

Korean Character Processing: Part I. Theoretical Foundation

(한글문자의 컴퓨터 처리 : I. 이론)

鄭 元 良^{*}
(Chung, Won Lyang)

要 約

2 부에 걸친 한글의 컴퓨터 처리에 관한 논문으로서, 제 1 부에서는 한글문자의 컴퓨터 처리의 제문제점을 확인하고 이를 위한 해결책으로 이론적 기반의 확립을 제안 하였다. 한글문자의 일차원적 문법과 이차원적 구조를 각각 BNF와 "patternal structure"를 사용하여 정의하고 이로부터 모아쓰기를 위한 lexical analysis 및 parsing algorithm을 이론적으로 토론 하였다. 모아쓰기 algorithm은 한글의 입력과 출력 모두에 응용이 가능하며, device-independence를 위해 "cardinal symbol set"의 개념을 소개하였다. 제 2 부에서는 한글 컴퓨터 처리의 역사적 개론과 상기 algorithm의 implementation 문제점들을 토론하고자 한다.

Abstract

This is Part I of a two-part article on Korean character processing by a computer. In part I, the problems in Korean character processing are identified and the theoretical foundation is laid out as a viable solution to them. The one- and two-dimensional syntactic structures of Korean characters are formally defined by means of BNF and "Patternal structure" respectively. Formal discussion of lexical and syntactic algorithms is given for character conversion. This character conversion algorithm is applicable to both input and output. For device-independence and implementation-independence, the concept of "cardinal symbol set" is introduced. We will present a historical survey of Korean character processing and discussion of implementation problems for the above algorithm in Part II.

1. Introduction

The history of efforts in computerizing the Korean character processing is relatively long considering the short history of computing in Korea. These efforts have been directed to two main areas: (i) discovery of Korean character processing algorithms and underlying theory, and (ii) mass production of input/output devices capable of handling Korean characters.

Naturally, any research and development pro-

ject would try to be successful in both areas (i) and (ii). In spite of the many attempts, the reason for the long delay of the introduction of commercially available Korean input/output terminals is largely due to the status of computing in Korea during the last decade. However, we have to admit that the scientists and engineers involved in the field are also responsible for the present status of Korean processing by computer.

A number of key problems in Korean character processing can be identified: (i) approaches to the algorithms for Korean I/O are more or less adhoc, (ii) these algorithms are too much device-dependent, and (iii) the lack of stand-

* 正會員, 韓國科學技術研究所
(Korea Institute of Science and Technology)
接受日字: 1979年 4月 23日

ards in keyboard set or coding scheme causes difficulties in communication and interfacing. High cost of terminals currently available is a serious obstacle to meeting the upsurging demands for Korean I/O terminals.

This article is a sincere attempt to provide a meaningful solution to the above problems (i) and (ii), with a high hope that this will lay a ground stone for future endeavors in the area (iii). It is believed to be a theoretical underpinning of Korean character processing, especially in relation to input/output and the linear to two-dimensional conversion of Korean characters. The syntactic structures of both linear and two-dimensional character representations have been analyzed (in Section 2). The theory of "patternal structure" is presented for the latter. The conversion algorithm is formally discussed (in Section 3) and the notion of the "cardinal symbol set" is introduced to preserve the device-independence. The theoretical background for the discussion is provided by the formal language and language translation theory of computer science.

A final remark should be made regarding the discovery of the algorithm and theory of this article. The discovery occurred independently without the knowledge of other works in late October of 1978. Considering the simplicity (and the accompanying beauty) of Korean character syntax, many independent works should have been carried out without being known. Thus, the claim is not a big one but being made in good faith for the full disclosure of a scholarly work.

2. Structure of Korean Characters

A. Formal Syntax of Korean Characters.

As any other language, the Korean language is generated by an alphabet (generating set of symbols). The "Korean alphabet" is divided into two distinct classes, "consonants" and "vowels", according to the phonetic role of a symbol.

There are 30 consonants and 21 vowels be-

ing used in modern Korean. "Consonants" are further broken down to three groups: simple (there are 14 of them), double (5), and compound (11) consonants. Similarly, vowels are grouped into two classes: simple (10) and compound (11). The alphabet symbols are classified and listed in table 1.

A proper combination of these 51 symbols produces a Korean "character" which is a unit of pronunciation and meaning, i.e., syllabus. Formation of a character is governed by a fixed set of rules, the subset of the syntax of Korean.

Each character must begin with a consonant, called "initial consonant or sound", which is followed by a vowel called "middle sound". Optional, third component of a character is a consonant which is called "terminal consonant or sound". Simple and double consonants can serve as both initial and terminal sounds.

Table 1.

Symbol class	alphabet symbols
Simple(14)	ㄱ, ㅋ, ㆁ, ㆅ, ㄷ, ㅌ, ㅁ, ㅂ, ㅅ, ㅇ,
Consoant Double(5)	ㅈ, ㅊ, ㅋ, ㆁ, ㆅ
Compound (11)	ㄴㅈ, ㄴㅊ, ㄴㄱ, ㄴㅅ, ㄴㅂ,
	ㄴㅇ, ㅁㅅ, ㅅㅈ, ㅅㅅ, ㄷㅅ, ㅁㅁ
Simple(10)	ㅏ, ㅑ, ㅓ, ㅕ, ㅗ, ㅛ, ㅜ,
Vowel	ㅠ, ㅡ, ㅟ
Compound (11)	ㅝ, ㅞ, ㅟ, ㅠ, ㅡ, ㅢ, ㅣ,
	ㅤ, ㅥ, ㅦ, ㅧ

Compound consonants can be a terminal sound only. The formal syntax of Korean characters is described by means of BNF notations in Figure 1.

There are two ways that Korean (sentence) can be written down. One way is to spell one symbol after another except for the necessary blanks or sentential marks (commas, colons, for example). That is, each position of a

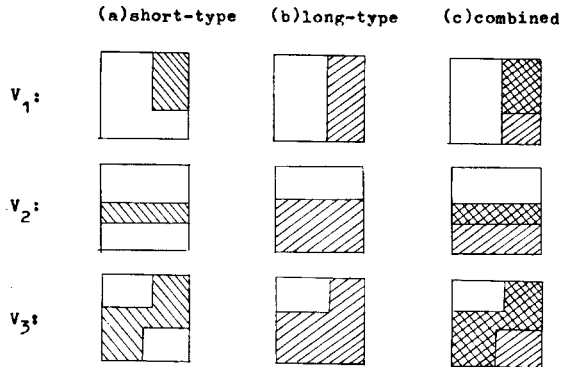


Fig. 2. Elementary patterns for vowels.

- (iii) type-3 consonants(C_3): initial consonants followed by mixed-shape vowels.
- (iv) type-4 consonants(C_4): terminal consonants preceded by short vertical or mixed-shape vowels.
- (v) type-5 consonants(C_5): terminal consonants preceded by short vertical or mixed-shape vowels.

Fig. 3 shows the window patterns for these consonant type. Shaded portions represent the areas occupied by consonants. We will call the above window patterns for vowels and consonants "elementary patterns."

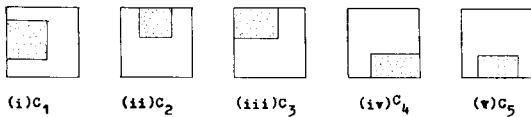


Fig. 3. Elementary patterns for consonants.

We now recall that there are only a finite number of ways that consonants and vowels are combined to comprise a meaningful character. This suggests that patternal structure of Korean characters should be described in terms of finite combination of elementary patterns. And this is exactly the case. The rules and codes for the patternal structure ("two-dimensional syntax" if you like) are succinctly summarized in table 2. As we may see, there are six combined patterns for legal Korean characters. These patterns are shown in fig. 4 together with examples of characters corre-

sponding to each pattern.

Table 2.

Vowel Types \ Consonant Types	Long-type			Short-type		
	V_1	V_2	V_3	V_1	V_2	V_3
Long-type	C_1	C_1V_1	-	-	-	-
Short-type	C_2	-	C_2V_2	-	-	$C_2V_2C_5$
	C_3	-	-	C_3V_3	$C_3V_1C_4$	-
	C_4	-	-	-	-	-
	C_5	-	-	-	-	-

Code : Pattern : Examples :

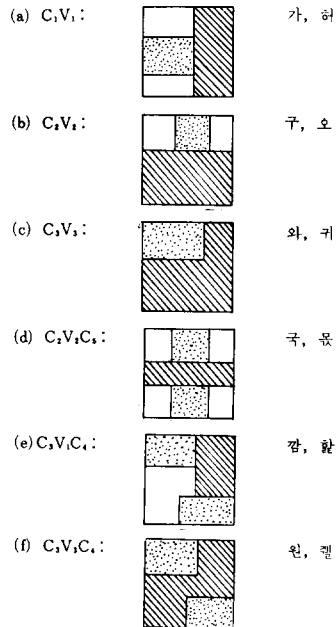


Fig. 4. Patterns for Korean characters.

3. Korean Character Combination / Conversion Algorithm

We have seen that the differences between the linear and two-dimensional character representations are rather superficial. By this we mean that the basic (abstract) syntax of two representations is identical as described by fig. 1.

This is so even though different concrete syntax based on patternal structure can be derived for two-dimensional character representation.

- (1) Lexical analysis : identification of chunks and characters from the linear input string.
- (2) Syntactic analysis : determination of symbol combination, i. e., patternal structure, of a 2-D character representation.

For implementation purposes, it is not essential that two parts are organized as separate but connected modules(e.g., coroutines). They may be integrated into a single program module.

A. Lexical Analysis

The most fundamental to lexical analysis is the property of the syntax of Korean that each character must begin with a (simple or double) consonant followed by one or two cardinal vowels.

Lexical analysis algorithm is best described by a finite automaton. It seems most convenient to

get two input symbols at a time , which also results in the simplest automaton.

The automaton is described in fig. 5.

The START state is entered for every character and , for all output states , it is necessary to re-set all pointers being used with an input string.

All "error states" are suppressed but can easily be added to the figure. Of course , any input symbol $a_i \in \emptyset$.

We note that the basis of the automaton is in fig. 2 - 1.

Three nonterminal states $S_1, S_2,$ and S_3 admit interesting interpretations:

S_1 means that initial and middle sounds have been identified (for noncardinal vowels, only the first symbol); S_2 means that a compound vowel has been detected; S_3 means hat a cardinal vowel is followed by a terminal consonant .

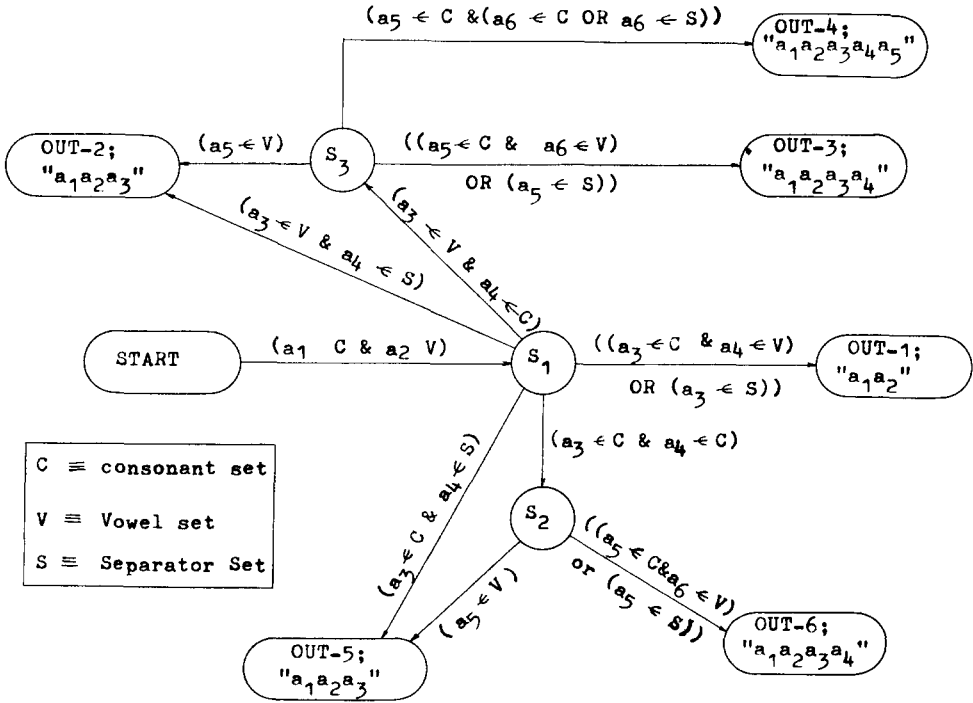


Fig. 5. Automaton for lexical analysis.

B. Syntactic Analysis (Parsing)

Syntactic analysis for Korean characters is extremely simple, to the degree of triviality. Once lexical analysis is done, the vowel (i.e., middle sound) of each character determines the patternal structure uniquely. The theoretical basis of this comes from table 2:

each column contains only a single entry. Discrimination between long and short vowel types is trivially done by the absence or presence of terminal consonant. The fact that choice of a column (from the type of a vowel) fixes the types of consonant (s) is useful for syntax error detection. That is, if the input consonant (s) type does not match that of a selected column, it signals the occurrence of a syntax error. From the viewpoint of table 2, syntactic analysis may be thought as the process of determining the patternal structure code, $c_i v_j$ or $c_i v_j c_k$, where $j \in \{1, 2, 3\}$ and $i, k \in \{1, 2, 3, 4, 5\}$

As a supplementary remark we add that proper encoding of alphabet symbols should be seriously considered for efficient lexical and syntactic analysis.

Fig. 6 describes the syntactic analysis algorithm, "SYNA", by means of a pyramid-like symmetric flowchart. Labels in EXIT nodes represent the patternal structure names as in Fig. 4.

(We note that the flowchart of fig. 6 is obtained by straight forward translation of table 2.

4. Conclusion

In this article, we discussed (i) the theory of Korean character syntax and patternal structure (for two-dimensional composition), and (ii) that of lexical and syntactic analysis algorithms for Korean character conversion. We were able to observe

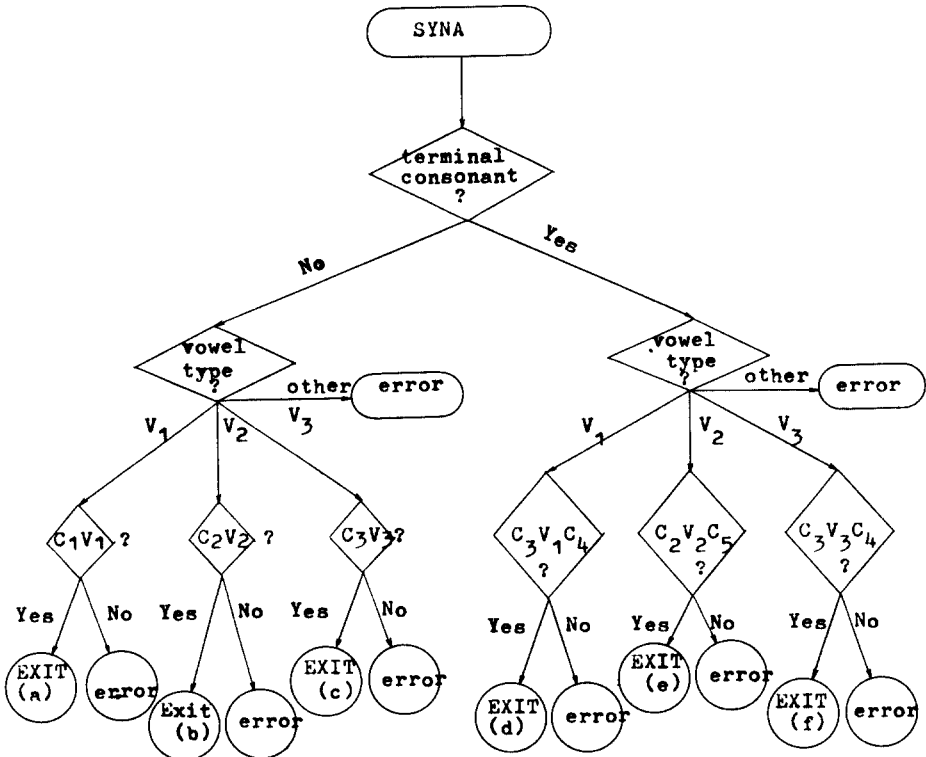


Fig. 6. Flowchart for syntactic analysis.

the utmost simplicity and beauty of Korean characters as the object of formal language and computer processing.

The discourse has been presented in a manner independent of device or implementation specifics. Thus any modification of the discussed algorithms will be straightforward.

We suggest the following applications or extension of the present theory :

(1) application of the patternal structure to on-line character recognition (hand-

written or printed).

(2) extension of the conversion algorithm to the lexical and syntactic analyses for a Korean (programming) language processor, and

(3) incorporation of formal and computer processing properties in determining the Korean standards for keyboards and coding scheme of Korean character processing.

Some of these problems will be discussed in the sequel.

