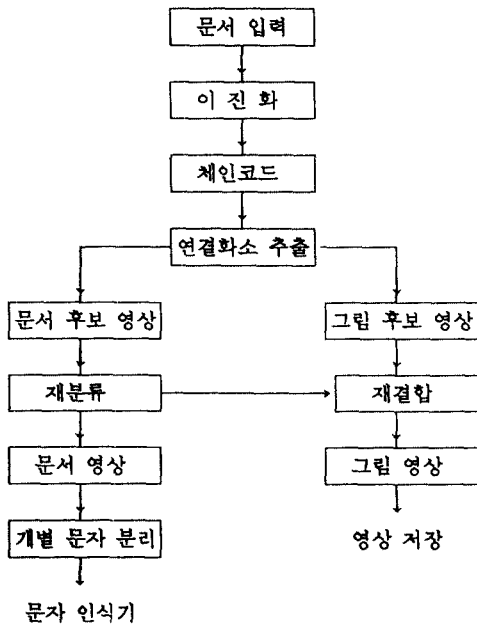


그림 2. 제안한 전처리 과정 흐름도  
Fig. 2. Flowchart of the proposed preprocessing step.

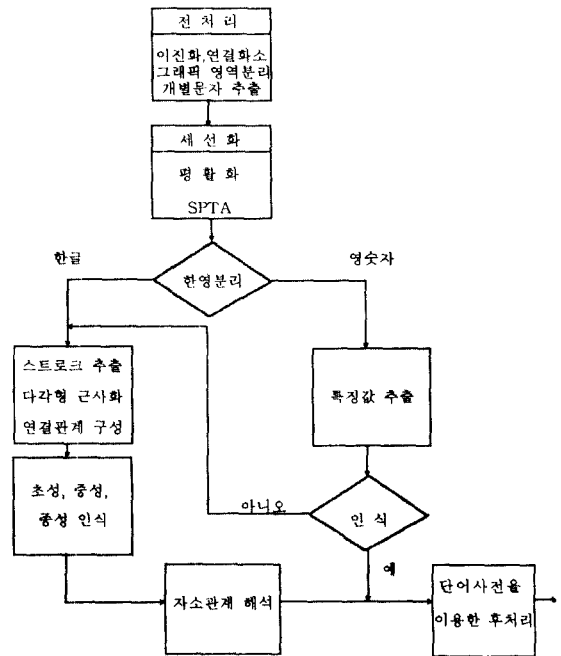


그림 3. 문자인식 흐름도  
Fig. 3. Flowchart of a character recognition step.

인식하였다. 여기서 단어사전 정합은 자소의 유사도와 후보문자를 이용하여 수행하였다. 그림 3은 본 시스템에서의 문자인식 흐름도이다. 먼저 세선화 과정을 거친 후 한글과 영문으로 분리하고 영문의 경우 특징값을 이용하여 계층적으로 자소를 분류한다. 만약 영문으로 인식되지 않을 경우 한글로 인식을 다시 시도한다. 한글 인식부분에선 세선화된 결과로부터 세그먼트를 추출한다. 세그먼트를 추출한 후 세그먼트간의 위치와 접촉관계를 그래프 구조로 구성한다. 다음은 DP정합 비용함수값에 따라 branch and bound 알고리즘을 이용하여 기본자소를 추출하고 추출된 자소에서 구조해석법을 이용하여 확장자소를 인식한다. 이러한 방법으로 초성을 인식하고 중성과 종성은 같은 방법으로 자소분류도에 바탕을 두어 인식된다. 각 자소를 추출한 후 자소사이의 위치관계를 이용하여 후보문자를 선택한다. 문자인식을 수행한 후 후보문자와 문자의 유사도를 고려한 단어사전을 이용하므로써 오인식 문자를 수정하여 인식률을 높인다.

### 1. 세선화

본 연구에서 사용한 문자인식법은 획분석을 기본으로 하므로 세선화 과정이 필요하다. 세선화는 본래 패턴의 모양과 연결성을 유지하는 골격을 추출하는 과정을 말한다. 보통 세선화 과정은 반복적으로 경계점을 지워나가는 방법인데 경계점을 결정하는 판단 기준에 따라 여러가지 세선화 방법이 있다.<sup>13)</sup> 본 과제에선 세선화 알고리즘으로 빠르고 좋은 성능을 가진 것으로 알려진 SPTA(safe point thinning algorithm)를 사용하였다.<sup>14)</sup> SPTA 알고리즘은 패턴의 연결성을 깨지 않고 심한 평활화를 일으키지 않는 것이 특징이다.

### 2. 한글과 영숫자 분리

본 시스템에선 한글과 영숫자를 인식하기 위하여 이들의 특성이 다르므로 각각 다른 방법을 사용하였다. 따라서 이들의 분리과정이 필요하게 된다. 기존의 분리 방법에선 한글의 수직모음과 수평모음을 인식함으로써 이들을 분리하나 필기체의 경우 이러한 방법을 사용할 경우 오류가 발생할 수 있다.

본 방법에선 한글과 영숫자를 분리하기 위하여 끝점의 개수와 연결요소의 개수를 이용하였다. 영문은 'l', 'j'를 제외한 모든 문자가 한개의 연결요소로 구성되며 끝점의 개수가 4개 이하이나 실제로 잡음이나 필기형태의 변화 그리고 세선화과정에서 발생하는

삐침 등으로 4개 이상의 끝점이 존재할 수 있다. 또한 한글의 경우 자소의 접촉이 발생한다해도 한개의 연결요소로 구성되며, 끝점의 개수가 2개 이하인 경우는 발생하지 않는다. 본 방법에선 한개의 연결요소로 구성되고 끝점의 개수가 2개 이하일 경우는 영숫자로 인식하며, 한개의 연결요소로 구성되고 끝점의 개수가 2개보다 클 경우 먼저 영숫자 인식을 시도하며, 영숫자로 인식되지 않으면 한글로 인식하게 된다. 이외의 모든 경우는 한글로 인식하게 된다. 만일 영문의 경우 'l', 'j'를 제외한 문자가 한개 이상의 연결요소로 분리될 경우 오류가 발생하는 단점이 있다. 그리고 'l', 'j'는 문자의 폭을 이용하여 문자를 분리하였다.

### 3. 다각형 근사화와 세그먼트 추출

본 방법은 획분석에 바탕을 두므로 입력패턴의 기본 획으로 분리하는 과정이 필요하다. 세선화 과정을 거친 후 곡선 형태로 된 획을 직선의 조각으로 분리하는 다각형 근사화를 수행한다.<sup>15)</sup> 다각형 근사화 방법으로 획 위의 두 점과 그 사이의 점들 중에 가장 먼 점까지의 거리가 임계치보다 클 때 그 점을 할극점으로 하는 방법과 체인코드의 1차 미분값과 2차 미분값을 이용하는 방법이 있다. 본 연구에서는 구현이 용이하며 비교적 좋은 성능을 가진 2차 미분값을 이용하는 방법을 사용하였다. 이 방법은 체인코드와 2차 미분값을 부호가 바뀌는 점들 중에서 1차 미분값이 임계치보다 큰 점을 할극점으로 선택한다.

### 4. 구조해석을 위한 한글자소의 유형분리

한글의 구조적 특징에 의한 문자의 구조분석은 문자인식과 직결된다. 더구나 자소간의 연관관계는 물론 자소내에서도 기본자소에 몇 개의 획이 추가되어 새로운 자소가 만들어지므로 자소 구조분석은 매우 중요하다. 본 연구에서 사용하는 구조해석법은 세그먼트의 속성 그리고 세그먼트 사이의 연관관계를 이용하였다. 여기서 세그먼트의 속성으로 세그먼트의 길이와 직선성 그리고 방향각을 선택하였으며, 세그먼트간의 연관관계로는 연결성과 연결된 경우 각도의 변화량과 연결되지 않은 경우 세그먼트간의 위치관계를 이용하였다. 여기서 항상 연결성의 관계를 갖는 자소와 추가 세그먼트에 의해 다른 자소를 생성할 수 있는 기본이 되는 자소를 기본자소로 하고, 연결될 필요없이 세그먼트의 추가로 생성할 수 있는 자소를 확장자소로 정의한다. 이러한 방법으로 자소들을 기본자소를 바탕으로 하여 계층적으로 나누어 간다.

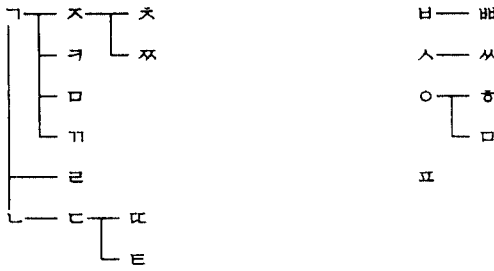


그림 4. 초성 자소 분류도  
Fig. 4. Classification diagram of head consonants.

그림 4는 구조해석법을 위한 초성의 분류도를 나타낸다. 그림에서 ‘ㄱ’, ‘ㄴ’, ‘ㅇ’은 연결성에 의해 표현될 수 있고, 이로부터 파생되는 자소들이 한글 자소의 대부분을 형성하므로 매우 중요한 자소이다. ‘ㅂ’, ‘ㅅ’, ‘ㅇ’은 직선의 세그먼트를 기본으로 방향이 다른 세그먼트가 다른 위치 연관관계를 가지고 자소를 형성한다. 이러한 기본자소를 바탕으로 접촉과 관계없는 몇 개의 세그먼트의 추가로 새로운 자소들이 된다. 이처럼 구조적으로 자소가 구성되므로 분류도에 바탕을 두어 구조해석법을 적용하였다.

5. PD정합을 이용한 branch and bound 알고리즘

본 연구에서는 구조해석법을 이용한 자소인식을 위하여 기본 자소를 추출한 후 주위 세그먼트의 존재와 속성 그리고 기본 자소와의 연관관계를 이용한다. 결국 기본자소 추출은 자소인식 및 문자인식에 있어서 가장 중요한 문제이다.

기본적인 구조 해석법에서 기본 패턴과 속성 그리고 연관 관계만을 이용함으로써 필기자의 필기왜곡에 대처할 수 없는 단점이 있다. 또한 통계적인 방법은 이러한 왜곡이나 부분적 획의 손실에는 적용할 수 있으나, 실제 필기체의 경우에는 특징을 추출하는 자체가 어려워 적용하기 어려운 단점이 있다. 따라서 본 연구에선 구조 해석법과 통계적인 특징을 함께 이용하는 방법을 채택하였다. 이것을 위하여 본 연구에선 branch and bound 알고리즘을 이용하여 자소 추출을 행하고, 여기에 필요한 비용함수로서 DP 정합을 사용하여 최적의 경로를 찾고 최소한의 비용을 갖는 세그먼트의 집합을 자소로 인식하는 알고리즘을 제안하였다. DP 정합은 1차원 탄력정합(elastic matching)으로서 보통 1차원 신호로 모델링이 가능한 음

성처리<sup>(6)</sup>나 온라인 문자인식<sup>(7,8)</sup>에 사용되어 왔다. 그러나 DP 정합의 우수성에도 불구하고 오프라인 문자 인식의 경우 입력영상 2차원이므로 적용에 어려움이 따른다. 본 연구에선 이러한 문제를 해결하기 위해 입력영상에서 세그먼트의 연결성이 존재하는 직선 세그먼트 및 기본자소 인식에 DP 정합 비용함수값에 따른 branch and bound 알고리즘을 적용하여 자소 추출 및 기준자소의 소속정도를 계산하여 문자인식을 수행하였다. 이러한 방법으로 자소간의 접촉이나 세선화 과정에서 발생하는 잡음에 강한 자소추출을 행하였다.

입력패턴과 기준패턴사이의 상이도는 정규화된 비율과 세그먼트의 방향각의 차를 선형결합으로 주어진다.  $A_i^k$ 과  $AR_i^k$ 는 각각  $k$ 번째 자소의 입력과 기준 패턴의  $i$ 번째 세그먼트의 방향을 나타낸다. 또한  $L_i^k$ 와  $LR_i^k$ 는 각각  $k$ 번째 자소의 모든 세그먼트의 길이로 정규화된 기준 패턴과 입력패턴의 획의 길이이다. 기준패턴의 한개의 세그먼트에 대응되는 입력패턴의 세그먼트들은  $C_j^k$ 개로 구성되며, 그것의 방향과 길이는 각각  $A_{ij}^k$ 와  $L_{ij}^k$ 로 나타낸다. 여기서  $j$ 는 1에서부터  $C_j^k$ 이다.  $M$ 개의 세그먼트로 구성되는  $k$ 번째 기준 패턴과 입력패턴 사이의 각도차는  $c_k$ 로 정의된다. 각도차는  $180^\circ$ 보다 작으며, 기준패턴의 한 세그먼트에 대응되는 입력패턴의 모든 세그먼트는  $L_{ij}^k$ 로 나타낸다. 방향각에 의한 상이도는

$$c_k = \sum_{i=1}^M \sum_{j=1}^{C_i^k} L_{ij}^k \times \min(|AR_i^k - A_{ij}^k|, 360 - |AR_i^k - A_{ij}^k|) \tag{1}$$

으로 주어지며, 정규화된 길이에 의한 상이도는

$$l_k = \sum_{i=1}^M |LR_i^k - \sum_{j=1}^{C_i^k} L_{ij}^k| \tag{2}$$

로 주어진다. 결국  $k$ 번째 자소에 대한 전체 상이도  $e_k$ 는  $d_k$ 와  $l_k$ 의 선형결합으로 주어지,

$$e_k = d_k + \alpha \times l_k \tag{3}$$

으로 표시된다. 여기서  $\alpha$ 는 실험적으로 결정되는 상수이다.

그림 5는 기준패턴 ‘ㄱ’과 입력패턴사이의 상이도를 계산하기 위하여 대응되는 세그먼트를 보여준다. 기준패턴 ‘ㄱ’은 2개의 세그먼트로 구성되어  $M$ 은 2가 된다. 기준패턴의  $LR_1^k$  세그먼트는  $L_{11}^k, L_{12}^k$  그리고

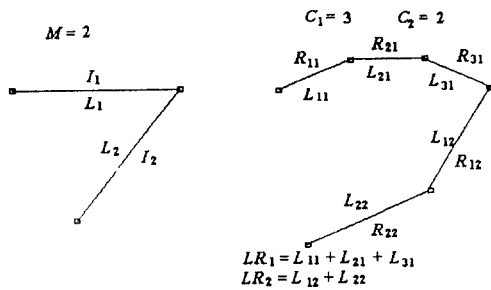


그림 5. 기준패턴과 입력패턴사이의 상이도  
Fig. 5. Dissimilarity between reference and input patterns.

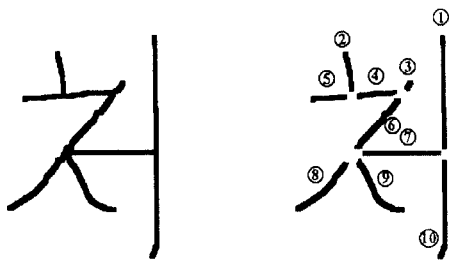


그림 6. 세그먼트 추출후 입력문자  
Fig. 6. Input character after segment extraction.

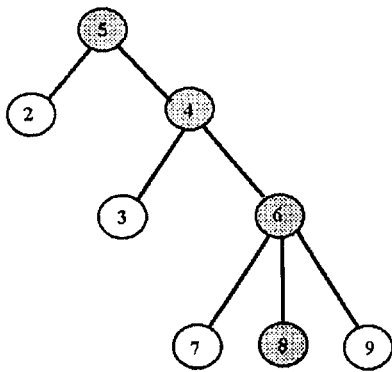


그림 7. Branch and bound 알고리즘에 의한 세그먼트 탐색과정  
Fig. 7. Segment search using the branch and bound algorithm.

$L_{11}^k$ 로 정합되고,  $LR_2^k$ 는  $L_{21}^k$ 와  $L_{22}^k$ 로 정합된다. 입력패턴과 기준패턴사이의 상이도는 정합된 세그먼트간의 특징값으로 계산된다. 이러한 방법으로 계산된 상이도를 branch and bound 알고리즘에 적용한다.

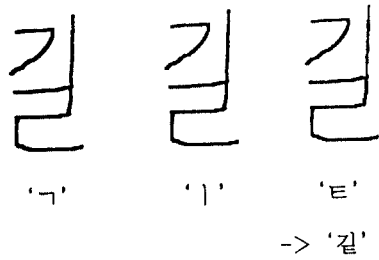
그림 6은 세션화 과정과 세그먼트 추출 후 세그먼트의 번호를 나타낸 것이다. 이와같이 세그먼트 추출 후 세그먼트간의 연결관계를 그래프 구조로 나타낸다. 초성 자소 'ㄱ'을 인식하기 위하여 최상단 왼쪽에서 'ㄱ'의 첫번째 획에 가장 가까운 수평획 5를 찾는다. 다음 이 획에 접속된 2번획과 4번획중 자소 'ㄱ'을 형성할 수 있는 정도를 DP정합 비용함수를 이용하여 계산한 후 이 값이 임계치보다 큰 비용을 가지면 제거하여 나간다. 결국 세그먼트 2는 제거되고 세그먼트 4가 추출된다. 다음 4번 세그먼트에 접속된 5번과 6번 세그먼트의 비용이 계속 계산된다. 여기서 사용된 DP정합은 기준자소의 방향각과 기본자소의 길이를 이용하여 비용함수를 계산하였다. 이것은 기존의 온라인 문자인식 방법에서 사용한 방법이다.

그림 7은 branch and bound 알고리즘에 따른 세그먼트의 탐색과정을 보여주고 있다. 여기서 색칠된 세그먼트가 기준자소 'ㄱ'에 가장 잘 정합되는 세그먼트를 보여준다. 또한 DP정합은 정합도뿐만 아니라 기준자소의 한 획과 정합되는 입력패턴들도 추출할 수 있다. 이 예의 경우는 기준자소 첫번째 획에 입력패턴의 5번과 4번 획이 정합되고, 기준자소 두번째 획에 입력패턴의 6번째와 8번째획에 정합됨을 보여준다. 이러한 제안한 방법으로 자소간의 접촉이나 세션화 과정에서 발생하는 잡음에 강한 자소추출을 행하였다.

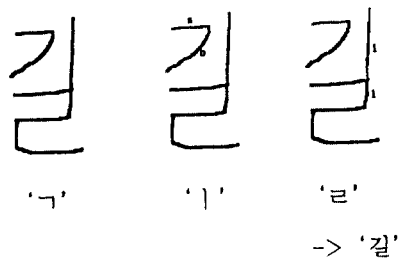
### 6. 자소간의 접촉관계를 고려한 자소추출

DP정합 방법은 기본자소의 인식뿐만 아니라 기준자소에 대응되는 세그먼트의 집합을 찾을 수 있다. 이러한 특성을 이용하여 기존의 인식된 획을 제거해가는 방법에서 발생하는 획간의 접촉으로 문자내 한 자소의 추출과정에서 다른 자소까지 연결되어 추출되어 다음 자소 추출이 실패하는 문제를 해결할 수 있다. 이러한 문제는 자소간의 접촉관계를 고려하여 줌으로써 자소 인식시 초기의 자소추출과 다음 자소추출후, 자소의 접촉 가능성을 고려하여 초기에 추출된 모든 세그먼트를 제거하지 않으므로써 해결하였다.

그림 8은 이러한 예를 보여 준다. 그림 8(a)는 기존의 방법으로 중성의 자소 추출과정에 의하여 종성 자소인식에 영향을 미치는 것을 보여 주고 있다. 그림



(a)



(b)

그림 8. 자소 접촉관계를 고려한 자소 추출

- (a) 기존의 자소추출 방법
- (b) 제안한 자소추출 방법

Fig. 8. Phoneme extraction by considering touching relationship

- (a) Conventional phoneme extraction method,
- (b) Proposed phoneme extraction method.

8(b)는 본 방법에서 이러한 문제를 해결하기 위하여 종속의 인식시 중성과의 접촉관계를 고려하여 자소를 추출한 예이다. 그림 8(a)에서 자소 'ㅣ'이 중성 'ㄱ'과 접촉될 경우 기존의 방법<sup>9)</sup>은 'ㅣ'의 추출 과정에 의하여 'ㄱ'의 추출시 문제가 생김을 보여준다. 그림 8(b)는 자소 'ㅣ'를 인식하고 다음 자소 'ㄱ'을 인식할 때 전 자소가 'ㅣ'였으므로 접촉의 가능성을 고려하여, 중성으로 인식된 모든 세그먼트를 제거하지 않고, 중성에 포함될 가능성이 있는 세그먼트까지 함께 고려하여 자소추출을 행하는 과정을 보여준다.

이러한 자소간의 접촉 가능성을 표 1에 나타내었다. 표 1에서 인용부호 안의 자소는 중성을 의미한다. 이 표에 보인 경우 이외에도 더욱 다양한 자소의 접촉이 있을 수 있으나, 본 연구에서는 실험문서에서 자주 발생하는 경우만을 고려하였다.

표 1. 자소간의 접촉관계

Table 1. Touching relationship between phonemes.

자 소	자 소
ㄱ, ㄲ	ㄷ, ㄸ, ㅌ, ㄱ, ㄲ
ㄴ	ㄷ, ㅌ, ㅍ, ㅍ, ㅅ, ㅆ
ㄷ, ㄸ	ㅌ
ㄹ	ㅌ, ㄱ, ㅋ
ㅋ	ㄱ, ㅋ, ㄷ, ㄸ, ㅌ, ㄱ, ㄲ
ㅌ, ㅍ, ㅍ, ㅌ, ㄱ, ㅋ, ㅋ, ㅅ, ㅆ, ㅅ, ㅆ, ㅅ, ㅆ, ㅅ, ㅆ, ㅅ, ㅆ, ㅅ, ㅆ, ㅅ, ㅆ	'ㄱ', 'ㄲ', 'ㄷ', 'ㄸ', 'ㄹ', 'ㄱ', 'ㅋ', 'ㅌ', 'ㅍ', 'ㅆ', 'ㅅ', 'ㅆ', 'ㅅ, ㅆ', 'ㅅ, ㅆ'

7. 구조해석법

문자가 몇개의 기본적인 패턴의 조합으로 구성된다면, 이러한 부패턴의 계층구조로서 문자를 나타낼 수 있다.<sup>(10)</sup> 한글의 경우 문자의 구성이 자소의 조합으로 구성되어 있으며 각 자소는 더 작은 기본 자소와 획들로 구성된다. 특히 한글의 경우는 문자의 유형이 여러 개 존재하는 특징을 가지고 있다. 본 논문에선 한글 유형으로 기존의 6가지 유형에 바탕을 두었다.<sup>(11)</sup> 본 방법에선 구조적 방법에 바탕을 두어 자소인식을 하며, 또한 자소 인식에 있어서도 기본자소와 몇개의 획을 인식하는 방법을 사용하였다. 이러한 과정을 위하여 문자의 계층구조를 패턴분법으로 기술하여야 한다.

본 논문에선 문자를 조성, 중성 그리고 중성으로 나누고 각 자소는 자소분류도에 따라 계층적으로 구성하였다. 예를 들어 '착'의 경우 조성 'ㅊ', 중성 'ㅌ' 그리고 중성 'ㄱ'으로 구성되며 이들간의 연관관계는 유형 4에 속한다. 조성 'ㅊ'은 기본자소 'ㄱ'에 2개의 직선 세그먼트를 합하여 구성되며 이들은 적당한 위치관계를 가지고 구성되어 있다. 이러한 문자구성 계층에서 한 단계 더 내려가며, 기본자소 'ㄱ'은 기본 패턴인 직선 세그먼트로 구성된다. 같은 방법으로 중성과 중성도 표현된다. 이러한 경우 속성문법의 S는 한 개의 문자를 말하며, 비종단기호는 자소와 기본자소를, 종단기호는 기본 패턴인 직선 세그먼트를 말한다. 구문생성 규칙은 문자유형에 따라 6가지 규칙이 있으며 또한 각 자소는 기본자소와 몇개의 세그먼트로 구성되고 그들간의 상하, 좌우 또는 크기비율 등의 관계로 구성된다.

8. 중성 및 중성자소 인식

본 방법은 구조해석법에 DP정합을 사용함으로써 구조적 특성과 통계적 특성을 모두 이용하였다. 이러



한 방법으로 사소의 조성은 초성자소 분류도에 따라 기본 자소를 인식하고 구조해석법에 바탕을 두어 확장자소를 인식하였다.

한글의 경우 중성의 인식은 중성의 인식과 함께 고려되어야 하며, 한글문자인식에 있어서 가장 중요한 문제이다. 이러한 특성은 한글이 중성 중심의 문자라는 점과 중성의 형태에 따라 중성과 다른 인관관계를 가지며 중성이 존재하지 않을 수도 있기 때문이다. 중성은 인식과정에서 중성과 접촉을 발생할 수 있으며 중성이 있는지 없는지 판단할 수 없기 때문에 수평모음과 중성의 분리는 어려운 문제이다. 본 연구에선 이러한 오류의 가능성을 줄이기 위하여 수직모음 인식을 먼저 수행하고 다음 수평모음의 인식을 시도함으로써 오류의 가능성을 줄였다. 이것은 수직모음의 인식과정에는 오류의 가능성이 적고, 수직모음의 인식되면 수평모음의 존재 가능성이 줄어들기 때문이다. 수직모음의 존재를 확인하고 존재한다면 수직모음 'ㄱ, ㅋ, ㆁ, ㆁ, ㆁ, ㆁ, ㆁ, ㆁ'를 인식한다. 이러한 수직모음이 인식되면 수평모음의 존재 가능성은 표위치에서처럼 줄어든다. 예를 들면 수직모음 'ㄱ'이 인식되었을 경우 수평모음이 존재하지 않거나 수평모음은 'ㅇ'일 수 밖에 없다는 것이다. 특히 'ㅋ, ㆁ, ㆁ, ㆁ'의 경우에는 수평모음이 존재하지 않는 경우이므로 중성과의 혼동 가능성이 줄어들게 되는 것이다. 이처럼 탐색영역이 줄어들기 때문에 모음추출시 오류를 줄일 수 있다.

9. 영숫자 인식

본 연구에서는 영숫자 인식을 위한 특징으로 흑백화소 변화수와 거리 그리고 부영값을 이용하였다.<sup>11)</sup> 기본 문자인식 과정은 위의 특징에 의하여 자소를 분리하고 인식하는 방법을 채택하였다. 그러나 기존의 특징 추출단계에서는 이러한 특징을 문자의 특정위치에서 구하였으나, 실제 필기체 문자의 경우 특징의 위치가 변할 수 있으므로 특징 추출의 안정성을 위하여 본 방법에선 영역을 크게 두고 영역내에서 특징값의 최빈값과 상위 3개의 평균값 그리고 하위 3개의 평균값을 이용하여 자소의 굴곡이나 잘못된 빼침 등에 의한 영향이 적도록 하였다. 거리 특징값을 찾을 때 그림 9(a)와 9(b)에서와 같이 일정한 위치를 고정시키면 특징값이 크게 변하게 된다. 또한 평균값이나 중간값도 불안정하고 특징으로 사용하기에 곤란하다. 이러한 경우 그림 9(c)와 같이 영역을 크게 주고

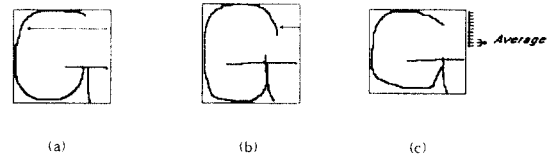


그림 9. 거리특징 추출 방법

- (a) 고정된 위치에서 한 입력패턴의 거리특징
- (b) 고정된 위치에서 다른 입력패턴의 거리특징
- (c) 제안한 특징추출 방법에 의한 거리특징

Fig. 9. Methods for extracting distance features.

- (a) Distance feature for one input pattern at fixed position.
- (b) Distance feature for another input pattern at fixed position.
- (c) Distance feature by the proposed feature extraction method.

여기서 거리가 큰 3개를 평균낸 특징이 'G'를 인식하기에 더욱 적합하다.

본 연구에서 영숫자인식은 여러가지 특징들을 이용하여 계층적으로 분류하는 방법을 사용하였다. 먼저 흑백화소 변화수에 따라 입력 패턴을 분류하고 분류된 패턴을 흑화소까지의 거리를 이용하여 분류하였다. 다음 계층에선 부분 부영값을 이용하여 자소를 분류하였다. 그리고 최종으로 끝점이나 분기점의 위치를 이용하여 패턴을 인식하였다.

10. 단어사전을 이용한 후처리과정

본 방법에선 문자인식기에서 만들어진 후보문자들과 단어사전의 각각의 단어들 사이에 유사도<sup>(13)</sup>를 계산하고 이들 중에 임계치를 넘으면서 가장 잘 정합되는 단어로 오류를 정정한다. 한글에 있어서 유사문자에 관한 연구는 주로 획의 개수와 형태에 따라 유사도를 판단한다. 본 방법에선 자소분류도에 기초를 둔 유사자소 행렬을 만들었다. 이러한 유사도 판단은 결국 패턴의 획의 숫자와 배치관계에 바탕을 둔 것이다. 이것은 통계적 문맥표현법에선 문자의 혼동확률과 같은 개념으로 이것을 한글의 경우 자소 혼동확률로 적용한 것이다. 이러한 자소혼동 확률과 함께 문자인식기에서 발생하는 후보문자를 고려함으로써 문자인식기에 생성되는 국부 정보를 함께 고려하였다.

$X_i(i=1, \dots, M)$ 를 문자인식기에서 생성되는 후보단어라 하고,  $X_i = x_{i1} \dots x_{im(i)}$ 는  $m(i)$ 개의 문자로 구성되어 있다고 하자. 또한 각 문자  $x_{ij} = hi_{ij} vi_{ij} bi_{ij}$

로 초성, 중성 그리고 종성으로 구성되어 있다. 또한 단어사전의 단어들을  $W_i(i=1, \dots, N)$ 라 하고, 각 단어는  $W_i = w_{i1} w_{i2} \dots w_{im(i)}$ 이고,  $w_{ij} = hr_{ij} vr_{ij} br_{ij}$ 라 하자.

후처리에서 정합 단어를 찾기 위하여 단어 유사도는

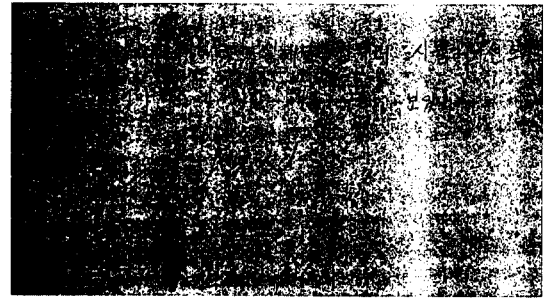
$$\begin{aligned}
 S(j) &= \text{Min}_{i=1, \dots, M} \{ S_{word}(X_i, W_j) \} \\
 &= \text{Min}_{i=1, \dots, M} \left\{ \sum_{k=1}^{m(i)} S_{char}(x_{ik}, w_{jk}) \right\} \\
 &= \text{Min}_{i=1, \dots, M} \left\{ \sum_{k=1}^{m(i)} (S_{head}(h_{iik}, hr_{jk}) + \right. \\
 &\quad \left. S_{vowel}(v_{iik}, vr_{jk}) + S_{bottom}(b_{iik}, br_{jk})) \right\} \quad (4)
 \end{aligned}$$

로 정의되며, 이 값을 최대화시키는  $S(j)$ 가 임계값 이상이면,  $W_j$ 는 정합단어로 선택된다. 여기서  $S_{head}$ ,  $S_{vowel}$  그리고  $S_{bottom}$ 은 자소간의 유사도이다. 이 값은 행렬로 주어지며 인식기에 발생하는 혼동정도를 나타내는 값으로 자소 분류도에서 인접된 자소는 유사도가 크며, 멀리 떨어져 있는 자소간에는 유사도가 적고 연결성이 없는 자소의 유사도는 최소의 값을 갖는다.

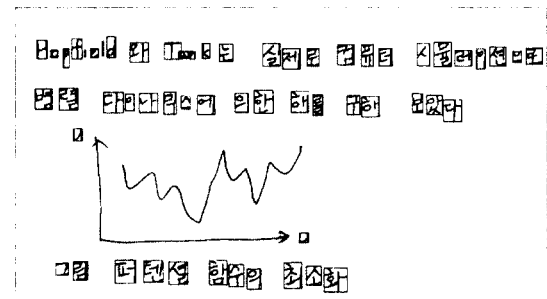
이러한 후처리를 사용함으로써 자소간의 유사도 측정으로 자소 추출에 의한 오류를 정정할 수 있었다. 그러나 문자 사이의 유사도는 자소 사이의 유사도의 합으로만 표현할 수 없는 경우도 있으며, 이러한 방법은 한글의 용언활용, 조사 등을 분석하는 방법이 아니어서 앞으로 많은 연구가 필요하다.

#### IV. 실험결과 및 분석

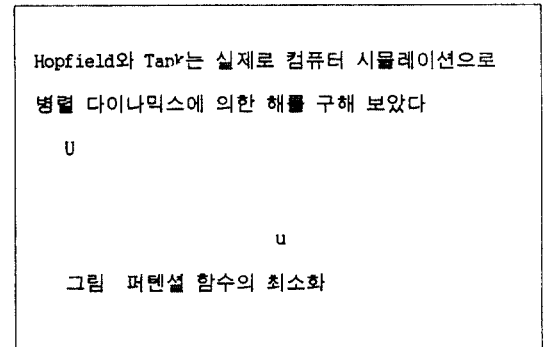
본 연구는 그래픽을 포함한 필기체 한글과 영숫자로 구성된 혼합문서 인식을 위한 알고리즘 구현에 대한 연구이다. 전처리 과정에서 체인코드를 이용하여 그래픽 영역을 분리하였고 몇개의 특징과 한글의 유형을 이용하여 개발문자영역을 추출하였다. 분리된 개별문자를 한글과 영숫자로 분리하고 각각 따로 인식하였다. 본 실험에선 문서에 50자에서 100자 정도의 문자를 포함한 실험문서 12장에 대하여 98%의 인식을 얻었다. 오인식은 주로 마침표나 영문의 대소문자 구별과 문서 상태의 불량으로 획이 손실되는 경우에 발생하였다. 본 시스템은 기존의 시스템<sup>(12)</sup>에 같은 필기체 문서를 사용하였을 경우보다 3~5%의 인



(a)



(b)



(c)

그림 10. 그래픽을 포함한 필기체 혼합문서 1

- (a) 입력영상
- (b) 이진화와 개별문자
- (c) 인식결과

Fig. 10. Handwritten mixed document 1 containing a graphic image.

- (a) Input image,
- (b) Binarization and character isolation,
- (c) Recognition result.

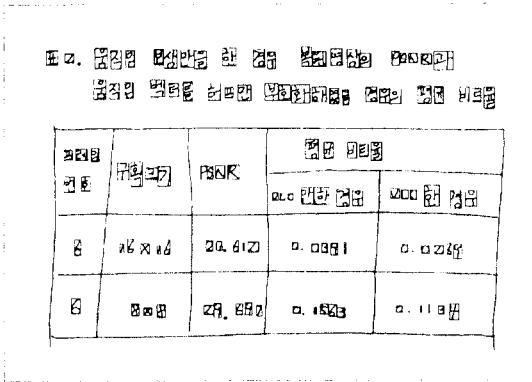


표 2. 움직임 보상을 한 경우 복원영상의 PSNR과 움직임 벡터를 허프만 부호화하였을 경우의 평균 비트율

프레임 번호	구획크기	PSNR	평균 비트율	
			VLC 안한 경우	VLC 한 경우
6	16x16	20.612	0.0391	0.0269
6	8x8	29.697	0.1563	0.1134

그림 11. 도표를 포함한 필기체 혼합문서 2  
(a) 입력영상  
(b) 이진화 추출결과  
(c) 인식결과

Fig. 11. Handwritten mixed document 2 containing a table.

- (a) Input Image,
- (b) Binarization and character isolation,
- (c) Recognition result.

식물 향상을 얻었다. 기존의 시스템은 인쇄체를 가정하여 특징을 추출하는 방법이므로 필기체의 경우 특징의 위치 등의 변화가 심하여 인식률이 저하되었다. 그러나 본 방법에선 획을 분석함으로써 더 좋은 성능을 보였다. 필기체의 경우 특징 추출의 어려움을 보여준 것이다.

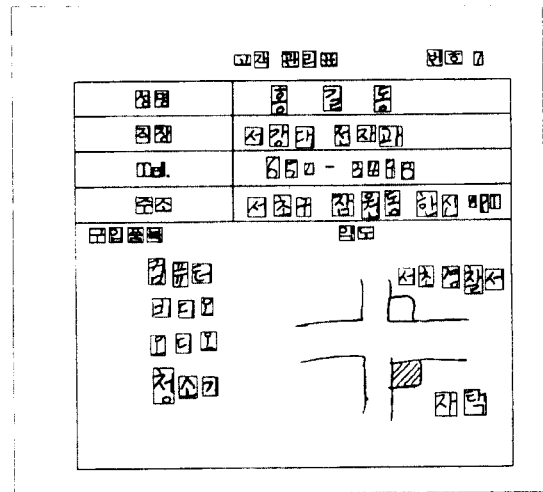
그림 10(a)는 실험에서 사용한 한글과 영문으로 구성된 입력영상으로 문서의 해상도는 300 dpi (dots per inch)이고 256 그레이레벨 값을 가진다. 영상 취득은 스캐너를 이용하였으며 TIFF 형식으로 저장하여 처리하였다. 실험은 Workstation R-4000상에서 C언어로 구현하였다. 이 영상을 국부적응 이진화 방법을 사용하여 이진화를 수행한 후 그래픽영역을 분리하고 개별문자를 추출한 결과를 그림 10(b)에 나타내었다. 이진화 결과가 잡음 등에 크게 영향을 받지 않고 바르게 이진화가 이루어졌음을 보여준다. 여기서 개별 문자로 추출된 문자는 최소 사각형으로 테두리를 씌웠으며, 그래픽 영역으로 분리된 것에는 테두리를 표시하지 않았다. 그림에서 보듯이 "Tank"에서 'T'와 'a'의 경우 투영에 의한 방법에서 개별문자 추출이 불가능하나 제안한 방법에서 바르게 추출되었다. 그러나 이 문서의 경우에는 마침표가 잡음제거 과정에서 제거되는 오류가 발생하였다. 그림 10(c)는 인식 결과로서 자소간의 접촉이 복잡한 '컴'자나 '물', '막', '화' 등과 같은 문자도 바르게 인식하였다. 이 문서에서 마지막 줄의 "함수의"에서 '의'사의 경우 자소의 접촉을 고려하여도 '의'와 '의'사로 인식될 수 있다. 이러한 경우는 실제로 문자 자체만의 인식으로는 인식할 수 없다. 이러한 경우 문맥이나 조사의 쓰임을 이용해야 하나 이러한 과정은 실제로 모든 경우에 사용하기에는 어려움이 있다. 본 방법에선 '의'의 경우 특별히 앞의 단어가 단어사전에 정합되며 '의'사 다음에 문자가 존재하지 않을 경우 '의'자로 인식한다.

그림 11(a)는 도표를 포함한 혼합문서이다. 이것의 이진화를 수행한 후 도표는 그래픽 부분으로 분리되고, 개별문자를 추출한 결과 그림 11(b)에 나타내었다. 개별문자의 추출결과를 보면 '상'과 '의'의 경우 투영을 이용한 경우 분리가 불가능한 경우이고, "보상만을"에서 '만'의 경우 한글문자의 유형과 중성자소의 인식을 통하여 개별문자로 추출된 경우이다. 여기서 일부 '이' 개별문자 추출중에 제거되었는데 이는 잡음제거시 임계치에 미치지 못해 제거되었다. 그림 11(c)는 문자인식의 결과이다. 여기서 'P', 'S', 'N', 'V', 'L', 'C' 및 그밖의 숫자들은 대부분 영숫자로

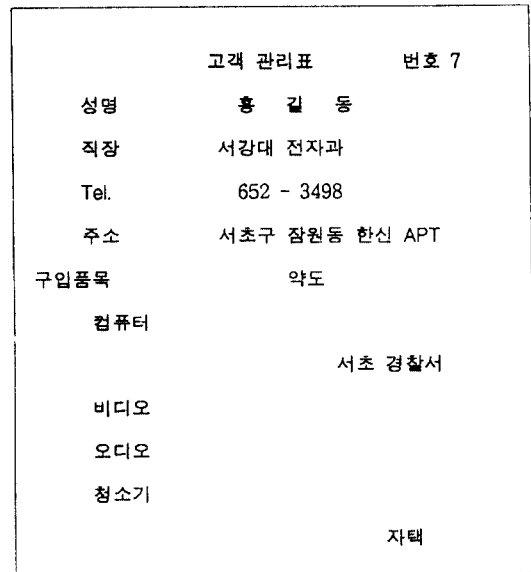
분리되어 옳게 인식되었으나, 'R'과 같은 경우 빼침 등에 의하여 불필요한 끝점 등이 발생하여 한글인지 영문인지 구별할 수 없다. 이러한 경우는 먼저 영문 인식 부분에서 'R'로 정인식되어 한글인식부분으로 넘어가지 않았다. 여기서 '임'자나 '관' 등은 자소의 접촉이 문제를 일으킬 수 있으나 자소의 접촉을 고려함으로써 문자를 옳게 인식하였다.

본 논문에서 제안한 알고리즘의 효율성을 보이기 위하여 인쇄체 및 필기체 혼합문서에 대하여 실험하였다. 이 문서는 인쇄체 및 필기체 혼합문서로 고딕체, 명조체 한글 그리고 영문을 모두 포함하고 있는 실험영상이다. 또한 도표부분과 그래픽부분을 가지고 있다. 이것을 국부적응 이진화 알고리즘을 사용하여 이진화한 결과를 수행하고 한글의 문자 유형과 부분적 수직모음의 인식을 통하여 개별문자를 추출한 결과를 그림 12(a)에 나타내었다. 인쇄체와 달리 필기체의 경우 연필로 필기되어 밝기특성이 다르나 바르게 이진화된 것을 볼 수 있다. 이진화 결과에서 잡음 제거 과정을 거침에도 불구하고 마침표와 크기가 비슷한 경우 완전히 잡음을 제거하지 못했으나, 개별 문자 추출단계에서 고수준(high-level) 지식을 사용함으로써 문자로 추출되지 않았다. 전처리 과정에서 영상의 밝기값만을 이용하지 않고 인간이 문서나 영상을 인식하는데 사용하는 지식을 사용함으로써 더욱 복잡한 문서를 인식할 수 있을 것이다. 앞으로 이러한 고수준 지식을 전처리과정 등에서도 이용할 수 있는 알고리즘 개발이 필요하다. 개별 문자 추출과정에서 영문과 한글 수직모음과 비슷하여 생길 수 있는 영문 'l'자도 주위의 문자 정보를 함께 이용함으로써 문자를 바르게 추출하였다.

그림 12(b)는 인쇄체 및 필기체 혼합문서의 인식결과이다. 먼저 한글과 영문의 분리과정이 수행되는데 영문의 경우 모두 영문으로 바르게 판단되나 문서에서 '오'의 경우 한개의 연결요소로 구성되어 영문으로 분리된다. 그러나 영문인식을 실패하여 한글로 바르게 인식되었다. 본 연구에서 제안한 알고리즘은 필기체를 위한 것이다. DP 정합을 이용함으로써 문자의 빼침에도 관계없이 바르게 인식함을 보여준다. 인쇄체의 경우는 필기체와 달리 문자의 획이 두꺼우며 획의 끝부분에서 빼침 등이 많아 문자인식에 어려움이 있다. 인쇄체의 경우도 획의 연속된 접촉이 '객'자나 '관'자 등 여러군데에서 나타나지만 이것을 획의 연결 가능성을 고려한 탐색을 수행하여 바르게 인식하였다. "구입품목"에서 '입'자의 경우 'ㅂ'자의 빼침이 작



(a)



(b)

그림 12. 인쇄체 및 필기체 포함한 혼합문서 3

- (a) 이진화의 개별문자
- (b) 인식결과

Fig. 12. Mixed document 3 containing printed/handwritten characters.

- (a) Binarization and character isolation,
- (b) Recognition result.

표 2. 수직모음에 따른 가능한 수평모음

Table 2. Possible horizontal vowels corresponding to each vertical vowel.

수직 모음	가능한 수평모음
ㅏ	ㅑ
ㅓ	ㅕ
ㅗ, ㅛ, ㅜ, ㅠ	ㅛ, ㅠ, ㅡ

아 'ㅏ'자로 인식되나 이것은 후처리과정을 통하여 빠르게 인식하였다. 또한 'ㅓ'자와 같이 'ㅇ'이 들어갈 경우 'ㅏ'자로 오인식하는 경우가 많은데 이러한 것도 후처리 과정을 통하여 빠르게 인식되었다.

V. 결 론

본 논문은 그래픽을 포함한 필기체 혼합문서의 인식에 관한 연구이다. 문서 인식을 위하여 먼저 제안한 이진화 알고리즘을 이용하여 이진화를 수행하였으며 한글의 유형과 크기 및 수직모음의 부분적 인식을 통하여 개별문자를 추출하였다. 이러한 방법을 사용함으로써 기존의 개별문자 추출이 영문과 한글을 모두 처리할 경우 생기는 문제점을 해결하였다. 또한 문자인식을 위하여 DP정합 비용함수값에 따른 branch and bound 알고리즘을 이용하여 자소의 접촉이 복잡하고 세선화 등으로 생기는 잡음에도 강한 인식 알고리즘을 개발하고 실험하였다. 또한 자소의 접촉 가능성을 고려함으로써 자소의 접촉이 복잡한 문자도 빠르게 인식하였다. 또한 간단한 후처리 과정을 통하여 오류를 정정하였다.

문자인식에 있어서의 중요한 점은 문자의 파손이나 문자의 심한 변형에 대처하는 것이다. 이러한 문제를 후처리과정에서 어느 정도는 보완할 수 있으나, 한글의 구조와 활용의 복잡성으로 문맥정보를 이용하는 방법을 개발하는데 어려움이 많아 앞으로 이에 대한 연구가 필요하다.

참 고 문 헌

1. 도정인, "한글문서인식 시스템의 개발," 정보과학회지, 제 9권, 제 1호, pp. 22-32, 1991년 2월.
2. D. Marr and E. Hildreth, "The detection of intensity changes by computer and biological

- vision system," *Computer Vision, Graphics, Image Processing*, vol. 22, no. 1, pp. 1-27, Apr. 1983.
3. L. Lam, S. W. Lee, and C. Y. Suen, "Thining methodologies-A comprehensive survey," *IEEE Trans. Anal. Mach. Intell.*, vol. PAMI-14, no. 9, pp. 869-885, Sep. 1992.
4. N. J. Naccache and R. Shinghal, "SPTA: A proposed algorithm for thinning binary patterns," *IEEE Trans. Syst. Man, Cybern.*, vol. SMC-19, no. 2, pp. 409-418, May/June 1984.
5. K. Wall and P. E. Danielsson, "A fast sequential method for polygonal approximation of digitized curves," *Computer Vision, Graphics, Image Processing*, vol. 28, no. 2, pp. 220-227, Nov. 1984.
6. H. Sakoe and S. Chiba, "Dynamic programming algorithm optimization for spoken word recognition," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-26, no. 1, pp. 43-49, Feb. 1978.
7. C.-K. Lin, K.-C. Fan, and F. T.-P. Lee, "On-line recognition by deviation-expansion model and dynamic programming matching," *Pattern Recognition*, vol. 26, no. 2, pp. 259-268, Feb. 1993.
8. 심동규, 함영국, 박래홍, "DP 매칭과 퍼지이론을 이용한 홀림체 온라인 한글인식," 전자공학회논문지, 제 30권 B편, 제 4호, pp. 116-129, 1993년 4월.
9. 류승필, 김태균, "속성분법을 이용한 필기체 한글 문서 내의 자모인식," 전자공학회논문지, 제 26권, 제 3호, pp. 85-94, 1989년 3월.
10. R. Schalkoff, *Pattern Recognition*, John Wiley & Sons, Inc., 1992.
11. 이주근, 남궁재찬, 김영진, "한글 pattern에서 sub-pattern분리와 인식에 관한 연구," 전자공학회논문지, 제 18권, 제 3호, pp. 1-9, 1981년 6월.
12. 함영국, 도상윤, 정홍규, 김우성, 박래홍, 이창범, 김상중, "문서 입력력 시스템 구성에 관한 연구," 전자공학회논문지, 제 29권 B편, 제 10호, pp. 100-112, 1992년 10월.
13. 민병우, 이성환, 김홍기, "문자인식을 위한 후처리 기법의 사례 연구," 제 1회 문자 인식워크샵 발표 논문집, pp. 91-103, 1993년 5월.



**沈 東 圭(Dong-Gyu Sim) 正會員**  
 1970년 7월 10일생  
 1993년 2월 : 서강대학교 전자공학과 졸업(공학사)  
 1993년 3월 ~ 현재 : 서강대학교 대학원 전자공학과 석사과정 재학중  
 ※주관심분야 : 영상처리



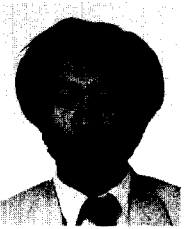
**金 仁 權(In Kwon Kim) 正會員**  
 1970년 2월 16일생  
 1992년 2월 : 서강대학교 물리학과 졸업(이학사)  
 1994년 2월 : 서강대학교 대학원 전자공학과 졸업(공학석사)  
 1994년 3월 ~ 현재 : 서강대학교 대학원 전자공학과 박사과정 재학중  
 ※주관심분야 : 영상처리



**咸 永 國(Young Kug Ham) 正會員**  
 1966년 7월 6일생  
 1990년 2월 : 서강대학교 전자공학과 졸업(공학사)  
 1992년 2월 : 서강대학교 대학원 전자공학과 졸업(공학석사)  
 1992년 3월 ~ 현재 : 서강대학교 대학원 전자공학과 박사과정 재학중  
 ※주관심분야 : 영상처리



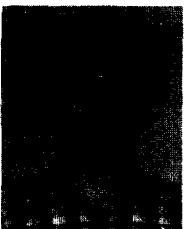
**朴 來 弘(Rae-Hong Park) 正會員**  
 1954년 1월 1일생  
 1976년 2월 : 서울대학교 전자공학과(공학사)  
 1979년 8월 : 서울대학교 전자공학과(공학석사)  
 1981년 6월 : Stanford Univ. 전기공학과(공학석사)  
 1984년 6월 : Stanford Univ. 전기공학과(공학박사)  
 1990년 1월 ~ 1990년 1월 : Univ. of Maryland, College Park, Center for Automation Research, Computer Vision Lab. 객원부교수  
 1984년 9월 ~ 현재 : 서강대학교 전자공학과 교수  
 ※주관심분야 : 영상통신, 컴퓨터 비전, 패턴인식



**李 昌 範(Chang Bum Lee) 正會員**  
 1955년 3월 23일생  
 1979년 2월 : 서강대학교 전자공학과 졸업(공학사)  
 1990년 8월 : 서강대학교 대학원 전자공학과 졸업(공학석사)  
 1983년 ~ 현재 : 한국전자통신 연구소 선임연구원  
 ※주관심분야 : 영상처리 통신시스템



**金 庠 仲(Sang Joong Kim) 正會員**  
 1949년 5월 15일생  
 1977년 2월 : 한양대학교 전자공학과 졸업(공학사)  
 1980년 2월 : 연세대학교 대학원 전자공학과 졸업(공학석사)  
 1977년 ~ 현재 : 한국전자통신연구소 통신접속연구실장  
 ※주관심분야 : 통신시스템 및 영상정보처리 등



**尹 炳 楠(Byeong Nam Yoon) 正會員**  
 1949년 11월 15일생  
 1975년 2월 : 한양대학교 전자공학과 졸업(공학사)  
 1989년 8월 : 청주대 대학원 전자공학과 졸업(공학석사)  
 1982년 ~ 현재 : 한국전자통신연구소 부장  
 ※주관심분야 : 통신시스템 및 영상정보처리 등