

技術文書に対する事象ラベル生成と付与の検討

深草理貴¹ 大八木悠聖¹ 喜多俊介¹ 森辰則¹

小野寺理恵² 伊藤拓海²

¹横浜国立大学大学院 ²株式会社 IHI

fukagusa-riki-wf@ynu.jp ooyagi-yuusei-kn@ynu.jp kita-shunsuke-yz@ynu.jp tmori@ynu.ac.jp

onodera3892@ihi-g.com ito9762@ihi-g.com

概要

品質規程の整備では、規程文書をプロセス単位で参照・監査できる形に再編する必要がある一方、業務文書は簡潔、規程文書は条件や例外を含む詳細記述になりやすく、粒度差のため単純な文字列突合が難しい。本研究は、文が表す状況を述語項構造に基づく JSON 形式の「構造ラベル」として表現し、述語概念と格付き項概念を付与して照合する枠組みを示す。複合名詞は内包事象を入れ子で保持し、焦点で解釈のぶれを抑える。特許請求項と明細を模擬データとして、一般化の有無による照合の挙動を評価した。

1 はじめに

品質規程の整備活動では、規程文書を理解・管理しやすい形に再編するため、記述をプロセス単位で束ね、プロセスごとに参照・監査・改定できる状態を目指すことが多い。ここでいうプロセスは業務の一部分を成す作業の系列として業務文書に記述され、実務上は短い一文、あるいは「対象 α の処理 β 」のような名詞句で表される。一方、規程文書は各プロセスに対して満たすべき規程を、付随作業や確認事項を含めて文章で詳細に記述するが、規則集合として厳密に整形されているとは限らず、「(作業 A) をする場合には、(付随作業 B) をし、(確認事項 C) を確認する」といった不定形の文を含むことが多い。このため適用条件の内部に埋め込まれたプロセス記述を、業務文書に現れるプロセス記述と対応付ける必要が生じる。

品質管理の観点では、業務文書における各プロセスが満たすべき規程を網羅的に把握できることが重要である。しかし、業務文書は簡潔であるのに対し、規程文書は条件・例外・確認手順などを含む詳細記述になりやすく、粒度差を無視して文字列類似度で

突合することは難しい。規程の適用条件に含まれるプロセス記述は、業務側と同一、または業務側を一般化した表現であることが期待されるため、粒度差を考慮しつつ一般化を介して対応関係を推定する仕組みが必要となる。

方法論として文間類似度、文間含意関係、文書分類などが考えられるが、粒度差が大きい場合には類似度や含意判定の頑健性が課題になりうる。また分類に基づく方法は固定カテゴリと学習事例を要し、業務文書の増加や規程改定を前提とすると、カテゴリの事前設計と維持が現実的でない場合がある。そこで本研究では、プロセス記述を述語—項構造に基づく構造化情報として表現し、さらにシソーラスに基づく語の分類情報を付与して一般化する戦略をとる。業務文書の短いプロセス記述から得られる小規模な構造化情報を出発点に、抽象化・一般化を行って規程文書側のプロセス記述と整合する表現に近づけることを狙う。加えて本研究で扱う「ラベル」は文字列のみではなく、「人間向け説明用の文字列」と「対応する構造化情報」を組として管理することを前提にする。構造化情報だけでは人間が理解・改定・再利用しづらく、文字列だけでは意味的な突合に耐えないためである。本稿ではまず、この組のうち構造化情報側を中心に、生成・付与の手順と評価を報告する。

2 関連研究と本校の貢献

述語項構造や格フレームに基づく意味解析は日本語処理の基盤技術であり、ゼロ代名詞や格の推定など、文内部の役割関係を用いた解析が長く研究されてきた [1]。また、述語や用言に対する意味属性付与は、語彙体系に基づく一般化の一形態として位置づけられ、運用で扱える上位カテゴリへ写像する設計が議論されている [2]。語の抽象化・上位概念化に関しては、日本語語彙大系のようなシソーラ

ス資源が整理されており、表層差を吸収して同一状況として扱う基盤を提供する [3] [4]。さらに、意味役割付与やその汎化は、述語項構造を用いた一般化の実装側の論点に関係し、意味役割の汎化に関する議論や深層学習モデルによる付与も報告されている [6] [7]。特許文書は定型表現と独特の言い回しを含むドメイン文書であり、ドメイン適応や構文解析の観点からも扱われてきた [5]。

以上を踏まえ、本稿では、固定カテゴリに一度で分類するのではなく、文が表す状況を述語—項構造として抽出し、構造を一般化して粒度差のある別文書表現と照合できるかを検証する。ラベルは「人間向け文字列」と「述語—項構造に基づく JSON」の組として管理するが、本稿では主に後者を扱い、複合名詞内の事象も入れ子で保持して、名詞句に内包された事象と文の主要事象の混同を抑える。さらに、一般化したラベルと、共通部分抽出で得た構造を比較し、どの程度「同じ状況」とみなせるかを評価に組み込む点を特徴とする。

3 課題設定

本研究の最終目標は、規程文書の文章をプロセス単位で束ね、プロセス体系の細粒度化とラベル付与を一体として支援することである。業務文書側にはプロセス系列が存在し、規程文書側には各プロセスに適用される規則・確認事項が文章として存在する。対応づけとは、業務文書のプロセス記述が、規程文書の適用条件に含まれるプロセス記述と同一または一般化関係にあることを判定し、規程を該当プロセスへ紐づける操作とみなせる。

実データの代替として、本稿では特許文書の請求項を「簡潔なプロセス記述」、明細を「条件・根拠・詳細の説明」と見立てる。請求項側のプロセス記述から得られる表現を、明細側のより説明的な表現と対応づけることを、粒度差を伴う対応づけの模擬課題として扱う。この設定では、請求項側は短い構造になりやすく、明細側は同一状況をより多くの修飾・条件・付随行為とともに表すため、請求項構造が明細構造の部分構造として現れるか、または一般化により整合し得るかが中心課題となる。

4 提案手法

4.1 構造ラベルの基本単位

本研究では、ラベルを人間向けの文字列ではなく、

述語項構造に基づく構造化意味表現 (JSON) として定義する。以降、この JSON 表現を構造ラベルと呼ぶ。構造ラベルは、文中の事象を「述語」と「項」の集合として表す。述語には表層 (原文の動詞・用言) と、運用で扱うための述語概念を付与する。項には表層 (名詞句)、表層格 (が・を・に等)、深層格 (主体・対象・場所等)、および項概念を付与する。本文では各フィールドの意味と、照合で利用する要素を説明する。構造ラベルの具体例となる JSON 本体は、付録の図 A1 としてまとめて掲載する。大規模言語モデル (LLM) にいくつかの事例 (表層の文字列と JSON による構造化意味表現の組) を与えると、このような構造を生成できる。

4.2 自動化の位置づけ

本研究の枠組みは、規程整備における対応付け作業を、機械可読な表現にもとづいて反復可能にすることを意図している。そのため、構造ラベルの作成は、人手での厳密なアノテーションを前提とせず、大規模言語モデル (LLM) を用いて半自動的に生成する設計を採る。ここでいう半自動とは、述語・項・格・概念付与までを LLM が一括で候補生成し、その出力を人手で確認して必要箇所のみ修正・再生成しながら確定する分担を指す。そのため、人手でゼロから構造を付与するのではなく、深層格の取り違えや複合名詞の入れ子範囲、焦点、概念付与の過不足といった解釈のぶれやすい点を中心に点検して確定する運用を想定する。具体的には、表層の文字列と JSON による構造化意味表現の組を少数例として与え、述語の抽出、項の抽出、表層格・深層格の付与、ならびに述語概念・項概念の付与を一貫して出力させる。さらに、複合名詞に対しては、名詞句内部に内包される事象を仮定した入れ子表現を生成させ、後段の照合で利用可能な構造として保持する。本稿では、このように生成された構造ラベルを入力として照合を行い、一般化の導入により粒度差のある記述間で対応付けが可能になるかを評価する。

4.3 述語項構造の構築と再帰的名詞句解析

入力が文である場合は、文の主要な述語を抽出し、述語が取る項 (対象や場所など) を抽出して構造ラベルを生成する。図 A1 では、文全体の事象として「敷き詰める」を述語に持つ事象 1 を立て、その述語概念を Connective Action として付与している。事象 1 の項には、「断熱性被覆フロート」を表層格「を」、

深層格「対象」として与えた項1と、「防液堤の底部」を表層格「に」、深層格「場所」として与えた項4を含める。

また、本研究では複合名詞を単なる名詞として扱わず、名詞句内部に内包される事象を仮定して分解し、入れ子の事象として保持する。図A1では、「断熱性被覆フロート」を、事象2(述語「被覆」)が取る項の組として表現し、その内部で「フロート」を概念「物品」として位置づけ、さらに「断熱性」を属性として付与している。同様に、「防液堤の底部」も、事象3(述語「底部」)として分解し、「防液堤」を全体として持つことで、場所表現がどの対象の部分であるかを明示する。このように、名詞句に含まれる事象(被覆、底部)と、文の述部として現れる事象(敷き詰める)を混同せず、入れ子構造として区別したまま表現できる。

提示するJSON形式では、各事象・項にIDを付与し、さらに名詞句として何を中心に解釈するかを「焦点」フィールドで明示する。図A1では、「断熱性被覆フロート」に対しては焦点を項2(フロート)に置き、「防液堤の底部」に対しては焦点を項4(底部)に置いている。焦点は照合時に「その名詞句が何を指しているか」とみなすかを固定する役割を持ち、複合名詞や部分表現における解釈のぶれを抑えるために導入する。さらに、全体としての焦点を事象1に置くことで、本構造ラベルが文全体の事象(敷き詰める)を代表することを明確化している。

4.4 述語のプロセス概念化

述語は、企業内運用で扱いやすい少数のプロセス概念へ写像する。本稿では「蓄積・保持」「移送・移動」「変換・変調」などのカテゴリを用い、述語表層をこれらの概念に対応づける設計を採る[2]。ここで重要なのは、文の述部としてのプロセスと、複合名詞内部に内包されるプロセスを区別したまま、同一の概念体系へ写像できる点である。

4.5 項のシソーラス一般化

項(名詞句)は表層表現に依存しすぎない一般化を行うため、シソーラス資源を参照して上位概念を付与する[3][4]。これにより、表層が異なる名詞句同士でも、より上位の概念階層で「同一状況の変種」として扱える可能性が高まる。

4.6 共通部分抽出とラベルの一般化

本研究では、二つの方向を区別して扱う。一つは、請求項ならびに明細からそれぞれから得られた生の二文(あるいは二つの構造)を照合し、共通部分だけを残した構造を抽出する方向である。これは「似た状況をまとめてラベル候補を得る」側の操作に近い。もう一つは、あらかじめ定めた規則に基づいて請求項から得られたラベルを一般化し、明細から得られた文の構造に当てはめて照合できるかを検証する方向である。本稿の評価では、この二つが同じ挙動を示すか、または違いがあるかを検討対象とし、ラベル側の一般化レベルと共通部分抽出結果のずれが、ラベルとしての妥当性判断にどう影響するかを議論する。

5 評価実験

本評価実験の目的は、請求項と明細の間に存在する記述粒度の差を前提としたとき、構造ラベルの一般化が照合の成立にどの程度寄与するかを確認することである。具体的には、人手で同一状況を表していると判断した対応関係を基準とし、構造ラベル照合がそれらをどの程度再現できるかを、再現率によって評価する。このとき再現率は、人手で正しいとされた対応のうち、照合により一致と判定できた割合として定義し、集計はマイクロ平均で行う。また、本稿では一般化の適用範囲の違いが照合結果に与える影響を調べるため、請求項側と明細側の双方を元構造のまま照合する手法、双方を一般化した構造同士で照合する手法、請求項側のみ一般化して明細側は元構造のまま照合する手法の三手法を同一の人手正解集合に対して適用する。

5.1 人手正解と一致判定単位

人手正解は、請求項側の構造と明細側の構造が同じ状況を表していると判断された対応に基づいて与える。判定単位は一律に固定しないが、述語概念と、項概念、および表層格・深層格の整合が中心となる。付録図A2のように、述語概念がConnective Actionであり、項概念が「物品」「建造物」で、それぞれが格情報と整合している場合には、この三点が正解判定に必要な情報となる。このように、一致判定に含める要素の選択が結果に影響するため、評価では照合条件ごとに同一の判定規則を適用し、マイクロ平均で集計する。

5.2 三手法とその再現率評価

再現率の評価では、人間が正しいと思っている対応の数を分母とし、システムが一致と判定できた数を分子として計算する。本稿では次の三条件を設定する。第一に、両側とも元の構造のまま照合する。第二に、両側とも一般化した構造同士を照合する。第三に、請求項側のみ一般化し、明細側は元構造のまま照合する。いずれもマイクロ平均で算出する。

6 結果

人手正解の総数を 31 とし、5.2 の三つの手法に対して再現率を算出した。表 1 に一致数と再現率を示す。手法①（元の構造同士）では 26 件が一致し、再現率は 0.84 であった。手法②（一般化した構造同士）では 28 件が一致し、再現率は 0.90 であった。手法③（請求項のみ一般化、明細は元構造）でも 29 件が一致し、再現率は 0.94 であった。手法①について一致した構造と不一致した構造をそれぞれ図 B1 と図 B2 に示す。

表 1 再現率の集計（マイクロ平均）

手法	一致数 (TP)	全体 (TP+FN)	再現率 (Recall)
①	26	31	0.84
②	28	31	0.90
③	29	31	0.94

7 議論

表 1 より、元構造同士の照合（手法①）は、一般化を導入した照合（手法②・③）に比べて再現率が低かった。本稿の一致判定は、文全体に対応する主要事象について、述語と項の整合が取れていることを重視する。そのため、一致したケースでは、共通構造が文全体の主要な事象の述語と項がすべて保持されている一方、不一致のケースでは、一方の構造に関して必須の述語と項が欠けていると言える場合と、共通構造が名詞句内部の入れ子事象に縮退した場合があった。

手法①では、両文の構造が詳細な入れ子まで含んだまま比較されるため、片方の文では名詞句内部の事象が強く現れ、もう片方では主要事象側に情報が寄るといった構造の偏りがあると、共通部分抽出が主要事象まで到達できず、入れ子側の部分構造だけが残る（＝縮退する）ことが起こり得る。これに

対して構造の一般化では、照合に寄与しにくい名詞句内部に対応する事象を除外し、文全体の主要な事象の述語と項の組を中心とした構造へ正規化する。その結果、入れ子の有無や内部の分解の深さに起因する差分が一致条件に入りにくくなり、粒度差をまたいでも事象レベルでの整合が成立しやすくなるため、手法②・③で再現率が改善したと考えられる。また、手法②と手法③が同程度であった点は重要である。請求項側をラベルとして一般化しておけば、明細側を元の構造のまま保持しても照合が成立し得ることを示しており、運用上、説明的な側まで一律に一般化しない設計でも網羅性を確保できる可能性がある。ただし、再現率の上昇がそのまま誤一致の増加を意味しないかどうかは、適合率側の検証が必要である。一般化は一致条件を緩めるため、過度に行うと誤一致が増える恐れがある。そこで、ラベル側の一般化レベルと共通部分抽出で得られた構造の差を比較し、ラベルが状況を十分表せているか／逆に二文に特化しすぎていないかを見極める必要がある。また一致判定でどの項まで必須とするかは運用目的（網羅性重視か精度重視か）によって変わる。今回の結果は、採用した判定規則の範囲では一般化が再現率を改善することを示しており、今後は適合率も含めて人手判断に近い一般化水準を検討する。

8 おわりに

本稿では、企業の品質規程整備における「規程文書と業務文書の粒度差を伴うプロセス対応づけ」という課題を、特許文書の請求項と明細の関係に擬して定式化し、述語項構造化とシソーラス一般化に基づく構造ラベル生成・照合の枠組みを示した。構造ラベルは、述語概念と格付きの項概念のセットとして表現し、複合名詞内部の事象を入れ子として保持することで、名詞句に内包された事象と文の述部としての事象の混同を抑える設計を採った。評価では、再現率を三条件で測定し、元構造同士の照合に比べて一般化を導入した照合で再現率が向上すること、また請求項側のみ一般化しても両側一般化と同程度の再現率が得られることを確認した。今後は、LLM による共通部分抽出を含む適合率評価を追加し、ラベル一般化と共通構造抽出のずれが、ラベルとしての妥当性判断に与える影響を定量・定性の両面から検証する。

参考文献

- [1] 中岩浩巳. 語用論的・意味論的制約を用いた日本語ゼロ代名詞の文内照応解析. 自然言語処理, 3(4), pp.49–66, 1996
- [2] 中岩浩巳, 白井諭, 池原悟. 用言意味属性を用いた日本語ゼロ代名詞の文内照応解析. 情報処理学会 第 49 回全国大会講演論文集, pp.3-173, 1994.
- [3] NTT コミュニケーション科学基礎研究所 (言語情報研究グループ). 日本語語彙大系 CD-ROM 版 (資源概要). 1999.
- [4] 池原悟 (編). 日本語語彙大系. 岩波書店, 1997.
- [5] 橋本力, 河原大輔, 吉田節行, 後藤広樹, 横山晶一. 特許文書の構文解析. 第 51 回自動制御連合講演会, 2008.
- [6] 松林優一郎. 自動意味役割付与における意味役割の汎化. 自然言語処理, 17(4), 2010.
- [7] タロク・カラム, 荒牧英治, 高村大也. 深層学習を利用した PropBank 形式の日本語意味役割付与モデル. 言語処理学会第 29 回年次大会 (NLP2023) 発表論文集, pp.2419–2421, 2023.
- [8] Kargupta, P., Zhang, N., Zhang, Y., Zhang, R., Mitra, P., Han, J. TaxoAdapt: Aligning LLM-Based Multidimensional Taxonomy Construction to Evolving Research Corpora. Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (ACL 2025), 2025.

付録

図 A1

```
{
  "表層": "防液堤の底部に断熱性被覆フロートを敷き詰める。",
  "意味": {
    "事象": {
      "ID": "事象1",
      "述語": { "表層": "敷き詰める", "意味": { "概念": "Connective Action" } },
      "項": [
        {
          "ID": "項1",
          "表層": "断熱性被覆フロート",
          "表層格": "を",
          "深層格": "対象",
          "意味": {
            "事象": {
              "ID": "事象2",
              "述語": { "表層": "被覆", "意味": { "概念": "Connective Action" } },
              "項": [
                { "ID": "項2", "表層": "フロート", "表層格": "を", "深層格": "対象", "意味": { "概念": "物品" } },
                { "ID": "項3", "表層": "断熱性", "表層格": "の", "深層格": "属性", "意味": { "概念": "性質" } }
              ]
            }
          },
          "焦点": "項2"
        }
      ]
    }
  },
  {
    "ID": "項4",
    "表層": "防液堤の底部",
    "表層格": "に",
    "深層格": "場所",
    "意味": {
      "事象": {
        "ID": "事象3",
        "述語": { "表層": "底部", "意味": { "概念": "場" } },
        "項": [
          { "ID": "項5", "表層": "防液堤", "表層格": "の", "深層格": "全体", "意味": { "概念": "建造物" } }
        ]
      }
    },
    "焦点": "項4"
  }
]
},
"焦点": "事象1"
}
```

図 A2

```
{
  "意味": {
    "事象": {
      "ID": "事象1",
      "述語": { "意味": { "概念": "Connective Action" } },
      "項": [
        { "ID": "項2", "表層格": "を", "深層格": "対象", "意味": { "概念": "物品" } },
        { "ID": "項5", "表層格": "に", "深層格": "場所", "意味": { "概念": "建造物" } }
      ]
    }
  },
  "焦点": "項2"
}
```

図 B1

```
◎一致
表層：漏出した液化ガス
{
  "共通構造": {
    "意味": {
      "事象": {
        "述語": { "意味": { "概念": "Physical Transfer" } },
        "項": [
          {
            "表層格": "が",
            "深層格": "対象",
            "意味": {
              "事象": {
                "述語": { "表層": "液化", "意味": { "概念": "変換・変調" } },
                "項": [
                  {
                    "表層格": "を",
                    "深層格": "対象",
                    "意味": { "概念": "気体" }
                  }
                ]
              }
            }
          }
        ]
      }
    },
    "焦点": "項2"
  }
}
```

図 B2

```
×不一致
表層：表層：漏出した液化ガス
{
  "共通構造": {
    "意味": {
      "事象": {
        "述語": { "表層": "液化", "意味": { "概念": "変換・変調" } },
        "項": [
          {
            "表層格": "を",
            "深層格": "対象",
            "意味": { "概念": "気体" }
          }
        ]
      }
    }
  }
}
```