

## 漢文の形態素解析・依存文法解析・ 直接構成鎖解析

安 岡 孝 一

筆者が班長を務める京都大學人文科學研究所共同研究班「東アジア古典文獻コーパスの實證研究」では、現在、古典中國語（漢文）における文法解析の自動化に全力で取り組んでいる。漢文の白文に對し、形態素解析・依存文法解析・直接構成鎖解析を順におこなうことで、白文の統語構造が解析可能となる、というのが、われわれの見通しである。形態素解析<sup>[1]</sup>によって、單語切りをおこなうと同時に、各單語の品詞を得る。依存文法解析<sup>[2]</sup>によって、單語と單語の間の係り受け關係を解析すると同時に、文の切れ目を得る。直接構成鎖解析<sup>[3]</sup>によって、各文の統語構造を解析木の形で得る。例として「孟子見梁惠王王曰叟不遠千里而來」という白文に對し、形態素解析・依存文法解析・直接構成鎖解析を順におこなう際の流れ（イメージ圖）を、圖 1 に示す。

なお、本研究は、科學研究費補助金基盤研究 (B)17H01835『古典漢文形態素コーパスにもとづく動詞の作用域の自動抽出』の研究助成を受けている。

### 漢文の形態素解析

漢文の形態素解析において、われわれは、MeCab という汎用の形態素解析ソフトウェアを用いている。MeCab は、もともとは日本語の形態素解析用だった<sup>[1]</sup>が、言語、辭書、コーパスに依存しない汎用的な設計がなされており、辭書とコーパスを準備すればいかなる言語にも對應できる。

MeCab の辭書には「品詞」（複数の階層が可能）が必要なことから、われわれは、日本語と漢文をつなぐ「構造」の一種である訓讀に着目し、返り點を「品詞」に反映させることを考えた。すなわち、訓讀における返り點が、漢文の動賓構造を表しているのみならず、動詞類に「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 漢文辭書の「第 1 階層の品詞」として定めた。次に「第 2 階層の品詞」を「名詞」「代名詞」「數詞」「動詞」「前置詞」「副詞」「助動詞」

孟子見梁惠王王曰叟不遠千里而來

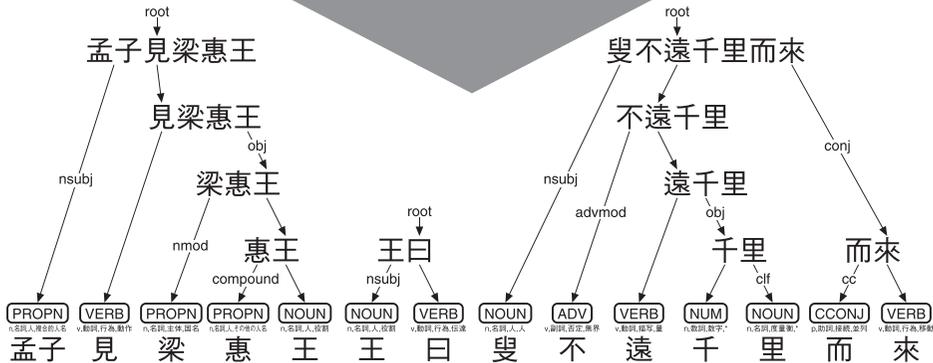
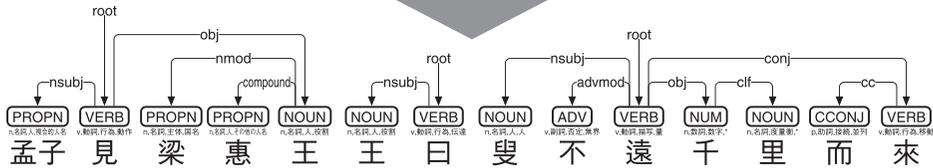
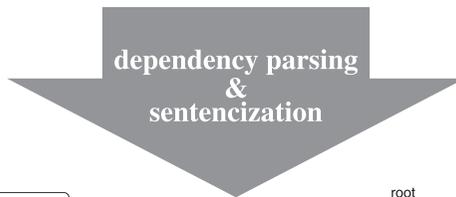


圖1 漢文の形態素解析・依存文法解析・直接構成鎖解析

「助詞」「感嘆詞」「接尾辭」の10種類とした。従來の漢文文法などで見られた「形容詞」を廢止して、「動詞」と統合している<sup>[4]</sup>のが特徴である。さらに「第3階層の品詞」として44種類の意味素性を、「第4階層の品詞」として88種類の小素性を定義し、形態素解析の結果として得られる各単語を、意味の面からも捉えやすいよう工夫した<sup>[5]</sup>。

MeCabを用いた形態素解析において、その中心となるアイデアは、CRF (Conditional

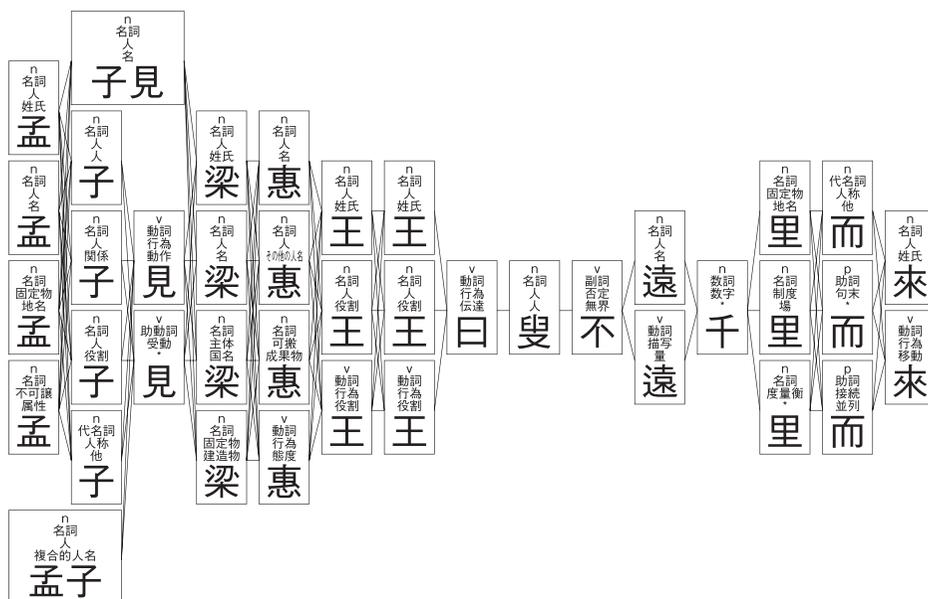


圖2 CRF を用いた漢文の形態素解析

Random Fields) である。われわれの漢文形態素解析に即して言えば、解析したい白文を MeCab 漢文辞書に基いて、可能性のある全ての単語（4階層の品詞を含む）の組み合わせの列に変換する。例として「孟子見梁惠王曰叟不遠千里而來」という白文に對する CRF（の一部）を、圖2に示す。「孟子」には「n, 名詞, 人, 複合的人名」という「品詞」が與えられており、その一方で、「孟」には4種類の「品詞」が、「子」には4種類の「品詞」が與えられうることから、これらの組み合わせが全て列挙されている。さらに「子見」に「n, 名詞, 人, 名」という「品詞」が與えられていて、「見」は「v, 動詞, 行為, 動作」と「v, 助動詞, 受動, \*」の可能性がある。これらの組み合わせに對し、MeCabは、各単語の出現確率と、隣り合う単語どうしの共起確率から、全ての組み合わせの中で最も確率が高くなるような単語列を抽出する。

抽出した単語列に對し、Universal Dependencies 向け品詞（表1）とグロス（近似的な逐語英譯）の付與を並行しておこない、形態素解析を完成する。グロスは基本的に<sup>[6]</sup>に依っているが、一人稱・二人稱・三人稱は [1PRON]・[2PRON]・[3PRON] で、受動の助動詞は [PASS] で、完了の助詞は [PFV] で、疑問の助詞は [Q] で表した。例として「孟子見梁惠王曰叟不遠千里而來」に對する形態素解析の結果を、圖3に示す。

表1 MeCab 向け漢文品詞體系と Universal Dependencies 向け品詞

大品詞, 品詞	UD 向け品詞	特 例
n, 名 詞	NOUN	n, 名詞, 人, 姓氏 n, 名詞, 人, 名 n, 名詞, 人, その他の人名 n, 名詞, 人, 複合的人名 n, 名詞, 主體, 國名 n, 名詞, 固定物, 地名 } PROPEN
n, 代名詞	PRON	
n, 數 詞	NUM	
v, 動 詞	VERB	
v, 前置詞	ADP	
v, 副 詞	ADV	
v, 助動詞	AUX	
p, 助 詞	PART	p, 助詞, 接續, 屬格 p, 助詞, 接續, 並列 } SCONJ } CCONJ
p, 感嘆詞	INTJ	
p, 接尾辭	PART	
s, 記 號	PUNCT	
s, 文 字	SYM	



圖3 「孟子見梁惠王曰叟不遠千里而來」の形態素解析結果

## 漢文の依存文法解析

漢文の依存文法解析<sup>[7]</sup>において、われわれは、Universal Dependencies<sup>[8]</sup>（以下「UD」）という、言語横断的な依存構造記述を用いている。UDは、品詞・形態素属性・依存構造情報を、言語に依存せず記述する手法である。句構造を考慮せずに係り受け関係を記述できるよう、全ての構文構造を「単語」間の依存関係で記述するのが特徴である。

「単語」間の係り受け関係に対しては、UD依存構造の「単語」間リンクを用いて表現し、各リンクに表2のUD依存構造タグ39種類を付与している。タグのうち32種類は、もともとUDで規定されているものであり、7種類（nsubj: pass・csubj: pass・obl: tmod・obl: lmod・discourse: sp・compound: redup・flat: vv）は、その派生形である。rootはリンク元を持たないが、他のタグによるリンクは、リンク元の「単語」とリンク先の「単語」を1つずつ有する。たとえば、漢文の動賓構造は、動詞をリンク元、賓語をリンク先、とするobjというリンクで表現する。リンクの本数は「単語」の個数に等しく、各リンクのリンク先は、全て互いに異なっている。すなわち、各「単語」から出るリンクは複数

表2 古典中國語（漢文）Universal Dependencies の依存構造タグ

	Nominals	Clauses	Modifier Words	Function Words
<b>Core arguments</b>	nsubj 主語 ↔nsubj: pass [受動文] obj 目的語 iobj 間接目的語	csubj 節主語 ↔csubj: pass [受動文] ccomp 節目的語 xcomp 節補語		
<b>Non-core dependents</b>	obl 斜格補語 ↔obl: tmod [時] ↔obl: lmod [場所] vocative 呼稱語 expl 形式語 dislocated 外置語	advcl 連用修飾節	advmod 連用修飾語 discourse 談話要素 ↔discourse: sp [文助詞]	aux 動詞補助成分 cop 繫辭 (copula) mark 標識 (marker)
<b>Nominal dependents</b>	nmod 體言による連體修飾語 nummod 數量による修飾語	acl 連體修飾節	amod 用言による連體修飾語	det 決定詞 clf 類別詞 case 格表示
<b>Coordination</b>	<b>MWE</b>	<b>Loose</b>	<b>Special</b>	<b>Other</b>
conj 接續 cc 接續詞	fixed 固着 compound 複合 (endocentric) ↔compound: redup [重疊] flat 並列 (exocentric) ↔flat: vv [動詞類]	list 細目 parataxis 隣接表現	orphan 親なし	root 親

ありうるが、各「単語」に入るリンクは1つだけである。また、リンクはループしない。さらに、われわれの古典中國語（漢文）UDでは、リンクどうしが交差しない、rootをまたぐリンクも存在しない、という制限も設けている。

依存文法解析のための手法は、これまでに数多く提案されているが、われわれの古典中國語（漢文）UDのように、複数のrootを持ち (dependency forest)、UD依存構造のリンクどうしが交差せず (planar)、rootをまたぐリンクがない (projective)、という条件においては、arc-planar<sup>[9]</sup>という (非決定性) アルゴリズムが、有効だと考えられる。arc-planarは、単語列の先頭から末尾に向かって「垣根」(stack-buffer boundary)を移動していく、というイメージで処理をおこなう。「垣根」がおこなう遷移は、**Shift・Reduce・Left-Arc・Right-Arc**の4種類である。

- ・ **Shift** 「垣根」を右に1単語分、移動する。
- ・ **Reduce** 「垣根」のすぐ左の単語を除去して、解析結果へ移す。
- ・ **Left-Arc** 「垣根」のすぐ右の単語から、すぐ左の単語へリンクを繋ぐ。
- ・ **Right-Arc** 「垣根」のすぐ左の単語から、すぐ右の単語へリンクを繋ぐ。

単語が全て **Reduce** されて、「垣根」がポツンと取り残された時点で、arc-planarは終了である。arc-planarによる「孟子見梁惠王王曰叟不遠千里而來」の依存文法解析の様子を、圖4に示す。あとは、リンクが入っていない「見」「曰」「遠」にrootを刺すこと

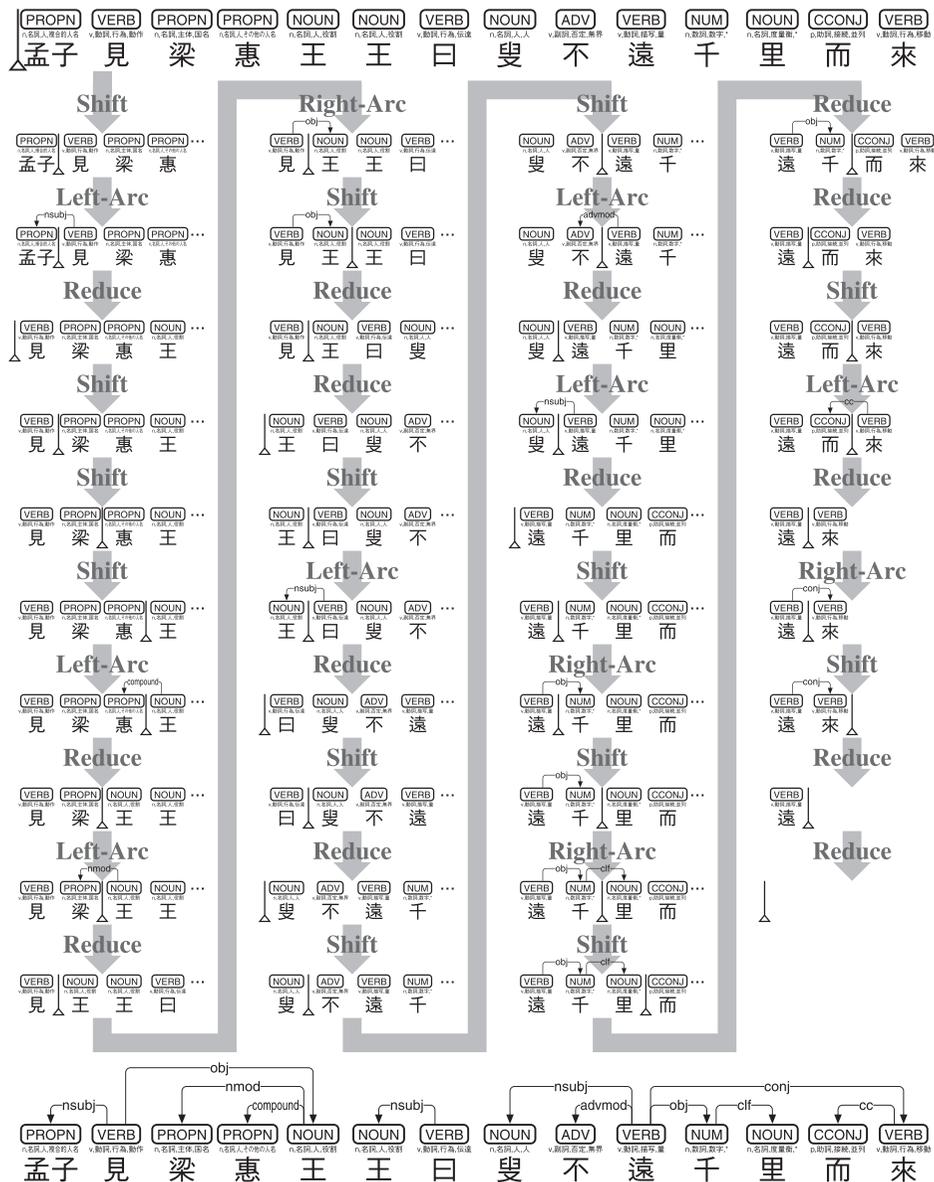


圖4 arc-planarによる漢文の依存文法解析

で、圖1の依存文法解析結果が得られるわけである。

ただし、arc-planarにおける「垣根」の遷移は、実際には非決定的\*である。圖4では

\* 「垣根」が単語列の先頭にある局面では **Shift** を、末尾にある局面では **Reduce** をおこなうしかなく、これらの場合だけは決定的である。また、**Left-Arc** の直後を **Reduce** に、

解析過程を一本道で示したが、現実には、各局面において複数の可能性が、枝分かれとして存在する。これら複数の可能性については、それぞれの遷移を選択した場合を、確率的に並行して解析することになる。われわれが解析に用いている UDPipe<sup>[10]</sup>においても、ほぼ、そのような形で実装がおこなわれている。

## 漢文の直接構成鎖解析

直接構成鎖解析 (immediate catena analysis)<sup>[3]</sup> は、構成素 (constituent)<sup>[11]</sup> による古典的な文法解析を、構成鎖 (catena)<sup>[12]</sup> を用いて依存文法へと擴張する手法である。われわれは、漢文の依存文法記述に UD を用いていることから、直接構成鎖解析においても、UD を擴張する形で解析手法の開発をおこなった。

構成鎖は、依存文法の文法木における連結な部分グラフである。たとえば、圖 5 左上の UD 「孟子見梁惠王」は、以下に示す 17 種類の構成鎖を含んでいる。

「孟子見梁惠王」「孟子見梁王」「孟子見惠王」「見梁惠王」「孟子見王」  
 「見梁王」「見惠王」「梁惠王」「孟子見」「見王」「梁王」「惠王」「孟子」  
 「見」「梁」「惠」「王」

直接構成鎖解析における基本コンセプトは、リンクの除去である。最初に root リンクの除去をおこない、そこから文法木を幅優先でリンクの除去をおこない、最終的には全てのリンクを除去する。圖 5 左を見てもよい。最初に root リンクを除去し、極大構成鎖「孟子見梁惠王」を得る。次に nsubj リンクを除去すると、「孟子見梁惠王」は、「孟子」と「見梁惠王」の 2 つの構成鎖に分割される。次に obj リンクを除去すると、「見梁惠王」は、「見」と「梁惠王」の 2 つの構成鎖に分割される。次に nmod リンクを除去すると、「梁惠王」は、「梁」「惠王」の 2 つの構成鎖に分割される。最後に compound リンクを除去すると、「惠王」は、「惠」と「王」の 2 つの構成鎖に分割される。この分割の様子を、再度、構文木として組み上げたのが圖 5 右であり、これが「孟子見梁惠王」に対する直接構成鎖解析の構文木である。

同様に、「王曰」「叟不遠千里而來」に対する直接構成鎖解析の構文木を、圖 1 下に示す。ただし、これらの直接構成鎖解析を適切におこなうためには、適切な順序でリンク

Right-Arc の直後を Shift に、それぞれ決め打ちすることは (局面に応じて) 可能である。

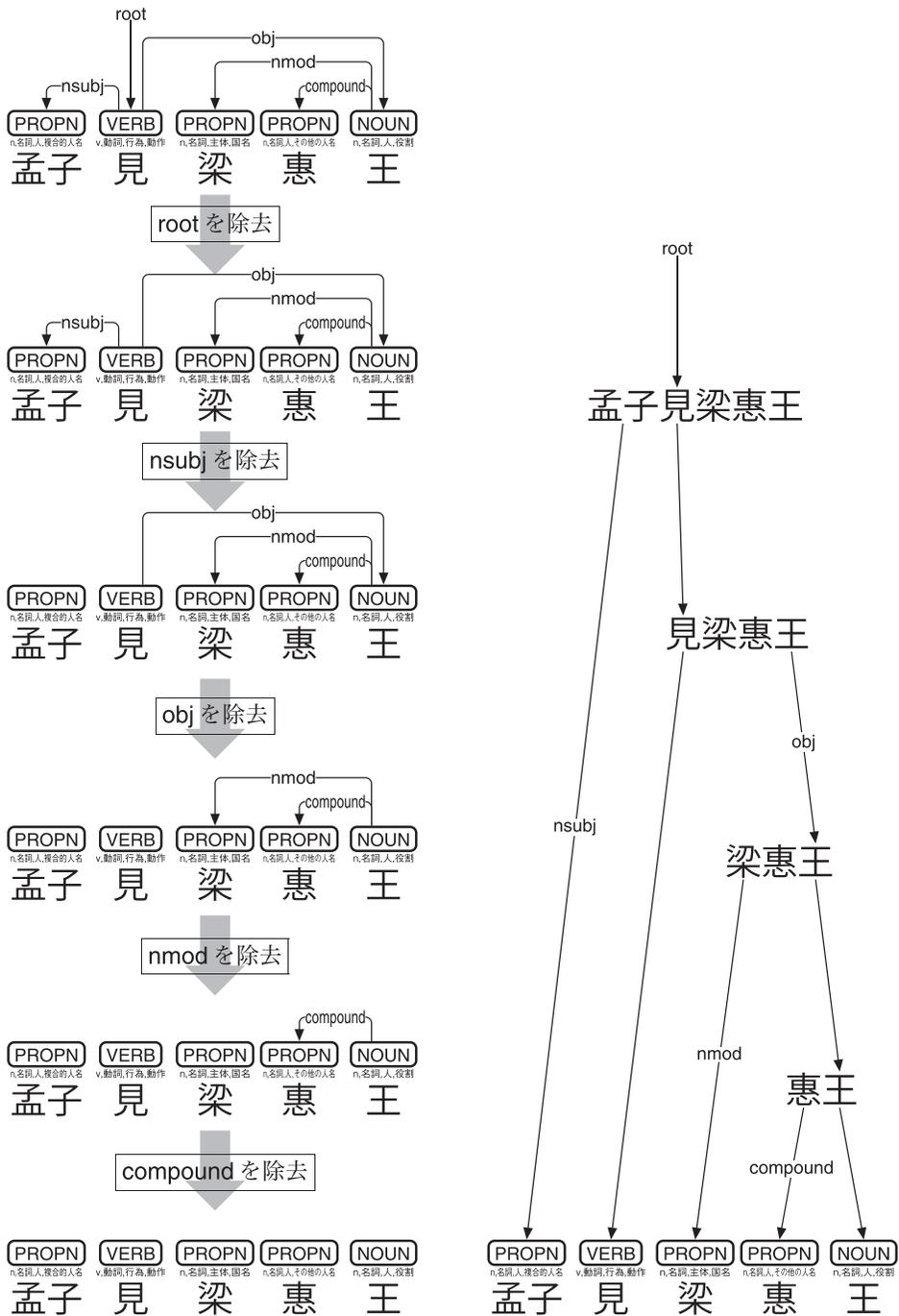


圖5 「孟子見梁惠王」の直接構成鎖解析

の除去をおこなう必要がある上に、複数のリンクの同時除去も考慮に入れなければならない。しかし、そのような除去順序を導出するようなアルゴリズムを、現時点のわれわれは開発しきれていない。われわれの今後の研究の進展に期待されたい。

## 参 考 文 献

- [1] Taku Kudo, Kaoru Yamamoto, Yuji Matsumoto: Applying Conditional Random Fields to Japanese Morphological Analysis, Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing (July 2004), pp.230-237.
- [2] Igor A. Mel'čuk: Dependency Syntax: Theory and Practice, New York: State University of New York Press (1988).
- [3] 安岡孝一: Universal Dependencies の拡張にもとづく古典中国語 (漢文) の直接構成鎖解析の試み, 情報処理学会研究報告, Vol.2019-CH-120 (2019年5月), No.1, pp.1-8.
- [4] 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體系再構築, 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.
- [5] 安岡孝一, ウィッテルンクリスティアン, 守岡知彦, 池田巧, 山崎直樹, 二階堂善弘, 鈴木慎吾, 師茂樹: 古典中国語 (漢文) の形態素解析とその應用, 情報処理学会論文誌, Vol.59, No.2 (2018年2月), pp.323-331.
- [6] Edwin G. Pulleyblank: Lexicon of Reconstructed Pronunciation in Early Middle Chinese, Late Middle Chinese, and Early Mandarin, Vancouver: UBC Press (1991).
- [7] 安岡孝一, ウィッテルンクリスティアン, 守岡知彦, 池田巧, 山崎直樹, 二階堂善弘, 鈴木慎吾, 師茂樹: 古典中国語 Universal Dependencies への挑戦, 情報処理学会研究報告, Vol.2018-CH-116 (2018年1月), No.20, pp.1-8.
- [8] Joakim Nivre: Towards a Universal Grammar for Natural Language Processing, CICLing 2015: 16th International Conference on Intelligent Text Processing and Computational Linguistics (April 2015), pp.3-16.
- [9] Carlos Gómez-Rodríguez, Joakim Nivre: A Transition-Based Parser for 2-Planar Dependency Structures, Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics (July 2010), pp.1492-1501.
- [10] Milan Straka and Jana Straková: Tokenizing, POS Tagging, Lemmatizing and Parsing UD 2.0 with UDPipe, Proceedings of the CoNLL 2017 Shared Task (August 2017), pp.88-99.
- [11] Rulon S. Wells: Immediate Constituents, Language, Vol.23, No.2 (April-June 1947), pp.81-117.
- [12] Timothy Osborne, Michael Putnam, Thomas Groß: Catenae: Introducing a Novel Unit of Syntactic Analysis, Syntax, Vol.15, No.4 (December 2012), pp.354-396.