

# 専門用語辞書管理システム および 用語の内部構造アノテーションについて

山田恵美子,

呂嘉, 松本裕治

(東京大学)

(奈良先端科学技術大学院大学)

概要: 医学生物学関連の専門用語の表記、読み、品詞、英訳、シソーラスコード等の意味クラス情報などの基本情報、および、同義語へのリンクや内部構造を記述できる辞書管理システムを開発している。現在は、ライフサイエンス辞書の日本語の用語が対象、複合語として複雑な構造をもつ専門用語の内部構造を単語単位の係り受け木構造としてアノテーションする機能を備えている。一部の用語の内部構造を本システムを用いてよってアノテーションし、それを学習データとして用いて、未解析の専門用語の内部構造の自動アノテーション実験を行った結果を報告する。

## 辞書管理システム「Cradle」

- ライフサイエンス辞書(京大金子研 2008.4版)を形態素解析用辞書として管理
- 見出し、読み、品詞などの用語の基本情報、および MeSHのシソーラスコード、および、英訳の情報を格納
- 専門用語(複合語)の内部構造アノテーション機能と 内部構造表示機能
- 検索機能(綴り、読み、品詞、構成語による検索)
- 同義語情報の付与

条件: 単語=>結核, 辞書=(WebLSD-200804*)									
ID	単語	読み	発音	BASE	ROOT	辞書	品詞	活用型	活用形
詳細 80772	結核	ケツカク	ケツカク	結核	■	■	名詞-一般		
詳細 198848	肺結核	ハイケツカク	ハイ	ケツカク	■	■	名詞-一般		
詳細 2197835	広範囲抗生物剤耐性結核	コハシキヤウザイ	コハシキヤウザイ	広範囲抗生物剤耐性	■	■	名詞-一般		
詳細 2198017	結核性腫瘍	ケツカキセイセキ	ケツカキセイセキ	結核性腫瘍	■	■	名詞-一般		
詳細 2198710	耐性結核菌	キヨウヤクサウキ	キヨウヤクサウキ	耐性結核菌	■	■	名詞-一般		
詳細 2198905	結核性	ケツカキセイ	ケツカキセイ	結核性	■	■	名詞-一般		
詳細 2200540	結核核	キヨウカクカク	キヨウカクカク	結核核	■	■	名詞-一般		
詳細 2204108	類結核性	ルイケッカクシ	ルイケッカクシ	類結核性	■	■	名詞-一般		

「結核」を含む語の検索結果画面

## 専門用語の内部構造解析

- 複合語の構造は単語からなる係り受け木によって記述可
  - D(通常の係り受け), R(逆向き係り受け), P(並列), U(その他の関係)によって表現
- 共通の文字が省略される(縮退する)現象がある
  - これを表現するため、文字間の係り受けによって解析
  - WB(語の構成要素の先頭として係る), WI(語の構成要素の内部として係る)の関係を解析用に新たに定義

## 専門用語の内部構造解析実験

- 一部の用語に単語間関係の内部構造アノテーションを行い、これを学習データとして自動解析実験を行った。
  - データ
    - ライフサイエンス辞書の掲載語から804語
  - 手法
    - Shift-Reduce法(MaltParser)[Nivre 2004]を利用
  - 結果: (5分割交差検定による精度)
    - 文字間係り受け 96.9% (ラベル無し: 97.8%)
    - 用語全体の解析 87.6% (ラベル無し: 89.9%)

[Nivre 2004] Nivre, J. and Scholz, M. "Deterministic Dependency Parsing of English Text," In Proceedings of COLING 2004, pp.64-70, Geneva, August 2004.

選択	更新	ID	単語	読み	品詞	構造	辞書	状態
①	○	2254710	コハシキヤウザイ	コハシキヤウザイ	名詞-一般	show	NEW	

選択	更新	ID	単語	読み	品詞	構造	辞書	状態
②	○	80772	結核	ケツカク	名詞-一般			CHECKED

複合語の内部構造アノテーション画面

ID	単語	読み	品詞	構造
2197835	広範囲抗生物剤耐性結核	コハシキヤウザイ	名詞-一般	

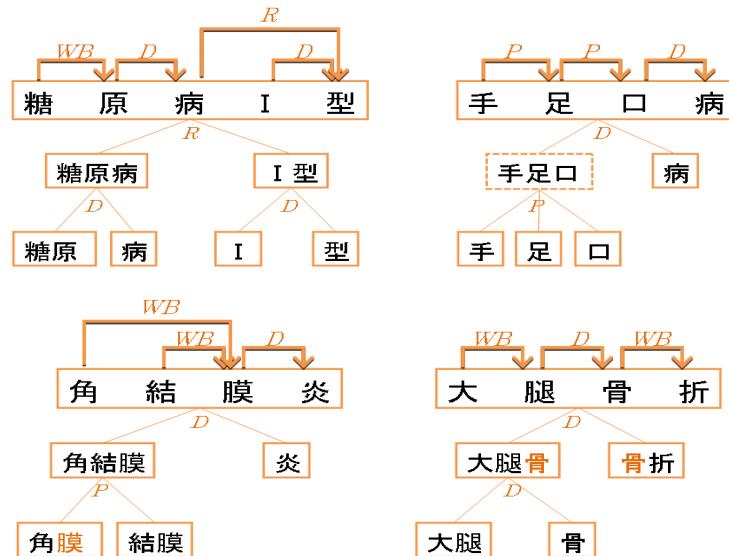
状態: NEW  
権限: 閲覧  
更新者: yamada  
更新時間: 2008-11-18 16:52:23

構造詳細

```

graph TD
    A[広範囲抗生物剤耐性結核] --> B[広範囲抗生物剤耐性]
    B --> C[結核]
    B --> D[抗生物剤耐性]
    C --> E[広範囲]
    C --> F[抗生物剤]
    C --> G[耐性]
    D --> H[耐性]
    D --> I[抗生物剤]
    E --> J[広範囲]
    F --> K[抗生物剤]
    G --> L[耐性]
    H --> M[耐性]
    I --> N[抗生物剤]
    J --> O[広範囲]
    K --> P[抗生物剤]
    L --> Q[耐性]
    M --> R[耐性]
    N --> S[抗生物剤]
    O --> T[広範囲]
    P --> U[抗生物剤]
    Q --> V[耐性]
    R --> W[耐性]
    S --> X[抗生物剤]
    T --> Y[広範囲]
    U --> Z[抗生物剤]
    V --> AA[耐性]
    W --> BB[耐性]
    X --> CC[抗生物剤]
    Y --> DD[広範囲]
    Z --> EE[抗生物剤]
    AA --> FF[耐性]
    BB --> GG[耐性]
    CC --> HH[抗生物剤]
    DD --> II[広範囲]
    EE --> JJ[抗生物剤]
    FF --> KK[耐性]
    GG --> LL[耐性]
    HH --> MM[抗生物剤]
    II --> NN[広範囲]
    JJ --> OO[抗生物剤]
    KK --> PP[耐性]
    LL --> QQ[耐性]
    MM --> RR[抗生物剤]
    
```

用語の詳細情報の提示画面(右は内部構造表示)



文字単位の係り受け解析(用語上部の矢印)と  
構成要素となる単語による木構造表現(下部の木構造)