

## 연결성분 자소를 이용한 문자 인식 연구

이경호<sup>o</sup>

<sup>o</sup>한라대학교 정보통신방송공학부

e-mail:khlee@halla.ac.kr<sup>o</sup>

## A Study on Character Recognition using Connected Components Grapheme

Kyong-Ho Lee<sup>o</sup>

OSchool of Information & Communication, Broadcasting Engineering, Halla University

### ● 요약 ●

본 연구에서는 한글 문자 인식을 수행하였다. 한글 인식을 수행하되 고딕 인쇄체 문자를 대상으로 하였고, 자소 단위 인식을 통한 인식을 수행하되 기존 한글 문자 인식 연구에서 사용하는 자음과 모음 단위의 자소가 아닌 연결성분을 이용하여 인식하는 새로운 자소를 이용하였다. 새로운 자소들은 끝점, 2선 모임점, 3선 모임점, 4선 모임점의 특징을 추출하고 특징에 의해 자소를 인식하는 데이터베이스를 구성하여 자소를 인식하게 하였다. 또한 연결 성분을 반영한 새로운 자소로 고딕 인쇄체 문자를 인식하므로 추출된 자소를 6가지로 분류하였고, 6가지 자소에 의해 구성되는 92가지 문자 구조를 제안하고 이에 따른 문자를 데이터베이스를 구축하였고, 자소의 무게 중심을 이용한 분포를 이용하여 제안된 구조를 통하여 데이터베이스를 이용한 문자인식을 수행하였다.

**키워드:** 문자 인식(character recognition), 자소 인식(grapheme recognition), 연결 성분(connected components)

### I. Introduction

본 연구에서는 한글 문자 인식을 수행하였다. 고딕 인쇄체 문자를 대상으로 하였고, 연결성분을 이용하여 인식하는 새로운 자소를 이용하였다. 새로운 자소들은 끝점, 2선 모임점, 3선 모임점, 4선 모임점의 특징을 추출하고 특징에 의해 자소를 인식하는 DB를 구성하여 자소를 인식하였다. 또한 연결 성분을 반영한 새로운 자소로 고딕 인쇄체 문자를 인식하므로 추출된 자소를 6가지로 분류하였고, 6가지 자소에 의해 구성되는 92가지 문자 구조를 제안하고 이에 따른 문자를 DB를 구축하였고, 자소의 무게 중심을 이용한 분포를 이용하여 제안된 구조를 통하여 DB를 이용한 문자인식을 수행하였다.

### II. Preliminaries

문자 인식은 많이 연구되어 온 분야이나, 지금도 여전히 많은 연구들이 나오고 있다. 문자열에 회전 알고리즘을 적용하여 인식을 향상시키거나, 휴리스틱 분할 알고리즘을 적용하여 향상을 꾀하고, 구글의 Open API를 이용하여 한글 문자 인식 향상을 위한 노력을 하며, 한글 문자 주변에 삼벌이 위치한 상황에서 문자 템플릿에 의존하지 않는 블립 투사를 제시하며, 한글 문자의 다양성 확대를 위해 다중 컬럼 딥 인공 신경망을 이용하여 인식을 시도하는 등 많은

연구들이 지속되고 있다. 연구가 지속되는 이유는 한글 문자의 특수성 때문으로 한글은 영문자 알파벳과 달리 초성과 중성 종성들이 2차원적 조합되어 다수의 문자가 구성되기 때문에 인식해야할 대상이 많으며, 작은 획 하나에 의해 구분될 정도로 유사성이 높아 처리 과정이 복잡하기 때문이다[1-7]. 그동안 연구된 한글 인식 방법들은 한글을 작성하는 방법에 따라 ‘인쇄체 인식’과 ‘필기체 인식’으로 분류할 수 있고 또 인식하는 단위에 따라 ‘문자 단위 인식’과 ‘자소 단위 인식’으로 분류할 수 있다. 두 가지 다 전처리 작업으로 입력된 영상에서 보정을 하고, 텍스트 영역 추출과 수평 분할, 수직 분할 등의 반복을 통해 문자 단위로 분할 한 후, 문자 단위에서는 각각의 방법으로 인식을 한다. 또한 대부분의 경우 추가되는 전 처리 작업으로 추출해 낸 문자의 크기를 자신들이 정한 규격으로 조절한 후 인식을 한다. 문자 단위 인식의 경우 정해진 규격 크기를 웨이블릿 변환을 하여 인식하거나, 신경망을 통한 인식에서는 문자를 구성하는 화소를 입력으로 하여 인식을 한다[8-11]. 자소 단위 인식에서 자소의 분할 연산은 매우 중요한 과정이다. 자소 단위 인식에서 자소 분리를 하다가 자소의 일부분이 소실될 수도 있고, 자소 구성 영상 이외의 잡영이 포함될 수도 있다. 자소의 범위를 작게 잡으면 자소의 일부가 소실되고 소실을 막기 위하여 자소의 영역을 크게 잡으면 잡영이 들어오는데 소실과 잡영을 처리하는데 애를 먹기도 한다[7].

### III. The Proposed Scheme

[12]에서는 기존 연구가 자소 분할에 어려움을 겪는 것을 보고 자소 분리에 유리한 연결 성분을 이용한 새로운 자소를 제안하고 분류하였으며, 또 제안된 자소를 인식하기 위한 특징 점들을 제안하고, 이를 데이터베이스화 한 후, 제안된 특징 점들을 이용하여 자소들을 인식하는 실험을 수행하였다. 본 연구에서는 [12]에서 제안한 자소를 이용하여 문자를 인식하는 과정을 연구 수행하였다. 수행한 알고리즘은 1) 문자 영상 입력 2) 문자 범위 확정 3) 입력된 컬러 영상을 회색 레벨 영상으로 변경 4) 회색 레벨 영상을 2진 영상으로 변경 5) 2진 영상으로 세선화 수행 6) 세선화 영상에서 특징점으로 끝점, 2선 모임점, 3선 모임점, 4선 모임점 추출 7) 2진 영상으로 역세선화로 자소별 영역 분류 고립점 추출 8) 영역별 범위 내 특징 점과 자소 특징 DB를 참고하여 자소 인식 9) 자소별 무게 중심 추출 10) 탐침에 의한 문자 구획 분류 11) 문자 구획과 자소별 무게 중심을 이용한 문자 조합을 거친다.

본 연구에서 이용된 자소와 분류는 그림 1과 같고, 제안된 자소에 의해 구성할 수 있는 문자 구조는 그림 2과 같고, 문자 구조에 따른 소속 문자는 그림 3과 같다.

분류	소속 자소
C	ㄱ/ㄴ/ㄷ/ㄹ/ㅁ/ㅂ/ㅅ/ㅇ/ㅈ/ㅊ/ㅋ/ㅋ/ㅌ/ㅍ/ㅍ/ㅍ in 싸 /ㅍ(고립점유) in 짜.../ㅍ(고립점무) in 쯤...
V	ㅏ/ㅑ/ㅓ/ㅕ/ㅗ/ㅛ/ㅜ/ㅠ/ㅡ/ㅣ/ㅣ/ㅣ
Cr	겨 in 계.../꺠/꺡/꺢/꺣/꺤/꺥/꺦/꺧/꺨/꺩/꺪/꺫/꺬/꺭/꺮/꺯/꺰/꺱/꺲/꺳/꺴/꺵/꺶/꺷/꺸/꺹/꺺/꺻/꺼/꺽/꺾/꺿 in 꺾.../꺾/어/여/오/요/ 쿠/꺠/꺡/꺢/꺣/꺤/꺥/꺦/꺧/꺨/꺩/꺪/꺫/꺬/꺭/꺮/꺯/꺰/꺱/꺲/꺳/꺴/꺵/꺶/꺷/꺸/꺹/꺺/꺻/꺼/꺽/꺾/꺿
VC	ㅅ in 녹.../ㅅ in 불.../ㅅ in 불.../ㅅ in 불.../ㅅ in 불.../ㅅ in 불... ㅅ in 녹.../ㅅ in 불.../ㅅ in 불.../ㅅ in 불.../ㅅ in 불.../ㅅ in 불... ㅅ in 불.../ㅅ in 불...
O	lower part of 포.../lower part of 표...
VCC	ㅅ in 뭉.../ㅅ in 뭉.../ㅅ in 뭉.../ㅅ in 뭉...

Fig. 1. Grapheme classification

단	조각	패턴 (보기)	패턴구조 코드
1	1	Cr	010101
	2	C V Cr V C Cr	010201 010202 010203
	3	C C V C V V C Cr V Cr V V	010301 010302 010303 010304
	4	C C V V	010401
2	2	C C V Cr V	020201 020202 020203 020204 020205
		C C V C C VC C Cr	020301 020302 020303 020304 020305
		Cr V Cr V Cr C V V V	020306 020307 020308 020309 020310
		V V V V V V	
		C C V C C VC C Cr	

4	020311 020312 020313	C V C Cr Cr V C V V C C C C C C C	
	020401 020402 020403 020404	C Cr V C C V C C V Cr V V C C C C VC C C C	
	020405 020406 020407 020408	C C V C V V C V V V V V V V C V V C V V C V V	
5	020409 020410 020411 020412	C C V V C V V Cr V V C C C C C C C	
	020501 020502 020503	C C V C C V V V C C C V V	
	020504 020505	C C V V C C C C	
6	020601	C C V V C C C	
	3	030301 030302 030303 030304 030305	C Cr V V V V V C C C C C C VC VCC V 030306 030307 030308 030309 030310
		030401 030402 030403 030404 030405	C C C V C V C V C V V V Cr V V C C C C C C C C 030406 030407 030408 030409 030410
030411 030412 030413 030414 030415		V V V V V V V O V C Cr VC C VC C V V V V 030416 030417	
030501 030502 030503 030504		C V V C C V C C V V V V C C C C C 030505 030506 030507 030508	
030509 030510 030511		V V V C V V C V C V V C C C V V	
5	030601	C C V V V C C C	
	4	040401 040402 040403	V C V C V V V V V C C C
		040501 040502 040503	V V V C C V V V C C C



040601	공룡
050601	흥

Fig. 3. Classified korean character

#### IV. Conclusions

본 연구에서는 고딕 인쇄체 문자 자소 단위 인식을 통한 한글 문자 인식을 수행하되 연결성분을 이용하는 새로운 자소를 이용하였다. 추출된 자소를 6가지로 분류하였고, 6가지 자소에 의해 구성되는 92가지 문자 구조를 제안하고 자소의 무게 중심을 이용한 분포를 이용하여 제안된 구조를 통하여 문자인식을 수행하였다.

#### References

[1] Hyuna Oh, EuGene Rhee, "Enhancement of Car Licence Plate and Security with Rotation Algorithm", Journal of Security Engineering, Vol. 13, No. 2, pp. 83~90, Apr. 2016

[2] Moon Yong Jin, Jong Bin Park, Dong Suk Lee, "Real-Time Vehicle Licence Plate Recognition System Using Adaptive Heuristic Segmentation Algorithm", KIPS Tr. Software and Data Eng., Vol. 3, No. 9, pp. 361-368, Mar. 2014

[3] Kang-San Kim, Seok-Cheon Park, Seok-Ho Oh, "Suggestion of Enhanced Korean Character Recognition Technique Using Google Tesseract Open API", Proceeding of Korean Society for Internet Information, Vol. 16, No. 1, Spring. 2015.

[4] Kyusoo Choung, "Text Area Detection of Road Sign Images based on IRBP Method", Journal of Intelligent Transportation System, Vol. 13, No. 6, Dec. 2014.

[5] Kyung-Wha Park, Byoung-Hee Kim, Dong-Sig Han, Seong-Ho Son, Woo-Yung Kang, Byoung-Tak Zhang, "Handwritten Hangul Recognition using Multi-column Deep Neural Networks", Proceeding of KIPS Spring Conference, Vol. 26, No. 1, 2016.

[6] Min-Soo Kim, Eun-Young Kang, Woo-Sung Kim, Sun-Hwa Han, Jin-Hyung Kim, "A Study on Implementation of Printed Character Recognition System And Performance Evaluation", Korea Information Processing Society, Vol. 7, No. 11, pp. 3584-3591, Nov. 2000.

[7] Kil Taek Lim, Ho Yon Kim, "A Study on Machine Printed Character Recognition Based on Character Type Classification", The Institute of Electronics and Information Engineers, Vol. 40, No. 5, pp. 26-39, Sep. 2003.

[8] Kil-Taek Lim, Gi-Seok Kim, "Reestimation of Recognition Result of MLP Classifier for Machine Printed Hangul - Feasibility Study on Softmax Method", Journal of Information & Electronic Technology, Vol. 6, pp. 93-105, 2007.

[9] Kil Taek Lim, Ho Yon Kim, "A Study on Machine Printed

Character Recognition Based on Character Type Classification", Journal of IEIE, Vol. 40, No. 5, pp. 266-279, Sep. 2003.

[10] Kil Taek Lim, Seon Hwa Jeong, Seung Ick Jang, Ho Yon Kim, "An Implementation Method of the Character Recognizer for the Sorting Rate Improvement of an Automatic Postal Envelope Sorting Machine", Journal of Korea Society of Industrial Information System, Vol. 12, No. 4, Dec. 2007.

[11] Duk-Ryong Lee, Woo-Youn Kim, Il-Seok Oh, "A Hangul Document Image Retrieval System Using Rank-based Recognition", Journal of The Korea Contents Association, Vol. 5, 2005.

[12] Kyong-Ho Lee, "A Study on Grapheme and Grapheme Recognition Using Connected Components Grapheme for Machine-Printed Korean Character Recognition", Journal of The Korea Society of Computer and Information, Vol. 20, No. 1, pp. 27-36, September, 2016