

Biologia Computacional
Um Espaço Desafiador para os
Profissionais de Computação

Katia S. Guimarães
katiag@cin.ufpe.br

Seminário CIn – 14/dezembro/2007

Roteiro

- Breve Histórico da Era **Genômica**
- Pós-Genoma, a Curta Era **Proteômica**
- A Essencialidade do **Interactoma**
- Nosso Trabalho de Pesquisa e Oportunidades no CIn

Bio-informática vs. Biologia Computacional

Há muita controvérsia quanto a estes termos.

Para alguns

Bio-informática é uma especialização da Informática que trata de desenvolver ferramentas para lidar com dados biológicos.

Biologia (Molecular) Computacional área de pesquisa que combina conhecimentos de Química, Física, Biologia, C. Computação, Matemática e Estatística para atacar problemas de Biologia Molecular.

Histórico da Era Genômica - 1990

Início: Outubro de 1990

Lançamento do Projeto Genoma Humano

- Seqüenciar o DNA humano ($3 \cdot 10^9$ pb) e
- Identificar os estimados 100 mil genes.

Atores Principais:

Consórcio envolvendo EUA, Inglaterra, França, Japão, Alemanha e China.

Prazo: 15 anos (terminaria em 2005)

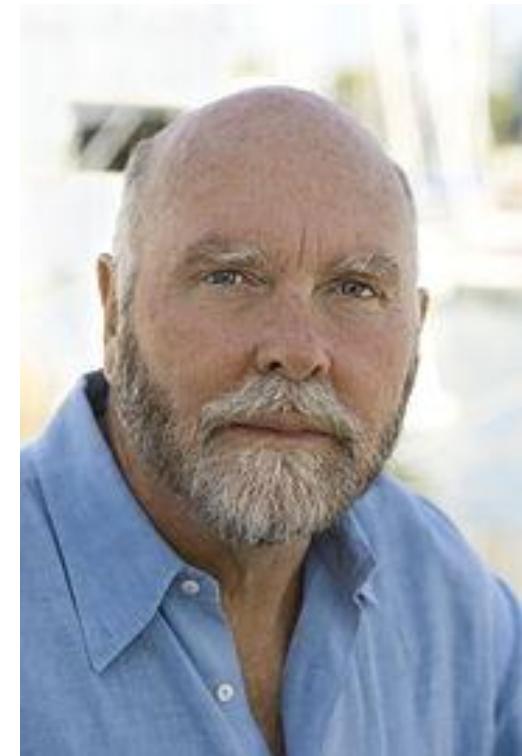
Orçamento: US\$3 bilhões de dólares

Histórico da Era Genômica - 1992

1992

→ Consórcio faz mapas dos cromossomos humanos

→ Craig Venter, pesquisador do NIH, funda *The Institute for Genomic Research (TIGR)*



Histórico da Era Genômica - 1995

Grupo de pesquisadores da TIGR publica na revista *Science* o artigo

Whole-genome random sequencing and assembly of *Haemophilus influenzae* Rd.
Fleischmann RD, Adams MD, et al.

com a seqüência de DNA da bactéria
Haemophilus influenzae (*otite, meningite*)

Tamanho: $2 \cdot 10^6$ bp

Técnica: *double-barrel shotgun sequencing*

- Mais custo computacional
- Muito menos tempo e custo em labs.

Strategies for Generating the Human Genome Sequence

Celera

HGP

Shotgun
fragmentation

Insert
fragment
into vehicle

Random
sequencing
phase

ACTGCA

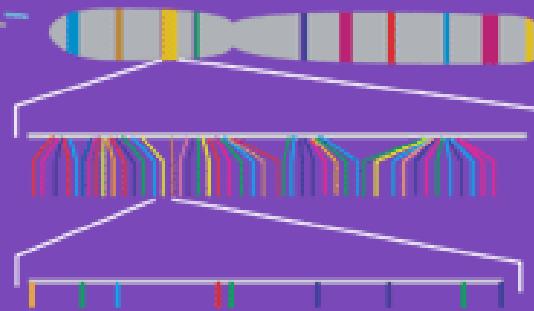
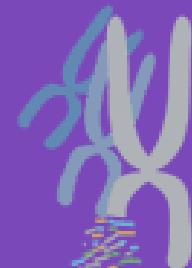
Insert large
fragments
into BAC
vector

Whole-
genome
assembly
(WGA)
approach

Regional
assembly
approach

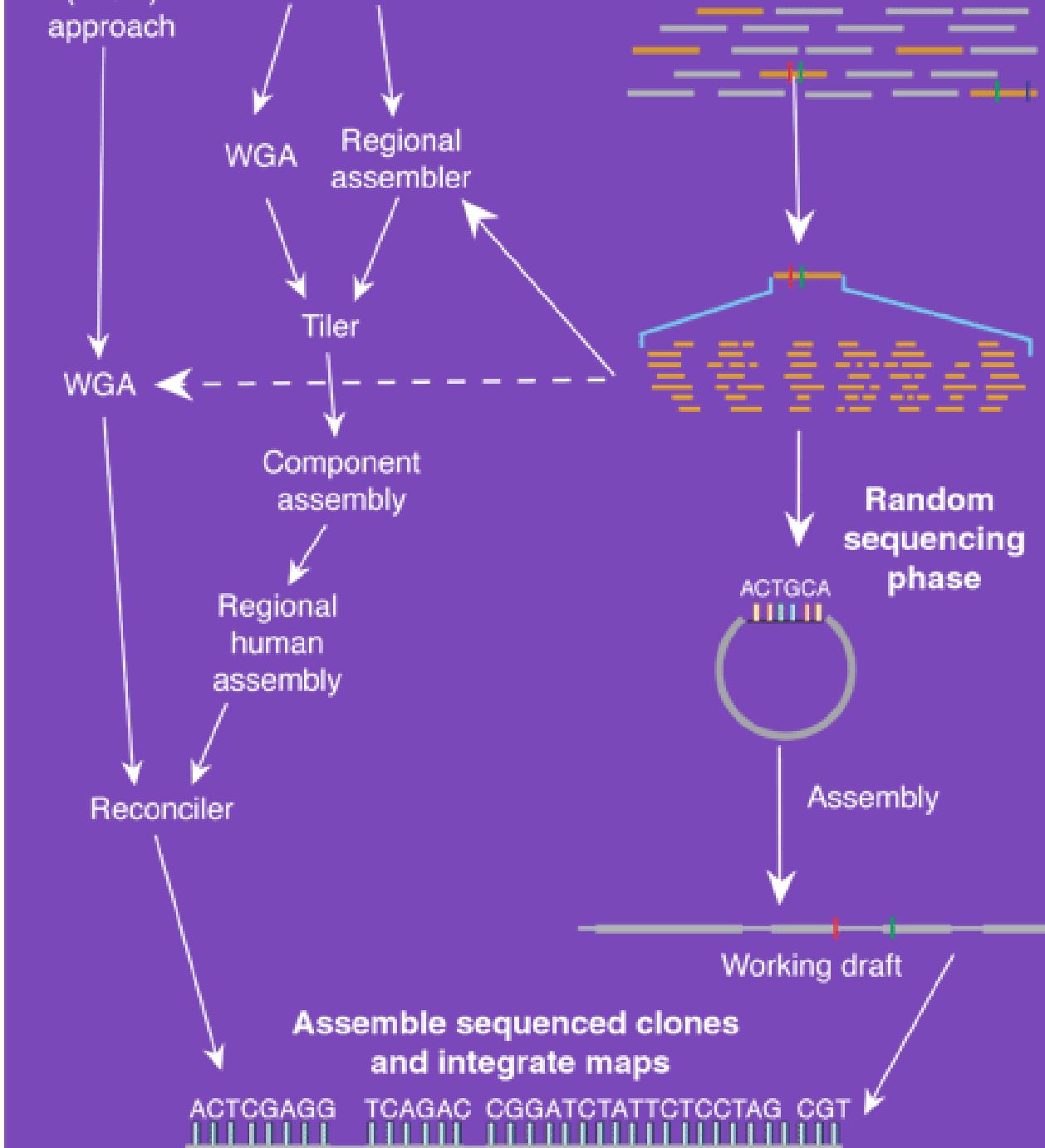
Regional

Landmark maps



Minimum
tiling path
from ordered
BACs





Histórico da Era Genômica – 96-98

Pesquisadores da TIGR publicam as seqüências de DNA de outras bactérias

Mycoplasma genitalium (1996) (menor bactéria)

Methanococcus jannaschii (1997)

1998

Craig Venter se associa com a Applied Biosystems para fundar a Celera Genomics Corp., com o objetivo de seqüenciar o genoma humano em 03 anos (2001), ao custo de US\$300 milhões (1/10 do orçamento do projeto do Consórcio).

Histórico da Era Genômica – 99-00

O Consórcio revê suas previsões, e anuncia a conclusão do seqüenciamento para 2003,

D

cromossomo

Ju

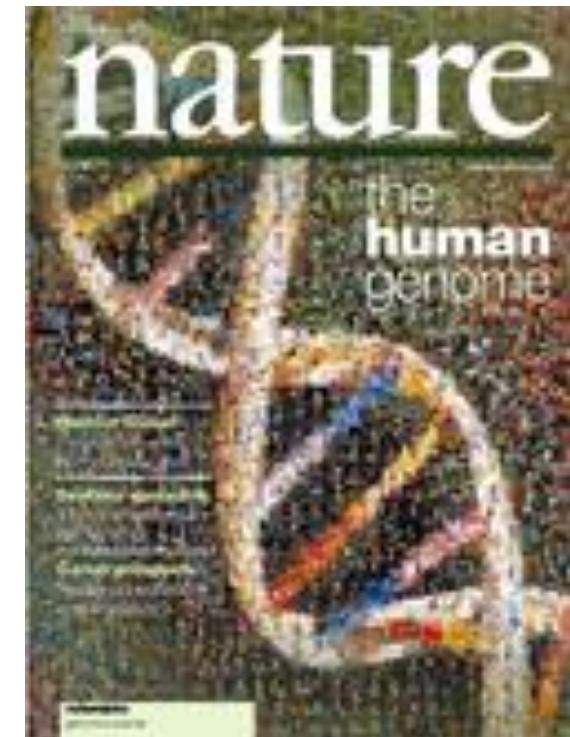
amento do primeiro consórcio mundial.

Anunciam a conclusão do genoma humano.

Histórico da Era Genômica – Fev 2001

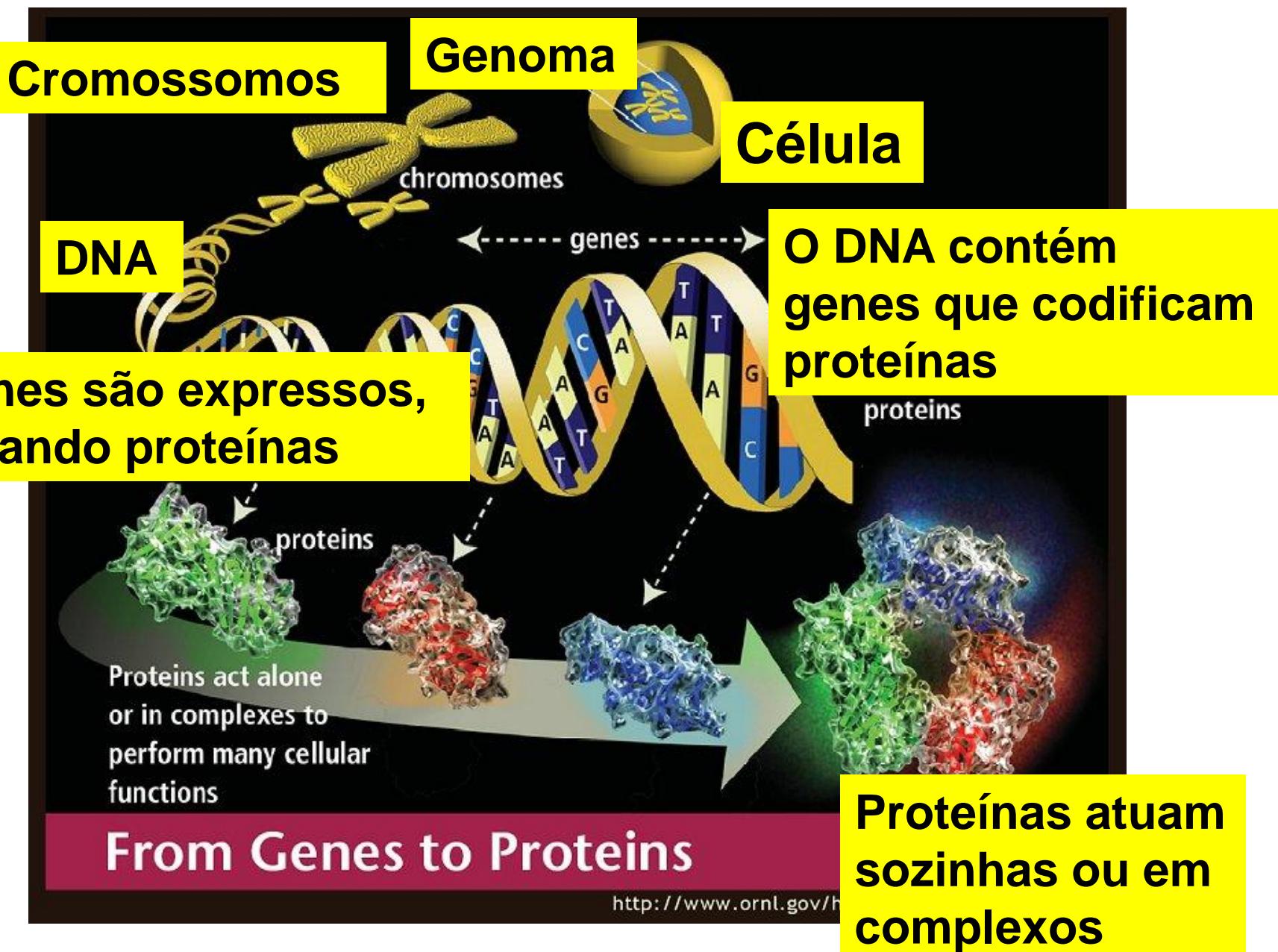


The Sequence of the Human Genome
J. C. Venter, M. D. Adams, E W. Myers, et al.



Initial Sequencing and Analysis of the Human Genome
Consórcio Mundial

Biologia Molecular 101 em 2 Minutos



O GenBank

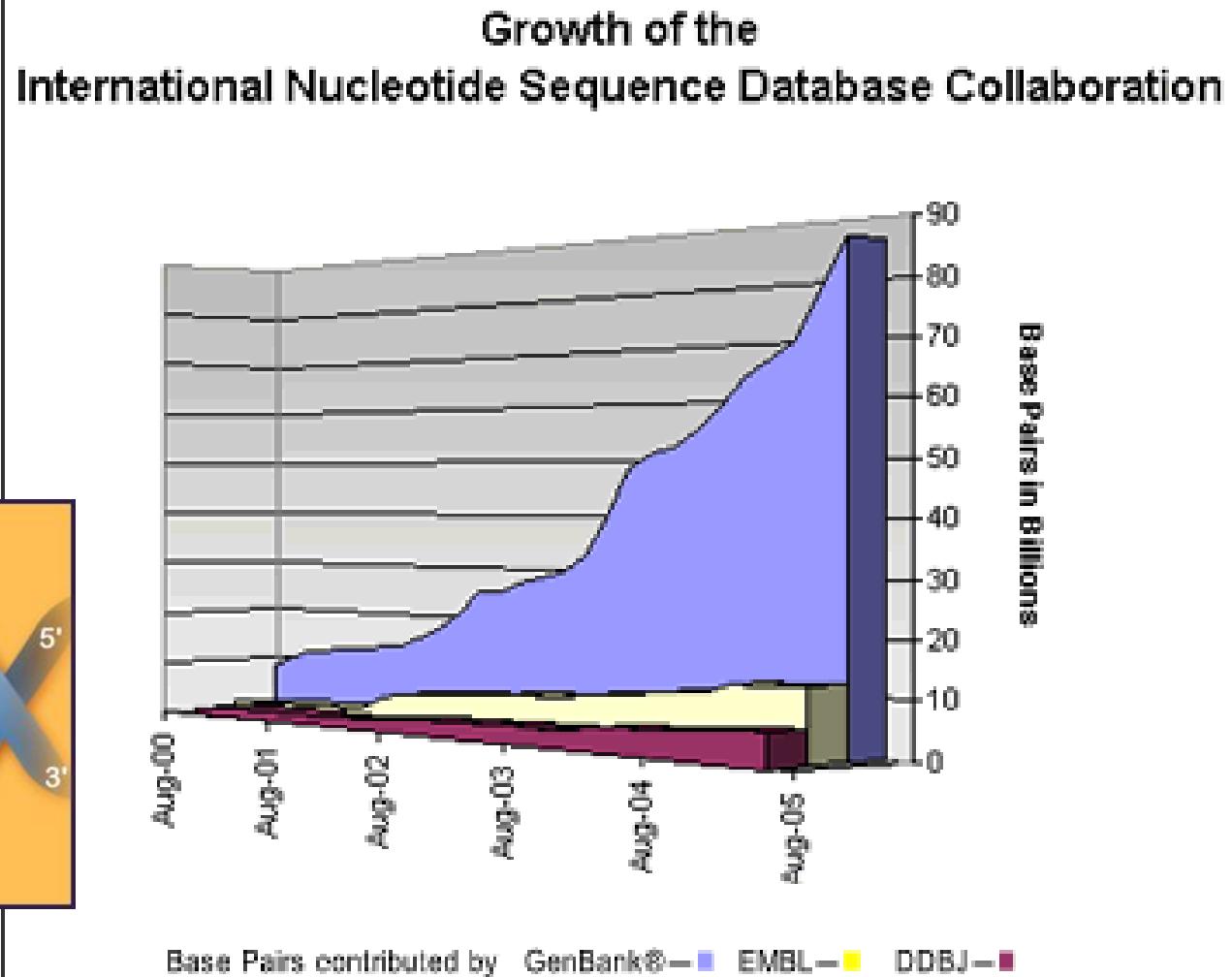
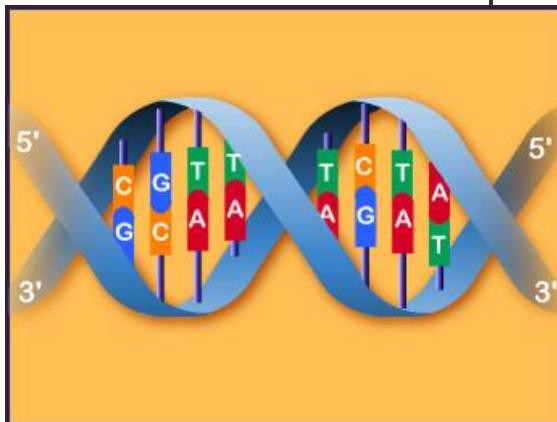
Do ponto de vista da computação, foram criadas inúmeras seqüências sobre o alfabeto {a, c, g, t}

1 gatcctccat atacaacggt atctccacct caggttaga tctcaacaac ggaaccattg
61 ccgacatgag acagtttagt atcgctgaga gttacaagct aaaacgagca gtatcgact
121 ctgcacatctga agccgctgaa gttctactaa gggtgataa catcatccgt gcaagaccaa
181 gaaccgccaa tagacaacat atgtaacata tttaggatat acctcgaaaa taataaaccg
241 ccacactgtc attattataa ttagaaacag aacgcaaaaa ttatccacta tataattcaa
301 agacgcgaaa aaaaaagaac aacgcgtcat agaactttg gcaattcgcg tcacaaataa
361 attttggcaa cttatgttc ctcttcgagc agtactcgag ccctgtctca agaatgtaat
421 aataacccatc gtaggtatgg ttaaagatag catctccaca acctcaaagc tccttgccga
481 gagtcgccct ctttgtcga gtaatttca ctttcatat gagaacttat ttcttatttc
541 ttactctca catcctgttag tgattgacac tgcaacagcc accatacta gaagaacaga
601 acaattactt aatagaaaaaa ttatatctc ctcgaaacga tttcctgctt ccaacatcta
661 cgtatatcaa gaagcattca cttaccatga cacagctca gattcatta ttgctgacag
721 ctactatatc actactccat ctagtagtgg ccacgcccata tgaggcatat cctatcgaa
781 aacaataaccc cccagtggca agagtcaatg aatcgttac atttcaaatt tccaatgata
841 cctataaaatc gtctgttagac aagacagctc aaataacata caattgcttc gacttaccga
901 gctggcttc gtttgactct agttcttagaa cgttctcagg tgaaccttct tctgacttac
961 tatctgatgc gaacaccacg ttgtatttca atgtaatact cgagggtacg gactctgccg

Link: <http://www.ncbi.nlm.nih.gov/Genbank/index.html>

O GenBank do NCBI

Aproxima-
damente
85 bilhões
de pares
de base



Link: <http://www.ncbi.nlm.nih.gov/Genbank/index.html>

O genoma
é como
um livro

dividido
em
capítulos
(cromos-
somos)

e
estes em
palavras
(genes)

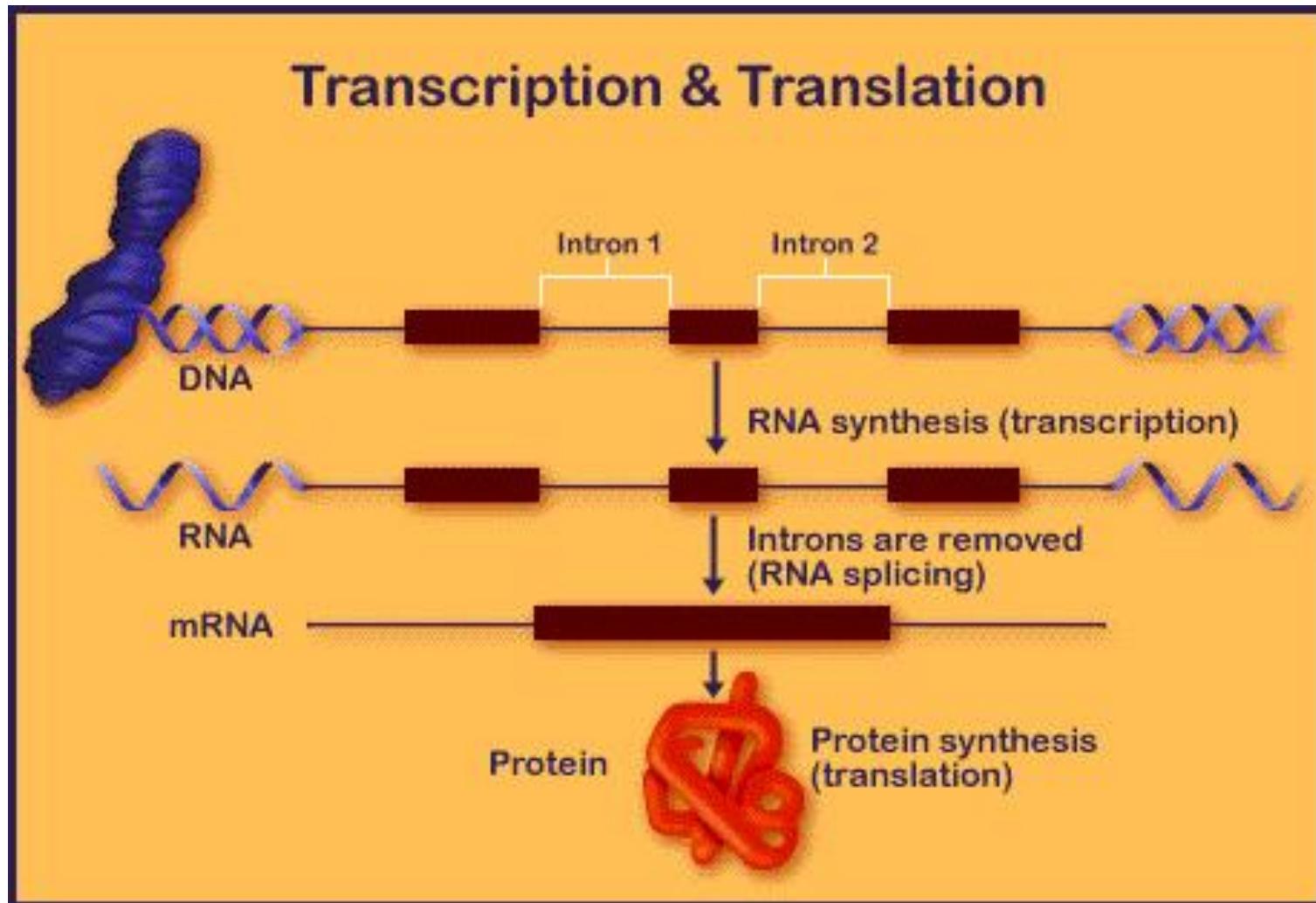
O que fazer com tanta Letrinha?

Rat prion gene D50092 upstream of exon 1: annotations

1 aagttttaa agccacttcc tggtgggaag agagoagtca gocagtaggt atttgattct
61 tagaggaaas agctgatata tttaaccaa goctttcaga gtcocccotgt gggagaggcc
121 agatggaggt gaggggaagg ctggccctt cccataagtga atctaccac acctctctcg
181 cgtccccctgc catctctgac agcaagttga gggccctctg gaaattccaaa gatggggc
241 ccctgatcac aatgtagact ctctgaggca ggaaggcaag gatctgaaag acgttcgt
301 tatttcttaa gaacacagag taatgttgc gcaaggccac caccactat ttaaacaggc
361 tgaaatctgg gtcatctcag cccaaaggcgt tgctccaggc ctctcatgc tgcccttc
421 cccaaacagct gtcaagagcgt tacccacata tgtagacacg cacacacaca cacaacac
481 attttatcat aactataatc ctctggctcc atgactactt cataacaact ttaatctcc
541 tggcccaat actctcttac caccctctcc tcctggatcc taatactggc cacaatatt
601 taatccaaat ccaatttgtt gtttgcatt aattttcaat gtccctccat cccatgtact
661 acaggagat acatttccat gtgggtttatg acccccttgc tcccaatgtact ggttgcatt
721 tatatacatgt goatggaaat atgttatatg tgatggggc tgtagccaca tatataatgt
781 ggcacatgtc atcctatgtc tgaggccac aggtcaatgt cccgtttcc cccatca
841 tccagggtgtc cctggattcc aaactcagggt cccatgtttt gggaaaccaac ccagccccc
901 gatccataac ccccttcttcc ttccaaacaag gtccttattt gttgccttgg tcagccctga
961 acttggatct cctgcctcga ctctctgtat tgccagggtttt cccctgtccaa agtccggagg
1021 catcttgcac aagagcatca ttccctttaa gtcgtttccat ggggtttca tgggggtct
1081 taaatgtatgtactttcc ttggacacaa gttaaaggccaa aacaaggatata aacggatgt
1141 acacgggttgt acatgtcaat aatgtgaaatc ctcaggaaatc tgaggccacaa gggttgc
1201 gagggttggg tcaatgttcaat cttacttagat ggttagaaacag gccaacctgg gctaagg
1261 totgttccaa aacataaaaga aaaaagggggggggggggggggggggggggggggggggggggg
1321 gagagatgtg acatcttaatg agcctgtatg ttgggttgt tgcccaacaa gactgtat
1381 aaaatgggtt atgttggaaa aggccgggggggggggggggggggggggggggggggggggggg
1441 actgtgggggggggggggggggggggggggggggggggggggggggggggggggggggggggggg
1501 aagacagcc aggttgcatttcc ttccatgttccatgttccatgttccatgttccatgttccat
1561 tgaggaggggggggggggggggggggggggggggggggggggggggggggggggggggggggggg
1621 aaaatgtatct tcaactgttccatgttccatgttccatgttccatgttccatgttccatgt
1681 aagggtgttca atggtatatttccatgttccatgttccatgttccatgttccatgttccatgt
1741 ttaatgtata ccataattttccatgttccatgttccatgttccatgttccatgttccatgt
1801 ttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttt
1861 ttaaaatgttccatgttccatgttccatgttccatgttccatgttccatgttccatgttccat
1921 ctggactttt ttagttatgtttttttccatgttccatgttccatgttccatgttccatgttcc
1981 tggtgttttt tataacttgcatgttccatgttccatgttccatgttccatgttccatgttcc
2041 aaatttttttccatgttccatgttccatgttccatgttccatgttccatgttccatgttccat
2101 gtcgttccatgttccatgttccatgttccatgttccatgttccatgttccatgttccatgttcc
2161 tagaattttccatgttccatgttccatgttccatgttccatgttccatgttccatgttccatgt
2221 ggccctgggtt ccgttccatgttccatgttccatgttccatgttccatgttccatgttccat
2281 tacatcaatgttccatgttccatgttccatgttccatgttccatgttccatgttccatgttcc
2341 caAATAAAAGa gtatccatgttccatgttccatgttccatgttccatgttccatgttccat
2401 ccgagaatgttccatgttccatgttccatgttccatgttccatgttccatgttccatgttccat
2461 tagtttgcatgttccatgttccatgttccatgttccatgttccatgttccatgttccatgttcc
2521 acgtatgttccatgttccatgttccatgttccatgttccatgttccatgttccatgttccatgt
2581 gggagccgtt ggttccatgttccatgttccatgttccatgttccatgttccatgttccatgt
2641 cgtcacatgttccatgttccatgttccatgttccatgttccatgttccatgttccatgttccat
2701 taacaatgttccatgttccatgttccatgttccatgttccatgttccatgttccatgttccat
2761 acaatcatgttccatgttccatgttccatgttccatgttccatgttccatgttccatgttccat
2821 cccgttccatgttccatgttccatgttccatgttccatgttccatgttccatgttccatgttccat

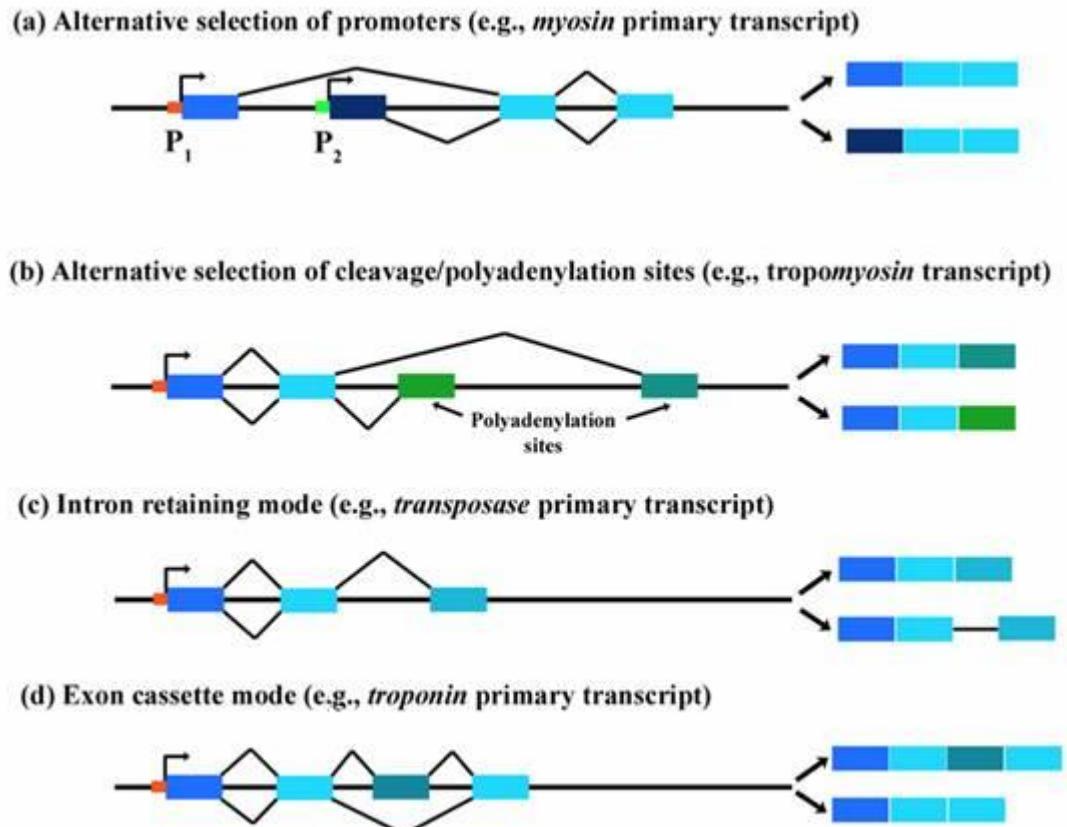
pre-gene region
+RSINE2 445-477
-RSINE2 742-881
+B1 1136-1268
-MER17A 1275-1342
5' UTR of cyt c; intron 1 out-splice
ATG start of cyt c pseudo ORF
intron 2 of cyt c out-splice
end of cyt c ORF
350 bp of 3' UTR homology
interrupting +RDRE1_RN 2174-2266
center of stem loop before poly A signal
poly A signal; end 1100 mRNA; new poly A
452 bp residual prion promoter
conserved motifs
unknown insert in rat with flanking gggg
SP-1, AP-1, AP-2 CCATT sites
exon 1 of rat prion; start of intron 2

Genes são transcritos e traduzidos dando origem a proteínas

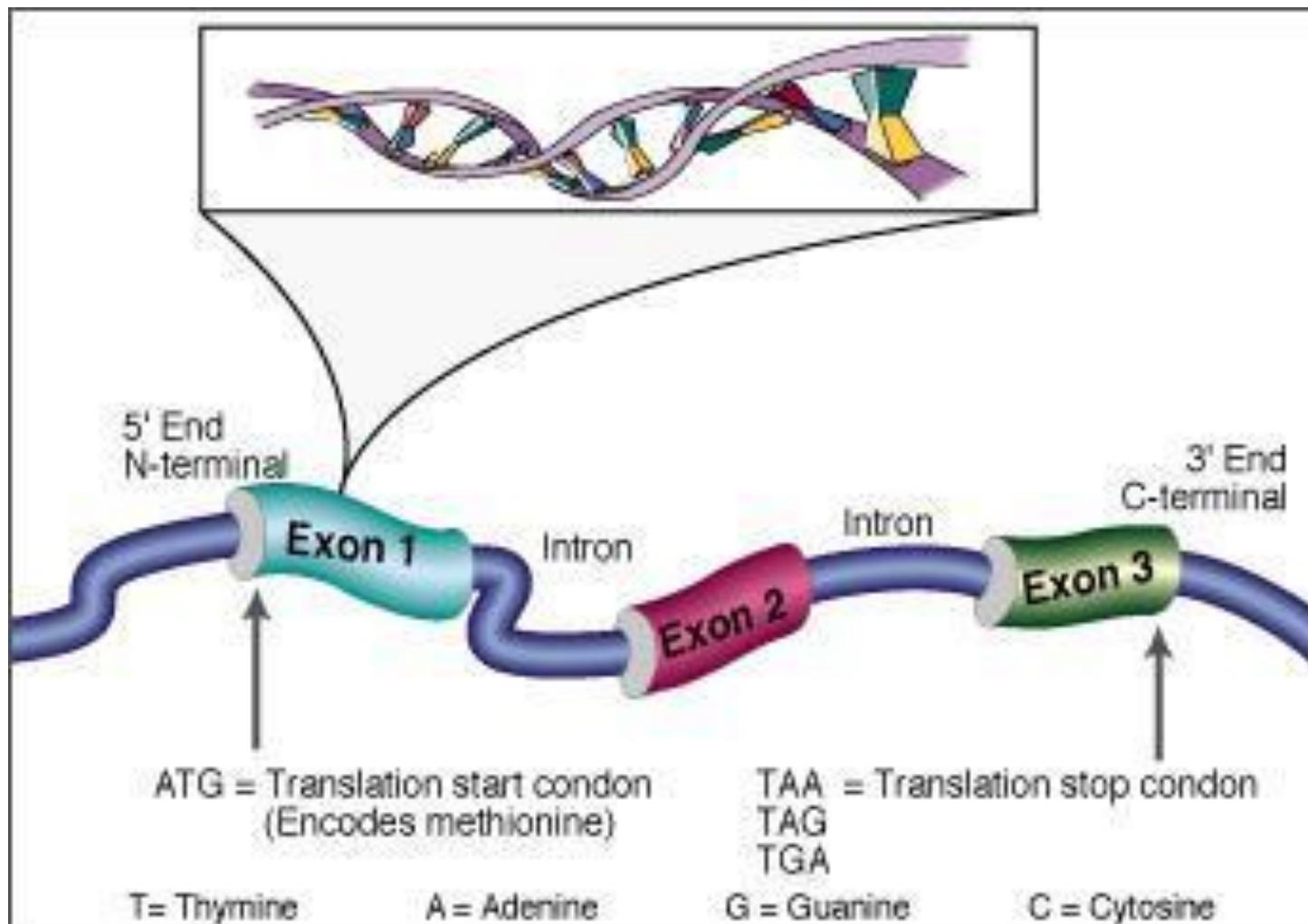


Mas se fosse muito simples não seria tão interessante

Foram encontrados apenas cerca de 25 mil genes, ao invés dos 60 a 100 mil genes estimados a princípio.

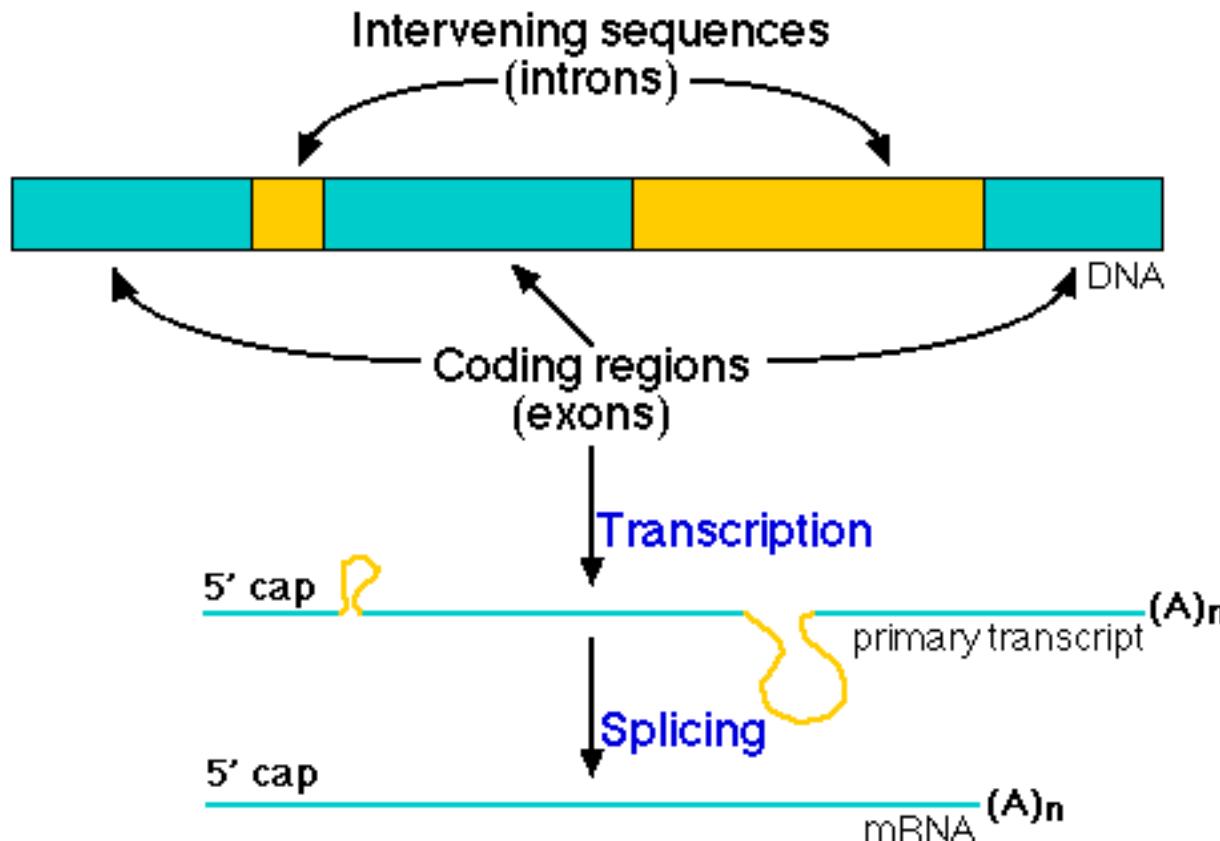


Os genes dos eucariotos são compostos por exons e introns ...



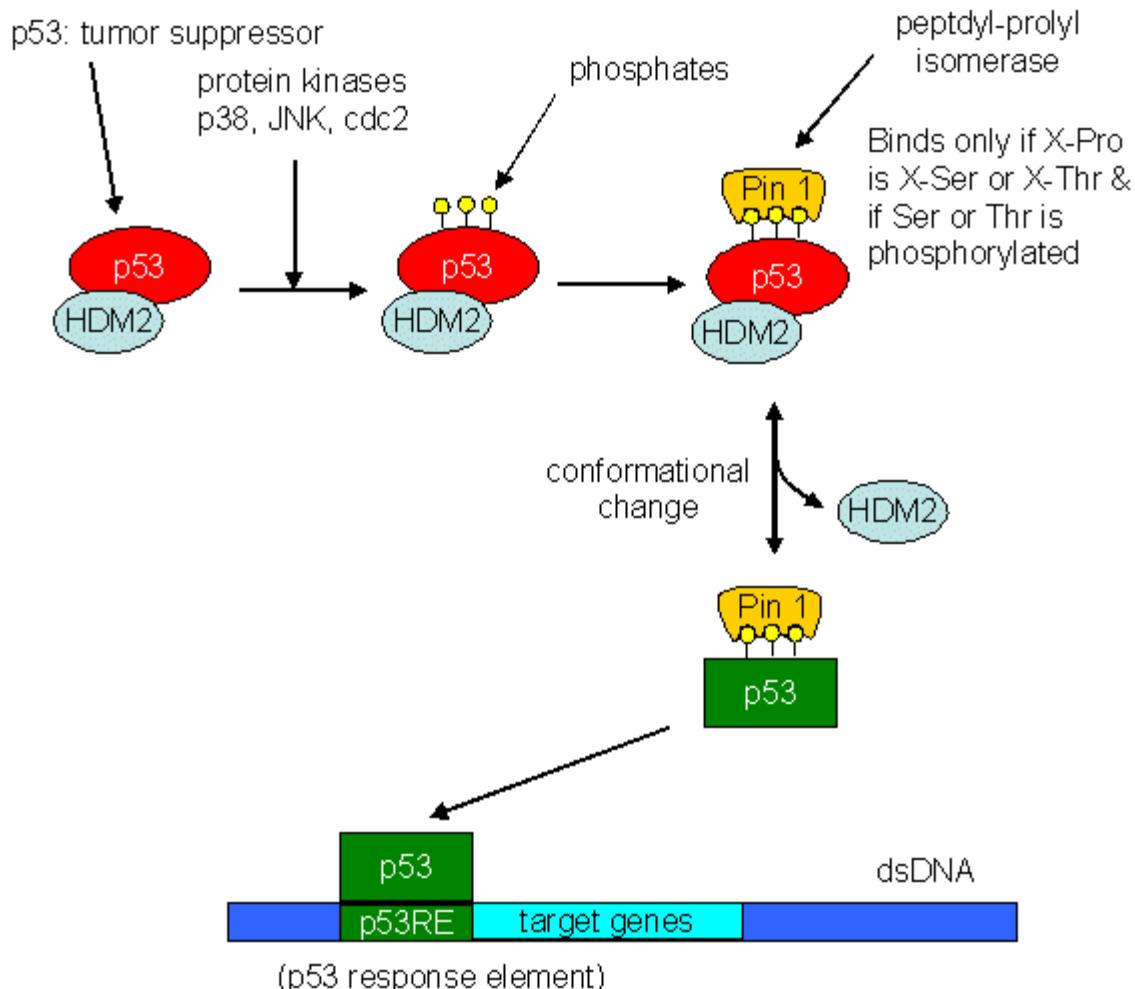
Os introns sempre desaparecem
por ocasião do splicing.

Os exons podem permanecer ou não.



Os genes têm a sua expressão controlada por Fatores de Transcrição

p53 as transcription factor: tumor suppressor activity



(after Ryan and Vousden, Nature, 419, pg 795, 2002)

O foco da
atenção
passou
então
para as
proteínas

Os níveis de expressão de um gene variam com os promotores que se ligam à região reguladora, imediatamente antes do gene

Promoter type	Expression vector promoter region	Expression level
Minimal ————— MCS — TATA — Reporter	Insert-dependent
Constitutive — AP1 — CAAT — AP1 — SP1 — SP1 — TATA — MCS —	High
Cell-specific — Cell-specific — CAAT — TATA — MCS —	Cell type-dependent
Regulated — HRE — HRE — HRE — HRE — TATA — MCS — hormone-response elements	Low → High

O que é Proteômica?

Proteômica é o estudo das proteínas, com suas estruturas e funções.

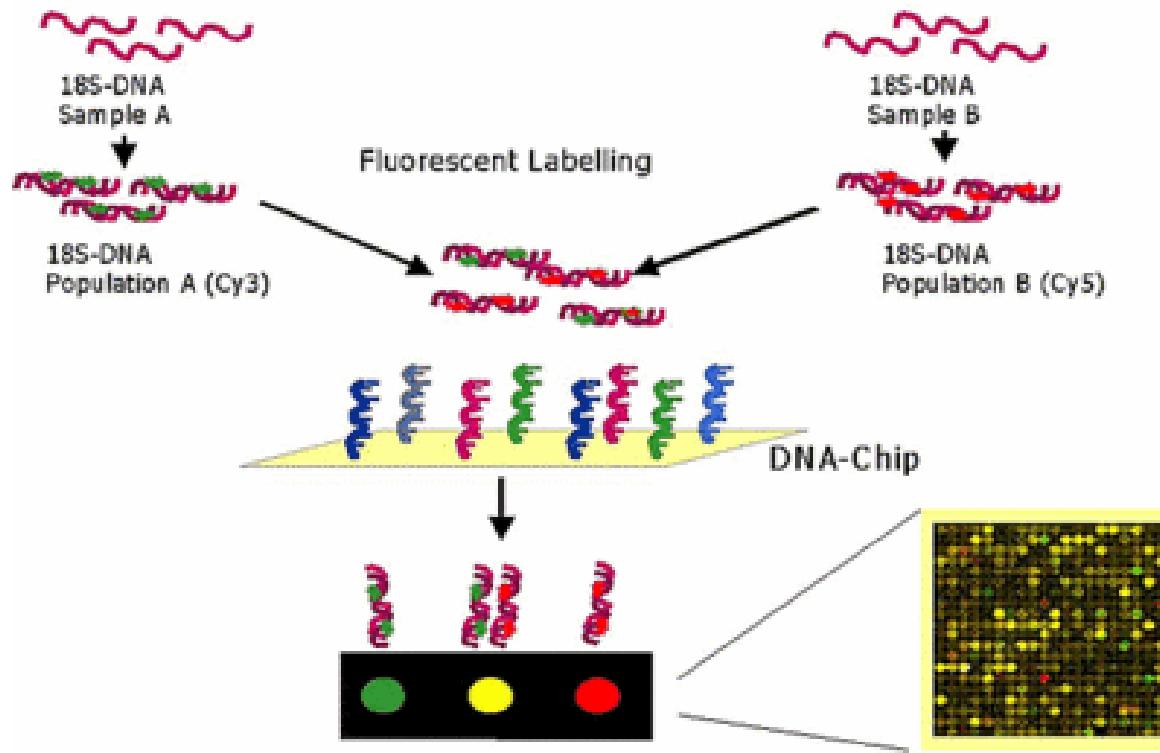
Link para Expasy:

<http://ca.expasy.org/cgi-bin/prosite/PSView.cgi?ac=>

Link para PDB:

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

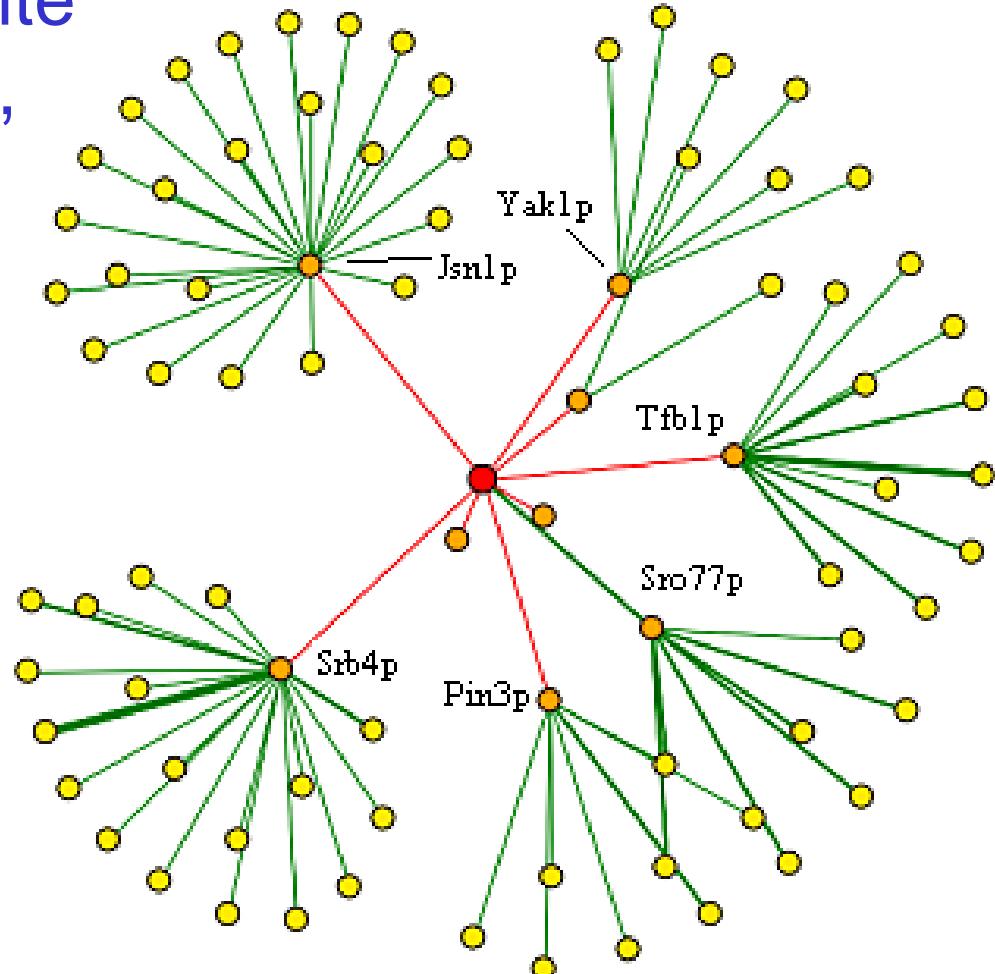
O nível de expressão dos genes em diferentes condições ou em intervalos de tempo pode ser medido



Em setembro de 2006, o [GEO \(Gene Expression Omnibus\)](#) do [NCBI](#) continha mais de 3.2 bilhões de medidas, tomadas sobre mais de 200 organismos.

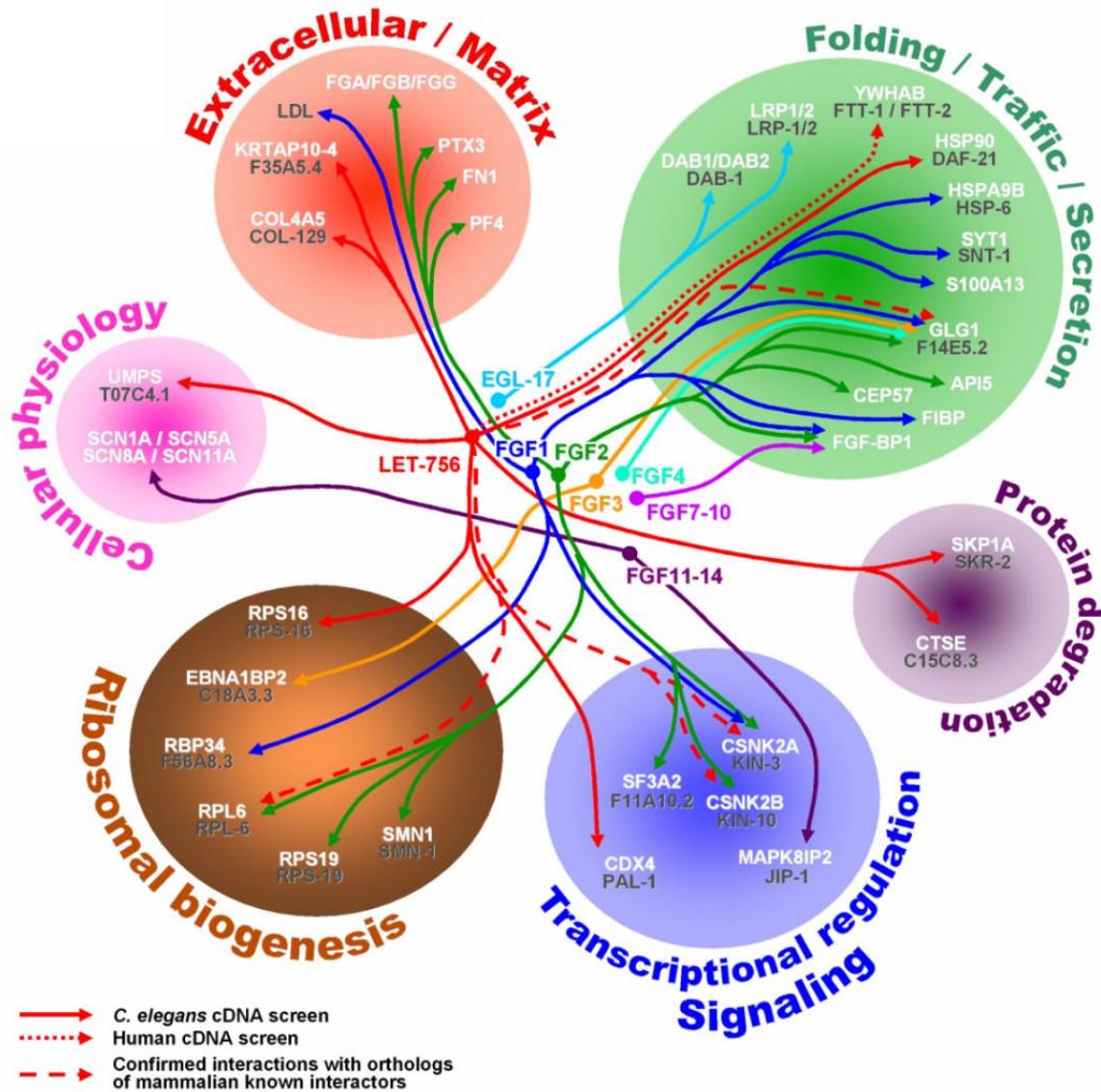
Interações entre Proteínas

Proteínas geralmente atuam em conjunto, e se organizam em redes do tipo *small world*, com muitos nós de grau baixo e poucos nós de grau alto (*hubs*)



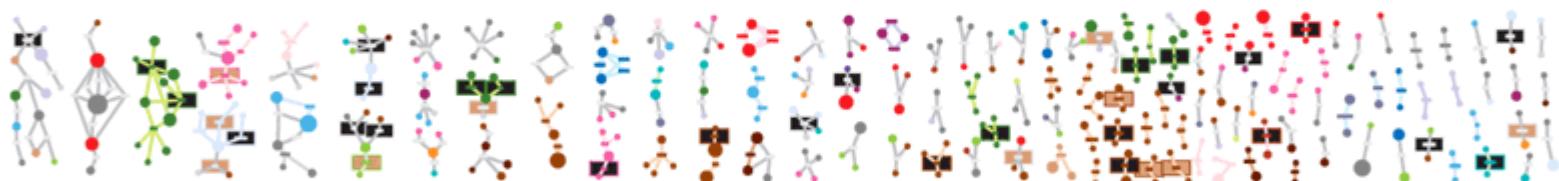
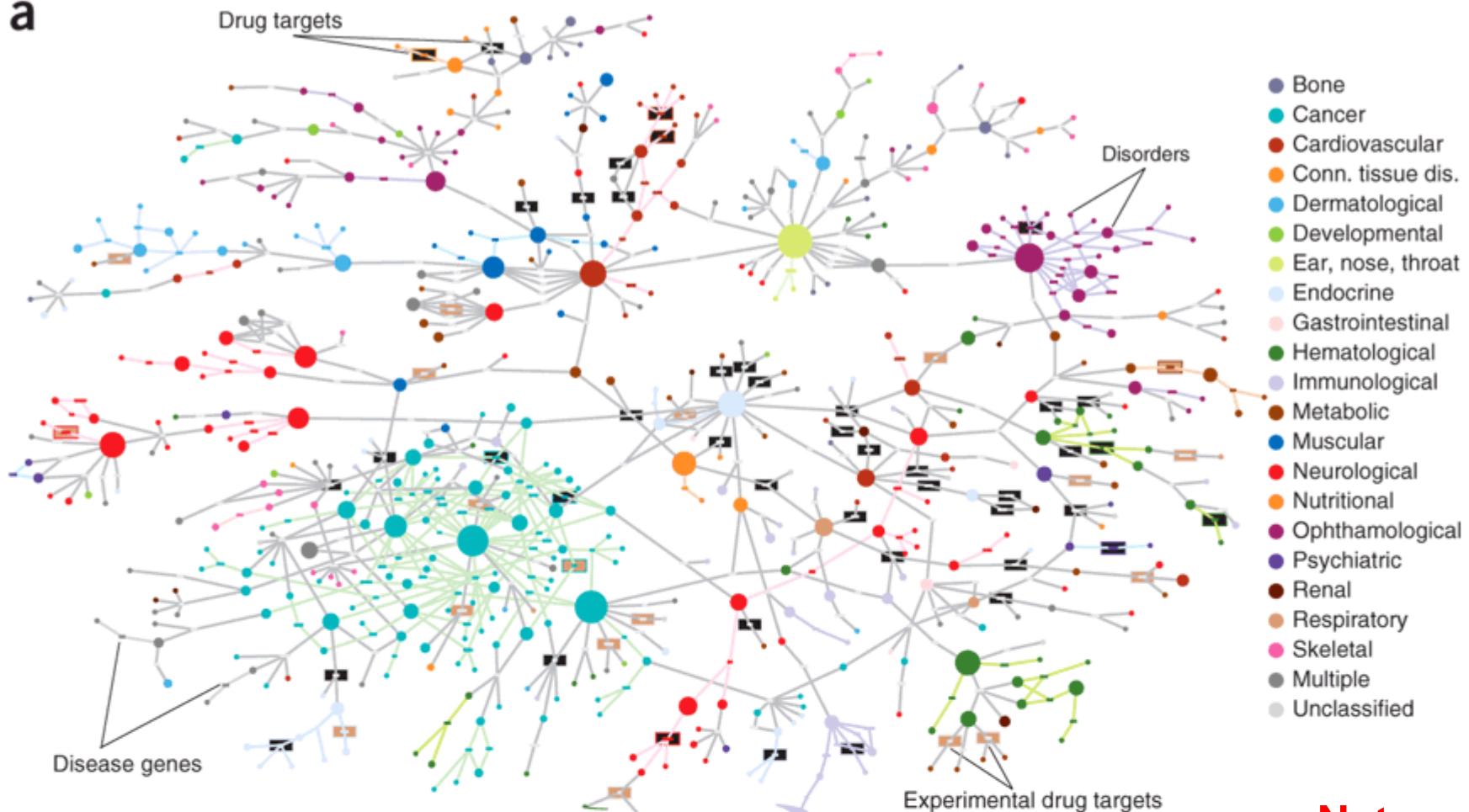
Interações entre Proteínas

Proteínas relacionadas funcionalmente encontram-se a uma distância muito pequena, e em geral são vizinhas nos mapas de interação.



Doenças Humanas e Alvos de Drogas

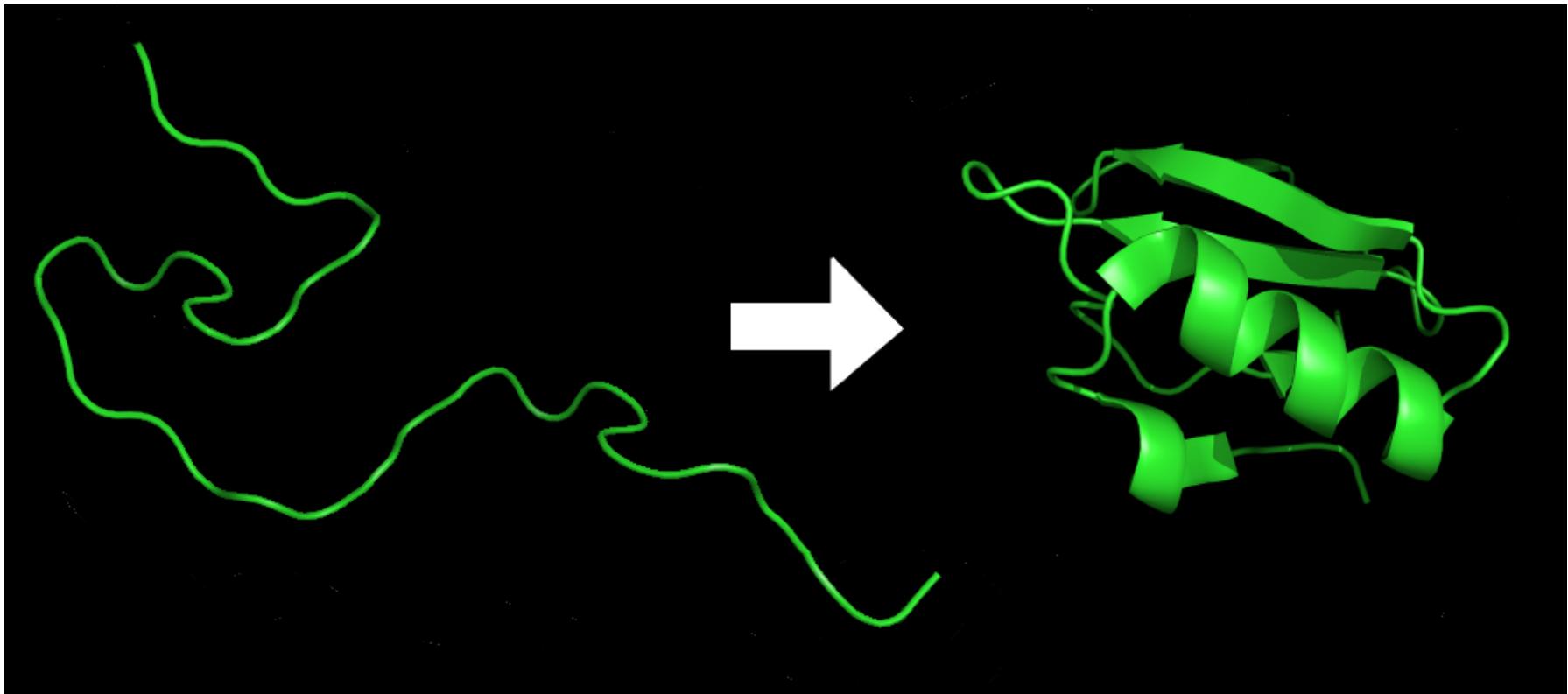
a



**Nature
Biotech,
Out 2007**

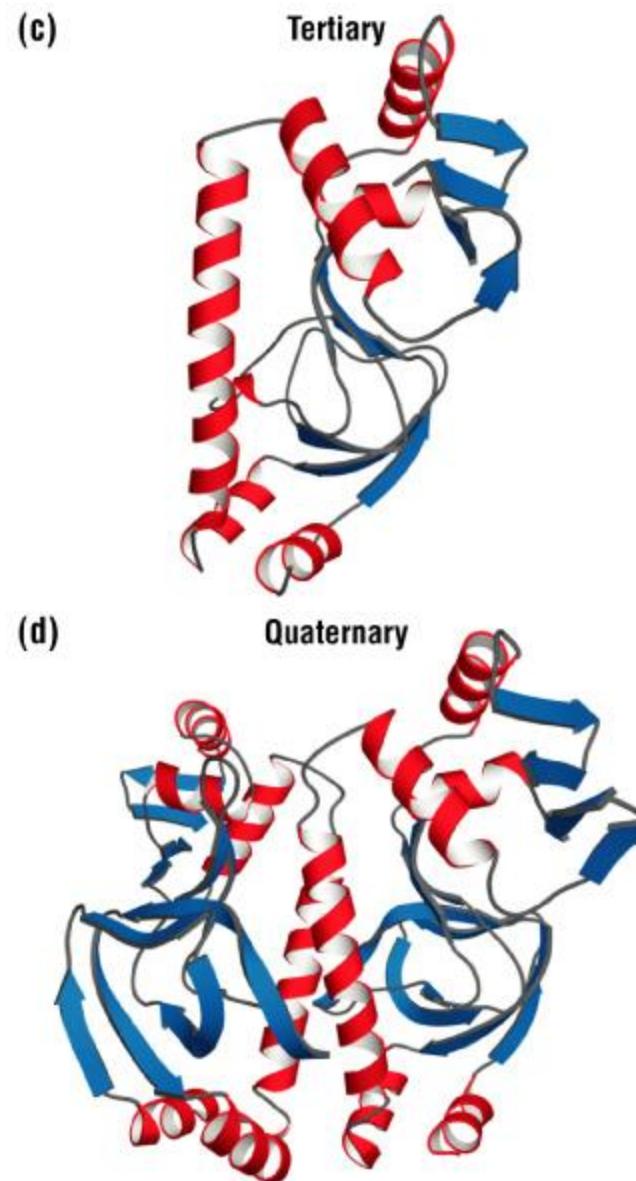
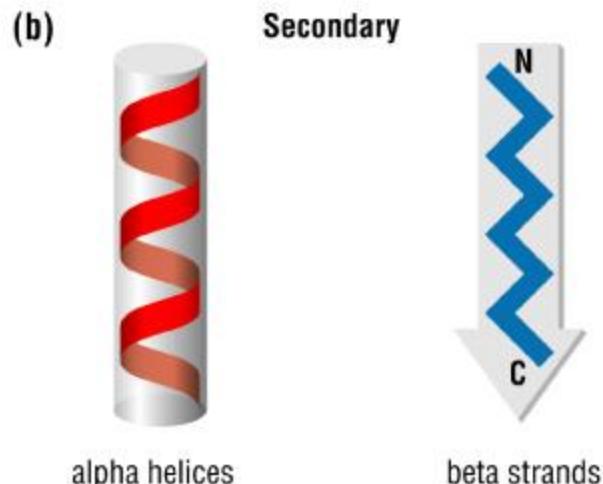
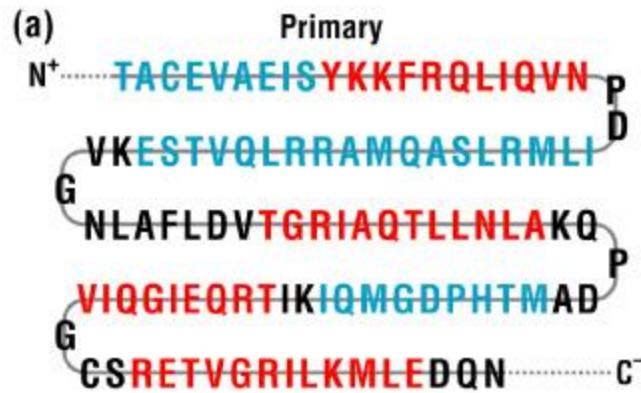
Dobramento de Proteínas

Um problema bem Difícil

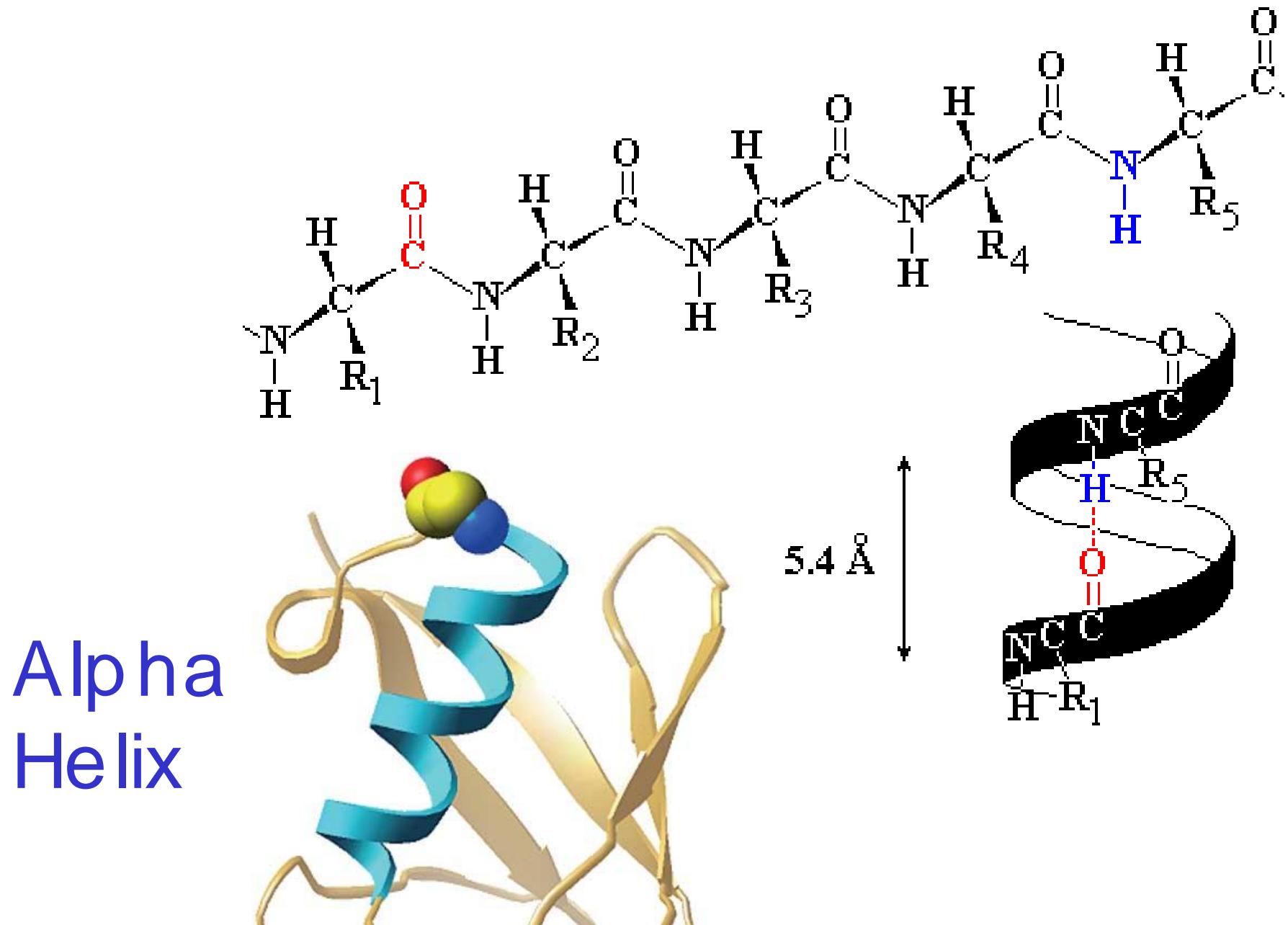


Predição de Estrutura de Proteínas

From **Protein Structure and Function**
by Gregory A Petsko and Dagmar Ringe

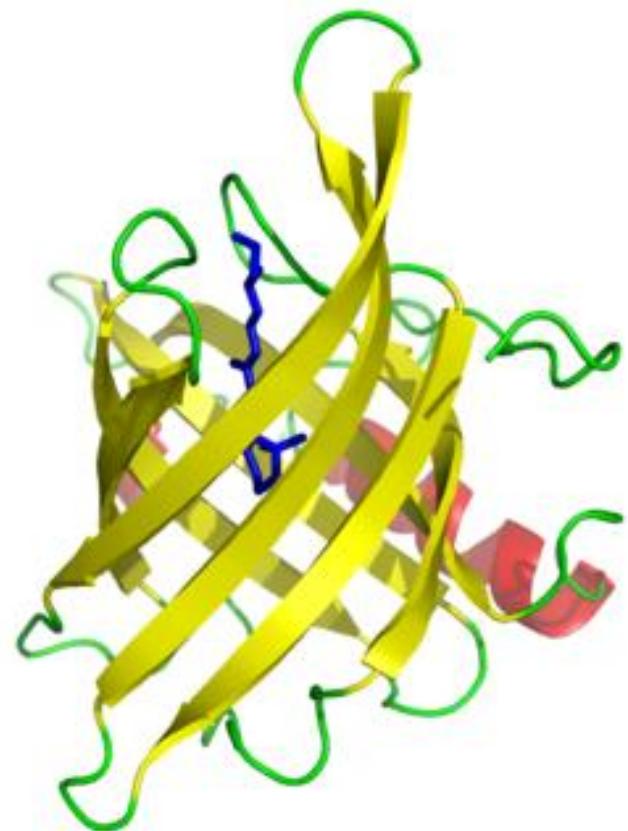
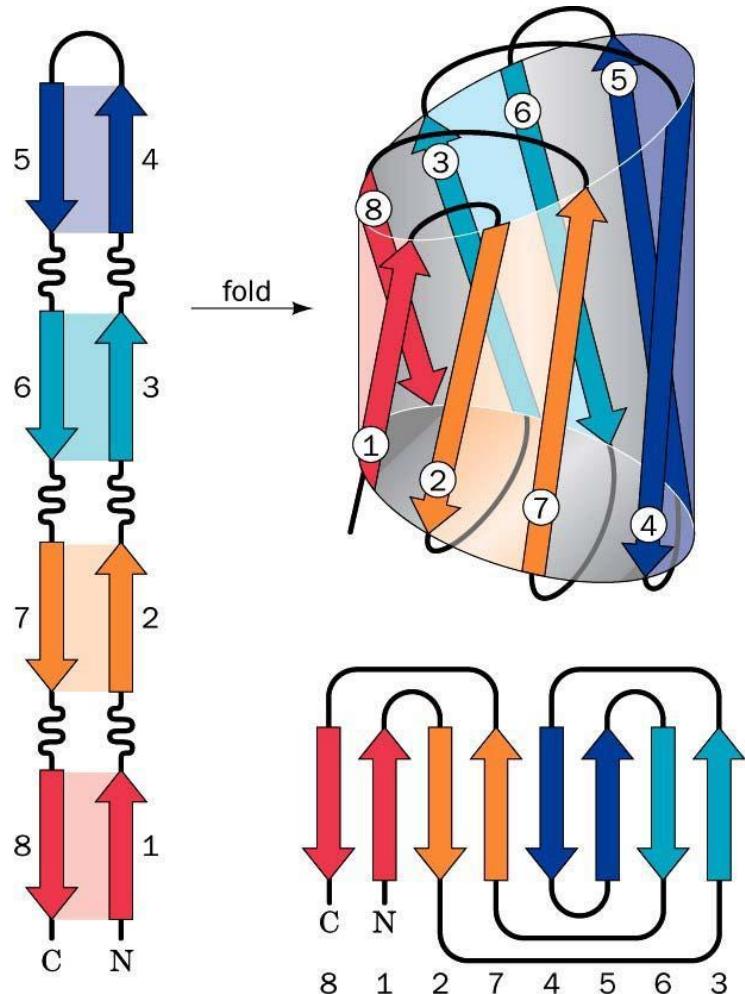


Dobramento de Proteínas



Dobramento de Proteínas

Beta Barrel



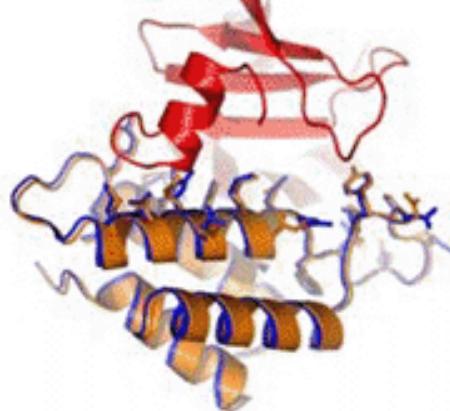
Docking e o Projeto de Drogas

Detalhes como orientação e ângulo de ligação de todos os resíduos do sítio ativo são essenciais.

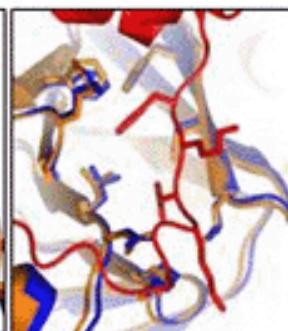
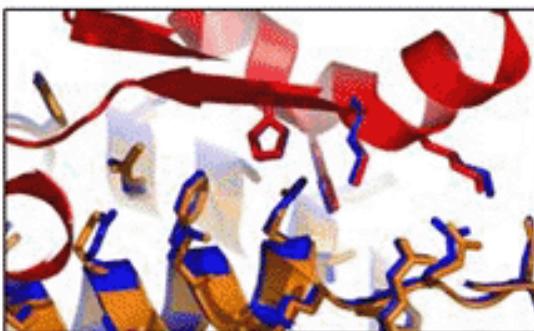
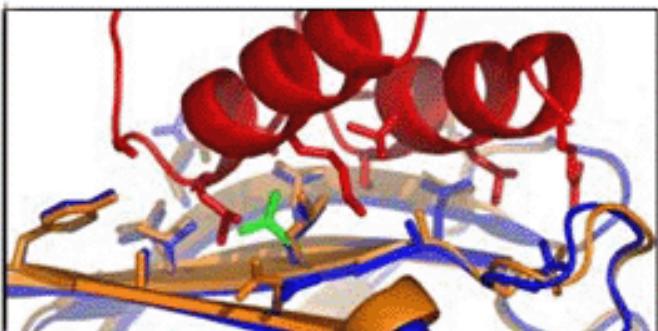
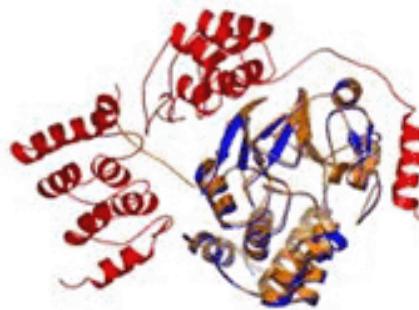
Target 12: Dockerin-cohesin complex



Target 15: Immunity protein-colicin tRNA_{Asp} complex



Target 14: Myosin phosphatase-targeting subunit—protein S/T phosphatase



The Blue Gene Project

Em dezembro de 1999, IBM anunciou um projeto orçado em \$100 milhões de dólares em 5 anos.

Objetivo: Construir um computador massivamente paralelo para ser aplicado no estudo de fenômenos biomoleculares, como protein folding.

IBM Research – Blue Gene supercomputadores que operam da ordem de 478 TFlops (continuado) e 596 TFlops no pique!

Top 500 parallel Computers, Nov/2007: <http://www.top500.org/>

Nosso Trabalho Nesta Área

Nos últimos anos temos trabalhado em:

- Predição de Estruturas de Proteínas
- Reconstrução de Redes de Genes
(em geral e relacionadas a doenças)
- Interação entre Proteínas e Genes
- SNPs (single nucleotide polymorphism)
e Haplotyping *Novo!*

Predição de Estruturas de Proteínas



**Combining Few Neural Networks for
Effective Secondary Structure Prediction**

*K. Guimarães, J. Melo e G. Cavalcanti
Bethesda, USA, Março 2003*



**Protein Secondary Structure Prediction:
Efficient Neural Network and Feature
Extraction Approaches**

*J. Melo, G. Cavalcanti e K. Guimarães
IEE Electronics Letters, 2004*

Reconstrução de Redes de Genes

IJCNN
2005

**A Simpler Bayesian Network Model for
Genetic Regulatory Network Inference**
Gustavo Bastos and Katia S. Guimarães
Montreal, CA, Agosto 2005



**Analyzing the Effect of Prior Knowledge in
Genetic Regulatory Network Inference**
Gustavo Bastos and Katia S. Guimarães
Dezembro 2005

Interação entre Proteína e Genes



ALGORITHMS FOR
MOLECULAR BIOLOGY

**Decomposition of overlapping protein complexes:
A graph theoretical method for analyzing static and
dynamic protein associations**

E Zotenko, K S Guimaraes, R Jothi, T Przytycka

Abril 2006



**Predicting domain-domain interactions
using a parsimony approach**

K S Guimaraes, R Jothi, E Zotenko, T Przytycka

Novembro 2006

Interação entre Proteína e Genes



Interrogating domain-domain interactions with parsimony based approaches

K S Guimarães and T Przytycka

Março 2008

SNPs e Haplotyping

... ataggtcc**C**tatttcgcgc**C**gtatacacacggg**A**ctata ... → **CCA**
... ataggtcc**G**tatttcgcgc**C**gtatacacacggg**T**ctata ... → **GCT**
... ataggtcc**C**tatttcgcgc**C**gtatacacacggg**T**ctata ... → **CCT**

0.1% diferença de um indivíduo para outro.

80% das variações em SNPs

Ponto frequente de caracterização de doenças

Abordagens Combinatórias e Estatísticos

SNPs e Haplotyping

Araújo, FRB ; GUIMARÃES, K. S. . A Case-Control Study of Non-parametric Approaches for Detecting SNP-SNP Interactions. In: XXX International Conference of the Chilean Computer Science Society, 2011, Curicó, Chile. Proc. of the XXX International Conference of the Chilean Computer Science Society, 2011.

SNPs e Haplotyping

Rosa, Rogério S. ; GUIMARÃES, K. S.. Insights on Haplotype Inference on Large Genotype Datasets. In: Brazilian Symposium on Bioinformatics (BSB) 2010, 2010, Búzios, RJ, BRAZIL. Lecture Notes in Bioinformatics. Berlin, Alemanha : Springer, 2010. v. 6268. p. 47-58.

Rosa, Rogério S. ; Santos, R.H.S. ; GUIMARÃES, K. S.. Accurate Prediction of Error in Haplotype Inference Methods through Neural Networks. In: IJCNN - Int. Joint Conference on Neural Nets, 2012, Brisbane. Proc. of the IJCNN 2012. Piscataway, NJ, USA : IEEE Publishing, 2012.

Micro Array Clustering

Monteiro, Carla C.R.R. ; GUIMARÃES, K. S.. Logistic Biclustering Models for Protein Network Inference. In: IEEE International Conference on Bioinformatics and Bioengineering, 2009, Taichung, Taiwan. BIBE 2009 Proceedings, 2009.