難病・希少疾患症例コーパス構築の課題と効率化に向けた LLM・ソールの活用

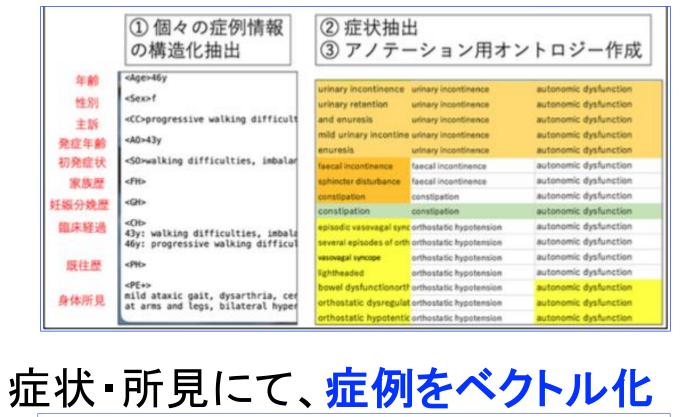
土肥栄祐¹⁾、金進東²⁾、早川格³⁾、松旗知康⁴⁾、高月照江²⁾、建石由佳²⁾、藤原豊史²⁾、山本泰智²⁾

神経内科、4) 広島大学神経内科 1) NCNP 神経研究所 疾病研究第三部、2) DBCLS、3) NCCHD

難病・希少疾患は約1万種存在し、症例数の少なさから診断まで時間を要し、特に 非典型例では時間がかかるため、質の高い症例コーパス構築やその構築の手法 開発が望まれている。本研究では、日本語症例報告に疾患名・症状名をタグ付けし た高品質コーパス構築において、大規模言語モデル(LLM)とWebベースのアノテー ション管理・編集ツールを組み合わせ、効率化を図った。J-STAGEに収載されている 医学文献から、NANDO辞書と形態素解析を用い、難病・希少疾患の疾患名に基づ き症例報告を抽出した。テキスト正規化後、LLMによる症状・所見のアノテーション を実施し。JSON形式出力、入力のチャンク化、テキスト部省略によるトークン節約に より、安定した結果を得た。アノテーションはPubAnnotationで管理し、GUIベースの TextAEで専門家が評価・修正を行った。Human-in-the-loopにより精度向上と作業負 担軽減を実現し、専門家参画を促進した。今後は用語オントロジーの適合性向上な ど、さらなる課題解決が必要である。

これまでの可視化の検証

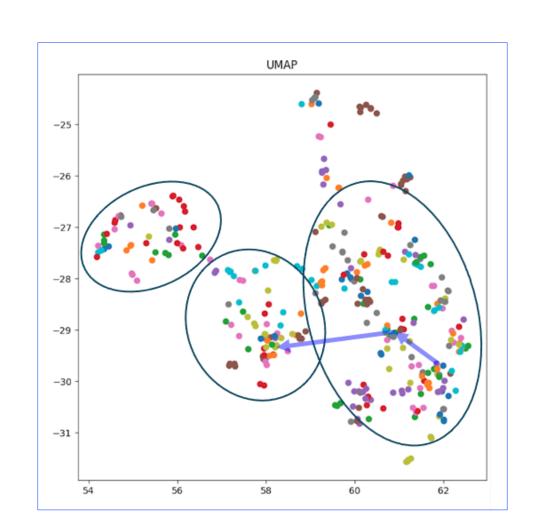
- ・アレキサンダー病の症例情報(200例)の構造化
- 938個の症状・所見 → 59個の症状所見にアノテーション



・クラスタリング

- 時系列の可視化
- 非典型例の推定



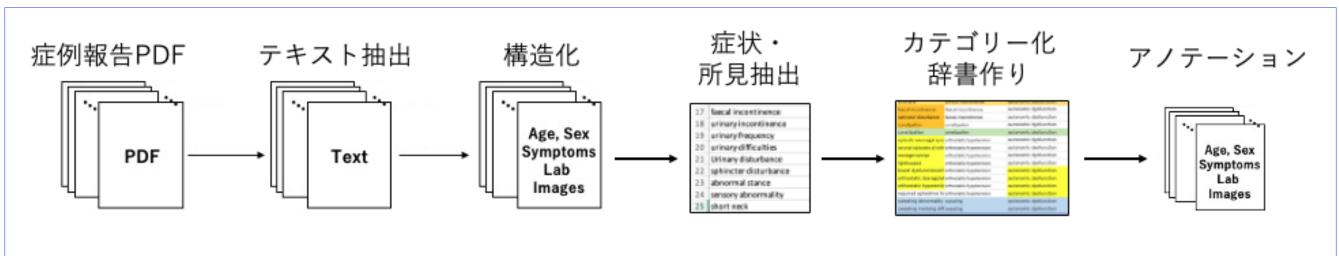


しかし、これらの工程から明らかとなった課題として、

1: 工程に時間と労力がかかる

2: ツールや、オントロジーは理解が難しく、コーパス作成に携 わるエキスパート(医療者)が、参入し難い

各工程の自動化やツールの導入を試みた



1:テキスト抽出

• 症例報告PDFを、LLM(ここではCloud版のChatGPT)にuploadし、プロンプト エンジニアリングにてテキスト抽出し、マニュアルで確認を行った。

2:構造化

・症例報告における、年齢・性別・主訴・現病歴などのカテゴリーを、大きく分 類しするためのPythonコードを作成。抽出したテキストを活用し、内部の構 造を把握し統計情報を取得したのちに、構造化されたテキスト形式で出力。

【既往歴】30 歳時に小麦アレルギー。

複数症例がある場合→仕分け カテゴリーの表記揺れ→標準化

テキストの特徴を解析

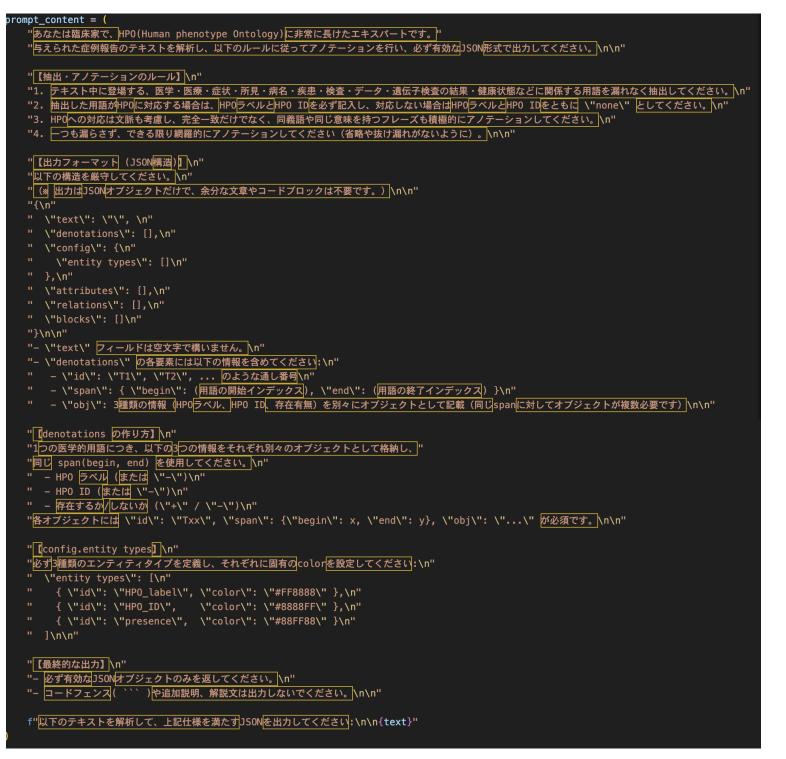
<除外基準>

•100文字以内 ・標準カテゴリーが2つ以下

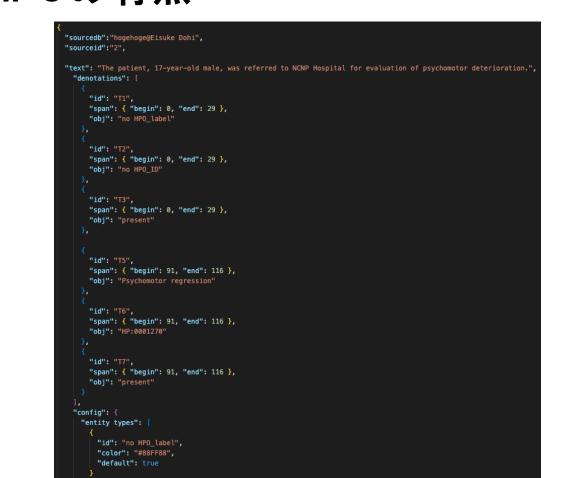
段階を踏むことで構造化

3:LLMによるドラフトアノテーション

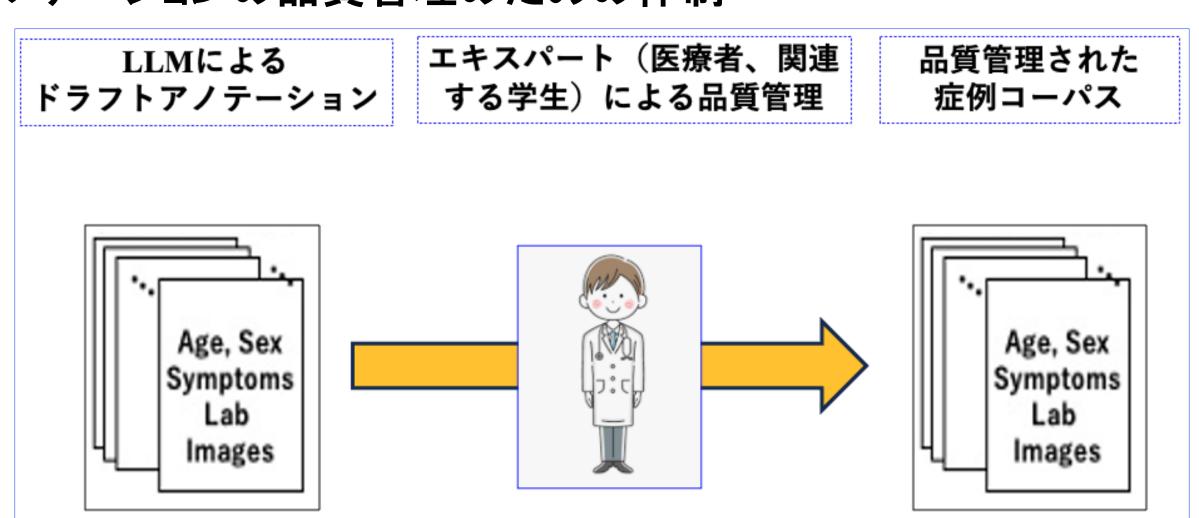
プロンプトエンジニアリングを用い、ドラフトアノテーション



- チャンク化(カテゴリー、句読点) ○ JSON出力(PubAnnotation対応)
- 症状・所見の有無
- 時間 ○ HPOの有無



4:アノテーションの品質管理のための体制



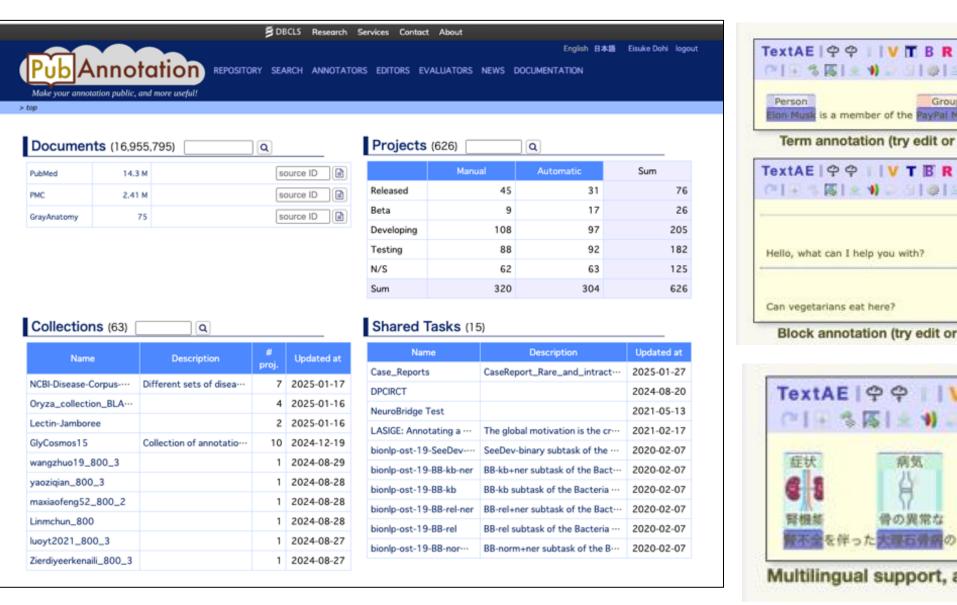
PubAnnotation:

(https://pubannotation.org/)

TextAE:

(https://textae.pubannotation.org/)

Attribute of annotation

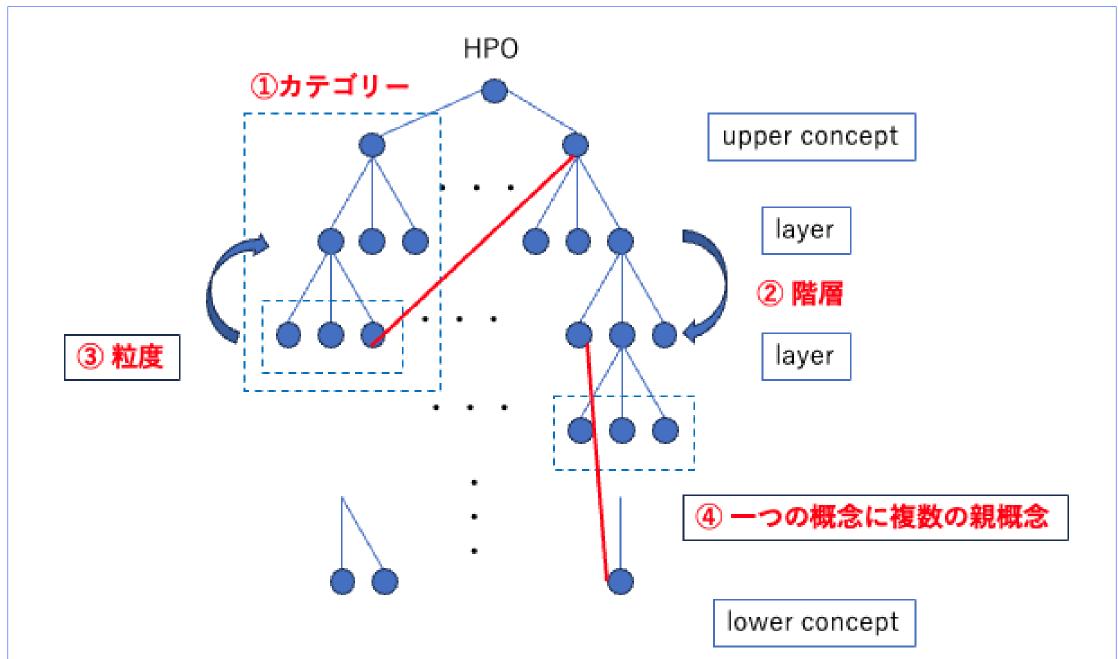


- アノテーション管理・編集が可能なウェブツール
- 複数のアノテーションの管理が可能、アノテーター間での突合が可
- ○JSON形式で、アノテーションデータのアップロード
- GUIベースでアノテーションの確認・編集が可能(TextAE)

LLMにて、PubAnnotationに対応したJSON形式のデータ出力。ドラフトアノテー ションと、ツールの導入により、エキスパートの参画はしやすい体制構築。

しかし、実際にアノテーションチェックの際に、 アノテーションの不備が、オントロジーに無いためか? オントロジーの中で最適なものが選択されているのか? この確認が困難であることが明らかとなった

5:オントロジーの視認性について



- ① カテゴリー:解剖(臓器)、システム(免疫、代謝etc)、時点など、... etc
- ② 階層:カテゴリーの対象とするものの複雑性に依存して深さが異なる。
- ③ 粒度: エンドフェノタイプを示すものが、別の階層に存在する。
- ④ 複数継承:一つの概念に複数の親がある時
- ⇒これが、俯瞰的な見方や全体像を掴むことを難しくしている

このDAGをわかり易く検索・利活用する可視化ツールを開発し、エキス パート(臨床医や医学生)のチェックの効率化を進めている

6: 達成できた課題と、残る課題

- ・実症例の蓄積+生成AI
- 元データの質の底上げ
- ・病名の表記揺れ・誤記
- ・情報の構造化(自動化) アノテーションの効率化
- · HPOの整備・可視化 or 独自基準
- ・時系列データの表現
- ・ローカルでの構築

7: まとめと展望

プロセスは技術で進歩してゆくため、生データの質、データの使い易さ・連結し やすさが、一層価値を増してくると考えている。







今後は、初期に行ったテキスト抽出などの自動化を再検討に加え、 DAGの可視化・編集ツールを横展開し下記の可能性を展開してゆく。

① 医療情報の新規可視化法による、医療-患者間コミュニケーションの改善

バックキャストによるデータ収集法の最適化

<u>③ 施工過程の可視化による医療者の学習の効率化</u>