

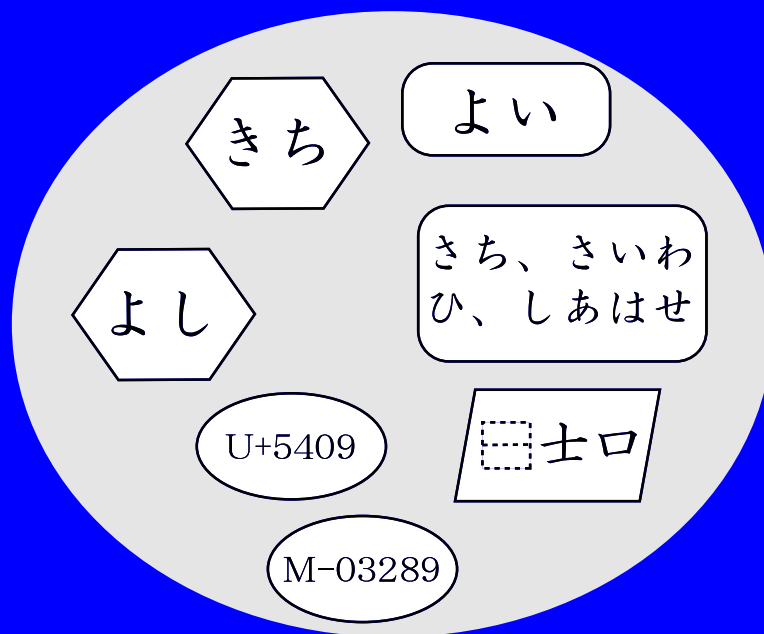
Character Processing based on Character Ontology

MORIOKA Tomohiko

2005-01-22

CHISE Project

Chaon model

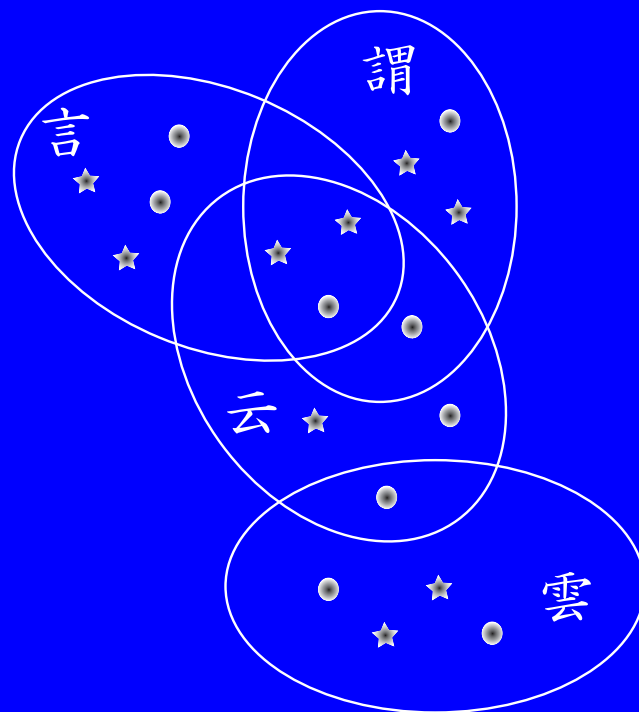


Each character is represented by its features

Comparison between character representation methods

type	easiness	expressivity
integer		×
set		
tree		
network (graph)	×	

Set operations for character features



Character operations in Chaon model

- Define/put characters • character features
- Get character features
- Search characters by character features

Implementations: by CHISE project

Fundamental library: libchise

Editing system: XEmacs CHISE

Scripting languages: Ruby/CHISE, Perl/CHISE

Multilingual Typesetter: /CHISE

Automatic Ideographic Glyph Generator: KAGE

Character databases for CHISE

1. General database distributed as a part of XEmacs CHISE (CHISE basic character database)
2. Database about structure information of Ideographs (CHISE IDS database)

CHISE basic character database

- Basic database for multi-scripts
- 'define-char' format
 - Character features are represented by Lisp (XEmacs CHISE) function 'define-char'
- About 100000 characters (in terms of character objects in CHISE database)
- UCS Unified Ideographs may be separated (now introducing multilevel unification rules)

Category of character features

1. General properties (like items in dictionaries)
2. Mapping to ID of character
 - two-way mapping (=CCS)
 - one-way mapping (=>CCS)
3. Relations between characters ($->$ *active-direction*, $<-$ *passive-direction*)

CHISE IDS database

Machine readable expression about combinations of components of Ideographs (解字)





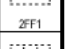









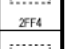









- *Ideographic Description Sequence (IDS)* of ISO/IEC 10646 is used as the format
- non-UCS characters are also used as components (if it is not found in UCS)

File Edit View Cmds Tools Options Buffers



U-00020020	甘	<input type="checkbox"/> 甘	×
U-00020021	禿	<input type="checkbox"/> 一	先
U-00020022	庇	<input type="checkbox"/> 一	此
U-00020023	雁	<input type="checkbox"/> 一	<input type="checkbox"/> 几 <input type="checkbox"/> 日 七 一
U-00020024	丛	<input type="checkbox"/> <input type="checkbox"/> 大	大 一
U-00020025	亘	<input type="checkbox"/> 一	<input type="checkbox"/> 日 一
U-00020026	州	<input type="checkbox"/> 一	州
U-00020027	巴	<input type="checkbox"/> 巴	三
U-00020028	汲	<input type="checkbox"/> 下	及
U-00020029	劊	<input type="checkbox"/> 几	丁
U-0002002A	其	<input type="checkbox"/> 甘	\
U -0002002B	其	<input type="checkbox"/> 甘	丿
U-0002002C	因	<input type="checkbox"/> 因	人
U-0002002D	考	<input type="checkbox"/> 采	丁
U-0002002E	裔	<input type="checkbox"/> 一	八 <input type="checkbox"/> <input type="checkbox"/> ×
U-0002002F	函	<input type="checkbox"/> 一	困

IDC

2FF		Ideographic description characters	
0		2FF0	 IDEOGRAPHIC DESCRIPTION CHARACTER LEFT TO RIGHT
1		2FF1	 IDEOGRAPHIC DESCRIPTION CHARACTER ABOVE TO BELOW
2		2FF2	 IDEOGRAPHIC DESCRIPTION CHARACTER LEFT TO MIDDLE AND RIGHT
3		2FF3	 IDEOGRAPHIC DESCRIPTION CHARACTER ABOVE TO MIDDLE AND BELOW
4		2FF4	 IDEOGRAPHIC DESCRIPTION CHARACTER FULL SURROUND
5		2FF5	 IDEOGRAPHIC DESCRIPTION CHARACTER SURROUND FROM ABOVE
6		2FF6	 IDEOGRAPHIC DESCRIPTION CHARACTER SURROUND FROM BELOW
7		2FF7	 IDEOGRAPHIC DESCRIPTION CHARACTER SURROUND FROM LEFT
8		2FF8	 IDEOGRAPHIC DESCRIPTION CHARACTER SURROUND FROM UPPER LEFT
9		2FF9	 IDEOGRAPHIC DESCRIPTION CHARACTER SURROUND FROM UPPER RIGHT
A		2FFA	 IDEOGRAPHIC DESCRIPTION CHARACTER SURROUND FROM LOWER LEFT
B		2FFB	 IDEOGRAPHIC DESCRIPTION CHARACTER OVERLAID

Demo

- XEmacs CHISE
- IDS

Conclusion

- Chaon model: each character is represented by a set of character features
- Chaon mode has enough expressivity and not so complex (easy to implement)
- CHISE Project are developing implementations based on Chaon model and some results are available
- CHISE character databases