

TEXT PROCESSING OF THAI LANGUAGE
=THE THREE SEALS LAW=

Shigeharu Sugita

National Museum of Ethnology
Expo Park, Senri, Suita
OSAKA 565, Japan

Abstract

Computer softwares for processing Thai language are developed at National Museum of Ethnology, Osaka, Japan. We use a popular intelligent terminal TEKTRONIX 4051 for inputting and editing, IBM 370 model 138 for KWIC making and sorting, and CANON's laser beam printer for final output.

Using these systems, "Kotmai Tra Sam Duang" (the Three Seals Law) which contains many kind of laws and ordinances proclaimed in Thai between 1350-1805 A.D. is computerized. This text has 1700 pages and about 1400000 letters. KWIC index becomes 200000 lines.

Some statistical data for this text are obtained. They are occurrence frequency data of single letter, group vowel, and letter combination (digram), etc.

Aknowledgements

This report is a result of joint project at National Museum of Ethnology. The member are Y. Ishii, I. Akagi, S. Tanabe, Y. Sakamoto, S. Uemura, A. Ishizawa, M. Sawamura, K. Sasaki, Y. Kurita, and S. Sugita. Their research field are ethnology, linguistics, computer science, and sociology etc.

We thanks Mr. Sophon Chitthasatcha, Miss Sumalee Maungpaisaln and Miss Hiroe Matsumoto for their help in segmentation, inputting and correction.

We also thanks Prof. K. Nakayama and A. Oikawa of Tsukuba University for their support on making Thai letter patterns and output software for laser beam printer.

Introduction

In the field of ethnology or cultural anthropology, ethnographies are very important information sources for comparative study of many different societies. Not only bibliographic data but also contents of text are necessary.

HRAF (Human Relations Area Files), which was developed by Dr. Murdock and now managed by HRAF Inc. at Yale University, is a unique retrieval system.

They use about 800 category codes by which analysts classify the contents of each pages of books.

Though HRAF system is an elaborate work, it is not easy to search necessary data by user terms, that is, natural words. If whole text are fed into computer, it is very easy to retrieve any part of text by the same natural words used in the text.

On-line retrieval system is smart and effective. But sometimes researcher wants printed index like as KWIC which is usable at any time and place. Combining KWIC index and thesaurus dictionary, it gives us a very powerful tools for searching special expression hidden in the text.

Till quite recently, at least in Japan, most cases of computer processing of natural language are distorted to indo-european language or Japanese. In the ethnological studies, we must treat many areas in the world. We need computer softwares which process unfamiliar languages for us, such as Arabic, Korean, Sumerian, Mongolian, Devanagari, Thai, etc.

National Museum of Ethnology at Osaka has introduced several computer systems to encourage humanity study, and now is developing many application softwares which are usable by any researchers who do not know computer programming or how to use computer.

This report describes one of such application softwares which treats Thai letters. The points of our work are as follows;

- 1) A popular computer terminal is used for Thai letter inputting and editing. It is easy to use because dead key operation is not necessary.
- 2) KWIC making and sorting software are implemented using FORTRAN language which can be transferred to any other computer system. The algorithm is not so complex but it was not implemented only because they are not popular language.
- 3) Statistical data of the text are obtained. They are occurrence frequency of single letter, group vowel, and letter combination. These data will help us as a contextual data in case of OCR.

Segmentation

There is no segmentation problems in case of indo-european languages, because they have clear separator for word unit such as space or comma. There are, however, many languages in Asia which have no clear separator. They are Korean (Hangul), Chinese, Japanese, and Thai, etc. Examples shown below mean that there exist several different segmentation. Segmentation affects to the meaning of sentence and retrieval efficiency.

国立民族学博物馆
 国立 / 民族学博物馆
 国立 / 民族学 / 博物馆
 国立 / 民族 / 学 / 博物 / 馆

오 불 밤 나 무 사 온 다

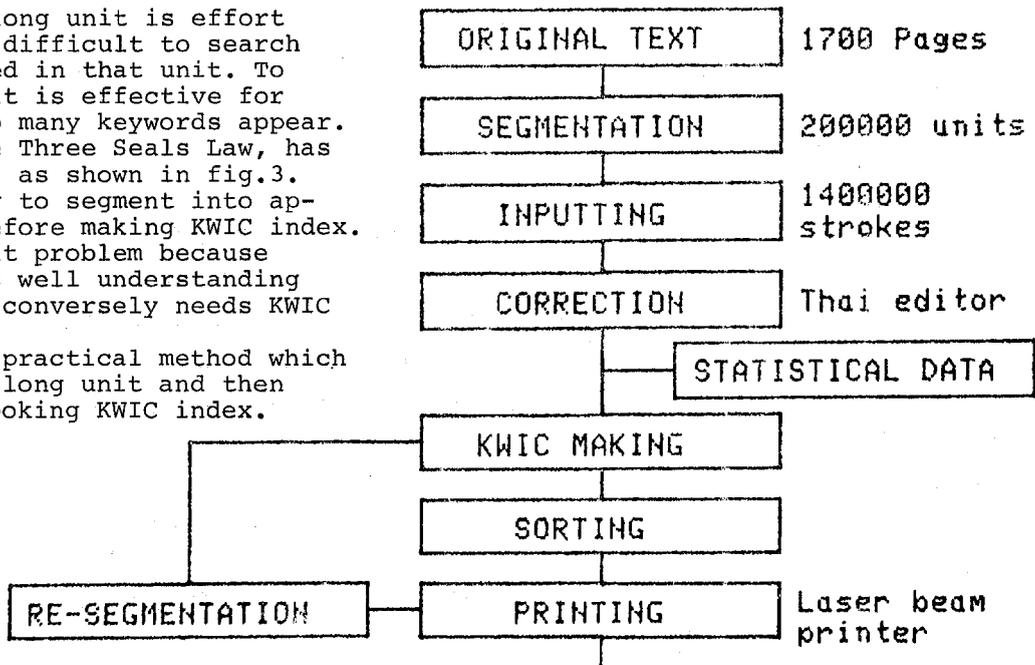
오	불	밤	나	무	사	온	다
오	불	밤	나	무	사	온	다
오	불	밤	나	무	사	온	다
오	불	밤	나	무	사	온	다
오	불	밤	나	무	사	온	다

Fig.1 Examples of different segmentations

To cut into long unit is effort saving, but it is difficult to search the string included in that unit. To cut into short unit is effective for searching, but too many keywords appear.

The text, the Three Seals Law, has no word separator, as shown in fig.3. So it is necessary to segment into appropriate units before making KWIC index. But it is difficult problem because segmentation needs well understanding of meaning, which conversely needs KWIC index.

We adopted a practical method which at first cut into long unit and then cut again after looking KWIC index.



Terminal

We use a popular intelligent graphic terminal TEKTRONIX 4051 which has usual alphabet keyboard. We stucked Thai letter labels on the side of each key as if it looks like Thai typewriter. A code table of Thai letters and corresponding english alphabets is shown in Table 1.

- The characteristics of this terminal are;
- 1) It generates Thai letter pattern by BASIC program in graphic mode. User can affirm the letter he typed.
 - 2) It has local casset memory, so that user can input and edit data anytime, even when host computer is not working.
 - 3) By way of communication line, stored data can be transmitted to host computer for time consuming work.
 - 4) It is easy to implement a flexible Thai language editor, which accept alphabet commands and display Thai letters.
 - 5) Copy of screen can be taken by the hard copy unit attached to it.

Rules for text inputting

The text has many irregular expressions. So following expediciencies are adopted.

- 1) Quotated words or phrases from Pali language are skipped by inserting special symbol to indicate there are skipped words.

Fig.2 Flow diagram of KWIC making

KWIC making

The most obvious complication is the fact that in Thai writing as many as three separate characters can appear at the same horizontal position in four different vertical positions. Therefore number of letters to take as before or after context must be carefully counted.

As a index of every unit, volume number, page number and line number are attached to the left side.

Computer algorithm

- 1) Every occurrence of pre-positioned vowel (ใ แ ไ้) is moved to a position immediately following consonant it precedes.
- 2) Diacritic symbols are moved to the end of word with the indication of position counted from the end of word.
- 3) Each letter is replaced by the code given in Table 1.
- 4) Then two words are compared as if they are numerals.

กะไล่ กะไล่ 0001' 08567146000103

งั้ง งั้ง 0002' 15571500020300

Sorting

Sorting algorithm of Thai words is not so simple as English.

We ignored algorithm 2), because our segmentation units are not necessarily words so that it does not work effectively.

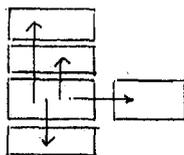
Table 1 Code table of Thai letter

Thai	ASCII	code												
ก	g	01	ข	G	21	ค	4	41	๐	[61	ด)	81
ค	H	02	ข	P	22	ค	,	42	๑	e	62	ด)	82
ค	j	03	ค	E	23	ค	p	43	๒	b	63	ด	*	83
ค	h	04	ค	D	24	ค	i	44	๓	u	64	ด	=	84
ค	U	05	ค	,	25	ค	A	45	๔	7	65			85
ค	J	06	ค	R	26	ค	l	46	๕	n	66			86
ค	H	07	ค	<	27	ค	?	47	๖	6	67			87
ค	d	08	ค	I	28	ค	;	48	๗	&	68			88
ค	!	09	ค	f	29	ค	L	49	๘	g	69			89
ค	:	10	ค	9	30	ค	K	50	๙	c	70	sp		90
ค	B	11	ค	5	31	ค	l	51	๐	F	71			91
ค	8	12	ค	m	32	ค	s	52	๑	.	72	-	2	92
ค	Y	13	ค	T	33	ค	>	53	๒	w	73	(2	93
ค	S	14	ค	o	34	ค	v	54	๓	0	74)	X	94
ค	\	15	ค	@	35	ค	U	55	๔	Q	75	,	M	95
ค	0	16	ค	x	36	ค	t	56	๕	"	76	*	↑	96
ค	C	17	ค	z	37	ค	y	57	๖	#	77	/	3	97
ค	-	18	ค	/	38	ค	<	58	๗	\$	78			98
ค	W	19	ค	r	39	ค	k	59	๘	%	79			99
ค	+	20	ค	a	40	ค	<	60	๙	<	80			

Statistical data

Total number of letters in the machine readable text is 1362602 which include special symbols such as separator, skip symbol, comma, etc. Total line number is 29582. In Table 2 is shown letter occurrence frequency for each letter. Table 3 shows occurrence frequency of compound vowels. Combination frequency of two letters are listed in Table 4. They are taken in order from the highest frequency. The combination is taken as shown below.

Fig.5 show a distribution of the ratio of upper and lower letters to the total number of letters in a line. Average ratio is 19%. A simple calculation give a ratio of 23% which is number of upper and lower letters among the horizontal positions. This means that in a line of Thai letter upper and lower letters is about 23% of normal horizontal positions.



T=total number of letters in one line

S=total number of upper and lower letters in the line

M=T-S=number of horizontal positions in the line

$$Q1=(S/T) \times 100$$

$$Q2=(S/M) \times 100$$

mean value of Q1=19%

" Q2=23%

Fig. 5

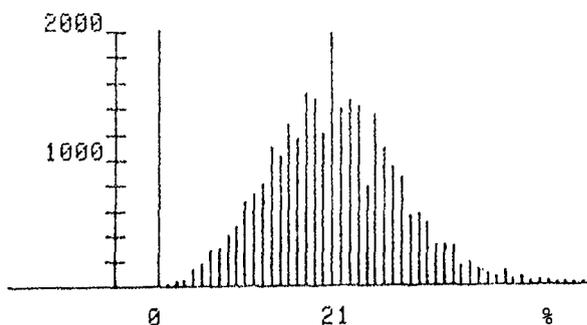


Table 2 Occurrence frequency of single letter

ก	92754	ค	27053	น	14502	ด	4233	ง	369
ข	70137	ท	22768	ก	13916	ต	3169	ช	365
ฃ	62913	ฒ	21840	ป	11407	ถ	3006	ฌ	276
ฅ	55392	า	20112	ย	10844	ท	2938	ฎ	150
ฉ	41798	น	19316	ม	10739	ช	2652	ฬ	125
ค	41407	พ	18866	ป	10018	จ	1900	ค	104
จ	39532	เ	18116	ก	8279	ร	1708	ฌ	63
ฉ	38624	โ	18049	พ	8069	พ	1493	ท	52
ช	37497	ค	17658	ก	7989	พ	1455	ฌ	28
ด	37310	แ	17549	ม	7810	ย	1186	ก	18
ด	37019	ล	17421	โ	6303	ด	865	+	17
ท	33185	ป	16903	ช	5661	ป	774	ฃ	7
ก	29376	ป	15405	ณ	4485	จ	698	ก	6
ย	28653	โ	15403	ณ	4241	ก	691	ด	5
น	27657								

Table 3 Occurrence frequency of compound vowel

- : consonant position

เ-ก	9947	-กย	2672	เ-อ	545	เ- ^๒	55	เ- ^๒ อ	0
-กย	9268	-กว	2622	เ- ^๒ ว	412	ย	49	เ- ^๒ ว	0
เ-อ	5020	เ- ^๒ ก	2134	เ- ^๒ ว	406	เ- ^๒ อ	22	เ- ^๒ ย	0
เ-ย	3617	เ- ^๒	1885	เ- ^๒ ว	339	ก	8	เ- ^๒	0
-อ	3434	เ-ย	1067	เ- ^๒ ว	235	-อ	6	เ- ^๒	0
-อ	3228	เ- ^๒ ย	1056	ก	107	เ- ^๒	2	- ^๒	0
- ^๒ ก	3085	เ- ^๒ ย	955	เ- ^๒	90	เ- ^๒ อ	0		

Table 4 Occurrence frequency of connected letters

/ : segmentation symbol, ท^๒ means ท^๒ , SP : space

น/	33472	แ	8248	น	6034	ต	4488	อ	3599
/ SP	28370	ล/	8173	ม	5943	/ค	4482	ข	3480
ก/	24672	/ท	8025	/พ	5930	น	4428	ม ^๒	3479
ง /	21972	อ	7856	น	5928	ด	4427	น	3470
/เ	18666	ร	7854	ก	5798	ล	4325	บ/	3447
ย/	14322	ว	7670	/ต	5606	ม ^๒	4286	ไป	3442
ด/	13239	/น	7584	น	5370	ท/	4162	ปร	3430
ท	11695	/พ	7498	/ท	5310	อ/	4118	/ค	3400
เ	11636	/เ	7363	/อ	5176	/ข	4094	น ^๒	3373
ร	11508	ว	7330	/จ	5086	น'	4066	ล	3368
ม/	10842	ว/	7110	พ	5016	ท	4034	ข	3350
ก/	10511	น ^๒	7021	จ	4997	ว'	3998	ม ^๒	3339
/แ	9924	/ส	6829	ก	4973	ก'	3966	แ	3339
/เ	9692	น	6756	น	4918	/,	3950	SPก	3337
ว	9305	ค	6676	ก	4860	ข	3925	บ/	3249
ย	9088	ว	6484	เป	4732	SPแ	3908	/บ	3204
น	8836	อ	6398	ท	4705	ก	3810	ล'	3179
/ม	8827	พ	6395	อ	4611	ค ^๒	3774	ท	3173
/ก	8752	พ	6299	ร ^๒	4601	ง	3663	SPต	3160
ร/	8494	น	6285	ท'	4508	เ /	3651	ม	3148

Printing

Laser beam printer

CANON LBP-3500 is a laser beam printer which can print out any kind of figure and characters. In a character mode, character must be defined as a dot matrix of 8X8,16X16,24X24,32X32,etc.

We use 16X16 matrix as a minimum module of Thai letter pattern. Thai characters are classified into fifteen types from the size of dot matrix. The largest pattern has 48X32 matrix which uses 6 modules.

One text line is printed by five horizontal zone. Each zone has 16 dot vertical width. The horizontal width of each letter can be changed character by character. But in a same zone, vertical size can not be changed.

Control of different letter width

The complex part of output program is to control the width so that heading part of KWIC index come in a line vertically.

An example of KWIC index is shown in Fig.6. We have printed about 200000 lines.

Reference

- 1) Ishii, Yoneo
1969 "Introductory remarks on the Law of Three Seals", East Asian Study, Vol.6, No.4, Kyoto University.
- 2) Murdock, George P.
1971 "Outline of cultural materials" Human Relations Area Files, Inc.
- 3) Oikawa, Akifumi & Nakayama, Kazuhiko & Sugita, Shigeharu
1979 "Printing of Thai letters by laser beam printer", the 20th annual meeting of information processing society of Japan
- 4) Sugita, Shigeharu
1979 "Computer use in ethnological studies", Bulletin of the National Museum of Ethnology, Vol.4, No.1
- 5) Udom Warotamasikkhadit & David Londe
"Computerized Alphabetization of Thai"

๔/๓๓๘/๑๓	ถ้า/ระชะบ้านทางไกลกัน/อยู่/เป็น /ทางชั่วเที่ยง	/ก็ดี/วันหนึ่ง/ก็ดี/จึง/จะ/ถึงบ้านค่างหน้า/นั้น/แล้ว
๒/๐๕๒/๐๖	กความ/ว่า/แก่/พญาน/ให้/ว่าตาม/รู้ตาม/แทน	/ก็ดี/ว่า/แต่/ตามจริง/ก็ดี/ ท่านว่า/จะ/เอา/เป็น/เจ
๒/๐๕๒/๐๗	า/พญาน/กล่าว/คำ/แก่/ผู้อ้าง/ก็ดี / ผู้คู่ความ	/ก็ดี/ว่า/มี/ผู้รู้/แทน/ผู้อ้าง/กล่าว/ข้อ/เนื้อความ/ซึ่ง/เจ
๒/๑๓๔/๑๓	ร้อมกัน/จะ/ว่าความ/ต่อกัน/โจท/ก็ดี/จำเลย	/ก็ดี/ว่า/หา/สมุด/คืน/สอ/มิ/ได้/ใช้/ ท่านว่า/ให้/เอา/
๒/๑๐๓/๑๕	๖ มาตราหนึ่ง/ วิวาท/กัน/ใน/สถาน /แห่ง/มีค	/ก็ดี/วิวาท/กัน/ใน/ทาง/สาม/เพรง/ก็ดี/ หา/สัก/จับ/พิ
๕/๓๓๘/๐๕	ง/ด้วย/วจา/ก็ดี/พะยัก/เอา/ด้วย/กาย/หน้าคา	/ก็ดี/ศีล/มุสา/ขาด/ ถ้า/แต่/คิด/ใน/ใจ/มิ/ได้/ขวล/ขว
๕/๓๓๑/๑๘	าง/ตาย/ โดย/อัน/น้อย/แต่/เรือค/แล/เสน/ตาย	/ก็ดี/ศีล/นั้น/ขาด/ แล/ศีล/ปานาคิ/ปาตฯ/จะ/ขาด/ต่
๓/๒๗๒/๐๕	กกล่าว/ให้/เป็น/คำ/นับ/ แล/เอา/ของ /ไม่/ขาย	/ก็ดี/ส่ง/ไป/ไกล/ก็ดี/ อยู่มา/เจ้า/ของ/มา/ภบ/แทน/จี
๑/๑๗๕/๐๖	ขมกัน/ให้/ตาย/ก็ดี/เป็น/กระสือ/กิน /กัน/ให้/ตาย	/ก็ดี/สรรพ/กระทำ/ให้/ตาย/ก็ดี/ หา/กัน/ว่า/จับ/ห
๓/๑๐๕/๐๓	/ใด/ลง/ก็ดี/ /ใด/เหนือ/ทรุข/ก็ดี /เหนือ/เชิง/ข้ง	/ก็ดี/สวน/แท้/แพ้ง/ริง/ใช้/ ให้/ใหม่ ๓ รอย/ๆ ละ ๒

volume page line

Fig. 6 Example of KWIC index of the Three Seals Law