

中日機械翻訳における中国語複合語の自動合成について†

范 莉 馨^{††} 任 福 繼^{†††}
宮 永 喜 一^{††} 柄 内 香 次^{††}

中国語科学技術文献の計算機処理において、近年の科学技術の発展に伴い、辞書に登録されていない新しい複合語が専門用語を中心に多数出現するという問題がある。したがって中国語文の解析において、このような複合語を正しく処理することが極めて重要である。本論文では、中国語文の形態素解析に際し、複合語になりうる文字列を抽出して複合語として扱う、複合語の自動合成手法を提案する。これにより、形態素解析の誤りを避け、構文解析における曖昧性を減少または解消することができる。特に原言語文の解析から目的言語文の合成に至る多段の処理を必要とする機械翻訳において、原言語文の形態素解析における曖昧性が後段の処理に大きく影響するので、本手法は有効であると考えられる。著者らは大量の教科書、科学技術文献中の複合語の調査に基づいて、複合語合成ルールをまとめ、これを組み込んで中日機械翻訳実験システムを構築した。実験により、本論文で提案した手法の有効性を確認することができた。

1. はじめに

現在、世界各地で機械翻訳の研究が精力的に行われ、いくつかの機械翻訳システムが実用化されているが、まだ多数の問題が残されている¹⁾⁻³⁾。また中日両言語間の機械翻訳に関しては、本格的な研究が開始されたばかりであり、未開拓の部分が極めて多い⁴⁾⁻⁶⁾。

中日機械翻訳では中国語複合語の処理が極めて重要である。その理由は、下記のとおりである。

中国語文ではいくつかの漢字が結合して複合語となっている場合が多く、また、特に科学技術文献では次々と新しい複合語が生まれている。一方、中国語文の解析という立場からは、なるべく複合語を単位とする方が曖昧性が減少し、都合がよい。それで、中国語文の解析に際し、すでに明らかになっている複合語だけでなく、複合語になりうる文字列を抽出し、複合語として扱うことが必要と考えられる。本論文では、この処理を複合語の自動合成という。日本語の複合語については既に多くの研究があるが⁷⁾⁻⁹⁾、中国語の複合語について、特に上記「自動合成」を念頭においた研究は非常に少ない。

われわれは、以下の2点に着目して中日機械翻訳のための中国語複合語の自動合成に関する研究を行っている。

(1) 複合語合成による曖昧性の解消

機械翻訳において文の解析を高精度で、かつ効率よく行うという観点から、できるだけ複合語単位で処理することが望ましく¹²⁾、個別の単語として扱うと、構文解析結果に曖昧性を発生しやすく、以後の翻訳過程に大きな影響を与える。次にその例を示す。

例①：機器翻訳は自然言語処理の一部。

(Jiqi fanyi shi ziran yuyan chuli de yi-bufen)

この文を形態素解析すると表1のような結果が得られる¹³⁾。この中で「翻訳 (fanyi)」、「是 (shi)」、「処理 (chuli)」の3語は動詞の属性をもつ品詞*なので、いずれもこの文の述語になりうる。したがって、このままでは構文解析結果に曖昧性が発生することになる。一般に、このような曖昧性を解消する方法としては、文の意味解析の段階で品詞間の接続頻度などによって、各品詞間の修飾・被修飾関係を確立することが考えられる。しかし、次に示すような中国語文の特徴により、この方法は中国語文に対してあまり有効でないと考えられる¹³⁾⁻¹⁶⁾。

a. 中国語文は単語間の切れ目のない「べた書き」形式で表記される。したがって形態素解析結果は日本語の場合と同様、多数の候補を含む、曖昧性の大きいものとなる。

b. 中国語の単語には多義性、多属性などをもって

† Automatic Composition of Chinese Compound Words for Chinese-Japanese Machine Translation by LIXIN FAN (Faculty of Engineering, Hokkaido University), FUJI REN (Machine Translation Development Dept., CSK Corporation), YOSHIKAZU MIYANAGA and KOJI TOCHINAI (Faculty of Engineering, Hokkaido University).

†† 北海道大学工学部

††† (株)CSK 技術開発本部

* 中国語の品詞は実詞と虚詞に大別される。実詞は文成分になることができる。一般に実質的に語彙的な意味を備えている。実詞には7種類があり、それらは名詞、動詞、形容詞、数詞、量詞、代名詞、副詞である。本論文ではそれらの属性記号をそれぞれn, v, a, m, q, r, dで表す。そして各々の品詞の下位区分をそれぞれ1, 2, 3...という細分類記号で表す。方位詞は名詞の一種であるが、中日両言語の文で特殊な役割をもっているため、通常の名詞と区別するため、その属性記号はhで表す。

表 1 形態素解析結果

Table 1 The result of morphological analysis.

No.	単語	属性記号
1	機器	n 4
2	翻訳	vn
3	是	v 7
4	自然	a 2
5	語言	n 4
6	処理	vn
7	的	u 1
8	一	m 1
9	部分	n 0

いるものが多い。さらに兼用品詞**では形態的な区別は全くない。例えば、名詞“学習”と動詞“学習する”は中国語では同じ表記“学习 (xuexi)”である。

c. 日本語と異なり、助詞、助動詞などの機能語がすべて漢字で書かれる。また動詞、助動詞、形容詞などの語尾変化、すなわち活用形もない。

ここで、「多単語が接続して同一の文成分になることができる」という中国語文法規則によって、形態素解析処理に際し接続可能な語を複合語***化することにより、曖昧性を解消、または減少させることができ、構文解析が容易になる¹⁷⁾⁻¹⁹⁾。例えば、前記表1の結果に対して複合語の合成を行うと、「機器翻訳 (jiqui fanyi)」、「自然語言処理 (ziran yuyan chuli)」、「一部分 (yi bufen)」などの新しい名詞性の複合語が合成され、表2に示すような結果が得られる。したがって、その後の構文解析において曖昧性を避けることができる。

(2) 複合語合成による形態素解析の誤りの回避

(1)からわかるように、中国語文の解析の際に、単文を構成する最小単位を複合語とする処理が有効である。そのために、すべての複合語を固定化して辞書に登録することが考えられるが、これは不相当である。その理由は、

- a. すべての複合語を登録することが実際に不可能であること；
- b. 常に複合語として処理することができない場合

** 兼用品詞とは2つの属性をもっている品詞である。兼用品詞の記号は相応品詞の属性記号としている。例えば、動詞と名詞の属性を同時にもっている品詞を動名詞と言ひ、“vn”で表記する。形容詞と動詞の属性を同時にもっている品詞を形容動詞と言ひ、“av”で表記する。

*** 1) 複合語は必ず2つ以上の単語からなるものでなければならない；2) 構文からみると、他の実詞と同じく、複合語も単文を構成する最小単位でなければならない。すなわちその間に他の単語を入れることができない。

表 2 複合語合成の結果

Table 2 The result of compound words composition.

No.	複合語	属性記号
1	機器翻訳	n 9
2	是	v 7
3	自然語言処理	n 9
4	的	u 1
5	一部分	n 9

がある；

という2つの理由からである。後者について次にその例を示す。

例②：他用機器翻訳技術文献。

(ta yong jiqui fanyi jishu wenxian)

ここで、“機器翻訳”は例①に示すように複合語である。一方、“翻訳”は動名詞である。もし“機器翻訳”という複合語を辞書に登録して常に用いることにすると、形態素解析の結果は[他, 用, 機器翻訳, 技術文献] (彼, で, 機械翻訳, 技術文献) になり、これは誤りである。正しい結果は[他, 用, 機器, 翻訳, 技術文献] (彼, で, 機械, 翻訳する, 技術文献) である。すなわち、この例文では“翻訳”は動詞で、文の述語になっている。それゆえこのような誤りを避けるためには、動名詞からなる複合語を固定的に辞書に登録せず、必要な場合にのみ複合語化することが望ましい。すなわち例文①の中の“機器”と“翻訳”は複合語“機器翻訳”を合成するが、例文②では合成せず、そのまま2つの単語として解析することが必要である。

われわれは表3に示す3種の文献から無作為に選んだ49篇のテキスト(1,862文, 76,342漢字)について、複合語の数量および種類に関する調査を行った。抽出された複合語は1,614語であった。表4にその種類を示す。表の中に、連続表記項目中の“*”はすべての細分類を意味する。“/”は接続を意味する。例えば、“n*”は名詞nのすべての細分類を表す。

上述のような観点から、本論文では中日機械翻訳のための中国語文の複合語自動合成手法を提案している。本手法の特徴として、①形態素解析の誤りを避けるため、動詞性の兼用品詞からなる複合語は辞書に登録せず、その都度合成すること、②構文解析の曖昧性を減少または解消するため、できるだけ早い段階で複合語を合成すること、③複合語合成ルールを用意して複合語合成を容易に実現し、処理時間を短縮すること、などがあげられる。

表 3 複合語の調査用文献

Table 3 Collected data for the investigation of compound words.

書名	No.	テ キ ス ト	文数	著 者	出版社(年)
科学 与 人 類	1	把日語訳成漢語的方法	53	張西祥, 孫麗英	機械工業出版社 (1987)
	2	科学与人的精神	34		
	3	新能源的探索	31		
	4	機械工業的發展基礎	26		
	5	科学進步的主要原因	48		
	6	地下資源的開發和使用	42		
	7	文明的基礎靠的是什麼	38		
	8	記憶“金屬”与宇宙工学	28		
	9	能源開發的方向	39		
	10	科学从鍊金術得到恩惠	109		
	11	未来的科学与技術	88		
	12	科学不是一件容易的事	157		
科 技 日 語 · 第 一 冊	1	X射線的發現	28	孫久明	科学普及出版社 (1984)
	2	从塑料到火箭燃料	16		
	3	侵蝕人体的黑霧	8		
	4	風箏為什麼能飛起来	14		
	5	交通工具及其顏色	22		
	6	大氣上層的離子發光	10		
	7	為什麼冬天煙霧多	13		
	8	電氣分解的發見	12		
	9	發明者是古代中国人	13		
	10	既使裝在耳朵里的收音機也能製造	12		
	11	具有結晶性質的奇妙液体	11		
	12	温泉	12		
	13	一火柴那麼多就抵得上一千噸	10		
	14	从感冒到青霉素休克	15		
	15	電話機的發明	18		
	16	電發光	13		
	17	新的能源和明天的技術	15		
	18	把平均寿命提高十年的主要角色	14		
	19	電子計算機和電子学	19		
	20	火山上開放的黄色花	13		
	21	愛迪生發明的電灯	50		
	22	合成洗滌剂	16		
	23	蜜蜂的生活	27		
	24	原子能發電	11		
	25	通古思大暴炸之迷	42		
	26	可以与石油相匹敵的資源	13		
	27	機器人	23		
	28	開拓利用海水的道路	14		
	29	鉄道の建設者	37		
	30	掌管生命和遺伝的神秘蛋白質	12		
語言 病 例 分 析	1	誤解詞義	74	蘇培成	南開大学出版社 (1987)
	2	成語使用不当	54		
	3	生造詞語	85		
	4	名詞動詞形容詞使用不当	61		
	5	数詞量詞使用不当	89		
	6	代詞使用不当	53		
	7	虛詞使用不当	227		
合計	49		1,862		

表 4 複合語調査結果の一覧
Table 4 The result of the compound words investigation.

品詞名	No.	表 記	例	数 量	%
四字 成語	1	n*/a1/n*/a1	山清水秀, 山高水險	4	7.9
	2	n*/n*/n*/n*	春風秋雨, 各種各様	5	
	3	m1/n*/m1/n*	一心一意, 一摸一樣	7	
	4	m1/a1/m1/a1	一清二白, 三長兩短	3	
	5	a1/n*/a1/n*	多種多様,	2	
	6	m1/q*/m1/q*	三番五次, 三天兩頭	2	
	7	無規則性	弄假成真, 哭笑不得	105	
名詞	8	n*/n*/n*	科学文化水平	151	47.9
	9	n*/n*	春夜, 大型計算機	534	
	10	n*/vn	機器翻訳, 社会生活	62	
	11	n*/s*	実験室里	26	
動詞	12	vn/n*	翻訳文献, 調査結果	150	9.8
	13	vn/vn	研究論証, 実験分析	8	
形容詞	14	a1/n*	大雪, 軽工業	42	10.2
	15	a2/n*	主要問題, 自然経済	79	
	16	ad/n*	整個国家, 最大值	4	
	17	an/n*	实用価値, 重要科学	29	
	18	a*/vn	共同研究, 重大発明	8	
	19	av/n*	相对関係,	2	
数詞	20	m*/n*	一九四九年	4	23.4
	21	m*/m*	三百六十五,	297	
	22	m*/q*	一回, 1.72 米	77	
代名詞	23	r*/r*	我們大家	7	0.8
	24	r*/q*	那幅, 這首	4	
	25	r*/s2	那一帶	2	
合計	25			1,614	100%

さらに、われわれは本手法に基づく複合語自動合成システムを構築し、実験により本手法の有効性を確認した。以下、2章で複合語の自動合成、3章で特殊品詞の解析、4章でこの手法の実験と考察などについて述べる。

2. 複合語の自動合成

2.1 複合語合成ルール

中国語は日本語と同様、単語間に明確な区切りのない表記方式である。そして、中国語には日本語の格助詞（ガ、ヲ、ニ、デ…）に相当する言葉、あるいは日本語やヨーロッパの言語にみられる語尾変化がない。それゆえ、中国語では主語、述語などの文法的関係は、主として単語の並べ方（語順）によって表される。したがって語順が極めて重要である¹³⁾⁻²⁰⁾。

本論文では、大量の中国語テキストから抽出した単語間接続関係に基づき、以下のことを考慮して複合語

合成ルールを定めている。

a. 前述のように、複合語合成ルールは形態素解析および構文解析の間における中間処理過程であることから、接続される単語は必ず形態素解析から見出された品詞とする。すなわちこれらの単語の属性は既知であるとする。

b. 新しい複合語を合成する最終的な目的は構文解析を簡単化し、曖昧性や誤りを減少あるいは解消するためである。それゆえ、複合語合成ルールを考えるにあたっては、単独で文の成分となりうる実詞間の接続規則しか考えないこととする。

c. 機械翻訳では処理時間をできるだけ短縮することが重要である。そこで、単語間の連続順序に右方向接続という制限を加え、構文解析の際に多重合成処理(2.2節で述べる)を行うことを可能とする。

以上のことを考慮し、品詞間の可能な接続関係を求めた^{10),11)}。

語は通常の複合語（属性記号 n9）と異なり，多重合成をしないので別な属性記号（N9）を与えている。

3. 特殊品詞の解析

中国語には，以下に示すような特別な品詞があり，特殊な文型を構成することができる。本章では，これらの品詞の特徴およびその文型について述べ，これらの文に対して中間処理を行う際に注意しなければならない問題点を検討する。

3.1 三向動詞

3.1.1 定義

中国語文法において，三向動詞は次のように定義されている^{17),19)}。

[三向動詞]: 同時に3個の名詞性の構成要素と連係される動詞を三向動詞と呼ぶ。

例③: 我 問 他 話 (wo wen ta hua).



私は 彼に 話を 聞く。

常用の三向動詞には表7に示す16個がある¹¹⁾。

3.1.2 文型および処理規則

中国語には，このような三向動詞（v1）を含む文型として以下の二種類がある。

(1) SVOO 型:

$s1 + v1 + o1 + o2$ (1)

例④, ⑤は SVOO 型の例文であり，アンダラインを付けた単語が三向動詞である。

表7 常用三向動詞一覧表

Table 7 The common SAN XIANG verbs.

No.	三向動詞	ピンイン	意味
1	報告	baogao	とどける, 報告する
2	称	cheng	と称する
3	称呼	chenghu	と呼ぶ
4	給	gei	与える
5	告訴	gaosong	告げる
6	還	huan	返す
7	教	jiao	教える
8	交	jiao	わたす
9	叫	jiao	という, ...と呼ぶ
10	借	jie	借りる, 貸す
11	留	liu	残す, 残る
12	賠	pei	賠償する, 弁償する
13	送	song	送る
14	通知	tongzhi	知らせる, 通知する
15	問	wen	たずねる, 問う, 聞く
16	贈	zeng	贈る, プレゼントする

例④: 老師送我翻譯文獻。

(laoshi song wo fanyi wenxian)

意味: 先生は 私に 翻訳の文獻を 送る。

ここで, s1は“老師”で, v1は“送”, o1は“我”, o2は“翻譯文獻”である。

例⑤: 我送給他學習用品。

(wo song gei ta xuexi yongpin)

意味: 私は 彼に 學習用品を 送る。

ここで, s1は“我”, v1は“送給”, o1は“他”, o2は“學習用品”である。

(2) 兼語式:

$s1 + v1 + os2 + v2 + (o1)$ (2)

上式の中で, v1は三向動詞である。v2は動詞また動詞の属性をもっている兼用品詞である。

例⑥, ⑦は兼語式の例文である。

例⑥: 老師叫我們翻譯文獻。

(laoshi jiao women fanyi wenxian)

意味: 先生は 私たちに 文獻を 翻訳させる。

ここで, s1は“老師”, v1は“叫”, os2は“我們”, v2は“翻譯”, o1は“文獻”である。

例⑦: 我教他包餃子。

(wo jiao ta bao jiaozi)

意味: 私は 彼に 餃子を 作るのを 教える。

ここで, s1は“我”, v1は“教”, os2は“他”, v2は“包”, o1は“餃子”である。

上例からわかるように, 三向動詞が存在する場合は, 上述の複合語合成ルールが適用できる場合と適用できない場合がある。例えば, 同じパターン“翻譯文獻”について, 例文④では複合語として合成できるが, 例文⑥では複合語になれない。

そこで, 三向動詞をもつ文における複合語合成規則を用意しなければならない。われわれは一般によく用いられる三向動詞をもつ中国語文の分析結果に基づき, 以下の3個の規則を抽出した。

規則I: 三向動詞の前に置かれている兼用品詞と他の単語との合成は複合語合成ルールが適用できる。

規則II: 「教(jiao)/教える」, 「叫(jiao)/という」, 「通知(tongzhi)/知らせる」, 「問(wen)/聞く」という4種の三向動詞の後に置かれている動名, 形動兼用品詞は文の述語になるので, 複合語合成ルールを適用できない。その他の三向動詞の文には複合語合成ルールが適用できる。

規則III: 「三向動詞+給/てあげる, てくれる, してもらう」という文には複合語合成ルールが適用できる。

3.1.3 解析の例

本項では、三向動詞のある例文2種を用いて、複合語処理過程を述べる。

例⑥：這個翻譯程序叫我們翻譯文獻。

(zhege fanyi chengxu jiao women fanyi wenxian)

(1) 形態素解析：

- CHO(1)=這個 CTY(1)=r 2
- CHO(2)=翻譯 CTY(2)=vn
- CHO(3)=程序 CTY(3)=n 4
- CHO(4)=叫 CTY(4)=v 1
- CHO(5)=我們 CTY(5)=r 1
- CHO(6)=翻譯 CTY(6)=vn
- CHO(7)=文獻 CTY(7)=n 4

ここで、CHO は中国語単語の配列であり、CTY は単語の属性の配列である。

この文には3つの動詞があり、三向動詞「叫(jiao)」と2つの動名詞「翻譯(fanyi)」である。

(2) 複合語合成：

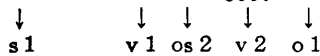
まず規則Iによって三向動詞「叫」の前に置かれている兼用品詞「翻譯」と名詞「程序(chengxu)」は複合語合成ルール No. 4 {vn/n*→n 9} によって合成し、名詞性の複合語 {翻譯程序} を合成することができる。その結果を以下に示す。

這個 {翻譯程序} 叫 我們 翻譯 文獻。

次に、規則IIによって三向動詞「叫」の後に置かれている兼用品詞「翻譯」は、文の述語であると判定される。したがって、その直後の名詞「文獻(wenxian)」と合成することはできない。

本例文は典型的な兼語式文である。以上により得られた結果を用いることにより、以後の処理は容易であり、以下に示す日本語訳文を合成することができる。

中国語構文：這個/翻譯程序/叫/我們/翻譯/文獻。



日本語訳文：この翻譯プログラムは 我々に 文獻を 翻譯させる。

例⑥：他送我學習用品。

(ta song wo xuexi yongpin)

(1) 形態素解析：

- CHO(1)=他 CTY(1)=r 1
- CHO(2)=送 CTY(2)=v 1
- CHO(3)=我 CTY(3)=r 1
- CHO(4)=學習 CTY(4)=vn
- CHO(5)=用品 CTY(5)=n 4

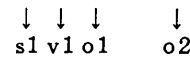
本例文では「送(song)」は三向動詞、「學習(xuexi)」は動名詞である。この文は一見兼語式の文型に見えるが、それによって翻訳すると以下に示すように誤った訳文を合成する。

誤り訳文：彼は 私に 用品を 學習するのを 送る。

実際はこの文の“學習”は動詞(學習する、勉強する)ではなく、“用品”を修飾する名詞で、複合語“學習用品”の一部である。

(2) 複合語合成：

規則IIによって三向動詞「送」の後に置かれている「學習」と「用品(yongpin)」は複合語合成ルール No. 4 によって合成できる。これにより、得られた結果はSV OO型によく対応できるので、訳文は容易に合成する。中国語構文：他/送/我/{學習用品}。



日本語訳文：彼は 私に 學習用品を 送る。

3.2 兼用品詞

3.2.1 特徴および処理規則

兼用品詞、特に動詞の属性をもつ兼用品詞の処理は中国語文の解析および機械翻訳における困難点の1つであると考えられる。それで、この兼用品詞についてさらに検討する必要がある。

中国語兼用品詞の構文解析は、活用による変化が全くないので、非常に面倒である。例えば：

例⑦：他的 生活 經歷 艱辛，曲折。

(ta de shenghuo jingli jianxin quzhe)

例⑧：他的 生活 經歷 了 兩個 時代。

(ta de shenghuo jingli le liangge shidai)

上の例文では「生活(shenghuo)」と「經歷(jingli)」という2つの単語は動名詞であり、同時に2つの属性をもっている：1つは名詞、もう1つは動詞である。

原形	名詞の意味	動詞の意味
生活	生活	暮らす
經歷	經歷	経る

形態上ではこの両者間の区別は全くないので、もしこのまま表5中の複合語合成ルールを利用し、名詞性の複合語を合成すれば、構文解析の誤りが発生することがある。

われわれの調査によれば、中国語文において常用される兼用品詞は356個あり、その中で動名詞は238個、形容動詞は118個である。そこで、中国語の論文、教科書からこれらの兼用品詞および関連情報を抽出し、複合語合成を行う際に考慮しなければならない

い、以下に示す3つの制限条件を得た。

規則IV：1つの単文には1つの動詞（あるいは動名詞，形容動詞）しかない，かつ文末の単語は形容詞でなければ，この動詞は文の述語であり，他の語と合成して複合語になることができない。

規則V：1つの単文には2つの連続している動名詞しかない，かつ文末の単語は形容詞でなければ，この2つの動名詞は合成した複合語になることができない。

規則VI：1つの単文には動詞付き型動詞*があれば，この動詞およびその後置かれている動詞（あるいは動名詞，形容動詞）は他の単語との複合をしない。

次に兼用品詞をもつ例文を示す。例⑨は複合語に合併できない例であり，例⑩，⑪は合成できる例である。

例⑨：那 是 翻译文献。

(na shi fanyi-wenxian)←合成できる

意味：それは 翻译文献で ある。

例⑩：他 正在 翻译 文献。

(ta zhengzai fanyi wenxian)←合成できない

意味：彼は 文献を 翻译し ている。

例⑪：他 的 翻译水平 很高。

(ta de fanyi-shuiping hengao)←合成できる

意味：彼の 翻译レベルは かなり高い。

3.2.2 解析の例

本項では例⑨，例⑩，⑪を用いて，兼用品詞の解析手順の概要を述べる。

まず，例文⑨では“他用機器翻译技术文献”の中の“機器 (n4)”と“翻译 (vn)”は合成ルール No. 3 (n*/vn→n9) により複合語“機器翻译”を合成できるが，この文には“用”は動詞付き型動詞なので，規則VIによって両者間の複合は行わない。

前述の例文⑩には連続して並んでいる兼用品詞が2個ある。規則Vによって，例⑩では文末の単語は形容詞なので，この2つの動名詞は合成して複合語になることができる：「生活経歴」{No. 5 の vn/vn→n9}。

逆に例⑪では文末の単語は形容詞ではないので，合成することができない。

上述の3例の訳文を以下に示す。

訳文⑨：彼は 機械で 技术文献を 翻译する。

訳文⑩：彼の 生活経歴は 苦しく，冗余曲折している。

訳文⑪：彼の 生活は 二つの 時代を 経てきた。

*中国語文にはある動詞（例えば“用”，“動”，“準許”など）は必ず後ろに他の動詞の出現を要求する。本論文ではこの種の特性をもつ動詞を動詞付き型動詞と呼ぶ。

4. 実験と考察

上述の手法に基づく実験プログラムを作成し，既に作成した形態素解析システムに組み込んで，実験により本手法の有効性を確認した。本章ではこのシステムを用いた実験結果および考察を述べる。

4.1 実験

以下に示す2回の実験を行った。1回目の実験では，本システムに組み込んだ複合語合成ルールおよび規則を抽出する際に用いた文献から163文を抽出して実験対象とした。2回目の実験では，上記以外の表8に示す科学技術に関連する論文，教科書から無作為に400文を抽出して実験対象文とした。それぞれのデータおよび翻訳結果を表9に示す。

表9中の正解率およびノイズの定義は以下に示すとおりである。

$$\text{正解率} = \frac{\text{正しく抽出された複合語数}}{\text{抽出されるべき複合語数}} \times 100\%$$

$$\text{ノイズ} = \frac{\text{誤って抽出された複合語数}}{\text{抽出されるべき複合語数}} \times 100\%$$

ここで，抽出されるべき複合語数とは，システムに組み込んだルールによって抽出されるべき（理論的）語数を示す。正しく抽出された複合語とは，複合語合成によって個々の単語として処理する場合に発生する構文的曖昧性を解消できるものの個数を意味する（表9中の構文的曖昧性を解消できる複合語の欄）。また，表9中の抑制された複合語とは，正しい構文解析結果が得られなくなるような誤った複合語合成が制限条件によって抑制された個数を意味する。

なお，合成した複合語のパターン別の出現頻度および語長別の出現頻度をそれぞれ表10に示す。

4.2 考察

実験結果から，複合語の合成は極めて良好であると判断される。合成した複合語の中に個々の単語として処理する場合の構文的曖昧性を解消できるものが167個であった。また，本システムには制限条件を用いて複合語の合成による曖昧性の発生を避けられると考えられる。以上により，本論文で提案した複合語の自動合成手法の有効性が確認された。

もちろん，本手法によりすべての複合語が完全に正しく合成できるわけではない。また，複合語の範疇は解析の用途および解析の方法により異なると考えられる。表9から，広範な分野にわたる大量のデータの場合に，複合語合成の正解率が若干低下することがわか

表 8 2 回目の実験用文献および教科書一覧
Table 8 The documents for the second experiment.

No.	書名	著者	出版社	出版年月
1	实用現代漢語語法	劉月華, 他	中国外語教学与研究出版社	1986. 10
2	現代漢語語法教程	陳国梁	中国西安市西安交大出版社	1986. 05
3	JAPANESE IN THIRTY HOURS BY EIICHI KIYOOKA 中訳本	清岡暎一著, 周炎輝訳	中国湖南科学技術出版社	1981. 05
4	中日交流標準日本語 2 中級	(中国)人民教育出版社 (日本)光村図書株式会社 合作編著	中国人民教育出版社	1990. 01
5	中国語文 3 月号	中国語文雑誌社編集部	中国国際図書貿易総公司	1983. 03
6	中国語文 4 月号	中国語文雑誌社編集部	中国国際図書貿易総公司	1983. 04

表 9 実験結果一覧
Table 9 The experimental results.

項目	回数	
	第 1 回目	第 2 回目
文数	163 (文)	400 (文)
文字数	6,683 (漢字)	15,598 (漢字)
平均文長	41 (字/文)	39 (字/文)
複合語総数	172 (個)	437 (個)
正しく合成された複合語	170 (個)	415 (個)
構文的曖昧性を解消できる複合語	71 (個)	157 (個)
抑制された複合語	35 (個)	87 (個)
誤って合成された複合語	3 (個)	9 (個)
正解率	98.8 (%)	95.0 (%)
ノイズ	1.7 (%)	2.1 (%)

った。これは複合語合成ルールおよび諸規則は大量の実験に基づきさらに充実させる必要があると意味する。次に誤りを生じた例文を示す。

例文⑨: 他 給 我 翻 訳 文 献.

(ta gei wo fanyi wenxian)

この例文は、構文的にも意味的にも曖昧性をもっており、以下の2つの構造が成立する。

構造 1 (兼語式): 他 給 我 翻 訳 文 献.

↓ ↓ ↓ ↓ ↓
s1 v1 cs2 v2 o1

(日本語訳文): 彼は 私に 文献を 翻訳してくれる。

構造 2 (SVOO 型): 他 給 我 {翻 訳 文 献}.

↓ ↓ ↓ ↓
s1 v1 o1 o2

(日本語訳文): 彼は 私に 翻訳文献を 与える。

すなわち、構造 2 に対し、“翻訳”と“文献”は複合

表 10(a) 合成した複合語のパターン一覧
Table 10(a) The patterns of composited compound words.

No.	パターン	コード	比率 (%)
①	複合名詞	n9	53.73
②	四字熟語	N9	21.20
③	数量詞	mq	9.64
④	複合代名詞	r9	8.68
⑤	複合数詞	m9	6.75

表 10(b) 合成した複合語の語長一覧
Table 10(b) The length of composited compound words.

No.	語長 (字)	比率 (%)
①	2	21.20
②	3	18.32
③	4	49.88
④	5	5.78
⑤	6	3.86
⑥	8	0.96

されなければならないが、構造 1 に対し、これらは複合できない。この構造を判断するためには、この文に関連する文脈情報が必要であると考えられる。

本論文で提案した複合語合成手法は中日機械翻訳システムの構成、とくに中国語文の解析のために考えられた方法であるが、実験結果から、その有効性は十分であると判断される。この結果を活用した中国語文の解析方式および中日翻訳アルゴリズムについては、改めて報告する予定である。

5. おわりに

中国語は単語間の切れ目のない「べた書き」形式で表記され、多属性をもつ兼用品詞が極めて多く、これらの兼用品詞は形態的に区別が全くない。また中国語

では語と語との文法関係は語順によって示される。さらに、中国語文には複合語が多数存在する。それゆえ、中国語から他の言語への機械翻訳において、これらを正しく処理しなければ、形態素解析の誤り、構文解析の曖昧性などを発生しやすく、その結果正しい訳文が得られないことになる。

本論文では、中日機械翻訳システムの構築のため、これらの中国語文の特徴を利用した中国語複合語の自動合成の手法を提案した。さらに、この手法に基づく実験システムを構築し、科学技術に関連する論文、人物伝記、教科書から無作為に抽出した400文を対象とした実験を行った。その結果、複合語合成の正解率は95%であった。これにより、本論文で提案した手法の有効性を確認することができた。今後、この結果を活用して中日機械翻訳システムの構築を予定している。

謝辞 日ごろ有益なご討論、ご助言をいただく研究室各位に深謝いたします。

参 考 文 献

1) 田中穂積, 野村浩郷: 機械翻訳, ビット別冊(1988).
 2) 長尾 真: 機械翻訳サミット, オーム社(1989).
 3) 牧野武則: 機械翻訳, オーム社(1989).

4) 寺下陽一, 二口邦夫, 鈴木 悟: 並列型パーサによる中国語文解析, 情報処理学会自然言語処理研究会報告, NL74-12 (1989.9).
 5) 任 福継, 范 莉馨, 枋内香次, 宮永喜一: 家族モデルを用いた文の分解に基づく日中機械翻訳システム, 情報処理学会論文誌, Vol. 32, No. 10, pp. 1249-1258 (1991).
 6) 寺田栄男, 孫 東恢, 田町常夫: 簡易型中日機械翻訳実験システムについて, 情報処理学会自然言語処理研究会報告, NL75-8 (1990.1).
 7) 長尾 真, 辻井潤一: 国語辞書の記憶と日本語文の自動分割, 情報処理学会論文誌, Vol. 19, No. 6, pp. 514-521 (1978).
 8) 武田浩一, 藤崎哲之助: 統計的手法を用いた漢字複合語の短単位分割, 自然言語処理, Vol. 48, No. 2, pp. 1-8 (1985.3).
 9) 宮崎正弘: 係り受け解析を用いた複合語の自動分割法, 情報処理学会論文誌, Vol. 25, No. 6, pp. 970-979 (1984).
 10) 范 莉馨, 任 福継, 宮永喜一, 枋内香次: 中国語文中の複合語の生成について, 電子情報通信学会言語理解とコミュニケーション研究会報告, NLC 90, AI 90-29 (1990.5).
 11) 范 莉馨: 中国語文解析システムにおける複合語の自動抽出に関する研究, 修士論文(1991.3).
 12) 辻井潤一: 辞書の構成と機械翻訳, 情報処理学会自然言語処理研究会報告, Vol. 26, No. 10,

付録 1 常用動名詞一覽表

No.	原形	ピンイン	No.	原形	ピンイン	No.	原形	ピンイン	No.	原形	ピンイン	No.	原形	ピンイン	No.	原形	ピンイン
1	愛好	aihao	41	代表	daibiao	81	号召	haozhao	121	考查	kaocha	181	失败	shibai	201	精通	xiaoqian
2	愛護	aihu	42	導演	daoyan	82	合理	hechangan	122	考慮	kaolu	182	承擔	shifan	202	兼做	xiangzheng
3	安排	anpai	43	点	dian	83	合舞	hezou	123	考驗	kaoyan	183	尖舞	shisuan	203	學生	xiesheng
4	捕	bai	44	調查	diaocha	84	備推	huailao	124	口誤	kouyi	184	試探	shitan	204	信任	xinren
5	幫助	bangzhu	45	調查	diaocha	85	幻想	huaxiang	125	拉扯	lailong	185	突擊	shiyang	205	行動	xingdong
6	包	bao	46	覆蓋	diaosu	86	國報	guobao	126	練習	lianxi	186	試驗	shiyang	206	休息	xiuxi
7	報道	baodao	47	離散	lisan	87	回春	huichun	127	領導	lingdao	187	收獲	shouhuo	207	宣傳	xuanchuan
8	報價	baofu	48	對比	duibi	88	金戰	jinzhan	128	埋伏	malfu	188	手術	shoushu	208	演說	yanshuo
9	報告	baogao	49	法	fa	89	活動	huodong	129	萌芽	mengya	189	說明	shuoming	209	演習	yanshi
10	保管	baogan	50	堆	dui	90	禍害	huohai	130	命令	mingling	170	束縛	sufu	210	要求	yaqiu
11	保障	baozhang	51	對面	duimian	91	計劃	jihua	131	烏黑	wuhei	171	遺記	yiji	211	議論	yilun
12	保証	baozheng	52	發明	fa ming	92	記帳	jizhang	132	捏造	niezao	172	隨從	suicong	212	除障	chuzhang
13	比喻	biyu	53	發現	faxian	93	紀念	jilian	133	判斷	panduan	173	探案	tansuo	213	聽心	yixin
14	變化	bianhua	54	反覆	fanzu	94	記載	jizai	134	判決	pianjue	174	提高	tigao	214	感傷	ganshang
15	表示	biaoshi	55	覆職	fanzhi	95	寄托	jituo	135	隨禮	peichen	175	体会	tihui	215	予感	yugan
16	表演	biaoyan	56	反映	fanying	96	假定	jiading	136	陪同	peitong	176	体现	tixian	216	運動	yundong
17	標誌	biaozhi	57	反應	fanying	97	建議	jianyi	137	批判	pipan	177	體驗	tixian	217	運動	yundong
18	滅	biyi	58	飛躍	feiyue	98	鑑定	jianding	138	批評	piping	178	排斥	tiaoxun	218	折斷	zhesuan
19	滅	bing	59	分析	fenxi	99	閉隔	biange	139	批示	pishi	179	通稱	tongchun	219	診斷	zhenduan
20	部署	bushu	60	風刺	fengci	100	剪輯	jianji	140	批注	pizhu	180	通告	tonggao	220	運動	zhendong
21	談判	tanpan	61	負擔	fudan	101	檢討	jiantaoyan	141	偏向	pianxiang	181	統計	tongji	221	支援	zhiyuan
22	殘廢	cankui	62	俘虜	fulu	102	剪貼	jiantie	142	評語	pingyu	182	運命	yunming	222	折斷	zhesuan
23	參考	cankao	63	改革	gaige	103	難關	nanguan	143	評論	pinglun	183	評論	pinglun	223	作用	zuoyong
24	參謀	cannou	64	改革	gaige	104	獎勵	jiangli	144	迫害	peihai	184	通知	tongzhi	224	障礙	zhuangai
25	沈浸	chending	65	改進	gaijin	105	教導	jiadao	145	欺騙	qipian	185	突擊	tuji	225	左右	zuoyou
26	沈浸	chenshe	66	改善	gaishang	106	教訓	jiaxun	146	齊亮	qiliang	186	突擊	tujian	226	展覽	zhanlan
27	聽講	tingjiang	67	干預	ganyu	107	教育	jiayou	147	企圖	qitu	187	突擊	tuji	227	偵探	zhenlan
28	聽講	tingjiang	68	干涉	ganse	108	經典	jiandian	148	辻政	touzheng	188	推測	tuiace	228	證明	zhengming
29	成就	chengjiu	69	感受	ganshou	109	結合	jiehe	149	審判	shenpan	189	運步	yunbu	229	指示	zhishi
30	嘗試	changshi	70	革新	gexin	110	補遺	buwei	150	審判	shenpan	190	歪曲	waiqu	230	主演	zhuanyan
31	刺	ci	71	工作	gongzuo	111	辭職	cizhi	151	傾向	qingxiang	191	妄想	wangxiang	231	主編	zhuibian
32	刺	ciji	72	貢獻	gongxian	112	警告	jingao	152	區別	qubie	192	威脅	weixie	232	主演	zhuanyan
33	啟動	chudong	73	痛苦	goutu	113	警告	jingao	153	曲解	qujie	193	誤會	wuhui	233	注釋	zhushu
34	處分	chufen	74	估計	gustu	114	審判	shenpan	154	欠火	quehuan	194	殘廢	cankui	234	注釋	zhushu
35	創造	chuangzao	75	關係	guanxi	115	編歷	bingli	155	認識	renshi	195	巧奪	qiaoduo	235	備置	beizhi
36	打掃	dasao	76	關係	guanxi	116	編歷	bingli	156	設計	sheji	196	侮辱	wuru	236	備置	beizhi
37	打算	dasuan	77	關係	guanxi	117	決定	jueding	157	生活	shenghuo	197	習慣	xiguan	237	備置	beizhi
38	打掃	dajiao	78	關係	guanxi	118	開始	kaishi	158	聲明	shengming	198	希望	xiwang	238	備置	beizhi
39	答復	dafu	79	規則	guizui	119	抗議	kuangyi	159	聲援	shengyuan	199	限制	xianzhi	238	備置	beizhi
40	代辦	daiban	80	規則	guize	120	看守	kanshou	160	勝利	shengli	200	笑話	xiaohua			

pp. 1-7 (1985. 10).
 13) ソフローノフ, M. B. (橋本万太郎(訳)): 中国語機械翻訳の一般原理, 中国語学 (1961).
 14) ヴェ・エム・ソーフツェフ (望月八十吉(訳)): 現代中国語概論, 中国語学 (1964).
 15) 劉永泉: 機器翻訳浅説, 中国語文 12月号 (1958).
 16) 望月八十吉: 中国語と日本語, 光生館 (1974).
 17) 劉月華, 潘文娉ほか: 実用現代漢語語法, 中国外語教学与研究出版社 (1986. 10).
 18) 三野昭一: 中国語文法法の基礎, 三修社 (1987).

19) 陳国梁: 現代漢語語法教程, 中国西安市西安交通大学出版社 (1986. 5).
 20) 朱美英, 内田裕士: 連接関係に基づく中国語辞書データの推定, 情報処理学会自然言語処理研究会報告, NL 73-4 (1989. 6).
 21) 范莉馨, 任福継, 宮永喜一, 枅内香次: 中国語複合語生成の実験システム, 電気関係学会北連大講演論文集, pp. 335-336 (1990. 10).
 (平成3年12月12日受付)
 (平成4年6月12日採録)

付録2 常用形容動詞一覧表

No.	原形	ピンイン	No.	原形	ピンイン	No.	原形	ピンイン
1	安定	anding	41	激動	jidong	81	深入	shenru
2	安靜	anjing	42	集中	ji zhong	82	省	sheng
3	便利	bianli	43	堅定	jianding	83	舒服	shufu
4	敗壞	baihuai	44	健全	jianquan	84	疏遠	shuyuan
5	寶貴	baogui	45	講究	jiangjiu	85	松	song
6	差	cha	46	靜	jing	86	松白	songbai
7	充實	chongshi	47	勻	yun	87	坦白	tanbai
8	純潔	chunjie	48	開闊	kaikuo	88	討厭	taoyan
9	錯	cuo	49	可憐	kelian	89	燙	tang
10	短	duan	50	肯定	kending	90	通	long
11	端正	duanzheng	51	吝	lin	91	痛快	longkuai
12	對	dui	52	寬	kuai	92	誇一	longyi
13	多	duo	53	秀明	xiu ming	93	調和	tiaohe
14	嚴	yan	54	累	lei	94	突出	tuchu
15	反複	fanfu	55	冷淡	lengdan	95	寬	wan
16	繁榮	fanrong	56	愉快	yuankuai	96	溫	wen
17	方便	fangbian	57	亂	luan	97	溫暖	wennuan
18	富余	fuyu	58	靈	ling	98	穩定	wending
19	紛爭	fensheng	59	麻煩	mafan	99	斜	xie
20	豐富	fengfu	60	滿意	manyi	100	習慣	xiguang
21	平靜	pingjing	61	模範	mofan	101	向	xiang
22	高興	gaoxing	62	密切	miqu	102	相當	xiangdang
23	公認	gongren	63	勉強	miangqiang	103	相對	xiangdui
24	恭敬	gongjing	64	明白	mingbai	104	啞	ya
25	公開	gongkai	65	平	ping	105	嚴格	yange
26	固定	guding	66	平均	pingjun	106	嚴肅	yansu
27	孤立	guli	67	平均	pingjun	107	愉快	yuankuai
28	固	gu	68	平穩	pingwen	108	寬狂	wankuang
29	光	guang	69	平穩	pingwen	109	陰險	yinxian
30	黑	hei	70	便宜	pianyi	110	讚	zan
31	橫	heng	71	漂亮	piaoliang	111	讚	zan
32	紅	hong	72	普及	puji	112	讚	zan
33	厚	hou	73	親切	qinche	113	讚	zan
34	煥發	huanfa	74	精確	jingque	114	正	zheng
35	緩和	huanhe	75	輕鬆	qingong	115	忠誠	zhongcheng
36	滄涼	cangliang	76	精確	jingque	116	壯	zhuang
37	荒涼	huangliang	77	熱鬧	renao	117	壯大	zhuangda
38	活	huo	78	熱鬧	renao	118		
39	活躍	huoyue	79	死	si			
40	擠	ji	80	少	shao			

付録3 常用動詞付き型動詞

No.	原形	ピンイン	No.	原形	ピンイン	No.	原形	ピンイン
1	愛	ai	21	派	pai	41	喜歡	xihuan
2	逼	bi	22	迫使	peishi	42	嫌	xian
3	奔	ben	23	佩服	peifu	43	笑	xiao
4	促使	chushi	24	批評	piping	44	欣賞	xinshang
5	催	cui	25	批評	piping	45	許	xu
6	打	dafa	26	氣	qi	46	過	guo
7	督促	ducu	27	強迫	qiangpe	47	過	guo
8	動員	dongyuan	28	讓	rang	48	要求	yaoqiu
9	發動	fadong	29	勸	quan	49	要求	yaoqiu
10	吩咐	fenu	30	讓	rang	50	怨	yan
11	号召	haozhao	31	認	ren	51	原諒	yuanliang
12	恨	hen	32	認為	renwei	52	用	yong
13	鼓舞	gubuyu	33	容許	rongxu	53	阻止	zuzhi
14	叫	jiao	34	允許	yunxu	54	繼續	juzhi
15	教育	jiayou	35	使	shi	55	準	zhun
16	截止	jinzhi	36	使得	shide	56	準許	zhunxu
17	可憐	kelian	37	討厭	taoyan	57	祝	zhu
18	留	liu	38	推選	tuiquan	58	贊	zan
19	罵	ma	39	推	tui			
20	命令	mingling	40	謝	xiexie			



范莉馨 (正会員)

1984年中国北京郵電学院無線通信専攻卒業。同年、中国郵電部郵電科学研究院に勤務。助理工程師。1991年北海道大学大学院電子工学専攻修士課程修了。現在、同大学院博士後期課程在学中。自然言語処理、機械翻訳に関する研究に従事。電子情報通信学会会員。



任福継 (正会員)

1982年中国北京郵電学院電信工部卒業。1985年同大学院計算機応用専攻修士課程修了。同年同大計算機工部計算機言語主講教師。1986年中国科学院博士課程、1987年中退来日。1991年北海道大学工学研究科電子工学専攻博士課程修了。工学博士。現在CSK技術開発本部研究員。計算機科学、自然言語処理、多言語機械翻訳の研究に従事。人工知能学会会員。



宮永喜一 (正会員)

1956年生。1981年北海道大学工学部電子工学専攻修士修了。工学博士。現在、北海道大学工学部電子助教授。並列計算機システム、デジタル信号処理等の研究に従事。電子情報通信学会、日本音響学会、IEEE各会員。



枅内香次 (正会員)

昭和14年生。昭和37年北海道大学工学部電気工学科卒業。昭和39年同大学院工学研究科修士課程修了。現在同工学部電子工学科教授。工学博士。自然言語処理、音声情報処理および信号処理プロセッサなどの研究に従事。電子情報通信学会、日本音響学会各会員。