

漢字構造変換の試み

守岡 知彦 (京都大学 人文科学研究所)

多くの漢字は複数の部品の組合せからなるため、その構造の機械可読記述は漢字の検索や処理において有用な情報である。特に、なるべく音価や意味のカテゴリーを保持した部品（機能的部品）を使って構造を記述した「機能的漢字構造」は漢字のオントロジーを記述する上で有用であるが、検索においては見掛け上の部品（皮相部品）や見掛け上の漢字構造（皮相漢字構造）も必要になることがある。そのための一つの方法は機能的漢字構造と皮相漢字構造の双方を記述することであるが、もし機能的漢字構造を皮相漢字構造に機械的に変換できれば機能的漢字構造だけを記述すれば良くなりデータ作成やデータ管理のコストを削減する上で効果的であるといえる。そこで、実際に機能的漢字構造を皮相漢字構造に変換するための書き換えシステムを試作した。本稿ではその概要について述べる。

An Attempt to Convert Structural Descriptions of Chinese Characters

MORIOKA, Tomohiko (Institute for Research in Humanities, Kyoto University)

Since many Chinese characters consist of combinations of multiple components, the machine-readable description of their structure is useful information in Chinese character retrieval and processing. In particular, the “functional structure” of Chinese character, which describes the structure using “functional components” that retain as much as possible the phonetic value and/or semantic categories, is useful in describing the ontology of Chinese characters, however apparent components and apparent structures are also necessary for retrieval of Chinese characters. A simple way to deal with this is to describe both the functional structure and the apparent structure, however if we can convert the functional structure to the apparent structure automatically, we can describe only the functional structure. If this is possible, it can be said that this is effective in reducing the cost of data creation and data management. Therefore, we developed a rewriting system to convert a functional structure into the apparent structure. This paper describes the outline of the method.

1 はじめに

多くの漢字は複数の部品の組合せからなる。例えば、

林=木木 雲=雨云 広=广ム
といった具合である。このような複数の部品の組合せからなる漢字の構造（部品をどのように組み合わせるか）のことを「漢字構造」[4]といい、その記述を「漢字構造記述」という。[5] 漢字構造

記述の標準形式としては ISO/IEC 10646 [1] で定義された Ideographic Description Sequence (IDS) 形式が普及している。

上記の例では部品や漢字構造は自明だが、「旗」の場合、方其 と書くか 𠂔其 と書くかの曖昧性がある。同様に、「嬴」の場合も 彔月女 𠂔女 と書くか 嬴女 と書くかの曖昧性がある。

ここで、「旗」の部品〈𠂔〉は単体字「𠂔」の意味のカテゴリーを保持した部品（意符）であり、

部品〈其〉は単体字「其」の音価的カテゴリーを保持した部品（声符）であるのに対し、見掛け上の部品〈方〉、〈箕〉は元となった漢字とのつながりを持たない。このように、元となった漢字の音義等の情報の一部を保持した部品のことを「機能的部品」といい、機能的部品の組合せで構成される漢字構造を「機能的漢字構造」と呼ぶことにする。また、「旗」における〈方〉や〈箕〉、あるいは、「羸」における〈月女𠂔〉のように元になった漢字が存在しなかったり同形の漢字が存在したとしてもそれとの音義等の機能的なつながりが存在しないような見掛け上の部品のことを「皮相部品」と呼び、「皮相部品」の組合せで構成される見掛け上の漢字構造のことを「皮相漢字構造」と呼ぶことにする。

漢字のオントロジーを記述するという観点では機能的漢字構造を記述する方が望ましいといえるが、検索においては皮相部品も利用できる方が望ましいといえる。そのための一つの方法は機能的漢字構造と皮相漢字構造の双方を記述することであるが、もし機能的漢字構造を皮相漢字構造に機械的に変換できれば機能的漢字構造だけを記述すれば良くなりデータ作成やデータ管理のコストを削減する上で効果的であるといえる。そこで、実際に機能的漢字構造を皮相漢字構造に変換するための書き換えシステムを試作した。本稿ではその概要について述べる。

2 機能的漢字構造と部品の生産性

もし全ての漢字の字源が明らかであればその字源に従って部品や構造を確定し記述することができるだろうが、実際には現在 ISO/IEC 10646 に収録されている統合漢字の多数の漢字の『正体』を明らかにすることは必ずしも容易ではない。漢字の音義が不明なこともあり、字源知識がなくてもある程度妥当な構造記述が可能で、かつ、字源が明らかな場合には字源的漢字構造になるべく一致するような漢字構造記述が

(可能な限り機械的に) 行えることが望ましい。

こうした漢字構造記述を行う上でのヒントとなるものは部品の生産性（部品の造字力）である。^{*1}そもそも多くの漢字が部品の組合せでできているのは比較的少数^{*2}の部品で多数の漢字を表現するためだといえ、ある部品を持つ漢字が多い程部品らしさが高いといえることができる。なお、この観点から見た機能的漢字構造や機能的部品、及び、それらの確からしさの計算については別稿で述べる予定である。^{*3}

3 漢字構造の書き換え

漢字構造記述の表現形式として IDS の枠組を用いることにする。IDS は部品の組み合わせ方を示す IDC (Ideographic Description Characters) と呼ばれる特殊な文字の後に2つもしくは3つの部品を並べた形式で、Lisp の S 式と同様な前置記法の一つである。但し、IDC の後に後置される部品の数は IDC の種類ごとに決まっているために括弧なしで表現できる。IDS では部品として漢字・部品文字と IDS を用いたものが利用可能であり、部品として IDS を用いることにより入れ子構造を表現することができる。

なお、ISO/IEC 10646 での定義では漢字部

^{*1} 漢字の部品は必ず意符か声符のどちらかという訳ではなく、両者を兼ね備える場合もあれば、「月」のように複数字源のものが合流してある漢字字形を見ただけではそのどれであるかがすぐには判別できず、その形状自体に意味があるようになったものもある。[3] では漢字が部品（文字）になる際の『記号化』という表意・表音作用を薄めることで音価や字義が異なるより多くの漢字を表現可能にする作用を想定している。部品の生産性という概念はこれを漢字構造記述のデータセットというある種のコーパスに対する統計的性質という視点からいわば逆向きに見たものとみなすことができる。

^{*2} 少なくとも、部品の数は漢字の数よりも少ないはずである。実際、CHISE 文字オントロジーの場合、収録文字数約 40 万オブジェクト（字種・字体・字形等各粒度のオブジェクトの延べ数）に対して部品数は 1 万オブジェクト程度である。

^{*3} “Viewpoints on the Structural Description of Chinese Characters”; cf. <https://grafematik2020.sciencesconf.org/>

品文字として使えるのは UCS に収録された漢字と部品用文字だけであるが、原理的にはそれ以外のものを用いることも可能であり、実際、CHISE 文字オントロジー [2] では UCS に収録されていない漢字・部品や UCS の抽象文字と異なる包摂粒度・包摂範囲を表現するために UCS 非収録文字も併用している。また、IDC に関しても拡張を行っている。[6] ここではこのような拡張を行った IDS（これを「拡張 IDS」と呼ぶことにする）を対象とするが、本稿での議論は通常の IDS でも適用可能である。

このように（拡張）IDS は木構造形式の一種であり、IDC を関数記号、部品として使用される漢字・部品用文字を定数とした項と見なすことができる。そして、[5] で述べたように、項書き換え系 (Term Rewriting System; TRS) における書き換え規則を用いて（拡張）IDS の構文木（の部分木）に対する書き換えを定式化することができる。[5] では包摂規準の書き換え規則化について議論したが、本稿では機能的漢字構造を皮相漢字構造に変換するための書き換え規則を導入する。

網羅的なリストはまだできていないが、経験則に基づき次のような書き換えを導入した。^{*4}

$$\square\square LRB \rightarrow \square L \square RB \quad (111)$$

例：旗 \square 於其 \rightarrow \square 方箕

$$\square\square\square \langle i \rangle \langle | \rangle RB \rightarrow \square \langle i \rangle \square RB \quad (112)$$

例：修 \square 攸多 \rightarrow \square 𠂇 \square 攸多

$$\square\square ATR \rightarrow \square A \square TR \quad (121)$$

但し、T は垂れ。

$$T = \{ \langle 丩 \rangle, \langle 丂 \rangle, \langle 丅 \rangle, \langle 尸 \rangle, \langle 凵 \rangle, \langle 冂 \rangle, \langle 廌/廐 \rangle, \langle 艸 \rangle, \langle 巛 \rangle, \langle 巛 \rangle, \langle xy \rangle \}$$

例：巛 \square 𠂇 \rightarrow \square 𠂇

$$\square\square A\bar{T}R \rightarrow \square A \square \bar{T}R \quad (122)$$

但し、 $\bar{T} = \{ T \text{ 以外} \}$

例：寢 \square 宀 \rightarrow \square 宀 \square 𠂇

$$\square\square E \langle i \rangle R \rightarrow \square E \square \langle i \rangle R \quad (131)$$

例：屨 \square 屮 \rightarrow \square 尸 \square 𠂇

$$\square\square EAB \rightarrow \square E \square AB \quad (132)$$

但し、 $A = \{ \langle i \rangle \text{ 以外} \}$

例：曆 \square 麻日 \rightarrow \square 厂替

$$\square\square LRB \rightarrow \square\square LBR \quad (210)$$

例：穀 \square 𠂇木 \rightarrow \square 𠂇

$$\square\square LRA \rightarrow \square\square ALR \quad (310)$$

例：𠂇 \square 𠂇 \rightarrow \square 𠂇 豕生

$$\square\square ABL \rightarrow \square\square LAB \quad (320)$$

例：染 \square 𠂇 \rightarrow \square 冫木

$$\square\square LRA \rightarrow \square\square ALR \quad (330)$$

例：𠂇 \square 放白 \rightarrow \square 𠂇 文

$$\square\square AEM \rightarrow \square A \square EM \quad (411)$$

但し、E は構え。

$$E = \{ \langle 口 \rangle, \langle 互 \rangle, \langle 𠂇 \rangle, \langle 西 \rangle \}$$

例：𠂇 \square 𠂇 \rightarrow \square 𠂇

$$\square\square A \langle 凵 \rangle M \rightarrow \square A \square \langle 凵 \rangle M \quad (412)$$

例：画 \square 𠂇 凵由 \rightarrow \square 𠂇 凵由

$$\square\square \langle 凵 \rangle BM \rightarrow \square\square \langle 凵 \rangle MB \quad (413)$$

例：𠂇 \square 𠂇 日 \rightarrow \square 𠂇 日 𠂇

$$\square\square ABM \rightarrow \square A \square MB \quad (414)$$

^{*4} 一部、略記のため、項から項への書き換え規則として記述していないが、変数になっている所を定数として展開することで、変数の制約として記述している部分も項から項への書き換え規則に展開できる。

但し、 $A = \{ \langle \text{白} \rangle \text{ 以外} \}$

かつ $B = \{ E \text{ と } \langle \text{口} \rangle \text{ 以外} \}$

例：橐 𠄎 橐石 → 𠄎 𠄎 𠄎

𠄎𠄎 $A \langle \text{口} \rangle BC \rightarrow 𠄎𠄎 A 𠄎 \langle \text{口} \rangle CB$ (415)

例：藁 𠄎 藁缶 → 𠄎 藁木

𠄎𠄎 $A \langle \text{宀} \rangle BC \rightarrow 𠄎𠄎𠄎 A \langle \text{宀} \rangle CB$ (416)

例：囊 𠄎 囊缶 → 𠄎𠄎𠄎 𠄎 𠄎 𠄎

𠄎𠄎 $A \langle \text{ノ} \rangle BC \rightarrow 𠄎 A 𠄎𠄎 C \langle \text{ノ} \rangle B$ (417)

例：翬 𠄎 翬 𠄎 → 𠄎 𠄎 𠄎 𠄎 𠄎 𠄎

𠄎𠄎 $AAC \rightarrow 𠄎 ACA$ (421)

例：楸 𠄎 林矛 → 𠄎 木矛木

𠄎𠄎 $LRC \rightarrow 𠄎 L 𠄎 CR$ (422)

例：衡 𠄎 行𠄎 → 𠄎 𠄎 𠄎

𠄎𠄎𠄎 $LRBA \rightarrow 𠄎𠄎 LARB$ (511)

例：嚮 𠄎 嚮束 → 𠄎 嚮 𠄎

𠄎𠄎𠄎 $LRMBA \rightarrow 𠄎𠄎 LARMC$ (512)

例：𠄎 𠄎 𠄎 女 → 𠄎𠄎 𠄎 女 𠄎 𠄎

𠄎𠄎𠄎 $LRC A \rightarrow 𠄎 L 𠄎 ACR$ (530)

例：𠄎 𠄎 𠄎 古 → 𠄎 𠄎 𠄎 𠄎 𠄎 𠄎

𠄎𠄎 $AKC \rightarrow 𠄎 A 𠄎 KCR$ (601)

但し、 K は構え。

$K = \{ \langle \text{門} \rangle, \langle \text{鬥} \rangle \}$

例：𠄎 𠄎 𠄎 蟲 → 𠄎 𠄎 𠄎

𠄎𠄎 $A 𠄎 LRC \rightarrow 𠄎 A 𠄎 LCR$ (611)

例：𠄎 𠄎 𠄎 女 → 𠄎 𠄎 𠄎 𠄎 𠄎 𠄎

𠄎𠄎 $AM 𠄎 LRC \rightarrow 𠄎 AM 𠄎 LCR$ (612)

例：𠄎 𠄎 𠄎 女 → 𠄎 𠄎 𠄎 𠄎 𠄎 𠄎

𠄎𠄎 $AMRC \rightarrow 𠄎 A 𠄎 MCR$ (620)

例：𠄎 𠄎 𠄎 言 → 𠄎 𠄎 𠄎 𠄎 𠄎

𠄎𠄎 $NAC \rightarrow 𠄎 N 𠄎 AC$ (710)

例：𠄎 𠄎 𠄎 一 → 𠄎 𠄎 𠄎

𠄎𠄎 $L BC \rightarrow 𠄎𠄎 L CB$ (811)

例：𠄎 𠄎 𠄎 米 → 𠄎𠄎 L 米大

𠄎𠄎 $LBC \rightarrow 𠄎𠄎 LCB$ (813)

但し、 B は繞 (にょう) 以外。

$B = \{ \langle \text{厶} \rangle, \langle \text{日} \rangle, \langle \text{日} \rangle, \langle \text{心} \rangle, \langle \text{壬} \rangle, \langle \text{壬} \rangle, \langle \text{王} \rangle, \langle \text{玉} \rangle \}$

例：望 𠄎 室月 → 𠄎 𠄎 王

4 実験結果

前節で提案した書き換え規則を用いて、CHISE 漢字構造情報データベースを取り込んだ状態の CHISE 文字オントロジー中の機能的漢字構造 2231 件に対し、式 111~813 を用いて皮相漢字構造を生成した。このうち 2227 件に対しては適切な変換が行われた。

適切に変換できなかったのは次のものである：

- 111 ● 拜：𠄎 拜 𠄎 → 𠄎 𠄎 𠄎 𠄎
「𠄎 𠄎 𠄎 𠄎」となるべきであるが、そもそも「𠄎 拜 𠄎」が「拜」の機能的漢字構造として妥当かどうかという問題がある。
- 122 ● 命：𠄎 𠄎 𠄎 中、 → 𠄎 𠄎 𠄎 中、
「𠄎 𠄎 𠄎 中、」となるべきであるが、前項と同様に機能的漢字構造の妥当性の問題がある。
- 210 ● 疆：𠄎 𠄎 𠄎 土 → 𠄎 𠄎 𠄎 𠄎 𠄎
「𠄎 𠄎 𠄎 𠄎 𠄎」となるべきである。「土」が意符で「疆」が声符であり、この機能的漢字構造は妥当である。

規則番号	適用件数	成功	失敗
111	341	340	1
112	41	41	0
121	138	138	0
122	71	70	1
131	9	9	0
132	781	781	0
210	122	121	1
310	2	2	0
320	17	17	0
330	9	9	0
411	30	30	0
412	3	3	0
413	1	1	0
414	256	255	1
415	6	6	0
416	5	5	0
417	3	3	0
421	14	14	0
422	164	164	0
511	85	85	0
512	6	6	0
530	7	7	0
601	7	7	0
611	76	76	0
612	10	10	0
620	6	6	0
710	5	5	0
811	2	2	0
813	14	14	0

表1 書き換え規則の適用結果

414 ● 年: 𠂔午二 → 𠂔𠂔二十

今の所、UCS には「𠂔」が収録されていないため、UCS に基づく狭義の IDS では適切な皮相漢字構造を決定しづらい。なお、「年」は「年」の異体字であるが、この現在の字体は説文小篆の「季」

や甲骨文字の「𠂔」とは別字源のものと考えられ、「午」も「二」も純粋な記号部品*5である。即ち、これらは声符でも意符でもないが、部品の生産性に着目した場合、十分な機能性を有しているといえる。

これらの例は「弓」や「十」、「中」のように、部品の矩形領域中に小さな部品を配置可能な空間があり、その上下左右に「丶」や「二」のような小さな部品を配置した場合に食い込まれる形になることによる。特に、「十」や「中」のように複数の領域に配置可能な形状の場合、この手法では決定困難だと思われる。同様に、𠂔を含んだ複合部品が現れる場合も問題が生じ得る。このような場合、IDC を拡張するなどしてどういう交差関係になっているかを記述するためのヒントを示す必要があるといえる。なお、今回はそうした複合部品の内部構造を皮相漢字構造とすることによって問題を回避した。

5 おわりに

前節で提案した書き換え規則を用いて、CHISE 漢字構造情報データベースを取り込んだ状態の CHISE 文字オントロジー中の機能的漢字構造 2231 件に対し、式 111~813 を用いて皮相漢字構造を生成した。このうち 2227 件に対しては適切な変換が行われた。

今回は機能的漢字構造から皮相漢字構造への変換を試みたが、皮相漢字構造から機能的漢字構造への変換や、任意の漢字構造から機能的漢字構造への変換についても考える必要があるだろう。

また、機能的漢字構造から皮相漢字構造への書き換え規則に限定しても、今回提案したものはこれまでに著者が発見した書き換えパターンの内、CHISE 漢字構造情報データベースに収録された文字での適用例が存在するものである。

*5 裘錫圭のいうところの「記号字」

潜在的にはまだ未発見の書き換えパターンが存在し得る。

また、縦もしくは横に連続して複数の部品が並ぶパターンに関しては今回は取り扱わなかった。こうしたケースの場合、より身近な（≒生産性の高い）部品の組合せにする方が判りやすいといえるが、検索に際しては可能な全パターンを生成することも考慮に値するだろう。この問題に関しては今後の課題としたい。未知の書き換え規則を発見するという観点でも、分割可能な全パターンから部品の生産力の総和を最大化するようなものを見つけ出すという『機能的漢字構造解析』という問題系を提案したい。

今回提案した書き換え規則を逆に用いた場合、項書き換え系の観点では完備な書き換え系を構成可能といえ、なんらかの正規形が得られるといえるが、その書き換え結果の妥当性は保証できない。意味のある機能的漢字構造を得るためにはこの部品の生産性の観点に基づいた機能的漢字構造解析が必要になるとと思われる。

参考文献

- [1] International Organization for Standardization (ISO). *Information technology — Universal Coded Character Set (UCS)*, 2014年9月. ISO/IEC 10646:2014.
- [2] Tomohiko Morioka. Multiple-policy character annotation based on CHISE. *Journal of the Japanese Association for Digital Humanities*, Vol. 1, No. 1, pp. 86–106, 2015年11月.
- [3] 浅原達郎. 漢字の字符, 1996年12月. 戎肆庵読裘記之一.
- [4] 守岡知彦. CHISE 漢字構造情報データベース. 東洋学へのコンピューター利用 第17回研究セミナー, 全国文献・情報センター人文社会科学学術セミナーシリーズ、京都大学学術情報メディアセンター 第78回研究セミナー, pp. 93–103, 2006年3月.
- [5] 守岡知彦. 項書き換え系を用いた漢字字体の包摂規準の形式化の試み. *情報処理学会論文誌*, Vol. 59, No. 2, pp. 332–340, 2018年2月.
- [6] 守岡知彦. 内容アドレッシングを用いた多粒度漢字構造情報表現の試み. *情報処理学会論文誌*, Vol. 61, No. 2, pp. 171–178, 2020年2月.