

遺伝子-蛋白質解析システム GENAS について

久原 哲¹⁾, 高木利久²⁾, ニ村祥一³⁾, 神佳え⁴⁾, 林 勝哉¹⁾, 松尾文碩³⁾

1) 九州大学農学部農芸化学科, 2) 工学部情報工学科, 3) 大型計算機センター, 4) 遺伝情報実験施設

遺伝子-蛋白質解析システム GENAS は Adbis E DBMSとして EMBL Nucleotide Sequence Data Library, GenBank, NBRF Protein DataBankをデータベースとし応用プログラムの実行環境をもつデータベースシステムである。検索機能だけでなく、情報処理機能をもつことにより配列解析に有効であり、同時に個人の配列に対しても同様の解析ができることによりその応用範囲は広い。

GENAS に登録されている応用プログラムは挿入欠出を含む文字列の高速検索や長大配列の二次構造推定などがあり有効な解析ができるようになっている。

霊長類の反復配列である L1ファミリーの遺伝子、蛋白質の構造解析を行なった。その結果、相同性の高い領域が存在することが明らかとなり、同時にコードしている蛋白質がウィルスの RNA依存性 DNAポリメラーゼと相同性があった。

"GENAS: A database system for nucleic acid and protein sequence analysis" (in Japanese)

Satoru KUHARA¹⁾, Toshihisa TAKAGI²⁾, Syouchi FUTAMURA³⁾, Yoshiyuki SAKAKI⁴⁾, Katuya HAYASHI¹⁾, Fumihiko MATSUO³⁾

1) Department of Agricultural Chemistry, 2) Department of Computer Science and Communication Engineering, 3) Computer Center and

4) Research Laboratory for Genetic Information, Kyushu University, 46-02, Fukuoka 812, JAPAN

A database system named GENAS, (GENE and protein Analyzing System) for computer analysis of sequence was constructed using a deductive DBMS called Adbis which manages integrated database GENEDB and application programs. GENAS enables us to retrieve any sequence data from EMBL nucleotide sequence data library, GenBank and NBRF protein sequence databank, and readily to analyze them (if necessary, together with private data).

Homology search program in application uses Wilber and Lipman's fast algorithm which allows insertion and deletion between a private sequence and database sequence. Second structure of large RNA sequence can be estimated by dynamic programming algorithm.

L1 family of repetitive DNA sequences in primate were analyzed homology among sequences using GENAS. Members showed a high homology among each other and long open reading frame in L1 family sequence has significant sequence homology to several RNA-dependent polymerase of virus.

1 はじめに

遺伝子、蛋白質は生体機能をつかさどる重要な生体高分子であり、古くから研究が行われ多くの研究成果が発表されている。その構成要素は、遺伝子ではアデニン(A)、グアニン(G)、チミン(T)、シトシン(C)の4種であり、蛋白質では20種のアミノ酸である。1950年代に蛋白質ではポーリングによって α -ヘリックスという2次構造のモデルが、遺伝子ではワトソン-クリックによって遺伝子(DNA)の2重らせん構造のモデルが出され、構造と機能についての研究がさかんとなるきっかけを作った。特にDNAの2重らせん構造は遺伝子の複製つまり遺伝に重要な知見を与え遺伝から分子生物学への橋渡しを行なった。蛋白質については構造解析の手法が1950年代に確立され、蛋白質のアミノ酸配列が次々に決定されていた。しかしながら当時は遺伝子の構造解析の有効な手法がなく解析はあまり進まなかった。1970年代になり、サンガー法やマキサム-ギルバート法などの遺伝子構造解析手法の確立と遺伝子工学と呼ばれる遺伝子操作技術の普及とが相まって容易にかつ短時間で遺伝子の構造が決定できるようになった。これ以降は遺伝子の塩基配列の発表数は飛躍的に増加し、現在では蛋白質のアミノ酸配列として発表されるもののうち約90%が直接蛋白質から配列を決定したものでなく遺伝子の上で決定されたものという状況になっている。

この遺伝子、蛋白質の配列の急激な増加に際して、配列情報のデータベース化が提案された。蛋白質に関しては1950年代からジョージタウン大学のNational Biomedical Research Foundation (NBRF)の故Dayhoff女史らが中心となりアミノ酸配列を収集しAtlas of Protein Sequence and Structure⁽¹⁾を出版し、1971年からはProtein Data Bank (NBRF蛋白質データベース)となり計算機可読な形となっている。これに対し、遺伝子では1979年にデータベース化の動きが始まり、欧州では欧州分子生物学研究所(EMBL)、米国ではNational Institute of Health (NIH)が中心となり遺伝子データベースを作り始めた。1982年には、EMBLがEMBL Nucleotide Sequence Data Library Release 1.0を米国ではLos Alamos国立研究所とBolt Beranek and Newman (BBN)社がGenBank Release 1.0を配布し始めた。

1980年代に入り、オゾンジーンと呼ばれる発癌蛋白質をコードしている遺伝子⁽²⁾が発見されたが蛋白質の機能は不明であった。しかしデータベースとの相同性の検索により、それまでに知られていた蛋白質リノ酸化酵素、増殖因子あるいはDNA結合蛋白質などに非常に類似していることが明らかとなり発癌機構の解明に一步近づくことになった。この結果、分子生物学の分野でもデータベースが盛んに利用されるようになってきている。

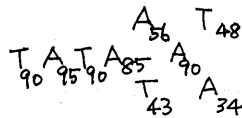
しかしながら、我国においては検索と処理機能を持つデータベースシステムが存在していなかった。そこで著者は九州大学大型計算機センターのFACOM M-382 OS IV/F4上に遺伝子-蛋白質解析システムGENAS⁽³⁾を構築した。GENASは情報検索と情報処理の2つの機能を持ち、GenBank, EMBL Nucleotide Sequence Data LibraryとNBRFのProtein Data Bankをデータとし、データベース管理システムとして九州大学大型計算機センターが開発したデータベース統合支援システムAdbis⁽⁴⁾を使用した関係型データベースシステムである。

2 遺伝子の構造と蛋白質の発現機構

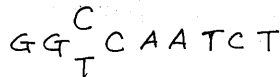
ここで遺伝子の構造と遺伝子にコードされている蛋白質の発現機構について簡単に述べる。この発現機構は図1に示すように転写、プロセッシング、翻訳の3過程に分けられ、それぞれの過程は配列と密接な関係にある。

2-1 転写過程

発現の最初はRNAポリメラーゼと呼ばれる酵素により、DNAの情報がRNAへ転写される過程である。転写が起る部位の周辺には各遺伝子間で共通性の高い2つの配列がある。1つはHogness box あるいはTATA box と呼ばれるもので転写開始点の上流30塩基付近にあり次のような配列を持っている。



ここでA,T,C,Gの下に付けた数字は出現頻度を示している⁽⁵⁾。このHogness boxのさらに約50塩基上流にCAAT box と呼ばれる



の共通な配列⁽⁶⁾があり、これらの共通配列は転写効率に影響することが知られている。これに対して転写の終了点の配列については不明であり、わずかに原核生物の転写終了の配列が明らかになっているに過ぎない。

2-2 プロセッシング過程

ポリメラーゼによる転写産物であるRNAは5'末端に7メチルグアノシン(CAP)、3'末端にはポリAが付加される。CAPが付加される位置はポリメラーゼによる転写開始に対応し、ポリAが付加される位置の10~20塩基上流にはAATAAAのポリAシグナルと呼ばれる共通配列が存在している。

真核生物では蛋白質をコードしている遺伝子はいくつかの部分に分かれて存在しているという特徴がある。蛋白質をコードしている部分をEXON(エキソン)、していない部分をINTRON(イントロン)と呼ぶ。このような構造をとっているRNAを前駆体RNAと呼び、この前駆体RNAからINTRONが切りとられEXONのみから成る成熟メッセンジャーRNA(mRNA)が出来る。この過程をスプライシングと呼ぶ。このスプライシング過程は蛋白質の配列に関係するものだけに正確さが要求され、その為にはINTRONとEXONの連結部位に何らかのシグナルが存在していると考えられ図1に示すような共通構造⁽⁷⁾が見つけられた。しかしながら

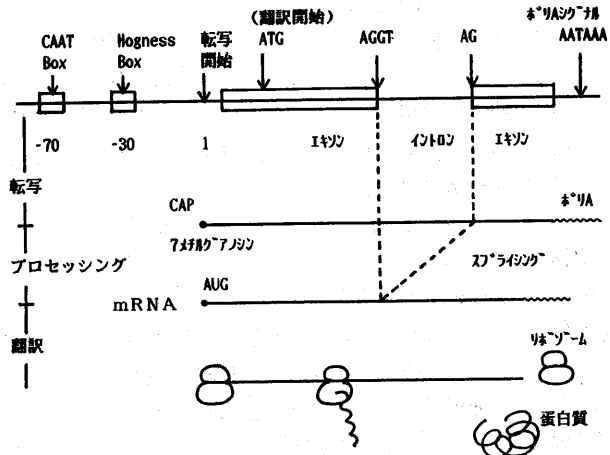


図1 蛋白質の発現機構

スプライシングの詳細な機構はいまだ不明な点が多い。

2-3 翻訳過程

プロセッシングを受けたできた mRNA は次にリボゾームと呼ばれる装置により蛋白質に翻訳される。その過程はまずメチオニンを示す AUG の mRNA 配列にリボゾームが結合し、次に表 1 に示すヌクレオチド (コドン) に従ってアミノ酸に翻訳され、アミノ酸間でペプチド結合が形成され蛋白質の一次構造が完成する。

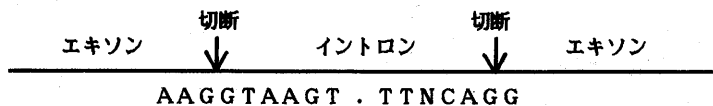


図2 スプライシング部位の共通配列

表1 遺伝コード表

First position (5'-end)	Second position				Third position (3'-end)
	U	C	A	G	
U	Phe	Ser	Tyr	Cys	U
	Phe	Ser	Tyr	Cys	C
	Leu	Ser	Stop	Stop	A
	Leu	Ser	Stop	Trp	G
C	Leu	Pro	His	Arg	U
	Leu	Pro	His	Arg	C
	Leu	Pro	Gln	Arg	A
	Leu	Pro	Gln	Arg	G
A	Ile	Thr	Asn	Ser	U
	Ile	Thr	Asn	Ser	C
	Ile	Thr	Lys	Arg	A
	Met	Thr	Lys	Arg	G
G	Val	Ala	Asp	Gly	U
	Val	Ala	Asp	Gly	C
	Val	Ala	Glu	Gly	A
	Val	Ala	Glu	Gly	G

3 遺伝子-蛋白質解析システム GENAS

3-1 データベース

遺伝子の塩基配列データベースとしては欧州では EMBL がデータ収集配布を行なっている EMBL Nucleotide Sequence Data Library と米国では NIH により Los Alamos 国立研究所がデータ収集を行ない BBN 社がデータ配布を行なっている GenBank, および NBRF がデータ収集配布を行なっている NBRF Nucleotide Data Bank の3つがある。これらのデータベースは1982年から公共データベースとして無料で磁気テープあるいはフロッピーディスクによる配布が開始されており、GenBank, NBRF は年々回また EMBL は年2回更新されている。

蛋白質のアミノ酸配列のデータベースとしては NBRF の Dayhoff らが Atlas of Protein Sequence and Structure を出版しており、1983年からは公共データベース NBRF Protein Data Bank として無料で配布されている。これらのデータベースの増加は塩基配列に関しては年間150万塩基以上になっている。

これら3つのデータベースはそれぞれ異なる書式を用いて書かれており、その例を図3に示す。EMBL を例にとり構造を説明すると1つの遺伝子の情報は1つのエントリーとして格納されている。各エントリーはまず ID 行で遺伝子を表わすコード名, AC 行でアセション番号, DT 行で入力された日付, DE 行で遺伝子の定義, KW 行でキーワード, OS 行で生物種, OC 行で分類, RN 行で文献番号, RT 行で文献題目, RA 行で著者名, RL 行で文献名, CC 行で注釈, FH 行で feature table の見出し, FT 行で feature table, SQ 行で配列の見出し, その次に配列, 最後に // 行で終了する順になっている。特に分子生物学の研究者が必要とする情報は FT 行に示されている feature table であり、各情報の位置とコードが key, From, To, Description の順で示されている。

著者らのシステムでは、これらのデータベースの中で塩基配列として EMBL, GenBank および蛋白質のアミノ酸配列として NBRF Protein Data Bank を基にして新たに関係型データベース GENEDB を構築している。GENAS を構築する上


```

DATABASE (GENEDB)
  L('EMBL + GENBANK + PROTEIN DATABASE')
*
DEF (REF < RF#ID)
  L('REFERENCE CODE')
*
DEF (RNA < RF)
  L('RNA OR CDNA')
*
DEF (FR < RF)
  L('FRAGMENT')
*
DEF (DE : RF -> DEF)
  L('DEFINITION')
*
DEF (OS < RF#GEN#SPEC)
  L('ORGANISM')
*
DEF (OC < RF#LVL#TAX)
  L('ORGANISM CLASSIFICATION')
*
DEF (HS < RF#GEN#SPEC)
  L('THE MOST COMMON HOST SPECIES')
*
DEF (HC < RF#LVL#TAX)
  L('HOST CLASSIFICATION')
*
DEF (BIB : RF#SN -> BI)
  L('BIBLIOGRAPHIC INFORMATION')
*
DEF (KY < RF#KH)
  L('KEYWORD')
*
DEF (AUTH < RF#AU)
  L('AUTHOR')
*
DEF (DC < RF#DS)
  L('DATA COMPILE')
*
DEF (YEAR < RF#YR)
  L('YEAR')
*
DEF (FN : RF -> NFTE)
  L('NUMBER OF FEATURE TABLE ENTRIES')
*
DEF (FT : RF#FTEN -> KEY#FROM#TO#DESCR)
  L('FEATURE TABLE')
  I(KEY, FROM, TO)
*
DEF (SH : RF -> APC#TG#MU#Y#R#AB)
  L('SEQUENCE HEADER')
*
DEF (LN : RF -> LEN)
  L('LENGTH OF SEQUENCE')
*
DEF (SQ : RF#SN -> SEQ)
  L('SEQUENCE')
*
DEF (RF)
  KEY
  T(14) L('REFERENCE CODE')
*
*** PROTEIN IF RF > 1,000,000,000; OTHERWISE NUCLEOTIDE.
*
DEF (ID)
  T(C10) L('SHORT IDENTIFICATION')
*
DEF (GEN)
  T(C100) L('GENUS')
DEF (SPEC)
  T(C50) L('SPECIES')
*
DEF (LVL)
  T(12) L('LEVEL IN TAXONOMIC TREE')
DEF (TAX)
  T(C50) L('ENTRY IN TAXONOMIC CLASSIFICATION')
*
DEF (DEF)
  T(C600) L('DEFINITION WRITTEN IN ENGLISH')
DEF (BI)
  T(C200) L('BIBLIOGRAPHIC INFORMATION')
DEF (KH)
  T(C50) L('KEY WORD')
DEF (AU)
  T(C50) L('AUTHOR')
DEF (DS)
  T(C8) L('DATA SOURCE')
DEF (YR)
  T(12) L('YEAR')
*
DEF (NFTE)
  T(12) L('NUMBER OF FEATURE TABLE ENTRIES')
DEF (FTEN)
  T(12) L('FEATURE TABLE ENTRY NUMBER')
DEF (KEY)
  T(C10) L('FEATURE TABLE KEY NAME')
DEF (FROM)
  T(14) L('ENDPOINT SPECIFICATION')
DEF (TO)
  T(14) L('ENDPOINT SPECIFICATION')
DEF (DESCR)
  T(C250) L('MORE INFORMATION ABOUT FEATURE')
*
DEF (A)
  T(14) L('NUMBER OF ADENINE BASES')
DEF (C)
  T(14) L('NUMBER OF CYTOSINE BASES')
DEF (T)
  T(14) L('NUMBER OF THYMINE BASES')
DEF (G)
  T(14) L('NUMBER OF GUANINES')
DEF (U)
  T(14) L('NUMBER OF URACIL BASES')
DEF (Y)
  T(14) L('NUMBER OF PYRIMIDINE BASES')
DEF (R)
  T(14) L('NUMBER OF PURINE BASES')
DEF (AB)
  T(14) L('NUMBER OF ANY BASES')
*
DEF (LEN)
  T(14) L('LENGTH OF SEQUENCE')
*
DEF (SN)
  T(14) L('SEQ NUMBER')
DEF (SEQ)
  T(C500) L('NUCLEOTIDE/PROTEIN SEQUENCE')

```

図4 データ定義

た。最も重要である FT (feature table) については NFTE (feature table のエントリ数) を作り RF X FTEN の直積ドメインとして KEY X FROM X TO X DESCR の直積をコードインとする関係とした。最後に3つのデータを区別するため DC (データコンパイル) の関係を作った。これらの関係の中で載置ファイルを作るものは FT のみとした。この定義に従って3つのデータより関係型データベース GENEDB を作った。

3-2 データベースシステム

前述の GENEDB をデータベースとして GENAS を作る上で次の2点に重点をおいた。

①情報の検索機能を持つ。担し検索手法は簡単であることが望ましい。

②情報の処理機能を持つ。特にデータが一次構造(配列)であり、共通配列とはいうものの変異があり配列の情報を得るには単純な文字列の検索だけでは不十分であり種々の処理を行なう必要がある。この2点を検討し GENAS のモードを図5に示すように作った。

3-2-1 情報検索

まず情報の検索機能については、分子生物学の分野の研究者は計算機を使った経験があまりないので、Adbis の推薦機能を直接利用する検索より、Adbis の情報検索型コマンドを使用することを考えた。次に検索を考えると、その検索のほと

んどがキーワードによる検索でありその他の検索、例えば詳細な分類による検索などはほとんど行なわないと考えられるので GENAS では情報検索型コマンドで検索可能な領域は KW(キーワード)、AU(著者名)、ID(コード名)とし、キーワードに重点をおくことにした。

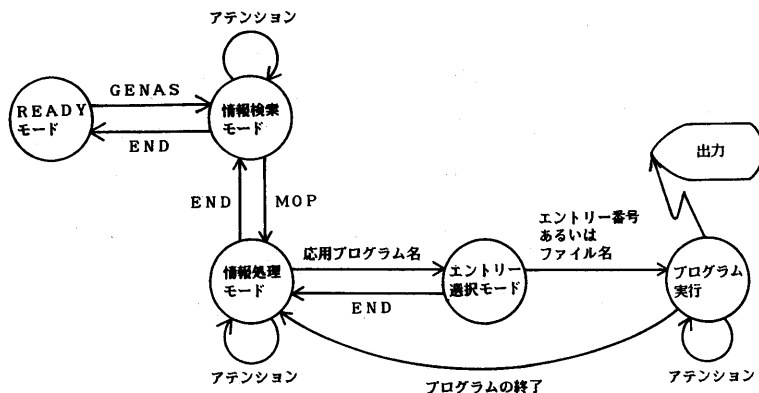


図5 GENASのモード遷移図

基となるデータのキーワードでは不足なので前述の DE 行、KW 行、RT 行、OS 行から単語を切り出しすべてキーワードとした。

3-2-2 情報処理機能

遺伝子あるいは蛋白質解析のデータベースシステムでは一般の学術情報データベースとは異なり、前述したように構造と機能の関係が一意的に決定されておらず不確定な要素が多いため、定形の処理だけではなく多様な処理ができる環境を持たなければならない。

アプリケーションとのインタフェース： MOP と呼ばれるコマンドプロシジャーを経由してアプリケーションを起動している。ただしアプリケーションの多様性および各アプリケーションのデータ要求の多様性から MOP の中でアプリケーションの選択を行ない、必要なデータあるいは環境の設定を行なっている。

アプリケーション： 現在 GENAS に登録されているアプリケーションの一覧を表すに示した。アプリケーションの実行環境では特に個人が所有している配列に対してもデータベース中の配列と同様な処理ができるようにした。処理のダイアグラムを図6に示した。

分子生物学の分野では、先に紹介したようにオンコジーンの見出し以来、実験で得た配列についてはデータベースの配列と相同性を検索することが一般的になってきている。この相同性の検索では文字列の完全一致ではなく文字列の不一致あるいは挿入欠失を含んだ一致を探さなければならない。不完全一致の相同性のプログラムについては Wilber, Lipman らが開発したアルゴリズム⁽⁸⁾がよく用いられている。遺伝子の2次構造の推定にはダイナミックプログラミングの手法が用

Adbisのコマンド呼び出し機能を使用

表2 GENASのアプリケーション

統計処理
塩基の出現頻度 コドンの出現頻度 類似性の統計処理
一次構造の文字列解析
相同性の検索 遺伝子からアミノ酸への翻訳 アミノ酸から遺伝子への逆翻訳 任意の文字列の検索
二次構造の予測
蛋白質の二次構造の予測 遺伝子の二次構造の予測

いられ計算時間の短縮を行なっているものがある(9)。

3-3 実際例(霊長類に存在する反復配列L1ファミリーの一次構造の比較)(10)

著者らのグループは高等動物の遺伝子の特徴の一つである反復配列の一つであるL1ファミリーの6個(L1Nc1~L1Nc6)の塩基配列を決定しGENASによる一次構造の解析を行なった。その結果図7に示すようにL1ファミリーの間では90%以上の相同性があり、その中のL1Nc4に

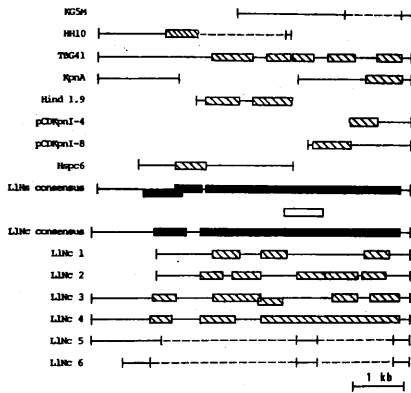
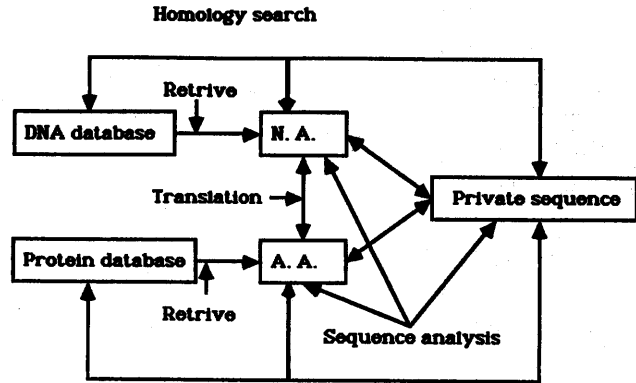


図7 L1ファミリーの構造と共通配列



Homology search

図6 処理のダイアグラム

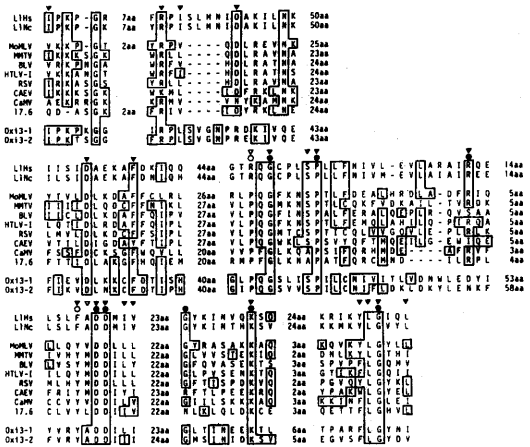


図8 L1のアミノ酸配列とウイルスのRNA依存性DNAポリメラーゼの相同性

は927アミノ酸残基から成る蛋白質がコードされている領域があることが明らかとなった。特に蛋白質をコードしている領域については他の種のL1ファミリーと65%以上の相同性があった。蛋白質をコードしている領域の相同性からみてL1ファミリーは進化的に保存されている蛋白質から誘導されてきたものであると考えられる。

次に、この保存されている蛋白質について検討した結果、図8に示すように、ウイルスのRNA依存性DNAポリメラーゼと部分的に強い相同性があることが明らかとなった。もし、コードしている蛋白質に逆転写活性があればL1ファミリーの分散の原因を考える上で重要な知見となるものである。

参考文献

- 1) Dayhoff, M.O., et al. (ed.) : Atlas of protein sequence and structure, Vol.1, National Biomedical Research Foundation, Silver Spring, Md. (1965)
- 2) Czernilofsky, A., et al. : Corrections to the nucleotide sequence of src gene of Rous sarcoma virus, Nature, Vol.301, pp. 736-738 (1981)
- 3) Kuhara, S., et al. : GENAS: A database System for Nucleic Acid Sequence Analysis, Nucleic Acids Res., Vol.12, pp.89-99 (1984)
- 4) 松尾文碩, 二村祥一, 高木利久 : 推論関係型データベース管理システム Adbis, 情報処理学論文誌, Vol.124, pp.249-255 (1983)
- 5) Lewin, B. : "GENES", Jhon Wiley & Sons, New York (1983)
- 6) Proudfoot, N.J., et al. : Structure of the Genes that do not Rearrange, Science, Vol.209, pp.1329-1336 (1980)
- 7) Lewin, B. : Alternatives for splicing: Recognizing the End of Introns, Cell, Vol.22, pp.324-326 (1980)
- 8) Wilber, W.J. and Lipman, O.J. : Rapid Similarity Searches of Nucleic acid and Protein Data Bank, Proc. Natl. Acad. Sci. U.S.A., Vol.80, pp.726-730 (1983)
- 9) Zuker, M. and Stiegler, P. : Optimal Computer Folding of Large RNA Sequence using Thermodynamics and Auxiliary Information, Nucleic Acids Res., Vol.9, pp.133-148 (1981)
- 10) Hattori, M. et al. : L1 family of repetitive DNA sequences in primates may be derived from a sequence encoding a reverse transcriptase related protein, Nature, Vol.321, pp.625-628 (1986)