

平成 28 年度  
「中国特許文献の解析及びデータ作成事業」

報 告 書

平成 29 年 3 月 24 日

一般財団法人 日本特許情報機構



## 目次

1. 中韓文献翻訳・検索システムにおける機械翻訳文の翻訳品質評価 .....	3
1.1 本評価の概要.....	3
1.1.1 技術分野（中 AU）ごとの翻訳品質評価（一次評価） .....	3
1.1.2 テーマコードごとの翻訳品質評価（二次評価） .....	4
1.2 一次評価の実施.....	4
1.2.1 評価対象文献の選定 .....	4
1.2.2 評価対象文の選定.....	5
1.2.3 重要技術用語の選定 .....	9
1.2.4 評価基準 .....	10
1.2.5 一次評価結果（評価対象文単位） .....	11
1.2.6 一次評価結果（中 AU 単位） .....	15
1.3 二次評価の実施.....	17
1.3.1 二次評価対象テーマコード及び対象文献の選定 .....	17
1.3.2 評価対象文及び重要技術用語の選定.....	19
1.3.3 評価基準 .....	19
1.3.4 二次評価結果（評価対象文単位） .....	20
1.3.5 二次評価結果（テーマコード単位） .....	20
1.4 和文抄録作成対象案件の選定 .....	23
1.5 結果の分析 .....	24
1.5.1 一次評価結果の分析 .....	24
1.5.2 二次評価結果の分析 .....	27
1.5.3 補足的な分析 .....	35
2. 特定テーマコードの中国特許公報の和文抄録の作成 .....	45
2.1 和文抄録の作成対象.....	45
2.2 和文抄録の作成結果.....	45
2.3 特殊案件.....	45
3. 国際調査報告で引用された文献の全文翻訳文の作成 .....	46
3.1 全文翻訳文の作成結果 .....	46
3.2 原文が 4 万文字以上の文献.....	46
4. 対訳コーパスの作成 .....	47
4.1 対訳コーパスの単位.....	47
4.2 対訳コーパスの作成対象 .....	49
4.3 対訳コーパスの作成件数 .....	50

5. 辞書データの作成.....	55
5.1 辞書データの作成手順.....	55
5.2 辞書データの作成件数.....	56
5.2.1 辞書データ作成件数.....	56
5.1.2 データフォーマット.....	56
6. テーマコード情報の作成.....	59
6.1 「テーマコード情報リスト」更新内容.....	59

# 1. 中韓文献翻訳・検索システムにおける機械翻訳文の翻訳品質評価

## 1.1 本評価の概要

平成 28 年度「中国特許文献の解析及びデータ作成事業」の一環として、「中韓文献翻訳・検索システム」を通じて提供している中国特許の機械翻訳データについて、二段階の翻訳品質評価を実施した。本評価の目的は、「中韓文献翻訳・検索システム」において相対的に機械翻訳精度の低い技術分野を特定することであり、上記事業は、この評価結果に基づき当該分野の中国和文抄録を作成することで、当該分野の言語資源を充実させ、そこから得られる中日対訳コーパスや中日辞書データによる当該分野の機械翻訳精度の向上を図ることをその目的の一つとしている。

本品質評価は、以下の二段階にて実施する。

- (1) 技術分野（中 AU）ごとの翻訳品質評価（一次評価）
- (2) テーマコードごとの翻訳品質評価（二次評価）

以下、それぞれの評価の実施要件について記す。

### 1.1.1 技術分野（中 AU）ごとの翻訳品質評価（一次評価）

特許庁より貸与された「中国公開特許公報の機械付与テーマコードリスト」及び「技術分野（中 AU）とテーマコードの対応表」に基づき、326 の技術分野（中 AU）から、平成 27 年発行の中国公開特許公報を各技術分野（中 AU）につき 10 文献ずつ、計 3,260 文献を選択し、「要約」または「発明の詳細な説明」に該当する部分から 1 文献につき 1 文ずつ中国語の文を抽出する。

抽出した中国語の文に対応する機械翻訳文を「中韓文献翻訳・検索システム」から用意するとともに、正解とみなせる基準翻訳文を用意し、評価対象文の属する各技術分野に精通している評価者が「内容の伝達レベル」及び「重要技術用語の翻訳精度」の評価を行う。

「内容の伝達レベル」の評価は、基準翻訳文を作成した者ではない二人以上の評価者により重複評価を行い、その平均値を評価結果とし、「重要技術用語の翻訳精度」の評価対象となる技術用語は、技術分野ごとに抽出した 10 文のセットそれぞれから 10 語ずつ、計 3,260 語を特許庁担当者と相談の上選定する。

なお、翻訳品質評価は、技術担当者が、選択した文献に付与されている当該技術分野（中

AU) が妥当であることを確認した上で実施する。

### 1.1.2 テーマコードごとの翻訳品質評価（二次評価）

「1.1.1 技術分野（中 AU）ごとの翻訳品質評価」の結果に基づき、「技術分野（中 AU）とテーマコードの対応表」及び「中国特許文献のテーマコードと件数の対応表」を用いて、一定以上の件数規模があるテーマコード（約 1,600 テーマ）のうち、翻訳品質の低い技術分野に対応する 3 割程度のテーマコード（約 480 テーマ）を選定し、該当するテーマコードごとに、平成 27 年発行の中国公開特許公報を 10 文献程度ずつ選定する（約 4,800 文献）。

さらに、上記(1)と同様に、選定した文献の「要約」または「発明の詳細な説明」に相当する部分の中から 1 文献につき 1 文ずつ中国語の文を抽出し、抽出した中国語の文に対応する機械翻訳文を「中韓文献翻訳・検索システム」から取得するとともに、正解となる基準翻訳文と重要技術用語を用意した上で、評価者が「内容の伝達レベル」及び「重要技術用語の翻訳精度」の評価を行い、翻訳品質が低いと判断された技術分野の中でも特に翻訳品質が低いテーマコードを特定する。

## 1.2 一次評価の実施

以下、一次評価の実施内容について記す。

### 1.2.1 評価対象文献の選定

一次評価の対象文献は、特許庁から貸与された「中国公開特許公報の機械付与テーマコードリスト（平成 15～27 年発行の中国公開特許公報に対しテーマコードを機械的に付与したテーマコードリスト）」及び「技術分野（中 AU）とテーマコードの対応表」に基づき、326 の技術分野（中 AU）から、平成 27 年発行の中国公開特許公報を各技術分野（中 AU）において 10 文献ずつ選択した。

一次評価では、326 の中 AU それぞれについて評価対象文献を 10 件ずつ選定したが、その際、この 10 文献おけるテーマコード（中 AU の下層分類）を、実際の中国文献における分布を反映した構成とした。具体的には、以下の基準で 10 文献を選定した。

- ① 中 AU ごとの中国公開特許文献のテーマコード別分布を調査
- ② 各中 AU におけるテーマコードの分布率に基づき 10 文献の配分を決定

つまり、ある中 AU に属する中国公開公報全件におけるテーマコード構成が、テーマコード A が 50%、B が 30%、C と D が 10%であれば、当該中 AU のために選定する 10 文献の

構成は、A に属する文献を 5 文献、B を 3 文献、C と D を 1 文献ずつとした。このような基準でサンプリングを行うことで、それぞれの中 AU に属する主なテーマコードを全てカバーし、かつ各中 AU におけるテーマコードの分布構成も反映した形で評価対象文献を選定することができる。

下表に、上記構成の実例として中 AU 「BHM7-1」における 10 文献の構成を示す。

表 1-1 中 AU 「BHM7-1」の評価対象 10 文献の構成

中 AU	テーマコード	中国文献数 (構成比)	10 文献の構成
BHM7-1	2G188	629 (30.1%)	3 文献
	2G084	594 (28.4%)	3 文献
	2G082	321 (15.4%)	2 文献
	2G085	196 (9.4%)	1 文献
	2G081	191 (9.1%)	1 文献
	2G078	111 (5.3%)	0 文献

## 1.2.2 評価対象文の選定

上記(1)にて選定した一次評価対象文献 3,260 文献 (326 中 AU×10 文献) について、1 文献あたり 1 文ずつ、一次評価対象文を選定した。以下、評価対象文の選定基準について記す。

### 1.2.2.1 選定基準

本品質評価の目的は、「中韓文献翻訳・検索システム」を通じて提供している中国特許の機械翻訳データについて、相対的に翻訳精度の低い技術分野を特定することである。

通常、機械翻訳の品質評価は、同一の入力文を用いて、複数の異なるエンジンや、バージョンアップ前後のエンジンにおける出力結果を比較するという形式が一般的であるが、本品質評価はこれと異なり、「単一のエンジンに対し、技術分野ごとに異なる入力文を用いてその機械翻訳結果を比較し、技術分野間の優劣をつける」必要がある。したがって、評価を正確に行うポイントは、各技術分野の評価対象文を、いかに「本来、同等の翻訳難易度と考えられる文」に揃え、各分野の条件を公平にできるかにある。

とはいえ、自然文の「本来の翻訳難易度」を測定することは極めて困難であり、明確な基準は存在しない。このため、本評価においては、各中 AU にて評価対象となる 10 文について、以下の厳密な選定基準を設け、原則、全ての文をこの条件で一律に揃えることで、技術

分野間の難易度のばらつきを極力抑えることとした。

- A. 対象文の文長を一定範囲（30～50 文字）に揃える
- B. 対象文の抽出箇所を「発明の詳細な説明」に揃える
- C. 難易度を左右する特殊な表現を含む文は除外する
- D. 所定文字数（3～5 文字）の重要技術用語候補を含む文のみを対象とする

以下、各条件について説明する。

#### 1.2.2.2 選定基準 A：対象文の文長を一定範囲（30～50 文字）に揃える

機械翻訳の精度は、基本的に入力文が長くなるにつれ悪化する。したがって、技術分野ごとの評価対象文の文長がばらついてしまうと、文長の差の影響が強くなり、各分野の翻訳品質の公平な比較にならない。このため、一次評価においては、中 AU ごとに抽出する 10 文の文長を一定の範囲内に限定し、各技術分野間の比較条件を極力揃えることとした。

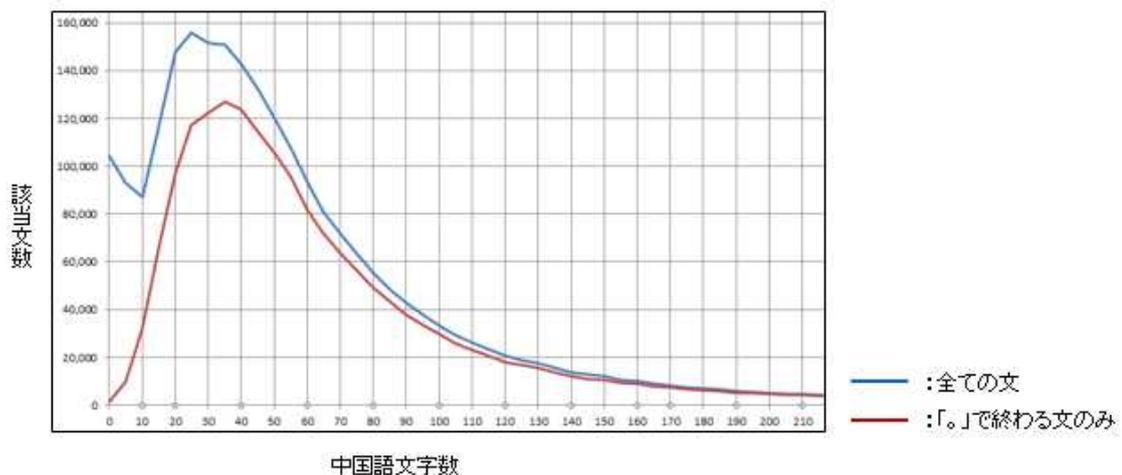
なお、分野間の条件を一律にするという意味では、複数の文長範囲を設定し、例えば 3 文、4 文、3 文など、10 文をそれぞれに割り振る手法も考えられる。しかしながら、本評価で最も重要なポイントは「各分野の『本来の翻訳難易度』を極力同等に揃えること」であり、かつ、この条件には明確な基準が存在せず、文長などコントロール可能な諸条件を可能な限り均一としたとしても、依然として偶然の要素が少なからず残留することは避けられないという事情がある。こうした要素を極力軽減するためには、同一条件のサンプル数をなるべく増やすことが非常に重要となってくる。このことから、本評価においては、複数の異なる文長範囲の文を少しずつ比較するよりも、10 文すべてを同一の文長範囲とし、サンプル数を多くするほうが、より平均的なデータ集合となり、公平な評価となる確率を高めることができるかと判断するに至った。

なお、文長の違いによる翻訳品質の差の把握は、異なるエンジン間の優劣を測る場合には非常に重要であるが、本調査は同一のエンジンでの比較であり、かつ「異なる文を（本来）同程度の難易度に揃える」という難しい条件を達成する必要があることから、サンプル数を増やしてこの条件を達成することを優先した。

10 文に適用する文長範囲については、極端に短いものや長いものは、翻訳の難易度も極端に低く／高くなり、結果として分野間の評価に差が出にくくなるため避けるべきであるが、それ以外は理論上、特段の適不適はない。ただし、中国文献において最も出現頻度の多い文長範囲を採用すれば、豊富な候補の中から本調査に最適の文を選定することができるため、分野間の条件のばらつきを最小限に抑えるという観点から最良と考えた。

この考えのもと、まずは実際の中国特許文献における文長の分布（5文字単位）を、中国公開特許明細書1ヶ月分（2015年12月公開分）全件、約180万文について調査し、中国文献の文長構成に関する分野横断的な最新の傾向を把握した。下記グラフ（図1-2）にその結果を示す。

図1-2：中国公開特許公報（2015.12）における文長分布



このグラフでは、1ヶ月分の中国公開公報の文長別の分布を、全ての文（青線）、および句点（。）で終わる文のみ（赤線）とで示しており、句点で終わる文（＝自然文とみなせる文）の文字数のピークは35～40文字の範囲であることがわかる。この結果から、中国特許文献全般においてはこのピーク周辺のものが最も文数が多く、したがってサンプルが豊富であるといえる。この文字範囲周辺は機械翻訳の難易度も適度な印象もあることから、本調査では、このピーク文字範囲35～40文字をコアとした30～50文字を評価対象文の文長の条件として、全ての評価対象文を極力この範囲のものに揃えることとした。

### 1.2.2.3 選定基準B：対象文の抽出箇所を「発明の詳細な説明」に揃える

本品質評価では、対象文の抽出箇所が「要約」または「発明の詳細な説明」に相当する部分であることが要件となっているが、本評価では、各中AUの評価対象文10文の全てを「発明の詳細な説明」から抽出した。

中国特許文献の「要約」では長文が多用される傾向にあり、本来であれば文を区切るべき話題の転換時にも、そのまま文が継続されがちである。中国文献の「要約」の平均文字数は260文字前後であるが、その多くは1文～3文程度で構成されており、「文長が極端に長く、（翻訳難易度が極端に高くなるため）技術分野間の差が出にくい」文が大半であると考えら

れる。こうした限られた数文から本評価の基準である「30～50文字」の文を搜索することは効率が悪く、条件が揃わない分野が多発するおそれがある。

一方、「発明の詳細な説明」は、候補となる文の数が圧倒的に多く、豊富なサンプル群から各分野の代表的な文を選定することができ、かつ文長をはじめとする選定文の諸条件を各分野で公平に揃えられる可能性がより高い。また、「発明の詳細な説明」は権利請求の解釈の対象範囲であるため、「要約」よりも正確な文で記載されている可能性が高く、文法的・内容的に辻褃の合わない文を誤って対象文としてしまうリスクも低くなる。

本評価では、これらの理由により、評価対象文は全て「発明の詳細な説明」から抽出すべきと結論した。

#### 1.2.2.4 選定基準C：難易度を左右する特殊な表現を含む文は除外する

自然文の中には、機械翻訳が苦手とする特殊な表現を含むものがある。例えば、文の途中にカッコ書きで長い補記が挿入されているような文は、それだけで他の文よりも機械翻訳の難易度は高くなる。こうした文が、無作為に一部の技術分野のみで評価対象文として選定されると、その分野は他の分野よりも不当に不利となる。

また、当然ながら原文に誤記や、その出願人しか使わないような特殊な用語・表現が含まれるような文は、これが原因で翻訳精度が低下することは明白であり、対象文に含めるべきではない。

こうした観点から、本評価では、評価対象文からこのような特殊な文を排除することとした。以下、その例を列挙する。

- ・誤字やタイプミスを含むもの
- ・長いカッコ書き補記を含むもの
- ・英数字記号等を多く含むもの
- ・上付・下付文字などの特殊文字を含むもの
- ・出願人特有の造語や独特の表現を含むもの
- ・固有名詞（会社名や商品名など）を含むもの

#### 1.2.2.5 選定基準D：所定文字数(3～5)の重要技術用語候補を含む文のみを対象とする

評価対象文と同様、その中からピックアップする重要技術用語についても、技術分野間で不公平な評価とならぬように配慮する必要がある。

技術用語を機械翻訳する際、概して単語よりも複合語（複数の単語からなる語）のほうが難易度が高くなることは言うまでもない。したがって、技術分野によって選定される用語の長さが不揃いでは、分野間の条件に根本的な有利不利が生じ、正確な分野間比較とならない。

本評価では、重要技術用語の翻訳精度評価においてこうした事態となることを回避し、分野間の翻訳精度の優劣が公平な条件で評価されるよう、選定する重要技術用語の文字数範囲を全中 AU において一律に揃えることとした。具体的には、各技術分野にて選定する全ての重要技術用語候補を、原則として、あらゆる分野で候補が豊富に存在し、かつ技術分野ごとの特徴が出やすいと考えられる、中国語 3～5 文字の用語（2～3 語からなる複合語が該当することが多いと思われる）に限定し、各分野の有利不利が極力解消されるようにした。

重要技術用語は、評価対象文のそれぞれから 1 語ずつ抽出することになる。このため、評価対象文の選定にあたっては、あらかじめ上記条件に該当する重要技術用語の候補を含んでおり、かつ技術分野で同一の語が重複しないことを文選定の条件とした。

#### 1.2.2.6 一次評価対象文の詳細

上記の各基準に従い選定した一次評価対象文 3,260 文の詳細は下記電子ファイルを参照されたい。

『01 一次評価（中 AU）文単位評価結果.xlsx』

#### 1.2.3 重要技術用語の選定

一次評価における「重要技術用語の翻訳精度」の評価のための重要技術用語の選定は、上記(2)にて選定した一次評価対象文 3,260 文より 1 語ずつ、計 3,260 語を選定した。一次評価対象文の選定にあたっては、前記(2)-5 にて述べたとおり、あらかじめ所定の文字数範囲の技術用語を含む文を選んでおり、各文からこの技術用語を採用した形となる。

重要技術用語についても、各中 AU で難易度を揃える必要があることから、前述のとおり、本評価では全ての中 AU で選定する全ての重要技術用語を「中国語 3～5 文字」のものに揃えた。この条件により、本評価で選定される重要技術用語は、ほぼ全てが 2～3 語で構成される複合語に統一される。

なお、重要技術用語は、同一分野で同一の用語を重複して選定することは不可としたが、分野間では重複を認めた。これは、一つの用語は必ずしも一つの分野のみで使われるものではなく、各分野で無作為に抽出された用語が分野間で重複した場合は、その語はそれぞれの分野で重要な技術用語であるため、無用な調整はせず、それぞれの分野にて評価対象にすべ

きである、との考えに基づく。

一次評価における重要技術用語 3,260 語の詳細は下記電子ファイルを参照されたい。

『01 一次評価（中 AU）文単位評価結果.xlsx』

#### 1.2.4 評価基準

一次評価では、各評価対象文に対して「内容の伝達レベル」と「重要技術用語の翻訳精度」の二つの観点から人手評価を行った。人手評価にあたっては、評価精度を高めるため、両観点とも、評価対象の分野に精通した二人の評価者が独立して評価した。すなわち、評価者 A が「内容の伝達レベル」と「重要技術用語の翻訳精度」を評価し、それと独立して評価者 B も同様の評価を実施した。そして「内容の伝達レベル」、「重要技術用語の翻訳精度」とも、評価者 A と B の平均値を各文の評価スコアとして採用した。

「内容の伝達レベル」および「重要技術用語の翻訳精度」の評価は、特許庁の「特許文献機械翻訳の品質評価手順 Ver.1.0<sup>1</sup>」に記載の下記評価基準に則り実施した。

##### 1.2.4.1 「内容の伝達レベル」の評価

「内容の伝達レベル」は、下記の 5 段階で評価した。

- 5：すべての重要情報が正確に伝達されている。(100%)
- 4：ほとんどの重要情報は正確に伝達されている。(80%～)
- 3：半分以上の重要情報は正確に伝達されている。(50%～)
- 2：いくつかの重要情報は正確に伝達されている。(20%～)
- 1：文意がわからない、もしくは正確に伝達されている重要情報はほとんどない。(～20%)

##### 1.2.4.2 「重要技術用語の翻訳精度」の評価

「重要技術用語の翻訳精度」は、下記の 4 段階で評価した。

- A（適訳語）：人手翻訳に照らし、技術的に同義かつ一般的に用いられる訳語である。
- B（可訳語）：技術用語として一般的に用いられる訳語ではないが、意味はおおむね正しい。
- C（誤訳語）：誤訳である。

---

<sup>1</sup> [http://www.jpo.go.jp/shiryoutoushin/chousa/tokkyohonyaku\\_hyouka.htm](http://www.jpo.go.jp/shiryoutoushin/chousa/tokkyohonyaku_hyouka.htm)

D（不訳語）：未知語、訳漏れである。

なお、「重要技術用語の翻訳精度」の評価スコア A、B、C、D は、集計の際、A を 5 点、B を 3 点、そして内容理解の役に立たない C と D はいずれも 0 点として数値化した。

#### 1.2.4.3 評価者の基準合わせの方法

本評価は、技術分野ごとに異なる評価者を起用し、また各文について二名の評価者による重複評価を行った。このため、各評価者間の評価基準をあらかじめ整合させて評価に臨む必要があった。

本評価においては、一次評価の実施に先立ち、全ての評価者に対し、以下の各手法にて評価基準の整合を行った。

##### ①過去の評価結果を利用した実例集の作成・配布

「特許文献機械翻訳の品質評価手順」の「内容の伝達レベル」と「重要技術用語の翻訳精度」に基づく既存の評価結果を用いて、ここからそれぞれの観点における各評点に該当する実例集を作成し、あらかじめ全評価者に配布した。

##### ②練習問題による事前の基準整合

同様に、既存の評価結果から適切な実例を採取し、これを練習問題として事前に全ての評価者に評価させた。その結果、担当者の定める正解評点とズレがある評価者に対しては、あらかじめ評価基準を調整し、全評価者の評価基準を整合させた。

##### ③評価基準の説明

評価に先立ち、本評価の手順や評価基準、注意点についての説明会を開催し、原則として全評価者を参加させた。

#### 1.2.5 一次評価結果（評価対象文単位）

一次評価では、全ての評価対象文に対し、「内容の伝達レベル」「重要技術用語の翻訳精度」とも二名の評価者にて重複評価を実施した。各文の評価スコアの算出にあたっては、下記の手順で複数の評価結果の平均値を取得し、採用した。

- ① 各文について二名の「内容の伝達レベル」の評価スコア（5～1点）の平均を算出
- ② 同じく二名の「重要技術用語の翻訳精度」の評価スコアを数値換算（A→5点、B→3点、C&D→0点）し、平均を算出
- ③ 前記①、②の平均値を合算して当該文の総合評価スコア（10～1点）を算出

一次評価における個々の評価対象文の「内容の伝達レベル（5～1 の 5 段階評価）」及び「重要技術用語の翻訳精度（A～D の 4 段階評価）」の評価結果は、下記電子ファイルを参照されたい。

[『01 一次評価（中 AU）文単位評価結果.xlsx』](#)

なお上記ファイルには、各文について二名の評価者それぞれの評価結果が併記されている。

評価結果のサンプルとして、一次評価で最も平均評価スコアが高かった中 AU である「JJJ3-1」と、最も低かった中 AU である「DGT1-1」の、各文の評価結果を表 1-3、1-4 に示す。

表 1-3 : 中 AU 「JJJ3-1」 10 文の評価結果 (総合評価最高位)

中AU	テーマ	原文	機械翻訳文	基準翻訳文	重要技術用語(中)	重要技術用語(MT)	重要技術用語(日)	評価1(内容) 5~1	評価1(用語) A~D	評価1のコメント(任意) 評価根拠等	評価2(内容) 5~1	評価2(用語) A~D	評価2のコメント(任意) 評価根拠等
JJUB-1	4.J026	上述の马来酸酐、烯基聚氧乙烯醚、丙烯酸酯的摩尔比为1:0.008~0.10:0.5~3;	上記の無水マレイン酸、アルケニルポリオキシエチレン、アクリル酸のモル比は1:0.008~0.10:0.5~3とする;	上記の無水マレイン酸、アシルポリオキシエチレンオレイルエーテル、アクリル酸のモル比は1:0.008~0.10:0.5~3とする;	马来酸酐	無水マレイン酸	無水マレイン酸	5	A		5	A	
JJUB-1	4.J026	作为优选,步骤(2)中,对于基体树脂为聚苯醚的接枝共聚物,温度为140至160°C。	好ましく、ステップ(2)にマトリックス樹脂に対してポリフェニレンエーテルのグラフト共重合体とし、温度は140~160°Cとする。	好ましくは、ステップ(2)において、マトリックス樹脂はポリフェニレンエーテルのグラフト共重合体であり、温度は140~160°Cである。	聚苯醚	ポリフェニレンエーテル	ポリフェニレンエーテル	5	A		5	A	
JJUB-1	4.J026	图1是本发明实施例1聚吡咯与亲水性聚苯异腈二嵌段共聚物的红外吸收光谱图。	図1は本発明実施例1ポリピロールと親水性ポリフェニレンイソシアニドジブロックコポリマーの赤外吸収スペクトルグラフである。	図1は本発明実施例1のポリピロールとポリフェニレンイソシアニドのジブロック共重合体の赤外吸収スペクトルグラフである。	聚吡咯	ポリピロール	ポリピロール	5	A		5	A	
JJUB-1	4.J026	在另一优选例中,在所述步骤(ii)中,所述的极性溶剂选自下组:四氢呋喃、乙醚,或其组合。	もう一つであり例が好ましいことに、前記工程(ii)に、前記極性溶媒は次のグループから選ぶ:テトラヒドロフラン、ジエチルエーテル、あるいはその組み合わせ。	他の好ましい例では、前記ステップ(ii)において、前記の極性溶媒は次のグループから選択される:テトラヒドロフラン、ジエチルエーテル、あるいはその組み合わせ。	四氢呋喃	テトラヒドロフラン	テトラヒドロフラン	5	A		5	A	
JJUB-1	4.J026	另外,由于存在异氰酸酯基,可能导致亚克力压敏胶与离型材料的粘附,无法分离。	また、イソシアネート基が存在するので、アクリル感圧接着剤と離形材料の粘着を引き起こす可能性があり、分離不可能。	更に、イソシアネート基が存在するために、アクリル感圧接着剤と離形材料の粘着は、分離不可能になる可能性がある。	异氰酸酯	イソシアネート	イソシアネート	4	A	「可能性がある」と「分離不可能」の位置が逆転している。	3	A	
JJUB-1	4.J041	本专利涉及一种天然胶乳快速凝固剂的合成方法,能快速高效地凝固天然胶乳,同时降低成本。	本特許は天然ラテックス急速凝固剤の合成方法に関し、素早く効率的に天然ラテックスを固化することができ、同時に低コスト化。	本特許は天然ラテックスの急速凝固剤の合成方法に関し、迅速で効率的に天然ラテックスを凝固させることができ、同時にコストを低減させる。	凝固剂	凝固剤	凝固剤	4	A	日本語文として不自然な部部が若干あるが、技術用語はほぼ正確。	4	A	
JJUB-1	4.J127	任选地,低聚物R可进一步包含至少一种不同于其它两种共聚单体的附加共聚单体。	任意に、オリゴマーRは少なくとも1種をさらに含んで他の2種類のモノマーの付加モノマーと異なることができる。	任意選択的に、オリゴマーRは、さらに、他の2種類のモノマーの付加モノマーと異なる少なくとも1種を含むことができる。	共聚单体	モノマー	モノマー	4	A		5	A	
JJUB-1	4.J127	上述混合单体由丙烯酸、甲基丙烯酸、丙烯酸钠和水按照重量份10:2:1:5混合而成。	上記の混合単量体はアクリル酸、メタクリル酸、アクリルスルホン酸ナトリウムと水重量分率の分10:2:1:5から混合される。	上記の混合単量体はアクリル酸、メタクリル酸、アクリルスルホン酸ナトリウムと水を重量部で10:2:1:5に従って混合される。	丙烯酸钠	アクリルスルホン酸ナトリウム	アクリルスルホン酸ナトリウム	5	A		4	A	「水重量分率の分~」が理解しにくい。
JJUB-1	4.J127	作为脂环族系多元酸,例如,可使用以下所举的脂环族二羧酸及其衍生物等。	脂環族系多塩基酸として、例えば、以下が挙げられた脂環式ジカルボン酸およびその誘導体などを用いることができる。	脂環族多塩基酸としては、例えば、以下に例示する脂環式ジカルボン酸及びその誘導体などを用いることができる。	衍生物	誘導体	誘導体	5	A		5	A	
JJUB-1	4.J127	其中,所述阻聚剂为4-甲氧基酚、对苯二酚或4-甲氧基酚与对苯二酚组成的混合物;	そのうち、前記重合防止剤は4-メトキシフェノール、ヒドロキノンあるいは4-メトキシフェノールとヒドロキノンが組成した混合物とする;	そのうち、前記重合防止剤は4-メトキシフェノール、ヒドロキノンあるいは4-メトキシフェノールとヒドロキノンで構成された混合物である;	阻聚剂	重合防止剤	重合禁止剤	5	A		5	A	

表 1-4 : 中 AU 「DGT1-1」 10 文の評価結果 (総合評価最低位)

中AU	テーマ	原文	機械翻訳文	基準翻訳文	重要技術用語(中)	重要技術用語(MT)	重要技術用語(日)	評価1(内容) 5~1	評価1(用語) A~D	評価1のコメント(任意) 評価根拠等	評価2(内容) 5~1	評価2(用語) A~D	評価2のコメント(任意) 評価根拠等
DGT1-1	3G070	燃气涡轮发动机还包括过滤器部件、至少一排可变入口导叶(VIGV)和入口空气流控制布置。	ガスタービンエンジンはまたフィルタ材、少なくとも1列の可变性入口案内翼(VIGV)と入口空気流制御配置を含んだことがある。	ガスタービンエンジンはまたフィルタ部品、少なくとも一列の可变入口ガイドブレード(VIGV)及び入口空気流制御装置を含む。	制御布置	制御配置	制御装置	2	C	「制御配置」では意味不明	2	C	エンジン入口の各要素を制御するための配置なのか、制御装置自体を指す
DGT1-1	3G070	排气系统是燃气轮机的重要辅助系统, 对其效率及正常工作有很大影响。	排気系はガスタービンの重要な支援システムであり、その効率および通常動作にとても大きい影響がある。	排気システムはガスタービンの重要な補助システムであり、効率及び通常動作に対して大きく影響する。	燃气轮机	ガスタービン	ガスタービン	3	A	排気システムの影響度には正負両面が有り得るので「支援」は不適当	3	A	「補助」と「支援」ではシステムの役割は異なるので、正しく理解することは困難である。「太陽熱」は「太陽エネルギー」との区別が不明確
DGT1-1	3G080	本发明涉及用于存储和利用来自间歇热源特别是来自太阳能收集器的能量的系统和装置。	本発明は記憶し利用して間欠の熱源の特に太陽エネルギーコレクタから来るエネルギーのシステムと装置とに関する。	本発明は間欠熱源の貯蔵及び利用に、特に太陽熱収集器からのエネルギーを利用するシステム及び装置に関する。	太阳能	太陽エネルギー	太陽熱	1	C	「太陽エネルギー」と「太陽熱」の区別が不明確	1	B	「太陽熱」は「太陽エネルギー」の部分集合
DGT1-1	3G081	本发明涉及余热发电技术领域, 尤其是涉及一种优化的有机朗肯循环低温余热发电系统。	本発明は余熱発電技術分野に、特に最適化に関する有機ランキンサイクル低温余熱発電システムである。	本発明は廃熱発電の技術分野に、特に最適化された有機ランキンサイクルによる低温廃熱発電システムに関する。	余热发电	余熱発電	廃熱発電	2	C	「余熱発電」は一般的ではない	2	C	「余熱発電」の訳が不適であるため、技術を正しく解釈することは難しい。
DGT1-1	3G104	本发明中, 所谓的“余热加热器”是指利用余热或大气对所述液体氧化剂加热的热交换器。	本発明中に、“余り蓄熱ヒーター”というものは余熱あるいは大気を利用することを指して前記液体酸化剤に加熱する熱交換器である。	本発明で、所謂“廃熱加熱器”は前記液体酸化剤に対し廃熱或いは大気を利用して加熱する熱交換器を指す。	余热加热器	余り蓄熱ヒーター	廃熱加熱器	2	C	余り、余熱、の訳語が不適切、蓄熱の訳語は誤解を招く	1	C	蓄熱機能があるかの誤解を与えるおそれあり不適切。
DGT1-1	3G104	此外, 将返回水再引入给水系统中减少了由于燃料预热而在水侧上引起的总压力损失。	また、水に戻ることは再び給水系统中に減少したことを導入して燃料予熱によって水側が引き起こす全圧に損失する。	また、復水を再度給水システムに戻し、燃料予熱によって水側に引き起こされる総圧力損失を軽減する。 ※水側は本件において「蒸気側」と対応する	燃料预热	燃料予熱	燃料予熱	1	A	訳として、まるで意味をなさない。	1	A	全体像が不明で理解できない。
DGT1-1	3G104	同时, 回热器也有一个特点, 燃气的进口压力小, 接近于正常大气压, 密封要求不高, 容许少量泄漏。	同時に、蓄冷器は1個の特徴が同じあり、燃焼ガスの入口圧は小さく、標準大気圧に接近し、密封要求は高くなく、少量が漏れることを許容する。	同時に、熱回収装置もガスの吸気圧が小さく、通常の大気圧に近く、シール性の要件は高くなく、少量の漏れは許容されるという特徴がある。	回热器	蓄冷器	熱回収装置	1	C	蓄冷器の訳語がまるで不適切	1	C	全体像が不明で理解できない。
DGT1-1	3G106	当矩形出口有主射流流出时, 将对两侧的气体进行引射抽吸使两侧产生低压区。	角型の出口にメインラジフロが流出がある時、両側のガスにエジェクタをすることを両側で低圧領域を発生することを吸引する。	矩形出口から主ジェット流が流出すると、両側の気体が吸引され、両側に低圧領域が形成される。	主射流	メインラジフロ	主ジェット流	1	C	誤訳および文法の間違い	1	C	訳語が不適切であり、機械翻訳結果から基準翻訳文を類推することは困難
DGT1-1	3G202	激光熔覆是目前较为先进的激光表面强化技术之一, 熔覆金属与母材为冶金结合, 强度高于传统钎焊。	レーザー溶着は現在比較的先進的レーザー表面強化技術の一つであり、溶着金属と母材は冶金の結合とし、強度は従来のろう付けより高い。	レーザーラッピングは現在の先進的レーザーによる表面強化技術の一つで、レーザーラッピング金属は母材と冶金結合し、従来のろう付けより強度が高い。	激光熔覆	レーザー溶着	レーザーラッピング	2	C	レーザー面強化技術の訳語が不適切で、文章の意味を阻害している	2	B	一般的説明ならば許容できる。
DGT1-1	3G202	优选地, 多孔材料中的单独的孔的孔尺寸包括大约5密耳与大约60密耳之间。	好ましくは、多孔質材料中の単独のホールの孔サイズは約5個のミルと約60個のミルの間を含む。	好ましくは、多孔質材料の単独の孔寸法は約5ミル〜約60ミルの範囲である。	多孔材料	多孔質材料	多孔質材料	2	A	5個のミル、60個のミルの訳語が不適切	1	A	ミルという用語を理解して翻訳すること。

### 1.2.6 一次評価結果（中 AU 単位）

上記 1.2.5 に示した評価対象文単位の評価結果を集計し、中 AU 単位での平均評価スコアを算出した。

算出は以下の手順で実施した。

#### [1. 文単位の平均値を算出]

- ① 各評価対象文について、二名の評価者の「内容の伝達レベル」の評価スコアの平均値を算出し、これを当該文の「内容の伝達レベル」の評価スコアとした。
- ② 各評価対象文について、二名の評価者の「重要技術用語の翻訳精度」の評価スコアを、A=5 点、B=3 点、C=0 点、D=0 点に換算のうえ平均値を算出し、これを当該文の「重要技術用語の翻訳精度」の評価スコアとした。

#### [2. 中 AU 単位の平均値を算出]

- ③ 各中 AU に属する 10 文について、①で算出した「内容の伝達レベル」の評価スコアの平均値を算出し、これを当該中 AU の「内容の伝達レベル」の評価スコアとした。
- ④ 各中 AU に属する 10 文について、②で算出した「重要技術用語の翻訳精度」の評価スコアの平均値を算出し、これを当該中 AU の「重要技術用語の翻訳精度」の評価スコアとした。
- ⑤ 各中 AU について、③で算出した「内容の伝達レベル」の評価スコアと、④で算出した「重要技術用語の翻訳精度」の評価スコアとを 1 : 1 の比率で合算し、これを当該中 AU の機械翻訳精度の総合評価スコア（10 点満点）とした。

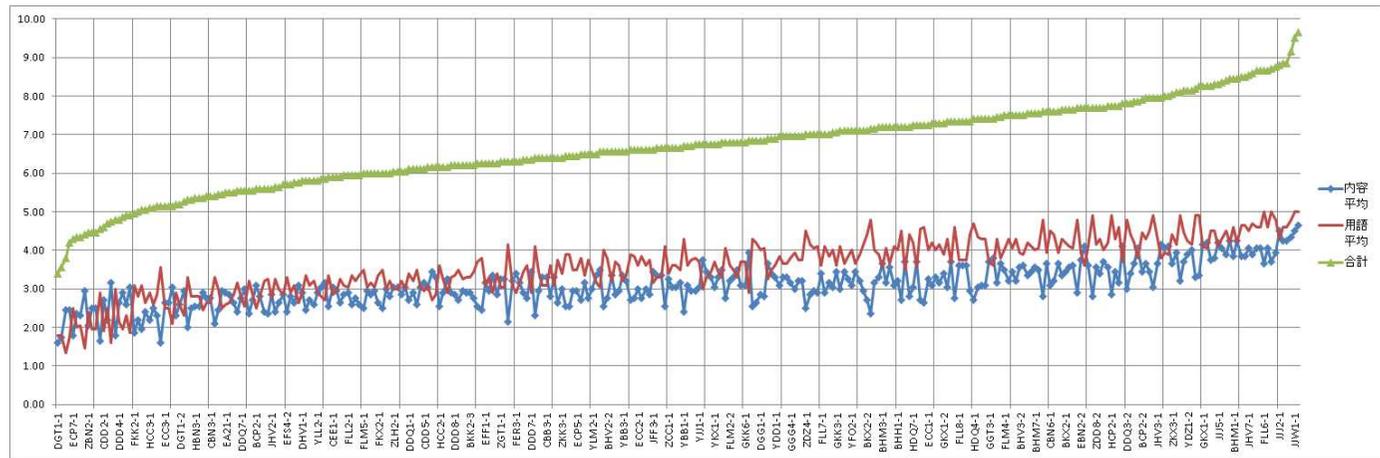
二つの評価観点の重みづけは、⑤に記載のとおり、1 : 1 とした。本事業の成果物のうち、「中韓文献翻訳・検索システム」の翻訳精度向上のために即時に利用できるのが「中日対訳辞書データ」であることから、同等レベルの翻訳精度であれば、辞書データでは直接改善できない「エンジンのロジックに起因する誤訳」よりも、辞書データで改善できる「技術用語の誤訳」が生じている中 AU を優先的に選定すべきと考えた結果、特定の一語の評価である「重要技術用語の翻訳精度」に、「内容の伝達レベル」と同等の重みを与える（すなわち、1 : 1）ことが適当と判断した。

上記手順にて算出した中 AU ごとの総合評価スコアの低い順（品質が悪いと評価された順）に下位 42 位までの一覧表と、全中 AU を平均スコア順に並べたグラフを表 1-5、1-6 に示す。全中 AU（326 分類）の総合評価スコアは別添 1-1「一次評価（中 AU）分類単位評価結果」を参照されたい。なお、本評価結果の分析は「1.5 結果の分析」の項にて述べる。

中AU コード (暫定)	内容 平均	用語 平均	合計	悪い 順
DGT1-1	1.60	1.80	3.40	1
EC2-1	1.75	1.80	3.55	2
DHH3-1	2.45	1.35	3.80	3
CBN5-1	2.45	1.75	4.20	4
EC7-1	1.80	2.50	4.30	5
DDD6-1	2.35	2.00	4.35	7
DJW4-1	2.30	2.05	4.35	7
FKK1-1	2.95	1.45	4.40	8
ZBN2-1	2.05	2.40	4.45	11
CBN1-1	2.50	1.95	4.45	11
EFF3-1	2.50	1.95	4.45	11
EC4-1	1.65	2.90	4.55	12
CDD2-1	2.70	1.90	4.60	13
ZCP2-1	2.20	2.50	4.70	14
DGT4-2	3.15	1.60	4.75	15
AFF2-1	1.80	3.00	4.80	17
DDD4-1	2.65	2.15	4.80	17
CER1-1	2.90	1.95	4.85	18
FLL8-2	2.60	2.30	4.90	20
EFS4-1	3.05	1.85	4.90	20
FKK2-1	1.85	3.10	4.95	21
HBN1-1	2.20	2.80	5.00	22
HDD1-1	1.95	3.10	5.05	24
EC5-1	2.40	2.65	5.05	24
HCC3-1	2.20	2.90	5.10	26
DJW7-1	2.50	2.60	5.10	26
DDQ2-1	2.30	2.85	5.15	31
DJJ8-1	1.60	3.55	5.15	31
EC3-1	2.65	2.50	5.15	31
YLM1-1	2.65	2.50	5.15	31
CDD6-1	3.05	2.10	5.15	31
CDD3-1	2.30	2.90	5.20	33
DGT1-2	2.60	2.60	5.20	33
DGT4-4	2.95	2.30	5.25	34
EC4-1	2.00	3.30	5.30	36
DDD5-1	2.50	2.80	5.30	36
HBN3-1	2.55	2.80	5.35	39
DDD2-1	2.55	2.80	5.35	39
ZDZ2-1	2.90	2.45	5.35	39
CBB4-2	2.75	2.65	5.40	42

(左表) 表 1-5 : 中 AU ごとの総合評価スコア一覧 (低評価順 42 位まで)

(下図) 図 1-6 : 一次評価 (中 AU) 分類単位評価結果グラフ



## 1.3 二次評価の実施

続いて、二次評価の実施内容について記す。

### 1.3.1 二次評価対象テーマコード及び対象文献の選定

中 AU 単位で実施した一次評価は、翻訳精度が相対的に低いテーマコードを特定するための「粗ぶるい」に相当する。したがって、一次評価の結果から、原則として品質評価の低かった中 AU を対象に、二次評価の対象とするテーマコードを選定することになる。

#### 1.3.1.1 二次評価対象テーマコードの選定基準

具体的には、以下の手順により、一次評価の技術分野（中 AU）ごとの翻訳品質評価の結果を基に、再び「技術情報分野（中 AU）とテーマコードの対応表」及び「中国特許文献のテーマコードと件数の対応表」を参照し、一定以上の件数規模があるテーマ（約 1,600 テーマ）を特定したうえで、ここから相対的に翻訳品質の低いと見なされる技術分野に対応する 480 テーマを選定した。

- ① 一次評価で最も評価スコアの低かった中 AU から順番に、その中 AU に属し、かつ件数規模（中国文献数）が上位 1,600 位（※）に入っている全てのテーマコードを二次評価の候補とした。

例えば中 AU 「AFS2-1」には 8 つのテーマコードが属するが、二次評価は、このうち件数上位 1,600 位までに入る 5 つのテーマコードのみを対象とした。

※詳細は電子ファイル『02 テーマコード別件数規模上位 1600.xlsx』参照。

表 1-7：AFS2-1 における二次評価対象テーマ候補の選定結果

中 AU	テーマ	文献数	件数順位	二次評価対象
AFS2-1	5J062	5,856	78	○
	5J070	3,014	211	○
	5J083	466	1,181	○
	5J058	248	1,523	○
	5J084	245	1,531	○
	5J061	9	2,268	×
	5J057	2	2,366	×
	5J063	0	2,483	×

- ② また、件数規模が上位 1,600 位に入り、かつファミリー率が 1%以下のテーマコードについては無条件で選定し、①で下位のテーマと差し替えた（通番 #450～480 が該当）。詳しくは、下記 1.3.1.2 ファミリー率の低いテーマの補完にて述べる。

③ ただし、当該テーマコードが機械付与された平成 27 年発行の中国公開特許公報が 10 件に満たない場合（つまり評価対象文献が 10 件選べない場合）、及び、当該テーマコードが機械付与された案件が実際にそのテーマに属する案件であるかを技術担当者によりランダムに 50 文献人手チェックした結果、実際にテーマに合致する文献が 10 件に満たなかった場合（つまり機械付与されたテーマコードの精度が低く、評価対象とする文献が必要数取得できない場合）は、そのテーマは対象外とし<sup>2</sup>、下位のテーマを繰り上げた。

上記手順により選定した 480 テーマの一覧は別添 1-2「二次評価対象 480 テーマ一覧」を参照されたい。

### 1.3.1.2 ファミリー率の低いテーマの補完

現状、中日特許文献の機械翻訳の精度向上に用いられる主要な言語資源は、中国文献と日本語特許ファミリー文献の対訳データとなる。したがって、中国文献の件数に比して中日ファミリー案件が少ない分野（テーマ）は、分野の規模に比して言語資源が少ないと考えられ、したがって機械翻訳精度が他の分野よりも低くなっている可能性がある。事実、平成 27 年度に特許庁が実施した品質調査の結果を分析したところ、重要技術用語の評価とファミリー率との間に一定の相関性がみられた。また、一次評価においても、下記のとおり、上記仮説を裏付ける全体的な傾向が見られた。

#### 【分析結果】

二次評価の対象範囲である 1,600 テーマ全体における低ファミリー率テーマの割合と、このうち一次評価にて評価の低かった下位 480 テーマにおける割合とを対比した結果、ファミリー率が 1%以下のテーマは、1,600 テーマ中に 94 件存在するが、このうち 52 件が品質下位 480 テーマに偏っていた。ファミリー率が 2%以下、3%以下でも同様の傾向がみられ、かつ、徐々に相関性が薄れる結果となった。

ファミリー率	1,600 テーマ中	品質下位 480 テーマ中
1%以下	94 件	52 件
2%以下	228 件	111 件
3%以下	379 件	168 件

<sup>2</sup> 別添 1-2「二次評価対象 480 テーマ一覧」の備考欄に「対象外」とある 40 テーマが該当。

二次評価対象テーマの選定においては、この傾向を踏まえ、ファミリー率が著しく低いテーマについては、一次評価の結果によらず、二次評価の対象に加えて再度精査することとした。なぜなら、こうしたテーマは言語資源が乏しいことは確実であり、仮に一次評価の結果が良好であったとしても、それは限られたサンプルによる誤差や、同じ中 AU に属する他のテーマの品質に影響された結果である可能性も考えられるからである。

補完対象とするテーマの具体的な条件は以下とした。

#### ・対象テーマ数

二次評価の対象に無条件で加えるファミリー率の閾値は、1%以下と定めた。これにより、42 テーマ (480 テーマの 8.8%) が補完対象となる。

### 1.3.1.3 評価対象文献の選定

二次評価の対象文献も、一次評価と同様、特許庁から貸与された「中国公開特許公報の機械付与テーマコードリスト」及び「技術分野 (中 AU) とテーマコードの対応表」に基づき、選定した 480 のテーマコードに属する平成 27 年発行の中国公開特許公報を各テーマコードあたり 10 文献ずつ選択した。

### 1.3.2 評価対象文及び重要技術用語の選定

二次評価における評価対象文 4,800 文及び重要技術用語 4,800 語は、一次評価と同様の基準にて選定した。選定した評価対象文並びに重要技術用語の詳細は下記電子ファイルを参照されたい。

[『03 二次評価 \(テーマ\) 文単位評価結果.xlsx』](#)

### 1.3.3 評価基準

二次評価は、一次評価と同一の基準及び手法、すなわち「内容の伝達レベル」と「重要技術用語の翻訳精度」という二つの観点から二名の評価者による重複評価を行った。評価基準の詳細は前記 1.2.4 を参照されたい。

#### 1.3.4 二次評価結果（評価対象文単位）

二次評価における各評価対象文の「内容の伝達レベル（5～1の5段階評価）」及び「重要技術用語の翻訳精度（A～Dの4段階評価）」の評価結果は、下記電子ファイルを参照されたい。

[『03 二次評価（テーマ）文単位評価結果.xlsx』](#)

#### 1.3.5 二次評価結果（テーマコード単位）

上記(4)に示した評価対象文単位の評価結果を集計し、テーマコード単位での平均評価スコアを算出した。算出方法は一次評価と同様であり、詳細は前記 1.2.6 を参照されたい。

上記手順にて算出したテーマコード別の総合評価スコアを別添 1-3「二次評価（テーマ）分類単位評価結果」に示す。なお、本評価結果の分析は「1.5 結果の分析」の項を参照されたい。

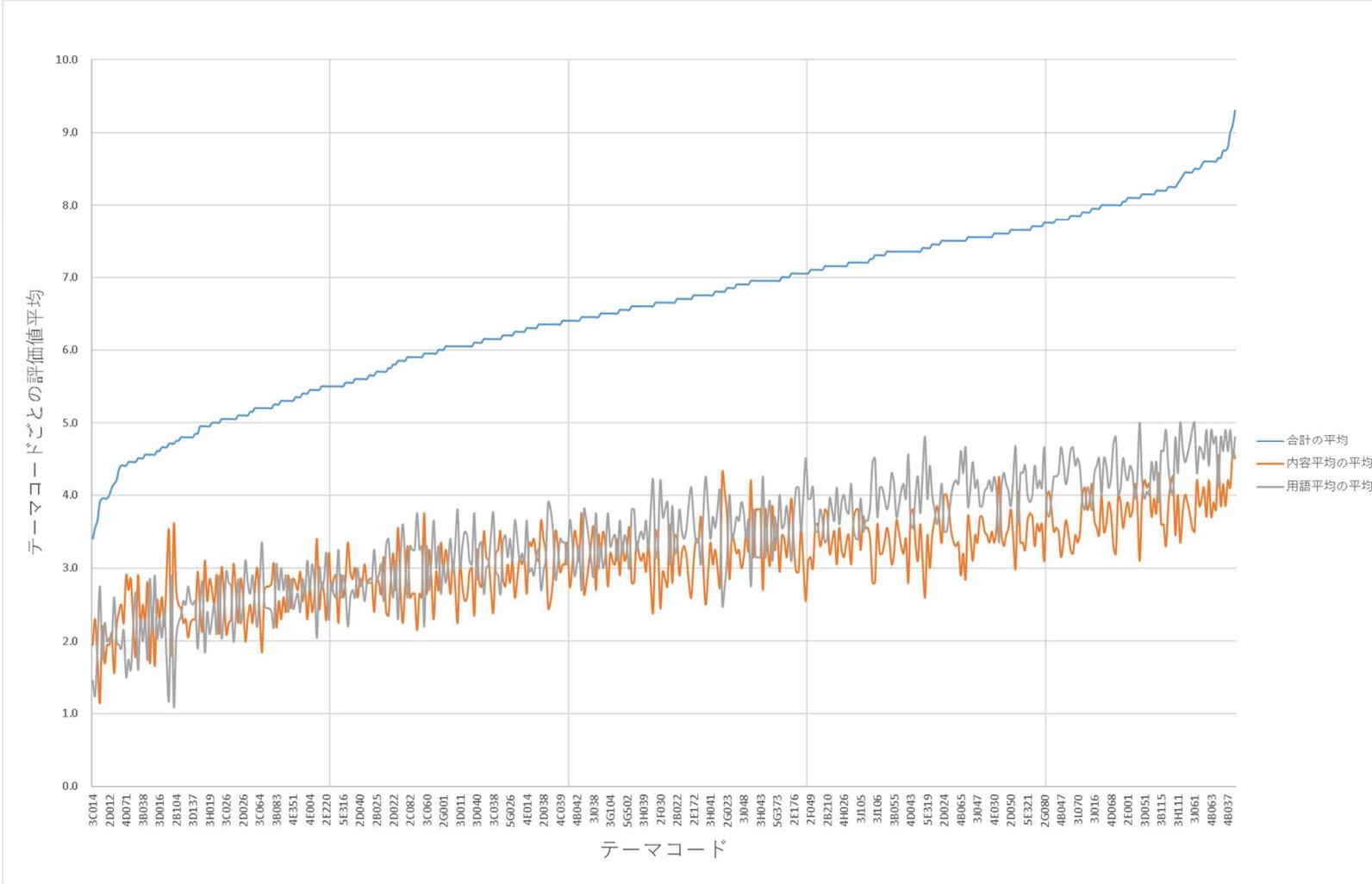
表 1-8 にその平均評価スコアが低い順に 30 テーマを掲載する。また、図 1-9 に、二次評価対象全 480 テーマを平均評価スコアの低い順に並べたグラフを示す。

なお、表 1-8 にて「該当案件数」が橙色・緑色にハイライトされたセルはそれぞれ案件数が 1,000 件以上、2,000 件以上のテーマを示している。また「上限 2,000 件」欄が緑色のセルは、案件数が 2,000 件以上につき上限の 2,000 件を和文抄録作成対象としてカウントしたことを表している。

表 1-8：二次評価の平均評価スコア下位 30 テーマ一覧

通番	中AUコード	テーマコード	中AU	テーマ名	件数規模 上位	中国 文献数	ファミリー率	平均 評価スコア	(内容)	(用語)	該当 案件数	上限 2000件	件数 累計
1	ECP2-1	3C014	非金属の加工	木材用鋸盤の構成部品, 付属品	1592	209		3.40	1.95	1.45	50	50	50
2	CDD4-1	2D013	掘削	ショベル系を除いた土砂掘削機及び施工法	1470	280	2%以下	3.55	2.30	1.25	135	135	185
3	ECC5-1	3C065	携帯工具	はさみ・ニッパ	860	773		3.65	1.75	1.90	177	177	362
4	CDD4-1	2D065	掘削	さく岩、採鉱及び採鉱機械とその方法	809	857	2%以下	3.90	1.15	2.75	830	830	1192
5	ECP7-1	3B221	装飾技術	装飾技術	184	3258	2%以下	3.95	1.95	2.00	1267	1267	2459
6	CEE2-1	2E108	仕上げ	屋根ふき・それに関連する装置または器具	466	1546	2%以下	3.95	2.20	1.75	483	483	2942
7	CDD4-1	2D003	掘削	掘削機械の作業制御	699	1026		3.95	1.70	2.25	384	384	3326
8	CDD4-1	2D012	掘削	ショベル系（制御を除く）	891	743		4.00	1.95	2.05	235	235	3561
9	ECC5-1	3C061	携帯工具	ナイフ	1036	586	3%以下	4.10	2.10	2.00	221	221	3782
10	ECC5-1	3C068	携帯工具	可搬型釘打ち機およびステーパー	787	883		4.15	1.55	2.60	239	239	4021
11	FLM5-1	3L102	製氷・冷蔵庫	冷蔵庫の箱体（壁体）2	708	1010		4.20	2.20	2.00	328	328	4349
12	FLM5-1	3L020	製氷・冷蔵庫	被冷蔵物の充填、照明装置	1002	624		4.35	2.40	1.95	202	202	4551
13	EFF3-1	4D063	破砕・粉砕	破砕・粉砕（1）	856	780	2%以下	4.40	2.25	2.15	475	475	5026
14	HDQ8-1	4D071	分離	液体又は風力による固体相互の分離	1594	207	2%以下	4.40	2.90	1.50	466	466	5492
15	EFF2-1	3F312	昇降・揚重	ジャッキ	1432	299	2%以下	4.40	2.50	1.90	340	340	5832
16	FKK2-1	4L029	生活家電	アイロン	920	708	2%以下	4.45	1.80	2.65	137	137	5969
17	CDD6-1	2D054	トンネル	立坑・トンネルの掘削技術	130	4120	1%以下	4.45	2.85	1.60	1354	1354	7323
18	CDD3-1	2D049	基礎工	基礎工事に適用される隔壁	1507	257	2%以下	4.45	2.70	1.75	202	202	7525
19	DDQ3-1	3D017	事故防止	乗員・歩行者の保護	661	1085		4.45	2.25	2.20	195	195	7720
20	FKK2-1	3B038	生活家電	ヘアカール	1451	290		4.50	2.50	2.00	26	26	7746
21	CDD3-1	2D043	基礎工	地盤の調査及び圧密・排水による地盤強化	740	960	1%以下	4.50	2.20	2.30	718	718	8464
22	CBB4-3	2B150	バイオ・その他	飼料（2）（一般）	68	6722	1%以下	4.50	2.90	1.60	2459	2000	10464
23	DGT1-1	3G106	タービン	ジェット推進設備	796	876		4.55	1.65	2.90	279	279	10743
24	DGT1-1	3G070	タービン	タービンの細部・装置	1463	284		4.55	2.45	2.10	168	168	10911
25	CBN3-1	2C030	教習具	練習用教習具	1044	579		4.55	2.25	2.30	92	92	11003
26	EFF2-1	3F205	昇降・揚重	ジブクレーン（門形、ケーブルクレーン）	295	2345		4.55	1.70	2.85	1239	1239	12242
27	CDD5-1	2D059	陸路	橋または陸橋	173	3413	1%以下	4.55	2.80	1.75	2107	2000	14242
28	CEE2-1	2E301	仕上げ	階段・手すり	1569	223	2%以下	4.60	2.55	2.05	97	97	14339
29	DDD1-1	3D016	車体構造	車両用パンパ	791	878		4.60	2.30	2.30	133	133	14472
30	CBN3-1	2C032	教習具	教示用装置	103	4867		4.65	2.60	2.05	961	961	15433

図 1-9 : 二次評価 480 テーマ 評価スコア下位順グラフ



## 1.4 和文抄録作成対象案件の選定

二次評価の結果から、480 テーマの中から翻訳精度が相対的に低いテーマコードを特定し、このテーマに属する中国特許文献を和文抄録作成対象として選定した。本評価では、以下の基準により、本事業にて和文抄録の作成対象とする 8 万件を選定した。

- ① 二次評価の平均スコアの低かったテーマから、原則として、抄録作成範囲として指定された「平成 15～27 年発行の中国公開特許公報（A 公報）に対応する中国特許公報（B 公報）」からそのテーマに属する全ての案件を対象とした。
- ② ただし、上記①で対象となるテーマの中には、案件数が数千件を超えるようなものも存在する。こうしたテーマについては、案件数の上限を 2,000 件と定め、これを超えるテーマについては登録日の新しい順に 2,000 件のみを対象とすることで、極力多数のテーマがカバーされるようにした。
- ③ 上記①～②にてカウントした対象文献数を平均スコア下位テーマから順に累計し、総数が 8 万件に達したところで終了した。結果、平均スコア 6.20 までが和文抄録作成対象となり、このうち 5H181 は、登録日の新しい順に 544 文献を対象とした。
- ④ なお 3B040（下位順で 47 位）、2D118（同 80 位）、5G016（同 174 位）は、上記①に該当する案件が 1 件も存在しなかったため、結果的に対象外となった。

上記手順により和文抄録作成対象とした全テーマコードを別添 1-3「二次評価（テーマ）分類単位評価結果」に示す。なお、選定された 8 万件の詳細は電子ファイル『04 和文抄録作成対象 8 万件.xlsx』を参照されたい。

## 1.5 結果の分析

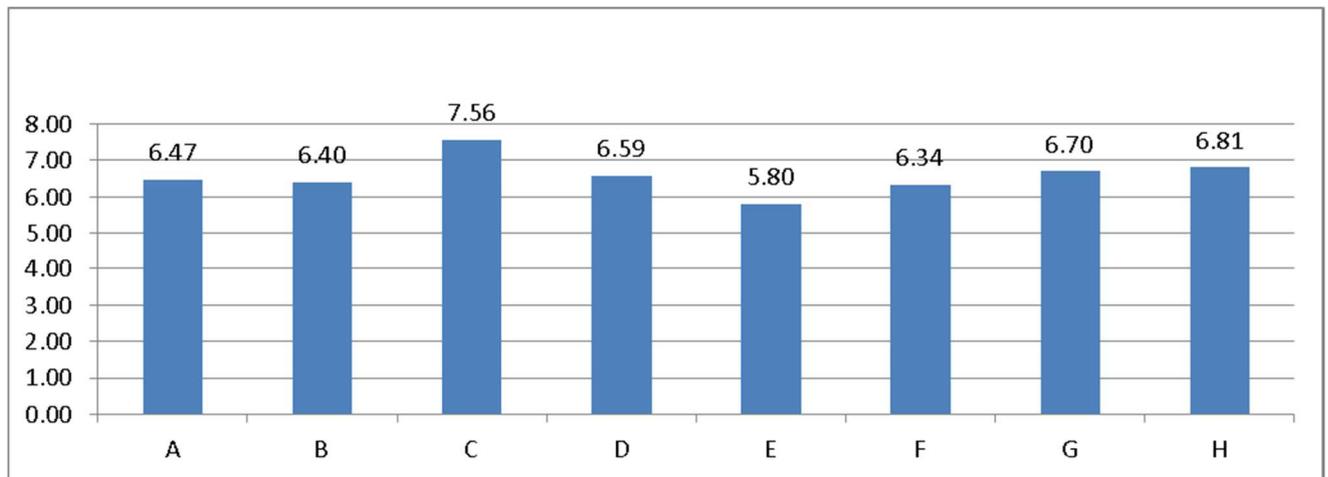
一次評価、二次評価で得られた人手評価結果を分析した。

### 1.5.1 一次評価結果の分析

#### 1.5.1.1 セクション単位の集計

各中 AU に対応する筆頭 FI の先頭文字を IPC のセクションとみなし、下記グラフのとおりセクション単位の評価スコアの平均を算出した。

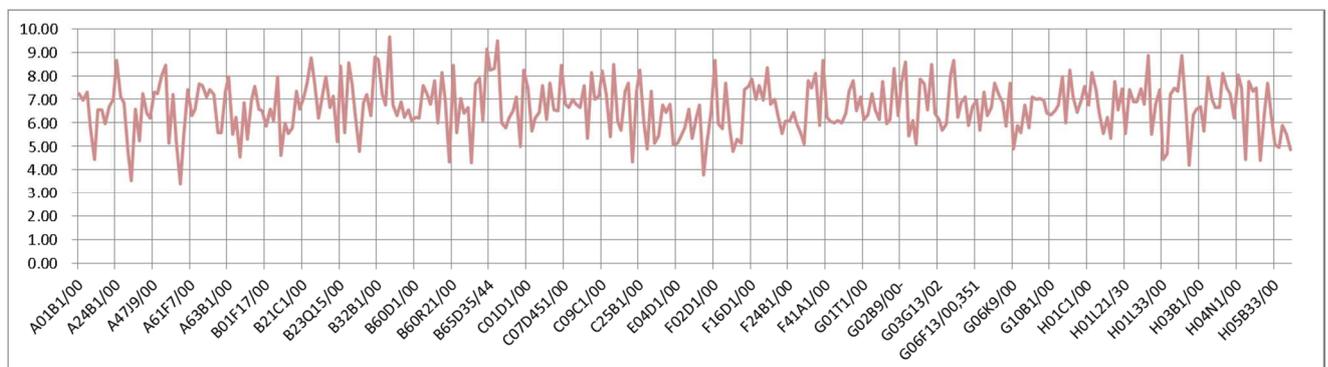
図 1-10：セクション単位の評価スコアの平均



グラフを見ると、平均スコア 7.56 の C セクションが最も評価が高く、5.80 の E セクションが最も低い。

この傾向を確認するため、各セクションを平均にまとめずに下図「筆頭 FI 順の評価スコア」を作成した。しかしながら、こちらの分析では、特段 C セクションのスコアが高く E セクションが低いという傾向は示されず、本評価においてはセクション単位での翻訳精度の良し悪しについて、明確な傾向は見られなかったといえる。

図 1-11：筆頭 FI 順の評価スコア



### 1.5.1.2 最も評価スコアが低い中 AU (DGT1-1)

一次評価の対象とした中 AU 全 362 分類で最も平均評価スコアが低かった分野は「DGT1-1 (タービン)」で、スコアは 3.40 (内容 1.60+用語 1.80) であった。なお、内容伝達の評価に大きく影響する中国語の平均文長は 39.1 文字であった。

表 1-12 : 中 AU 「DGT1-1」 一次評価対象文及び評価スコア (平均スコア最低)

	A	B	C	D	E	F	G	H	I	J	K	L	M	N
	中AU	テーマ	原文	機械翻訳文	基準翻訳文	重要技術用語(中)	重要技術用語(MT)	重要技術用語(日)	評価1(内容) 5~1	評価1(用語) A~D	評価1のコメント(任意) 評価根拠等	評価2(内容) 5~1	評価2(用語) A~D	評価2のコメント(任意) 評価根拠等
1	DGT1-1	3G070	燃气涡轮发动机还包括过滤器部件、至少一排可变入口导叶(VIGV)和入口空气流控制布置。	ガスタービンエンジンはまたフィルタ材、少なくとも1列の可変性入口案内翼(VIGV)と入口空気流制御配置を含んだことがある。	ガスタービンエンジンはまたフィルタ材、少なくとも1列の可変性入口ガイドブレード(VIGV)及び入口空気流制御装置を含む。	制御布置	制御配置	制御装置	2	C	「制御配置」では意味不明	2	C	エンジン入口の各要素を制御するための配置なのか、制御装置自体を指す
2	DGT1-1	3G070	排气系统是燃气轮机的重要辅助系统, 对其效率及正常工作有很大影响。	排気系はガスタービンの重要な支援システムであり、その効率および通常動作に対しても大きく影響がある。	排気システムはガスタービンの重要な補助システムであり、効率及び正常動作に対して大きく影響する。	燃气轮机	ガスタービン	ガスタービン	3	A	排気システムの影響度には正負両面が有り得るので「支援」は不適当	3	A	「補助」と「支援」ではシステムの役割は異なるので、正しく理解することは
3	DGT1-1	3G080	本发明涉及用于存储和利用来自间歇热源特别是来自太阳能收集器的能量的系统和装置。	本発明は記憶し利用して間欠の熱源の特に太陽エネルギーコレクタから来るエネルギーのシステムと装置とに関する。	本発明は間欠熱源の貯蔵及び利用に関し、特に太陽熱収集器からのエネルギーを利用するシステム及び装置に関する。	太阳能	太陽エネルギー	太陽熱	1	C	太陽光エネルギーと太陽熱エネルギーとの区別が不明確	1	B	誤訳のため類推が困難である。「太陽熱」は「太陽エネルギー」の部分集合
4	DGT1-1	3G081	本发明涉及余热发电技术领域, 尤其是涉及一种优化的有机朗肯循环低温余热发电系统。	本発明は余熱発電技術分野に関し、特に最適化に関する有機ランキンサイクル低温余熱発電システムである。	本発明は廃熱発電の技術分野に関し、特に最適化された有機ランキンサイクルによる低温廃熱発電システムに関する。	余热发电	余熱発電	廃熱発電	2	C	「余熱発電」は一般的ではない	2	C	「余熱発電」の訳が不適であるため、技術を正しく解釈することは難しい。
5	DGT1-1	3G104	本发明中, 所谓的“余热加热器”是指利用余热或大气对所述液体氧化剂加热的热交换器。	本発明中に、“余り蓄熱ヒーター”というものは余熱あるいは大気を利用することを指して前記液体酸化剤に加熱する熱交換器である。	本発明で、所謂“廃熱加熱器”は前記液体酸化剤に対し廃熱あるいは大気を利用して加熱する熱交換器を指す。	余热加热器	余り蓄熱ヒーター	廃熱加熱器	2	C	余り、余熱、の訳語が不適切、蓄熱の訳語は誤解を招く	1	C	蓄熱機能があるかの誤解を与えるおそれあり不適切。
6	DGT1-1	3G104	此外, 将返回水再引入给水系统中减少了由于燃料预热而在水侧上引起的总压力损失。	また、水に戻ることは再び給水系中に減少したことを導入して燃料予熱によって水側が引き起こす全圧に損失する。	また、復水を再度給水システムに戻し、燃料予熱によって水側に引き起こされる総圧力損失を軽減する。 ※水側は本件において「蒸気側」に対応することを許容する。	燃料预热	燃料予熱	燃料予熱	1	A	訳として、まるで意味をなさない。	1	A	全体像が不明で理解できない。
7	DGT1-1	3G104	同时, 回热器也有一个特点, 燃气的进口压力小, 接近于正常大气压, 密封要求不高, 容许少量泄漏。	同時に、蓄冷器は1個の特徴が同じくあり、燃焼ガスの入口圧は小さく、標準大気圧に接近し、密封要求は高くなく、少量が漏れることを許容する。	同時に、熱回収装置もガスの吸気圧が小さく、通常の大気圧に近く、シール性の要件は高くなく、少量の漏れは許容されるという特徴がある。	回热器	蓄冷器	熱回収装置	1	C	蓄冷器の訳語がまるで不適切	1	C	全体像が不明で理解できない。
8	DGT1-1	3G106	当矩形出口有主射流流出时, 将对两侧的气体进行引射抽吸使两侧产生低压区。	角型の出口にメインラジフロが流出がある時、両側のガスにエジェクタをすることを両側で低圧領域を発生することを吸引する。	矩形出口から主ジェット流が流出すると、両側の気体が吸引され、両側に低圧領域が形成される。	主射流	メインラジフロ	主ジェット流	1	C	誤訳および文法の間違い	1	C	訳語が不適切であり、機械翻訳結果から基準翻訳文を類推することは困難
9	DGT1-1	3G202	激光熔覆是目前较为先进的激光表面强化技术之一, 熔覆金属与母材为冶金结合, 强度高于传统钎焊。	レーザ溶着は現在比較的に先進的レーザ表面強化技術の一つであり、溶着金属と母材は冶金の結合とし、強度は従来のろう付けより高い。	レーザクラッピングは現在の先進的レーザによる表面強化技術の一つで、レーザクラッピング金属は母材と冶金結合し、従来のろう付けより強度が高い。	激光熔覆	レーザ溶着	レーザクラッピング	2	C	レーザ面強化技術の訳語が不適切で、文章の意味を阻害している	2	B	一般的説明ならば許容できる。
10	DGT1-1	3G202	优选地, 多孔材料中的单独的孔的孔尺寸包括大约5密耳与大约60密耳之间。	好ましくは、多孔質材料中の単独のホールの孔サイズは約5個のミルと約60個のミルの間を含む。	好ましくは、多孔質材料の単独の孔寸法は約5ミル〜約60ミルの範囲である。	多孔材料	多孔質材料	多孔質材料	2	A	5個のミル、60個のミルの訳語が不適切	1	A	ミルという用語を理解して翻訳すること
11														

### 1.5.1.3 最も評価スコアが高い中 AU (JJJ3-1)

一方、全中 AU で最も平均評価スコアが高かった分野は「JJJ3-1 (付加系高分子 (特殊))」で、スコアは 9.65 (内容 4.65+用語 5.00) である。内容伝達の評価に大きく影響する中国語の平均文長は 39.6 文字であり、最低平均スコアの DGT1-1 より 0.5 文字長かった。このことから、本評価で採用した文長範囲 (30~50 文字) の範囲内であれば、平均文長の違いによる有利不利はさほど考慮する必要はないと考えられる。

表 1-13 : 中 AU 「JJJ3-1」一次評価対象文及び評価スコア (平均スコア最高)

	A	B	C	D	E	F	G	H	I	J	K	L	M	N
	中AU	テーマ	原文	機械翻訳文	基準翻訳文	重要技術用語(中)	重要技術用語(MT)	重要技術用語(日)	評価1(内容) 5~1	評価1(用語) A~D	評価1のコメント(任意) 評価根拠等	評価2(内容) 5~1	評価2(用語) A~D	評価2のコメント(任意) 評価根拠等
1	JJJB-1	4.J026	上述的马来酸酐、烯基聚氧乙烯醚、丙烯酸酯的摩尔比为1:0.008~0.10:0.05~3;	上記の無水マレイン酸、アルケニルポリオキシエチレン、アクリル酸のモル比は1:0.008~0.10:0.05~3とする;	上記の無水マレイン酸、アリルポリオキシエチレンオレイルエーテル、アクリル酸のモル比は1:0.008~0.10:0.05~3とする;	马来酸酐	無水マレイン酸	無水マレイン酸	5	A		5	A	
2	JJJB-1	4.J026	作为优选,步骤(2)中,对于基体树脂为聚苯醚的接枝共聚物,温度为140至160℃。	好ましく、ステップ(2)にマトリックス樹脂に対してポリフェニレンエーテルのグラフト共重合体とし、温度は140~160℃とする。	好ましくは、ステップ(2)において、マトリックス樹脂はポリフェニレンエーテルのグラフト共重合体であり、温度は140~160℃である。	聚苯醚	ポリフェニレンエーテル	ポリフェニレンエーテル	5	A		5	A	
3	JJJB-1	4.J026	图1是本发明实施例1聚吡咯与亲水性聚苯醚二嵌段共聚物的红外吸收光谱图。	図1は本発明実施例1ポリピロールと親水性ポリフェニレンイソシアニドジブロックコポリマーの赤外吸収スペクトルグラフである。	図1は本発明実施例1のポリピロールとポリフェニレンイソシアニドのジブロック共重合体の赤外吸収スペクトルグラフである。	聚吡咯	ポリピロール	ポリピロール	5	A		5	A	
4	JJJB-1	4.J026	在另一优选例中,在所述步骤(ii)中,所述的极性溶剂选自下组:四氢呋喃、乙醚,或其组合。	もう一つであり例が好ましいことに、前記工程(ii)に、前記極性溶媒は次のグループから選ぶ:テトラヒドロフラン、ジエチルエーテル、あるいはその組み合わせ。	他の好ましい例では、前記ステップ(ii)において、前記の極性溶媒は次のグループから選択される:テトラヒドロフラン、ジエチルエーテル、あるいはその組み合わせ。	四氢呋喃	テトラヒドロフラン	テトラヒドロフラン	5	A		5	A	
5	JJJB-1	4.J026	另外,由于存在异氰酸酯基,可能导致亚克力压敏胶与离型材料的粘结,无法分离。	また、イソシアネート基が存在するので、アクリル感圧接着剤と離形材料の接着を引き起こす可能性があり、分離不可能。	更に、イソシアネート基が存在するために、アクリル感圧接着剤と離形材料の接着は、分離不可能になる可能性がある。	异氰酸酯	イソシアネート	イソシアネート	4	A	「可能性がある」と「分離不可能」の位置が逆転している。	3	A	
6	JJJB-1	4.J041	本专利涉及一种天然胶乳快速凝固剂的合成方法,能快速高效地凝固天然胶乳,同时降低成本。	本特許は天然ラテックス急速凝固剤の合成方法に關し、素早く効率的に天然ラテックスを固化することができ、同時に低コスト化。	本特許は天然ラテックスの急速凝固剤の合成方法に關し、迅速で効率的に天然ラテックスを凝固させることができ、同時にコストを低減させる。	凝固剂	凝固剤	凝固剤	4	A	日本語文として不自然な部部が若干あるが、技術用語はほぼ正確。	4	A	
7	JJJB-1	4.J127	任选地,低聚物R可进一步包含至少一种不同于其它两种共聚单体的附加共聚单体。	任意に、オリゴマーRは少なくとも1種をさらに含んで他の2種類のモノマーの付加モノマーと異なることができる。	任意選択的に、オリゴマーRは、さらに、他の2種類のモノマーの付加モノマーと異なる少なくとも1種を含むことができる。	共聚单体	モノマー	モノマー	4	A		5	A	
8	JJJB-1	4.J127	上述述混合单体由丙烯酸、甲基丙烯酸、丙烯酸酯和水按照重量份10.2:1.5混合而成。	上記の混合単量体はアクリル酸、メタクリル酸、アクリルスルホン酸ナトリウムと水重量分率の分10:2:1:5から混合される。	上記の混合単量体はアクリル酸、メタクリル酸、アクリルスルホン酸ナトリウムと水を重量部で10:2:1:5に従って混合される。	丙烯酸酯	アクリルスルホン酸ナトリウム	アクリルスルホン酸ナトリウム	5	A		4	A	「水重量分率の分~」が理解しにくい。
9	JJJB-1	4.J127	作为脂环族多元酸,例如,可使用以下所举的脂环族二羧酸及其衍生物等。	脂環族系多塩基酸として、例えば、以下が挙げられた脂環式カルボン酸およびその誘導体などを用いることができる。	脂環族多塩基酸としては、例えば、以下に例示する脂環式カルボン酸及びその誘導体などを用いることができる。	衍生物	誘導体	誘導体	5	A		5	A	
10	JJJB-1	4.J127	其中,所述阻聚剂为4-甲氧基酚、对苯二酚或4-甲氧基酚与对苯二酚组成的混合物;	そのうち、前記重合防止剤は4-メトキシフェノール、ヒドロキノンあるいは4-メトキシフェノールとヒドロキノンが組成した混合物とする;	そのうち、前記重合禁止剤は4-メトキシフェノール、ヒドロキノンあるいは4-メトキシフェノールとヒドロキノンで構成された混合物である;	阻聚剂	重合防止剤	重合禁止剤	5	A		5	A	
11	JJJB-1	4.J127												

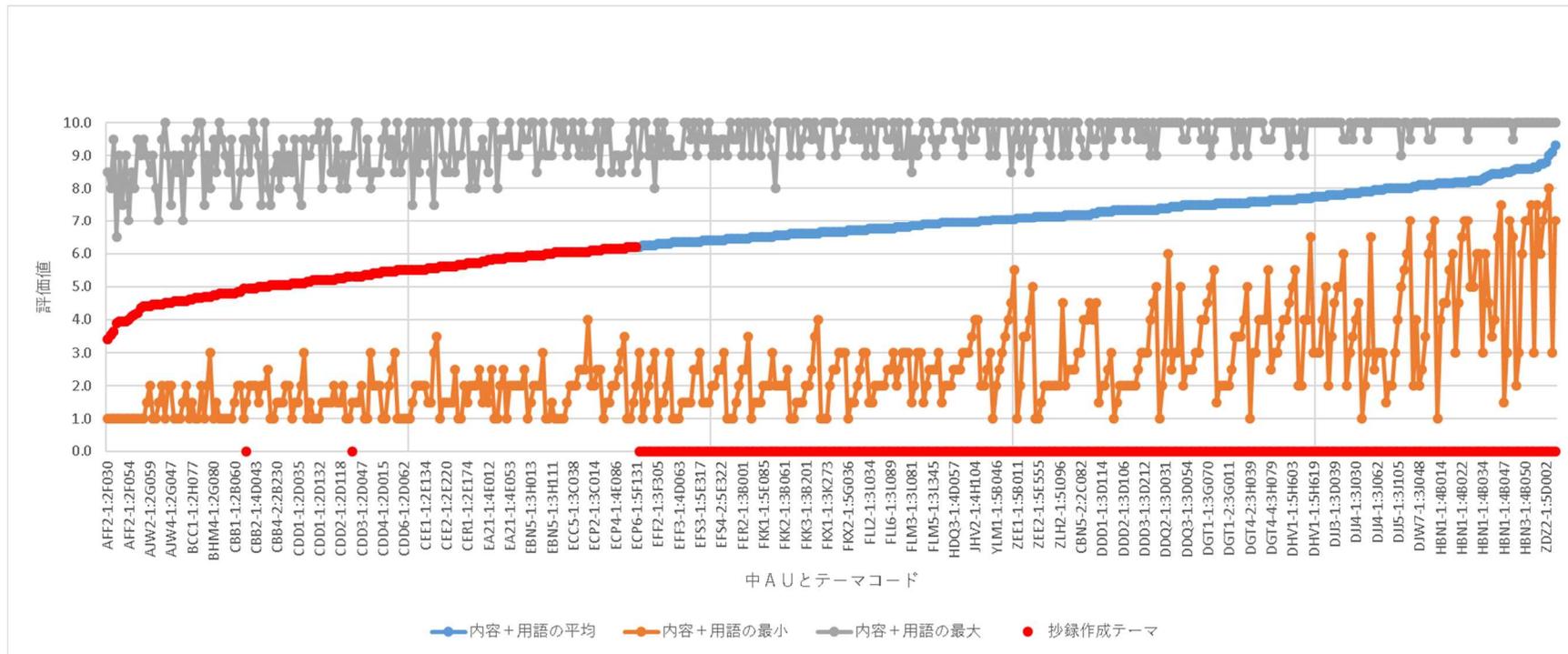
## 1.5.2 二次評価結果の分析

### 1.5.2.1 テーマコードの平均評価スコアと抄録作成対象（174テーマ）の関係

二次評価における 480 テーマの平均評価スコアと抄録作成対象テーマ（174テーマ）との関係を示す。平均評価スコアが赤色で示されているものが和文抄録作成対象となるテーマである。

このグラフには各テーマの平均スコアに加えて最大と最小スコアを示した。全体的に最大スコアと最小スコアには非常に大きな開きがある。一つのテーマの評価文数は 10 文であるため、採用された評価対象文における難易度の偏りが、平均評価スコアに大きく影響している可能性がある。

図 1-14 二次評価テーマ 480 テーマの品質下位順グラフ



### 1.5.2.2 翻訳品質が低いテーマコードの分析① 重要技術用語の翻訳精度

本評価では、「内容の伝達レベル」と「重要技術用語の翻訳精度」の評価スコアを意図的に1:1の重みづけで合算して総合評価としている。このため、各文で「重要技術用語」として選定された語の翻訳精度は、それだけで総合評価の50%を占め、さらに当然ながら「内容の伝達レベル」の評価にも加味されるため、双方の観点でダブルカウントされることとなる。このため、本評価において「重要技術用語」の翻訳精度は、翻訳品質の優劣評価に非常に大きな影響力をもつ。

これを踏まえ、二次評価の対象とした480テーマ中、まずは「重要技術用語」の平均スコアが1点台の16テーマについて、その傾向を分析した。分析対象となったテーマ及び各テーマの平均スコアを下表1-15に示す。

表 1-15 「重要技術用語の翻訳精度」2点未満の16テーマ

順位	中 AU	テーマ	用語	内容	総合
1	ECP2-1	3C014	木材用鋸盤の構成部品、付属品	1.95	3.40
2	CDD4-1	2D013	ショベル系を除いた土砂掘削機及び施工法	2.30	3.55
3	ECC5-1	3C065	はさみ・ニッパ	1.75	3.65
6	CEE2-1	2E108	屋根ふき・それに関連する装置または器具	2.20	3.95
12	FLM5-1	3L020	被冷蔵物の充填、照明装置	2.40	4.35
14	HDQ8-1	4D071	液体又は風力による固体相互の分離	2.90	4.40
15	EFF2-1	3F312	ジャッキ	2.50	4.40
17	CDD6-1	2D054	立坑・トンネルの掘削技術	2.85	4.45
18	CDD3-1	2D049	基礎工事に適用される隔壁	2.70	4.45
22	CBB4-3	2B150	飼料（2）（一般）	2.90	4.50
27	CDD5-1	2D059	橋または陸橋	2.80	4.55
31	EFS1-1	3C030	自動組立	2.90	4.65
33	CER1-1	2E044	はしご	3.50	4.70
35	FER2-1	3K099	衣服及び関連具のホルダー	3.60	4.70
45	DDD2-1	3D106	鉄道車両の車体細部・付帯設備	2.95	4.85
49	DDD1-1	3D038	推進装置の冷却、吸排気、燃料タンクの配置	3.10	4.95

これら16テーマは当然ながら総合評価も低く、上表に示したとおり、下位1位～3位を占めるほか、残りの各テーマも全て下位50位以内となっている。

これら 16 テーマにおける「重要技術用語」の訳語（機械翻訳結果）を概観した印象として、以下が挙げられる。

#### ①名詞と解釈されなかった重要技術用語を含む

16 テーマの多くは、重要技術用語が複合名詞と解釈されなかったケースが複数件含まれていた。以下、そのケースの実例をいくつか挙げる。

表 1-16：重要技術用語が複合名詞と解釈されなかった実例

テーマ	中国語	基準翻訳	機械翻訳
3C014	多片出条机	チップソー切断機	高く数枚を転送して条機
2D013	挖泥深度	浚渫深さ	削り泥の深さ
3C065	碰触块	接触ブロック	ブロックに当接する
2E108	绿色建筑	グリーンビルディング	緑色に・・・建築する

このように、複合語が複合名詞と解釈されず、単語レベルで個別に異なる品詞に解釈されてしまうような場合、その悪影響はその複合名詞のみにとどまらず、文全体の翻訳精度を著しく損なうことは言うまでもない。たとえば上記 4 例は「重要技術用語の翻訳精度」は当然 C 評価（0 点）であるが、これらの語を含む文全体の「内容の伝達レベル」の評価もそれぞれ 1 点、2 点、1 点、1 点と非常に低い。そして、このような用語が重要技術用語として選ばれた場合、二つの観点（「重要技術用語の翻訳精度」と「内容の伝達レベル」）において厳しく減点されるため、こうしたケースが複数件発生したようなテーマは、おのずと評価順位が下位に偏ってくるものと考えられる。上記結果は、これを裏付けるものである。

なお、評価対象文には技術用語が複数個含まれることも多く、また重要技術用語の選定は恣意的に「機械翻訳の難しそうな」用語を選んだわけではないため、このように複合名詞と解釈されない用語が、評価対象文中に含まれるものの、重要技術用語としては選定されていないというケースも当然想定される。ただしこうした場合も、「内容の伝達レベル」での大幅な減点は免れないため、10 文中にこのような用語を多く含むテーマは、自然と評価下位に集まるものと考えられる。

## ②重要技術用語の単純な誤訳

重要技術用語を複合名詞と正しく解釈はしている（つまり、品詞自体は正しく解析できている）ものの、訳語が不適切なため C 評価というケースも、当然ながら下位テーマには多数存在している。以下、その一例を示す。このパターンの誤訳の場合、少なくともこれが直接の原因となって文全体の構文解析を崩すことはないため、理論上は、「内容の伝達レベル」を含めた総合評価への悪影響は上記①のパターンより小さいといえる。

表 1-17：重要技術用語が複合名詞のまま誤訳された実例

テーマ	中国語	基準翻訳	機械翻訳
3C014	锯齿结构	鋸齒構造	ジグザグ構造
2D013	反冲管路	逆洗配管	キック管路
3C065	紧固齿轮	締め付けギヤ	固着歯車
2E108	木牌科	木造組物	木札科

ただし、文の主題となる重要技術用語が大きく誤訳されると、その時点で（文法的に成立している文でも）内容の理解は難しくなることも多く、結果的には「内容伝達レベル」が①のパターンと同等程度に低評価となる可能性は高い。実際のところ、「重要技術用語のみが誤訳でそれ以外の部分は正確」というケースはほとんど存在せず、全体的には「重要技術用語」が低評価の文は「内容伝達レベル」も低い評価に止まる傾向にあった。

## ③重要技術用語が未知語であるケースは少ない

機械翻訳において、未知語は翻訳精度を大きく損なう主要な要因の一つであり、これは中日翻訳においても原則として当てはまる。ただし中日翻訳の場合、漢字から漢字への翻訳が大半であるため、未知語と誤訳との境界が曖昧で区別しにくいこと（例えば上記②に示した「木札科」などは、未知語なのか誤訳なのか微妙）、また、同じ漢字なので未知語であっても意味がかなり把握できる場合もあるなど、未知語のみに限定した悪影響を抽出することはかなり困難である。

ただし本調査の場合、二次評価結果にて「未知語（D 評価）」と判定された重要技術用語は 4,800 語中 50 語程度（C 評価は 1,000 語以上）ときわめて少なく、「中韓文献翻訳・検索システム」においては未知語の発生率はかなり低い。このため、評価結果全体への影響も局所的なものにとどまっていると見なせる。

### 1.5.2.3 翻訳品質が低いテーマコードの分析② 内容伝達レベル

内容伝達レベルが2点未満の評価となったテーマは下表に示す12テーマである。このうち「重要技術用語」の評価も2点未満であったテーマは2テーマのみであるが、他の10テーマも大半は「重要技術用語」は2点台に止まっており、全体平均(3.43)に比して決して高くはない。

表 1-18: 「内容の伝達レベル」2点未満の12テーマ

順位	中 AU	テーマ	内容	用語	総合	
1	ECP2-1	3C014	木材用鋸盤の構成部品, 付属品	1.95	1.45	3.40
3	ECC5-1	3C065	はさみ・ニッパ	1.75	1.90	3.65
4	CDD4-1	2D065	さく岩、採鉱及び採鉱機械とその方法	1.15	2.75	3.90
5	ECP7-1	3B221	装飾技術	1.95	2.00	3.95
7	CDD4-1	2D003	掘削機械の作業制御	1.70	2.25	3.95
8	CDD4-1	2D012	ショベル系 (制御を除く)	1.95	2.05	4.00
10	ECC5-1	3C068	可搬型釘打ち機およびステープラー	1.55	2.60	4.15
16	FKK2-1	4L029	アイロン	1.80	2.65	4.45
23	DGT1-1	3G106	ジェット推進設備	1.65	2.90	4.55
26	EFF2-1	3F205	ジブクレーン (門形、ケーブルクレーン)	1.70	2.85	4.55
34	ECC5-1	3C066	湿式かみそり	1.80	2.90	4.70
76	EFF2-1	3F204	クレーンの細部 (制御、安全)	1.85	3.35	5.20

内容伝達レベルの評価に影響する要素はさまざまであり、10文のサンプルからその傾向を特定することは非常に難しい。ただし、本調査の場合、各テーマとも同じ翻訳エンジンを使用しており、翻訳のロジック自体は共通しているため、各テーマ同等程度の難易度の文がサンプリングされているという前提に立てば、各テーマの評価の差は、主にそこに含まれる技術用語を正しく解析・翻訳できているかの差と考えることができる。上記12テーマにおける「重要技術用語」の評価スコアはおおむね1~2点台であり、480テーマ全体の平均スコア3.43よりも有意に低いことから、この考えが裏付けられる。また、評価文中において技術用語は「重要技術用語」として選ばれた一語以外にも存在することが大半であり、こうした「重要技術用語以外の技術用語」の翻訳精度の良否も、「内容の伝達レベル」の評価を左右する主要な要因となっていることは確実といえる。

例えば、内容伝達レベルの評価が最も低かったテーマ2D065(スコア1.15)は、10文全てが1点~2点という低評価であった。このうち4文は「重要技術用語」はA評価を得ているが、例えば下の例のように、それ以外の技術用語に訳語が不適切(例えば「プラウ

による採炭法」⇒「ホーベル採炭法」、「ドリル式採炭機による採炭法」⇒「式採炭機採炭法…を穿孔する」) なものがあり、これが低評価の一因となっている。

表 1-19: 「重要技術用語」が A 評価に対し「内容伝達レベル」が低評価の実例

原文	薄煤层综合机械化开采方法可分为三种：刨煤机采煤法、钻式采煤机采煤法和滚筒采煤机采煤法。
基準翻訳	薄炭層の完全機械化による採掘方法は：プラウによる採炭法、ドリル式採炭機による採炭法及びドラム式採炭機による採炭法の三種類に分類される。
機械翻訳	薄炭層総合機械化採掘方法は3種類可分であってなる：ホーベル採炭法は、式採炭機採炭法とドラム採炭機採炭法を穿孔する。
重要技術用語	薄煤层
基準翻訳	薄炭層
機械翻訳	薄炭層

もちろん、用語が原因ではなく、その分野特有の文体の癖などによって翻訳精度が低くなっているというケースも考えられる（上例では「可分であってなる」の箇所）。ただし、文体や文構造が原因の誤訳は、原則としてエンジンのロジックを改善しない限り辞書登録では解決不可能なことが多い。このため本評価では、評価スコアの集計において重要技術用語の比重を相当に重くしており、その分、構文解析の比重は比較的軽くなっている。よって、同じ程度の内容伝達レベルの文であれば、用語の訳質が原因であるもののほうが、それ以外の原因のものよりも下位か、少なくとも同等に扱われる仕組みとなっており、全体的に見れば、評価スコア下位のテーマは、文中に含まれる技術用語の翻訳精度が相対的に低いテーマであると見なせる。

#### 1.5.2.4 翻訳品質が低いテーマコードの分析③ 全体的な傾向

ここまで「重要技術用語の翻訳精度」と「内容の伝達レベル」それぞれの観点で特に低い評価を受けたテーマを中心に、その傾向を分析した。その結果、「重要技術用語の翻訳精度」は当然として、「内容の伝達レベル」の評価が低い文についても、「重要技術用語」をはじめとする文中の技術用語の翻訳精度が低いケースが多いことが示された。

実際のところ、出現する技術用語がすべて正しく翻訳されている文が「内容の伝達レベル」において著しく低い評価となることは考えにくく、全体的な傾向として、評価の低いテーマは概して技術用語の翻訳精度に難あり、と考えるのが妥当であろう。

特に、「重要技術用語の翻訳精度」の評価が顕著に低かったテーマでは、複合名詞である「重要技術用語」がそう解釈されず、単語に分解されそれぞれが誤った品詞に解釈され

ているケースが目立った。このようなケースでは「内容の伝達レベル」も連動して低い評価となっており、このようなケースが多発したテーマが総合評価下位となる傾向が見て取れた。

「中韓文献翻訳・検索システム」に用いられるルールベース方式の機械翻訳エンジンは、文を構成する各語の品詞解釈を文法規則に則って行うため、この「複合名詞を誤って他の品詞と解釈してしまう」タイプの誤訳が発生すると、それは局所的な誤りにとどまらず、文全体の構文解析に悪影響を及ぼすことが多くなる。なお、このような品詞誤りは、重要技術用語に選ばれなかった用語にも同様に発生する可能性があり、「内容の伝達レベル」が低評価のテーマには、こうしたケースも数多く含まれるものと考えられる。

中国語は、同じ用語が複数の品詞に解釈可能なことが多く、これが上記パターンの誤訳を多発させる要因の一つである。したがって、「中韓文献翻訳・検索システム」において相対的に翻訳品質が低いテーマに対し、その分野で多用される複合名詞の辞書登録を進めることは、技術用語の翻訳精度を向上させるという局所的な改善効果のみならず、その用語に正しい品詞解釈を与える（複合語の場合、単語に比べて品詞解釈の選択肢は限定される）という意味においても、きわめて有効な品質改善策である。

#### 1.5.2.5 技術分野ごとの傾向の分析

本評価では、一次評価では中 AU 単位、二次評価ではテーマコード単位で、技術分野ごとの機械翻訳精度の傾向を把握したが、ここでは、より広い技術分野ごとの傾向を分析した。

各テーマコードは中 AU に属し、各中 AU は 2B～5Z に区分された「技術単位」に属する。さらに、この「技術単位」は、上 1 ケタ (2～5) で集約することで、特許庁内審査部門ごとの 4 つの「担当技術」に区分することができる。

分析では、下記①～③の 3 つの条件それぞれにおいて、上記 4 つの「担当技術」別の構成比を比較することで、各「担当技術」にあたる技術分野の機械翻訳品質の傾向の調査を試みた。

- ① 件数規模（中国文献数）上位 1,600 テーマにおける構成比
- ② 二次評価対象 480 テーマにおける構成比
- ③ 評価スコア 6.20 以下（和文抄録作成対象の閾値）の 177 テーマにおける構成比

以下、4 つの担当技術単位の①～③の構成比を示す。

表 1-20: 「担当技術単位 (1~4)」別テーマコード構成比

担当技術	①1,600 テーマ		②480 テーマ		③177 テーマ	
1	378	23.6%	132	27.5%	58	32.8%
2	602	37.6%	254	52.9%	101	57.1%
3	332	20.8%	61	12.7%	10	5.6%
4	288	18.0%	33	6.9%	8	4.5%
合計	1,600		480		177	

上表を見ると、①の件数上位 1,600 テーマにおいては、担当技術 1~4 の構成比は、2 が約 4 割、他の 3 つはそれぞれ 2 割前後と、比較的均一な構成比であるのに対し、一次評価結果を踏まえた②、二次評価結果を踏まえた③と、対象範囲が品質下位に絞られていくにつれ、担当技術 1 と 2 が比率を大きく上げ、その一方で担当技術 3 と 4 は比率を大きく下げている。結果、和文抄録作成対象テーマとなる③の時点では、「担当技術」1 と 2 で 177 テーマ中 159 テーマと、その大部分を占める結果となった。

この集計結果からは、現状の「中韓文献翻訳・検索システム」においては、担当技術 1 と 2 は、担当技術 3 と 4 よりも相対的に機械翻訳品質が低い傾向にあると言える。

なお、「技術単位 (2B~5Z)」ごとの該当テーマ数や比率の集計結果については、下記電子ファイルを参照されたい。

[『05 二次評価の技術単位別集計結果.xlsx』](#)

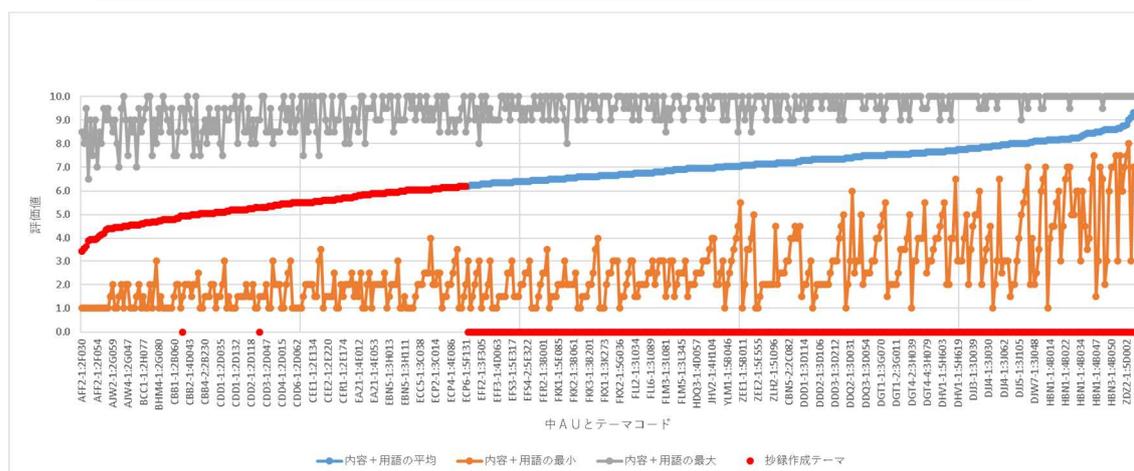
### 1.5.3 補足的な分析

#### 1.5.3.1 評価対象文における特異値の影響

本評価は一次評価、二次評価とも一つの単位の評価文数が10文と少ないため、採用された評価対象文における難易度の偏りが評価結果に影響している可能性がある。そこで、各テーマの10文から特異値を排除する目的で評価スコアが最高と最低のものを2文ずつ除外した6文で判定した場合と10文全て使った場合を比較して、抄録作成対象と判断されるテーマにどの程度、違いがあるのかを調査した。

下記は各テーマ10文で各テーマの評価スコア平均を算出して低い順に並べたグラフである。中央の曲線が平均スコア、上方のグレーの折れ線がテーマごとの最高スコア、下方のオレンジの折れ線が最低スコアを示す。平均スコアの曲線が赤い部分に相当するテーマが、和文抄録作成対象となる174テーマである。

図 1-21 二次評価テーマ 480 テーマの品質下位順グラフ (図 1-14 の再掲)



次に各テーマの10文から評価スコアが最高と最低のものを2文ずつ除外した6文での平均評価スコアを算出し、評価スコアの低い順に並べて和文抄録作成対象となる174テーマがどう変化するかを調査した。

その結果、評価文数10文と6文とで和文抄録作成対象となるテーマの入れ替わりは9テーマ(10文での対象テーマ3H036、3C038、5E344、3K014、2G001、2B033、3L046、4D021、2G070が、6文では5G016、3D050、2B230、3G105、3C017、3L059、3L081、2B260、2E002に入れ替わった)にとどまった。入れ替わり率は174テーマ中の9テーマで5.2%となり、94.8%は不変であった。また、入れ替わりのテーマは何れも閾値近辺に分布するもの

であった。

この結果から、480 テーマから翻訳精度の相対的に低い 174 テーマを選定するという今回の二次評価では、各テーマとも特異値の影響は許容範囲内であり、10 文のサンプリングにて妥当な精度が得られたと考えられる。

### 1.5.3.2 重複評価の相関について

本評価では、一次評価、二次評価とも、また「内容の伝達レベル」「重要技術用語の翻訳精度」の二つの観点ともに、二名の評価者による重複評価を行った。

なお、評価の実施にあたっては、1.2.4.3 に記載したとおりの各種手段にて事前の基準合わせを実施し、評価者間の基準を極力整合させるよう努めた。

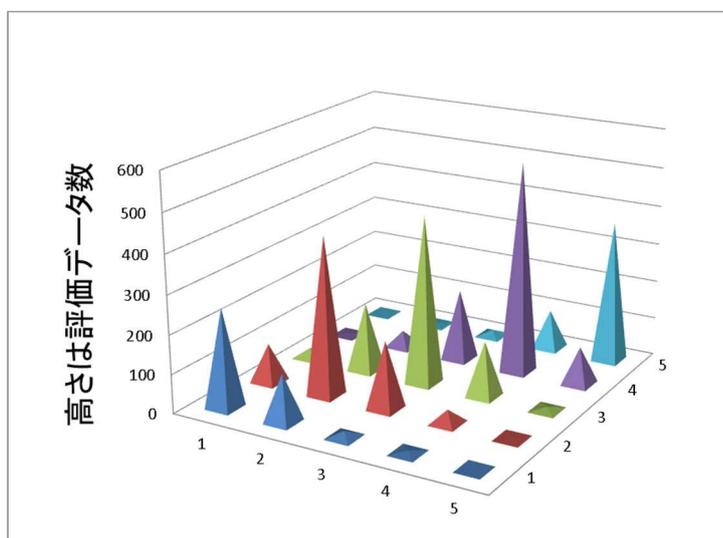
以下、本評価にて実施した重複評価における二名の評価者間の相関性について分析した。

#### ①内容伝達レベルの評価者二名の相関（相関係数：0.82）

内容伝達レベルの評価の二名の評価者の相関係数は 0.82 で強い相関が認められる。なお、全体 3,260 文の内 2,008 文(61.6%)は両者が一致した。二名の評価者の評価スコアの差が 2 以上の行数は 155 文(4.75%)にとどまった。

図 1-22：内容伝達レベルの評価 評価者二名の相関

評価者1→ 評価者2↓	1	2	3	4	5
1	255	123	16	6	0
2	99	409	174	35	1
3	25	177	433	142	8
4	3	42	184	545	97
5	0	3	16	101	366

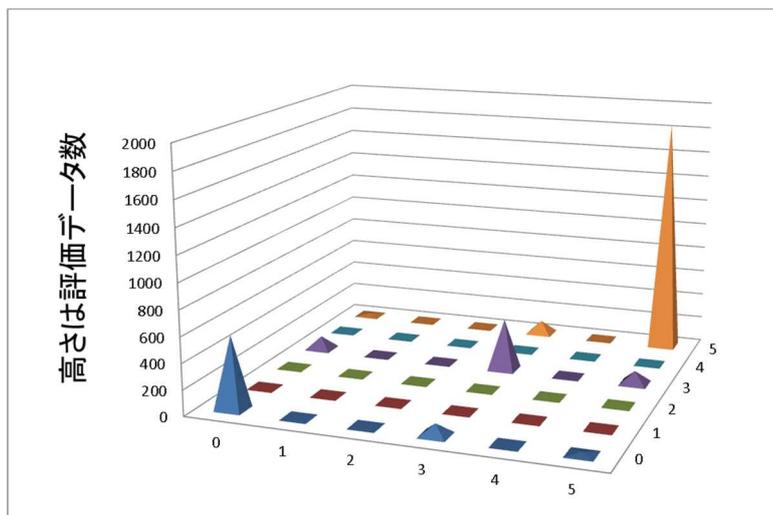


②重要技術用語の評価の評価者二名の相関（相関係数：0.87）

重要技術用語の評価の二名の評価者の相関係数は 0.87 で強い相関が認められる。全体 3,260 語の内 2,827 語(86.7%)は両者の評価が一致した。二名の評価者の評価スコアの差が A と C,D、B と CD のように大きく異なるものは 229 文(7.02%)であった。なお、評価スコアは A を 5 点、B を 3 点、C と D を共に 0 点に換算している。

図 1-23：重要技術用語の評価 評価者二名の相関

評価者1→ 評価者2↓	0	1	2	3	4	5
0	568	0	0	89	0	18
1	0	0	0	0	0	0
2	0	0	0	0	0	0
3	104	0	0	407	0	101
4	0	0	0	0	0	0
5	18	0	0	103	0	1852



本評価では、「内容の伝達レベル」、「重要技術用語の翻訳精度」の双方を二名による重複評価にて実施したが、いずれも高い相関が見られた。このことは、本評価の評価結果が、事前の基準合わせの効果により信頼性の高いものとなったことを示すと同時に、本評価に用いた評価基準が、評価者によるばらつきが出にくく、高い客観性を有するものであることを示している。

### 1.5.3.3 ファミリー文献数と評価スコアの関係

本調査の評価対象である機械翻訳結果はルールベースの機械翻訳結果である。ルールベース翻訳の重要な要素である中日用語辞書は中国と日本のパテントファミリーをソースとして作成されることが多い。このため、パテントファミリーが多い分野は中日辞書が充実しているため翻訳品質が高く、パテントファミリーが少ない分野は相対的に翻訳品質が低いという仮説が成り立つ。

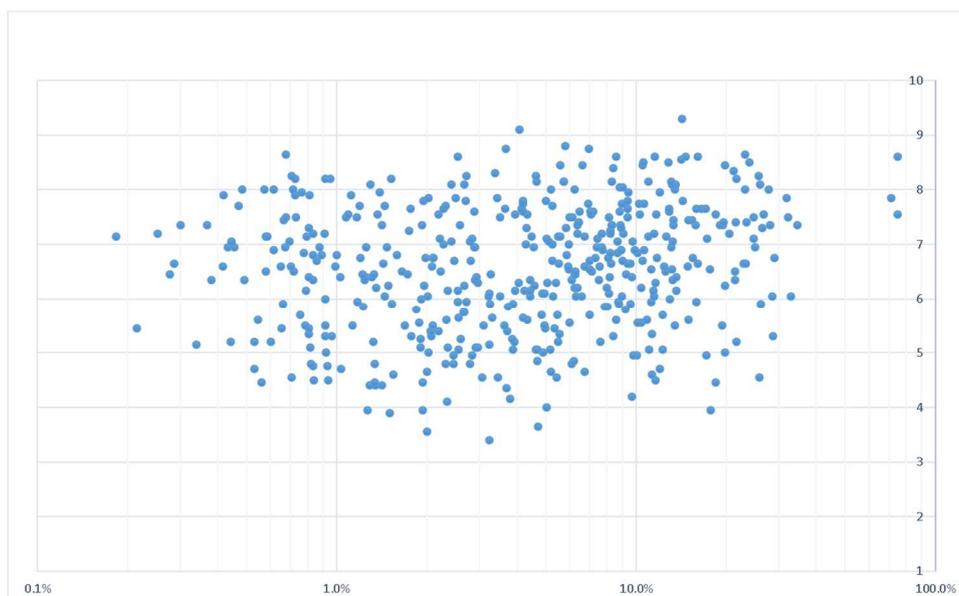
そこで、二次評価対象の 480 テーマの評価スコアと中日ファミリー文献との関係を調査した。

#### ①480 テーマの評価スコアとファミリー率の関係（相関係数 0.21）

二次評価対象 480 テーマの評価スコア（内容+用語）の平均とパテントファミリーの関係を図 1-24 にてグラフ化した。

縦軸が評価スコア、横軸がパテントファミリー率（そのテーマに属する全中国文献中、日本のファミリー文献を有するものの割合）である。当該テーマの全文献が日本特許のファミリーを持つ場合が 100%となる。大部分が 10%以下であるため、横軸を対数目盛とした。二つの値のピアソンの相関係数は 0.21 で、相関は弱い。

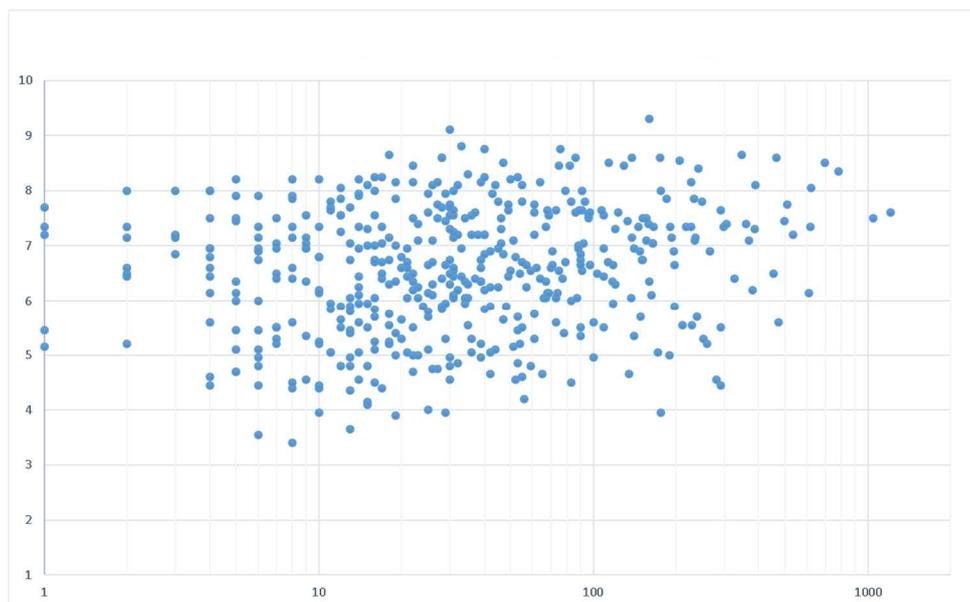
図 1-24：二次評価対象 480 テーマの平均スコアとファミリー率の相関



## ②480 テーマの評価スコアとファミリー文献数の関係 (相関係数 0.181)

縦軸が平均評価スコア、横軸が日本のパテントファミリー文献数である。日本のファミリー文献数が 100 件以下のテーマが多かったため、横軸を対数目盛とした。二つの値のピアソンの相関係数は 0.181 であり、やはり相関は認められない。

図 1-25 : 二次評価 480 テーマの平均スコアとファミリー文献数の相関 (対数目盛)

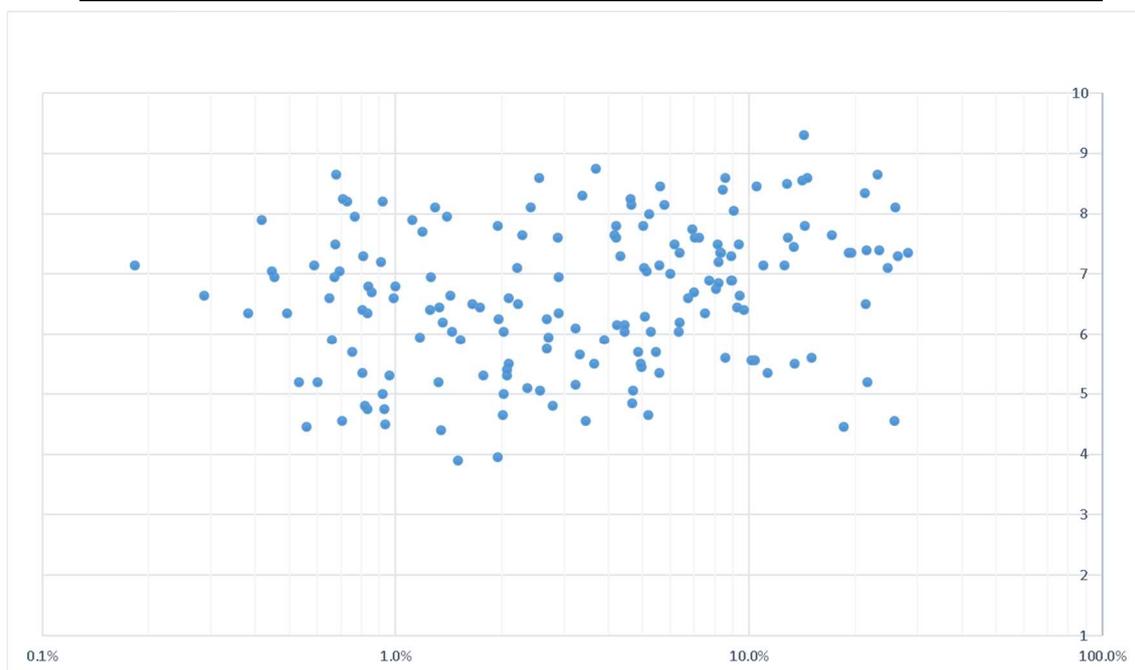


③中国文献数 $\geq 1,000$ の173テーマの評価スコアとファミリー率の関係（相関係数 0.23）

テーマの中国文献数（ファミリー文献数と非ファミリー文献数の合計）が1,000件以上の173テーマに限定して二次調査の評価スコア（用語＋内容）の平均とパテントファミリー率の関係をグラフ化した。

二つの値のピアソンの相関係数は0.23であり、弱い相関ではあるが①の全体での相関に比べて若干強くなった。

図 1-25：中国文献数 1,000 件超の 173 テーマの評価スコア／ファミリー率の相関

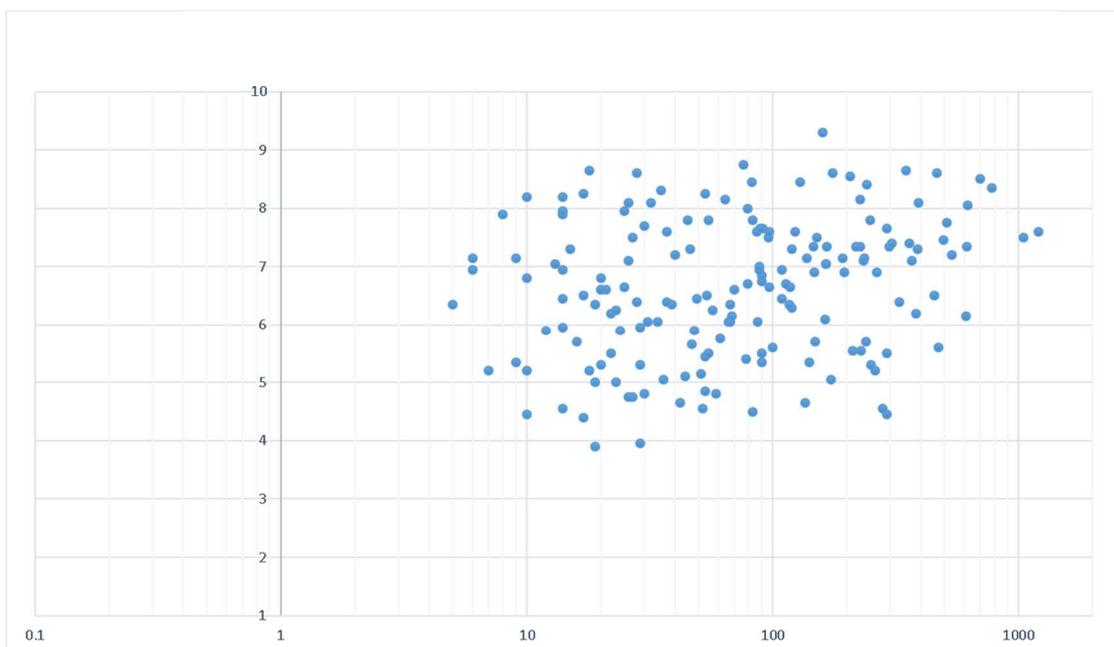


④中国文献数 $\geq 1000$ の173テーマの評価スコアとファミリー文献数の関係(相関係数0.25)

テーマの中国文献数が1,000件以上の173テーマに限定して二次調査の平均評価スコアと、パテントファミリーの文献数の関係をグラフ化した。

結果として、相関係数が0.18から0.25に向上した。このことから、ファミリー文献数がある程度あれば、既存の中日用語辞書の充実が進んでいるとみなせる可能性がある。

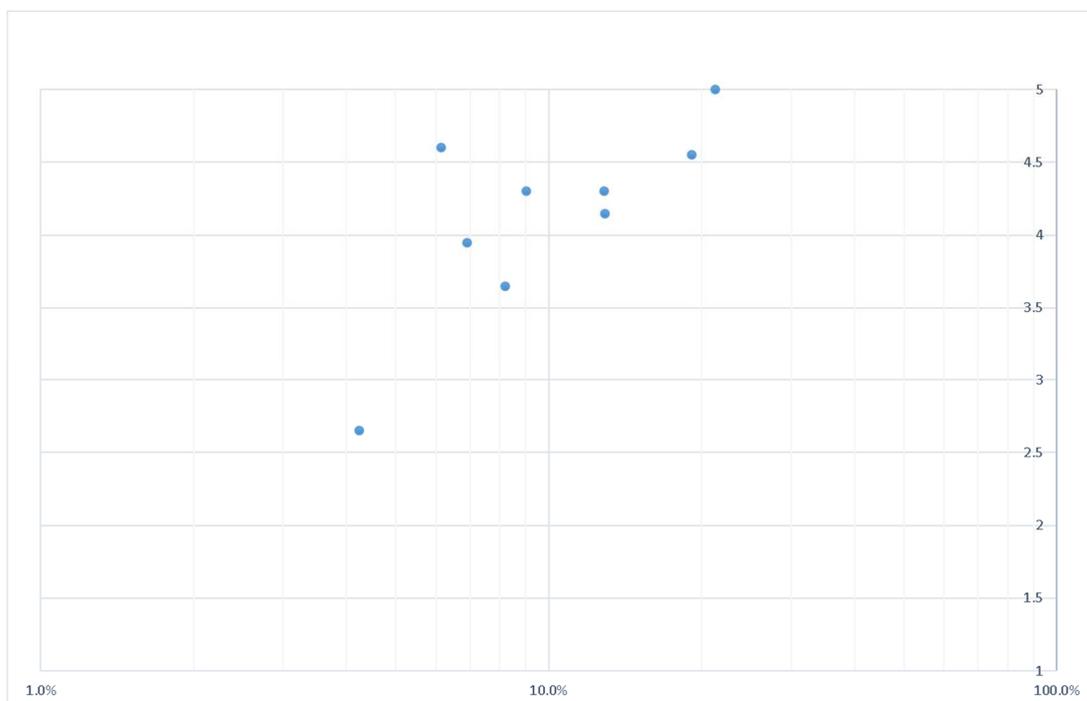
図 1-26 : 中国文献数 1,000 件超の 173 テーマの評価スコア / ファミリー文献数の相関



⑤ファミリー文献数 $\geq 500$  の 9 テーマの重要技術用語の評価スコアとファミリー文献数の関係 (相関係数 0.70)

ファミリー文献数と、既存の用語辞書の充実との関連性を確かめるため、ファミリー文献数が 500 以上の 9 テーマについて「重要技術用語」の評価スコアとファミリー文献数の割合の関係を調査したところ、相関関係は 0.70 に高まった。また、その平均評価スコアは 9 テーマ中 6 テーマにおいて 4.0 以上と高かった。このことから、ファミリー文献数が 500 以上あるテーマは、用語辞書はある程度充実していると思わせる可能性があるが、サンプル数が極端に少ないため、十分に証明されたとはまではいえない。

図 1-27 : ファミリー文献数 500 件超の 9 テーマの用語評価スコアとファミリー文献数の相関



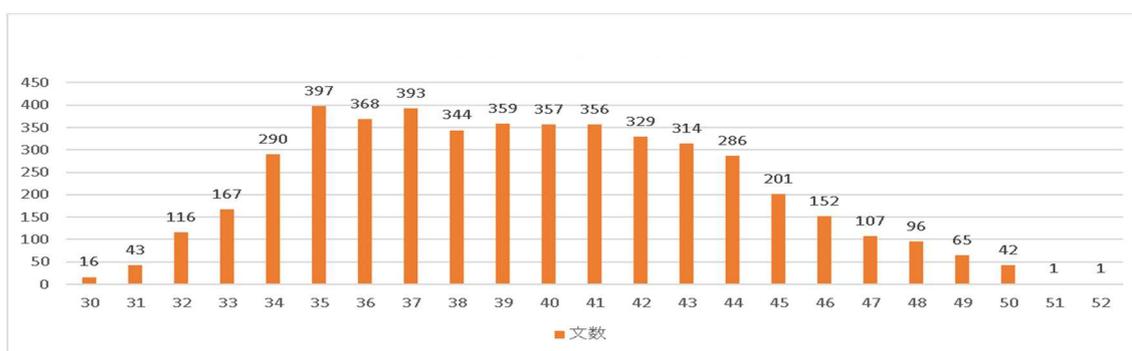
上記①～⑤について考察を行ってきたが、ファミリー率やファミリー文献数と翻訳精度との間に十分な相関があることは見い出せず、結論として、ファミリー率やファミリー文献数だけでは翻訳精度に課題のあるテーマを特定することは困難であることがわかった。

### 1.5.3.4 評価対象文の文長と評価スコアの関係

本調査では、評価対象文（原文）の文長は原則 30 文字～50 文字の範囲に統一した。この範囲における文長と評価スコアとの相関性を見るべく、二次評価の対象全 4,800 文の文の長さ「内容の伝達レベル」の平均評価スコアとの関係をグラフ化した。

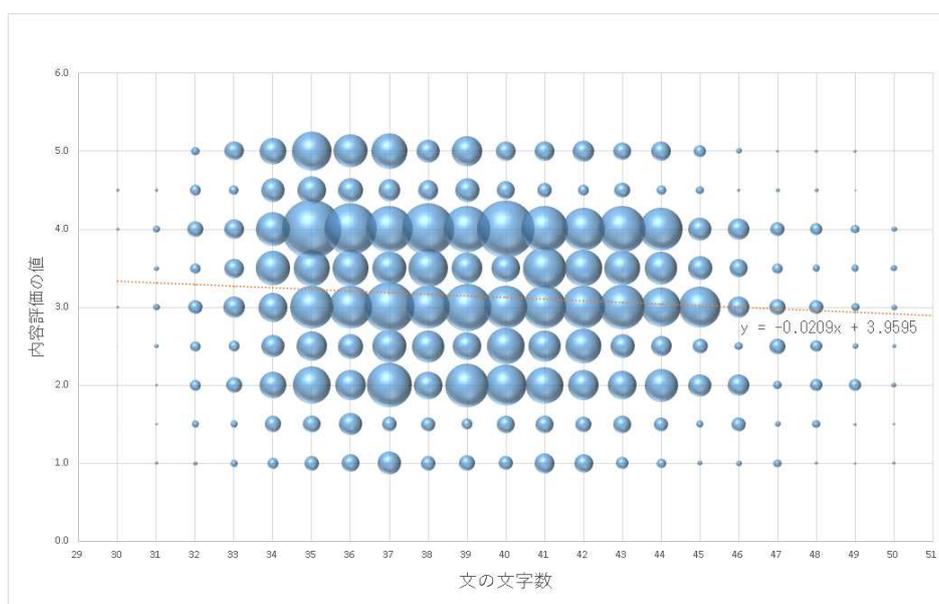
下のグラフは縦軸が評価スコア、横軸が原文の文字数（文長）である。棒グラフは各文字数に該当する文数を示している。

図 1-28：評価対象文の中国語文字数と内容伝達レベルの評価スコアの関係



一方、下記のバブルチャートは縦軸が「内容伝達レベル」の平均評価スコア、横軸が文の長さ、バブルの直径が当該ドットに対応する文数を示している。橙色の点線は近似曲線  $y = -0.0209x + 3.9595$  である。

図 1-29：評価対象文の中国語文字数と内容伝達レベルの評価スコアの関係（バブルチャート）



二つの値のピアソンの相関係数は-0.08 で、文長が長くなるにしたがって評価スコアが低くなることを示しているが、その傾斜はごくわずかである。また、上記バブルチャートでは大半の文が文字範囲のコア（35～45 文字）に集中しており、かつどの文字数においても概ね同様の評価スコア分布を得ていることが見て取れる。

これらのことから、少なくとも今回の評価においては、各テーマの翻訳品質の比較に文長は大きく影響していないことがわかった。したがって、本評価に用いた 30～50 文字という文長範囲設定は妥当であったとみなせる。

## 2. 特定テーマコードの中国特許公報の和文抄録の作成

機械翻訳文の翻訳品質評価に基づき、174 のテーマコードに属する中国特許文献 8 万件を選定し、和文抄録を作成した。以下、その概要を記す。

### 2.1 和文抄録の作成対象

和文抄録の作成対象テーマおよび作成対象文献の選定基準については、「1. 中韓文献翻訳・検索システムにおける機械翻訳文の翻訳品質評価」の「1.4 和文抄録作成対象案件の選定」に記載したとおりである。

この基準に基づき和文抄録作成対象とした全テーマコード（174 テーマ）を別添 1-3 「二次評価（テーマ）分類単位評価結果」に示す。なお、選定された 8 万件の詳細は電子ファイル『04 和文抄録作成対象 8 万件.xlsx』を参照のこと。

### 2.2 和文抄録の作成結果

和文抄録の作成結果を以下に示す。

作成件数		80,000 件
平均文字数（日）		2,203 文字
（請求項）		736 文字
（発明の目的又は効果）		240 文字
（実施例）		1,227 文字
最大文字数（日）	CN103741640B	20,994 文字
最小文字数（日）	CN1325318B	301 文字

### 2.3 特殊案件

和文抄録作成にあたり、中国特許文献に実施例が存在しない等、和文抄録作成にあたり定常の翻訳箇所選定が行えない案件については、その都度対応案を策定のうえ、特許庁担当者の承認を得て所定の対応を取った。

本年度は、別添 2-1 「特殊案件対応一覧」に掲載した 38 文献について、作業員より報告を受け、別添 2-1 に記載の対応を取った。

### 3. 国際調査報告で引用された文献の全文翻訳文の作成

特許庁より貸与された「全文翻訳文作成対象文献リスト」に掲載された中国文献（国際調査報告で引用された中国文献のうち日本語及び英語で公開されたファミリー文献を有しない中国文献）、全 19,541 件について、全文翻訳文を作成した。以下、その概要を記す。

#### 3.1 全文翻訳文の作成結果

全文翻訳文の作成結果を以下に示す。

作成件数		19,541 件
平均文字数（日）		12,973 文字
最大文字数（日）	CN1852974A	183,600 文字
最小文字数（日）	CN1819857A	333 文字

#### 3.2 原文が 4 万文字以上の文献

全文翻訳文の対象となる中国特許文献の文字数が 4 万文字を超えている場合、原文を 4 万文字程度に短縮したうえで全文翻訳文を作成した。これに該当する案件は 129 件存在した。一覧を電子ファイル『[06 4 万文字超過文献一覧.xlsx](#)』に示す。

## 4. 対訳コーパスの作成

本事業にて翻訳対象とした中国特許原文と、作成した翻訳文（和文抄録、全文翻訳文）とを用いて中日対訳コーパスデータを作成した。以下、その詳細を示す。

### 4.1 対訳コーパスの単位

対訳コーパスは、中国特許原文データを以下の基準で分割した文単位にて作成した。

- ・文中の句点「。」で分割
- ・文中の改行で分割

なお、上記基準により、出願人が一文を意図的に改行している場合は、その改行ごとに1レコードとなる。出願人の意図的な改行は、請求項における要素列挙や発明の効果の列挙にて多用されるものであるが、改行の単位で事実上一文として扱うべき場合が多いこと、さらには、意図的改行を連結して一文とすると文長が長大となり、対訳コーパスとしての有用性が低下することなどから、上記基準を採用した。以下、出願人による意図的な改行の一例を示す。

本発明的优点在于：	←改行
1、制备的阳离子淀粉改性AKD为固体，便于储存和运输；	←改行
2、产品熔融后通过高速剪切乳化，……其施胶效果高于传统的AKD乳液。	←句点



連結せず、意図的改行の単位で対訳コーパス化

本発明的优点在于：
1、制备的阳离子淀粉改性AKD为固体，便于储存和运输；
2、产品熔融后通过高速剪切乳化，……其施胶效果高于传统的AKD乳液。

なお、出願人の意図的な改行ではなく、原文データ上の不備等により一文が複数の改行によって分割されているものも存在する（例えば、公報レイアウトの改行位置にて逐一改行マークが挿入されている）。こうした場合は一文に連結した。

步骤 7:支持虚拟 SIM 卡的新移动终端开机后, 通过短消息连接支持虚拟 SIM 卡的呼	←改行
叫服务中心的服务号码, 支持虚拟 SIM 卡的呼叫中心向移动终端发送短消息提示,	←改行
请用户通过新移动终端输入获取原移动终端的账户信息的功能代码。	←句点



連結して本来の文単位で対訳コーパス化

步骤 7:支持虚拟 SIM 卡的新移动终端开机后, 通过短消息连接支持虚拟 SIM 卡的呼叫中心的服务号码, 支持虚拟 SIM 卡的呼叫中心向移动终端发送短消息提示, 请用户通过新移动终端输入获取原移动终端的账户信息的功能代码。

なお、こうした原文の分割・連結は、翻訳文の作成前に行い、全ての翻訳作業はこの単位で実施された。翻訳・校閲にあたっては、専用のインタフェースを用い、原文と翻訳文との対訳関係のズレを物理的に防止した。以下、その一例を示す。

所述方法包括以下步骤：	前記方法は以下のステップを含む：伝送する必要があるデータ量又はサービスタイプに基づいてアップリンク及びダウンリンクの帯域幅を決める。
根据需要传输的数据量或业务类型决定上下行的带宽。	

文区切りのズレによる空欄のアラートや、原文と翻訳文の句点の個数チェック等、原文と翻訳文の正確な文対応を作業インタフェース上で確保。

こうした措置により、翻訳が完了した時点で原文と翻訳文とが一文対一文の単位で完璧に対応づけられた対訳コーパスデータが自動的に完成することとなる。したがって本事業の対訳コーパス作成においては原文と翻訳文との対応付けのためのデータ解析処理（文アライメントツールによる）は一切行っておらず、原則として本事業で翻訳対象とした全ての文を余すことなく対訳コーパス化している。

## 4.2 対訳コーパスの作成対象

対訳コーパスは、原則として、本事業にて作成した全ての翻訳文(和文抄録、全文翻訳文)を対象とした<sup>3</sup>。ただし、以下に該当する文については、中日対訳コーパスという観点からは有用な対訳文とならないため、あらかじめ除外した。

### [(i) 原文がイメージタグ、テーブルタグ、数式タグを含むデータ]

原文中に化学式や表などを示すイメージタグ、テーブルタグ、また、数式を示す数式タグが存在する場合は、翻訳文では、原則として該当箇所を【IMG】に置換している。(テキスト化可能な文字列はテキスト化する場合あり)。

こうした場合、【IMG】に置換するにせよテキスト化するにせよ、原文がタグデータである以上、翻訳文と原文とは正確な対訳データとはなりえない。このようなデータが対訳コーパスに含まれると、コーパスの使用時に悪影響を与えるおそれがある。このため、原文がイメージタグ、テーブルタグ、数式タグを含むレコードは、あらかじめ除外した。

※上付・下付文字タグはあらかじめ翻訳文から除去する仕様につき除外対象とはならない。

### [(ii) 原文、翻訳文が英数字のみで構成されるデータ]

数式や元素記号の列挙など、原文、翻訳文ともに中国語、日本語を含まない文は、中日の対訳コーパスとしては不要であるため、あらかじめデータから除外した。

### [(iii) 原文に不備があり翻訳文と対応していないデータ]

たとえば下例のように、原文の不備により翻訳文と正確に対応しないデータがまれに存在する。このようなデータについて、典型的なパターンによる検出や、翻訳者・校閲者からの指摘に基づく特定を行い、データから除外した。

以示例的方式、 <b>错误！未找到引用源。</b> 给出了本发明的真空紫外光谱测定装置示意图。	実施例により、 <b>図1</b> は本発明の真空紫外スペクトル測定装置の概略図を示す。
---	--

「**エラー！引用元が見つかりません**」という意味のアラート。翻訳文はPDFを参照して「**図1**」と正しく訳されているが、対訳コーパスとしては原文と対応しておらず、除外すべきと判断。

<sup>3</sup> 全文翻訳対象文献に既存の「日本語要約」及び「発明の名称」が存在する場合、当該文献の「要約」及び「発明の名称」についてはコーパス作成の対象外とした。一文対一文での翻訳が必ずしもされている保証がないためコーパスの精度を低下させるリスクがあり、また要約及び発明の名称と同一／類似の文は明細書本文中にほぼ確実に再度出現するため、明細書本文中の記載から得られる対訳コーパスによりカバー可能であると判断した。

### 4.3 対訳コーパスの作成件数

本事業にて作成した対訳コーパスデータの内訳を以下に示す。

和文抄録より	1,556,688 件
全文翻訳文より	2,673,928 件
合 計	4,230,616 件

なお、完成した対訳コーパスデータに対し、以下のダブルチェックを実施し、全てのレコードが一定の文字数比率の範囲内であることを確認した。

#### [対訳文の文字数比較によるチェック]

何らかの理由で文対応不良が発生した場合に備え、原文と翻訳文の文字数比較に基づくチェックを実施した。

具体的には、コーパスの各レコードについて、日本語と中国語の文字数の比率を算出し、両者が大きく乖離する文対や、双方の文字数が非常に多い文対の存在の有無と、存在した場合はその内容を確認した。

以下、文字数比率が最も甚だしかったレコードと、文字数が最も多かったレコードを示す。いずれも、内容的には正しいものとみなされた。

#### ・[文字数比率が最も甚だしかったレコード (比率が小さいもの)]

最も文字数比率が小さいレコードは、日本語文字数と中国語文字数の比率 0.2 (つまり、中国語文字数が日本語文字数の 5 倍) であった。比率の小さいレコードについて抽出した一部の例を以下に示す。同様に比率の低いパターンをサンプル的にチェックしたが、いずれも下例同様、中国語原文と日本語文の対応に問題はなかった。

日/中比率	日本語文	中国語文
0.2	畳	一种榻榻米
0.2	—●詳細はチェックのこと	----- -----●查看详情
0.3	鋳型	一种铸造模型
0.4	<比誘電率>	&lt;相对介电常数&gt;

・[文字数比率が最も甚だしかったレコード (比率が大きいもの)]

最も文字数比率が大きいレコードは、日本語文字数と中国語文字数の比率 19.0 (日本語文字数が中国語文字数の 19 倍) であった。このレコードの原文「则」は、前後に存在するイメージについての接続詞的に用いられているため、翻訳文では意味の通る意識がなされている。同様に、「以下のとおりである」のような文についても比率が大きくなる傾向が見られたが、類似するパターンにおいて特に問題のないことを確認した。その他、比較的比率の大きくなる傾向のあるレコードは、EUC 文字での置き換えが不可能な文字を含んでいる文が該当した。このような文は、翻訳、校閲時の基準によって当該文字を「\* (Unicode U+文字コード)」で置き換えることとしているため比率が大きくなっている。

日/中比率	日本語文	中国語文
19.0	そうすると、すなわち以下になる：	则
13.0	それは以下のとおりである：	有
12.5	a n 按安暗岸俺案鞍* (Unicode U+6C28)* (Unicode U+80FA) 厂广庵・・・(中略)・・・* (Unicode U+9D95)	an 按安暗岸俺案鞍氨胺厂广庵 … (中略) …鵲

・[文字数が最も多かったレコード]

文献番号 CN102383759B の請求項 1 の翻訳文において、日本語文字数が 10,415 文字、中国語原文文字数が 6,657 文字であった。この文献の当該請求項は一文として存在しており、途中に意図的な改行も含まれていないため問題がないと判断した。その他、一定以上の文字数を含むレコードについてサンプル的にチェックを実施し、問題のないことを確認した。

[NM (日本語文数－中国語文数) が 1-1 でないもの]

本事業の仕様に従い、対訳コーパスの各レコードには、日本語文数と中国語文数とを N-M の形で表示している。対訳コーパスは原則として 1 文対 1 文で作成しているため、大半のレコードは 1-1 となるが、そうならないものが若干数存在するため、代表的なパターンを以下に示す。

なお、この文数カウントは、「1+各文の文字列中に存在する句点の数」で機械的に計算した。但し、次の句点はカウントしない。

- ・文末尾の句点
- ・閉じ括弧<sup>4</sup>直前の句点

### ・「n-1」パターン①

原文は 1 文であるが、長大のため日本語の読みやすさ優先や要素列举の明確化のため、日本語が複数文で訳されている。

<p>1. インバータエアコンの室外機であって、それは以下を含む：室外機の底部に設置され且つその上に室外機の各部品を取り付けることができる底板；室外機の前に設置されるフロントパネル；底板上部の一側に設置される圧縮機；フロントパネル下部の一側に設置される制御ボックス；制御ボックス内で、PCBの後に設置される放熱板及び放熱板の後の一側に設置される放熱ファン。///その特徴は以下のとおりである：前記圧縮機の入口管にバイパス管が接続され、バイパス管は放熱板の後に蛇状に巻き付けられる。</p>	<p>1.一种变频式空调器的室外机，包括：设置在室外机的底部且其上能够安装室外机各部件的底盘；设置在室外机前面的前面板；设置在底盘上部一侧的压缩机；设置在前面板下部一侧的控制盒；设置在控制盒内、PCB后面的散热片以及设置在散热片后面一侧的散热风扇，其特征在于：在所述的压缩机的入口管上连接有一旁通管，旁通管在散热片的后面盘成蛇形。</p>
--	---

<p>これに鑑み、本発明は有機発光ダイオードディスプレイを提供し、以下を含む：ガラス基板；及び前記ガラス基板の片側に設けられた回路層、アノード層、有機発光層、カソード層。///更に以下を含む：タッチ誘導層及びタッチ検出ユニット。</p>	<p>有鉴于此，本发明提供一种有机发光二极管显示屏，包括：玻璃基板；以及设于所述玻璃基板一侧的线路层、阳极层、有机发光层、阴极层，还包括：触控感应层以及触控检测单元。</p>
--	---

### ・「n-1」パターン②

原文特有の「；及び」「在于：」「步骤：」「其中：」「如下：」等での意図的な改行により、日本語が複数文となっている。

<p>38. 請求項37に記載のシステムであって、その特徴は、前記BS内の制御モジュールが以下を含むことである：解析部、制御コマンド生成部及び制御部。///ここで、</p>	<p>38.如权利要求37所述的系统，其特征在于，所述BS中的控制模块包括：分析单元、控制命令产生单元以及控制单元，其中，</p>
--	---

<sup>4</sup> 丸括弧、かぎ括弧、角括弧、黒墨括弧、ダブルクォーテーション

・「1-n」パターン①

原文は誤記により複数文となるが、本来は1文のため結合し、日本語は1文としている。

<p>第四級アンモニウム塩の生産プロセスについて多くの研究や開発が行われ、中国特許文献では多数の第四級アンモニウム塩の製造プロセスが開示され、そのうち特許番号が200510061094.4の文献では、炭酸ジエステル（炭酸塩）とアミン（アンモニア）塩を触媒の作用下で第四級アンモニウム塩をワンステップで合成する方法が開示されている。</p>	<p>人们对季铵盐的生产工艺进行了大量的研究与开发，已由中国专利文献公开的季铵盐制备工艺数量众多，其中专利号为200510061094.4的文献公开了从碳酸二酯(脂)与胺(氨)盐在催化剂作用下一步合成季铵盐的方法。</p>
---	---

・「1-n」パターン②

原文は複数文であるが日本語の読みやすさを優先にし、結合し、日本語は1文としている。

<p>Z1及びZ2バイトの使用が標準化されていないため、IOSL装置はこれらのバイトを全0で充填する。</p>	<p>由于 Z1 和 Z2 字节的使用没有标准化。//IOSL 装置将这些字节用全 0 填充。</p>
<p>一部の未重合の小分子又は小分子量の重合分子は、高温下の溶媒抽出方式によって除去する。</p>	<p>部分未聚合的小分子或小分子量的聚合分子。//通过高温下溶剂萃取的方式除去。</p>

・「n-n」パターン

括弧<sup>5</sup>内の句点は文の区切りとはみなすべきでなく、原則、括弧内の句点での分割は行なわれないが、機械的カウントでは文末尾の句点と閉じ括弧直前の句点以外の句点は一律でカウント対象となるため、n が 2 以上とされるケースが発生する。

<p>適切な閾値 (閾値を経験的方法で確定する。///本実施例において円周サンプリングを480にすると、隣接する3点の夾角は <math>\pi - 2\pi/480</math> であり、測定誤差の影響を考慮して、閾値を <math>\pi - 5 * 2\pi/480 \sim \pi - 10 * 2\pi/480</math> に設定することができる。///本実施例における閾値を3.05ラジアンに設定する) を設定することで、どの点が切り欠き上の点であるかを判断できる。</p>	<p>设定合适 的閾値 (閾値采用经验的方法确定。///本实施例中圆周采样为480, 则相邻3点的夹角为 <math>\pi - 2\pi/480</math>, 考虑到测量误差的影响, 閾値可以设定为 <math>\pi - 5 * 2\pi/480 \sim \pi - 10 * 2\pi/480</math>。///本实施例中的閾値设定为3.05弧度), 可以判断哪些点是缺口上的点。</p>
---	--

<sup>5</sup> 丸括弧、かぎ括弧、角括弧、黒墨括弧、ダブルクォーテーション

## 5. 辞書データの作成

中国特許公報の和文抄録の作成及び全文翻訳文の作成で作成した翻訳文（和文抄録、全文翻訳文）と原文からなる対訳コーパスから対訳辞書データを作成した。以下、その概要を記す。

### 5.1 辞書データの作成手順

対訳辞書データの作成手順として、まず対訳コーパスを解析し、対訳辞書候補データを作成した。その後、特許庁より貸与された約 230 万件語の対訳辞書データとの突き合わせを行い、重複しないデータのみを人手確認用辞書候補データとして作成した（146,970 件）。この人手確認用辞書候補データを頻度順に、採用語数（下表「対訳辞書として適切」に該当）が 10 万語に達するまで人手確認を行い、対訳辞書データを作成した。人手確認結果を以下に示す。

表 5-1：対訳辞書データ候補 人手確認結果

対訳辞書として適切	100,000 件
見出し語自体が不適切	13,339 件
訳語自体が不適切	4,220 件
見出し語と訳語とが対応していない	2,435 件
辞書に登録する用語として不適切	2,163 件

また、上記 10 万語に加え、本事業にて実施した翻訳品質評価の対象とした重要技術用語 8,000 語のうち、機械翻訳結果が基準翻訳文の訳語と異なるものについても、同様に人手確認・重複削除を行い、辞書データとして適切と判断された 1,680 件を対訳辞書データに追加した。なお、これらの語は対訳コーパスからの採用ではないため、出現頻度が 0 の場合がある。

## 5.2 辞書データの作成件数

### 5.2.1 辞書データ作成件数

上記の対訳辞書として適切と判断されたデータ（101,680 件）と特許庁より貸与された対訳辞書データを統合し、「見出し語（中国語）、訳語（日本語）、品詞（見出し語（中国語））、品詞（見出し語（日本語））」が付された形式とした上で、UTX 形式に変換し、対訳辞書納入データ（納入物④）を作成した(2,369,120 件)。また、人手確認で不採用となった候補データ全件（22,157 件）についても、別途「人手確認により対訳辞書から除外したデータ」として納入した。

①対訳辞書データ（新規作成）	101,680 件
②特許庁辞書データ（既存）	2,267,440 件
対訳辞書納入データ（①+②）	2,369,120 件

### 5.1.2 データフォーマット

以下、「対訳辞書納入データ」及び「人手確認により対訳辞書から除外したデータ」のフォーマットを示す。

#### 5.1.2.1 「対訳辞書納入データ」フォーマット

ファイル名：JPO-CJ-DICT-28fy.utx

フォーマット：冒頭の#で始まる行はヘッダ部（UTX 1.11 に準拠）であり、その後のデータ部の各行は一つの見出し語；訳語のセットとなっている。

- ・ヘッダ部1

```
#UTX 1.11; zh-CN/ja-JP; yyyy-mm-dd; JPO
```

上記 yyyy-mm-ddには作成日が設定される。

- ・ヘッダ部2

```
src <tab> tgt <tab> …
```

ヘッダ部2には下記データ部に示す記号項目名がタブ区切りで設定される。

・データ部

データ部には下記に示す項目がタブ区切りで設定される。

項番	記号項目名	項目
1	src	中国語・見出し語
2	tgt	日本語・訳語
3	src:pos	中国語・品詞
4	tgt:pos	日本語・訳語
5	comp-num	複合語の語数
6	freq1-all	頻度情報・全分野
7	freq1-chem	頻度情報・化学分野
8	freq1-elec	頻度情報・電気分野
9	freq1-mach	頻度情報・機械分野
10	freq1-phys	頻度情報・物理分野
11	freq2-all	補正された頻度情報・全分野
12	freq2-chem	補正された頻度情報・化学分野
13	freq2-elec	補正された頻度情報・電気分野
14	freq2-mach	補正された頻度情報・機械分野
15	freq2-phys	補正された頻度情報・物理分野

なお、上記項目の複合語の語数、頻度情報、補正された頻度情報は「平成27年度 中国特許文献の機械翻訳の品質評価及び辞書整備に関する調査<sup>6</sup>」の項番4.3.3 「統合された対訳辞書作成」に記載の作成方法に準拠し作成した。

<sup>6</sup> [https://www.jpo.go.jp/shiryou/toushin/chousa/pdf/kikai\\_honyaku/h27\\_01.pdf](https://www.jpo.go.jp/shiryou/toushin/chousa/pdf/kikai_honyaku/h27_01.pdf)

### 5.1.2.2 「人手確認により対訳辞書から除外したデータ」フォーマット

ファイル名 : cj\_dict\_del\_1\_4.txt

データ項目

下記に示す項目がタブ区切りで設定される。

項番	項目	備考
1	人手確認結果	記号D1, D2, D3, or D4が設定される。 D1(見出し語自体が不適切) D2(訳語自体が不適切) D3(見出し語と訳語とが対応していない) D4(辞書に登録する用語として不適切)
2	中国語	
3	日本語	

## 6. テーマコード情報の作成

本事業にて全文翻訳文を作成した中国特許文献 19,541 件の全件について、特許庁により機械的に付与されたテーマコードの妥当性チェックを実施した。チェックの結果に基づき、特許庁より貸与された「テーマコード情報リスト」の内容を更新した。以下、作業結果について報告する。

### 6.1 「テーマコード情報リスト」更新内容

テーマコードの妥当性チェックの対象となった中国文献 19,541 件における、機械付与テーマコードの総数と、このうち妥当（正解）であったテーマコード件数、誤りであり削除したテーマコード件数、不足であり追加したテーマコード件数について以下に示す。

	機械付与 テーマ数	正解 テーマ数	削除 テーマ数	追加 テーマ数	修正後 テーマ数
合計	46,476	32,276	14,200	8,480	40,756
平均	2. <sup>378</sup>	1. <sup>652</sup>	0. <sup>727</sup>	0. <sup>434</sup>	2. <sup>086</sup>
		69%	31%	18%	

上表に示したとおり、対象案件 19,541 件における機械付与テーマコード数は平均 2.378 であったのに対し、人手確認後の付与テーマ数は平均 2.086 と、若干（平均 0.292）の減少が見られた。

機械付与テーマコードの正解率は約 69%であった。なお、一案件に対して機械付与されたテーマが全ての正解テーマ（＝確認後付与テーマ）を含んでいた案件は 19,541 件中 9,860 件（約 50%）であり、このうち 7,515 件（約 38%）は、確認後付与テーマと完全に一致した（つまり削除も追加も不要であった）。