

บทที่ 3

การเข้าและถอดรหัสข้อความไทยในเอกสารพีดีเอฟ

จากปัญหาของการค้นข้อความไทยในเอกสารพีดีเอฟ ที่ไม่สามารถค้นข้อความได้อย่างถูกต้อง เนื่องจากเครื่องมือการค้นในโปรแกรมแสดงเอกสารพีดีเอฟไม่เข้าใจการเข้ารหัสที่แบบอักษรไทยใช้ ในการแก้ไขปัญหา ผู้วิจัยทำการสร้างเครื่องมือการค้นข้อความในเอกสารพีดีเอฟที่เข้าใจการเข้ารหัสของชุดแบบไทยที่ใช้ในเอกสารพีดีเอฟ โดยผู้วิจัยทำการศึกษาว่าข้อความไทยในเอกสารพีดีเอฟมีการเข้ารหัสอย่างไรและควรใช้วิธีการถอดรหัสอย่างไร ในบทนี้จะกล่าวถึงการเข้ารหัสตัวอักษรไทยในเอกสารพีดีเอฟและวิธีในการที่จะถอดรหัสตัวอักษรเหล่านั้น

3.1 การเข้ารหัสตัวอักษรของแบบอักษรไทยในเอกสารพีดีเอฟ

เอกสารทั่วไปเมื่อทำการแปลงมาอยู่ในรูปแบบเอกสารพีดีเอฟ แบบอักษรที่ใช้ในเอกสารเดิมอาจจะถูกเปลี่ยนชนิดหรือเปลี่ยนการเข้ารหัสต่างไปจากข้อมูลเดิม ขึ้นอยู่กับข้อกำหนดและสภาพแวดล้อมที่ใช้ในการสร้างเอกสารพีดีเอฟ เอกสารพีดีเอฟภาษาไทยก็เช่นกัน แบบอักษรไทยที่ใช้ในเอกสารทั่วไปจะเข้ารหัสตามข้อกำหนดมาตรฐาน มอก.620 เมื่อแปลงมาอยู่ในรูปแบบเอกสารพีดีเอฟ พบว่า แบบอักษรไทยในเอกสารพีดีเอฟมีการเข้ารหัสตัวอักษรอยู่ 3 แบบด้วยกัน คือ

1. เข้ารหัสตามข้อกำหนดมาตรฐาน (Standard)
2. เข้ารหัสแบบผู้กำหนดการเข้ารหัสตัวอักษรเอง (Custom)
3. ใช้ข้อมูลการเข้ารหัสที่มีอยู่ในแฟ้มแบบอักษร (Built-in)

ข้อความไทยในเอกสารพีดีเอฟที่ใช้แบบอักษรไทย ที่มีการเข้ารหัสตัวอักษรแบบผู้กำหนดการเข้ารหัสตัวอักษรเอง และแบบใช้ข้อกำหนดการเข้ารหัสในแฟ้มแบบอักษร จะเป็นปัญหาในการที่เครื่องมือการค้นข้อความไม่สามารถค้นข้อความได้ เนื่องจากไม่เข้าใจการเข้ารหัสที่แบบอักษรไทยใช้ จึงทำให้ถอดรหัสข้อความออกมาไม่ถูกต้อง

3.1.1 การเข้ารหัสของแบบอักษรไทยแบบผู้กำหนดการเข้ารหัสตัวอักษรเอง

จากการศึกษาพบว่า แบบอักษรไทยที่มีการเข้ารหัสโดยผู้กำหนดการเข้ารหัสตัวอักษรเอง เมื่อพิจารณาตามรหัสตัวอักษรและชื่อของตัวอักษรที่ใช้ในการเข้ารหัสแล้ว จะมีอยู่ด้วยกัน 4 ลักษณะ คือ

ลักษณะที่ 1 การเข้ารหัสของแบบอักษร มีรหัสตัวอักษรตรงกับมาตรฐาน มอก.620

การเข้ารหัสในลักษณะนี้ รหัสของตัวอักษรแต่ละตัวอักษรในเอกสารพีดีเอฟมีค่าคงที่ และรหัสตัวอักษรตรงตามข้อกำหนดมาตรฐาน มอก.620 แต่ชื่อตัวอักษรในแบบอักษรเดิมเปลี่ยนไป ตัวอย่างเช่น ข้อความ "ABCDกขคกงคขกABCD" ใช้แบบอักษร AngsanaUPC ขนาด 14,18 ตามลำดับ ภายในเอกสารพีดีเอฟ จะใช้รหัสตัวอักษรและการเข้ารหัสตัวอักษร ดังนี้

```
F1(ABCD\241\242\244\247)TjTf
Encoding << 1/G01 /G02 /G03 .../GFE /GFF >>
F2(\241\242\244\247ABCD)TjTf
Encoding << 1/G01 /G02 /G03 .../GFE /GFF >>
```

จากตัวอย่างจะเห็นว่า ข้อความในเอกสาร ที่เป็นข้อความภาษาอังกฤษจะยังคงตัวอักษรเดิมที่ใช้ ส่วนข้อความไทยจะใช้เป็นรหัสตัวอักษรฐานแปด โดยรหัสตัวอักษร 'ก' คือ 241₈ รหัสตัวอักษร 'ข' คือ 242₈ รหัสตัวอักษร 'ค' คือ 244₈ และรหัสตัวอักษร 'ง' คือ 247₈

จะได้ว่า การเข้ารหัสในลักษณะนี้ รหัสตัวอักษรตรงตามข้อกำหนดมาตรฐาน มอก.620 (ดู ตารางที่ 6 ประกอบ)

ตารางที่ 6 ตารางการเข้ารหัสแบบอักษรไทยในเอกสารพีดีเอฟ ที่รหัสตัวอักษรตรงตามมาตรฐาน มอก.620

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
B	๐200	๐201	๐202	๐203	๐204	๐205	๐206	๐207	๐210	๐211	๐212	๐213	๐214	๐215	๐216	๐217
๐	๑	๒	๓	๔	๕	๖	๗	๘	๙	๐	๑	๒	๓	๔	๕	๖
9	๐220	๐221	๐222	๐223	๐224	๐225	๐226	๐227	๐240	๐241	๐242	๐243	๐244	๐245	๐246	๐247
๐	๑	๒	๓	๔	๕	๖	๗	๘	๙	๐	๑	๒	๓	๔	๕	๖
A	๐240	๐241	๐242	๐243	๐244	๐245	๐246	๐247	๐250	๐251	๐252	๐253	๐254	๐255	๐256	๐257
	ก	ข	ค	ด	ค	ข	ง	จ	ฉ	ช	ฅ	ญ	ฎ	ฏ		
B	๐260	๐261	๐262	๐263	๐264	๐265	๐266	๐267	๐270	๐271	๐272	๐273	๐274	๐275	๐276	๐277
๐	๑	๒	๓	๔	๕	๖	๗	๘	๙	๐	๑	๒	๓	๔	๕	
C	๐300	๐301	๐302	๐303	๐304	๐305	๐306	๐307	๐310	๐311	๐312	๐313	๐314	๐315	๐316	๐317
๐	๑	๒	๓	๔	๕	๖	๗	๘	๙	๐	๑	๒	๓	๔	๕	
D	๐320	๐321	๐322	๐323	๐324	๐325	๐326	๐327	๐330	๐331	๐332	๐333	๐334	๐335	๐336	๐337
๐	๑	๒	๓	๔	๕	๖	๗	๘	๙	๐	๑	๒	๓	๔	๕	
E	๐340	๐341	๐342	๐343	๐344	๐345	๐346	๐347	๐350	๐351	๐352	๐353	๐354	๐355	๐356	๐357
๐	๑	๒	๓	๔	๕	๖	๗	๘	๙	๐	๑	๒	๓	๔	๕	
F	๐360	๐361	๐362	๐363	๐364	๐365	๐366	๐367	๐370	๐371	๐372	๐373	๐374	๐375	๐376	๐377
๐	๑	๒	๓	๔	๕	๖	๗	๘	๙	๐	๑	๒	๓	๔	๕	

ลักษณะที่ 2 การเข้ารหัสของแบบอักษรมีรหัสตัวอักษรลดลงจากมาตรฐาน มอก.620 ไปด้วยค่า 35_๘

การเข้ารหัสในลักษณะนี้ รหัสตัวอักษรมีค่าคงที่ แต่เมื่อเปรียบเทียบกับมาตรฐาน มอก.620 จะลดลงไปด้วยค่า 35_๘ และชื่อของตัวอักษรในแบบอักษรเดิมเปลี่ยนไป

ตัวอย่างเช่น ข้อความ "ABCDกขคงคขกABCD" ใช้แบบอักษร AngsanaUPC ขนาด14,18 ตามลำดับ ภายในเอกสารพีดีเอฟ จะใช้รหัสตัวอักษรและการเข้ารหัสตัวอักษร ดังนี้

F1(144\45\46\47\204\205\207\212)TjTf

Encoding << 1/G36 /G37 /G38 /G39 /G132 /G133 /G135 /G138 >>

F2(212\207\205\204\44\45\46\47)TjTf

Encoding << 1/G36 /G37 /G38 /G39 /G132 /G133 /G135 /G138 >>

จะเห็นว่า รหัสตัวอักษรของตัวอักษร 'A' ตามมาตรฐานมอก.620 คือ 101_๘

การเข้ารหัสในลักษณะนี้ ตัวอักษร 'A' ใช้รหัสตัวอักษร 44_๘

จะได้ผลต่างในการเข้ารหัสใหม่คือ (101_๘ - 44_๘) = ลดลงจากเดิม 35_๘

รหัสตัวอักษรของตัวอักษร 'ก' ตามมาตรฐานมอก.620 คือ 241_๘

การเข้ารหัสในลักษณะนี้ ตัวอักษร 'ก' ใช้รหัสตัวอักษร 204_๘

จะได้ผลต่างในการเข้ารหัสใหม่คือ (241_๘ - 204_๘) = ลดลงจากเดิม 35_๘

การเข้ารหัสแบบนี้ รหัสตัวอักษรจะลดลงไปด้วยค่า 35_๘ (ดูตารางที่ 7 ประกอบ)

ตารางที่ 7 ตารางการเข้ารหัสแบบอักษรไทยในเอกสารพีดีเอฟ ที่รหัสตัวอักษรลดลงจากมาตรฐาน มอก.620

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
6	๐140	๐141	๐142	๐143	๐144	๐145	๐146	๐147	๐1๕0	๐151	๐152	๐153	๐154	๐1๕๖	๐15๗	
7	๐160	๐161	๐162	๐16๓	๐164	๐165	๐16๖	๐167	๐170	๐171	๐172	๐173	๐174	๐175	๐176	๐177
8	๐200	๐201	๐202	๐203	๐204	๐20๕	๐20๖	๐207	๐210	๐211	๐212	๐213	๐214	๐215	๐21๖	๐217
9	๐220	๐221	๐222	๐223	๐224	๐22๕	๐22๖	๐227	๐2๓0	๐2๓1	๐2๓2	๐2๓3	๐2๓4	๐2๓5	๐2๓6	๐2๓7
A	๐240	๐241	๐242	๐243	๐244	๐245	๐246	๐247	๐250	๐251	๐252	๐253	๐254	๐255	๐256	๐257
B	๐260	๐261	๐262	๐263	๐264	๐265	๐26๖	๐267	๐270	๐271	๐272	๐273	๐274	๐275	๐27๖	๐277
C	๐300	๐301	๐302	๐303	๐304	๐30๕	๐306	๐307	๐310	๐311	๐312	๐313	๐314	๐315	๐316	๐317
D	๐320	๐321	๐322	๐323	๐324	๐325	๐326	๐327	๐3๓0	๐3๓1	๐3๓2	๐3๓3	๐3๓4	๐335	๐336	๐337
E	๐340	๐341	๐342	๐343	๐344	๐34๕	๐346	๐347	๐350	๐3๕1	๐3๕2	๐3๕3	๐3๕4	๐355	๐356	๐357

ลักษณะที่ 3 การเข้ารหัสแบบอักษรมีการเข้ารหัสตัวอักษรแต่ละครั้งไม่คงที่ แต่ชื่อตัวอักษรคงที่

การเข้ารหัสในลักษณะนี้ รหัสตัวอักษรในเอกสารพีดีเอฟไม่คงที่ ขึ้นอยู่กับลำดับของตัวอักษรในเอกสารของแบบอักษรนั้น แต่ชื่อตัวอักษรในการเข้ารหัสคงที่

ตัวอย่างเช่น ข้อความ "ABCDกขคกงคขกABCD" ใช้แบบอักษร AngsanaUPC ขนาด 14, 18 ตามลำดับ ภายในเอกสารพีดีเอฟ จะใช้รหัสตัวอักษรและการเข้ารหัสตัวอักษร ดังนี้

$F1(ABCD \setminus 001 \setminus 002 \setminus 003 \setminus 004)TjTf$

Encoding << 1/afii59681 /afii59682; /afii59684; /afii59687; >>

$F2(\setminus 001 \setminus 002 \setminus 003 \setminus 004 \text{ ABCD})TjTf$

Encoding << 1/afii59687; /afii59684; /afii59682; /afii59681; >>

จากข้อมูลการเข้ารหัสข้างต้น จะได้ว่า

รหัสตัวอักษรไทยจะขึ้นอยู่กับลำดับตัวอักษรที่ปรากฏในเอกสารนั้นๆ พิจารณาที่ ตัวอักษร 'ก' เป็นอักษรตัวแรกของแบบอักษรขนาด 14 ในเอกสารพีดีเอฟ ตัวอักษร 'ก' จะใช้รหัสตัวอักษร 001 การเข้ารหัสตัวอักษร "afii59681" ในแบบอักษรขนาด 18 ตัวอักษร 'ก' เป็นตัวอักษรลำดับที่ 4 ในเอกสารพีดีเอฟ ตัวอักษร 'ก' จะใช้รหัสตัวอักษร 004 การเข้ารหัสตัวอักษร "afii59681"

จะเห็นว่า การเข้ารหัสแต่ละครั้ง รหัสตัวอักษรของ 'ก' ไม่คงที่ แต่ชื่อของตัวอักษรที่ใช้ในการเข้ารหัสคงที่ (ดูตารางที่ 8 ประกอบ)

ตารางที่ 8 ตารางชื่อตัวอักษรที่แบบอักษรไทยใช้ในเอกสารพีดีเอฟ

	0	1	2	3	4	5	6	7
128	f700 ฐ	f701 ฏ	f702 ฏ	f703 ฏ	f704 ฏ	ellipsis; ...	f705 ฏ	f706 ฏ
136	f707 ฏ	f708 ฏ	f709 ฏ	f70a ฏ	f70b ฏ	f70c ฏ	f70d ฏ	f70e ฏ
144	f70f ฏ							
152	f710 ฏ	f711 ฏ	f712 ฏ		f714 ฏ	f715 ฏ	f716 ฏ	f717 ฏ
160		afii59681 ก	afii59682 ข	afii59683 ข	afii59684 ก	afii59685 ข	afii59686 ขง	afii59687 ง
168	afii59688 จ	afii59689 ค	afii59690 ช	afii59691 ช	afii59692 ฉ	afii59693 ฉ	afii59694 ฉ	afii59695 ฉ
176	afii59696 ฉ	afii59697 ชา	afii59698 ฉง	afii59699 ฉง	afii59700 ฉ	afii59701 ฉ	afii59702 ฉ	afii59703 ก
184	afii59704 ฐ	afii59705 ณ	afii59706 บ	afii59707 ป	afii59708 ต	afii59709 ต	afii59710 พ	afii59711 ฟ
192	afii59712 ภ	afii59713 ม	afii59714 ย	afii59715 ว	afii59716 ฉ	afii59717 ฉ	afii59718 ภ	afii59719 ว
200	afii59720 ศ	afii59721 ช	afii59722 ส	afii59723 ค	afii59724 ห	afii59725 อ	afii59726 อ	afii59727 ย
208	afii59729 ค	afii59728 ค	afii59730 ก	afii59731 ก	afii59732 ค	afii59733 ค	afii59734 ค	afii59735 ค
216	afii59736 ค	afii59737 ค	afii59738 ค	afii59739 ค	afii59740 ค	afii59741 ค	afii59742 ค	afii59743 ค
224	afii59744 ค	afii59745 ค	afii59746 ค	afii59747 ค	afii59748 ค	afii59749 ค	afii59750 ค	afii59751 ค
232	afii59752 ค	afii59753 ค	afii59754 ค	afii59755 ค	afii59756 ค	afii59757 ค	afii59758 ค	afii59759 ค
240	afii59760 ค	afii59761 ค	afii59762 ค	afii59763 ค	afii59764 ค	afii59765 ค	afii59766 ค	afii59767 ค
248	afii59768 ค	afii59769 ค	afii59770 ค	afii59771 ค	f718; ค	f719; ค	f71a; ค	

ลักษณะที่ 4 การเข้ารหัสของแบบอักษรที่มีการเข้ารหัสตัวอักษรแต่ละครั้ง รหัสตัวอักษรและชื่อตัวอักษรไม่คงที่

การเข้ารหัสในลักษณะนี้ รหัสตัวอักษรในเอกสารพีดีเอฟไม่คงที่ขึ้นอยู่กับลำดับของตัวอักษรในเอกสารของแบบอักษรนั้น และชื่อตัวอักษรในการเข้ารหัสถูกเปลี่ยนไปจากชื่อที่ใช้ในแบบอักษรเดิมและชื่อที่ถูกเปลี่ยนไปนั้นไม่คงที่

ตัวอย่างเช่น ข้อความ "ABCDกขคดงคขกABCD" ใช้แบบอักษร AngsanaUPC ขนาด 14, 18 ตามลำดับ

ภายในเอกสารพีดีเอฟ จะใช้รหัสตัวอักษรและการเข้ารหัสตัวอักษร ดังนี้

F1(10011002100310041005100610071010)TjTf

Encoding << 1/G01 /G02 /G03 /G04 /G05 /G06 /G07 /G08 >>

F2(10011002100310041005100610071010)TjTf

Encoding << 1/G01 /G02 /G03 /G04 /G05 /G06 /G07 /G08 >>

จากตัวอย่างข้างต้น จะได้ว่า

ทั้งข้อความภาษาไทยและภาษาอังกฤษ รหัสตัวอักษรจะถูกเปลี่ยนไปขึ้นอยู่กับลำดับตัวอักษรที่ปรากฏในเอกสารนั้นๆ เช่น ตัวอักษร 'A' เป็นอักษรตัวแรกของแบบอักษรขนาด 14 เมื่อแปลงมาอยู่ในรูปเอกสารพีดีเอฟ ตัวอักษร 'A' จะใช้รหัสตัวอักษร 001 เข้ารหัสตัวอักษร /G01 ในแบบอักษรขนาด 18 ตัวอักษร 'A' เป็นตัวอักษรลำดับที่ 5 ในแบบอักษร ตัวอักษร 'A' ในแบบอักษรนี้ จะใช้รหัสตัวอักษร 005 เข้ารหัสตัวอักษร /G05 ตัวอักษรภาษาไทยก็เช่นกัน ตัวอักษร 'ก' เป็นอักษรลำดับที่ 5 ของแบบอักษรขนาด 14 เมื่อแปลงมาอยู่ในรูปเอกสารพีดีเอฟ ตัวอักษร 'ก' จะใช้รหัสตัวอักษร 005 เข้ารหัสตัวอักษร /G05 ในแบบอักษรขนาด 18 ตัวอักษร 'ก' เป็นตัวอักษรลำดับที่ 4 ในแบบอักษร ตัวอักษร 'ก' ในแบบอักษรนี้ จะใช้รหัสตัวอักษร 004 เข้ารหัสตัวอักษร /G04 จะเห็นว่า การเข้ารหัสแต่ละครั้งรหัสตัวอักษรและชื่อของตัวอักษรไม่คงที่

3.1.2 การเข้ารหัสแบบอักษรไทยที่ใช้ข้อกำหนดการเข้ารหัสในแฟ้มแบบอักษร

โดยปกติเอกสารพีดีเอฟ จะต้องมีข้อมูลการเข้ารหัสระบุไว้ในพจนานุกรมแบบอักษร ถ้าไม่มีข้อมูลการเข้ารหัสในพจนานุกรมแบบอักษร ให้ใช้การเข้ารหัสที่อยู่ในแฟ้มแบบอักษรนั้น แบบอักษรไทยที่มีการเข้ารหัสแบบนี้ จะใช้การเข้ารหัสตามข้อกำหนดมาตรฐานยูนิโค้ด (ดูตารางที่ 9 ประกอบ)

ตารางที่ 9 ตารางรหัสยูนิโค้ดที่แบบอักษรไทยใช้ในเอกสารพีดีเอฟ

	0	1	2	3	4	5	6	7
128	1700 จ	1701; ๑	1702; ๒	1703; ๓	1704; ๔		1705; ๕	1706; ๖
136	1707; ๗	1708; ๘	1709; ๙	170a; ๐	170b; ๑	170c; ๒	170d; ๓	170e; ๔
144	170f; ๕							
152	1710; ๖	1711; ๗	1712; ๘	1713; ๙	1714; ๐	1715; ๑	1716; ๒	1717; ๓
160		OE01 ก	OE02 ข	OE03 ฃ	OE04 ค	OE05 ด	OE06 ฉ	OE07 ง
168	OE08 จ	OE09 ฉ	OE0A ช	OE0B ซ	OE0C ฅ	OE0D ฉ	OE0E ฐ	OE0F ฎ
176	OE10 ฬ	OE11 ท	OE12 ฑ	OE13 ฒ	OE14 ด	OE15 ต	OE16 ถ	OE17 ท
184	OE18 ธ	OE19 น	OE1A บ	OE1B ป	OE1C ผ	OE1D ฝ	OE1E พ	OE1F ฟ
192	OE20 ภ	OE21 ม	OE22 ย	OE23 ร	OE24 ฤ	OE25 ล	OE26 ฬ	OE27 ว
200	OE28 ศ	OE29 ษ	OE2A ส	OE2B ห	OE2C ฬ	OE2D อ	OE2E ฮ	OE2F ย
208	OE30 ๐	OE31 ๑	OE32 ๒	OE33 ๓	OE34 ๔	OE35 ๕	OE36 ๖	OE37 ๗
216	OE38 ๘	OE39 ๙	OE3A ๐	OE3B ๑	OE3C ๒	OE3D ๓	OE3E ๔	OE3F ๕
224	OE40 ๖	OE41 ๗	OE42 ๘	OE43 ๙	OE44 ๐	OE45 ๑	OE46 ๒	OE47 ๓
232	OE48 ๔	OE49 ๕	OE4A ๖	OE4B ๗	OE4C ๘	OE4D ๙	OE4E ๐	OE4F ๑
240	OE50 ๒	OE51 ๓	OE52 ๔	OE53 ๕	OE54 ๖	OE55 ๗	OE56 ๘	OE57 ๙
248	OE58 ๐	OE59 ๑	OE5A ๒	OE5B ๓	1718, ๔	1719; ๕	171a, ๖	

3.2 การถอดรหัสข้อความไทยในเอกสารพีดีเอฟ

จากหัวข้อที่ผ่านมาได้นำเสนอลักษณะการเข้ารหัสของแบบอักษรไทย ที่เป็นปัญหาในการค้นข้อความ ทำอย่างไรจึงจะถอดรหัสข้อความไทยที่ไม่เป็นไปตามข้อกำหนดมาตรฐานเหล่านั้น ให้ได้ข้อความไทยที่ถูกต้องตรงตามข้อกำหนดมาตรฐาน มอก.620

จึงแบ่งปัญหาที่จะต้องแก้ไขในกระบวนการถอดรหัสข้อความไทยในเอกสารพีดีเอฟออกเป็นขั้นตอนดังนี้

- ทำอย่างไรจึงจะทราบว่า แบบอักษรไทยมีการเข้ารหัสลักษณะใดตามที่ได้กล่าวมา
- ทำการถอดรหัสข้อความอย่างไร จึงจะได้รหัสอักษรไทยที่ถูกต้องตามมาตรฐาน มอก.620

จากปัญหาข้างต้น ออกแบบกระบวนการถอดรหัสข้อความไทยในเอกสารพีดีเอฟ โดยให้มีขั้นตอนดังนี้

1. วิเคราะห์ลักษณะการเข้ารหัสแบบอักษรไทยในเอกสารพีดีเอฟ
2. คัดแยกข้อความในเอกสารพีดีเอฟแล้วแปลงให้เป็นรหัสที่ตรงตามมาตรฐาน มอก.620

3.2.1 การวิเคราะห์ลักษณะการเข้ารหัสแบบอักษรไทยในเอกสารพีดีเอฟ

เนื่องจาก การเข้ารหัสของแบบอักษรไทยตามที่ได้กล่าวมาในหัวข้อการเข้ารหัสไม่มีระบุไว้ในแฟ้มเอกสารพีดีเอฟ ทำอย่างไรจึงจะทราบว่าเอกสารพีดีเอฟนี้มีการเข้ารหัสตามลักษณะใด จากการศึกษาเอกสารพีดีเอฟภาษาไทยที่ใช้กันทั่วไปในระบบเครือข่ายอินเทอร์เน็ต ได้แก่ วารสารตัวอย่างเพื่อใช้ในการโฆษณาของอ.ส.ม.ท (www.mcot.or.th/v_bookworld1) จำนวน 79 เอกสาร ตัวอย่างนิตยสารขวัญเรือน (www.thaimag.com/kwanruen) จำนวน 15 เอกสาร รายงานของ นิสิตภาควิชาวิศวกรรมคอมพิวเตอร์ มหาวิทยาลัยเกษตรศาสตร์ ประกอบวิชาสัมมนา (www.cpc.ku.ac.th/~semina) จำนวน 35 เอกสาร จดหมายข่าวของศูนย์คอมพิวเตอร์ คณะวิศวกรรมศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย (www.eng.chula.ac.th) จำนวน 8 เอกสาร เอกสารของสมาคมวิศวกรรมยานยนต์ประเทศไทย (www.eng.chula.ac.th/~mech/tsae) จำนวน 10 เอกสาร บทความทางการแพทย์ของ นายสุรชัย อัญเชิญ ภาควิชาเภสัชวิทยา คณะเภสัชศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย (www.pharm.chula.ac.th/surachai) จำนวน 2 เอกสาร บทความสอนการใช้งานคอมพิวเตอร์ของ นายสรวิชัย พลสิทธิ์ ชุมนุมคอมพิวเตอร์ สถาบันเทคโนโลยีพระจอมเกล้าลาดกระบัง (www.kmitl.ac.th/~oom) จำนวน 2 เอกสาร เอกสารคู่มือการใช้งานเครือข่ายอินเทอร์เน็ตของศูนย์คอมพิวเตอร์ จุฬาลงกรณ์มหาวิทยาลัย (www.it.chula.ac.th) จำนวน 1 เอกสาร และจากผู้วิจัยดำเนินการทดลองสร้าง จะได้ว่า รายละเอียดในข้อมูลแบบอักษรและข้อมูลสภาพแวดล้อมบางส่วนในเอกสาร สามารถใช้ระบุลักษณะการเข้ารหัสของแบบอักษรไทยได้ แต่วิธีการนี้ไม่สามารถนำมาวิเคราะห์การเข้ารหัสแบบอักษรไทยในเอกสารพีดีเอฟได้ทุกเอกสาร ในบางเอกสารก็มีรายละเอียดข้อมูลไม่เพียงพอที่จะใช้วิธีนี้ทำการวิเคราะห์ จึงแบ่งกระบวนการวิเคราะห์ลักษณะการเข้ารหัสแบบอักษรไทยในเอกสารพีดีเอฟ ออกเป็น 2 ขั้นตอน คือ

ขั้นตอนที่ 1 กระบวนการวิเคราะห์การเข้ารหัสแบบอักษรไทยในเอกสารที่มีข้อมูลเพียงพอในการวิเคราะห์การเข้ารหัส

ขั้นตอนนี้อาศัย ข้อมูลชื่อแบบอักษร ประเภทแบบอักษร พจนานุกรมแบบอักษร เครื่องมือในการสร้างเอกสารพีดีเอฟ มาวิเคราะห์การเข้ารหัสแบบอักษร ได้ดังนี้ (ดูตารางที่ 10 ประกอบ)

ตารางที่ 10 ตารางข้อมูลที่ใช้ในการวิเคราะห์การเข้ารหัสชุดแบบอักษรไทยในเอกสารพีดีเอฟ

ข้อมูลที่ใช้ในการวิเคราะห์การเข้ารหัสชุดแบบอักษรไทยในเอกสารพีดีเอฟ				ลักษณะการเข้ารหัส				
ชื่อชุดแบบอักษร	ประเภทชุดแบบอักษร	สภาพแวดล้อมในการสร้าง	องค์ประกอบอื่นๆ	มอก.620	ยูนิโค้ด	ลดลงด้วยค่าคงที่	ชื่ออักษรวงศ์ที่	รูปแบบไม่คงที่
ลงท้ายด้วย "UPC"	ทรูไทย	ดิสทิลเลอร์รุ่นที่ 3		X				
ลงท้ายด้วย "UPC"	ทรูไทย	ดิสทิลเลอร์รุ่นที่ 4			X			
		ดิสทิลเลอร์รุ่นที่ 3	ที่ข้อมูลผู้ผลิต(Producer) ระบุว่าสร้างโดย ผลิตภัณฑ์ของอะโดบี เช่น Adobe PageMaker	X				
			ในตารางแสดงความกว้างลำดับตัวอักษรตัวแรกเท่ากับ 32 และตัวสุดท้ายเท่ากับ 255	X				
ขึ้นต้นด้วย "MSTT"	ทรูไทย	ดิสทิลเลอร์รุ่นที่ 3				X		
ลงท้ายด้วย "UPC"	ประเภทที่ 1	ดิสทิลเลอร์รุ่นที่ 3				X		
			ชื่ออักษรในการเข้ารหัสเริ่มต้นด้วย "afii59"				X	
ลงท้ายด้วย "UPC"	ประเภทที่ 1		รหัสตัวอักษรเริ่มต้นที่ 001					X

I 199910214

จาก ตารางที่ 10 จะได้ว่า

- เอกสารที่สร้างโดยผลิตภัณฑ์ ของบริษัท Adobe Inc เช่น PageMaker เข้ารหัสแบบผู้
ใช้กำหนดเอง จะได้การเข้ารหัสชุดแบบไทยที่มีรหัสตัวอักษรตรงมาตรฐาน มอก.620
- ตารางแสดงความกว้างในแบบอักษร ลำดับตัวอักษรตัวแรกเริ่มต้นที่ ลำดับตัวที่ 32
และลำดับตัวอักษรตัวสุดท้ายเท่ากับ 255 จะได้ การเข้ารหัสชุดแบบไทยที่มีรหัสตัว
อักษรตรงมาตรฐาน มอก.620
- ชื่อแบบอักษรคงชื่อเดิม แบบอักษรเป็นประเภททรูไทป์ สร้างโดย Distiller รุ่นที่ 4 ไม่
พบการเข้ารหัสในเอกสารพีดีเอฟ จะได้ การเข้ารหัสชุดแบบไทยที่ใช้ข้อกำหนดการเข้า
รหัสตามมาตรฐานยูนิโค้ด
- ชื่อแบบอักษรคงชื่อเดิม แบบอักษรเป็นประเภทที่ 1 เข้ารหัสแบบผู้กำหนดเอง สร้าง
โดย Distiller รุ่นที่ 3 จะได้ การเข้ารหัสชุดแบบไทยที่มีรหัสตัวอักษรลดลงจากข้อ
กำหนดในมาตรฐาน มอก.620 ไปด้วยค่า 35_๕
- ชื่อแบบอักษรขึ้นต้นด้วย MSTT แบบอักษรเป็นประเภททรูไทป์ สร้างโดย Distiller รุ่นที่
3 จะได้ การเข้ารหัสชุดแบบไทยที่มีรหัสตัวอักษรลดลงจากข้อกำหนดในมาตรฐาน
มอก.620 ไปด้วยค่า 35_๕
- ชื่อตัวอักษรในข้อมูลการเข้ารหัส ใช้ชื่อตัวอักษรที่เริ่มต้นด้วย "afii59" จะได้ การเข้า
รหัสชุดแบบไทยที่รหัสตัวอักษรไม่คงที่ แต่ชื่อตัวอักษรในข้อมูลการเข้ารหัสคงที่
- รหัสตัวอักษรตัวแรกเริ่มต้นด้วย รหัสตัวอักษร 001 และใช้ชื่ออักขระ "G01" จะได้ การ
เข้ารหัสชุดแบบไทยที่รหัสตัวอักษรและชื่อตัวอักษรในข้อมูลการเข้ารหัสไม่คงที่

ขั้นตอนที่2 กระบวนการวิเคราะห์การเข้ารหัสแบบอักษรไทยในเอกสารที่มีข้อมูลไม่เพียงพอในการ
วิเคราะห์การเข้ารหัส

จะทำการวิเคราะห์โดยให้ผู้ใช้ช่วยทำการวิเคราะห์ โดยจะนำข้อความในเอกสารพีดีเอฟมา
แสดงเป็นข้อความ 2 ข้อความ ข้อความแรกจะใช้รหัสตัวอักษรตามที่ใช้ในเอกสารพีดีเอฟ ข้อความ
ที่สองจะบวกรหัสตัวอักษรด้วยค่า 35_๕ ให้ผู้ใช้เลือกว่าสามารถอ่านข้อความใดได้ เพื่อตรวจสอบดู
ว่าเอกสารพีดีเอฟนี้ใช้การเข้ารหัสที่ตรงมาตรฐาน มอก.620 หรือลดลงจากมาตรฐาน มอก. 620

ตัวอย่างเช่น เอกสารพีดีเอฟแสดงข้อความ "สมรวย" เอกสารเข้ารหัสแบบอักษรในลักษณะที่รหัส
ตัวอักษรลดลงจากมาตรฐาน มอก.620 ด้วยค่า 35_๕

ภายในเอกสารพีดีเอฟ จะใช้รหัสตัวอักษรและการเข้ารหัสตัวอักษร ดังนี้

F1(1255\244\246\252\245)TjTf

Encoding << 1/G164 /G167 /G168 /G170 /G173 >>

เมื่อวิเคราะห์การเข้ารหัสโดยใช้สภาพแวดล้อมในเอกสารพีดีเอฟแล้วพบว่าข้อมูลในเอกสารพีดีเอฟไม่เพียงพอในการวิเคราะห์การเข้ารหัส จะนำรหัสตัวอักษรในเอกสารพีดีเอฟมาแสดงใน 2 ลักษณะ คือ

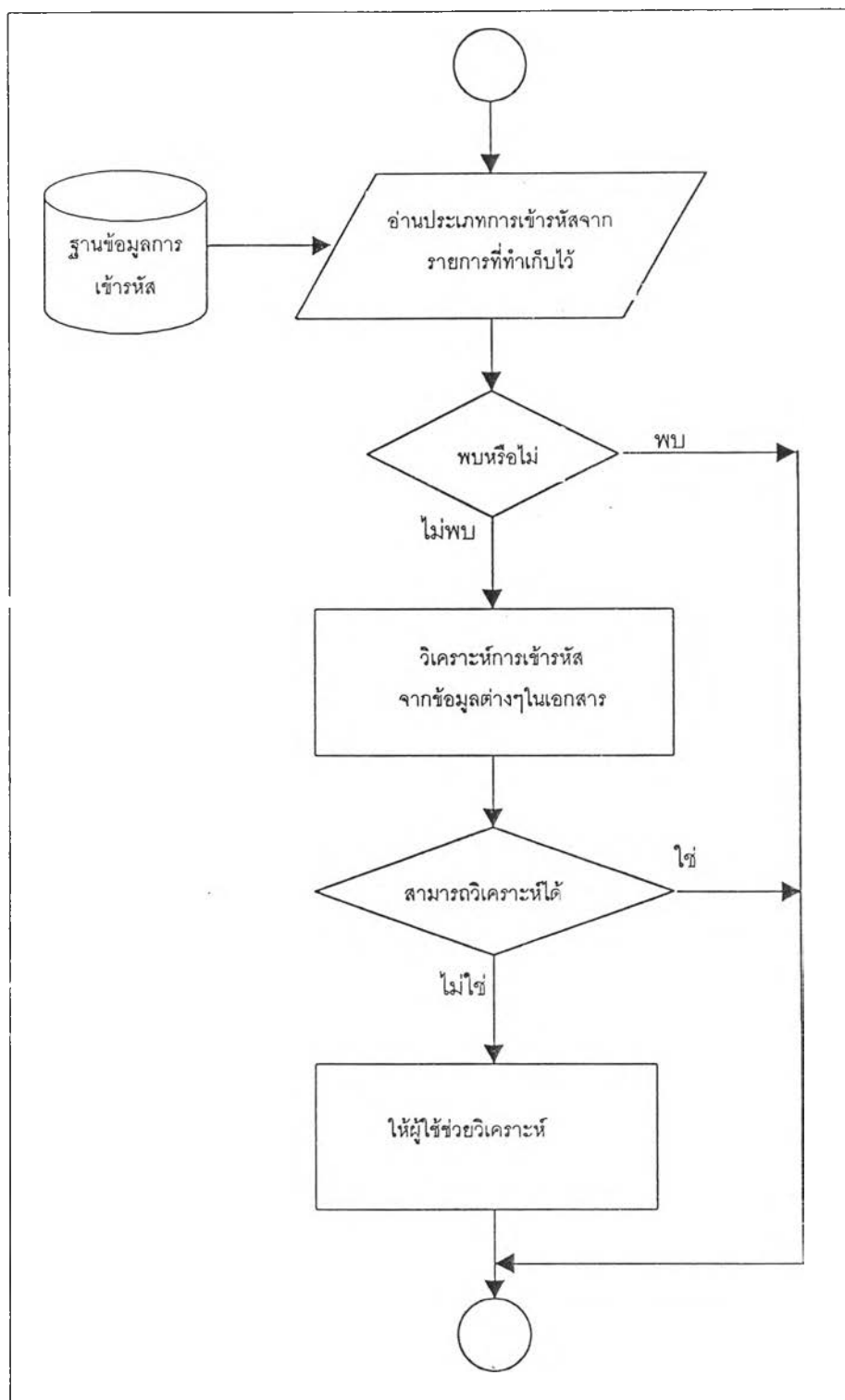
ข้อความที่ใช้รหัสตัวอักษรที่ใช้ในเอกสารพีดีเอฟ คือ " ญคจขง "

และข้อความที่บวกรหัสตัวอักษรด้วยค่า 35_{๑๖} คือ " สมรวย "

จะเห็นว่าข้อความที่สองสามารถอ่านได้ เอกสารนี้มีการเข้ารหัสแบบที่รหัสตัวอักษรลดลงจากมาตรฐาน มอก.620 ด้วยค่า 35_{๑๖}

สาเหตุที่ใช้ข้อความเพียง 2 ข้อความ ทั้งนี้เนื่องจาก การเข้ารหัสแบบอักษรที่รหัสตัวอักษรไม่คงที่แต่ชื่อตัวอักษรในการเข้ารหัสคงที่ ชื่อตัวอักษรสามารถใช้ระบุการเข้ารหัสแบบนี้ได้ ถ้าไม่พบชื่อตัวอักษรที่ใช้ตรวจสอบก็จะไม่ใช่การเข้ารหัสแบบนี้ เอกสารที่เข้ารหัสโดยใช้การเข้ารหัสที่มีอยู่ในแฟ้มแบบอักษร ตรวจสอบก็จะไม่พบข้อมูลการเข้ารหัสในเอกสารพีดีเอฟ ทำให้การเข้ารหัสที่ยังไม่ทราบแน่ชัดจะเป็น การเข้ารหัสที่รหัสตัวอักษรตรงมาตรฐาน มอก.620 การเข้ารหัสที่รหัสตัวอักษรลดลงมาตรฐาน มอก.620 และการเข้ารหัสที่รหัสตัวอักษรและชื่อตัวอักษรไม่คงที่ จึงนำข้อความในเอกสารพีดีเอฟ มาแสดงโดยการเข้ารหัสใน 2 ลักษณะ คือ ใช้รหัสตัวอักษรที่ใช้ในเอกสารพีดีเอฟ ถ้าสามารถอ่านได้จะเป็นการเข้ารหัสที่รหัสตัวอักษรตรงมาตรฐาน มอก.620 แต่ถ้าข้อความที่บวกรหัสตัวอักษรด้วยค่า 35_{๑๖} สามารถอ่านได้จะเป็นการเข้ารหัสที่รหัสตัวอักษรลดลงจากมาตรฐาน มอก.620 ด้วยค่า 35_{๑๖}

กระบวนการวิเคราะห์การเข้ารหัสแบบอักษรไทยในเอกสารพีดีเอฟ มีการทำงานดังนี้



รูปที่ 19 ผังงานการวิเคราะห์การเข้ารหัสแบบอักษรไทยในเอกสารพีดีเอฟ

การถอดรหัสข้อความไทยในเอกสารพีดีเอฟ ที่เข้ารหัสตัวอักษรในลักษณะที่รหัสตัวอักษรไม่คงที่ แต่ชื่อตัวอักษรในการเข้ารหัสคงที่

เอกสารที่มีการเข้ารหัสในลักษณะนี้ จะมีชื่อตัวอักษรคงที่ตลอดการเข้ารหัส การถอดรหัสจะใช้ชื่อตัวอักษรในการแปลงค่ารหัสตัวอักษรที่คัดแยกออกมาจากเอกสารพีดีเอฟ ไปใช้ค่ารหัสตัวอักษรที่ตรงตามมาตรฐาน มอก.620 เมื่อพิจารณาการเข้ารหัสด้วยชื่อตัวอักษรตามที่แสดงในตารางที่ 8 จะเห็นว่า ชื่อตัวอักษรที่ใช้ในแบบอักษรไทยจะแบ่งออกได้เป็น 2 กลุ่มด้วยกัน คือ กลุ่มที่ชื่อตัวอักษรขึ้นต้นด้วย "afii" และกลุ่มที่ชื่อตัวอักษรขึ้นต้นด้วย "f7"

กลุ่มที่ชื่อตัวอักษรขึ้นต้นด้วย "afii" จะเริ่มต้นที่รหัสตัวอักษร 161 ใช้ชื่อตัวอักษร "afii59681" เป็นลำดับ ไปจนถึงรหัสตัวอักษร 251 ใช้ชื่อตัวอักษร "afii59771" การถอดรหัสตัวอักษรของกลุ่มที่ชื่อตัวอักษรขึ้นต้นด้วย "afii" สามารถทำได้ดังนี้ นำค่าตัวเลข 3 ตัวหลังของชื่อ ลบด้วย 520 จะได้รหัสตัวอักษรที่ตรงตามมาตรฐาน มอก.620

กลุ่มที่ชื่อตัวอักษรขึ้นต้นด้วย "f7" จะนำค่าตัวเลขที่ต่อท้าย "f7" สร้างเป็นตารางอ้างอิงเพื่อใช้สำหรับในการแปลงค่ารหัสตัวอักษร โดยค่าตัวเลขที่ต่อท้าย "f7" จะเป็นตำแหน่งที่เก็บค่ารหัสตัวอักษร และในตารางจะเก็บค่ารหัสตัวอักษรที่ตรงตามมาตรฐาน มอก.620 เมื่อทำการคัดแยกข้อความออกมาจากเอกสารพีดีเอฟ จะดูว่ารหัสตัวอักษรนี้เข้ารหัสด้วยชื่ออักษรว่าอย่างไร ถ้าเป็นชื่อที่เริ่มต้นด้วย "f7" จะนำค่า 2 ตัวเลขที่อยู่ด้านหลัง "f7" เป็นค่าตำแหน่งของตารางอ้างอิงไปหาค่ารหัสตัวอักษรในตาราง (ดูตารางที่ 12 ประกอบ)

ตารางที่ 12 ตารางถอดข้อความไทยในเอกสารพีดีเอฟ
ที่แบบอักษรเข้ารหัสด้วยชื่อตัวอักษร "f7"

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	00	01	02	03	04	05	06	07	08	09	10	11	12	13	14	15
	176	212	213	214	215	232	233	234	235	236	232	233	234	235	236	173
1	16	17	18	19	20	21	22	23	24	25	26	27				
	209	254	231	232	233	234	235	236	216	217	254	219				

ตัวอย่างเช่น เอกสารพีดีเอฟแสดงข้อความ "กินดีมีเงิน" เข้ารหัสโดยชื่อตัวอักษรคงที่

ภายในเอกสารพีดีเอฟ จะใช้รหัสตัวอักษรและการเข้ารหัสตัวอักษร ดังนี้

F1(001\002\003\004\005\006\005\007\0\10\002\003)TjTf

Encoding << 1/afii59681 /f701; /afii59705; /afii59700; /f702; /afii59713;

/afii59744; /afii59687; >>

จะทำการถอดรหัสดังนี้ นำค่ารหัสตัวอักษร 001 ไปหาชื่อตัวอักษรในข้อมูลการเข้ารหัสของ พจนานุกรมแบบอักษร จะได้ว่า รหัสตัวอักษรเข้ารหัสด้วยชื่อ afii59681 นำค่าตัวเลข 681 ลบ ด้วย 520 จะได้รหัสตัวอักษร 161 ('ก') ตัวอักษรถัดไป 002 เข้ารหัสด้วยชื่อ f701 เจอว่าชื่อตัว อักษรคือ f7 นำค่าตัวเลขที่ตามหลัง f7 คือ 01 ไปหาค่ารหัสตัวอักษรในตาราง จะได้รหัสตัวอักษร 212 (สระอ) ทำการถอดรหัสต่อไปนี้ต่อไปเรื่อยๆจนหมดข้อความ

จากที่ทราบแล้วว่าการเข้ารหัสตัวอักษรในลักษณะนี้จะทำการเข้ารหัสตัวอักษรเริ่มต้นที่ 001 002 เป็นลำดับไปเรื่อย ทำให้มีตัวอักษรบางตัวใช้ค่ารหัสอักขระ 13 หรือ 32 ซึ่งในความหมาย ของการเข้ารหัสตามมาตรฐานทั่วไป จะเป็นรหัสตัวอักษรควบคุม ค่ารหัสอักขระ 13 หมายถึง การ ขึ้นบรรทัดใหม่และค่ารหัสอักขระ 32 หมายถึง การเว้นวรรค ทำให้เมื่อต้องเข้าไปในเอกสารพีดีเอฟ ทำการคัดแยกข้อความออกมาจากเอกสารพีดีเอฟ ตัวอักษรที่ใช้ค่ารหัสอักขระ 13 และ 32 ซึ่งมีความหมายว่า การเว้นวรรค หรือการขึ้นบรรทัดใหม่ จะไม่ถูกทำการคัดแยกออกมาด้วย จึงทำให้ การถอดรหัสข้อความในเอกสารที่มีการเข้ารหัสในลักษณะนี้ จะมีข้อจำกัดที่ไม่สามารถคัดตัว อักษรตัวที่ 13 และ ตัวที่ 32 ออกมาได้

การถอดรหัสข้อความไทยในเอกสารพีดีเอฟที่ไม่มีข้อมูลการเข้ารหัสในพจนานุกรมแบบอักษร

การถอดรหัสข้อความไทยในเอกสารพีดีเอฟที่มีลักษณะเช่นนี้ จะทำการถอดรหัสโดยคัด แยกข้อความในเอกสารพีดีเอฟออกมาเป็นรหัสยูนีโค้ด รหัสยูนีโค้ดจะใช้รหัสขนาด 2 ไบต์ เช่น OE01, OE02 เป็นต้น การถอดรหัสข้อความจะใช้ค่ารหัสไบต์หลังของรหัสยูนีโค้ดในการแปลงเป็น รหัส มอก.620 วิธีการถอดรหัสยูนีโค้ดที่เริ่มต้นด้วย E0 จะนำค่ารหัสไบต์หลังของรหัสยูนีโค้ดที่พบ แปลงจากตัวเลขฐานสิบหกเป็นตัวเลขฐานสิบ บวกด้วย 160 จะได้รหัสตัวอักษรที่ตรงตามมาตร ฐาน มอก.620 วิธีการถอดรหัสยูนีโค้ดที่เริ่มต้นด้วย f7 จะทำเช่นเดียวกับการถอดรหัสด้วยชื่อตัว อักษรที่ขึ้นต้นด้วย "f7" คือ จะนำค่ารหัสไบต์หลังของยูนีโค้ดที่เริ่มต้นด้วย f7 ไปหาค่ารหัสตัวอักษร ในตาราง

ตัวอย่าง เอกสารพีดีเอฟแสดงข้อความ "กินดีมีเงิน"

ไม่มีข้อมูลการเข้ารหัสในพจนานุกรมแบบอักษร ทำการคัดแยกข้อความจะได้รหัสตัวอักษร

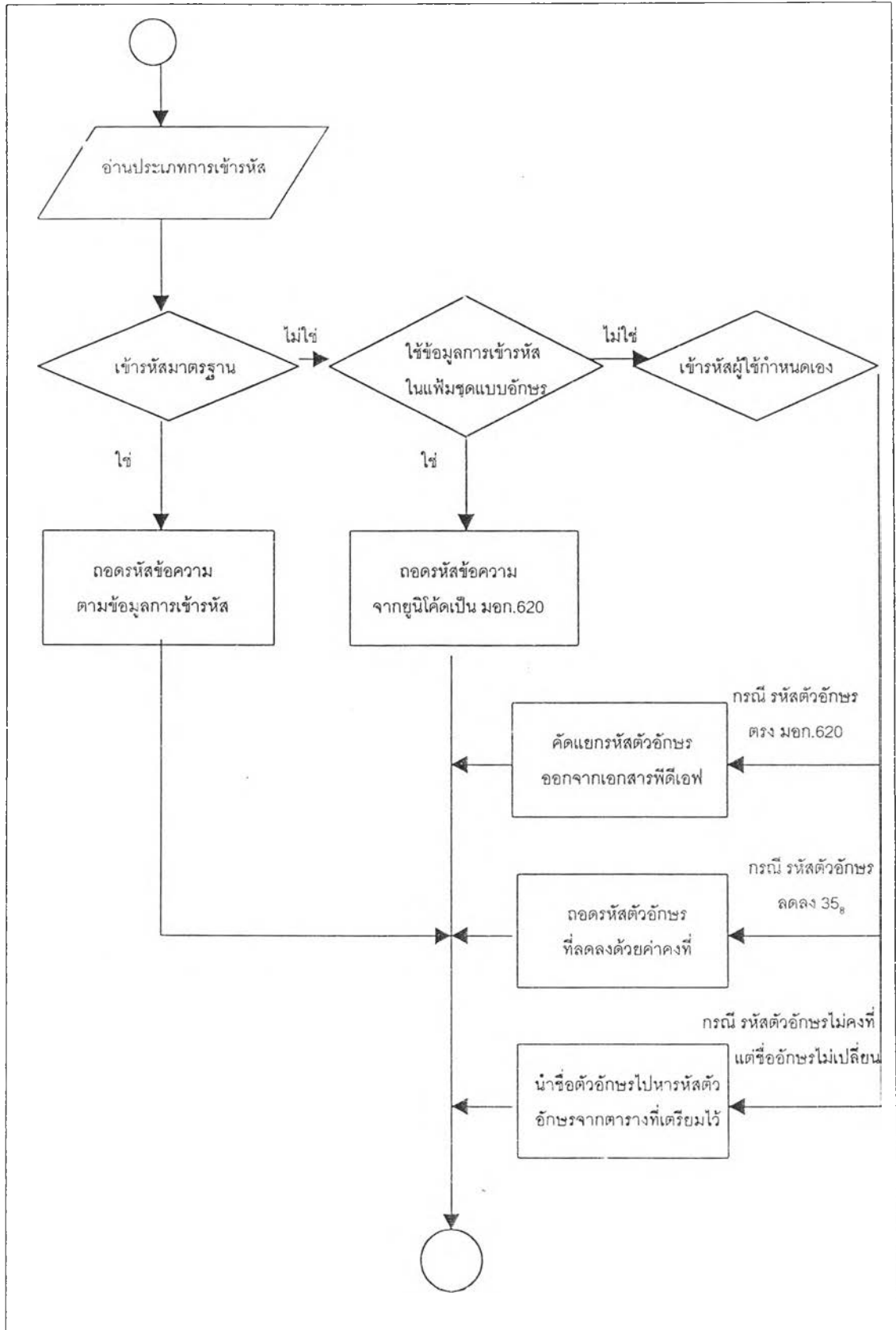
[OE01f701E0190E14f7020E21f7020E400E07f7010E19]

รหัสตัวอักษรตัวแรกคือ OE ใช้รหัสตัวอักษรไบต์ถัดไป คือ 01 บวกด้วย 160 จะได้ 161 ('ก')

รหัสตัวอักษรตัวที่ 3 คือ f7 ใช้รหัสตัวอักษรไบต์ถัดไป คือ 01 หาค่ารหัสตัวอักษรจากตารางจะได้ รหัสตัวอักษร 212 (สระอ)

รหัสตัวอักษรตัวที่ 5 คือ OE ใช้รหัสตัวอักษรไบต์ถัดไป คือ 25(19) บวกด้วย 160 จะได้ 185 ('น') ทำการถอดรหัสต่อไปนี้ต่อไปเรื่อยๆจนหมดข้อความ

กระบวนการตัดแยกข้อความแล้วแปลงให้เป็นรหัสอักขระที่ตรงตามข้อกำหนด มอก.620



รูปที่ 20 ผังงานการตัดแยกข้อความแล้วแปลงให้เป็นรหัสอักขระที่ตรงตามข้อกำหนด มอก.620