

統合データベース開発:

専門用語辞書管理システムと専門用語 解析技術の開発

奈良先端大

松本裕治, 新保仁, 浅原正幸,
原一夫, 鈴木郁美, 呂嘉

- 専門用語解析技術
 - 専門用語辞書システムの開発
 - 専門用語解析技術の開発
- 専門用語抽出ツールの設計と開発
 - 専門用語辞書拡張支援ツールの設計と開発



今年度の成果目標

(1) 専門用語解析技術

① 専門用語辞書システムの開発

- 9月末まで:
 - 10万語以上の規模の生命科学用語の辞書を格納(ライフサイエンス辞書, 病名マスターなど)
 - Webブラウザ上で用語の検索や用語のもつシソーラスコードなどの情報を表示・修正機能
- 22年度末:
 - 用語間の意味関係(シソーラスに基づく概念的な上位下位関係や類似度)の表示機能
 - 全体的な専門用語辞書システムを完成させる。

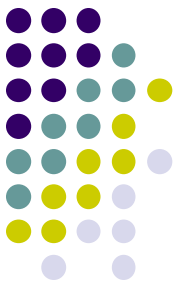


今年度の成果目標

(1) 専門用語解析技術

② 専門用語解析技術の開発

- 9月末まで
 - 学習データとして内部構造解析済みのデータを2000語以上に拡大
- 22年度末:
 - 前年度に設計した一般的な統語解析法に基づくアルゴリズムをコンピュータ上に実装し、90%以上の内部構造解析精度を達成する



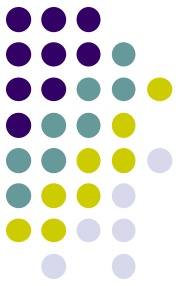
今年度の成果目標

(2) 専門用語抽出ツールの設計と開発

① 専門用語辞書拡張支援ツールの設計と開発

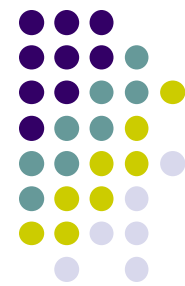
平成21年度までに、専門分野のテキストに現れる新規の専門用語(新規語)と既知語との意味的な類似度計算する手法を実装し、新規語と意味的に類似性の高い専門用語辞書(シソーラス)内の既知語候補を、シソーラスの構造とともに提示するインターフェースの設計を行った。

- 今年度9月末まで:
 - 用語の内部構造の情報など、種々の情報を用いて類似度判定の精度を向上を図る
 - 新規語に対してシソーラスコードを付与するための支援機能
- 22年度末:
 - 専門用語辞書システムの一機能としてWebブラウザ上で利用可能な機能として統合する。



専門用語解析システム

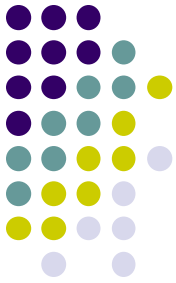
1. 専門用語辞書システム



辞書管理システムCradle

- 形態素解析用辞書の管理ツール(現状)
 - 現在登録している辞書, 用語集
 - ライフサイエンス辞書(京大金子研究室)
 - 標準病名マスターv2.80を格納
 - 仲里さんよりいただいた専門語候補(18万語)
 - 複合語に対する内部構造付与
 - 現在約1800語について人手により内部構造付与
 - 辞書の標準項目による検索以外に, 同義語検索, 内部構造に基づく検索を実装
 - 表示項目のカスタマイズ機能を実装をより一般化

検索画面



Cradle--ChaSen Dictionary Management System - Mozilla Firefox
ファイル(E) 編集(E) 表示(V) 履歴(S) ブックマーク(B) ツール(I) ヘルプ(H)
http://dahlia.naist.jp/cradle/

Cradle--ChaSen Dictionary Man...

CRADLE--茶筌辞書管理システム

日本語辞書 中文辞典 matsu | [Preference](#) | [User list](#) | [Logout](#)

単語属性

ID	=	<input type="text"/>	単語	=	<input type="text"/>
読み	=	<input type="text"/>	発音	=	<input type="text"/>
品詞	=	<input type="text"/>	活用型	=	<input type="text"/>
活用形	=	<input type="text"/>	Base	=	<input type="text"/>
辞書	or	<input type="text" value="NAIST-jdic-20080707"/> <input type="text" value="WebLSD-200804*"/> <input type="text" value="標準病名マスター-V2.80*"/> <input type="text" value="pne_kw*"/> <input type="text" value="techterm*"/>	文字数	=	<input type="text"/>
更新時間	<=	<input type="text"/>	状態	=	<input type="text"/>
親概念日本語表記	=	<input type="text"/>	新規者	=	<input type="text"/>
手動参照先の日本語コード	=	<input type="text"/>	更新者	=	<input type="text"/>
階層の深さ	=	<input type="text"/>	親概念英語表記	=	<input type="text"/>
自動参照先ID	=	<input type="text"/>	日本語コード	=	<input type="text"/>
自動参照先表記	=	<input type="text"/>	手動参照先の日本語表記	=	<input type="text"/>
親概念ID	=	<input type="text"/>	ツリー番地	=	<input type="text"/>
ICD10	=	<input type="text"/>	ツリー日本語	=	<input type="text"/>
確信度	=	<input type="text"/>	ツリー英語	=	<input type="text"/>
			頻度(文数)	=	<input type="text"/>
			頻度(文献数)	=	<input type="text"/>

複合語属性

内部表記	include	<input type="text"/>	内部読み	include	<input type="text"/>
内部PCS	include	<input type="text"/>	状態	=	<input type="text"/>

単語情報の表示



Cradle--ChaSen Dictionary Management System - Mozilla Firefox

ファイル(E) 編集(E) 表示(V) 履歴(S) ブックマーク(B) ツール(I) ヘルプ(H)

http://dahlia.naist.jp/cradle/jp/show/2223089

Cradle--ChaSen Dictionary Man...

INDIGENOUS LAB COMPUTATIONAL LINGUISTICS

CRADLE--茶筌辞書管理システム

matsu | [Preference](#) | [User list](#) | [Logout](#)

日本語辞書 中文辞典

単語詳細

ID	2223089	
単語	遺伝子発現量	
読み	イデンシハツゲンリョウ	
発音		
品詞	名詞一般	
活用型		
活用形		
BASE	遺伝子発現量	系列
ROOT		
辞書	WebLSD-200804*, pne_kw*, techterm*	
親概念日本語表記		
親概念英語表記		
手動参照先の日本語コード		
日本語コード	J058351	
階層の深さ		

構造詳細

状態	NEW
備考	
更新者	matsu
更新時間	2010-01-27 13:17:17

遺伝子発現量

構成	遺伝子, 発現量
枝の種類	D
縮退文字の位置	
省略文字の位置	none

ツリー構造

```

graph TD
    A[遺伝子発現量] --> B[遺伝子]
    A --> C[発現量]
    B --> D[遺伝]
    B --> E[子]
    C --> F[発現]
    C --> G[量]
  
```


専門用語辞書システムの開発項目

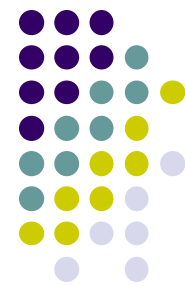


- 辞書システムの機能拡張
 - 表示項目の追加・削除, 値の属性定義のカスタマイズ
 - 辞書によって表示項目を自由に定義できるようにする
 - DBシステムをMySQLからMongoDBへ変更予定
 - 管理者権限の分類, グループ化
 - 一部のDBにのみ権限をもつ管理者
 - 一部のDB修正機能にのみ権限をもつ管理者
- など, 最上位の管理者以下に様々な異なる管理者グループを定義できるようにする



専門用語解析技術

2. 専門用語解析技術の開発



専門用語解析技術の開発項目

- 内部構造解析データの拡張
 - 現状の1800語を2000語以上に拡張
 - できれば3000語以上の用語の内部構造タグ付けを行いたい
 - 内部構造の自動解析システムの実装・評価
 - 昨年度までに設計した統語解析アルゴリズムの一般化に基づく方法を実装
 - 半教師付き手法による精度向上実験
 - これまでの実装(文字ベースの決定性アルゴリズム)との比較評価



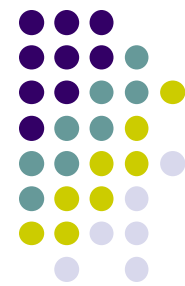
(2) 専門用語抽出ツールの 設計と開発

専門用語辞書拡張支援ツール



研究の目的＝シソーラス拡張

- 新規の専門用語に対して，類似度が高い順に登録済の専門用語をランク付けし，提示するシステムの構築
- シソーラス辞書の編集者は，システムが提示するランキング上位語を参考に，新しい専門用語をシソーラス辞書に格納する



専門用語辞書拡張支援ツール

- 専門文書から対象とする語の文脈情報を抽出して用語の隣接グラフを作成し，グラフ構造を用いて用語間の類似度を算出する手法を提案
- 雑誌「蛋白質・核酸・酵素」を実験データとして用い，そこに登場する新規の専門用語と類似度の高い語をライフサイエンス辞書から検索する
- 検索した類似度上位の語のシソーラス内の位置を表示するインタフェースを構築

類義語検索ツールの初期画面



Synonym Acquisition - Mozilla Firefox

ファイル(E) 編集(E) 表示(V) 履歴(S) ブックマーク(B) ツール(I) ヘルプ(H)

http://cl.naist.jp/~naonori-a/synonym/request.html

Synonym Acquisition

類義語マッピング

新規専門用語をシソーラスにマッピングする支援を行うシステム

クエリ: 1 件

(クエリ例: ニワトリ, 培養細胞, インターロイキン)

- システムの概要
 - シソーラスに登録されていない新規の医療、バイオの専門用語をシソーラス上にマッピング支援をする
 - 新規専門用語は「蛋白質・核酸・酵素 (PNE)」の文脈情報を使って、既に登録されている専門用語と類似性を比較する
 - 専門用語間の類似度はコサイン類似度を使用している

完了

類似の上位語とシソーラスの部分表示



Result - Mozilla Firefox

ファイル(E) 編集(E) 表示(V) 履歴(S) ブックマーク(B) ツール(I) ヘルプ(H)

http://cl.naist.jp/~naonori-a/cgi-bin/show_result.cgi?op=show&id=N94CG4sj

Result

アルブミンの類義語検索結果

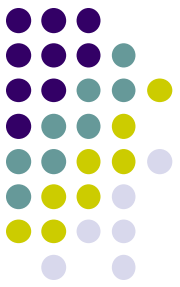
ランキング

アルブミン: [\[D12.776.034\]](#)

1位 ラミン	[D12.776.660.650.875]
2位 プラスミン	[D08.811.277.656.300.760.625]
3位 ヒドロキシルアミン	[D01.625.075.525] , [D02.092.570]
4位 ウシ血清アルブミン	[D12.776.034.841.540] , [D12.776.124.727.540]
5位 グルコサミン	[D09.067.342.531]

- 解剖学[\[A\]](#)+
- 生物[\[B\]](#)+
- 病気[\[C\]](#)+
- 化学物質と薬物[\[D\]](#)+ 1位 2位 3位 4位 5位
 - 酵素および補酵素[\[D08\]](#)+ 2位
 - 酵素[\[D08.811\]](#)+ 2位
 - 加水分解酵素[\[D08.811.277\]](#)+ 2位
 - ペプチド加水分解酵素[\[D08.811.277.656\]](#)+ 2位
 - エンドペプチダーゼ[\[D08.811.277.656.300\]](#)+ 2位
 - セリンエンドペプチダーゼ[\[D08.811.277.656.300.760\]](#)+ 2位
 - プラスミン[\[D08.811.277.656.300.760.625\]](#) 2位
- アミノ酸・ペプチド・タンパク質[\[D12\]](#)+ 1位 4位
 - タンパク質[\[D12.776\]](#)+ 1位 4位
 - 核タンパク質[\[D12.776.660\]](#)+ 1位
 - 核マトリクス結合タンパク質[\[D12.776.660.650\]](#)+ 1位
 - ラミン[\[D12.776.660.650.875\]](#)+ 1位

専門用語辞書拡張支援ツールの 今後の予定



- 9月末まで：
 - 用語の内部構造の情報などの情報も利用することにより類似度判定の精度を向上
 - 新規語に対してシソーラスコードを付与するための支援機能
- 今年度末まで：
 - 専門用語辞書システムと連携させ、辞書システムから呼び出すことにより、辞書システムの一機能として利用可能にする