

新素材研究開発分野における文献からの研究支援情報抽出技術の提案 Technology to Extract Research Clue from Scientific Papers for New Materials Developments

照屋絵理[†] 小野寛太[‡] 竹内理[†] 森田秀和[†] 林貴之[†]
Eri Teruya Kanta Ono Tadashi Takeuchi Hidekazu Morita Takayuki Hayashi

1. はじめに

材料開発分野では国際的な開発競争の激化に伴い、短時間・低コストでの材料開発が必須となっていることを背景に、マテリアルズインフォマティクス (MI) が活発化してきた[1]。MI は物質・材料の物理的・化学的性質に関する膨大なデータと AI などの IT 技術を駆使して、材料開発を効率化する取り組みである。近年では、過去の実験データ等の数値データの解析により求める特性を持つ新規材料の原料候補を絞り込むことで、新規材料開発にかかる実験回数を削減する試みが盛んに行われており、実験回数を最大 10 万分の 1 に削減したという事例も発表されている[2]。

上記のように、従来は実験データ等の数値データに対する MI が主流であった。しかし、MI の有用性を示す事例が増えるにつれ、更なる材料開発の効率化に向け、数値データのみならずテキストデータを活用した MI への期待が高まっている。特に、テキストマイニング技術により、科学文献から材料開発を支援する情報 (研究支援情報) を抽出したいというニーズが増大している。例えば、文献から物質の未知の性質やある物質の類似物質を抽出する、過去に特定の物質に対してよく適用された技術 (実験や解析手法など) を抽出する、過去に特定の物質の研究を実施した研究者や所属の情報を抽出するなどである。これらのニーズに対し、Tshitoyan ら [3] は、科学論文に対して word embedding 技術 [4] を用いて単語をベクトル化し、単語間の距離を見ることで、材料研究の支援となる物質とその物質の性質との対応関係が抽出可能であることを示している。

彼らの試みにより、材料分野における word embedding 技術への期待が高まっている。しかし、材料研究実業務での本技術の活用に向けては、彼らの研究内容に加え、材料研究者が着目する観点に関する情報の検索や可視化の仕組みが求められる。加えて、材料研究時に研究者がよく参考とする、過去に特定物質を研究した研究者やその所属等に関する関係性の抽出も求められている。

そこで、本研究では、材料研究では物質名等の特定のキーワードや、著者や所属等の科学文献に関する特定の情報およびそれらの関係性が研究を支援する上で特に重要であることに着目し、これらの情報および関係性を研究支援情報として文献から自動抽出し、さらに関係性の視認性が良いグラフ構造で可視化、検索を可能とする技術を提案する。以降では、2章で提案技術の概要を説明する。3章で提案技術の初期評価および考察を実施し、4章でまとめを述べる。

2. 研究支援情報抽出・可視化技術の提案

本研究では、研究支援情報を科学文献から半自動で抽出、さらに情報の絞込み、検索、可視化を可能とする研究支援情報抽出・可視化技術を提案する。本技術の開発に先立ち、材料研究者と議論する中で、研究支援情報として、物質名とその物性値の詳細情報、技術名、物性値名に関するキーワード、および論文情報 (journal や出版年などを含む)、著者名、所属名とそれらの関係性を求めていることが分かった。そこで、これらの情報を科学文献から固有表現抽出技術 [5]、word embedding 技術、科学文献出版社が提供する Web からの論文情報取得の仕組み (文献取得 API) を用いて半自動で抽出する。さらに、材料研究者は関係する複数の情報から物質の性質等を連想することがあり、関係性を一目で把握できることを求めていることに着目し、視認性のグラフ形式で情報を可視化する。

図 1 に提案技術の概要を示す。提案技術では情報抽出部にて研究支援情報を抽出し、可視化部にて抽出した情報を可視化する。下記にてそれぞれについて説明する。

2.1 情報抽出部

情報抽出部では、文献からキーワードとして物質名、技術名、物性値名を固有表現抽出技術を用いてキーワード種別毎に抽出する。さらに word embedding 技術を用いてキーワード間の関係性を抽出する。また、論文情報、著者名、所属名の文献のメタデータとそれらの関係性を抽出する。

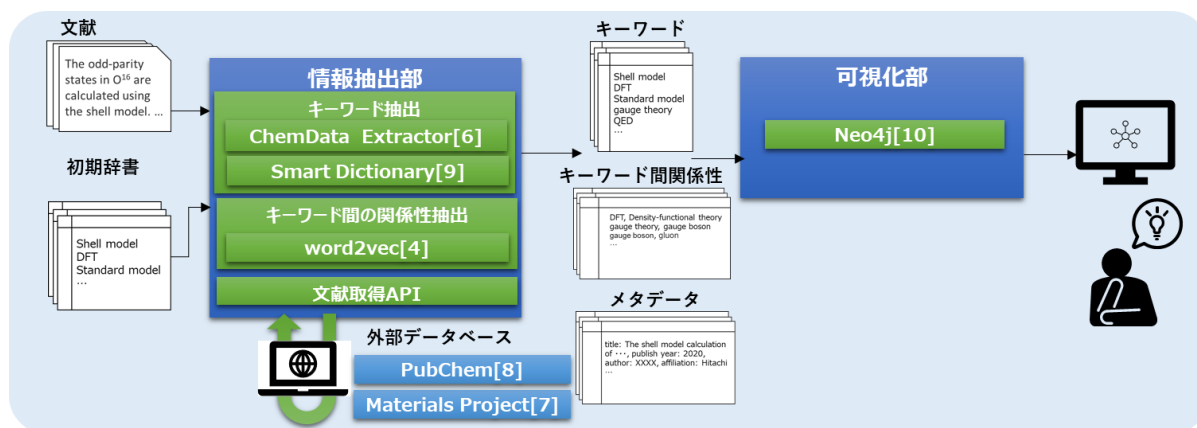


図 1 提案技術概要

[†] (株) 日立製作所 Hitachi, Ltd.

[‡] 高エネルギー加速器研究機構 KEK

