

「ドメイン依存の固有表現抽出技術の現状」付録：ドメイン依存の固有表現抽出に関連する論文リスト

- ・本リストは、ドメイン依存の固有表現抽出技術の現状を調査するにあたり、以下3つの学会の論文を参照し、その中からドメイン依存の固有表現抽出をテーマとしたものを抽出したものです。
 - ・言語処理学会年次大会（2019年3月～2022年3月）
 - ・言語処理学会論文誌（2018年1月～2021年12月）
 - ・情報処理学会NL研究会（2018年5月～2021年9月）
 - ・電子情報通信学会テキストアナリティクス・シンポジウム（2011年7月～2021年11月）
- ・論文の本文内で紹介していない研究も含みます。
- ・論文内で紹介した研究はタイトル部分をハイライトしています。
- ・順序は、論文内でのドメインの出現順→論文内で紹介したものは本文内の出現順・紹介していないものは学会別・発表時期の新しい順となっています。
- ・本リストは論文執筆作業中の参考材料として、理解可能な範囲で整理したものです。

ドメイン	発表学会／掲載誌	論文タイトル (ハイライトは本文に記載したもの)	著者	概要	課題設定／着眼点	手法	手法の一部となっているデータ (学習データ、辞書など)	入力データ	出力	URL
化学	言語処理学会年次大会（2019年）	外部知識源を使用したWikipediaからの化合物情報抽出	渡土名朝飛, 野中尋史, 小林皖雄, 関根聡	日本語版Wikipediaの化合物記事から、原材料・製造方法を抽出する。化合物の属性6つのうち、原材料・製造方法は抽出が難しいとされるが、その抽出精度を高めるために深層学習（Bi-LSTM+CRFモデル）と外部知識源から作成した化合物辞書を組み合わせる。	化合物の属性の内、抽出が難しい原材料・製造方法について精度高く抽出したい。化合物の情報抽出に深層学習モデルを適用する場合、化合物の種類が非常に多く、トレーニングデータ中での出現頻度が低い（または出現しない）化合物名が大量にあることで精度が低下しやすい。	ベースをBi-LSTM+CRFモデルとし、外部知識源から作成した化合物辞書を組み合わせることで出現頻度の低い表現による精度低下を防ぐ。辞書はWikiData、PubChem、ChEBIから構築したものと、日化辞を併用。	・WikiData、PubChem、ChEBIから構築した化合物名辞書と日化辞 ・上記辞書で化合物名を置換した文→Bi-LSTM+CRFモデルの入力とする ・学習データは、森羅プロジェクトで公開されている構造化データとそれに対応するWikipedia化合物記事	Wikipediaの化合物記事	化合物の原材料と製造方法	https://www.aclweb.org/anthology/P19-1002
化学	言語処理学会年次大会（2021年）	構文情報とラベルなしデータを用いた化学分野の関係抽出	新城大希, 徳永健伸, 牧野拓哉, 岩倉友哉	化学論文から、化学物質間の相互作用等の関係を自動抽出する。BioBERTの関係抽出精度をさらに高めるために、OpenIEを補助タスクとしたマルチタスク学習を行う。またPubMedにラベルを付与し学習に利用する。	BioBERTは大量のラベル付きデータを必要とすること、構文情報を利用していないことに改善の余地がある。	BioBERTによる関係抽出と、OpenIEによる補助タスクを組み合わせる。またPubMedにラベルを付与し学習に利用する。	・CHEMPROTなどのラベル付きデータとPubMedにラベルを付与したデータで固有表現抽出器・関係抽出器を学習 ・PubMedのアブストラクトにラベルを付与したものを、既存のラベル付きデータとともに学習に利用	化学関係の論文や特許など。（実験はChemprot、GAD、EU-ADRで行った）	化学物質間の関係	https://www.aclweb.org/anthology/P21-1002
化学	情報処理学会 NL研究発表会 第241回（2019年）	辞書を用いたコーパス拡張による化学ドメインのDistantly Supervised固有表現認識	辰巳守祐, 後藤啓介, 進藤裕之, 松本裕治	Distant Supervisionを使った自動アノテーションのノイズの除去手法と、既存の辞書を使ったコーパス拡張によるRecall向上手法	専門分野の固有表現認識では学習データのアノテーションコストが高い。それを解決する1手法として、DSによる自動アノテーションがある。だがDS NERでは教師データのノイズが多く、またRecallが低いことが課題となっている。	ノイズに関しては、4 fold cross validationによるNERの予測を行い、自動アノテーションと予測が不一致の場合はノイズとみなす。Recallに関しては、DSで生成したデータ中の特定単語と辞書の単語を入れ替えることで辞書単語を含む疑似センテンスを作り、DSデータに加えてコーパスを拡張する。	化学物質名辞書（CTD、Mesh、PubChemからの収集）、論文アブストラクト（MedlineからPubMed経由収集）、人手アノテーションデータ（ChemDNERのTestデータ）	なんらかの化学系文書（実験では論文アブストラクト）	文書内の固有表現（化学物質名など）へのアノテーション	https://www.aclweb.org/anthology/P19-1002
化学	言語処理学会年次大会（2020年）	無機材料文献からの合成プロセス抽出のための関係抽出	牧野晃平, 國吉房貴, 小澤順, 三輪誠	無機材料の合成プロセスは複数文にわたり記述されるため、文間の関係を対象とした関係抽出手法が必要となる。深層学習モデルとルールベースモデルの2つを提案。	材料開発の分野では、膨大な文献に記述された合成プロセスの解析により探索・開発時間短縮する技術へのニーズがある。合成プロセスは複数文にわたり記述されるため、文間の関係を対象とした関係抽出手法が必要となる。	深層学習モデルとルールベースモデルの2つを提案。前者はBRANのTransformer部をBERTに置き換えたモデル。後者はデータセットの観測により定義。	下記（学術文献からの無機材料合成プロセス抽出のためのグラフ表現）で作成した、文献にタグ付けされたコーパスを拡張。材料を示すタグについてサブラベルを追加した。	タグ付きの無機材料文献	文献中の用語間の関係	https://www.aclweb.org/anthology/P20-1002
化学	言語処理学会年次大会（2021年）	文献抄録中の主題材料に着目した超伝導材料に関する情報抽出	山口京佑, 旭良司, 佐々木裕	文献抄録から超伝導材料に関する情報をスロット抽出する。	効率的な材料探索を目指すMaterials Informaticsにおいては材料・物質特性について大量のデータが必要だが、多くのデータが構造化されていない。よって、文献抄録から超伝導材料に関する情報をスロット抽出するシステムの構築を目指す。	固有表現・関係・イベントを抽出するモデルと、主題材料分類モデル（それぞれニューラルネットワーク）、これらを統合しスロット抽出するルールベースモジュール。	超伝導材料に関する文献抄録1,000件に対し固有表現クラスなどを人手でタグ付けしたものの。	超伝導材料に関する文献抄録	超伝導材料のElement（元素名、化合物名など）、Doping（ドーピング操作）、Value（45%などの定量表現）、SC（超伝導特性に関連する固有表現）など	https://www.aclweb.org/anthology/P21-1002
化学	言語処理学会年次大会（2021年）	Relation Extraction Task for Inorganic Material Synthesis Procedure	Shanshan liu, 松本裕治	無機材料合成手順の手順表現抽出から関係抽出までをバイバイラインとして(joinでなく)行う。	関係抽出タスクは現実的な課題における検証が不十分であるため、無機材料合成手順の抽出という課題を選択した。	Bi-LSTMとALTOPを比較する。言語表現モデルとしてSciBERTを利用。	熱電材料に関する論文241件に人手でエンティティと関係をタグ付けしたものの。	無機材料に関する論文など	無機材料に関する「Material」「Condition」「Method」「Process」	https://www.aclweb.org/anthology/P21-1002
化学	言語処理学会年次大会（2020年）	複数の事前学習モデルを併用した化学分野の関係抽出	配合智史, 嶋田和孝, 渡邊大貴, 三浦明波, 岩倉友哉	化学ドメインで固有表現間の関係抽出を行うためのより軽量なモデルの検討。	化学分野の関係抽出にはBERTモデルが成果を上げているが、サイズが大きくなりがちであり計算時間がかかる。そこでContextual String Embeddings（CSE）を使った手法を提案する。	まずBiLSTMを使った文字レベルの言語モデルを事前学習する。関係分類にはBiLSTM-Attentionモデルを用いる。その入力として、既存の二種類の分散表現（GloVe分散表現とElmo分散表現）に加えて、CSEを利用する。	PubMed、PMC、ChemRxivからの科学論文アブストラクトと本文。	化学論文などの文書	固有表現間の関係（遺伝子と疾患、タンパク質と化合物など）	https://www.aclweb.org/anthology/P20-1002

「ドメイン依存の固有表現抽出技術の現状」付録：ドメイン依存の固有表現抽出に関連する論文リスト

- ・本リストは、ドメイン依存の固有表現抽出技術の現状を調査するにあたり、以下3つの学会の論文を参照し、その中からドメイン依存の固有表現抽出をテーマとしたものを抽出したものです。
 - ・言語処理学会年次大会（2019年3月～2022年3月）
 - ・言語処理学会論文誌（2018年1月～2021年12月）
 - ・情報処理学会NLP研究会（2018年5月～2021年9月）
 - ・電子情報通信学会テキストアナリティクス・シンポジウム（2011年7月～2021年11月）
- ・論文の本文内で紹介していない研究も含まれます。
- ・論文内で紹介した研究はタイトル部分をハイライトしています。
- ・順序は、論文内でのドメインの出現順→論文内で紹介したものは本文内の出現順・紹介していないものは学会別・発表時期の新しい順となっています。
- ・本リストは論文執筆作業中の参考材料として、理解可能な範囲で整理したものです。

ドメイン	発表学会／掲載誌	論文タイトル (ハイライトは本文に記載したもの)	著者	概要	課題設定／着眼点	手法	手法の一部となっているデータ (学習データ、辞書など)	入力データ	出力	url
化学	言語処理学会年次大会（2020年）	Contextual Subword Embeddingsを考慮した 文書からの化合物名抽出実験	関根裕人, 浦澤合, 乾孝司, 岩倉友哉	化合物名をより細かいサブワードに分割した ものを利用した化合物名抽出。	化合物名抽出では極端に長い単語や 未知語の存在が課題となる。そこで 「methyl」「amino」などのサブ ワードを利用した抽出を行う。	単語を分割したサブワードを BiLSTMを使ったモデルで処理。そ れを単語系列のBiLSTM-CRFモデル の入力とする。	CHEMDNERコーパスを実験用デー タとする。（論文Abstract1万件に化 合物を手手でタグ付けしたもの）	化学論文などの文書	化合物名 (サブワード情報を使うことにより 未知語でも抽出可能に)	https://www.semanticscholar.org/paper/Contextual-Subword-Embeddings-for-Compound-Extraction-Kaneko-Urawa-Kadota/2020/abstract
化学	言語処理学会年次大会（2020年）	自動生成した学習データを用いたマルチタ スク学習によるタンパク質と化学物質間の関係 抽出	新城大希, 西川仁, 徳永健博, 牧野拓哉, 岩倉友哉	BioBERTの精度改善	BioBERTは ・ラベル付きデータの量が限定的 (作成コストがかかる) ・関係抽出では構文情報が有用だ が、BioBERTでは利用していない	・Open IEの手法（文内から「2つ のエンティティとその関係」を抽 出）で補助タスクの学習データ作成 ・主タスク（関係抽出）とOpen IE で抽出されたペアかどうかを分類す る補助タスクを同時学習して主タ スクの精度向上を図る	CHEMPROTのうち、タンパク質と 化学物質のペアを含む文	化学論文などの文書	タンパク質と化学物質の関係	https://www.semanticscholar.org/paper/Auto-generated-learning-data-for-multi-task-learning-on-protein-chemical-relation-extraction-Shinoda-Kaneko/2020/abstract
化学	言語処理学会年次大会（2020年）	Data Augmentation Technique for Process Extraction in Chemistry Publications	Yuni Susanti, Hikaru Yokono, Hiroaki Yoshida	化合物合成プロセス抽出のためのデータ拡張	化学ドメインでのラベル付きデータ は少なく、作成コストが高い。既存 のデータを最大限活用すべくデータ 拡張を行う。	・入力文と類似する文を学習データ から探す・ ・入力文中のエンティティを、類似 文の中で同じエンティティタイプを 持つエンティティに置き換える ことにより、元の文に意味の近い文 を新たに得る。 プロセス抽出はBiLSTM-CRFで行 う。学習に拡張データを使う割合を 変化させ、抽出性能を比較する。	合成プロセス235件に対し、化学ド メインの専門家が材料や環境条件な どをタグ付けしたもの。	化合物合成プロセスを含む文書	化合物合成プロセス	https://www.semanticscholar.org/paper/Data-Augmentation-Technique-for-Process-Extraction-in-Chemistry-Publications-Yuni-Susanti-Hikaru-Yokono-Hiroaki-Yoshida/2020/abstract
化学	言語処理学会年次大会（2020年）	無機化合物を対象とした論文に対する化学物 質名抽出システムの性能分析	町光二郎, 吉岡真治	生命医化学分野で学習した言語モデルを用 いた化学物質名抽出	機械学習ベースのシステムでは化学 物質名抽出の際、コーパスに存在し ない無機化合物の再現率が低い。 ニューラル言語モデルベースのシ ステムとサブワードによる単語分解の 枠組みを利用することで再現率を高 められるのではないかと。	・BioBERTを抽出に利用 ・WordPieceをサブワードに利用	・BioBERTをCHEMDNERで学習 ・評価データには、ナノ結晶デバイ ス開発分野の論文5件に化学物質名や 物質特性をタグ付けしたもの。	化学論文などの文書	化学物質名	https://www.semanticscholar.org/paper/Performance-analysis-of-chemical-named-entity-recognition-system-for-inorganic-compounds-Machino-Yoshioka/2020/abstract
化学	言語処理学会年次大会（2020年）	学術論文からのポリマー・溶媒の固有表現お よび溶解性の自動抽出	山口泰弘, 進藤裕之, 松本裕治	ポリマーと溶媒のスパン予測（固有表現抽 出）と、ポリマーと溶媒の間の溶解性の関係 抽出	物質化学論文において、ポリマーの データは数値として表にまとめられ ている事が多いのに対し、溶解性 に関する情報はテキストに記述され ることが多い。これを機械学習モデ ルで自動的に抽出したい。	・固有表現抽出はBiLSTM-CRF ・関係抽出はBiLSTM ・固有表現抽出・関係抽出ともに、 単語埋め込みにchar-CNNとBERT- Base、SciBERTを用いて性能比較	ポリマーと溶媒のスパンをタグ付 けた599文	物質化学論文などの文書	ポリマー、溶媒、溶解性の関係	https://www.semanticscholar.org/paper/Auto-extraction-of-polymer-solvent-expression-and-dissolution-property-from-scientific-paper-Yamaguchi-Suetsugu-Matsumoto/2020/abstract
化学	言語処理学会年次大会（2020年）	Extraction of Inorganic Material Synthesis Procedure from Literature	Liu Shanshan, Fusataka Kuniyoshi, Jun Ozawa, Masaki Kiyono, Yuji Matsumoto	無機化合物からの情報抽出において、従来 の抽出対象は固有表現およびアクショングラフ 程度にとどまっていた。これを一歩進め、合 成時の圧力など、素材合成工程全体の抽出に 取り組む。	無機化合物合成は複数の工程を経て 行われる。各工程は、何らかの素材 に対し、何らかの条件下で、何らか の加工をすることである。したがっ て、素材、条件、加工を抽出した 後、それらの関係を抽出すること で、化合物合成過程抽出ができると 考えた。	材料合成過程を記述したデータに、 エンティティと関係をタグ付けた もの（Kuniyoshi et al.2019）を使 う。固有表現抽出部分と関係抽出部 分で構成。前者はELMo、BiLSTM、 CRFを組み合わせたもの。加工エン ティティにはルールベースも併用。 加工表現は17種類に分類した。関係 抽出部分は複数手法を比較。	材料合成過程を記述したデータに、 エンティティと関係をタグ付けた もの（Kuniyoshi et al.2019）、素材 エンティティを含むデータ（Kim et al.2018）	無機化合物合成過程を記述した文書	無機化合物、その材料、加工時の環 境条件、加工内容（mix、meltなど）	https://www.semanticscholar.org/paper/Extraction-of-Inorganic-Material-Synthesis-Procedure-from-Literature-Liu-Shanshan-Fusataka-Kuniyoshi-Jun-Ozawa-Masaki-Kiyono-Yuji-Matsumoto/2020/abstract
化学	言語処理学会年次大会（2019年）	化合物の同義語辞書を用いた固有表現抽出	渡邊大貴, 田村晃裕, 二宮崇, 牧野拓哉, 岩倉友哉	化合物名抽出と化合物の言い換えをマルチ タスク学習することで、化合物名抽出の性能を 改善する。	化合物名には表記ゆれが多いため、 表現の同一性を学習する必要がある。	アテンションに基づくニューラル機 械翻訳（ANMT）で化合物を言い換 える。PubChem名称辞書の同一IDの 化合物ペアを学習して言い換えモデ ル	PubChem名称辞書（言い換えモデル の教師データ作成に利用）	化学ドメインの文書	化合物名	https://www.semanticscholar.org/paper/Compound-named-entity-recognition-using-synonym-dictionary-Watanabe-Tamura-Fukushima-Kaneko/2019/abstract
化学	言語処理学会年次大会（2019年）	学術論文からのポリマー-溶解性データの自動 抽出	岡博之, 吉澤篤志, 進藤裕之, 松本裕治, 石井真史	学術論文からのポリマーとその溶解性（どの 溶媒に対し可溶性を持つか）の関係抽出	現在ポリマーデータを手で抽出し DB化している。それを効率化した い。	・ポリマー名抽出はルールベース ・溶媒名抽出は辞書マッピング ・関係抽出はルールベース	・溶媒名辞書（135件）	化学ドメインの文書	ポリマー名-良溶媒名の関係	https://www.semanticscholar.org/paper/Auto-extraction-of-polymer-solubility-data-from-scientific-paper-Ogawa-Yoshizawa-Suetsugu-Matsumoto-Ishii/2019/abstract
化学	言語処理学会年次大会（2019年）	化学ドメインにおける教師無し固有表現抽出	辰巳守祐, 進藤裕之, 松本裕治	化学物質名の固有表現抽出において、以下を 明らかにする。 ・分散表現抽出において、文／単語／サブ ワード／文字ベースのうちの、どの処理単位が 最も有効か ・Distant Supervisionで生成された擬似コー パスからどのようにノイズを取り除くか	化学物質名の固有表現抽出では （1）新しい化学物質が日々生み出 される→文脈情報を用いた分散表現 を使う手法が考えられてきた。で は、分散表現抽出はどのような単位 で行うべきか。 （2）辞書やコーパス作成に専門知 識が必要 → Distant Supervision → ノイズが生じるのでいかに取り 除くか。	（1）Flairによる分散表現抽出を文 字ベースとサブワードベースで行っ たうえでBiLSTM-CRFでの固有表現 抽出を行い、性能を比較。 （2）学習済みの固有表現抽出器に 推論させてFalse Positiveとなる単語 を除去する。	実験では以下を使用。 ・分散表現の事前学習データとして Medline ・固有表現抽出器の学習データとし てChemdNERとMedline ・Distant Supervisionに使う化学物 質辞書としてCTDとMeSH（計40万 語）	化学ドメインの文書	化学物質名	https://www.semanticscholar.org/paper/Unsupervised-named-entity-recognition-in-chemical-domain-Matsumoto-Tsukagawa-Tsukagawa/2019/abstract

・本テストは調査執筆作業中の参考材料として、理解可能な範囲で重複したものです。

3 / 7 ページ

「ドメイン依存の固有表現抽出技術の現状」付録：ドメイン依存の固有表現抽出に関連する論文リスト

- ・本リストは、ドメイン依存の固有表現抽出技術の現状を調査するにあたり、以下3つの学会の論文を参照し、その中からドメイン依存の固有表現抽出をテーマとしたものを抽出したものです。
 - ・言語処理学会年次大会（2019年3月～2022年3月）
 - ・言語処理学会論文誌（2018年1月～2021年12月）
 - ・情報処理学会NLP研究会（2018年5月～2021年9月）
 - ・電子情報通信学会テキストアナリティクス・シンポジウム（2011年7月～2021年11月）
- ・論文の本文内で紹介していない研究も含まれます。
- ・論文内で紹介した研究はタイトル部分をハイライトしています。
- ・順序は、論文内でのドメインの出現順→論文内で紹介したものは本文内の出現順・紹介していないものは学会別・発表時期の新しい順となっています。
- ・本リストは論文執筆作業中の参考材料として、理解可能な範囲で整理したものです。

ドメイン	発表学会／掲載誌	論文タイトル (ハイライトは本文に記載したもの)	著者	概要	課題設定／着眼点	手法	手法の一部となっているデータ (学習データ、辞書など)	入力データ	出力	URL
医療・薬事	言語処理学会年次大会（2020年）	医薬品添付文書からの薬剤情報抽出システム	小島諒介, 岩田清明, 中津井雅彦, 奥野恭史	市販薬の添付文書のPDF群から、効能などの情報を抽出できるシステム。	市販薬の情報は新薬開発効率化に有用だが、企業横断の情報源が乏しい。医薬品添付文書のPDFはまとめてアクセスしやすいため、そこから情報抽出するシステムを構築する。	・PDFからのテキスト・レイアウト抽出 ・抽出対象の情報のアノテーション（添付文書をSGML化したデータを代替とした） ・「比較的規模の小さい」、 「BiLSTMを含む」ネットワークでの深層学習 ・クエリに対し回答を返すインターフェース	・添付文書をSGML化したデータ	市販薬添付文書（または薬剤関連文書一般）、クエリ	クエリへの回答（薬剤の効能、禁忌など）	https://www.researchgate.net/publication/350504644
医療・薬事	言語処理学会年次大会（2019年）	データベースの説明文を利用した薬物相互作用抽出	浅田真生, 三輪誠, 佐々木裕	薬物相互作用抽出に、薬物データベース（DrugBank）にある薬物説明文を利用する。	薬学論文等で報告される薬物相互作用を効率よくデータベースに登録するため、深層学習を使った自動抽出手法が検討されている。そこに既存の薬物データベースを活用したい。	DrugBankの薬物説明文と、相互作用抽出元となる文書それぞれをCNNで表現し、ふたつのCNNを同時に学習する。	DrugBankの薬物説明文（薬物名に対し説明文が対応）	薬学論文等	薬物間の相互作用（動態的作用、薬力学的作用、併用の際の推奨、相互作用有無）	https://www.researchgate.net/publication/346041444
医療・薬事	自然言語処理	病名アノテーションが付与された医療テキスト・コーパスの構築	荒牧 英治, 若宮 翔子, 矢野 憲, 永井 晋之, 岡久 太郎, 伊藤 薫	電子カルテに対し病名（および当該カルテの患者においてその症状が発生しているかどうか）をアノテーションした、45,000テキストの大規模コーパス構築。	電子カルテへの病名アノテーションの詳細な仕様を作成し、フィージビリティを検討する。また、作成したコーパスを使った病名抽出器を構築しアノテーションを検証。	電子カルテへの病名アノテーションの詳細な仕様を作成し、フィージビリティを検討する。また、作成したコーパスを使った病名抽出器を構築しアノテーションを検証。	日本内科学会に報告された44,761件の症例報告。あらかじめ少数のデータを作成し、それを機械学習することで完全に仮のタグを付け、効率化を図った。	症例報告	病名、およびその症状が発生しているかどうか	https://www.istage.ist.go.jp/article/plnp/25/1/25_119/pdf/-char/ja
医療・薬事	電子情報通信学会 テキスト・シンポジウム 第9回（2016年）	Twitterを用いた皮膚障害情報の抽出	阿部健一, 吉田博哉	消費者被害、なかでも皮膚障害被害の拡大を未然に防ぐべく、Twitterから皮膚障害情報を抽出したい。最終的には、被害を出している商品名や企業名を特定したいが、ここではその前段階として、消費者被害に関するツイートを特定するシステムを構築。	化粧品等による健康被害を早期に発見するため、Twitterから被害事例を抽出したい。ここでは被害事例を含むらしいツイートの特定を行う。	危険表現を含むツイートの中からノイズとなる文やTwitterユーザーを排除し、皮膚障害らしき（真朝度）を付与してコーパスを構築。システムのユーザーが入力する危険表現を含むツイートを出力する。	・皮膚障害に関する危険表現を取得するためのキーワードとして、「かぶれた」等6種類の語を使用。	Twitterデータ、クエリ（危険表現）	危険表現を含むツイート	https://www.researchgate.net/publication/302202624
企業情報・金融	言語処理学会年次大会（2020年）	ニュース記事からの企業キーワード抽出	奥田裕樹, 高橋寛治	Sansanでユーザー向けに配信しているニュース記事（新聞や通信社の記事+企業プレスリリース）から、「企業活動の中で生まれたモノやサービスを表す名称」を抽出。	Sansanでは、ユーザーが名刺に記載された企業の情報を検索する際、「名刺に書かれている内容以上の情報での検索」を可能にしている。そのために必要となる企業キーワードの収集を自動化したい。	ニュース記事の中で鉤括弧で囲まれた部分を抽出し、それが企業キーワードであるかどうかの二値分類を行う。判定には、BiLSTM-CRFのうち、Contextual String Embeddings（CSE）による系列モデルを適用した。	2019年に配信されたニュース記事のうち3,978件からカギカッコ部分7,225件抽出、企業キーワードかどうかを手で判定したもの。	Sansanでユーザー向けに配信しているニュース記事（新聞や通信社の記事+企業プレスリリース）	企業キーワード	https://www.researchgate.net/publication/350504644
金融・企業情報	言語処理学会年次大会（2020年）	StruAPを用いた金融分野の開示文書からの情報抽出	柳井孝介, 佐藤美沙, 十河泰弘, 山脇功一, 渋谷淳	有価証券報告書などから、売上、利益、キャッシュフローなど投資家が投資判断に使える情報を抽出。	有価証券報告書や有価証券届出書は膨大に開示されており、投資家が必要とする情報は1文書あたり50～80項目と、人手での理解にコストがかかる。これら文書はXBRLで公開されており、「〇〇をベンチマークとして」のような定型的表現も多いため、これらを活用できる。また提案手法であるStruAPを使うことで、表現そのものではなく、木構造のパターンを使った抽出が可能になっている。	抽出したいものの自体の辞書ではなく、木構造のパターンと関係を表す表現の辞書でマッチングする。	・木構造パターン474件 ・辞書80語	有価証券報告書などの金融文書	「売上高」「キャッシュフロー」等の経営指標名、および「8.0%」のような経営指標値。	https://www.anlp.jp/proceedings/annual_meeting/2020/pdf/-/B1-4.pdf
金融・企業情報	自然言語処理（2020年）	金融・経済ドメインを対象とした言語処理	坂地 泰紀, 和泉 潔, 酒井 浩之	金融・経済ドメインでの自然言語処理の状況調査	-	-	-	-	-	https://www.istage.ist.go.jp/article/plnp/27/4/27_951/pdf/-char/ja
金融・企業情報	自然言語処理（2021年）	クラウド名刺管理サービスに関連する自然言語処理の取り組み	高橋 寛治, 真鍋 友則	名刺管理サービスに関連して行われる自然言語処理の活用事例紹介（本表では企業名抽出について述べる）	名刺交換した相手の社名とニュースを紐付け、ビジネス上の付き合いを与えたい。	ルールベース（ \hookrightarrow 辞書マッチ）および機械学習（Transformerモデル）	企業名辞書、機械学習で間違いやすい項目のブラックリスト	ニュース記事	該当する企業に関連するニュース	https://www.istage.ist.go.jp/article/plnp/28/1/28_297/pdf/-char/ja

「ドメイン依存の固有表現抽出技術の現状」付録：ドメイン依存の固有表現抽出に関連する論文リスト

- ・本リストは、ドメイン依存の固有表現抽出技術の現状を調査するにあたり、以下3つの学会の論文を参照し、その中からドメイン依存の固有表現抽出をテーマとしたものを抽出したものです。
 - ・言語処理学会年次大会（2019年3月～2022年3月）
 - ・言語処理学会論文誌（2018年1月～2021年12月）
 - ・情報処理学会NLP研究会（2018年5月～2021年9月）
 - ・電子情報通信学会テキストアナリティクス・シンポジウム（2011年7月～2021年11月）
- ・論文の本文内で紹介していない研究も含まれます。
- ・論文内で紹介した研究はタイトル部分をハイライトしています。
- ・順序は、論文内でのドメインの出現順→論文内で紹介したものは本文内の出現順・紹介していないものは学会別・発表時期の新しい順となっています。
- ・本リストは論文執筆作業中の参考材料として、理解可能な範囲で整理したものです。

ドメイン	発表学会／掲載誌	論文タイトル (ハイライトは本文に記載したもの)	著者	概要	課題設定／着眼点	手法	手法の一部となっているデータ (学習データ、辞書など)	入力データ	出力	url
企業情報・金融	電子情報通信学会 テキスト・シンポジウム 第3回（2013年）	企業WEBページからの企業の事業に関連するキーワードの自動抽出	勝田研一・酒井浩之	企業の事業内容を素早く把握すべく、事業関係キーワードを抽出	企業に関する知識の浅い学生が就職活動をするにあたり、自身の関心や専門分野に近い企業を探すことは難しい。そこでたとえば「テキストマイニング」で検索した結果「野村総合研究所」や「マクロミル」といった関連性の高い企業がヒットするようなシステムを作りたい。	1) 企業のWebページから名詞を抽出し頻度順に並べる 2) IDF値の低い名詞を排除することで一般的すぎる語を除く 3) ただしIDFが低い語を人手でチェックし、重要なものであれば残す 4) 出現する企業数の少ない語（特異すぎる用語）を除く 5) IDF値の低い名詞を含む語も排除する ※たとえば「凸版印刷」において、「印刷」は1で抽出され、2で排除されるが3で復活し、5で「印刷ページ」「印刷画面」といった明らかに不適切なものを排除できる。		ニュース記事	企業に関連性の深いキーワード	https://www.researchgate.net/publication/260211126_Automated_Keyword_Extraction_from_Enterprise_Web_Pages
機械加工	言語処理学会年次大会（2021年）	機械加工文書における用語入れ子構造とトリガードを考慮した用語関係同時抽出	稲熊陸, 小島大, 東孝幸, 三輪誠, 古谷克司, 佐々木裕	機械加工技術文書内の「切削速度が増加すると切削温度が増す」といった因子（切削速度・切削温度など）とそれらの関係（Aが増すとBも増す など）を、トリガード（「物理量を示す二用語間の変化を表す単語」）を考慮して抽出	ものづくり現場の高齢化により、技術継承が難しくなっている。熟練を要する業務に、工程策定業務があり、それを行うには機械加工因子間の関係の知見が必要となる。そのため、機械加工関連の文書から、機械加工因子と、因子間の関係を同時に抽出できる技術が求められている。	入れ子構造を考慮した用語抽出と、「トリガード（物理量を示す二用語間の変化を表す単語）」を考慮した関係抽出を行う。	・切削加工に関する教科書文に対し、用語同士の関係ラベル（例：Aが大きくなるとBが大きくなる→Positive、など4パターン）を付与したもの ・各文のトークン分割のために、Sentencepieceを日本語版Wikipediaで事前学習する。	機械加工の切削加工に関する教科書文（文中に関係ラベルを付与）	機械加工用語、用語間の関係4種類（正の相関・負の相関・AはBの一種の関係・定性的な関係）	https://www.researchgate.net/publication/350204246_Automated_Keyword_Extraction_from_Enterprise_Web_Pages
機械加工	言語処理学会年次大会（2020年）	入れ子構造を考慮した機械加工用語抽出	稲熊陸, 小島大, 東孝幸, 三輪誠, 古谷克司, 佐々木裕	機械加工分野の技術者の判断支援や知見の継承支援のために知識ベースを作りたい。その第一歩として、機械加工文書から機械加工用語とその関係を抽出する。	機械加工用語は「切削加工」「切削」「加工」を含むように入れ子構造になっている事が多い。	BERT→畳み込みニューラルネットワーク→トークン数ごとの用語抽出を行う。	機械加工分野の教科書2,881文に対し人手でアノテーションしたものを。	機械加工文書	機械加工用語	https://www.researchgate.net/publication/350204246_Automated_Keyword_Extraction_from_Enterprise_Web_Pages
文学（小説）	言語処理学会年次大会（2021年）	小説あらすじを用いて学習した系列ラベリングモデルによる小説本文からの人物情報抽出の性能検証	岡裕二, 安藤一秋	あらすじのテキストで事前学習したモデルを使い、小説から人物の性別や年齢、職業などを抽出する。	ライトノベルなどの作品が増え、作品を探す労力が増大している。特に、小説の内容に読み込んだ検索機能が実装されていない。小説内の人物相関図やあらすじの生成を目指すことでそれに資する。	小説のあらすじデータにタグ付けし4つの深層学習モデルで学習、固有表現抽出。	・小説のあらすじデータ（NihのWebcat Plusから、Wikipediaの日本の小説家一覧の小説家名で検索したもののうち、同「日本のファンタジー作家一覧」にある作家の作品、または「ファンタジー」という単語を含むもの。 ・小説の本文データ（なろう小説API利用） ・人手によるタグ付け	小説の本文またはあらすじ	文中にある名前、性別、年齢、容姿、職業、所属、場所、人物関係など	https://www.researchgate.net/publication/350204246_Automated_Keyword_Extraction_from_Enterprise_Web_Pages
文学（小説）	言語処理学会年次大会（2020年）	系列ラベリングによる小説のあらすじからの人物情報抽出の検討	岡裕二, 安藤一秋	小説のあらすじから人物の属性などを抽出	同上。本稿ではあらすじからの人物情報・人物関係表現抽出手法を検討。	小説あらすじテキスト（1008件、約5000文）に対し、名前・性別・年齢表現をタグ付け。CRFでラベリングを行う。素性として、表記等とともに文字uni-gram、bi-gram等を加え、8パターンの素性組み合わせを作った性能比較。	タグ付けしたあらすじデータ	小説のあらすじ	登場人物の名前、性別表現、年齢表現、容姿や特性表現、職業や立場表現、組織・種族名、その他（異星人、神等）、地名や建物名、人物関係表現	https://www.researchgate.net/publication/350204246_Automated_Keyword_Extraction_from_Enterprise_Web_Pages
食	電子情報通信学会 テキスト・シンポジウム 第13回（2018年）	レストラン・レビューにおける食べ物・飲み物表現の抽出	新堂安孝, 友利涼, 富田祐平, 兼村厚範, 森信介	レストランレビューをマーケティングに利用するため、「そば」のような単純な表現ではなく「「香り高くのごとし抜群のおいしい十割そば」のような長く複雑なものを抽出したい。現状ではその難しさを定量的に表したデータもないため、その評価も含めて実験を行う。	レストランレビューをマーケティングに利用するため、固有表現抽出のRecall、Precisionともに改善したい。「そば」のような単純な表現ではなく「「香り高くのごとし抜群のおいしい十割そば」のような長く複雑なものを抽出したい。現状ではその難しさを定量的に表したデータもないため、その評価も含めて実験を行う。	飲み物・食べ物表現を人手でアノテーションしてモデルの学習、開発、評価を行う。モデルはナイーブなCRFベースと、DNNを用いたCRFベースを比較。	・2016年の食べログレビューデータ230万件を元に50ジャンルから一定数のレビューを抽出したもの ・上記から飲み物・食べ物表現を人手でアノテーションしたもの	食べログのレストランレビュー	「香り高くのごとし抜群のおいしい十割そば」のように、食べ物・飲み物＋その性質を表す表現。	https://www.researchgate.net/publication/350204246_Automated_Keyword_Extraction_from_Enterprise_Web_Pages

- 「ドメイン依存の固有表現抽出技術の現状」付録：ドメイン依存の固有表現抽出に関連する論文リスト
- ・本リストは、ドメイン依存の固有表現抽出技術の現状を調査するにあたり、以下3つの学会の論文を参照し、その中からドメイン依存の固有表現抽出をテーマとしたものを抽出したものです。
 - ・言語処理学会年次大会（2019年3月～2022年3月）
 - ・言語処理学会論文誌（2018年1月～2021年12月）
 - ・情報処理学会NL研究会（2018年5月～2021年9月）
 - ・電子情報通信学会テキストアナリティクス・シンポジウム（2011年7月～2021年11月）
 - ・論文の本文内で紹介していない研究も含まます。
 - ・論文内で紹介した研究はタイトル部分をハイライトしています。
 - ・順序は、論文内でドメインの出現順→論文内で紹介したものは本文内の出現順・紹介していないものは学会別・発表時期の新しい順となっています。
 - ・本リストは論文執筆作業中の参考材料として、理解可能な範囲で整理したものです。

ドメイン	発表学会／掲載誌	論文タイトル (ハイライトは本文に記載したもの)	著者	概要	課題設定／着眼点	手法	手法の一部となっているデータ (学習データ、辞書など)	入力データ	出力	URL
食	情報処理学会 NL研究発表会 第237回（2018年）	文字分散表現に基づく単語分類情報を用いた レシビ固有表現抽出	平松 淳, 若林 啓, 原島 純	ドメイン（この例では料理）に関連する言語 資源を使った固有表現抽出	固有表現抽出の教師データをドメイン ごとに構築するのはコストが大き い。よって、文中の単語をカテゴリ に分類し、分類情報を固有表現抽出 器の入力として利用する。	LampleらのBiLSTM-CRFの処理に、 単語分類器からの情報を追加する。 文中の単語について、オントロジー での属性ラベルを予測する分類器を 学習し、固有表現抽出器の特徴量に 組み込む。	・ レシビNEコーパス（笹田ら） （クックパッドの手順データに固有 表現を付与したもの） ・ 料理オントロジー（Nanbaら） （料理関係の用語に上位下位関係、 同義語などを付与したもの）	料理ドメインのテキスト （長期的には、より一般的なドメイ ンに適用したい）	料理関係の固有表現	https://www.proceedings.nlp.dcc.ac.jp/2018/papers/p101.pdf
食	人工知能									
みやげ品	情報処理学会 NL研究発表会 第243回（2019年）	固有表現抽出によるブログテキストからの品 名・店名抽出	池田 流弥, 安藤 一秋	「現地でしか購入できない土産」の情報をブ ログなどのUGCから抽出する手法	現地でしか購入できない土産情報を Webから収集するシステムを構築す るにあたり、既存の固有表現抽出手 法の多くはUGCや日本語テキストに 対し未評価である。よって、日本語 ブログから構築したデータを用い、 CRFと深層学習による固有表現抽出 の性能評価を行う。	CRFとBiLSTM-CRFの2モデル、 BiLSTM-CNN-CRF、Char- BiLSTM-CRF。	・ 日本の著名な土産をまとめた OMIYAIによる土産名をクエリとし、 Yahoo!ブログの菓子・デザート カテゴリでヒットしたブログ記事の うち680エントリに入手で固有表現 タグを付与したもの。	ブログ記事	土産名、土産／菓子店名	https://www.proceedings.nlp.dcc.ac.jp/2019/papers/p101.pdf
みやげ品	言語処理学会年次大会（2019年）	深層学習によるブログ記事からの土産の品 名・店名抽出	池田流弥, 安藤一秋	「現地でしか購入できない土産」の情報をブ ログなどのUGCから抽出する手法	多くの人が旅行の際に「現地でしか 購入できない土産」を求めている が、「現地でしか購入できない」と いう情報がほとんどない。そこで Webから収集した情報をもとに、土 産品情報を抽出したい。 （本研究では「現地でしか購入でき ないかどうか」は判定せず、土産品 情報全体を抽出する）	Yahoo!ブログから土産情報の書かれ た記事を収集し、土産の品名・店名 を人手でタグ付けする。BiLSTM- CRFでの抽出とCRFのみでの抽出を 比較する。	・ 土産サイトOMIYAIの土産名リス ト（ブログ記事収集時のクエリとし て使用） ・ ブログ記事（680件）への土産品 名・店名にタグ付けしたデータ（学 習用） ・ Wikipedia本文全体（分散表現構築 用）	ブログ記事	土産名、土産／菓子店名	https://www.proceedings.nlp.dcc.ac.jp/2019/papers/p101.pdf
交通	言語処理学会年次大会（2020年）	オントロジー形式アノテーションを対象とし た交通用語・関係抽出と正誤問題の回答	鈴木直樹, Savong Bou, 三輪誠, 佐々木裕	交通関係の文書からの用語関係抽出にあた り、タグ付けデータの形式を変えることで精 度を向上させた。	従来の交通関係用語抽出／関係抽出 は精度が不十分だった。用語分類に おいて似た意味のタグがあること、 関係の種類が細かすぎたことが要因 ではないか。	似た意味を持つタグを整理し、関係 を用語としてタグ付け（オントロ ジー形式でアノテーション）する。 関係の種類は47種類→5種類に削 減。AlanらによるFlairで用語抽出。	オントロジー形式でタグ付けした データ。	交通法規に関する文書	交通用語と用語間の関係	https://www.proceedings.nlp.dcc.ac.jp/2020/papers/p101.pdf
交通	言語処理学会年次大会（2019年）	CNNを用いた交通教則からの交通用語関係抽 出	八木智也, 三輪誠, 佐々木裕	自動運転での利用を目指した交通オントロ ジーのための、交通用語同士の関係抽出。	現在の自動運転プログラムには交通 法規やマナーなどの知識がコード内 に組み込まれており、更新が難し い。そのような知識を交通オントロ ジーとして整理し、自動運転プログ ラムから参照するという構造にすれ ば、法改正や異なる国での利用も容 易になると考えられる。オントロ ジー構築には膨大な作業が必要とな るため、それを補助するようなテキ ストからの情報自動抽出が求められ る。	交通文書（約2000文）内の用語間に 66種類のタグを付け、CNNを用いた モデルで用語ペア間の関係を学習す る。	・ 交通文書へのタグ付けデータ	交通文書	文書内の用語同士の関係	https://www.proceedings.nlp.dcc.ac.jp/2019/papers/p101.pdf
その他	言語処理学会年次大会（2020年）	会議録に含まれる法律名を対象としたEnd- to-Endのエンティティリンキングの性能評価	松森拓真, 木村泰知, 荒木健治	会議録から法律名（皮肉表現など多様な表記 ゆれ含む）を抽出。	法律名には「働き方改革関連法」を 「過労死促進法」と呼ぶなど、人手 でも難しい表記ゆれがある。その ため、メンションを「法律名を表 す語句」とし、法律名のメンション 抽出および曖昧性解消の両方（End- to-End Approach）を行う。	“国会・地方議会会議録上の法律名部 分にタグ付け（メンションとそれに ひもづくWikipedia記事のアノテ ーション）を行う。 deppavlov（BERT）で固有表現抽 出モデルを作成してメンションを抽 出。それをWikipedia記事と結びつけ て曖昧性解消を行う。”	“・国会・地方議会会議録10日分（5 つの検査語にそれぞれ2日分）に対し、 人手でタグ付けを行った。 ・タグ付けの中で、各メンションと 対応するWikipedia記事を結びつけ た。”	法律名を含む会議録	法律名と、それに対応するWikipedia 記事。	https://www.proceedings.nlp.dcc.ac.jp/2020/papers/p101.pdf
その他	人工知能学会全国大会（2019年）	ECサイトにおける商品タイトルからの商品名 抽出	張 瑞楠	ECサイトの商品タイトルから商品名だけを抽 出。	商品タイトルはSEO対策のために商 品名以外の要素が多く付随しており わかりにくい。一般的な文章と違い 名詞・名詞句の羅列であることが多 く、このタスクに特化した手法が必 要。	Term Weighting問題としてTF-IDF を使う手法と、系列ラベリング問題 としてCRFを使う方法、BiLSTM- CRFを使う手法を比較した。	Yahoo!ショッピングから抽出した商 品タイトル1万5000件にタグ付けし たもの。	商品タイトル （例：【純正品】HPインクカート リッジ【送料無料】）	商品名 （例：インクカートリッジ）	
その他	人工知能学会全国大会（2020年）	構文解析情報を用いたテキストからの数値情 報の抽出	黒土 健三, 森本 康嗣, 佐藤 美沙, 柳 井 孝介	論文から技術トレンドを把握するための数値 情報抽出。ここでは応用物理分野の論文を対 象とする。	科学論文中の数値情報からは、ムー アの法則のような技術トレンドを読 み取れるのではない。	StruAP（係り受け構造に基づくルー ルベースの抽出ツール）を用いる。	論文 （実験では半導体パッケージング技 術に関する論文）	数値と項目名のペア		

・本テストは調査執筆作業中の参考材料として、理解可能な範囲で重複したものです。

<http://www.informaworld.com/0000-0000>