

Seep による科学ファクトデータベースシステム

5V-9

佐藤 誉夫 竹田 正幸 松尾 文碩
九州大学工学部

1. まえがき

Seep は、著者らが開発した、ルールベースとデータベースを統合管理するシステムである。Seep は、UNIX と MS-DOS のもとで動作する。九州大学大型計算機センターでは、演繹 DBMS である Adbis¹⁾ を使って、結晶構造データベースシステム XDT²⁾、遺伝子情報データベースシステム GENAS³⁾ のサービスを行ってきたが、これらを Seep で再構築する作業を行なった。Seep を使うことにより、ルールベースに基づく規則推論機構を用いたデータ解析が可能となる。

2. データベースシステムの構成

Seep を用いた科学ファクトデータベースシステムの構成を、XDT を例に説明する。図 1 に XDT の構成を示す。図のように、Seep はルールベース管理機構

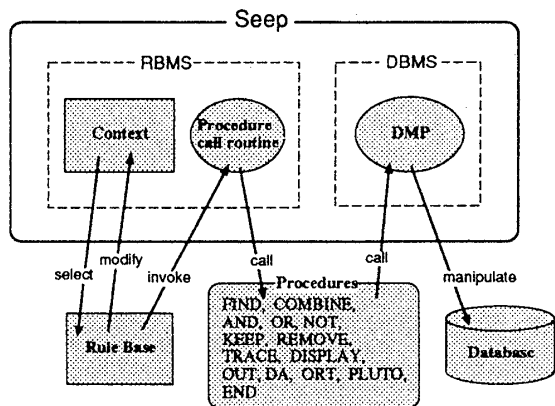


図 1 XDT の構成

(RBMS) とデータベース管理機構 (DBMS) から成っている。Seep の DBMS の外部へのインターフェースは、DMP (Data Manipulation Primitives) である。これは、find, insert, retrieve, create などのデータベースの基本操作の集合であり、C 言語で記述された手続きのライブラリである。この DMP をアプリケーションなどから呼び出すことにより、データベースを操作することができる。

通常のルールベースシステムと異なり、Seep ではルールから手続きを呼び出せるようにしている。Seep の DBMS は固有の問合せ言語をもたない。そこで、FIND, DISPLAY などの情報検索系のコマンドに対応する手続きを、DMP を用いて作成し、ルールから呼び出すことにした。また、アプリケーションに関しては、ケンブリッジ結晶データ付属の FORTRAN で書かれた解析プログラム等をもとに ORT, PLUTO などの手続きを作成した。すなわち、情報検索系のコマンドとアプリケーションを区別しない。図 2 に XDT の検索例を示す。

```

kyu-cc%xdt
.find triazole
1: 253 document(s) found
.find amino
2: 7676 document(s) found
.combine 1*2
3: 71 document(s) found
.display
3: combine 1*2
1/ 71
MF = C9 H10 N6 O1 S1
CN = 4-AMINO-3-(BETA-BENZOYLHYDRAZINO)-5-MERCAPTO-1,2,4-TRIAZOLE
BI = R.C.SECOMBE,C.H.L.KENNARD : J.CHEM.SOC.,PERKIN T., VOL. , PAGE
1973.
+
2/ 71
MF = C24 H30 N6 O2
CN = 5-AMINO-1-BUTYL-3-(N-BUTYL-N',N'-DIBENZOYL-...OLE-N$1,N)-
BI = G.RECK,M.CZUGLER,L.PARKANYI,E.SAUPE
PAGE 565, 1981.
- 2357, 1990.
+
3/ 71
MF = C10 H11 Cl2 N6 1+
CN = 4-AMINO-3-(2-...
HYDROCHI...
BI = P.-E... AMINO-3-METHYL-1,2,4-TRIAZOLE-N$1,N$2)-BIS(MU)2$-CHLORO
+
...LANFRANCHI,M.A.PELLINGHELLI : J.CHEM.RES., VOL. 214, PAGE
/12, 1990.
+
71/ 71
MF = C3 H6 N4 S1
CN = 4-AMINO-3-METHYL-1,4-DIHYDRO-1,2,4-TRIAZOLE-5-THIONE
BI = F.BIGOLI,M.LANFRANCHI,M.A.PELLINGHELLI : J.CHEM.RES., VOL. 214, PAGE
1712, 1990.
+

```

図 2 XDT の検索例

3. Seep のデータベース

Seep のデータベースは、関係の集合である。関係の特別な形として関数を区別した。例えば、関係 $ATMP$

$$ATMP \subseteq ID \times ATOM \times APN$$

は、三つの領域 $ID, ATOM, APN$ の直積の部分集合を意味する。これらの領域は関係データベースにおい

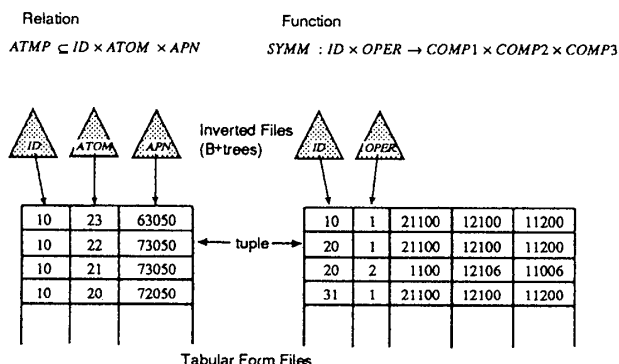


図3 関係・関数のデータファイル構成

て属性と呼ばれる。一方、関数 $SYMM$

$$SYMM : ID \times OPER \rightarrow COMP1 \times COMP2 \times COMP3$$

は、 $ID \times OPER$ から $COMP1 \times COMP2 \times COMP3$ への写像を意味する。関係および関数は、図3に示すように、表形式ファイルと索引転置ファイルとして実現される。関係ではすべての領域について索引転置形を作成するが、関数では関数値の領域に関して索引転置形を作成しない。

関係データベースのすべての関係がひとつの属性を含むとき、その属性をデータベースキーとよぶ。XDT や GENAS などの科学ファクトデータベースでは観測ごとに付与した参照番号をデータベースキーとすることができる。Seep では、データベースキーをもつデータベースについては、キー集合を単位とした検索や集合演算などの操作を行なうことができる。

キーワード検索のためには、キーワード KW から関連するデータベースキー ID の集合を求める必要がある。これは、関係

$$KY \subseteq KW \times ID$$

によって実現できるが、Seep では、時間・領域量の効率化を図るため、

$$KY : KW \rightarrow 2^{ID}$$

という特別な関数を用いることができる。この関数は、図3の表形式をもたない。

4. データベースの定義と構築

データベースを作成するためには、データ定義言語 (Data Definition Language; DDL) を用いて、データベースの定義をあたえなければならない。図4にXDTのデータ定義の一部分を示す。図において、

```

DATABASE(XDT)
ADMINISTRATOR(g70073a)
MASTER_PASSWORD(*****)
DATABASE_KEY(ID)
.....
DEF(ATMP < ID*ATOM*APN)
  L('ATOM PROPERTIES')
  PATH(/home/center/database/w_vf1/xdt/)
DEF(SYMM : ID*OPER -> COMP1*COMP2*COMP3)
  L('SYMMETRY POSITION')
  PATH(/home/center/database/w_vf1/xdt/)
.....
DEF(ID) T(L4) L('ID. CODE')
DEF(ATOM) T(I2) L('NUMBER OF ATOM')
DEF(APN) T(I4) L('10000*ELEM+1000*NCA+100*NH+NCH+50')
DEF(OPER) T(I2) L('NTH SYMMETRY OPERATOR')
.....

```

図4 XDTのデータ定義

DATABASE_KEY(ID)

は、領域 ID をデータベースキーとすることを意味する。また、

```

DEF(ATMP < ID*ATOM*APN)
DEF(SYMM : ID*OPER -> COMP1*COMP2*COMP3)

```

はそれぞれ、関係 $ATMP$ 、関数 $SYMM$ を定義している。

図4の形式によるデータ定義のファイルをもとに、DMP の define によってデータベースの定義が行なわれる。各関係・関数は、それぞれの入力形式データファイルから DMP の create によって創成が行なわれる。

5. むすび

Seep を用いた科学ファクトデータベースの構築について述べた。既存のデータ解析プログラム等は、ルールベースにより効果的にデータベースの応用プログラムとして統合することができる。今後は、ルールベースによる規則推論機構を用いた遺伝子情報の解析等を行なうことにより、Seep の規則推論機構の実用環境下での有効性を検証したいと考えている。

参考文献

- 1) 松尾, 二村, 高木: データベース統合支援システム Adbis(1)-Adbis の概要と DMP の仕様-, 九州大学大型計算機センター広報 16, 6, 172-217 (1981).
- 2) 河野, 高木, 松尾, 二村, 鬼塚: 結晶構造データベースシステム XDT の使用法, 九州大学大型計算機センター広報 16, 1, 32-53 (1983).
- 3) 久原, 榎, 高木, 松尾, 二村, 鬼塚: 核酸塩基配列データベースシステム GENAS の使用法 (1), 九州大学大型計算機センター広報 16, 5, 497-521 (1983).