

形態・構文的パターンを用いた言い換えコーパスの構築

藤田 篤[†] 乾 健太郎^{††}

語彙・構文的言い換えの中には、形態・構文的パターンに基づいて一括りにできるものの、表現を構成する語の統語・意味的な特性に依存して言い換える可否や言い換え方が決まる現象が少なくない。たとえば、複合語を構成語に分解するような言い換え、機能動詞構文の言い換え、態や格の交替、種々の動詞交替、語彙的派生などはこの語彙構成的言い換えるの範疇に含まれる。我々は現在、これら語彙構成的言い換えに関わる語の統語・意味的な特性を明らかにするため、および言い換え生成技術の定量的評価のために、個々の言い換えクラスごとに言い換え事例集（言い換えコーパス）を構築している。本稿では、言い換え前後の表現の形態・構文的パターンと既存の言い換え生成システムを用いて言い換え事例を半自動的に収集する手法について述べる。また、日本語の機能動詞構文の言い換え、動詞の自他交替を対象とした予備試行の結果を報告する。

キーワード：言い換えコーパス、語彙構成的言い換え、言い換えクラス、機能動詞構文、自他交替

Building a Paraphrase Corpus Using Morpho-syntactic Paraphrasing Patterns

ATSUSHI FUJITA[†] and KENTARO INUI^{††}

Several classes of paraphrases have a potential to be compositionally explained by combining syntactic and semantic properties of constituent words: e.g., composing/decomposing compounds, voice/case alternation, various verb alternation, and lexical derivation. Towards deep analysis of these compositional classes of paraphrases, we have examined a class-oriented framework for collecting paraphrase examples, in which sentential paraphrases are collected for each paraphrase class separately by means of automatic candidate generation using morpho-syntactic paraphrasing patterns, followed by manual judgement. Our preliminary experiments on building two paraphrase sub-corpora have so far been producing promising results with regard to cost-efficiency, exhaustiveness, and reliability.

Keywords: paraphrase corpus, compositional paraphrase, paraphrase class, light-verb construction, transitivity alternation

1. はじめに

意味が近似的に等価な言語表現の異形を言い換えと言う。言い換えの問題、すなわち同じ意味内容を伝達する言語表現がいくつも存在するという問題は、曖昧性の問題、すなわち同じ言語表現が文脈によって異なる意味を持つという問題と同様、自然言語処理における重要な問題である。

言い換えるの自動生成に関する工学的研究には、言い換えるを同一言語間の翻訳とみなし、異言語間機械翻訳（以下、単に機械翻訳）で培われてきた技術を応用する試みが多い。たとえば、構造変換方式による言い換え生成^{12),19)}、コーパスからの言い換え知識（同義表現

対や変換パターン）獲得^{1),14),16)}の諸手法は、機械翻訳向けの手法と本質的にはそれほど違わない。ただし、言い換えるは入出力が同一言語であるため、機械翻訳とは異なる性質も備えている。たとえば、平易な文章に変換する、音声合成の前処理として聴き取りやすいように変換するなど、ミドルウェアとしての応用可能性が高いことがあげられる。すなわち、言い換えるを生成する過程のどこかに、応用タスクに合わせた言い換え知識の使い分け、および目的適合性を評価する処理が必要になる。

事例集の位置付けも異なる。翻訳文書は日々生産・蓄積されており、大規模な対訳コーパスが比較的容易に利用可能である。これらは主に、翻訳知識の収集源あるいは統計モデルの学習データとして用いられる。一方、言い換え前後の文または文書の対が明示的かつ大規模に蓄えられることはほとんどない。言い換えコーパスを構築する試みはいくらも見られるものの、我々が知る限り、現在無償で公開されている言い換えコー

[†] 京都大学情報科学研究科

Graduate School of Informatics, Kyoto University

^{††} 奈良先端科学技術大学院大学情報科学研究科

Graduate School of Information Science, Nara Institute of Science and Technology

パスは Dolan ら⁵⁾ が開発したものしかない。さらに、言い換え知識の収集源として用いられるような言い換えコーパスはあっても、言い換えと呼べる現象の類型化、個々の種類の言い換える特性の分析、言い換え生成技術の開発段階における性能評価などの基礎研究に利用可能な言い換えコーパスはない。

我々は主に、言い換える実現に必要な情報を明らかにするため、および言い換え生成技術の定量的評価のために言い換えコーパスを構築している。本稿では、このような用途を想定して、

- どのような種類の言い換えるを集めるか
- どのようにしてコーパスのカバレッジと質を保証するか
- どのようにしてコーパス構築にかかる人的コストを減らすか
- 言い換え事例をどのように注釈付けて蓄えるかなどの課題について議論する。そして、コーパス構築の方法論、およびこれまでの予備調査において経験的に得られた知見について述べる。

以下、2 節では言い換えコーパス構築の先行研究について述べる。次に、我々が構築している言い換えコーパスの仕様について 3 節で、事例収集手法の詳細を 4 節で述べる。予備調査の設定を 5 節で述べ、構築したコーパスの性質について 6 節で議論する。最後に 7 節でまとめる。

2. 先行研究

言い換えコーパスの構築に関する先行研究は、内省に基づく生成、コーパスからの自動獲得の 2 種類に大別できる。いずれにおいても、コーパスを構成する個々の事例は、言い換えるの関係にある文対である。

2.1 内省に基づく言い換える生成

同じ原文に対して複数の翻訳がある場合、それらは言い換えるとみなすことができる。機械翻訳では、システムの評価方法として 1 つの原文に対して複数の正解翻訳例を用意することが一般的になってきており、そうした複数の翻訳例を含む対訳コーパスもいくつか整備されつつある^{11),15),17),18),20)}。

人間が内省に基づいて言い換えるを記述するアプローチは、大きな人的コストを要するにもかかわらずカバレッジを保証していない。すなわち、どのような種類の言い換えるを集めるのか、その範囲の言い換えるをどのようにして網羅的に集めるのかという課題に対する答えは明確でない。上述の先行研究がこのような観点での評価に至っていないのは、機械翻訳器の改善を主たる目的としているためであろう。

Web 上のニュース記事から抽出した 5,801 文対に対して 2 名の評価者が言い換えるか否かのラベルを付与したコーパス。
http://research.microsoft.com/research/nlp/msr_paraphrase.htm 文献 17), 18) は語の用法の網羅性について述べているが、ある範囲の言い換えるを網羅的に集めることには言及していない。

2.2 コーパスからの自動獲得

近年、同義表現対や変換パターンなどの言い換え知識の収集源として言い換え事例を自動的に収集する試みが報告されるようになってきた。とくにここ数年は、同じ出来事を報道している複数の新聞社の記事に対応付ける試みが多い^{2)~5),14),16)}。このアプローチでは、異なるコーパス中の文と文を、内容語や固有表現の重なり具合、構文構造の類似度、文の抽出元の記事の日付や記事中の文の位置などのメタ情報に基づいて照合し、言い換えるらしい文対を出力する。

言い換え文対の自動獲得手法には人的コストを必要としないという利点がある。収集された個々の言い換え文対には多くのノイズが含まれるが、これを人手で除外するとしても内省に基づいて事例を記述する手法に比べてコストは低い。また、未知の種類の言い換えるを発見できる可能性も秘めている。一方、文の照合における制約が収集可能な言い換え事例の種類を擬似的に限定してしまう。さらに、類似する文対を漠然と集めているに過ぎないため、複数の言い換えるが組み合わさった複雑な言い換えるの例が含まれる、潜在的に可能な言い換えるの種類や分布を反映していない、網羅性が低いなどの欠点もある。

3. 対象とねらい

言い換えると呼べる現象は多岐にわたる。その中には談話の状況に関する高度な推論を要するものもあり⁸⁾、現在の技術ですべてをカバーするのは難しい。そこで、まずはどのような種類の言い換えるの事例を集めるかについて議論する。

言い換えるに関する工学的研究のほとんどが、語あるいは言語表現の内包的意味が等価であるような現象、すなわち語彙・構文的言い換えるを対象としている。我々の焦点もこの例に洩れない。語彙・構文的言い換えるに限っても、純粋に統語論で扱えそうな言い換えるから語の詳細な意味に立ち入る必要のある言い換えるまで多岐にわたるが、実現に必要な知識の観点から以下のように 4 種類に分けて考えることができる。

統語的言い換える 個別の語の意味に立ち入らなくても、統語論の記述レベルで概ね説明できる言い換える

- (1) a. 最初に合格したのは高橋さんだ。
b. 高橋さんが最初に合格した。

語彙的言い換える 語の同義性だけで概ね説明できる、統語操作を伴わない局所的言い換える

- (2) a. 一層の苦境に陥る恐れがある。
b. 一層の窮地に陥る可能性がある。

語彙構成的言い換える 語の統語的特性と意味的特性に基づいて構成的に説明できると考えられる規則性の高い言い換える

- (3) a. 2 位が先頭との距離を縮めた。
b. 2 位と先頭の距離が縮まった。

推論的言い換え 世界知識や社会慣習に根ざし、統語論、意味論だけでは説明が難しい言い換え

- (4) a. 財政再建が急務の課題だ。
b. 緊急に財政再建する必要がある。

言い換への計算モデルが実用規模で機能するためには、大規模な言い換え知識が必要となるので、その開発および保守を効率化するための方法論が重要な研究課題になる。これには主として、人手で作成された既存の語彙資源を利用するアプローチと 2.2 項のような手法で得られた言い換えコーパスから言い換え知識を自動獲得するアプローチがある。詳細は文献 8) に譲るが、既存の語彙資源から抽出できるのは限定的な種類の言い換え知識だけであり、またコーパスから力任せに自動獲得する方法もこれまでのところ実用に耐える成果を挙げられていない。

語彙・構文的言い換えの中には、次に示す一連の例のように、上で語彙構成的言い換えと呼んだような構成的に計算できる可能性が高い言い換えも少なくない。

- (5) 態交替
a. 暴風雨が多くの地域を見舞った。
b. 多くの地域が暴風雨に見舞われた。
- (6) 場所格交替
a. 通りが群衆であふれた。
b. 群衆が通りにあふれた。
- (7) 自他交替
a. 洗濯物が風に揺れる。
b. 風が洗濯物を揺らす。
- (8) 動詞交替
a. うちの営業が S 社にコピー機を 5 台売った。
b. S 社がうちの営業からコピー機を 5 台買った。
- (9) 統語カテゴリ間の交替
a. 部屋は十分暖まっている。
b. 部屋は十分暖かい。
- (10) 複合動詞
a. 父が息子に土地を譲る。
b. 息子が父から土地を譲り受ける。
- (11) 補文構文
a. 太郎が犯人であると認める。
b. 太郎を犯人と認める。
- (12) 機能動詞構文
a. 息子が友人の活躍に刺激を受ける。
b. 息子が友人の活躍に刺激される。
- (13) 複合名詞
a. 財政再建が課題だ。
b. 財政を再建することが課題だ。
- (14) 名詞接尾辞の着脱
a. 新しい機材の必要性を議論する。
b. 新しい機材が必要かどうかを議論する。
- (15) 修飾語の係り先の交替
a. 厳密に審査基準を定める。
b. 厳密な審査基準を定める。

- (16) 主辞交替
a. リサイクルの効率化が求められる。
b. 効率的なりサイクルが求められる。

これらの例はそれぞれ異なる形態・構文的パターンによって特徴付けられる。このパターンに基づいて一括りにできる言い換え現象を本稿では言い換えクラスと呼ぶ。言い換えクラスの実在性は言語学的な分析においても示されており^{9),10),13)}、場所格交替や自他交替の構成性を言語学的に説明する試みもある。これを踏まえると、語彙構成的言い換えについては、個別の語の統語・意味的特性に関する知識と一般性の高い原理的な変換規則によって実現することが望ましい。語彙構成的言い換えが構成要素の語彙的知識から組み合わせ的に計算できるとすれば、少なくともそのクラスの言い換えについては、人手で開発・保守できる規模の語彙資源でカバーすることができる。

我々の言い換えコーパス構築の動機は、これら語彙構成的言い換えに関わる語の統語・意味的な特性を明らかにすること、その過程で言い換え生成における仮説を定量的に評価することにある。そこで、次に示すような設計で、個々の言い換えクラスごとに言い換えコーパスを構成する。

- 言い換えコーパスは言い換えクラスごとのサブコーパス群からなる。
- 各サブコーパスは所与の言い換えクラスに属する言い換え関係にある文対の集合からなる。
- 各サブコーパス中の言い換え事例は実世界における分布（密度、多様性）を十分に反映している。

4. 提案手法

3 節の議論を踏まえ、所与の言い換えクラス C 、事例収集源となる文集合 S に対して、 C に属する言い換え事例を S から (i) できるだけ網羅的に、(ii) できるだけ少ない人的コストで収集するという目標を設定する。当然、各事例の言い換えとしての適否の判定の (iii) 信頼性をできるだけ高く保たねばならない。

まずは、どのような方法論でどのような言い換えクラスをカバーできるかを経験的に調査する必要がある。その最初の試みとして、本稿では、2 節で述べた 2 種類のアプローチの中間にあたる、次の 3 ステップからなる半自動的な事例収集手法を検討する。

ステップ 1. 所与の言い換えクラス C について、形態・構文的変換パターン集合と辞書的な語彙知識を記述する。

ステップ 2. 既存の言い換え生成システムを用いて、所与の文集合 S に変換パターンを適用し、言い換え事例の候補集合を生成する。

ステップ 3. 言い換えクラスごとに適否判定ガイドラインを用意し、それに基づいて個々の言い換え候補を適格、不適格に分類する。

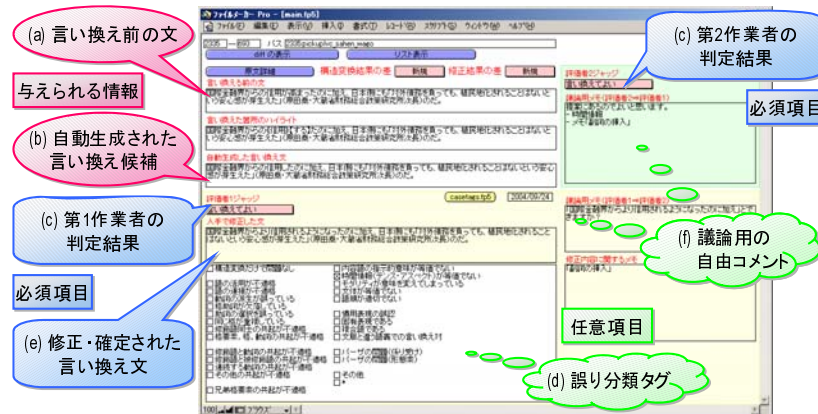


図 1 自動生成した事例の適否を判定するための作業環境

言い換え生成システムの利用により、(i) 言い換え事例収集における人的コストを低減すると同時に、(ii) 所与の言い換えクラスに対するカバレッジ、および (iii) 適否判定の質を保証しようというねらいである。

この手法では、ステップ 1 と 3 に人手を要する。まず、ステップ 1 において、所与の言い換えクラス C を定義するための形態・構文的パターンを記述する必要がある。これは、文献 6) における文法開発と同様に、少数の典型的な言い換え事例から帰納的に作成する。たとえば、(12) から、機能動詞構文の言い換えに関する (17) のようなパターン¹を記述する。

- (17) s. N を ($\Rightarrow V$) V
 t. $V(N)$

ここで、 N 、 V は各々名詞、動詞を表す変数、 $V(N)$ は N の動詞形を表す。係り受け関係に関する条件を右下に添えた矢印で表す(上の例では「 N を」が「 V 」に係ることを条件としている)。

所与の文集合 S に含まれる言い換え可能な文を網羅的に収集するために、形態・構文的パターンには過剰な制約を記述しないように配慮する。逆に、不適格な例を大量に生成してステップ 3 のコストを増やしてしまわないように、言い換えのクラスごとに語彙的な資源を用意する。たとえば、パターン (17) においては、 N と $V(N)$ を動作性名詞とその動詞形に限定する。形態素解析や係り受け解析の精度が十分実用的な精度であるため、形態・構文パターンと語彙的制約に基づいて言い換えクラスを定義するアプローチは現実的であると考えられる。

ステップ 3 では自動生成された言い換え候補の適否判定に人手を要する。ここでは言い換えクラスごとにどのような種類の誤りが生じるかをある程度予測できるという仮説に基づき、言い換えクラスごとに適否判定ガイドラインを作成しておく。このガイドラインは、文献 7) が示す言い換え誤りの分類に基づき、あらか

じめいくらかの作例に基づいて予測できる範囲で誤りの種類を列挙したものである。作業過程で未知の誤りが出現した場合、作業員間で判定結果が分かれた場合は、いくらか事例が溜まった時点で議論し、このガイドラインを更新する。

作業員は図 1 に示す作業環境で個々の言い換え候補の適否を判定する。(a) 言い換え前の文、(b) 自動生成した言い換え候補が与えられたときに、作業員は、(c) その言い換え候補の適否、(d) もしも不適格であればその原因の分類情報、(e) 修正することで言い換え可能、あるいは複数の言い換えが可能ならばそれらを記述する。判定に迷った候補については、(f) 議論用の自由コメントも記述する。

5. 言い換えコーパス構築の予備試行

前節で述べた言い換え事例の収集方法が、言い換えコーパス構築における種々の課題に対してどれだけ有効であるかを検証するため、事例分析に直接使える程度の規模のサブコーパスを構築した。今回は、機能動詞構文の言い換え、および動詞の自他交替の 2 つの言い換えクラスを取り上げた。この節では、共通の設定について述べた後、各言い換えクラスを対象としたサブコーパス構築の詳細について述べる。

5.1 共通の設定

我々の手法では、形態・構文的パターンと対象文との照合のためにいくつかのソフトウェアを必要とする。今回は、形態素解析器『茶筌』²、係り受け解析器『南瓜』³、言い換え生成モデル『KURA』⁴を用いた。

言い換え候補を収集する文のドメインは新聞記事中の文とした。具体的には日本経済新聞⁵(2000年、一文あたり平均 25.3 形態素)を用いた。茶筌、南瓜が

¹ s, t. は各々言い換え前後の文あるいはパターンを表す。

² <http://chasen.naist.jp/>

³ <http://chasen.org/taku/software/cabocha/>

⁴ <http://cl.naist.jp/kura/doc/>

⁵ <http://sub.nikkeish.co.jp/gengo/zenbun.htm>

- (22) s. N_1 が $(\Rightarrow V)$ N_2 {に, から, で} $(\Rightarrow V)$ V
 t. 変形なし.

ここで, N_1, N_2 は名詞を表す変数, V は動詞を表す変数である. 2つの格要素が動詞に係ることを条件としているが, これらの順序は問わない.

言い換え候補を1,000件程度生成することにし, LVCとのおおまかな頻度の比較から言い換え候補の収集源として25,000文を用いることにした. この文集合に(22)などのパターン群を適用したところ, V に対応する動詞800語が取り出された. そして, 各動詞に対して人手で自動詞, 他動詞を付与を記述したところ, 自動詞と他動詞の組を212組収集できた.

次に, 言い換え候補の自動生成のために, (23)のようなパターンを記述した.

- (23) s. N_1 が $(\Rightarrow V_i)$ N_2 に $(\Rightarrow V_i)$ V_i
 t. N_2 が $(\Rightarrow V_i(V_i))$ N_1 を $(\Rightarrow V_i(V_i))$ $V_i(V_i)$

N_1, N_2 はここでも名詞を表す変数である. 一方, $V_i, V_i(V_i)$ は自動詞とそれに対応する他動詞を表しており, 上の212組を用いて実現する. 動詞の自他交替には例(24), (25)のように様々な助詞が関わるが, どの要素を主格に据えるべきかは文脈に依存するため, すべての候補を別々に生成する. また, 例(26)のように他動詞文を自動詞文に言い換える例も同時に収集するため, 合計8種類のパターンを記述した.

- (24) s. 与党の法案に野党から反対意見が出る.
 t. 与党の法案に野党が反対意見を出す.
 (25) s. 戦火や迫害で難民が生まれる.
 t. 戦火や迫害が難民を生む.
 (26) s. 2位が先頭との距離を縮めた.
 t. 2位と先頭の距離が縮まった.

動詞の自他交替についても適否を判定するためのガイドラインを作成し, 修正の例を掲載した. 具体的には, (i) 活用形の修正, (ii) 格助詞の修正, (iii) ヴォイス表現の変更, の3種類の修正処理を許可した. たとえば, 例文(26s)のように他動詞を自動詞に置き換える場合, 「2位が」や「先頭との」をどのように残すべきかは形態・構文的な情報のみでは特定できない. ゆえに, 非決定のまま生成した候補を人手で修正する.

自動詞と他動詞の組を得る際に用いた25,000文に上述のパターン群を適用した結果, 985件の言い換え候補が生成された. これまでに964件の判定結果が確定しており, 484件の言い換え事例を得ている.

6. 結果と考察

前節で述べた2つの言い換えサブコーパスの仕様を表1に示す. また, 図3, 4に適否の判定結果が確定した言い換え候補の数を示す. 図中の横軸は2名の作業時間の合計であり, 言い換え候補の判定時間, 作業者間の議論の時間, 適否判定ガイドラインの更新後に再度判定する時間を含む. 以下, (i) 事例収集効率, (ii)

表1 構築した言い換えサブコーパスの仕様

言い換えクラス	LVC	TransAlt
言い換え候補の収集源の文数	10,000	25,000
言い換えパターンの数	4	8
語彙知識の種類	$\langle n, v_n \rangle$	$\langle v_i, v_t \rangle$
語彙知識の規模	20,155	212
言い換え候補の数	2,566	985
適否を判定した言い換え候補の数	1,067	964
収集した言い換え事例の数	591	484
作業時間(人時間)	118	169.5

収集した事例の網羅性, (iii) 判定結果の信頼性について述べ, (iv) 言い換えクラスの定義について議論する.

6.1 事例収集効率

現在までに2,031件の言い換え候補の判定結果が確定(5.1項で述べた通り不適格な候補の大半は1名だけの判定結果)しており, 1,075件の言い換え事例が収集できた. 図3, 4が示すように, 判定の速度は比較的安定していた. 一人時間あたりでは, 7.1件の言い換え候補の適否を確定, 3.7件の言い換え事例を収集できている. 先行研究では事例収集効率を定量的に評価していないため, 我々の手法がどれほど効率的であるかを比較によって示すことはできない. ただし, 同じ作業者が判定結果を見直すための時間, 作業者間の議論の時間も作業時間に含めていることを考慮すれば, 妥当な速度といえよう.

さらなる事例収集効率の向上のためには, どの作業に最も時間を要しているかの分析が必要である. 今回は各作業の時間を計測していなかったため, 作業者のヒアリングに基づいて次の2つの原因を取り上げる. 第一に, 言い換え候補を不適格とした場合にどのような誤りが原因で不適格としたかの記述(図1の(d))に時間を要していた. 誤り分類の体系は形態素情報や品詞体系, 係り受け構造の情報に基づいているため, 馴染みのない作業者には分類が難しかったようである. 第二に, 言語テストの難しさが作業効率を低下させる原因となっていた. これは, TransAltにおいてLVCよりも著しく(1.75倍)作業効率が悪かったことにも現れている. 6.4項で詳述する.

6.2 網羅性

どれだけ網羅的に言い換え事例を収集できているかを見積もるために, LVCで用いた文集合から無作為に750文取り出し, 人手で同じクラスの言い換えを試行した. 作成された206事例のうち獲得済みの事例は

文献3), 5)では, Web上のニュース記事から抽出した10,000文対を2名の作業者が独立に言い換えか否かに分類している. この試みでは, (i) 言い換えるクラスを限定せず, (ii) 適否に関する厳密なガイドラインなしに節の重複の度合いと作業者の直感に基づいて判定し, (iii) 判定結果が分かれた場合は議論なしに不適格としている. Chris Brockett氏とのパーソナルコミュニケーションによると2~3日(4~6人日)で作業を終えたとのことであるが, 少々速すぎる(粗すぎる)ように思える.

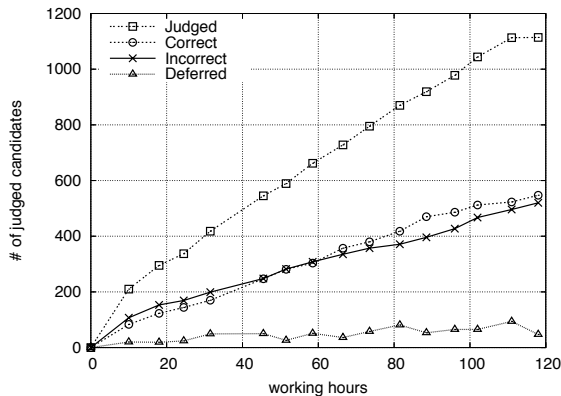


図3 適否を判定した言い換え候補の数およびその判定結果の内訳 (LVC)

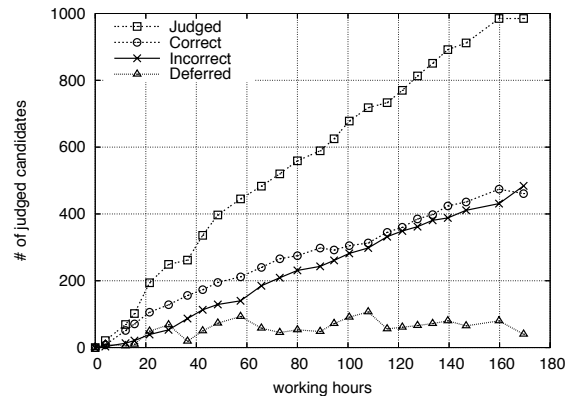


図4 適否を判定した言い換え候補の数およびその判定結果の内訳 (TransAlt)

158 事例であり、カバレッジは約 77% (158/206) となった。形態・構文的パターンでは収集できなかった 48 事例のうち、解析誤りによるものは 1 件のみであった。ゆえに、形態・構文的パターンを用いた候補生成は現実的なアプローチであると言える。34 件は形態・構文的パターンをいくつか追加することで自動的に収集できる。たとえば、(27) のようなパターンを追加すれば、(28) のような事例も収集できるようになる。

- (27) s. N 化{ が, を, に }($\Rightarrow v$) V
 t. $V(N$ 化)
- (28) s. これは市場の活性化にむけた規制緩和策だ。
 t. これは市場を活性化する規制緩和策だ。

残りの 13 件の取りこぼしは、〈ズレ、ズれる〉、〈伸び、伸びる〉のような動詞形を持つ語の品詞が IPADIC において一般名詞となっていたことに起因する。これらをあらかじめ辞書に記述しておくことはパターンの記述に比べると難しいが、形態素辞書の整備が進めばカバレッジを上げられると期待できる。

手持ちのパターンおよび語彙資源がどれだけのカバレッジを持っているか、制約としてどれだけ適切であるかを、言い換え生成および人手による適否判定の前に知ることはできない。ゆえに、上のような人手による分析は、我々がある言い換えクラスに対して持っている直感的な定義と自動的に収集できる範囲との違いを見極めるために欠かせない作業である。

6.3 判定結果の信頼性

判定結果の信頼性を保証するためにはより多くの作業者をを用いる必要がある。ただしそれは人的コストとのトレードオフになる。そこで我々は、作業者間の判定結果に揺れが生じないように言い換えクラスごとに適否判定ガイドラインを設け、適格な言い換え候補についてのみ多重判定を施した (図 2)。また、判定に

悩んだ場合は何日か後に見直す、作業者間で判定結果が分かれた場合は議論を通じて適否判定ガイドラインを更新するなどの工夫を施した。

適格と判定された言い換え候補に関する作業者間の一致率は、作業への習熟、および適否判定ガイドラインの更新に伴って上昇した。たとえば、LVC の場合の作業者間一致率は、74% (3 日目)、77% (6 日目)、88% (9 日目)、93% (11 日目) であった。このことは、作業者間の議論によって判断に悩むような言い換え候補、作業者間で判定結果が分かれるような言い換え候補に関する情報が整理され、ガイドラインが洗練されてきていることを示唆している。

図 2 の判定手順の信頼性をより正確に見積もるため、今後は第 1、第 2 作業者とは独立に言い換え候補の適否を判定する第 3 作業者を立てる予定である。

6.4 言い換えクラスの定義に関する議論

特定の言い換えクラスのみを考えるならば言い換えの適否の判定基準を明確に定義できると期待していた。しかし、LVC と TransAlt の作業効率の比較から、必ずしもその期待は満たされることが明らかになった。

TransAlt では他動詞を自動詞に言い換える際に格要素が欠落することをどこまで認めるかが議論になり、我々は言い換えによって生成された自動詞文の主格要素が意志性 (あるいは内在的コントロール¹⁰⁾) を持つか否かに着目した。すなわち、自動詞文に「自ら」、「勝手に」などの副詞を挿入した場合に文が成り立つ場合には、言い換え前の他動詞文の主格が含意されなかったため不適格とした。この言語テストに照らすと、例 (29) は適格、(30) は不適格と判定される。

- (29) s. 彼がスープを温めた。
 t. スープが温まった (*勝手に)。
 (30) s. 彼が氷を溶かした。
 t. 氷が溶けた (勝手に)。

しかしながら、言い換え前の文の主格が言い換えによって欠けるため、両例とも不適格だとの考えもある。

TransAlt の場合は、格が省略されている文を抽出できないため LVC よりもカバレッジが低いと予想される。

今回の試みによって蓄えられた多くの言い換え事例と適否判定ガイドラインには、今後このような問題を議論するための素材としての用途もある。

7. おわりに

言い換えという現象を工学的・言語学的側面の両方から解明するためには、様々な言い換えを漠然と扱うだけでなく、特定の言い換えクラスに焦点を絞った事例研究が欠かせない。本稿では、このような基礎研究の基盤となる言い換えコーパスを構築するための、言い換え前後の表現の形態・構文的パターンと既存の言い換え生成システムを用いる半自動的な事例収集手法について述べた。また、2つの言い換えクラスを取り上げた予備試行を通じ、この手法が比較的頑健に作用することを示した。

今後は、事例収集効率と適否判定の信頼性の改善をはかりながら、3節の例(5)~(16)に示したような言い換えクラスについてコーパスを構築し、言い換える構成性を裏付ける語の統語・意味的な特性を解明していきたい。なお、構築したコーパス、語彙資源、適否判定のガイドラインを公開に向けて準備中である。

参 考 文 献

- 1) Bannard, C. and Callison-Burch, C.: Paraphrasing with bilingual parallel corpora, *Proceedings of the 43th Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 597–604 (2005).
- 2) Barzilay, R. and Lee, L.: Learning to paraphrase: an unsupervised approach using multiple-sequence alignment, *Proceedings of the 2003 Human Language Technology Conference and the North American Chapter of the Association for Computational Linguistics (HLT-NAACL)*, pp. 16–23 (2003).
- 3) Brockett, C. and Dolan, W. B.: Support Vector Machines for paraphrase identification and corpus construction, *Proceedings of the 3rd International Workshop on Paraphrasing (IWP)*, pp. 1–8 (2005).
- 4) Dolan, B., Quirk, C. and Brockett, C.: Unsupervised construction of large paraphrase corpora: exploiting massively parallel news sources, *Proceedings of the 20th International Conference on Computational Linguistics (COLING)*, pp. 350–356 (2004).
- 5) Dolan, W. B. and Brockett, C.: Automatically constructing a corpus of sentential paraphrases, *Proceedings of the 3rd International Workshop on Paraphrasing (IWP)*, pp. 9–16 (2005).
- 6) Dras, M.: *Tree adjoining grammar and the reluctant paraphrasing of text*, PhD Thesis, Division of Information and Communication Science, Macquarie University (1999).
- 7) 藤田篤, 乾健太郎: 語彙・構文的言い換えにおける変換誤りの分析, *情報処理学会論文誌*, Vol. 44, No. 11, pp. 2826–2838 (2003).
- 8) 乾健太郎, 藤田篤: 言い換え技術に関する研究動向, *自然言語処理*, Vol. 11, No. 5, pp. 151–198 (2004).
- 9) Jackendoff, R.: *Semantic structures*, The MIT Press (1990).
- 10) 影山太郎: 動詞意味論—言語と認知の接点, くらしお出版 (1996).
- 11) 金城由美子, 青野邦夫, 安田圭志, 竹澤寿幸, 菊井玄一郎: 旅行会話基本表現に対する日本語パラフレーズデータの収集, *言語処理学会第9回年次大会発表論文集*, pp. 101–104 (2003).
- 12) Lavoie, B., Kittredge, R., Korelsky, T. and Rambow, O.: A framework for MT and multilingual NLG systems based on uniform lexico-structural processing, *Proceedings of the 6th Applied Natural Language Processing Conference and the 1st Meeting of the North American Chapter of the Association for Computational Linguistics (ANLP-NAACL)*, pp. 60–67 (2000).
- 13) Mel'čuk, I. and Polguère, A.: A formal lexicon in meaning-text theory (or how to do lexica with words), *Computational Linguistics*, Vol. 13, No. 3-4, pp. 261–275 (1987).
- 14) Quirk, C., Brockett, C. and Dolan, W.: Monolingual machine translation for paraphrase generation, *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 142–149 (2004).
- 15) 下畑光夫, 竹澤寿幸, 菊井玄一郎: 旅行会話における英語の同義表現コーパスの作成と分析, *情報科学技術レターズ*, pp. 83–85 (2003).
- 16) Shinyama, Y. and Sekine, S.: Paraphrase acquisition for information extraction, *Proceedings of the 2nd International Workshop on Paraphrasing: Paraphrase Acquisition and Applications (IWP)*, pp. 65–71 (2003).
- 17) 白井諭, 山本和英: 換言事例の収集—日英基本構文を対象として, *言語処理学会第7回年次大会発表論文集*, pp. 401–404 (2001).
- 18) 白井諭, 山本和英: 換言事例の収集—機械翻訳における多様性確保の観点から, *言語処理学会第7回年次大会ワークショップ論文集*, pp. 3–8 (2001).
- 19) Takahashi, T., Iwakura, T., Iida, R., Fujita, A. and Inui, K.: KURA: a transfer-based lexico-structural paraphrasing engine, *Proceedings of the 6th Natural Language Processing Pacific Rim Symposium (NLPRS) Workshop on Automatic Paraphrasing: Theories and Applications*, pp. 37–46 (2001).
- 20) Zhang, Y., Yamamoto, K. and Sakamoto, M.: Paraphrasing utterances by reordering words using semi-automatically acquired patterns, *Proceedings of the 6th Natural Language Processing Pacific Rim Symposium (NLPRS)*, pp. 195–202 (2001).