



NCBI

田中 俊典

I. はじめに

NCBI とは米国の国立衛生研究所 (NIH: National Institutes of Health) の中にある国立医学図書館 (NLM ; National Library of Medicine) の一部門で、国立生物工学情報センター (National Center for Biotechnology Information) のことです。1988 年の設立以後、生命科学、とくに分子生物学や生物情報科学 (bioinformatics) の研究に資するために、種々のデータベースの構築やそれらを利用するためのソフトウェアの開発をしています。これらのリソースの使用は基本的に無償で、全世界の研究者が利用しており、今や NCBI が無ければ全く研究は進まないと言っても過言ではないでしょう。

NCBI にあるたくさんの機能の中でも文献データベースである PubMed は医学生物学系のデータベースとして最も有名です。PubMed の元のデータはほとんどが MEDLINE (NLM のデータベース) 由来なのですが、MEDLINE がカバーしていないデータ、例えば full text 文献へのリンクや古い文献などが含まれています。したがって PubMed は MEDLINE のスーパーセットと行うことができるでしょう¹⁾。

拙稿ではたくさんある NCBI の活動の中から、PubMed 以外のいくつかを、あまり細かい所には立ち入れませんが、筆者の能力及ぶ範囲で紹介していきたいと思えます。実は NCBI にはチュートリアルやヘルプといったリソースも充実しており²⁾、この辺の話はそれらを見ていた

だいた方が間違いがないので、もし興味がおありの方はぜひ参照していただければと思います。

II. 概観

表 1 は The NCBI Handbook³⁾ の目次の項目部分だけを抽出したものです。Part. 1 のデータベースを見ていただくと、11 の項目が並んでいます。

例えば Part. 1 の最初に GenBank がありますが、これは核酸配列、すなわち DNA や RNA のデータベース (アミノ酸の配列情報も含まれます) で、ヨーロッパ (EMBL; the European Molecular Biology Laboratory) や日本 (DDBJ; the DNA Data Bank of Japan) との協力です。実際は、この 11 項目以外にも多くのデータベースがあります。

Part. 3 にはこれらのデータベースを用いるために有用なツールなどが書かれています。この中で Entrez というのは複数のデータベースに渡って検索できる非常に強力なサーチエンジンです。図 1 は NCBI の最初のページにある検索ボックスに Entrez で prion という語を入れ、「all databases」を選択して、検索した画面です (prion というのは蛋白質の一種で、ウシの狂牛病やヒトのクロイツフェルト・ヤコブ病に関係することで知られています)。

たくさんのアイコンがありますが、それらがすべてデータベースを表しており、これらのすべてのデータベースから prion という語を検索した結果が表示されています。左の一番上には PubMed がありますね。アイコンの左にある 13484 という数字は PubMed を検索して引っか

たなか としふみ: 藍野大学医療保健学部臨床工学科教授
t-tanaka@me-u.aino.ac.jp

表1 The NCBI Handbook の目次抜粋

Part 1. The Databases

1. GenBank: The Nucleotide Sequence Database
2. PubMed: The Bibliographic Database
3. Macromolecular Structure Databases
4. The Taxonomy Project
5. The Single Nucleotide Polymorphism Database (dbSNP) of Nucleotide Sequence Variation
6. The Gene Expression Omnibus (GEO): A Gene Expression and Hybridization Repository
7. Online Mendelian Inheritance in Man (OMIM): A Directory of Human Genes and Genetic Disorders
8. The NCBI BookShelf: Searchable Biomedical Books
9. PubMed Central (PMC): An Archive for Literature from Life Sciences Journals
10. The SKY/CGH Database for Spectral Karyotyping and Comparative Genomic Hybridization Data
11. The Major Histocompatibility Complex Database, dbMHC

Part 2. Data Flow and Processing

12. Sequin: A Sequence Submission and Editing Tool
13. The Processing of Biological Sequence Data at NCBI
14. Genome Assembly and Annotation Process

Part 3. Querying and Linking the Data

15. The Entrez Search and Retrieval System
16. The BLAST Sequence Analysis Tool
17. LinkOut: Linking to External Resources from Entrez Databases
18. The Reference Sequence (RefSeq) Database
19. Gene: A Directory of Genes
20. Using the Map Viewer to Explore Genomes
21. UniGene: A Unified View of the Transcriptome
22. The Clusters of Orthologous Groups (COGs) Database: Phylogenetic Classification of Proteins from Complete Genomes

Part 4. User Support

23. User Services: Helping You Find Your Way
24. Exercises: Using Map Viewer

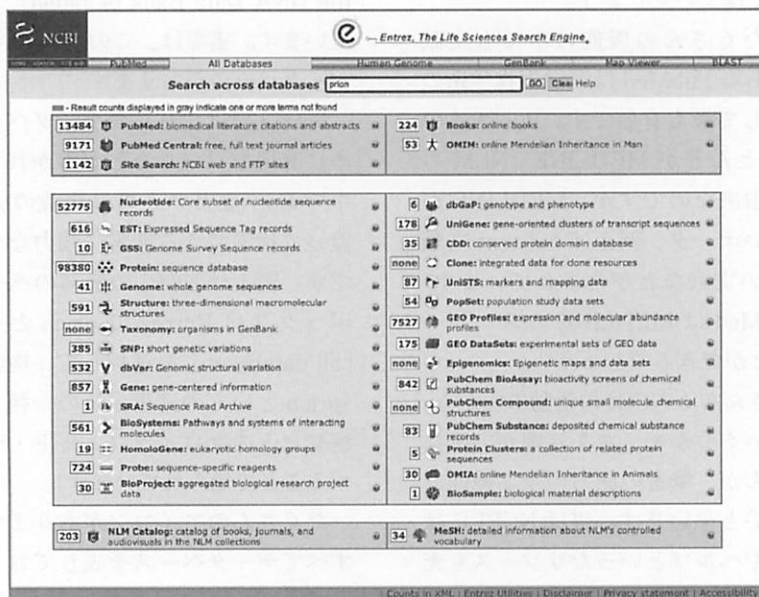


図1 Entrez で prion を検索した画面

かってきた件数を示しています(2012年9月20日現在の件数です)。ここをクリックすればその文献すべてにアクセスすることができるという訳です。このように Cross-database search ができるというのも NCBI データベースの大きな特徴の一つと言えるかもしれません。また、表1にはありませんが、比較的新しいサービスとして、MyNCBI といって検索利用時に利用者が自分で使用方法に合わせてカスタマイズできるツールなどもあり、どんどん便利に進化しているようです。

Ⅲ. 分子生物学の予備知識

データベースの解説をする前に、DNA や RNA といった用語について確認をしておきましょう。少しややこしいですがおつきあいください。

地球上のすべての生命はその遺伝情報を DNA あるいは RNA という高分子の形にして保持しています。DNA すなわちデオキシリボ核酸 (RNA はリボ核酸) はヌクレオチドという分子が鎖状につながった長い紐状の分子です。ヌクレオチドには5種類あり、記号で A,G,C,T,U と表されます。DNA はこのうち A,G,C,T の4種類、RNA は A,G,C,U の4種類のヌクレオチドから出来上がっています。この A,G,C,T (あるいは A,G,C,U) の並び方が、つまるところ遺伝情報という訳です。たとえば

・・・AAGCTTTATAGGGCGA・・・

といった具合に延々と繋がっているのが DNA (RNA) なのです。この並びを配列(シーケンス)と呼びます。

さてわれわれヒトのような真核細胞生物は、通常はこの長い DNA をいくつかの部分に切り分けて、くるくると巻き取って細胞の核の中にしまい込んでいます。この一塊を染色体と言います。ヒトの場合は46本の染色体を持っていますが、実はこのうち半分は母親から、残り半分は父親から受け継いだもので、ほとんど同じものなのです(全く同じという訳ではもちろんあ

りません)。つまりほとんど同じ DNA を2セット持っているということです。この1セットを「ゲノム」と呼びます。

ところで生物は DNA (言い換えると遺伝情報) を使って何をしているのでしょうか。実はその情報から蛋白質を作っているのです。ヒトの細胞を例にとると、まずゲノムの DNA のシーケンスを、RNA に写し取ります(この操作を「転写:transcript」と言い、できた RNA をメッセンジャー RNA と呼びます)。次にメッセンジャー RNA のシーケンスを元にして、今度はアミノ酸を並べてつないでいきます(これを「翻訳:translation」と言います)。詳細は省略しますが、3つのヌクレオチドに対して1つのアミノ酸が対応するようになっています。アミノ酸がつながった分子ということは、すなわち蛋白質、という訳です。比喩的に言うと、DNA の上に蛋白質の作り方が書いてある、ということになります。注意しておきたいのは、長いゲノムの DNA のうちで蛋白質の作り方が書かれているのは、ほんのごく一部の場所であり、ほとんどの場所は何のためにあるのかよくわかっていません。この蛋白質の作り方が書かれているゲノム DNA の部分を「遺伝子:gene」と言うのです。さらにややこしいことに、遺伝子は飛び飛びの状態で書かれていることが多いのです。そのため RNA に転写された後に、飛び飛び部分をくっつけて一つにする作業が必要です(これをスプライシングと言います)。

それでは細かい話はこの辺にして、Entrez を使って実際にいろいろなデータベースにアクセスしてみましょう。

Ⅳ. 核酸配列の検索

上でも触れましたが、NCBI における核酸の情報は GenBank というデータベースにまとめられています。GenBank には10万を越える生物種のデータが登録されています。そこにはゲノムのシーケンスもあれば、RNA のデータ、さらにアミノ酸の配列まで含まれています。

さきほどの Entrez の画面 (図 1) を見てください。ここから prion 遺伝子 (遺伝子名は PRNP です) を含む、ヒトのゲノムのシーケンスを探してみましょう。左の上から 4 つ目に nucleotide というデータベースがありますが、これを利用します。しかしこのままではヒット数は 52,777 件と多いので、「ゲノムの DNA」、「ヒトのデータ」条件を絞ってみます。そうすると 41 件が抽出されました。その中から説明文を手がかりにして、目的のデータを選び出した結果が図 2 です (ごく一部です)。

図の a, t, g, c の 4 つのアルファベットがずらっと並んでいる部分が DNA の配列、すなわちゲノムの遺伝情報を表現しています。chromosome="20" と表示されているのは、20 番染色体に属するゲノム DNA だという意味です。これで PRNP が、20 番染色体のゲノムの中

にあることはわかりましたが、実際にこのゲノムの配列のどの部分が PRNP の配列なのかはわかりづらいですね。そこで、ゲノムのデータではなく、遺伝子の配列そのもの (上で述べたように飛び飛びになっている情報も含めて) を知りたいときは、図 1 の左下にある、gene というデータベースを使えば簡単です。結果は省略しますが、これによって PRNP の DNA 配列を知ることができます。

V. 蛋白質の検索

IV. では核酸の配列を調べましたが、NCBI には蛋白質の配列 (この場合はアミノ酸配列ということになります) のデータベースもあります。ついでにこれも見ておきましょう。図 1 の protein というデータベースからヒトの prion 蛋白質の配列データを選び出した結果が図 3 です

```

FEATURES             Location/Qualifiers
     source            1..11079
                     /organism="Homo sapiens"
                     /mol_type="genomic DNA"
                     /db_xref="taxon:9606"
                     /chromosome="20"
ORIGIN
1  aatgacagtt  tgggtgaagc  agtatgtcca  tgtaaacca  cttttccaaa  ctttgcaaat
61  ccacagtggt  gaaatgtgat  gcagttatgg  tcaatgatat  taagcataat  atttttggtg
121  gggcatcttg  taagaacctt  aaaagagagg  aaaaaagtag  ttaagaggca  ttatccatta
181  gcctttgcct  ttctgctttc  ttctgcatga  agcatgggca  aaatgttga  gtcaatgcag
241  ccactttgca  aacattcagc  aaaagcacag  tctaggaatg  tattttaaaa  agactagaaa
301  ttgcctcata  cttgatatca  ttgtgaagct  attatagcat  tccttgattt  tatgttatat
361  aagaataata  aaccctaata  ttataaaagg  caacattttt  catgttttat  gtttcagcct
421  attgtacttg  tgaatagatt  atgtattata  gatcaaatgt  gatttagaaa  tatagccctt
481  tctccaatct  atcattgatg  ggcatttagg  ttgattctgt  cttttatatt  ggggaatagt
541  ctgcaatgaa  catacgcgtg  catagaggat  caggaaagat  aaacaatggg  taccaggctt
.....

```

図 2 nucleotide データベースの検索結果 (一部)

```

CDS             1..253
                /gene="PRNP"
                /coded_by="AB300823.1:11..772"
                /note="PrP 129M"
ORIGIN
1  manlgcwmiv  lfvatwsdlg  lckkrpkpgg  wntggsrypg  qgspggnryp  pqggggwgqp
61  hgggwgqphg  ggwgqphggg  wgqphggggg  qgggthsqwn  kpskpktnmk  hmagaaaaga
121  vvgglggyml  gsamsrpiih  fgsdyedryy  renmhrypnq  vyyrpmdeys  nqnnfvhdvc
181  nitikqhtvt  tttkgenfte  tdvkmervv  eqmcitqyer  esqayyqrqs  smvlfspppv
241  illisflifl  ivg
//

```

図 3 protein データベースの検索結果 (一部)

(一部です)。

下半分にアルファベットの暗号のような文字列がありますが、この文字一つ一つがアミノ酸を表しています。先ほどの核酸のように4種類ではなく、アミノ酸は20種類あるので、このような文字列になっているのです。

蛋白質は生物の体の中でいろいろな機能を持つのですが、そこで重要なのは立体構造です。立体構造は基本的にはアミノ酸の配列で決まります。つまり、例えば2つの蛋白質の一部分に共通のアミノ酸配列があると、少なくともその部分はよく似た立体構造となり、機能もよく似たものになる(だろう)と考えられるのです。このような共通部分をたくさん集めてデータベース化したものがCDD; conserved protein domain databaseです。図1の右の真ん中より少し上あたりにあります。このデータベースを利用すれば、機能はよくわからないけれどもアミノ酸配列はわかっている蛋白質の立体構造を推定することができ、さらには機能まで推測することができます。逆にある種のアミノ酸配列をもった蛋白質を人工的に作れば、ある特定の機能をもたせることができる(かもしれない)という利用の仕方も可能です。これは薬を設計したりするときに非常に有用です。

VI. BLAST

今まで見てきた検索というのは、prionという物質の遺伝子の核酸配列やアミノ酸配列を知るための方法でした。しかし逆に、例えば手元に何かわからない核酸の配列データがあるのだけれど、これが一体何の配列なのかわからないといった場合はどうすればよいのでしょうか。あるいは、ある詳細不明な蛋白質を分析して、アミノ酸配列は知ることができたけれども、この蛋白質はどういった蛋白質なのか、それとも全く新しい蛋白質なのかを知りたい。そういうときのためにNCBIではBLASTというツールを用意しています。BLASTとはBasic Local Alignment Search Toolのことです。

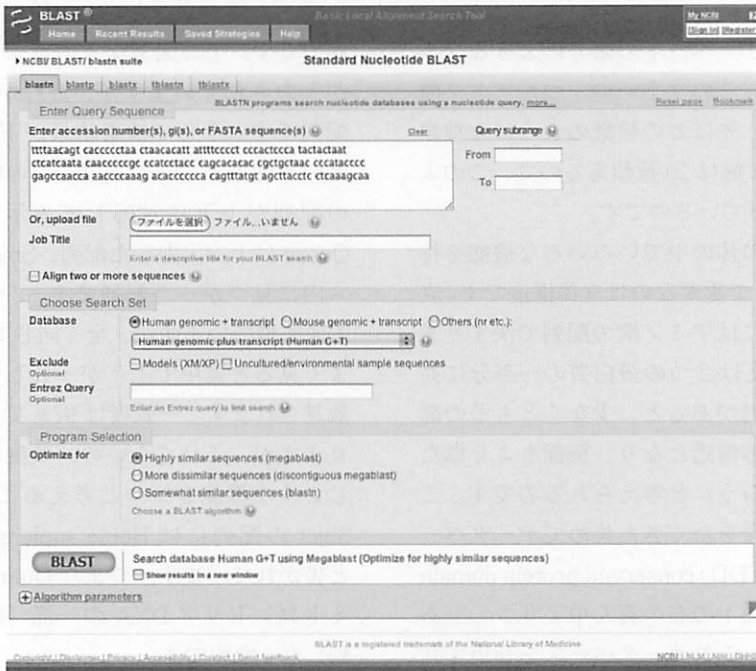
図4を見てください。上段がBLASTの検索画面です。上の大きな入力ボックスに核酸配列が入力されていますが、これが調べたいDNA配列です。そして検索結果が下段です。Queryという配列とSbjct (Subjectの略)という2つの配列が上下に対応して表示されています。Queryは上で入力した配列、Sbjctがデータベース内に見つかった配列です。2つの配列はとてもよく似ていますが、全く同じではありません。よく見ると途中で上下が一致していない場所が散見されます。一致度は97%です。ものにもよりますが、これぐらいの一致度ならば、ほぼ同じような機能を持つと考えることもできます。Sbjctの配列にはHomo sapiens mitochondrionと書かれています。つまりQuery配列はヒトのミトコンドリアDNAの一部(にとっても近い配列)だったのです。

実はこの例で用いたQuery配列もGenBankからとってきた配列でした。何の配列かと言いますと、Homo sapiens neanderthalensisつまりネアンデルタール人のミトコンドリアゲノムだったのです。そのため現代人の配列との間には微妙な違いがあったわけです。それにしても化石から採取されたデータまであるなんて、とても面白いですね。ちなみに、こういったデータの蓄積によって、ネアンデルタール人はわれわれ人類の直接祖先ではなく、同時代に生息していた別の種だったと考えられるようになったということです。

VII. 病気との関連を調べる

最後に遺伝疾患に関するデータベースを紹介します。OMIM: Online Mendelian Inheritance in Manというもので、図1の右上の方にあります。ここをクリックするとprionに関連した疾患についての詳細な情報を得ることができます(図5)。

図5の下の方にはCreutzfeldt-Jakob diseaseやGerstmann-Sträussler diseaseといったprion関連の疾患が見えます。病気そのものの臨床病



```
>ref|NC_012920.1| Homo sapiens mitochondrion, complete genome
Length=16569

Score = 303 bits (164), Expect = 4e-80
Identities = 178/184 (97%), Gaps = 4/184 (2%)
Strand=Plus/Plus

Query 1  TTTTAAAGTCACCCCTAACTAACACATTATTTCCCTCCCACTCCCATACTACTAAT 60
Sbjct 421 TTTTAAAGTCACCCCTAACTAACACATTATTTCCCTCCCACTCCCATACTACTAAT 60

Query 61  CTCATCAATAACAACCCCGCCATCCTACCCAG----CACACCCGCTGTAACCCATA 116
Sbjct 481 CTCATCAATAACAACCCCGCCATCCTACCCAGCACACACCCGCTGTAACCCATA 116

Query 117 CCCCCGACCAACCAACCCCAAGACACCCCCACAGTTTATGTAGCTTACCTCTCAA 176
Sbjct 541 CCCCCAACCAACCAACCCCAAGACACCCCCACAGTTTATGTAGCTTACCTCTCAA 176

Query 177 GCAA 180
Sbjct 601 GCAA 604
```

図4 BLASTによる検索

像など詳細が知りたければ、ここから簡単に参照することができますし、解説の文書中の引用文献はPubMedと、配列情報はGenBankとリンクしていますので効率よく知識を得ることができます。

VIII. おわりに

以上、たくさんあるNCBIの機能の中から、いくつか紹介させていただきましたが、これはほんの入り口に過ぎません。それぞれのデータベースを本当に使いこなすためには専門的知識

と専門的ニーズを持っていないと難しいと思います。とはいえ、データベースは誰にでも無料で利用できますし、仕事ではなくぶらりとNCBIを訪れて、気の向くままに遊んでみるというのなかなか楽しいものです。そんな楽しい世界の一端を拙稿で感じていただければ幸いです。

Home | About | Statistics | Downloads/API | Help | External Links | Terms of Use | Contact Us Select Language: ▼

prion Search Sort by: Relevance Date updated

Advanced Search: OMIM, Clinical Synopses, OMIM Gene Map Toggle search terms highlighted, [changes highlighted Search History: View, Clear

*176640

PRION PROTEIN; PRNP

Alternative titles; symbols

PRP

PRION-RELATED PROTEIN; PRIP

HGNC Approved Gene Symbol: PRNP

Cytogenetic location: 20p13

Genomic coordinates (GRCh37): 20-4,666,796 - 4,682,233 View NCBI

- Table of Contents - *176640
- External Links:
- Genome
- DNA
- Protein
- Gene Info
- Clinical Resources
- Variation
- Animal Models
- Cellular Pathways

Gene Phenotype Relationships

Location	Phenotype	Phenotype MIM number
20p13	Creutzfeldt-Jakob disease	123460
	Gerstmann-Sträussler disease	137440
	Huntington disease-like 1	603218
	Insomnia, fatal familial	603072
	Prion disease with protracted course (Kuru, susceptibility to)	606468 243300

図 5 OMIM (一部)

参考文献

- 1) PubMed と MEDLINE の違いについて. [引用 2012-09-20]
http://www.nlm.nih.gov/pubs/factsheets/dif_med_pub.html
- 2) NCBI Mini Courses. [引用 2012-09-20]
<http://www.ncbi.nlm.nih.gov/Class/minicourses/>

- 日本語訳を科学技術振興機構が公開している;
<http://www-bird.jst.go.jp/minicourses/>
- 3) The NCBI Handbook. [引用 2012-09-13]
<http://www.ncbi.nlm.nih.gov/books/NBK21101/>