

## THE UNIFIED PHONETIC TRANSCRIPTION FOR TEACHING AND LEARNING CHINESE LANGUAGES

Jiann-Cherng Shieh

Graduate Institute of Library and Information Studies  
National Taiwan Normal University, Taiwan  
jeshieh@ntnu.edu.tw

### ABSTRACT

In order to preserve distinctive cultures, people anxiously figure out writing systems of their languages as recording tools. Mandarin, Taiwanese and Hakka languages are three major and the most popular dialects of Han languages spoken in Chinese society. Their writing systems are all in Han characters. Various and independent phonetic transcriptions have been thus developed to be as the mapping mechanisms between Chinese mother tongue languages and Han characters. For teaching and learning facilitation purposes, we really require a convenient phonetic transcription system between daily Mandarin, Taiwanese and Hakka to speed Han characters data processing applications. The Roman spelling system is a universal tool that owns the one and only one spelling rule. By studying and analyzing the Roman spelling system, we have disclosed that 4135 Romanized phonetic transcriptions can be adequately applied to handle Han characters' mappings of Mandarin, Taiwanese and Hakka spoken dialects. In this paper, we propose a minimal perfect hashing function to process unified 4135 Mandarin, Taiwanese and Hakka Romanized phonetic transcriptions to their corresponding Han characters simultaneously. The unified phonetic transcription can be used to promote Chinese mother tongue languages applications and developments. Furthermore, it can be applied as a mechanism to popularize digital learning and teaching of Chinese mother tongue languages.

### INTRODUCTION

People generally recognize that it is valuable to teach and learn mother tongue languages in today societies. People anxiously figure out writing systems of their languages to record and preserve their distinctive cultures. Mandarin, Taiwanese and Hakka languages are the three major spoken dialects in Chinese society. There are many speakers of the languages in China, Malaysia, Singapore, Philippine, Thailand and Indonesia.

Mandarin is the widest spoken language in the world and there are about 1300 millions people worldwide. Pinyin, more formally Hanyu Pinyin, is the most common Standard Mandarin Romanization system in use. Hanyu means the Chinese language. pin means "together, connection, annotate" and yin means "sound". Pinyin uses the Latin alphabet to represent sounds in Standard Mandarin. Taiwan has adopted Tongyong Pinyin on the national level since October 2002. Tongyong Pinyin is a modified version of Hanyu Pinyin. Based on the Chinese remainder theorem, Chang and Wu (Chang & Wu, 1988) designed the hashing function to process 1303 distinct Mandarin phonetic transcriptions of Han characters.

Minnanyu refers to a family of Chinese languages which are spoken in southern Fujian and neighboring areas, and by descendants of emigrants from these areas in diasporas. It is usually called Taiwanese by residents of Taiwan, and Hokkien by residents of Southeast Asia. Taiwanese can be written with the Latin alphabet using a Romanized orthography which was developed first by Presbyterian missionaries in China and later by the indigenous Presbyterian Church in Taiwan; use of the orthography has been actively promoted since the late 19th century. Taiwanese is one of the most used dialects spoken in Taiwan, and evolved from the ancient languages of China, the Ho-Lo language family. According to the traditional but representative and authoritative Taiwanese dictionary (Shen, 2001), Shieh (2003) developed the hashing function of 3028 Taiwanese phonetic transcriptions of Han characters.

Hakka dialect is one of the seven major spoken dialects in Chinese Society. The Hakka language has numerous dialects spoken in southern provinces of China, Taiwan, Singapore, Philippine and Indonesia. It is the 32nd widest spoken language in the world and there are about 100 millions Hakka speakers worldwide. Hakka is not mutually intelligible with Mandarin, Cantonese, Minnan and most of the significant spoken variants of the Chinese language. The Hakka dialects across various China provinces differ phonologically, but the Meixian dialect of Hakka is considered the archetypal spoken form of the language. Shieh and Hsu (Shieh & Hsu, 2007) proposed a minimal perfect hashing function for the 1428 Hakka phonetic transcriptions of Han characters from authoritative Meixian Hakka dialect dictionary (Lee, 1995).

Many researchers have enthusiastically endeavored to study related the spoken languages subjects such as language curriculums in multicultural society (Kilimci, 2010), typologies of spoken language learning aids (Kartal, 2005), mappings between spoken language and its writing system (Chang & Wu, 1988; Shieh, 2003;

Shieh & Hsu, 2007) (as depicted in Figure 1), etc. They are striving to protect and promote their individual native cultures, and make them widespread utilization.

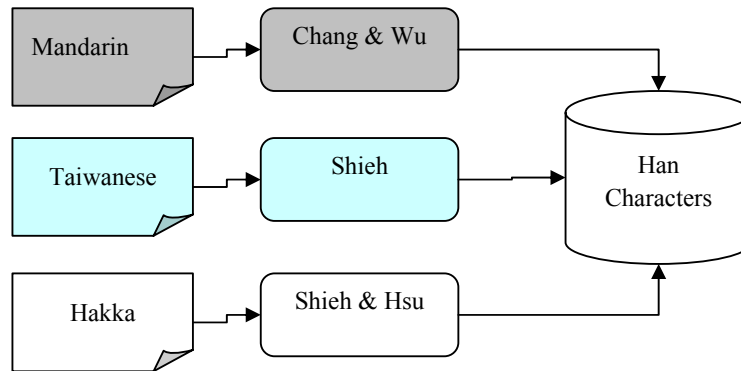


Figure 1: Various and Independent Han characters mappings for Chinese Languages

In Chinese societies, for teaching and learning facilitation purposes, we really require a convenient phonetic transcription system between daily Mandarin, Taiwanese and Hakka to speed Han characters data processing applications. These spoken languages are all with their respective Romanized phonetic transcriptions. Pleasantly surprised, the Roman spelling system is a universal tool that owns the one and only one spelling rule and can be generally and simultaneously applied to different languages applications. By studying and analyzing the Roman spelling system, we have disclosed that 4135 Romanized phonetic transcriptions can be adequately applied to handle Han characters' mappings of Mandarin, Taiwanese and Hakka spoken dialects. The 4135 integrated phonetic transcriptions are composed of 7 tones, 29 consonants, and 120 vowels at most. For language application purposes, it is much important for us to establish a mechanism to efficiently retrieve different Han characters and their corresponding pronunciations from its vocabulary repository, as illustrated in Figure 2. Many Chinese language learning applications, such as on-line or mobile dictionaries, translations, text-to-speech conversions, e-books, etc., can be further developed to help learners and teachers.

In this paper, we apply the Chinese remainder theorem to design a fast and efficient hashing function (Knuth, 1998) to map the unified 4135 phonetic transcriptions to corresponding Han characters of Mandarin, Taiwanese and Hakka languages. We also give a proof that the loading factor is more than 0.887, which is the best one when applying the Chinese remainder theorem to the design of hashing functions for the word sets.

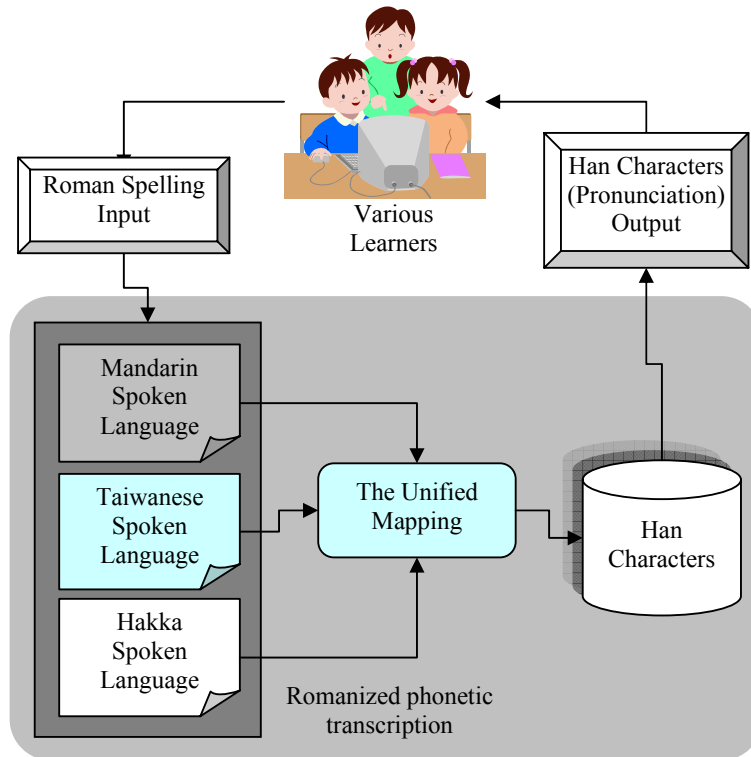


Figure 2: The Unified Mapping for Chinese Languages

### Hashing Functions Based on the Chinese Remainder Theorem

In this section, we first introduce the Chinese remainder theorem and its application to hashing function designs of character data sets. Then we review the hashing function designs of Mandarin, Taiwanese and Hakka phonetic transcriptions based on the theorem.

#### *The Chinese remainder theorem* (Chang & Lee, 1986)

Theorem 1. Let  $r_1, r_2, \dots, r_n$ , be  $n$  integers. There exists an integer  $C$  such that  $C=r_1 \pmod{m_1}$ ,  $C=r_2 \pmod{m_2}$ , ...,  $C=r_n \pmod{m_n}$ , if  $m_i$  and  $m_j$  are relatively prime to each other for all  $i \neq j$ .

For example, let  $r_1=1, r_2=2, r_3=3, r_4=4$  and  $m_1=4, m_2=5, m_3=7, m_4=9$ . Here  $m_i$  and  $m_j$  are relatively prime for  $i \neq j$ ,  $1 < i, j < 4$ . By the Chinese remainder theorem, there exists an integer  $C=157$  such that  $C \pmod{m_1}=157 \pmod{4}=1=r_1$ ,  $C \pmod{m_2}=157 \pmod{5}=2=r_2$ ,  $C \pmod{m_3}=157 \pmod{7}=3=r_3$ ,  $C \pmod{m_4}=157 \pmod{9}=4=r_4$ .

The following theorem results easily from the Chinese remainder theorem.

Theorem 2. Given a finite integer key set  $K=\{L_1, L_2, \dots, L_n\}$ . If  $L_i$  and  $L_j$  are relatively prime to each other for all  $i \neq j$ , there exists a constant  $C$  such that  $h(L_i)=C \pmod{L_i}$  is a minimal perfect hashing function (Chang & Lee, 1986).

#### *Hashing scheme based on the Chinese remainder theorem*

Based on the Chinese remainder theorem, Chang and Lee (1986) proposed a letter-oriented minimal perfect hashing scheme for a set of words. For a finite word set  $K=\{L_1, L_2, \dots, L_n\}$ , it is heuristically assumed that there exist  $s_1$  and  $s_2$  such that the extracted letter pairs  $(L_{i1}, L_{i2})$  are distinct, where  $L_{i1}$  and  $L_{i2}$  are the  $s_1$ -th and  $s_2$ -th characters of the word  $L_i$ ,  $i=1, 2, \dots, n$ . Chang and Lee's hashing function is defined as  $h(L_i) = H(L_{i1}, k_{i2}) = d(L_{i1}) + C(L_{i1}) \pmod{p(L_{i2})}$ , where  $d$  and  $C$  are integer value functions, and  $p$  is a prime number function. Chang and Lee's applied the hashing scheme to 12 months and 9 major planets with 0.154 and 0.103 loading factors respectively.

When applying the Chinese remainder theorem to the design of letter-oriented minimal perfect hashing functions, we often encounter the intractable issue of extracting letters from the word sets to form distinct letter pairs, especially from large data sets. Chang and Shieh (1985) used a zero value rehash index to resolve the problem. They successfully applied the technique to rehash the 59 reserved words for data-flow language VAL,

the 65 Z-80 commands, and the 256 frequently used words. Furthermore, Chang and Wu (1988) utilized the characteristics of Mandarin phonetic symbols to cluster the word set and then produced 1303 distinct letter pairs. The hashing scheme is introduced in the next section.

**Mandarin phonetic symbols hashing scheme** (Chang & Wu, 1988)

Chinese characters are constructed by 37 Mandarin phonetic symbols accompanied by one of the five tones. There are a total of 1303 distinct Mandarin phonetic transcriptions of Chinese characters. The phonetic symbols are divided into three categories: (1) the consonant, (2) the first vowel, and (3) the second vowel. For each symbol  $x$  in the symbol set, we have its order  $O(x)$ . In the hashing scheme, Chang and Wu translate all the phonetic transcriptions to letter pairs of two phonetic symbols.

Chang and Wu (1988) then cluster all letter pairs according to the five tones. In each equal-tone cluster, letter pairs with the same leading character are further grouped together. We see that the maximum number of character pairs in one group might go up to 33. From the experiment, as applying the Chinese remainder theorem, this would make the constant  $C$  quite large. By dividing the character pairs into three sets, they thus can assign the least 11 prime numbers to the corresponding characters in each group of the three sets. The minimal perfect hashing function is defined as  $H_j(L_{i1}, L_{i2}) = d_{jk}(L_{i1}) + C_{jk}(L_{i1}) \bmod p(L_{i2})$ , where  $d_{jk}$  and  $C_{jk}$  are integer value functions of each  $L_{i1}$  in the  $k$ -th set of each  $j$ -tone cluster, and  $p$  is a prime number function of each  $L_{i2}$ . The total size of space used is  $38 \cdot (3 \cdot (5 \cdot 2 + 1) + 3) + 1303 = 2671$ , where 38 stands for 37 phonetic symbols and 1 dummy symbol;  $3 \cdot (5 \cdot 2 + 1)$  stands for  $d_{jk}$  and  $C_{jk}$  of 5 clusters and index  $k$  in three sets. The number 3 is for the functions  $O$ ,  $p$ , and  $W$ . Thus, the loading factor is about 0.4878. If only the contiguous space is considered, the size of the space that is used becomes  $38 \cdot 3 \cdot 14 + 1303 = 2899$ ; the loading factor is about 0.45.

**Taiwanese phonetic transcriptions hashing scheme** (Shieh, 2003)

The Taiwanese phonetic transcription system, referred to a traditional but representative and authoritative Taiwanese dictionary (Shen, 2001), is composed of 7 tones (Table 1), 15 consonants (Table 2), and 45 vowels (Table 3). Each Taiwanese phonetic transcription consists of a vowel, a consonant, and a tone. Theoretically there are a total of 4725 transcriptions. However, only 3028 of the transcriptions are associated with Han characters. Shieh (2003) takes these 3028 transcriptions as study word set.

Table 1: Taiwanese Seven Tones

Code $k_{i1}$	1	2	3	4	5	6	7
Tone	kun	kún	kùn	kut	kûn	kün	kuť

Table 2: Taiwanese Fifteen Consonants

Code $k_{i2}$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
Consonant	liú	pinn	kiú	khi	tē	phó	thann	tsan	jíp	sí	ing	mng	gí	tshut	hí
Assigned Prime $P(k_{i2})$	2	3	5	7	11	13	17	19	23	29	31	37	41	43	47

Table 3: Taiwanese Forty-five Vowels

Code $k_{i3}$	Vowel	Code $k_{i3}$	Vowel	Code $k_{i3}$	Vowel	Code $k_{i3}$	Vowel	Code $k_{i3}$	Vowel
1	kun	10	kuan	19	kam	28	ka	37	kong
2	kian	11	ko	20	kue	29	ki	38	ń
3	kim	12	kiau	21	kang	30	kiu	39	muê
4	kui	13	ki	22	kiam	31	kenn	40	îng
5	ka	14	kiong	23	kau	32	kng	41	kiaunn
6	kan	15	kau	24	khia	33	kiô	42	tsim
7	kong	16	kai	25	kuè	34	kiunn	43	ngâu
8	kuai	17	kin	26	kam	35	kuan	44	kiann
9	king	18	khiong	27	ku	36	koo	45	kuan

Shieh handled 3028 distinct letter pairs of  $(k_{i1}, k_{i2}, k_{i3})$ 's, each with  $k_{i1}$  tones,  $k_{i2}$  consonants, and  $k_{i3}$  vowels. He sorted these letter pairs by their lexical orderings and then assigned each  $(k_{i1}, k_{i2}, k_{i3})$  a unique address. According to  $k_{i1}$ , Shieh got seven groups and computed their starting addresses  $d(k_{i1})$ . For each group  $k_{i1}$ , based on 15 consonants, he produced 15 tone/consonant subgroups and computed their corresponding relative

subgroup starting addresses  $d_{k_{i1}}(k_{i2})$ . For each subgroup, there are at most 45 letter pairs. He clustered the subgroup into 5 bunches by  $b(k_{i3})$  and also calculated each relative starting address  $d_{k_{i1},k_{i2}}(b(k_{i3}))$ , where each  $k_{i1}$  is associated with  $b(k_{i3})$ . Then he sequentially assigned the least 9 prime numbers  $P(k_{i3})$ 's to  $k_{i3}$  cyclically in each tone/consonant/vowel cluster. Finally, for every cluster, he applied the Chinese remainder theorem to compute constant  $C_{k_{i1},k_{i2}}(b(k_{i3}))$  such that  $C_{k_{i1},k_{i2}}(b(k_{i3})) \bmod P(k_{i3})$  equals the relative address of the cluster. The corresponding minimal perfect hashing function is defined as  $H(k_{i1}, k_{i2}, k_{i3}) = d(k_{i1}) + d_{k_{i1}}(k_{i2}) + d_{k_{i1},k_{i2}}(b(k_{i3})) + C_{k_{i1},k_{i2}}(b(k_{i3})) \bmod P(k_{i3})$ . Totally, it takes 4235 spaces: 3028 for key words,  $7 d(k_{i1})$ 's,  $7*15 = 105 d_{k_{i1}}(k_{i2})$ ,  $7*15*5 = 525 d_{k_{i1},k_{i2}}(b(k_{i3}))$ 's,  $7*15*5 = 525 C_{k_{i1},k_{i2}}(b(k_{i3}))$ 's and 45  $P(k_{i3})$ 's. The loading factor is  $3028/4235 = 0.715$ .

#### **Hakka phonetic transcriptions hashing scheme** (Shieh & Hsu, 2007)

According to the selected Meixian Hakka dialect dictionary, the Hakka phonetic transcription system is composed of 6 tones (Table 4), 17 consonants (Table 5), and 72 vowels (Table 6). Each Hakka phonetic transcription consists of a tone, a consonant, and a vowel. However, only 1428 of the transcriptions are associated with Han characters. Shieh and Hsu took these 1428 transcriptions as study word set.

Table 4: Hakka Six Tones

Tone	Yin Ping	Yang Ping	Shang	Qu	Yin Ru	Yang Ru
Code $k_{i1}$	1	2	3	4	5	6

Table 5: Hakka Seventeen Consonants

Consonant	p	p'	m	f	v	t	t'	n	l	ts	ts'	s	k	k'	ŋ	h	Ø
Code $k_{i2}$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17

Table 6: Seventy-three Vowels

Code $k_{i3}$	Vowel	Code $k_{i3}$	Vowel	Code $k_{i3}$	Vowel	Code $k_{i3}$	Vowel	Code $k_{i3}$	Vowel	Code $k_{i3}$	Vowel
1	i	13	eu	25	ok	37	e	49	ip	61	uo
2	au	14	uŋ	26	at	38	iap	50	əm	62	n
3	u	15	on	27	uk	39	iun	51	iai	63	uaŋ
4	oŋ	16	oi	28	ap	40	ot	52	ət	64	ep
5	a	17	en	29	it	41	iuk	53	ua	65	uat
6	o	18	iam	30	et	42	ɿ	54	uai	66	uet
7	ai	19	ui	31	ak	43	iet	55	uan	67	iut
8	un	20	aŋ	32	im	44	iak	56	ion	68	m
9	iau	21	iaŋ	33	iuŋ	45	iok	57	əp	69	iui
10	an	22	ien	34	ian	46	em	58	io	70	uen
11	in	23	ioŋ	35	ut	47	ən	59	uon	71	uak
12	am	24	iu	36	ia	48	iat	60	uoŋ	72	uok

They handled 1428 distinct letter pairs of  $(k_{i1}, k_{i2}, k_{i3})$ 's, each with  $k_{i1}$  tone,  $k_{i2}$  consonant, and  $k_{i3}$  vowel. Shieh and Hsu sorted these letter pairs by their lexical orderings and then assigned each  $(k_{i1}, k_{i2}, k_{i3})$  a unique address. According to  $(k_{i1}, k_{i2})$ , they had  $6*17$  groups and compute their starting addresses  $d(k_{i1}, k_{i2})$ 's. Then, they assigned appropriate prime numbers  $P(k_{i3})$ 's for  $k_{i3}$ . Finally, for every group, they applied the Chinese remainder theorem to compute constant  $C(k_{i1}, k_{i2})$  such that  $C(k_{i1}, k_{i2}) \bmod P(k_{i3})$  equals the relative address of character pair  $(k_{i1}, k_{i2}, k_{i3})$  in group headed with  $(k_{i1}, k_{i2})$ . The corresponding minimal perfect hashing function is defined as  $H(k_{i1}, k_{i2}, k_{i3}) = d(k_{i1}, k_{i2}) + C(k_{i1}, k_{i2}) \bmod P(k_{i3})$ . It takes 1704 spaces: 1428 key words,  $6*17 C(k_{i1}, k_{i2})$ 's,  $6*17 d(k_{i1}, k_{i2})$ 's, and 72  $P(k_{i3})$ 's. The loading factor is  $1428/1704=0.838$ .

#### **The Unified Phonetic Transcription Design**

##### **Hashing Function Design**

The unified Mandarin, Taiwanese and Hakka Romanized phonetic transcription is composed of 7 tones (Table 7), 29 consonants (Table 8), and 120 vowels (Table 9) associated with a prime number  $P(k_{i3})$ . Each phonetic transcription  $(k_{i1}, k_{i2}, k_{i3})$  consists of a tone  $k_{i1}$ , a consonant  $k_{i2}$ , and a vowel  $k_{i3}$ . There are totally 24360 combinations of  $(k_{i1}, k_{i2}, k_{i3})$ 's. According to our further analysis, we worked out that we can use exactly 4135 phonetic transcriptions to associate their corresponding Han characters.

Table 7: Tones

Tone	1	2	3	4	5	7	8
------	---	---	---	---	---	---	---

Code $k_{i1}$	1	2	3	4	5	6	7
---------------	---	---	---	---	---	---	---

Table 8: 29 Consonants

Consonants	b	c	ch	chi	d	f	g	h	j	ji	k	kh	l	m	n
Code $k_{i2}$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
Consonants	ng	p	ph	r	s	sh	shi	t	th	ts	tsh	tz	v	null	
Code $k_{i2}$	16	17	18	19	20	21	22	23	24	25	26	27	28	29	

Table 9: 120 Vowels

Vowel	Code $k_{i3}$	$P(k_{i3})$	Vowel	Code $k_{i3}$	$P(k_{i3})$	Vowel	Code $k_{i3}$	$P(k_{i3})$
a	1	2	iat	41	2	mh	81	2
ah	2	3	iau	42	3	ng	82	3
ai	3	5	iauh	43	5	nggh	83	5
ainn	4	7	iaunn	44	7	o	84	7
ak	5	11	ie	45	11	oh	85	11
am	6	13	iem	46	13	oi	86	13
an	7	17	ien	47	17	ok	87	17
ang	8	19	iet	48	19	on	88	19
ann	9	23	ieu	49	23	ong	89	23
annh	10	29	ih	50	29	onn	90	29
ap	11	31	ii	51	31	onnh	91	31
at	12	37	iim	52	37	oo	92	37
au	13	41	iin	53	41	ot	93	41
auh	14	43	iip	54	43	ou	94	43
aunn	15	47	iit	55	47	u	95	47
aunnh	16	53	ik	56	53	ua	96	53
e	17	59	im	57	59	uah	97	59
eh	18	61	in	58	61	uai	98	61
ei	19	67	ing	59	67	uainn	99	67
em	20	71	inn	60	71	uan	100	71
en	21	73	innh	61	73	uang	101	73
eng	22	79	io	62	79	uann	102	79
enn	23	83	ioh	63	83	uat	103	83
ennh	24	89	iok	64	89	ue	104	89
ep	25	97	ion	65	97	ueh	105	97
er	26	101	iong	66	101	uei	106	101
et	27	103	iou	67	103	uen	107	103
eu	28	107	ip	68	107	uenn	108	107
i	29	109	it	69	109	uet	109	109
ia	30	113	iu	70	113	uh	110	113
iah	31	127	iuann	71	127	ui	111	127
iai	32	131	iue	72	131	uih	112	131
iak	33	137	iuh	73	137	uinn	113	137
iam	34	139	iui	74	139	uk	114	139
ian	35	149	iuk	75	149	un	115	149
iang	36	151	iun	76	151	ung	116	151
iann	37	157	iung	77	157	uo	117	157
iannh	38	163	iunn	78	163	ut	118	163
iaong	39	167	iut	79	167	yu	119	167
iap	40	173	m	80	173	Null	120	173

We handle 4135 distinct letter pairs of  $(k_{i1}, k_{i2}, k_{i3})$ 's, each with  $k_{i1}$  tone,  $k_{i2}$  consonant, and  $k_{i3}$  vowel. We sort these letter pairs by their lexical orderings and then assign each  $(k_{i1}, k_{i2}, k_{i3})$  a unique address. According to  $(k_{i1}, k_{i2})$ , we have  $7*29$  groups and compute their starting addresses  $d(k_{i1}, k_{i2})$ 's. There are at most 120 characters  $k_{i3}$  in each group  $(k_{i1}, k_{i2})$ . Then, we heuristically assign appropriate prime numbers  $P(k_{i3})$ 's for  $k_{i3}$ . Finally, for every group, we apply the Chinese remainder theorem to compute constant  $C(k_{i1}, k_{i2})$  such that  $C(k_{i1}, k_{i2}) \bmod$

$P(k_{i3})$  equals the relative address of character pair  $(k_{i1}, k_{i2}, k_{i3})$  in group headed with  $(k_{i1}, k_{i2})$ . The corresponding minimal perfect hashing function is defined as  $H(k_{i1}, k_{i2}, k_{i3}) = d(k_{i1}, k_{i2}) + C(k_{i1}, k_{i2}) \bmod P(k_{i3})$ .

The hashing function design of the unified phonetic transcription is summarized as follows:

Step 1: Using tone  $k_{i1}$ , consonant  $k_{i2}$  and vowel  $k_{i3}$ , we can have 4135 distinct letter pairs  $(k_{i1}, k_{i2}, k_{i3})$ . We sort them in their lexical orders and assign each a unique address.

Step 2: We allocate each  $(k_{i1}, k_{i2})$  group a  $d(k_{i1}, k_{i2})$ , the first address of the letter pairs headed with  $(k_{i1}, k_{i2})$ 's.

Step 3: Associated with each group  $(k_{i1}, k_{i2})$ , we assign each  $(k_{i1}, k_{i2}, k_{i3})$  a relative address.

Step 4: We assign appropriately the prime numbers  $P(k_{i3})$ 's to  $k_{i3}$ .

Step 5: Consider the letter pairs  $(k_{i1}, k_{i2}, k_r)$ ,  $1 \leq r \leq n$ , with the same  $(k_{i1}, k_{i2})$ , that is they are in the same group  $(k_{i1}, k_{i2})$ , and the corresponding assigned prime numbers of  $k_r$ 's are  $P_1, P_2, \dots, P_n$ , where  $P_1 < P_2 < \dots < P_n$ .

Applying the Chinese remainder theorem to find a constant  $C(k_{i1}, k_{i2})$  such that  $C(k_{i1}, k_{i2}) \equiv 1 \pmod{P_1}$ ,  $C(k_{i1}, k_{i2}) \equiv 2 \pmod{P_2}$ , ..., and  $C(k_{i1}, k_{i2}) \equiv n \pmod{P_n}$ . Our proposed minimal perfect hashing function is simply defined as  $H(k_{i1}, k_{i2}, k_r) = d(k_{i1}, k_{i2}) + C(k_{i1}, k_{i2}) \bmod P(k_r)$ . The values of all  $C(k_{i1}, k_{i2})$ 's are illustrated in the Appendix.

### Loading Factor Comparisons

Loading factor is used to measure the efficiency of memory usage in hashing design. It is defined as a ration of the number of data and the total size of memory used. The loading factor of the hashing function designed in this paper is derived as follows: (1) Used spaces, 4135 key words,  $7 \times 29$   $C(k_{i1}, k_{i2})$ 's,  $7 \times 29$   $d(k_{i1}, k_{i2})$ 's, and 120  $P(k_{i3})$ 's. We take 4661 spaces in total. (2) The loading factor is  $4135/4661=0.887$ . The following table (Table 10) shows the loading factors of various minimal perfect hashing functions designed for diverse word sets by the Chinese remainder theorem. Obviously, it can be shown that our hashing function is superior to others.

Table 10: Comparisons of Loading Factors

Word Sets	Loading Factor	
Names of 12 months	0.154	Chang & Lee, 1986
59 VAL reserved words	0.312	Chang & Shieh, 1985
65 Z-80 commands	0.263	Chang & Shieh, 1985
256 frequently used words	0.472	Chang & Shieh, 1985
1303 Mandarin phonetic transcriptions	0.448	Chang & Wu, 1988
3028 Taiwanese phonetic transcriptions	0.715	Shieh, 2003
1428 Hakka dialect phonetic transcriptions	0.838	Shieh & Hsu, 2007
4135 unified Mandarin, Taiwanese and Hakka phonetic transcription	0.887	This paper

### Number C Analysis

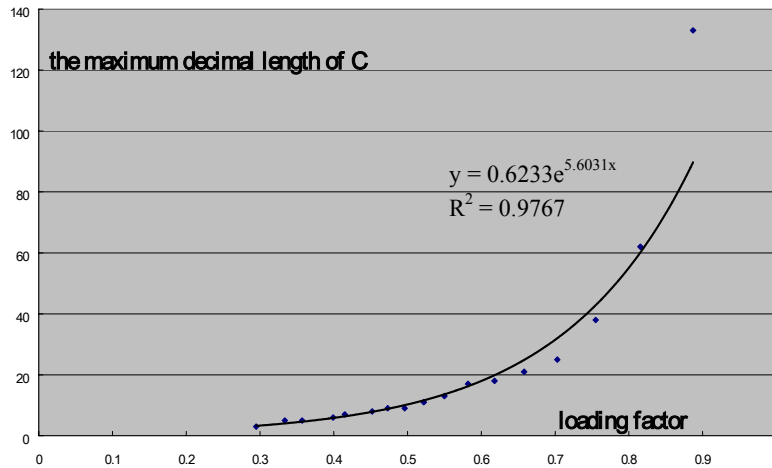
The numbers  $C$ 's are the most intractable ones as applying the Chinese remainder theorem to design hashing functions for data sets. On observing variations of  $C$ 's resulted by the experimental designs for the unified Romanized phonetic transcriptions, we have concluded that the size of  $C$  is dependent on the number of associated primes that we have used in each vowel group. In fact, during the hashing design, we can group  $k_{i3}$ 's vowels for each  $(k_{i1}, k_{i2})$  in different sizes to have alternate  $C$ 's results. The smaller size each vowel group has, the smaller constant  $C$  we result. However, what we should pay for is loading factor. There is a tradeoff between the size of constant  $C$  and the loading factor. The following table (Table 11) shows the experiments.

Table 11: The  $C$ 's and Loading Factors of Different Groups

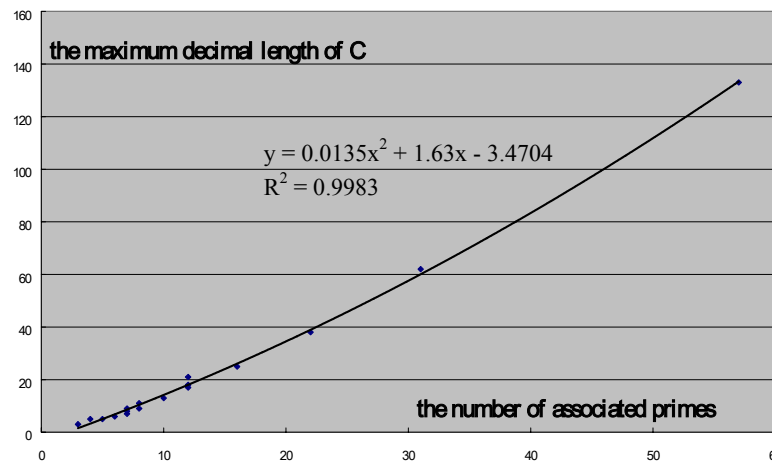
Number of vowel groups	Number of associated primes	Loading factors	The maximum length of $C$ 's
1	57	0.887	133
2	31	0.816	62
3	22	0.755	38
4	16	0.703	25
5	12	0.658	21
6	12	0.618	18
7	12	0.582	17
8	10	0.550	13
9	8	0.522	11
10	7	0.496	9
11	8	0.473	9
12	7	0.452	8
14	7	0.415	7
15	6	0.399	6

18	4	0.357	5
20	5	0.334	5
24	3	0.295	3

Next, we apply the statistical regression analysis to the experimental data to profile the correlations between the above parameters: loading factor vs. the decimal length of constant C, the number of associated primes vs. the decimal length of C and loading factor vs. the number of associated primes. The results are shown in the following Figure 3, where  $R^2$  is the coefficient of determination and its value is between 0 and 1.

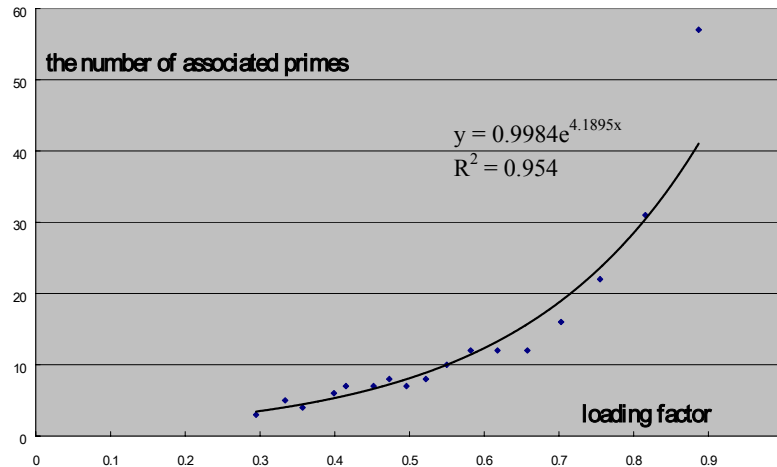


(a) Loading factor and the maximum decimal length of C



(b) The number of associated primes and the maximum decimal length of C





(c) Loading factor and the number of associated primes

Figure 3: The Regression Analysis of Effective Factors .

What the preferred situation is to have high loading factor with short decimal length of C. However, from the regression analysis, we have the fact that there is a tradeoff between two parameters. This will give us the concrete suggestion as we apply the unified phonetic transcription on diverse learning and teaching devices.

## CONCLUSIONS

With the unified phonetic transcription system for mapping Mandarin, Taiwanese and Hakka mother tongue languages to their Han characters, people will be convenient to promote their learning and teaching activities, as well as to record and preserve their particular cultures. In this research, we have successfully applied the Chinese remainder theorem to design a novel minimal perfect hashing function for 4135 integrated Mandarin, Taiwanese and Hakka Romanized phonetic transcriptions of Han characters. We have achieved significant results in terms of loading factors. We further give experimental investigation and mathematical regression analysis for considered factors of hashing effectiveness. We get the conclusion that the size of number C is dependent on the number of associated primes. For the unified Romanized phonetic transcriptions case, we propose the grouping technique to promote the effective applications as concerning the practicability of accessing constants C's. However, we have explored that there is a tradeoff between the loading factor and the size of C. The unified phonetic transcription can be used to promote Chinese mother tongue languages applications and can be applied as a tool to popularize digital learning of the languages further.

## REFERENCES

- Chang, C. C., & Lee, R. C. T. (1986). A letter oriented minimal perfect hashing scheme. *Computer Journal*, 29, 277-281.
- Chang, C. C., & Shieh, J. C. (1985). On the design of letter oriented minimal perfect hashing functions. *Journal of Chinese Institute Engineers*, 8, 285-297.
- Chang, C. C., & Wu, H. J. (1988). A fast Chinese characters accessing technique using mandarin phonetic transcriptions. *International Journal of Pattern Recognition and Artificial Intelligence*, 2, 105-137.
- Kartal, Erdoğan. (2005). The internet and autonomous language learning: A typology of suggested aids. *The Turkish Online Journal of Educational Technology*, 4(4), 54-58.
- Kilimci, Songül. (2010). Integration of the internet into a language curriculum in a multicultural society. *The Turkish Online Journal of Educational Technology*, 9(1), 107-113.
- Knuth, D. E. (1998). *The Art of Computer Programming III: Sorting and Searching 2nd ed.*, MA: Addison-Wesley, Reading.
- Lee, R. (1995). *Meixian Dialect Dictionary*, Jiangsu Education Publish.
- Shen, J. F. (2001). *Hui Im Po Kam. 46th ed.* Taiwan: Wen-I Publish.
- Shieh, J. C. (2003). An efficient accessing technique for Taiwanese phonetic transcriptions. *ACM Transaction on Asia Language Information Processing*, 2(1), 63-77.
- Shieh, J. C., & Hsu I. Y. (2007). The study of Hakka Han characters retrieval based on Chinese remainder theorem. In *Proceeding of The Sixth International Conference on Information and Management Sciences* (pp.1-6). Lhasa, Tibet, China.

## Appendix: The Values of C(k<sub>11</sub>, k<sub>12</sub>)'s

Tonek <sub>i1</sub>	Vowel k <sub>i2</sub>	d Address	C(k <sub>i1</sub> , k <sub>i2</sub> ) Value
1	1	0	64944529802520314615071967316445377
1	3	18	22072260041420641168670456121911469647357
1	4	37	2341262191069822373794511463164075
1	5	51	17826075993273201357220661350819824012907
1	6	70	36750251768899802683830396935
1	7	84	853423971975501060369045480800064457
1	8	102	3230888234697313835737580133054940934078475906950 6405134158841920335470218544507965940504191556559 11774608617063627
1	9	153	2811216111161579226740265389366581840581217
1	10	173	2341262191069822373794511463164075
1	11	187	5580647162611650774150618419340213976176344914817 6938660322244270242937251421132853009341113616024 428807332733392463007
1	12	239	5992871137556135789475657067386415924456169544983 3894341282590670541979389250807654745457
1	13	279	2419721735365016321777114959630980944915355085821 32838475657383488099190217
1	14	313	3977418278170894117865883174693348709678451922127 67
1	15	339	696097658238809533265458610083531288767
1	16	359	837645492417954537087209440
1	17	371	9256816887250690574107803504585337791212251648534 025093024420752547647
1	18	403	5302917695889361096059961272251746978804318170782 2176722777
1	19	430	1
1	20	431	9139333575553435344217573681380336055117755848887 4096007881464715928500115628213403688313249448958 96180387845109617
1	21	482	295741854374958593532533406802157
1	22	498	2341262191069822373794511463164075
1	23	512	7107816320470116249432012767652915896330712723618 0232942362905093585975451695572232441855417
1	24	554	1015974035642889592010734726566928011766396482742 4965345545544520122099986242077
1	25	590	2568083120669429978818728146297905012787131406591 3680434272401072566577907876566001440722639060786 2742466451591047
1	26	640	4964913407385885314546459313197740006643945653388 43097797541044913863962802531791645790796217
1	27	682	71711593239345115255790093017
1	28	696	12678327602740390681
1	29	705	2544866659299794137538047059881254367721733177260 5823322822721875158327900530386438590267244477609 41810631355723228823155204875566007
2	1	762	2659659253858485141086416953487398456564240347
2	3	784	4081499448220247076582253363034397

2	4	800	22492123753135128914677172498929
2	5	813	18918548605449
2	6	820	15656398053976404884240051025293
2	7	835	32478880744933868806045498990713283027
2	8	852	4514287264331388834844767949128952386082923444244 662398453514708904398465937552711849399021517
2	9	893	5576652203238090897311189140037
2	10	908	46991006860557
2	11	914	7839378807204054516147988363233949331261449480452 1543383818411446409126296404252400573506901816689 336887
2	12	960	7712718818069213264262424523835158900799430218587 906050597966600767581607
2	13	993	1514316853780982837710041954454136582163357047154 489327522740884312427908550466792087
2	14	1031	4659130194219555724727590242215279123970263126843 386987416675537
2	15	1062	1597654559044358957796421802051547923604424138759 4606870299138871861
2	16	1094	188459278322167127153310528405995558315107
2	17	1113	7703351501819123908547394978381198438297936947366 177839853089887
2	18	1143	33551547097592679892970230932194841
2	19	1159	2966945680576791803941513
2	20	1170	955226257940522272595537823100916958389536384672 0670887183208686293110648399890673471687
2	21	1210	80206105335638512818
2	22	1219	8100153285279498757088758568809
2	23	1232	3967019274860162899246388867812728434169084145156 22392718841084984610278516803037
2	24	1270	7291090621590647041439072207814174065327563872138 429998946
2	25	1297	2360804216443227435798667343668479052493606886105 2117344284790568803024813899597787348158445339764 18337
2	26	1344	1623532452270782186281332385066015582425401938240 4065514604339945261215768313481429641647
2	27	1383	879162072665
2	28	1390	1610599946459446681057
2	29	1400	1083206920889392469867123942521147945249568040525 7462415907115669137260295506321932051196389762439 74173060877007869777
3	1	1451	1749763638869608412023452687
3	3	1466	322241013048720509746348605415
3	4	1481	38996819222356748983769636
3	5	1492	1283654813299438293278779569397703907
3	6	1510	12944896164105491213843598330327
3	7	1526	115686602028041351850558986590963196427080136
3	8	1546	3357737135958534201527982731928588947077811758764 1764645865158324866160558099879928133577

3	9	1585	112519155483150811495060023100337397
3	10	1602	7375532636858408518398308344372
3	11	1615	1684783329030632754731377940271226718883299528812 2796893082348548895450031281698617653759518975954 179726981687
3	12	1663	7452694336649028941633590459945878578444992496635 454164446302943434372609536040359022872257447637
3	13	1706	1107987255140470453405591416742240701173456169477 31252124664935737261232456137
3	14	1741	1146716624527498163646366474429344250436465555139 7
3	15	1765	9378386571883512292646043084519325601841982717522 457
3	16	1790	129857631503867593810398312613400562127
3	17	1809	1324242754488800281714359122126851640079844770095 9730454708332410405727
3	18	1841	1300680649457591397293105543898150135112943781218 234973749353742267
3	19	1871	1938060110162566238
3	20	1880	1988671469739037638239403554225348077743556316160 6221940277390005968431648205240776763900096086435 69925366209153517
3	21	1930	266639658260349771068599529535437
3	22	1946	109377925939343038412220929
3	23	1957	5886821438505272931573958023368272745762577408572 173635907458823399172815340460914062560867
3	24	1998	1811038696058785063155490515767440481494605441286 52455346751735217536204765697
3	25	2033	2153966445706652547761825880466532229583531735956 2709056680833772470328455473953405921584279597819 927
3	26	2078	1215486198672120342411177158581553550434181396583 06461856626751777405797099461521832119009777
3	27	2119	37799600695088183854976164207
3	28	2133	16543624983193568638265
3	29	2144	2969146386276534554633176134274238100067395208131 7633140090885800771410499540926195098086014029692 4930108383645516211711941095757
4	1	2199	79913438279382459481489364265610711902572513
4	3	2222	724221717480452268740291021217
4	4	2237	6769896829431796783602782
4	5	2248	6980762466230080398861812412681587787
4	6	2266	17679299271350614000587
4	7	2278	1339920494387103920138721910610219766667857
4	8	2298	2951817521419370493764873433572225767662515619330 336139959221595448897864534760178582166524112417
4	9	2342	1794955116411917954709185512823943437
4	10	2360	5364090509388600007259758339069
4	11	2373	4780891390647215639980960268523434534887533915758 705362397056228729721990995951667915520218842

4	12	2415	26795197191227870856298511674858202168681232360269914
4	13	2439	1160120626706775692412796535591946154548304564859720398611688551537434225507192344678763
4	14	2480	973737779510096385392368583461073261924636037227
4	15	2504	1990947696902394643408173919694705315891614351764326875537
4	16	2532	28800703500190574616359876
4	17	2543	2610757393422535143644715869006825249716688643336441634608435164690063
4	18	2577	1255375958567247599798506332145332808
4	19	2594	55733162486284969083133
4	20	2604	43665913080555789728531326382132817474615494629781519356837546771964344981032202775906519620257
4	21	2647	155599977486564813010145580993355477
4	22	2664	2341262191069822373794511463164075
4	23	2678	56489372725693140581528568809536440159445884725284672118178475285043823045997390009233428753
4	24	2721	619401890263600802564735186852281568028143620
4	25	2742	15994428068795413828218144417715498169511362061528447687174623669850044300383819359940711176127
4	26	2784	117795152091889881210962711056161633914832844014897391
4	27	2808	124984874469926146257810472496
4	28	2822	5257164150129
4	29	2828	1623574972742378109983811417825200303692565072216639931617246225998695859273054766918434544702661591292579522920835441200893
5	1	2883	168878377699201632388054139240250981910241757
5	3	2904	1
5	4	2905	1
5	5	2906	1
5	6	2907	338826979773502888683565267
5	7	2921	546346856133854386461919403582789822821185162228595017
5	8	2946	2005397379617378022596671315414152052591005045578311624782663720962125478217362238156938751723125927
5	9	2991	2403837296262839863720324972
5	10	3003	20324972
5	11	3005	1526071016818703996330628204219975979083131452418117891912536510252620342305684332540851
5	12	3043	862323498696940745456525981650874688957556440144267190379151690016809838325631
5	13	3077	127673445075430236388007678348014265807775642697937591767821601530826359378197
5	14	3112	30539680348545329986649631833498718169991124565798288743276167
5	15	3141	9726901571118510230101687655349575127085799287

5	16	3165	1058728566088652243623586350911683250014000435416 23189878881177
5	17	3194	1669521193801301929676310284362728611750183160207
5	18	3217	3945647105453198330392078100279553534865048654928 0944903607
5	20	3244	4952424038687225802218686195906740998248491334207 651265315626967369909192024366546331790087
5	21	3285	31790087
5	22	3287	1
5	23	3288	1793530233601319628923995099342213372101796216563 70244686780813642822743121549977727
5	24	3326	1872459309563264049880522527678371796921339523547 477369476731761841687
5	25	3358	6260139171305869803523889775019808586114182015262 8497996105936522158303747192330346427
5	26	3398	5455540982003668924967943969509675551618840564446 19587690053613027588611174921537280818084507
5	27	3440	1
5	28	3441	6284602471488624036687
5	29	3452	1063512091338202606821365361821105909555754133304 111818326841617258522488818595962548625584587
6	1	3494	991291924043032694366940565068354227143
6	7	3512	189335613045177604054307401009276517
6	8	3529	1524817751545235818196122522779085103067082177245 19932343078540547
6	9	3559	77231424308070559
6	11	3566	5940175862658561808887053707488699515021403008687 1800571087
6	12	3592	35412161870767049414814
6	13	3602	1522525697674741080103176998268905023366602539287 170193797867
6	14	3630	1032697793415830176466938
6	15	3641	289578503341955866322657
6	16	3653	3006263
6	17	3657	1385974483316912168022583053437254148038608485598 7
6	18	3680	98334985008736815942514736877516469074908177
6	20	3701	8494617303933234208076030365057578308815484312989 98996
6	23	3724	3753403896537713415983677550634223700790719902095 20661495865142500727
6	24	3756	6337471000405276586237115312859281
6	25	3772	9627014677185898406282200050135665841754518693784 3951171786971032177
6	26	3803	801128368398242695123417
6	29	3816	5189813134507463236347957040810627329486403882309 551659472936601273
7	1	3846	214643791629891130652686735759237
7	6	3861	182287063653
7	7	3866	5068148460282900594363322574201

7	8	3879	73682591563608850472748865300829284731733
7	9	3898	3452032766467
7	11	3903	124403403267415881506160389626
7	12	3916	4935409942089541429388313310860735781
7	13	3933	4482628447566173767329199931470555272675772565355 70497949
7	14	3959	2562991393797764131711391
7	15	3971	1661275737352881
7	16	3978	1456363052715719456
7	17	3986	12172994738807119694929487674
7	18	3999	16222681315494410513594737121487
7	20	4013	3718892506217348358334646674903246426744707814296 385
7	23	4036	7166523818239668976161736537297480181038766848
7	24	4057	200218239419106128490513795004
7	25	4071	1613452437872808691961211787657095315483678643035 81037
7	26	4095	33287520344187096554000336156494463862
7	28	4112	175858486
7	29	4116	10966326177647918930460307313434677092766934