

「東アジア古典文献コーパスの研究」共同研究班報告

1 はじめに

2008年4月から2013年3月にかけて、われわれは「東アジア古典文献コーパスの研究」共同研究班(班長: 安岡孝一)を組織し、漢文解析に向けたコーパスの研究をおこなった。この共同研究班は、2008年3月に終了した「漢字情報学の構築」班[1]での、白文自動「点」打ちプロジェクトでの知見を受け、漢文の自動解析に特化しておこなわれた研究班である。

2 韻文の自動解析に向けて

漢文の自動解析に際し、先の「漢字情報学の構築」班の知見として、韻文と散文では全く文章構造が異なっており、それぞれに異なるアプローチを要することが、明らかになっていた。

韻文と散文の最大の違いは、もちろん、韻文が「韻」を踏んでいるという点である。すなわち韻文は、八文字、十文字、あるいは十二文字おきに脚韻を踏む、という形が一般的であり、しかもその脚韻も『廣韻』等に基づいている。このような条件から、白文の中から韻文を自動的に探し出して、文の単位で切り出すこと自体は、ほぼ100%の精度でおこなうことが可能となった[1]。

ただ、韻文では「韻」を優先するために、語順が崩壊してしまっており、一般的な意味での形態素解析は難しいというのが、われわれの感触だった。少なくとも、散文と同一の形態素解析手法は、韻文には適用できない。したがって、韻文の自動解析においては、散文とは異なるアプローチが必要だという結論になったのである。

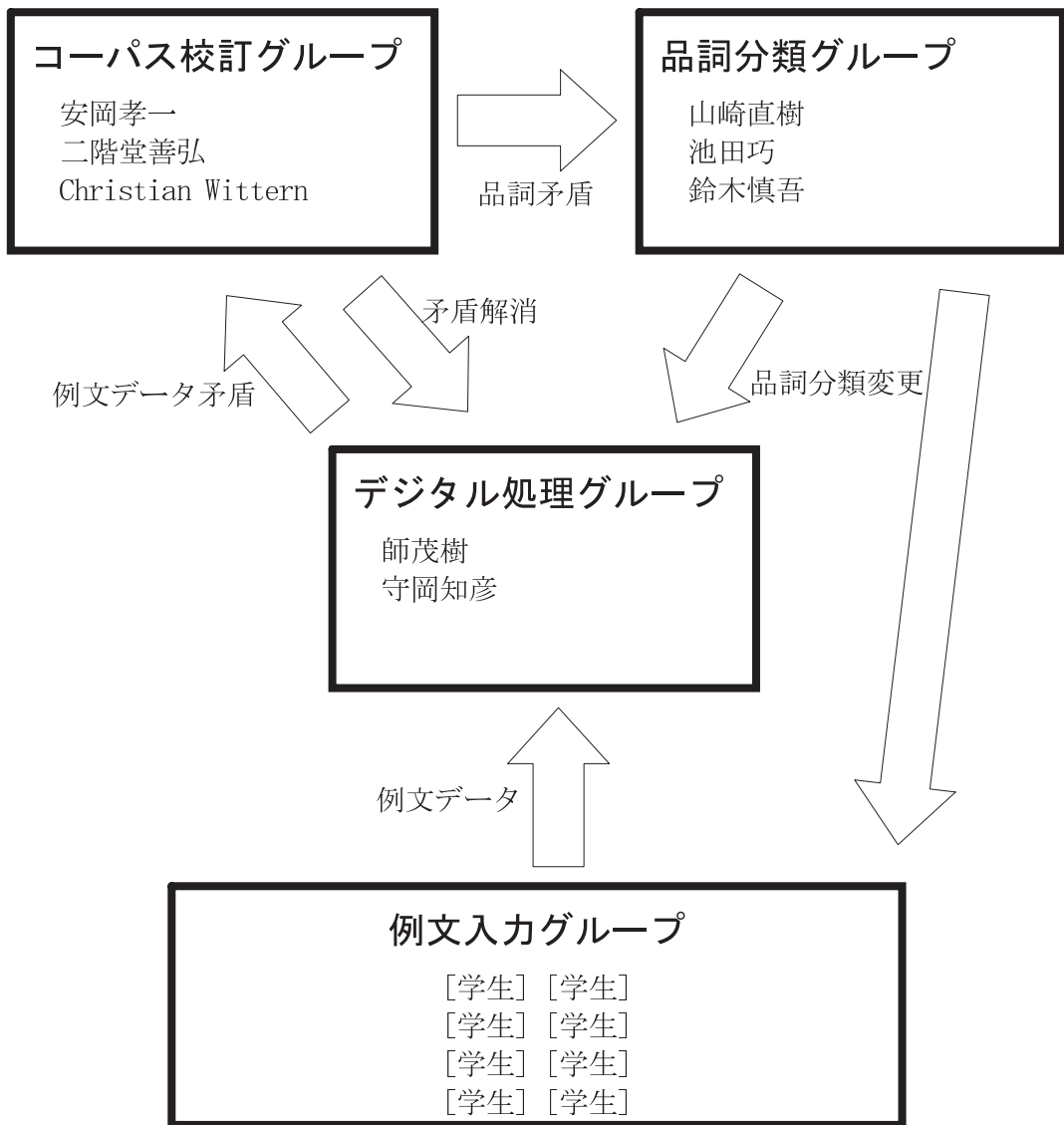
韻文の自動解析において、われわれは、対句構造に着目したアプローチを採ることにした。

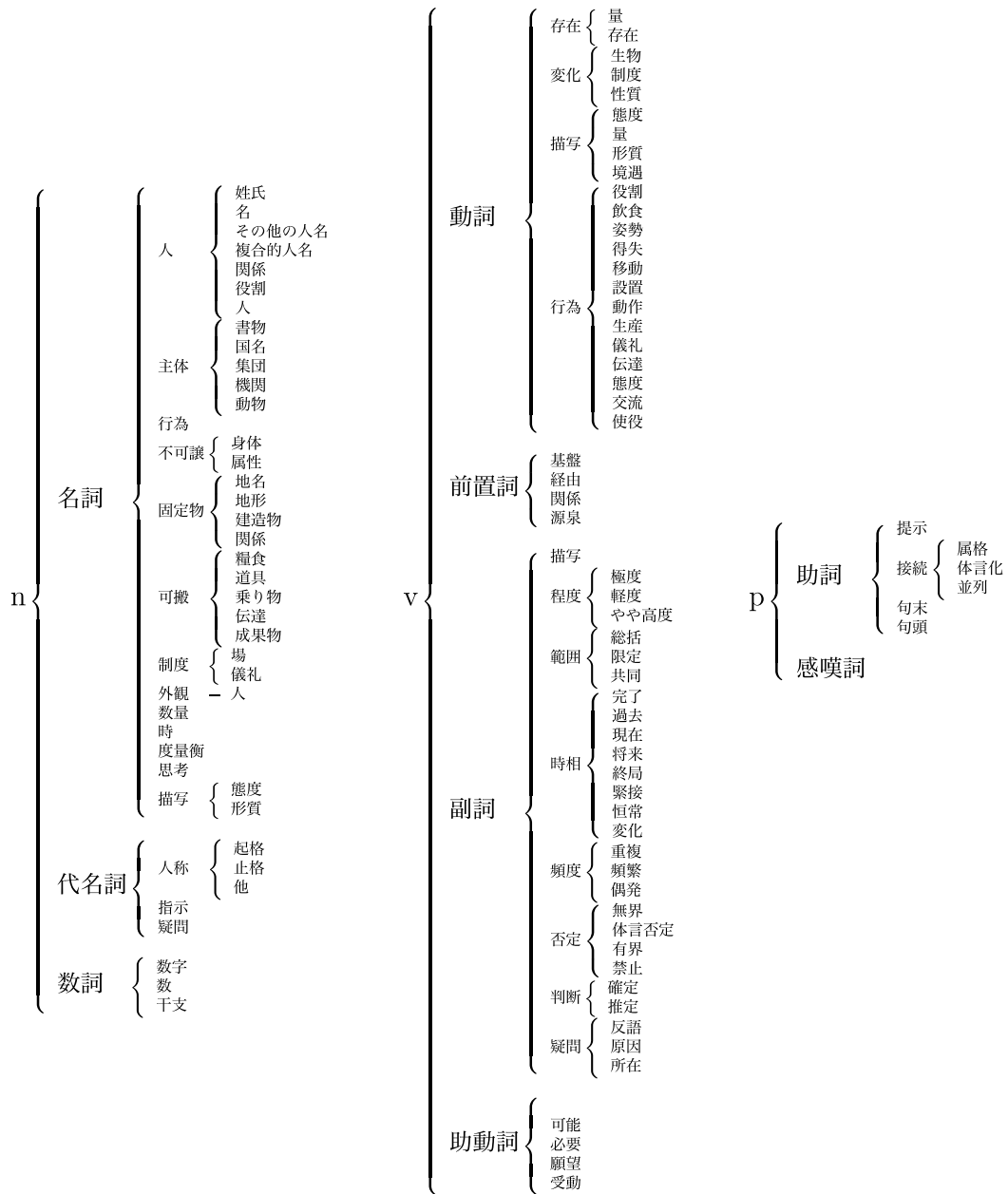
3 散文の自動解析に向けて

散文の自動解析において、われわれは、MeCabというソフトウェアを用いることにした[2]。MeCabはオープンソースの形態素解析エンジンで、言語、辞書、コーパスに依存しない汎用的な設計がなされており、辞書とコーパスを準備すればいかなる言語にも対応できる、というのが売りだった。ならば、漢文(の散文)にもMeCabを使用できるはずだ、というのが、われわれの直感だったが、われわれ以前には誰もそれを試したことがなかった。

MeCabの辞書には4階層の「品詞」が必要なことから、われわれは、日本語と漢文を繋ぐ「構造」の一種である訓読に着目し、返り点を「品詞」に反映させることを考えた。すなわち、訓読における返り点を、漢文の動賓構造を表しているものとみなし、動語に「v」という「品詞」を、賓語に「n」という「品詞」を、その他の語に「p」という「品詞」を、それぞれ、MeCab漢文辞書の「第1階層の品詞」(以下「大品詞」と呼ぶ)として定めることにしたのである。次に「第2階層の品詞」(以下「品詞」と呼ぶ)だが、これはIPAの日本語辞書から、デッチあげてみることにした[3]。「第3階層の品詞」(以下「意味索性」と呼ぶ)と「第4階層の品詞」(以下「小索性」と呼ぶ)に関しては、初期段階では付与しないことにしてみた。

このMeCab漢文辞書(IPA由来版)と、それに基づいて作ったMeCab漢文コーパスを用いて、高校教科書の漢文例や、三国志呉書列伝などの白文を、MeCabで形態素解析してみた。





4 おわりに

新規に「東アジア古典文献コーパスの応用研究」共同研究班を立ち上げる予定である。

参考文献

- [1] 「漢字情報学の構築」共同研究班報告, 東方學報, 第83冊 (2008年9月), pp.360-349.
- [2] 守岡知彦: MeCabを用いた古典中国語の形態素解析の試み, 情報処理学会研究報告, Vol.2008-CH-79 (2008年7月), pp.17-22.
- [3] 守岡知彦: MeCabを用いた古典中国語形態素解析器の改良, 情報処理学会研究報告, Vol.2009-CH-84 (2009年10月), No.3, pp.1-5.
- [4] Koichi Yasuoka: Toward a Syntactic Analysis of Classical Chinese Texts, Osaka Symposium on Digital Humanities 2011 (September 2011), p.34.
- [5] Naoki Yamazaki: Toward Syntactic Frame Retrieval of Classical Chinese Rhymes using Japanese 'kun' readings and Syntactic parallelism of couplets, Osaka Symposium on Digital Humanities 2011 (September 2011), p.35.
- [6] Tomohiko Morioka: A Prototype of a Classical Chinese Morphological Analyzer based on MeCab, Osaka Symposium on Digital Humanities 2011 (September 2011), p.36.
- [7] 守岡知彦: 古典中国語形態素コーパス編集システムの開発, 東洋学へのコンピュータ利用, 第23回研究セミナー (2012年3月), pp.75-83.
- [8] 山崎直樹, 守岡知彦, 安岡孝一: 古典中国語形態素解析のための品詞体系再構築, 人文科学とコンピュータシンポジウム「じんもんこん2012」論文集 (2012年11月), pp.39-46.