




## Banco de Dados aplicado a Sistemas Biológicos

Saulo França Amui  
 (sauloamui@gmail.com)

Departamento de Genética – FMRP/USP  
 Grupo de Bioinformática






## Dados Biológicos

- A internet mudou a maneira como os cientistas compartilham os dados e possibilitou que um depósito central de informações atendessem totalmente a uma comunidade de pesquisa.





Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui






## Dados Biológicos

- A tendência é armazenar dados biológicos brutos de todos os tipos em bancos de dados públicos, com acesso aberto pela comunidade de pesquisa.




Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

## Dados Biológicos

- Em vez de fazer pesquisas preliminar no laboratório, os cientistas vão primeiro aos bancos de dados

→ economia de tempo e recursos



Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui






## Dados Biológicos

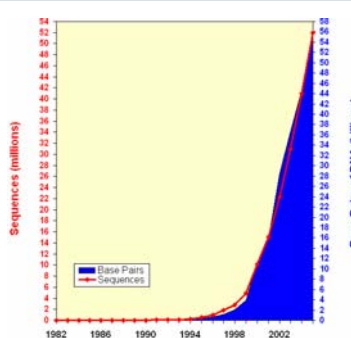
- Com o avanço da tecnologia, existem cada vez mais sequências e anotações e não é possível determinar a quantidade de informações que ainda será obtida de diversos organismos com o andamento do projeto genoma.
- Isso torna fundamental o uso de um banco de dados **bem estruturado** que permita o armazenamento, o acesso e o processamento destas informações de forma simples e eficiente.



Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

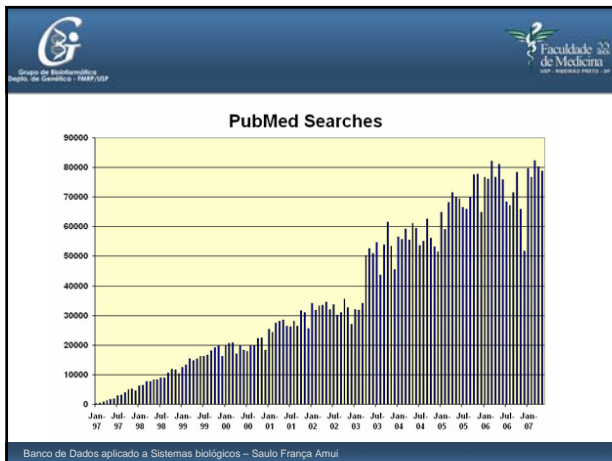



## Dados Biológicos



**Crescimento do Gen-Bank (1982 a 2005)**

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui



Grupo de Bioinformática  
Depto. de Genética - FMBP/USP

Faculdade de Medicina  
USP - Avenida Pasteur 301

### Banco de Dados

- Conjunto de informações relacionadas entre si, referentes ao mesmo assunto, organizadas prática e racionalmente, para que o usuário levante e **recupere informações**, tire conclusões e tome decisões.

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

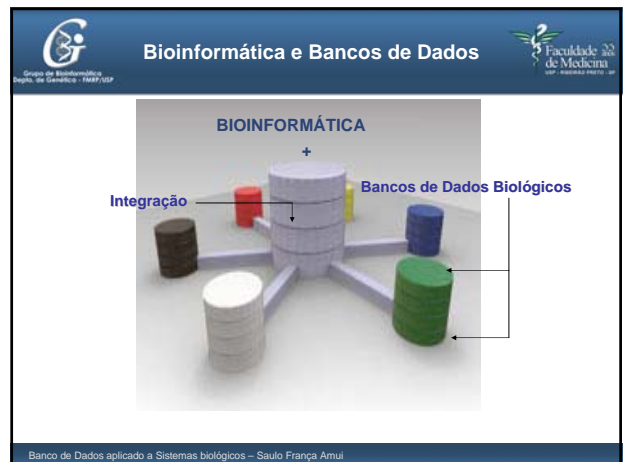
Grupo de Bioinformática  
Depto. de Genética - FMBP/USP

Faculdade de Medicina  
USP - Avenida Pasteur 301

### Bioinformática e Bancos de Dados

- Representação, armazenamento e a distribuição de dados → aspecto funcional da bioinformática.
- O desenvolvimento de ferramentas analíticas para revelar o conhecimento contido nos dados é o segundo, e mais científico, aspecto da bioinformática.

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui



Grupo de Bioinformática  
Depto. de Genética - FMBP/USP

Faculdade de Medicina  
USP - Avenida Pasteur 301

### Bancos de Dados Biológicos

- Todos os dados resultantes das análises de um projeto genoma são armazenados nos chamados bancos de dados biológicos.

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

Grupo de Bioinformática  
Depto. de Genética - FMBP/USP

Faculdade de Medicina  
USP - Avenida Pasteur 301

### Bancos de Dados Biológicos

- Constitui um grande conjunto de dados persistentes, geralmente associado a um software projetado para atualizar, consultar e recuperar componentes dos dados armazenados no sistema. (Bioinformatics FactSheet 2004)
- Eficácia:** Fácil acesso às informações.
- Objetivo:** Métodos para extrair somente as informações necessárias para responder uma específica pergunta biológica.

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

**Bancos de Dados Biológicos**

- Visão direcionada para a biologia molecular
- Base da biologia molecular: DNA - Nucleotídeos {A, C, G, T}
- Além do DNA, existem outros tipos de dados na linha de evolução (RNA e proteínas)
- Uma seqüência de DNA pode possuir milhares de pares de nucleotídeos.
- Cada seqüência possui uma identificação, funções biológicas e podem pertencer a vários organismos.

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

**Bancos de Dados Biológicos**

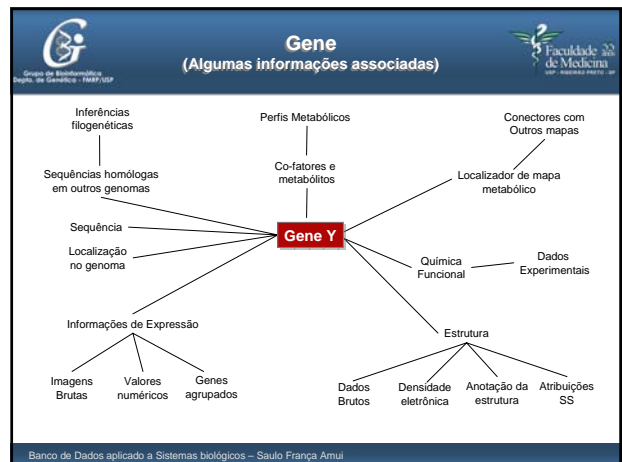
Organismos	Número de genes	Pares de base
Plantas	< 50.000	< 10 <sup>11</sup>
Mamíferos	100.000	3 x 10 <sup>8</sup>
Vermes	14.000	~ 10 <sup>8</sup>
Bactérias	2-4.000	< 10 <sup>7</sup>
<b>dsDNA vírus</b>		
Vacina	< 300	187.000
<b>ssRNA vírus</b>		
Influenza	12	13.500

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

**O que se pode descobrir sobre um gene por meio de uma busca a um Banco de Dados?**

- **Informação evolutiva:** genes homólogos, freqüências dos alelos, ...
- **Informação genômica:** localização no cromossomo, introns, regiões reguladoras, ...
- **Informação estrutural:** estruturas da proteína correspondente, tipos de folds (grande similaridade estrutural), domínios estruturais, ...
- **Informação de expressão:** expressão específica a um dado tecido, fenótipos, doenças, ...
- **Informação funcional:** função molecular/enzimática, papel em diferentes rotas, papel em doenças, ...

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui



**Bancos de Dados Biológicos**

- Disponibilizar dados biológicos para os cientistas
- Dados publicados podem ser difíceis de encontrar ou acessar
- Coleta-los da literatura consome muito tempo
- Disponibilizar dados em formato que possa ser lido por um computador

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui


**Estruturação dos Dados X Acessibilidade**

- O entendimento da diferença entre **dados estruturados** e **desestruturados**, e o projeto de um formato que se ajuste ao seu armazenamento de dados e às suas necessidades de acesso, é a chave para tornar os dados **úteis e acessíveis**.
- Formas de Armazenamentos
  - Diferentes bases de dados para armazenamento de dados biológicos

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

**Formas de Armazenamento**

- arquivos em formato de texto
- arquivos estruturados
- bancos de dados relacionais
  - dados estão guardados em tabelas
- bancos de dados objeto-relacionais
- bancos de dados orientado a objetos



Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

**Formas de Armazenamento**



ACTTCGCA  
TCCAGTCG  
CGGTACTA...

(arquivo texto)

(bases de dados disponíveis na Web)

o biólogo extrai os dados necessários para a sua pesquisa

(anotação da seqüência com base nos dados obtidos)

Iniciando com uma seqüência DNA em um arquivo texto ...

usa-se bases de dados na Web para pesquisar seqüências similares

o Sequência idêntica? funções, categorias, family, referências

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

**Realidades e Problemas comuns**

- Muitas bases de dados são construídas pelos próprios biólogos
- Não padronização da taxonomia
- Dificuldade na adoção de um vocabulário entre os grupos de pesquisa
- Termos diferentes para conceitos iguais
- Conceitos diferentes para termos iguais

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

**Realidades e Problemas comuns**

- **Qualidade dos dados disponíveis na Web**
  - Grupos de pesquisa submetem suas descobertas
  - Algumas bases aceitam de qualquer maneira
  - Muitas bases não verificam a qualidade dos dados
  - Outras bases preocupam-se com a qualidade dos dados, onde um comitê valida-os. Estas bases ganham destaque da comunidade científica.

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

**Qualidade da Informação**

- “Sua capacidade de julgar a qualidade da informação e do software encontrados na Web aumentará à medida que você continuar a atuar neste campo” (Gibas & Jambeck 2001)
  - Qual é a Fonte? Quem são os autores? Qual é o propósito da organização responsável pelas informações? É uma Organização acadêmica? Um Órgão Governamental? Uma empresa?
- Estima-se que qualquer seqüência do GenBank provavelmente contenha pelo menos um erro.
  - O software oferecido por órgãos públicos como o NCBI e o PDB pode ainda estar em desenvolvimento. Grande parte dele é de boa qualidade.
- Transparência (Documentação) e Atualização

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

**Realidades e Problemas comuns**

- **Versionamento dos dados**
  - Dados não acurados podem ser melhorados
  - Novas versões sobre a anotação de uma seqüência são submetidas às bases de dados
  - Muitos pesquisadores “carregam” toda uma base de dados para seu ambiente local, logo, não possuem as últimas atualizações.
  - O custo de armazenamento e do tratamento das versões pode levar algumas bases de dados à simples atualização da “versão” corrente (a única)

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

**Na prática ....**

- Para que estes bancos de dados possam ser realmente utilizados na prática é necessário tratar de vários pontos importantes:
  - a definição do modelo de dados mais adequado,
  - as necessidades de processamento,
  - as análises e controles semântico dos dados,
  - os meios de acesso e o problema da integração das bases de dados.

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

- Avanço da tecnologia → crescimento exponencial do volume de biosseqüências → dados submetidos aos bancos de dados através da **Internet**
- Grande facilidade na submissão de biosseqüências aos bancos de dados → muito importante para que os biólogos possam acessar e fazer suas análises em novos dados mais rapidamente.
- Existem **diversos bancos de dados**, cada um com um modelo de dados distinto e utilizando diferentes tecnologias, sobre os quais os usuários têm necessidade de interagir.

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

**Bancos de Dados Públicos**

- Bancos de Dados Públicos (mais de 348).
  - BD de seqüências de nucleotídeos
    - EMBL (<http://www.ebi.ac.uk/embl>)
    - GenBank (<http://www.ncbi.nlm.nih.gov/GenBank>)
    - DDBJ (<http://www.ddbj.nig.ac.jp>)
  - BD de seqüências de proteínas
    - SWISS-PROT, TrEMBL (<http://www.expasy.ch/sprot>)
    - PIR (<http://pir.georgetown.edu>)
  - BD de motivos
    - Pfam (<http://www.sanger.ac.uk/Software/Pfam>)
    - PROSITE (<http://www.expasy.ch/prosite>)
  - BD de estruturas macromoleculares 3D
    - PDB (<http://www.rcsb.org/pdb>)

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

**O que se deseja saber sobre uma seqüência?**

- Essa seqüência é similar a de algum gene conhecido? Quão semelhante? Qual é o seu significado?
- O que nós conhecemos sobre o **gene 1234**?
  - Genômica (localização, regiões regulatórias)
  - Estrutura (domínios?)
  - Funcional (doenças)
- Informações evolucionárias:
  - Esse gene pode ser encontrado em quais organismos?
  - Qual é a sua taxonomia?

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

**Classificação dos bancos de dados**

- Primários
  - Deposição direta de seqüências sem qualquer processamento ou análise (não curados).
  - GenBank, EMBL-Bank, DDBJ, etc.
- Secundários
  - Derivam dos primários porém com alguns tipos de análises (geralmente curados).
  - Swiss-prot, Uniprot, PROSITE, Blocks, PDB, etc.

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

**Bancos de Dados Compostos**


NRDB	OWL	MIPSX	SP + TrEMBL
PDB	SWISS-PROT	FIRL-4	SWISS-PROT
SWISS-PROT	PIR	MIPSOm	TrEMBL
PIR	GenBank	MIPSTn	
GenPept	NRL-3D	MIPSH	
SWISS-PROTupdate		FIRMOd	
GenPeptupdate		SWISS-PROT	
		EMTrans	
		GBTtrans	
		Kabat	
		PaaqIP	

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui

**Membros do Consórcio internacional**

**INSDC - International Nucleotide Sequence Database Collaboration**

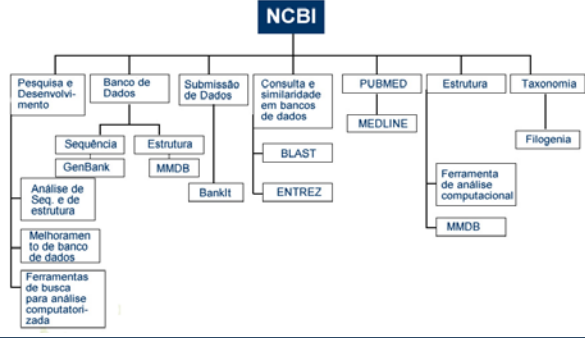
- Genbank (NCBI - National Center for Biotechnology Information)
- EMBL (European Molecular Biology Laboratory)
- DNA DataBank of Japan (DDBJ)



<http://www.insdc.org>

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui.

**NCBI**



Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui.

**Genbank**


Genbank (<http://www.ncbi.nlm.nih.gov/>)

- Responsável
  - National Center for Biotechnology Information (NCBI);
  - Desde 1982.
- Assunto
  - Dados sobre seqüências de ácidos nucléicos (DNA) e Proteínas derivadas;
  - Outros dados associados.
- Aproximadamente 61.132.599 seqüências

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui.


**Genbank - Formato**

- **HEADER**
  - Título
  - Taxonomia
  - Citação
- **FEATURES**
  - Gene
  - mRNA
  - AA
- **SEQÜÊNCIA**



Banco de Dados Biológicas - Rita Botelho San-Saúd (rbotelho@fmrp.usp.br)

**NCBI - Vários**



Banco de Dados Biológicas - Rita Botelho San-Saúd (rbotelho@fmrp.usp.br)

**NCBI - Entrez**



**ENTREZ**

- Interface de busca (keywords);
- Busca de genes, seqüências, proteínas, referências, etc.

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amui.

**NCBI - PubMed**

Banco de Dados Biológicos - SisRets - RodrigoSanCesari@fmp.usp.br

**NCBI - Blast**

Banco de Dados Biológicos - SisRets - RodrigoSanCesari@fmp.usp.br

**NCBI - Protein Clusters (novo)**

Banco de Dados aplicado a Sistemas biológicos - Saulo França Amui

**EBI**

Banco de Dados aplicado a Sistemas biológicos - Saulo França Amui

**EMBL**

EMBL (<http://www.ebi.ac.uk/>)

- Responsável
  - European Bioinformatics Institute (EBI);
  - Desde 1982.
- Assunto
  - Dados sobre seqüências de ácidos nucleicos;
  - Outros dados associados.
- Aproximadamente 98.271.924 seqüências

Banco de Dados aplicado a Sistemas biológicos - Saulo França Amui



**EMBL - Estatísticas**

Banco de Dados aplicado a Sistemas biológicos - Saulo França Amui










**PDB**  
 (Protein Data Bank)
 

[Home](#) | [Structure Summary](#) | [Biological & Chemical](#) | [Materials & Methods](#) | [Sequence Details](#) | [Statistics](#)

**4cha**  Learn more: [M] DOI: 10.2210/pdb/4cha.pdb

**Structure and Visualization**

**1a4a**  **Biological Molecule**

**Title:** STRUCTURE OF ALPHA-CHYMOTRYPSIN REFINED AT 1.68 ANGSTROMS RESOLUTION

**Authors:** Tsukada, H., Blom, D.M.

**Primary Citation:** Tsukada, H., Blom, D.M. Structure of alpha-chymotrypsin refined at 1.68 Å resolution. *J Mol Biol* pp.703-717, 1985 [Abstract] [PDF]

**History:** Deposition: 1984-11-26 Release: 1985-06-01

**Experimental Method:** Type: X-RAY DIFFRACTION Data [ [ EDC ] ]

**Parameters:**

Resolution	B-Value	R-Value	Space Group
1.68	0.179	0.19	P 2 <sub>1</sub> 2 <sub>1</sub> 2 <sub>1</sub>
	[ RMS ]		

**Unit Cell:** Length (Å) : 48.22 9 67.28 c : 65.88  
 Angles (°) : alpha 90.00 beta 120.80 gamma 90.00

**Molecular Description:** Polymer: 1 Molecule ALPHA-CHYMOTRYPSIN  
 Chain: A,B EC no: 3.4.21.1 [PDF]



**Display Options:**

- Stick
- Ball
- Model
- RMS Superposition\*
- RMS Fitting Workhorse
- QuickFit
- MS Straight

 \* Equivalents of interacting biological molecules.

<http://www.rcsb.org/pdb>

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amul.


**Nucleic Acids Research**


- Diversidade de bancos de dados (Nucleic Acids Research)

[Home](#) | [Contact Us](#) | [My Basket](#) | [My Account](#)

**Receive this page by email each issue** [See us for details]

**All articles in this issue are published under an open access license**

**Contents: Volume 35, Database Issue**

**Contents:** (continued) Volume 35 Database Issue, January 1, 2007

<p><b>1087R</b> A database of coiled genes is designed for supporting genome evolution</p> <p><b>1091R</b> The Eukaryote Repeat Database</p> <p><b>1093R</b> A database resource for analysis of non-coding RNA sequences</p> <p><b>1095R</b> The NCPD database: analysis and comparative genomics of alternative splicing in 37 animal species</p> <p><b>1097R</b> Alternative splicing database of complete genomes and manually annotated full-length cDNAs from 28 vertebrates</p> <p><b>1100R</b> A database of orthologous 132 gene orthologous genes</p> <p><b>1102R</b> GenomETool: a whole genome resource for the detection of transposon factor binding sites associated with conserved and non-coding elements across eukaryotic genomes and human 100 genomes</p> <p><b>1104R</b> A database of genome-wide mammalian orthologous nearby genes</p> <p><b>1106R</b> A database of genome-wide mammalian orthologous genes</p> <p><b>1108R</b> The mouse 5hmA database of regulatory sequences in genome-wide transcriptional enhancers</p>	<p>G. Wu, Y. Zhang, J. Qi, W. Li, T. Song, S. He and S. Zhang</p> <p>V. Gaidarov, A. Rodriguez and G. Brown</p> <p>A. Vior, J. G. Thompson, J. Drenth and L. A. Partridge</p> <p>S. Kim, A. V. Abkhovich, M. Bay and C. Lee</p> <p>T. Lee, Y. Lee, S. Kim, Y. Kim, Y. Kim, Y. Kim, W. H. Chung, H. Lee and S. Lee</p> <p>J. Gaidarov, Y. Gaidarov, M. Kim, Y. Kim, S. He and S. Zhang</p> <p>E. A. Aksoy</p> <p>A. El Jai, J. El Jai, S. Ghomriani, M. A. Doudouk, S. Ghomriani, S. Ghomriani and S. Ghomriani</p> <p>V. Gaidarov, C. Rodriguez, D. Bergman, R. J. Gaidarov, F. B. Baker and M. B. Drenth</p> <p>M. Pechkov, L. He, S. He and E. C. Minter</p> <p>A. G. Papanicolaou, V. B. Anagnostou, A. T. K. Vassiliadis and I. Gaidarov</p> <p>C. Zhang, Y. Zhang, Y. Zhang and M. Q. Zhang</p>	<p>076-079</p> <p>080-087</p> <p>088-091</p> <p>091-094</p> <p>094-0101</p> <p>0104-0109</p> <p>0109-0117</p> <p>0118-0123</p> <p>0124-0131</p> <p>0132-0138</p> <p>0139-0141</p> <p>0142-0148</p>
--	--	--

<http://nar.oxfordjournals.org/>

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amul.


**PDB**


**Referências**

B. BOECKMANN, A. BAIRICH, R. APWEILER, et al. "The SWISSPROT protein knowledgebase and its supplement TREMBL in 2003". *Nucleic Acids Research*, pp. 365-370, 2003.

D. A. Benson, I. Karsch-Mizrachi, D. J. Lipman, J. Ostell, B. A. Rapp, D. L. Wheeler. "GenBank". *Nucleic Acids Research* 28(1), pp. 15-18, 2000.

Gibas, C.; Jambek, P.; Desenvolvendo Bioinformática: ferramentas de software para aplicações e, biologia; Rio de Janeiro, Ed. Campus, 2001

H. M. Berman, J. Westbrook, Z. Feng, G. Gilli, T. N. Bhat, H. Weissig, I. N. Shindyalov, P. E. Bourne. "The Protein Data Bank". *Nucleic Acids Research* 28(1), pp. 235-242, 2000.

L. FALQUET, M. PAGNI, P. BUCHER, et al. "The PROSITE database, its status in 2002". *Nucleic Acids Research*, pp. 235-238, 2002.

W. Baker, A. van den Broek, E. Camon, P. Hingamp, P. Sterk, G. Stoesser, M. Ann Tuli. "The EMBL Nucleotide Sequence Database". *Nucleic Acids Research* 28(1), pp. 19-23, 2000.

W. C. Barker, J. S. Garavelli, H. Huang, P. B. McGarvey, B. C. Orcutt, G. Y. Srinivasarao, C. Xiao, L. L. Yeh, R. S. Ledley, J. F. Janda, F. Pfeiffer, H. Mewes, A. Tsugita, C. Wu. "The Protein Information Resource (PIR)". *Nucleic Acids Research* 28(1), pp. 41-44, 2000.

Y. TATENO, K. FUKAMI-KOBAYASHI, S. MIYAZAKI, et al. "DNA Data Bank of Japan at work on genome sequence data". *Nucleic Acids Research*, pp. 16-20, 1998.

Banco de Dados aplicado a Sistemas biológicos – Saulo França Amul.