

A CHINESE CHARACTERS CODING SCHEME FOR COMPUTER INPUT AND INTERNAL STORAGE

Chorkin Chan, Computer Centre, University of Hong Kong, Hong Kong

Abstract

A coding scheme for inputting Chinese characters by means of a conventional keyboard has been developed. The code for each Chinese character is composed of two strings of keys, one corresponds to the spelling and the other the ideographic property of the character. Each code requires no more than seven keys (average five and a half keys) and 99.5% of the ten thousand characters in a dictionary 'XianDai HanYu CiDian' have unique codes. Each input code can be packed into 32 bits for internal representation.

Introduction

Over the last few years, encoding Chinese characters has become a very active subject of research. Numerous papers have appeared, mainly written in Chinese (hence difficult to be referenced in English), proposing various kinds of inputting schemes. Unfortunately, most of these papers offered only the ideas without accompanying implementation and experimentation. This paper presents a coding scheme of Chinese characters based on their ideographic properties as well as their spellings so that a conventional typewriter keyboard can be used for inputting purposes. This scheme has been implemented at the University of Hong Kong using an IBM 3031 under VM/CMS. Without a proper output device to display the Chinese characters, when the code of a Chinese character is entered, the address of that character (where it can be found) in a dictionary 'XianDai HanYu CiDian' is displayed. This is awkward but still sufficient to prove the correctness of the code recognition procedure.

The Coding Scheme for Inputting

In this scheme, a code for a Chinese character consists of two strings of symbols concatenated together. One string of three symbols corresponds to the ideographic radicals the character is composed of. The other of no more than four symbols is the spelling of the character. Corresponding to each of the ten thousand characters in the dictionary 'XianDai HanYu CiDian', with the exception of twenty six pairs, there exists a unique code in this scheme. In other words, this coding scheme is 99.5% unique. Furthermore, among these pairs of characters sharing the same codes in a pair-wise manner, the non-uniqueness

of eighteen pairs can be easily removed by deleting one member of each pair from the vocabulary because they are either 'dead' characters appearing in ancient classics only or they can be replaced by other characters of equivalent meaning. The non-uniqueness of the remaining eight pairs can be removed also by either changing the ideographic pattern or the spelling of one of the members in each pair. Thus, by means of these remedial measures, this coding scheme offers a unique code to each of the ten thousand characters found in this dictionary. The list of characters sharing the same codes is in Table 1 together with the suggested remedies to overcome the problem of non-uniqueness.

The Spelling of Chinese Characters

There are two standard systems to spell Chinese characters, one in terms of the Latin alphabets and the other in terms of Mandarin Pin Yin symbols. By means of the former; a maximum of five alphabets are normally required to spell a Chinese character. However, since the alphabet 'G' (except when it is the leading alphabet) always appear with 'N' as 'NG', one can replace 'NG' with 'G' and reduce the maximum number of alphabets required from five to four. By means of the latter, no more than three symbols are required to spell a Chinese character. This can be an important saving but in this paper, spellings are in terms of Latin alphabets just because a conventional terminal keyboard does not have Mandarin Pin Yin keys.

It is not always obvious whether one should read certain Chinese characters with or without a curling tongue, i.e., whether one should spell with 'C' or 'CH', 'S' or 'SH' and 'Z' or 'ZH'. This is particularly difficult to those whose mother tongue is not Mandarin. In order to be more forgiving, this coding scheme allows one not to differentiate 'C' from 'CH', 'S' from 'SH' and 'Z' from 'ZH' so that, for example, 'SHAO' can be spelled as 'SAO'. As a consequence, there will be three additional pairs of characters sharing the same codes in a pair-wise manner as listed in Table 2. Fortunately, the non-uniqueness so engendered can be easily eliminated by deleting one member of each pair because of its rare occurrence. For the same reason, this coding scheme also allows one to confuse a leading 'N' with a leading 'L'. For example, 'LUAN' can be spelled as 'NUAN' and vice versa. No non-uniqueness is introduced as a result of this

Table 1: Pairs of Chinese Characters Sharing the Same Codes

Spelling	Radical Composition	Char. 1	Char. 2	Suggested Remedy	Justification
AN	5	厂	广	delete char. 2	same meaning
BI	9TE	蔽	弊	write char. 2 as 痲	it means a defect
BO	M	馗	跛	delete char. 1	same meaning
DIAO	V	鸟	吊	delete char. 2	replaced by 吊
DUN	KB	吨	囤	delete char. 1	uncommon
E	KGX	阿	婀	delete char. 1	uncommon
E	T2-	吡	囧	delete char. 2	uncommon
FU	VM	危	弗	spell char. 2 as FO	so is 佛
GU	K2	咕	固	delete char. 1	uncommon
JIA	FDK	架	枷	write char. 2 as 枷	metallic shackle
JIAN	JY	兼	检	delete char. 1	uncommon
JING	D6=	劦	劲	write char. 1 as 例	human activity
JUAN	PKL	賈	冼	delete char. 2	uncommon
LIAN	-Y	鍊	殮	delete char. 1	uncommon
LING	KYR	呤	囧	delete char. 1	uncommon
MAO	HOP	瑁	髦	delete char. 2	uncommon
PANG	87	旁	旁	write char. 1 as 旁	that's original
QI	I	屹	屹	write char. 2 as	that's original
SHAO	2DK	邵	邵	delete char. 2	replaced by 邵
SI	I	厶	私	write char. 1 as 私	that's original
XIAO	8EL	宵	霄	write char. 2 as 霄	being celestial
YI	F?	杙	棧	delete char. 1	uncommon
YI	?	乙	弋	delete char. 2	uncommon
YU	O	鱼	禺	spell char. 2 as OU	so is 禺
YUN	2K;	陨	郢	delete char. 2	uncommon
ZHANG	27	幛	障	write char. 1 as 幛	made of fabric
ZHEN	JPX	椹	斟	delete char. 1	uncommon
ZI	X.X	些	吡	delete char. 1	replaced by 些
ZI	;.X	些	些	delete char. 1	replaced by 些

relaxation because the complete code consists of the radical string as well as the spelling string. Over the ten thousand characters in 'XianDai HanYu CiDian', this coding scheme requires an average of 2.5 alphabets to spell a Chinese character.

The Radical Composition of Chinese Characters

One traditional method of looking up a Chinese character in a dictionary is first to identify a radical in the graphic representation of the character. There are hundreds of different standard radicals used in a dictionary and there are rigid rules to apply in order to identify one. The number of Chinese characters identified to a single radical is numerous. Even a combination of the spelling and the identifying radical together is not sufficient to yield a unique code for a Chinese character.

An experiment was conducted in which each of the ten thousand characters mentioned above was decomposed into a string of as many as eight radicals. In order to do so, a total of four hundred and fifty six radicals were employed. These radicals were grouped into fifty sets according to their common graphical properties. Each set is then associated with a key of a conventional keyboard. Table 3 lists all these radicals, their groupings and their associations with the keys of a keyboard. Human engineering aspects were considered when the set-key association was determined. The radical string for a Chinese character consists of the keys corresponding to the first three radicals composing the character. In case the character is decomposed into less than three radicals, blanks are used as fillers to make up a string of three keys. For instance, the character 将 is decomposed into 9LT and the radical string for 将 is I. In this coding scheme, the grouping of radicals into sets is of paramount importance. On the one hand, they are grouped according to their common graphic properties into as few sets as possible. On the other hand, care is exercised to assure the uniqueness (or almost uniqueness) of the code-character correspondence.

The Coding Scheme for Internal Representation

For data processing purposes, it is

necessary to arrange the Chinese characters into a collating sequence which is a direct result of their internal representation in computer memory. Hence, when one is designing the internal codes, besides minimizing the length of the codes, one should also observe that the collating sequence that follows is logical and practical. This paper attempts to derive the internal codes logically from the input codes which, in turn, are logically related to the spellings and graphical properties of the Chinese characters. When a new character is created in the future with a unique input code, this scheme guarantees that the internal code will also be unique and a logical place in the collating sequence for it is assured.

The maximum number of keys used for an input code is seven. Storing seven symbols, in general, requires seven bytes. We recall that three symbols out of the seven serve to indicate which sets of radicals the Chinese character is composed of. Since there are fifty sets of radicals altogether, there are a total of 125,000 possible combinations. Seventeen bits will be sufficient to represent these combinations. The remaining four alphabetic symbols used to represent the spelling have the following properties:- The first symbol can be any alphabet from A to Z (except V). Five bits would suffice to represent it. The second symbol can be a blank, A, E, H, I, M, N, O, R, U or V, a total of eleven possibilities. Four bits would suffice. The third symbol can be a blank, A, E, G, I, N, O, or U, a total of eight possibilities. Three bits would suffice. The fourth symbol can be a blank, A, G, I, N, O, or U, a total of seven possibilities. Three bits would suffice.

Thus the spelling can be packed into fifteen bits. Combining with the seventeen bits required for the radicals, a code in these scheme requires only thirty two bits of memory space.

As a consequence of this internal representation, the collating sequence would be such that where a character should appear in the sequence first depends on the spelling of the character. The order of two characters of the same spelling depends on the keys used in the radical strings for the two characters.

Key	Radicals in Sets
:	貝用典肉肉甫舟骨
:	彳亍
'	而丙兩面扇斷百頁兩西
"	...黑
,	ノ未各重垂重手白白白牛
<	車東東七
.	虫止止业丑豎翅业
>	子女聿隶肃肃尹尹古
/	讠(言)讠
?	乙弋弋戈戈戈弋弋气电毗我幾 幾(乙)

The Next Step

In order to evaluate the effectiveness of this coding scheme, the author plans to experiment with different users and measure their coding efficiencies as a function of training and experience as well as their reaction towards using this scheme. The acceptance of the users is the ultimate measure of success of any invention. The design of the set-key association in Table 3 is somewhat arbitrary. Since it has a subtle impact on the collating sequence, more research in this area is necessary.

Acknowledgement

The author is indebted to Professor T.C. Chen for his constructive suggestions and criticisms. The author is also grateful to Mr. T.H. Tse for his assistance and discussions.