

Unicode 기반 다국어 명함인식기 개발

장 동 협[†] · 이 재 홍^{††}

요 약

명함을 이용한 전세계적인 고객 관리 시스템을 구축하기 위해 다국어 명함인식기를 개발하였다. 먼저 다양한 언어의 문자인식 및 학습을 위해 Unicode 기반 문자 이미지 DB를 구축하였으며, 다양한 입력 장치를 통해 획득한 명함 영상에 대하여 정확한 데이터를 얻기 위한 다양한 컬러영상 처리 기술이 적용되었다. 다음에 다층 퍼셉트론 신경망, 언어 유형별 개별 문자인식, 각 언어별 명함에 사용된 필드별 키워드 DB를 이용한 후처리를 적용하여 명함 인식률을 향상시켰다.

키워드 : Unicode, 명함, 다국어, 신경망, 문자인식, 후처리

A Development of Unicode-based Multi-lingual Namecard Recognizer

Dong-Hyeub Jang[†] · Jae-Hong Lee^{††}

ABSTRACT

We developed a multi-lingual namecard recognizer for building up a global client management systems. At first, we created the Unicode-based character image database for character recognition and learning of multi languages, and applied many color image processing techniques to get more correct data for namecard images which were acquired by various input devices. And by applying multi-layer perceptron neural network, individual character recognition applied for language types, and post-processing utilizing keyword databases made for individual languages, we increased a recognition rate for multi-lingual namecards.

Keywords : Unicode, Namecard, Multi-Lingual, Neural Network, Character Recognition, Post-Processing

1. 서 론

현재 기업들은 인적채널, 텔레마케팅, DM, 인터넷 등의 커뮤니케이션 수단으로 고객정보를 수집 관리함으로써 고객정보의 한계적인 수밖에 없다. 반면, 명함을 통해 얻어지는 정보들은 고객과 기업 간의 신뢰를 바탕으로 이루어지므로 그 정보가 정확성을 가지고 있지만, 대부분 기업들이 명함정보를 제 가치대로 활용하지 못하고 있는 것이 현실이다. 또한 기업들이 구축하고 있는 고객관리 시스템은 각 국가 간 언어의 상이함으로 데이터 활용에 제한적이었다.

문자인식기술을 기반으로 한 명함인식은 국내에는 하이네임, 이르미, 국외제품으로는 World Card, BizReader 등이 상업용으로 출시되어 있다. 이들 중 World Card가 19개 언어를 지원하고 있으며, BizReader는 인식대상 언어를 선택하도록 하고 있는 실정이다. 따라서 본 연구에서는 Unicode를 채택하여 다국어 학습 및 인식을 지원 하며, 신경망을 이용하여 다양한 언어를 유형별로 분류하여 인식률과 인식속도를 향상시킬 수 있는 순수 국내기술의 명함인식기를 개

발하여 전세계적 고객관리 시스템 구축에 활용하고자 한다.

다국어 명함 인식 과정은 1) 개별 문자 학습과 인식을 위해 Unicode 문자 이미지 데이터베이스를 구축하고, 2) 스캐너, 웹-카메라 등의 다양한 입력장치로부터 획득한 명함 영상에 다양한 칼라 영상 처리 기술을 적용하여 정확한 명함 정보 영역을 추출하고, 3) 유형별 문자인식, 다양한 특징 추출을 기반으로 신경망을 이용한 다국어 문자인식 엔진을 거쳐 개별 문자로 인식된 후, 4) 각 언어별 키워드 DB 구축을 통한 지식기반 후처리로 구성된다. (그림 1)은 본 연구에서



(그림 1) 다국어 명함인식기의 구성도

† 정 회 원 : (주)한국인식기술 부설 연구소 소장
 †† 정 회 원 : 전남도립대학 보건의료과 부교수
 논문접수: 2008년 9월 17일
 수정일: 1차 2008년 10월 24일
 심사완료: 2008년 11월 1일

개발된 다국어 명함인식기의 구성도를 보여준다.

본 논문의 구성은 2장에서 다국어 명함인식을 위한 과정들을 기술하고, 3장에서 구현된 명함인식기에 대한 실험 및 평가를 한 후, 4장에서 결론을 맺는다.

2. 다국어 명함인식기

2.1 Unicode 문자 이미지 DB 구축

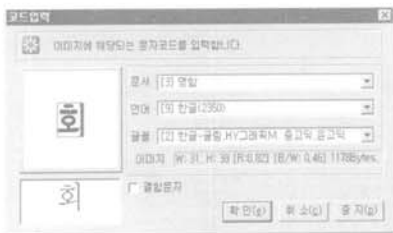
Unicode는 국제표준으로 제정된 2바이트계의 만국 공통의 국제 문자부호 체계(UCS: Universal Code System)이다 [1,2]. 따라서 개발하려는 다국어 명함인식기는 다양한 언어를 지원할 수 있도록 Unicode를 지원한다.

문자인식 및 학습단계에서는 개별문자 이미지에 대한 구조적, 통계적 정보를 추출하여 인식에 필요한 각종 특징정보를 추출한다. 다양한 언어의 문자나 글꼴(font)을 인식하기 위한 인식시스템을 개발하기 위해서는 개별문자 이미지 DB(UCDB: Unicode Character DB)구축이 필요하다. UCDB는 Unicode 기반 다국어 문자표를 언어·글꼴·크기·강조효과(기울임, 진하게)등을 구분하여 다양하게 입력한다.

다국어 문자 이미지를 입력하기 위해서는 언어·글꼴·문서유형(일반서적, 신문, 명함 등)을 지정했다. (그림 2)는 명함 이미지에서 개별문자를 UCDB에 추가하기 위한 화면을 보여준다.

<표 1> UCDB 테이블

필드	속성	길이	필드 설명
IDX	int	-	기본키
CODE	varchar	8글자	이미지에 해당되는 문자코드 값
SEPT	int	-	결합문자 표시 필드 (개별문자:0, 결합문자:1)
LANG	int	-	언어표시
FONT	int	-	글꼴표시
CX	int	-	이미지 폭 (pixel)
CY	int	-	이미지 높이 (pixel)
ISIZE	int	-	이미지 버퍼 크기 (bytes)
IMAGE	image	128K	이미지 버퍼



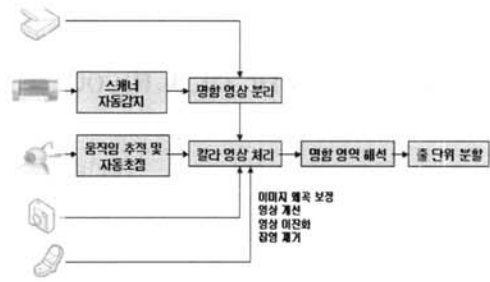
(그림 2) 개별 문자 이미지를 UCDB에 등록하는 화면

2.2 정보입력 모듈

다국어 명함인식 시스템에서는 언제 어디서나 정보를 입력할 수 있도록 다양한 디지털 기기로부터 영상을 획득할 수 있어야 한다. 대부분의 영상 입력 수단으로 사용되고 있

는 카메라는 명함정보 획득의 편리성이 있어 많이 사용되고 있다. 그러나 카메라로부터 입력받은 영상의 경우 조망의 상태나 방향, 주변 환경의 변화로 인하여 스캐너를 통해 입력된 영상에 비해서 영상의 질적 변화나 왜곡이 발생되기 쉬우므로 스캐너를 이용하여 입력한 영상에 비해 인식하기가 어려운 것으로 알려져 있다. 그러므로 디지털 기기로 입력받은 영상으로부터 보다 정확한 명함 정보를 추출할 수 있도록 칼라 영상처리에 대한 기술 개발이 필요하다.

정보입력모듈에서는 다양한 정보입력 기기 제어, 칼라 영상 처리, 명함 구조해석을 통해 명함인식 엔진 개발의 전처리 과정을 수행할 수 있도록 하였다. 정보입력 모듈에서 개발한 모듈은 (그림 3)과 같다.



(그림 3) 정보입력 모듈 구성도

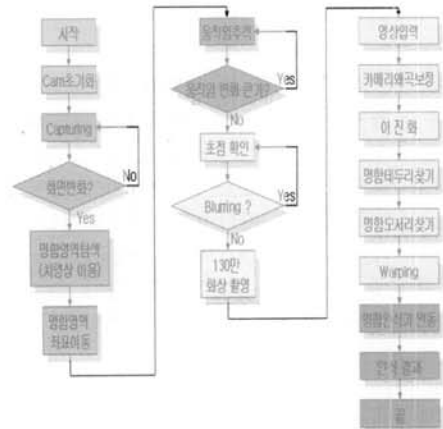
2.2.1 명함영상 획득

정보입력 모듈에서는 웹-카메라를 통해 명함을 인식하고자 할 때, 명함을 사용자가 카메라에 정 위치 시키지 못하는 점을 보완해야 한다.

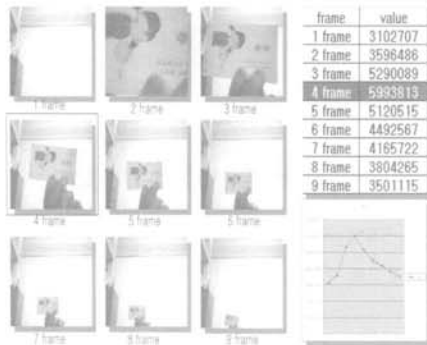
일정 범위 내에 명함을 위치시키면 명함 전체 영상을 얻기 위해 카메라의 움직임을 제어한다. 웹-카메라는 30만 화소 Pan/Tilt를 지원해야 한다. 명함의 움직임을 탐지하기 위해서 배경 프레임과 현재 프레임과의 차영상을 이용하여 명함 영역을 탐색한다.

(그림 4)는 웹-카메라를 통한 명함 영상 획득 과정을 보여준다.

각 영상으로부터 특징을 구하고 그 특징을 기준으로 시간에



(그림 4) 웹-카메라를 통한 명함영상 획득 과정



(그림 5) 최적의 정지 영상 획득 과정

따라 추적하게 되는데, 유용한 특징이 되는 것이 영상에서 꼭지점, 모서리 영역들이 해당된다. 이 특징 점을 기준으로 인접해 있는 픽셀들의 밝기 정보 차이를 통해 움직임을 추적하게 된다. 이 움직임이 거의 없을 때 명함이 카메라의 영상 획득 범위에 들어온 것으로 판단하고 자동으로 촬영한다.

(그림 5)는 웹 카메라에서 가장 최적의 영상을 획득하게 되는 과정을 보여 주고 있다.

2.2.2 칼라영상 처리

디지털 카메라로 찍은 영상은 스캐너와는 달리 렌즈의 특성에 따라 영상의 가장자리 부분이 휘어져 보이는 왜곡현상이 발생하게 된다. 이로 인해 명함인식 시에 정확한 정보를 추출하기 어렵다. 영상의 워핑(warping)을 이용하여 명함인식에 적합한 형태로 보정한다[3,4,5]. 카메라 렌즈에 의한 왜곡 보정은 원 영상에서 렌즈 왜곡 변수 값을 계산하여 inverse spatial mapping 방식을 사용하여 이미지 워핑을 하였다. 모양에 따른 왜곡 보정은 이진화-명함 외곽선 분리-명함의 4-모서리와 꼭지점 찾기-사변형에서 사각형으로의 양선형 변형(bilinear transformation) 과정을 거친다.

다양한 명함을 처리해야 하는 명함인식기에서 배경영상이 있는 명함이나, 역상으로 처리된 명함 등과 같은 경우 문자 영역과 배경 영역이 명도 레벨에서 확연히 구분되지 않는 경우가 존재하여 명함의 정보를 추출하기가 쉽지 않다. 따라서 명도 영상을 대상으로 영상개선을 수행하여 명도 대비가 더욱 분명하도록 하여 보다 정확한 정보를 추출할 수 있도록 명함영상을 획득하도록 한다. 명함 영상에는 잡영도 존재하고 명함을 스캔하면서 나타는 글자의 흐려짐을 개선하기 위해 필터링을 사용하여 보다 정확한 정보를 획득하도록 하였다.

명함 이미지는 다양한 형태의 디자인과 색깔 분포를 가지고 있어 이진화에 제약 사항이 많으며 양질의 글자 이미지를 획득하기가 힘든 경우가 많다. 명함이미지의 경우 전체 반전과 부분반전 등을 포함하고 있어서 반전된 영상의 부분은 정보를 추출하지 못하는 경우가 발생하므로 이를 위한 영상개선이 필요하다[6]. (그림 6)에서 필터링을 적용하였을 때 바탕색이 있는 영역의 정보를 획득할 수 있음을 볼 수 있다.

명함 영상의 해석과 인식이 고속으로 정확하게 처리되기 위해서 컬러 영상의 이진화가 필수적으로 요구된다. 특히



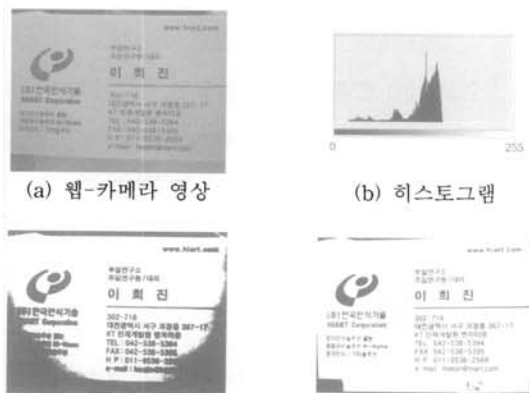
(그림 6) 필터링을 이용한 영상 개선

카메라로 입력받은 명함영상의 경우 조도와 빛의 방향에 따라 문자 부분이 훼손되어 나타나는 경우가 많다. 그러므로 빛의 영향에 둔감한 영상 이진화 알고리즘 및 이진화 과정이 고속으로 이루어져야 될 필요가 있다[6,7].

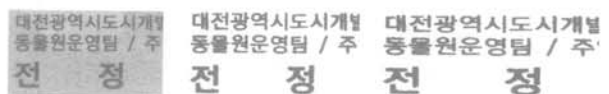
본 연구에서는 전역적 이진화 방법과 각 입력 영상에 따른 최적의 임계값을 결정하는 방법을 혼합한 방법을 사용하고 있다. 임계값은 영상의 그레이 레벨의 평균을 반복 알고리즘을 이용하여 최적의 임계값을 결정하였다. 반복 알고리즘의 기본적인 연산식은 아래와 같이 배경 T_b 와 물체 T_c 의 각 그레이 레벨의 평균을 나타낸다.

$$T = \frac{(T_0 + T_b)}{2}$$

위 식의 연산에 의해 한계값 T 를 구하면, T 를 이용하여 T_c 와 T_b 의 값을 다시 설정한다. 전체적인 처리 흐름은 동일한 T 가 나올 때까지 T_c 와 T_b 의 값을 다시 설정하게 되며, 결국 T 는 영상에서 적절한 한계치로 나타낼 수 있다. 본 연구에서는 전역적 이진화 방법을 사용함으로써 지역적 이진화 방법보다 더 빠르게 처리할 수 있었으며, 영상에 따라 최적의 임계값을 추출하므로 효율적인 영상의 이진화가 가능하였다.



(a) 웹-카메라 영상 (b) 히스토그램
(c) 히스토그램의 평균값을 이용한 이진화 영상 (d) 반복알고리즘을 이용한 이진화 영상
(그림 7) 이진화 과정



(a) 원영상 (b) 잡영 영상 (c) 잡영 제거된 영상
(그림 8) 잡영 제거 과정

2.2.3 명함 영역 해석

명함 영상 구조 해석은 획득된 명함 영상 데이터의 구성을 분석하여 그림·사진 등이 포함되어 있는 인식과 관계없는 영역과, 키워드·본문 등의 문자열이 존재하는 인식해야할 영역으로 영역 사이의 결합 구조를 밝혀내는 과정이다[8,9].

본 연구에서는 우선 명함의 일반적인 유형들을 분류한 후, 명함의 글자 영역들을 추출하였다.

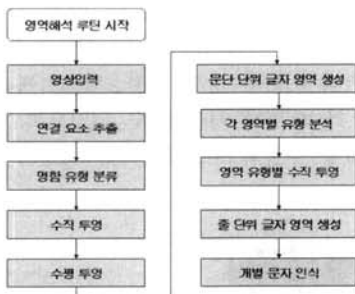
(그림 9)는 분류한 명함 유형을 보여준다.

영역해석 및 라인단위 분리 처리 과정은 (그림 10)의 순서로 개발되었다. 처리 과정은 우선 줄 단위로 추출하고 각 줄에서 개별 글자들을 추출하였다. (그림 11)은 줄 단위 개별 문자 영역을 추출한 결과를 보여주고 있다.

(그림 11)과 같이 각 언어별에 맞게 분리된 글자 영역을 가지고 분리 병합을 거치며 최적의 인식결과를 검색한다.



(그림 9) 명함 유형



(그림 10) 영역 해석 및 줄 단위 분리 과정



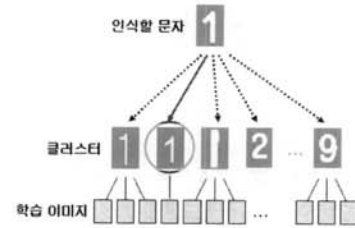
(그림 11) 줄 단위 글자영역 추출 결과

2.3 다국어 문자인식 모듈

다국어 명함인식 모듈에서는 다양한 특징추출을 기반으로 신경망을 이용한 학습 알고리즘을 이용하였다.

2.3.1 다층 퍼셉트론

명함인식엔진에서는 역전파(backpropagation) 학습 알고리즘을 이용한 다층 퍼셉트론(multi-layer perceptron)을 이



(그림 12) 클러스터링을 이용한 학습

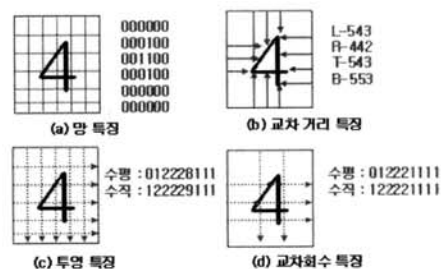
용하였다. 일반적인 다층 퍼셉트론 신경망을 이용하였으며, 입력층은 1개, 은닉층은 1개, 출력층은 1개로 구성하였다.

문자 인식을 위한 학습 단계에서는 다양한 형태의 입력 문자 패턴들을 클러스터링을 이용하여 유형별로 분류하여 학습시켰다. 본 연구에서는 비계층적 클러스터링 방법 중 문서 클러스터링에 많이 적용되는 k-means 클러스터링을 적용했다[10,11]. 클러스터링을 이용함으로써 인식대상 문자가 속한 언어권과 관계없이 문자인식을 용이하게 하도록 하였다.

2.3.2 특징 추출

신경망을 이용한 특징 추출은 전처리 과정을 거친 영상을 다양한 특징을 추출하여 신경망의 입력에 사용하는 방법으로 인식을 향상에 많은 도움을 준다. 특징을 추출하기 위해서는 문자영상에 대한 정규화 과정이 필요하다. 다양한 글꼴과 다 크기의 문자들을 인식하기 위해서 학습시킨 표본 문자와 비교하기 위해서는 입력 문자 영역을 학습시킨 표본 문자의 크기로 정규화하였다[6,12]. 정규화 된 문자 패턴에 대해 다른 문자 패턴과 구분되는 주요 특징을 추출하여 조합함으로써 효율적인 인식 성능을 구현하기 위한 입력 벡터를 생성하였다. 본 연구에서는 특징 추출을 위해 기존의 망(mesh)과 교차 거리(cross distance), 교차 회수(cross number) 등의 방법을 사용하였다.

입력 영상은 32×32로 정규화하여 8×8의 부 영역으로 나누고, 각 부 영역의 검은 화소 값을 기준으로 0.0~1.0 사이의 값으로 망 특징 값이 추출된다. 또한, 문자 영상의 주어진 경계로부터 좌우로 주사선을 주사하였을 때 획과 만나는 거리를 특징으로 이용했다. 투영 특징은 분할 결과를 정규화시킨 문자 영역에서 수평 방향 및 수직방향의 흑화소들에 대한 누적합 분포를 계산한 것이다. 교차 회수 특징은 문자 영역의 수평 방향 및 수직 방향으로 주사선을 통과시켰을



(그림 13) 격자 및 거리 특징 추출

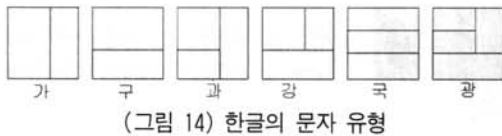
때 문자 영역과 주사선이 만나는 회수를 구한 것이다.

2.3.3 유형별 문자인식

인식기를 구성하여 문자를 인식하는 방법은 하나의 인식기를 사용하여 전체 인식 대상 문자들 가운데서 한 문자를 인식하는 방법과 문자들의 구조적 특징에 따라 분류된 각 유형별로 인식기를 구성하여 문자를 인식하는 방법으로 나눌 수 있다. 유형별로 문자를 분류하여 각 유형별로 인식하는 방법은 인식에 앞서 인식 대상 문자들의 특징에 따라 유형을 분류하고 각 유형별로 다시 세부 분류하여 인식함으로써 분류기별 부담도 적고 비슷한 문자들이 많은 한글과 같은 문자 인식에 매우 효과적이다.

예를 들어 한글의 경우, 두 개 또는 그 이상의 자음과 모음들이 수직 또는 수평 방향의 이차원 형태로 결합되어 있다. 이러한 한글의 구조적 특징을 고려하여 (그림 14)와 같이 여섯 개의 문자 유형으로 분류하였다.

명함인식엔진은 유형분류와 문자인식의 2단계 신경망으로 구성하였다. 인식대상 7318자(한글 2350자, 한자 4888자, 알파벳 52자, 기호·숫자 32자)에 대해 28개의 유형으로 분류한 후, 다시 각 유형에 대해 80~700개의 출력을 갖도록 신경망을 구성했다.



(그림 14) 한글의 문자 유형

2.3.4 문자인식

(그림 11)의 각 언어별에 맞게 분리된 글자 영역을 가지고 분리 병합을 거친 후 명함인식엔진에 의해 최종 인식된 결과는 (그림 15)와 같다.

글자 단위로 인식된 결과는 원영상 아래에 표시하였다.



(그림 15) 최종 인식된 결과

2.4 후처리 모듈

다국어 문자인식 모듈에서 나온 인식결과를 가지고 각 필드 별로 분류함으로써 정보가 되고 이를 이용함으로써 더욱 더 가치 있는 정보가 된다[13]. 인식 모듈에서 나온 인식결과로부터 필드를 추출하기 위한 데이터를 생성한 후, 각 필드 별 키워드 DB와 각 필드 별 형태적 특징들을 이용하여 의미 있는 정보들 즉, 각 나라별 이름, 회사, 직위, 부서, 전

화, 팩스, 휴대폰, 홈페이지, e-메일, 주소, 기타 항목(이차 필드)으로 분류한다. 또한 분류 결과를 보정하거나 해당 형식으로 맞추어 보기에 용이하게 해준다.

중국어(간체)·일본어 후처리 키워드 DB와 중국의 명함의 특징을 찾아 각 각의 필드로 보정하고 분류하였으며, 라틴계 언어에 대해서도 후처리 키워드 DB와 각 나라의 명함의 공통적인 특징을 찾아 각 각의 필드로 보정하고 분류하였다.

3. 실험 및 고찰

본 연구에서 개발된 다국어 명함인식기의 성능을 <표 2>에 보였다. 성능 시험에는 인텔 펜티엄IV 2.4GHz, 2GB RAM, 80GB HDD의 PC를 사용했다. <표 2>에서 인식속도는 명함 1장을 인식하는 데 소요되는 시간이며, 개발된 명함인식기는 한글·일본어·중국어·라틴계 언어를 포함하여 17개 언어를 인식한다. 명함인식 속도는 3~6초이었으며, 순수한 문자인식 속도는 초당 160자 이내이었다. 일본어 명함의 인식률이 다른 언어에 비해 다소 낮은 이유는 일본어가 문자 구조상 탁음(위첨자)이 존재하고, 문자열의 띄어쓰기가 없으므로 개별문자를 분리하는데 어려움이 있었기 때문이다.

<표 2> 개발된 다국어 명함인식기의 성능

	인식속도 (초)	인식률 (%)	용량 (MByte)	실험대상 명함(장)
한글	3~4	98.2	10	500
라틴계	2	98.5	1	100
중국어	5~6	96.5	15	200
일본어	3~4	96	8	302



(a) 중국어 명함



(b) 독일어 명함

(그림 16) 후처리를 거친 다국어 명함인식 결과

4. 결 론

본 연구에서는 다국어 명함인식기를 개발하기 위해 다양한 시도를 하였다. 개별문자 인식 및 학습을 위해 Unicode 문자 이미지 DB를 구축하였고, 명함 인식에 가장 적합한 영상을 획득하기 위해 스캐너와 웹-카메라로부터 획득한 명함 영상에 다양한 컬러 영상 처리 기술을 적용하였으며, 명함 영상 구조 해석을 통해 명함에서 글자영역과 그림 영역을 구별하여 명함인식률을 향상시키고자 하였다. 줄 단위로 분리된 영역에서 개별 문자를 분리한 후 다양한 특징점을 추출하고 신경망을 이용하여 개별 문자를 인식하였다. 마지막으로 인식된 개별 문자들에 각 국가별로 구축된 지식 기반 DB를 적용한 후처리를 통해 개별 문자로부터 유용한 정보를 획득할 수 있도록 하였다.

다국어 명함 인식기의 개발로 이를 활용한 다양한 제품 개발이 가능하게 되어 국내 시장에만 국한되어있던 문자인식기술 시장에 활력을 불러일으킬 것으로 기대된다. 또한 다국어 지원으로 여러 나라의 명함을 동시에 등록, 검색이 가능하게 됨에 따라서 전 세계적 고객관리가 가능하여 기업에 구축되어 있는 CRM, ERP 시스템에 적극 활용 될 수 있을 것이다.

참 고 문 헌

[1] ISO/IEC 10646-1, Universal Multiple-Octet Coded Character Set, 1995.
 [2] 전상훈, "C로 구현한 한글 코드 마스터", 도서출판골드, 2005.
 [3] Rafael C. Conzalez, Richard E. Woods, "Digital Image Processing", Addison Wesley, pp.443-458, 1992.
 [4] 손화정 김수형, "문자인식 분야의 기계학습 기법", 정보과학회지, 25권 3호, pp.12-20, 2007.
 [5] 유원필, 정연구, "내장형 렌즈 왜곡 보정 알고리즘 구현을 위한 이미지 워핑 방법", 한국정보처리학회논문지B, 10B권 4호, pp.373-380, 2003.
 [6] 장동혁, "디지털 영상처리의 구현", 정보게이트, 2001.
 [7] 김계경 김진호, "모바일 폰 카메라로 획득한 문서영상에 대한 국소 적응적 이진화 알고리즘", 한국화상학회지, 10권 1호, pp.17-26, 2004.
 [8] 이인동, 김태권, 권오석, "문서 인식을 위한 전처리 기술의 소개", 한국정보과학회지, 19권, 1호, pp.14-21, 1991.
 [9] 이인동, "문서영상에서 문자와 비문자의 분리 추출방법", 한국정보과학회 논문지, 17권, 3호, pp.247-258, 1990.
 [10] Tapas Kanung, "The Analysis of a Simple k-Means Clustering Algorithm," Proc. of ACM Symposium on Computational Geometry, June, 2000.
 [11] 오형진, "클러스터 중심 결정 방법을 개선한 K-Mean Algorithm의 구현", 전북대학교 대학원 석사학위논문, Aug., 2002.
 [12] 최형일, 이근수, 이양원 공역, "영상처리 이론과 실제", 홍릉과학출판사, pp.256-261, 1997.

[13] 민병우, 이성환, 김흥기, "문자 인식을 위한 후처리 기법의 사례 연구", 제1회 문자 인식 워크샵 발표 논문집, pp.91-104, 1992.



장 동 협

e-mail : dhjang@hiart.com

1994년 한밭대학교 전자계산학과(학사)

1996년 한밭대학교 전자계산과(공학석사)

2002년~현 재 (주)한국인식기술 부설 연구소 소장

관심분야 : 문자인식, 음성인식, HCI, 네트웍 망관리 등



이 재 홍

e-mail : jhlee@dorip.ac.kr

1986년 충남대학교 전자공학과(학사)

1988년 충남대학교 전자공학과(공학석사)

1999년 충남대학교 컴퓨터공학과(공학박사)

1988년~1994년 국방과학연구소 연구원

1994년~1995년, 1999년 (주)한국인식기술 연구원

2000년~현 재 전남도립대학 보건의료과 부교수

관심분야 : 문자인식, 멀티미디어, 센서네트워크