

# Viewpoints on structure description of Chinese character

Morioka Tomohiko

Center for Informatics in East Asian Studies  
Institute for Research in Humanities, Kyoto University

June 18th, 2020

# Introduction

Many Chinese characters (漢字) are complex characters composed of multiple components. So we can describe their structures: e.g.

林 = 木木      雲 = 雨云      広 = 广厶

But in some cases, there are ambiguity to analyze their structures and components: e.g.

- 旗 = 方其 or 广其
- 羸 = 𠂇月女𠂇 or 羸女

# Who am I?

## Works:

- CHISE (CHaracter Information Service Environment) <http://www.chise.org/>
- Bibliography of Oriental Studies on the Web <http://ruimoku.zinbun.kyoto-u.ac.jp/>
- MeCab-Kanbun (Morpheme Analyzer for classical Chinese; Joint research) <https://corpus.kanji.zinbun.kyoto-u.ac.jp/gitlab/Kanbun/mecab-kanbun>
- etc.

# CHISE IDS database

- <https://gitlab.chise.org/CHISE/ids>
  - one of the most comprehensive IDS dataset with a large number of characters that supports almost all CJKV Unified Ideographs coded in UCS.
- CHISE character ontology
  - CHISE IDS database is a part of CHISE character ontology. Each components are defined in the ontology.
- CHISE IDS Find  
<http://www.chise.org/ids-find>
  - a Web service for searching Chinese characters that contains specified components. It is also an entrance to the CHISE character ontology.

# Structural description requirements

There are a lot of Chinese characters, so it is not easy to maintain data quality.

- Versatility: Write once, use anywhere
- Consistency
- Coverage of components: describe all Chinese characters with as few components as possible
- Intelligibility (especially for native users and classical Chinese scholars)

→ We need models

# Description based on apparent structure

Components are a visible objects

- 林 = 𣏟木木
- 雲 = 𩇛雨云
- Then, if 羸 = 𩇛亡口𣏟月女卂,  
is 𣏟月女卂 a component?

# Description based on functional structure

Component is an interface to associate phonetic and/or semantic values and shapes

→ In this view, 𠄎月女𠄎 is not a component

- If you do not know the target character, you will not know the functional components (maybe it is the goal)

# Description based on glyph design variation of component

Component is a unit to describe glyph variations of Chinese characters. cf. unification rules

- 「習」「翊」「習」: 「羽」「羽」「羽」

*If an abstract component  $\langle \text{羽} \rangle = \{ \text{羽}, \text{羽}, \text{羽} \}$  is defined, it is possible to describe abstract character  $\langle \text{習} \rangle = \text{習} \langle \text{羽} \rangle \text{白}$*



# Description based on productivity

Components are objects that combine them to create Chinese characters

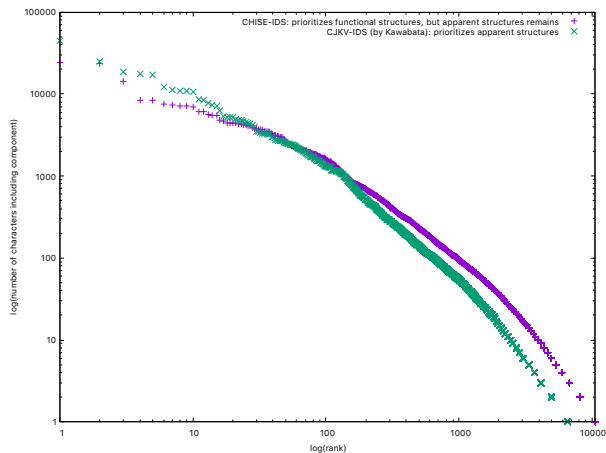
- Components that can produce many Chinese characters have high “componentness”.
- If a component is included in only one Chinese character, it is meaningless to regard it as a component (inappropriate decomposition?)
- Mechanical analysis is possible using the CHISE IDS database

# In case 羸

- 羸: 「羸」 「羸」 「羸 (羸, 羸, 羸)」 「羸」 「羸」  
「羸」 「羸」 「羸」 「羸」 「羸」 「羸」 「羸」 「羸」  
「羸」 「羸」 「羸」 「羸」 「羸」 「羸」 「羸 (羸)」  
「羸」 「羸」 「羸」 ...
- 𠃉月女卂: 「羸」 「羸」 「瀛 (羸)」 「羸」 「羸」



# Occurrence of components



This distribution seems to follow the Zipf's law

# Equivalence

In many cases, descriptions based on apparent structure and descriptions based on functional structure have equivalent information.

We can write rewriting rules: e.g.

- $\square\square\square ABC \rightarrow \square A \square BC$   
(旗 :  $\square$  旂 其  $\rightarrow$   $\square$  方 箕)
- $\square\square\square ABC \rightarrow \square\square\square ABC$   
(類 :  $\square$  須 女  $\rightarrow$   $\square\square\square$  多 女 頁)

Term Rewriting Systems (TRS) can also normalize glyph variants with unification rules.

# Ambiguity of apparent structure

虛：𠄎虍业 → 𠄎虍日七业 → 日卜𠄎厂日七业

Apparent component is also depended on knowledge.

# Conclusion

- Structural description of Chinese character should be based on Chinese character analysis (Chinese character studies), like grammatical analysis of natural language.
- It depends on knowledge, but statistical analysis for CHISE-IDS database helps discover this knowledge.
  - productivity of components
- Grapholinguistic model and algebraic model (such as Term Rewriting System) are the two wheels to describe structure of Chinese characters.