

6S-6

大規模電子化辞書開発における高機能辞書エディタ
 — そのデータ構造 —

三輪和弘、中沢正幸、鈴木美穂、安達久博*
 株式会社 日本電子化辞書研究所
 株式会社 東芝

はじめに

(株)日本電子化辞書研究所(EDR)では、実際の文章データに基づいて辞書開発を行っている。つまり、辞書のすべての情報は少なくとも一つ以上の文章によって保証された情報である。この保証を常に維持することは、辞書情報あるいは文章データの変更を他のデータへ反映しなければならないことを意味する。EDR辞書は現段階においても約400メガバイトを有し、コーパスデータ及びKWICデータは約30ギガバイト近くになる予定である。高機能辞書エディタは、上記したがってこの保証を実現するためにEDRでは、高機能辞書エディタを開発している。高機能辞書エディタは、上記データをスムーズに参照しエディティングするために開発されている。また、設計思想としては、Hyper Text及びオブジェクト思考を導入している。さらに、概念辞書の開発において、概念間のリンクを行う概念統合、概念の抽象化を行う概念体系の開発にも高機能辞書エディタが不可欠である。本論文では、高機能辞書エディタの対象となるそれぞれのデータの構造とデータ間の関係、及びそのアクセス方法について報告する。

1. EDR辞書データとコーパスデータの関係について

EDR辞書とコーパスは、単語や概念IDによってリンクが張られている。したがって、辞書データあるいはコーパスデータの修正はそのデータのみでなく、他のデータをも修正したことになる。また、それぞれのデータの内部においても相互に依存しあっているため1カ所の修正は他と矛盾を起こすことがある。これらの管理も行わなければならない。データ間の関係は、図-1のデータを辿ることにより参照される。それぞれのデータは、基本的な辞書情報としてのデータと高機能辞書エディタ内の情報としての修正履歴及び他のユーザとの整合性を取るためのユーザ辞書(一時的な修正部分の格納場所)とから構成されている。

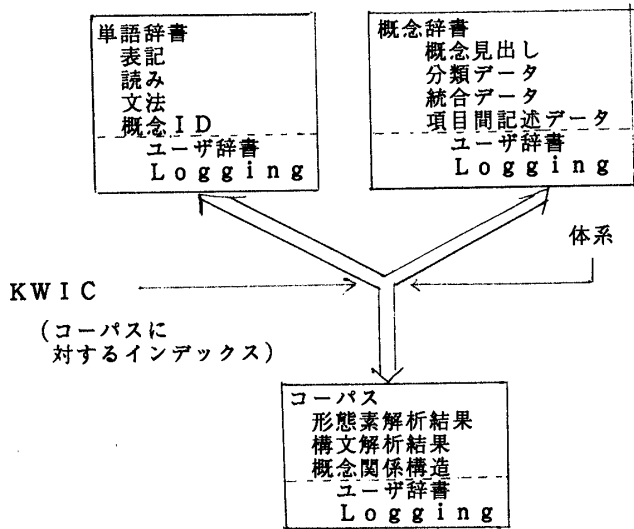


図-1 高機能辞書エディタが対象とするデータの構成図

2. EDR辞書データとコーパスデータのデータ構造について

EDR辞書データのデータ構造
 EDR辞書は、単語辞書と概念辞書の2つから構成されている。単語辞書は形態素解析、構文解析に利用できる文法情報が格納されている。また単語辞書のつのエントリーは、1つの概念に対応している。その概念にはユニークなIDと他の概念と区別できる程度の文章表現が書かれておりそれを概念見出しと呼んでいる。ただし、概念見出しは単語辞書ではなく概念辞書に記述されている。概念辞書は、概念見出しの他に体系の実データとなる分類データ、概念の発散を抑えるための統合データ、概念間の関係を記述した関係データ、概念項目間の関係を記述した項目間記述データといった概念に関係したデータが概念IDによって格納されている。それぞれのデータに対応するクラスとその参照されるデータを以下に示す。

High Quality Editor for Developing Electronic Dictionary
 Hisahiro ADACHI, Kazuhiro MIWA, Masayuki NAKAZAWA and Miho SUZUKI
 Japan Electronic Dictionary Research Institute, Ltd.

単語辞書エントリに対応するクラスWordClass、概念辞書エントリに対応するクラスConceptClass、コーパスデータのクラスとそのデータ構成以下のようにになっている。

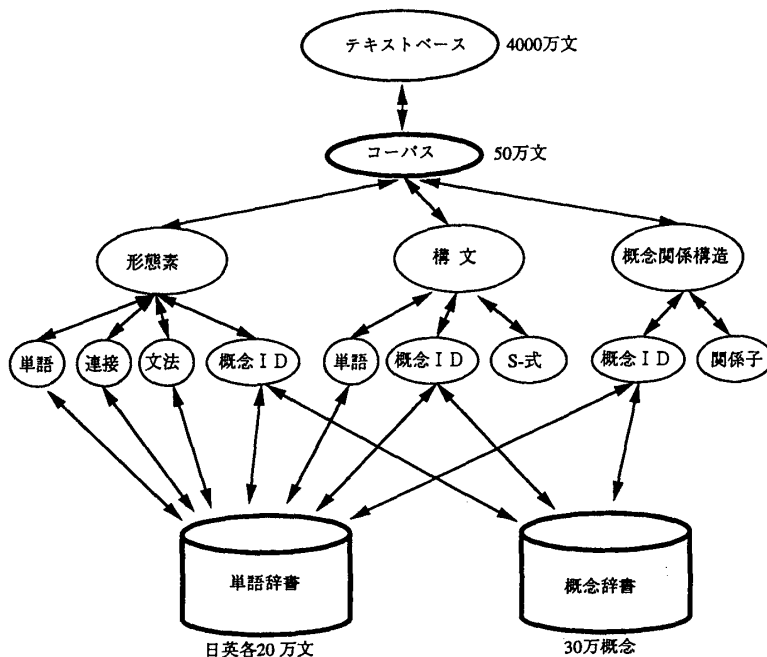
WordClass	:	単語の表記
単語表記	:	単語の読み
読み	:	形態素解析で利用するための文法属性
接続属性	:	構文解析で利用される文法属性
文法情報	:	固有の概念に付けられるID
概念ID	:	概念ID
ConceptRelClass	:	他の概念と明確に区別できる程度の自然言語表記
概念見出し	:	概念見出しに対応したID
ConceptID1	:	概念ID1と概念ID2の間に成り立つ関係
関係	:	他の概念見出しに対応したID
ConceptID2	:	概念IDリスト
ConceptCPS1Class	:	形態素解析結果
ConceptCPS2elClass	:	構文解析結果
ConceptCPS3elClass	:	係り受け構造データ
概念関係構造	:	概念関係構造

対訳リンクや分類データ及び統合データ、概念記述データはすべてこの関係として表現されている。

コーパスデータは、EDR辞書により解析され人手で検証され、その文分析結果として得られた形態素解析結果と構文解析結果、概念関係構造から構成されている。内部データとしては、単語表記、概念IDのリスト及び関係構造を持つクラスのデータとして登録される。

3. EDR辞書データとコーパスデータのアクセスについて

図-1で示したようにそれぞれのデータは、自分自身の情報とリンク情報を持っている。ユーザは常にすべてのデータを参照したいわけではない。それどころか参照したいのは部分情報の方が多い。そこで我々は、データの持つ情報を属性とし、エントリごとのデータの取出しをクラスメソッドとして定義した。これによってユーザは、自分の欲しい情報を持つインスタンスを生成しデータにアクセスすることができる。それぞれのクラス内の情報の取出しは、それぞれのクラスへのメッセージによって行われる。以下にその関係を示す。



ある概念の分類項目が知りたい場合は以下のようなクラスの組み合わせになる。

```
WordClass(単語, ConceptID);
ConceptRelClass(3, ConceptID, Label);
```

まとめ

種々の観点から辞書及び参照文章をエディティングするためのデータの構造とアクセス方法を実現できた。今後は1つのデータの変更が他にどのような影響を及ぼすかユーザが把握できるシミュレーション機能を検討中である。

謝辞

最後に、この研究の機会を与えて下さったEDRの横井所長、第四研究室吉田室長、及び終始貴重なご意見を頂いた第一研究室内田室長並びに研究員の皆様に感謝致します。

参考文献:

- [1] EDR電子化辞書: EDR TR-016 November 1989
- [2] 辞書開発支援システム: EDR TR-015 November 1989
- [3] Joan Peckham, Fred Maryanski: "Semantic Data Models" ACM Computing Surveys Vol. 20, No3