

研究報告書

申請者：宮下博幸（関西学院大学）

助成対象者：Julian Michael Stawecki, PhD（デュッセルドルフ大学）

研究課題「ドイツ語の項構造同定を目的とする解析システム（パーサ）の開発」

現在、ドイツ語に関しては大規模な電子データが存在している。しかしそのデータを言語学的に重要な句の単位に分け、またその句の単位から構成される文の項構造を自動的に同定するようなツールがないのが状況である。生成AIは有用なツールとなりうるが、この作業に特化したものでないため現段階ではなお実用的ではない。このような背景から、本共同研究はドイツ語の大規模コーパスにおいて名詞句、前置詞句などの句を認定し、さらにそれらと動詞との組み合わせを自動的に表示することを可能とする解析システム（パーサ）の開発することを目的とした。パーサの開発によりこのような解析の自動化が可能になれば、ある動詞がどのような句要素とともに出現するのかを大規模コーパスのデータをもとに実証的に示すことが可能となり、さらには動詞と句要素との共起パターンを頻度として把握できることで、今後の辞書記述にも多大な貢献となることが期待される。

以上の目的を達成すべく、助成対象者のStawecki氏とともに研究を進めた。その際、必要な解析内容の検討を含む解析システムの立案は主に申請者が行い、Stawecki氏はその立案に基づいてプログラミングを行う形で進めていった。解析システムは実際に解析データ上で走らせながらさまざまな改善を積み重ねていく必要があったが、日独間のメールならびにオンライン会議による検討を重ね、昨年末までに試行版のパーサが完成了。続けて試行版パーサの精度の検証を行った。検証を行うにあたっては新聞、フィクション、ノンフィクション、学術のそれぞれの比較的短いテキストを用意し、それらのテキストをパーサに分析させ、手作業でその結果を分析した。その結果、新聞では0.91、ノンフィクションでは0.92の精度を得ることができた。一方、フィクション、学術テキストに関してはこれより精度が低く、それぞれ0.81、0.84であった。これは新聞、ノンフィクションのテキストに比べてフィクションや学術テキストの複雑性が高いことに起因すると考えられる。4つのテキストの平均の精度は0.87であった。

以上の結果は2025年3月に長崎の出島メッセで開催された第31回言語処理学会（NLP2025, 2025年3月10日～14日）で報告した。Stawecki氏はこの共同の報告のために来日した。滞在中は学会報告で得られたフィードバックについて氏と対面で話し合うことができ、大変有意義な機会となった。氏の帰国後も、引き続きパーサのさらなる改良のための議論を継続している。

今後の課題としては、開発したパーサにとって苦手なテキスト種であることが判明したフィクションや学術テキストに関し、さらに精度を上げるための工夫を行っていくこと、またこれまで精度検証のために非常に小規模なサンプルテキストを用いていたが、今後例えれば平均で0.95程度の精度に達した時点で、本来の目的である大規模のデータの処理へと移行し、またそこから冒頭に目的として挙げたドイツ語の項構造の量的分析につなげていくことが求められる。Stawecki氏とともに、これらの課題を引き続き進展させてきたい。

なお本奨学金の費用明細は以下の通りとなっている。

往復渡航費（フランクフルト～日本）	20万円
日本滞在費（宿泊費、交通費等）	26万円
言語処理学会会費・参加費	2万円
研究関連の書籍購入	2万円

本奨学金を与えていただいたことに改めて感謝申し上げたい。

2025年11月4日