



# Alinhamento de Seqüências Biológicas

Profª Drª Silvana Giuliatti  
 Departamento de Genética – FMRP/USP  
 silvana@rge.fmrp.usp.br

Grupo de Bioinformática  
 Depto. de Genética - FMRP/USP

1

## O que se compara?

- A comparação de seqüências de **DNA**, RNA e proteínas é uma das bases da bioinformática.

NC1=NC=CC(=O)N1  
**Citosina**

CC1=CNC(=O)NC1=O  
**Timina**

NC1=NC2=C(N1)N=CN=C2N  
**Guanina**



NC1=NC=NC2=C1N=CN2  
**Adenina**

A C G T

nucleotídeos

Grupo de Bioinformática  
 Depto. de Genética - FMRP/USP

2

## O que se compara?

- A comparação de seqüências de DNA, **RNA** e proteínas é uma das bases da bioinformática.

NC1=NC=CC(=O)N1  
**Citosina**

O=C1NC=CC(=O)N1  
**Uracila**



NC1=NC2=C(N1)N=CN=C2N  
**Guanina**

NC1=NC=NC2=C1N=CN2  
**Adenina**

A C G U

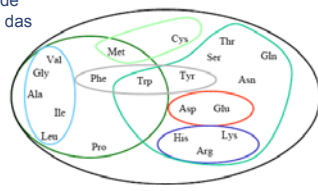
Grupo de Bioinformática  
 Depto. de Genética - FMRP/USP

3

## O que se compara?

- A comparação de seqüências de DNA, RNA e **proteínas** é uma das bases da bioinformática.



G A S T C V I L P F Y M W N Q H D E K R

Aminoácidos

Grupo de Bioinformática  
 Depto. de Genética - FMRP/USP

4



**Histórico**

- 1970**
  - Matrizes de Pontos
  - Programação Dinâmica
    - Alinhamento Global
- 1981**
  - Programação Dinâmica
    - Alinhamento local
- 1988**
  - Alinhamento com Banco de dados
    - FASTA
- 1990**
  - BLAST
- 1994**
  - Alinhamento múltiplo
    - CLUSTAL

Profª Drª Silvana Giulitti 9

**Matrizes de Pontos**

- Descrito pela primeira vez por Gibbs e McIntyre (1970).
- Método usado para alinhar duas seqüências
- Todas as regiões possíveis de serem alinhadas são encontradas
- Não permite a inclusão de gaps

Profª Drª Silvana Giulitti 10

**Matrizes de Pontos**

- Colocar uma seqüência em uma linha e a outra seqüência em uma coluna
- Colocar um ponto em todas as posições onde houver similaridade
- Diagonais revelam a similaridade entre as duas seqüências

Profª Drª Silvana Giulitti 11

**Matrizes de Pontos**

	A	C	G	T	A	C	G	T	A	C
G			o				o			
G			o				o			
T				o				o		
T				o				o		
A	o				o				o	
C		o				o				o
G			o				o			
G			o				o			
T				o				o		
C		o				o				o

Profª Drª Silvana Giulitti 12

**Programas para Matrizes de Pontos**

- **Dotlet**
  - [www.isrec.isb-sib.ch/java/dotlet/Dotlet.html](http://www.isrec.isb-sib.ch/java/dotlet/Dotlet.html)
  - Sequências curtas: até 10.000 caracteres
- **Dotter**
  - [www.cgr.ki.se/cgr/groups/sonnhammer/Dotter.html](http://www.cgr.ki.se/cgr/groups/sonnhammer/Dotter.html)
  - Sequências até 100.000 caracteres
- **EMBOSS Dottup, Dotmatcher**
  - [www.emboss.org](http://www.emboss.org)
  - Sequências maiores de 100.000 caracteres

Profª Drª Silvana Giulitti 13

**Programação Dinâmica**

- “Problema do caixeiro viajante”
- Procura por todas as soluções possíveis
  - Encontra a solução ótima



Profª Drª Silvana Giulitti 14

**Alinhamento com Programação Dinâmica**

- Alinhamento de pares de seqüências
  - Global e Local
  - Pode considerar lacunas (gaps) ao longo do alinhamento
- Encontrar o melhor alinhamento possível – alinhamento ótimo
  - Pode existir mais de um alinhamento ótimo
- **Limitação:** pode se tornar lento dependendo do tamanho das seqüências

Profª Drª Silvana Giulitti 15

**Score**

- Considerar as seqüências ACGGACT e ATCGGATCT

A	-	C	G	G	A	-	C	T	A	-	C	-	G	G	-	A	C	T
A	T	C	G	G	A	T	C	T	A	T	C	G	G	A	T	-	C	T

Qual destes é o melhor alinhamento?

Profª Drª Silvana Giulitti 16

**Score (Pontuação)**

- **Score ou Pontuação:** Medida pela qual os alinhamentos são quantificados
- Considere o seguinte esquema simples de pontuação

+1 para igualdade (match)
-1 para desigualdade (mismatch)
-2 para lacunas (gap)

Profª Drª Silvana Giuliatti 17

**Score (Pontuação)**

**Qual é o melhor alinhamento?**

Alinhamento 1

A	-	C	G	G	A	-	C	T
A	T	C	G	G	A	T	C	T

+1 -2 +1 +1 +1 -2 +1 +1 = +2

Alinhamento 2

A	-	C	-	G	G	-	A	C	T
A	T	C	G	G	A	T	-	C	T

+1 -2 +1 -2 +1 -1 -2 -2 +1 +1 = -4

**Melhor Alinhamento: Alinhamento 1**

Profª Drª Silvana Giuliatti 18

**Programação Dinâmica**

- **Alinhamento Global** - Algoritmo de Needleman-Wunsch (1970).
- **Alinhamento Local** – Smith-Waterman (1981)
  - Modificação do Algoritmo de Needleman-Wunsch
- Gerar uma matriz
- Encontrar todos os possíveis alinhamentos

Profª Drª Silvana Giuliatti 19

**Programação Dinâmica**

**Alinhamento Global**

	0	1	2	3	4
0 -	0	-2	-4	-6	-8
1 C	-2	-1	-3	-5	-7
2 T	-4	-3	0	-2	-4
3 T	-6	-5	-2	-1	-3
4 A	-8	-7	-4	-1	0
5 G	-10	-7	-6	-3	-2
6 A	-12	-9	-8	-5	-2

**- G T A - A**  
| | |  
**CTTAGA**  
-2 -1 +1 -2 +1 = -2

Profª Drª Silvana Giuliatti 20

**Programação Dinâmica**

### Alinhamento Local

		0	1	2	3	4
	-		G	T	A	A
0	-	0	0	0	0	0
1	C	0	0	0	0	0
2	T	0	0	1	0	0
3	T	0	0	1	0	0
4	A	0	0	0	2	1
5	G	0	1	0	0	1
6	A	0	0	0	1	1

Profª Drª Silvana Giulitti 21

**Matrizes de Substituição**

- Sistema de pontuação biologicamente relevantes
- Para produzir alinhamentos biologicamente significativos
- Matrizes
  - PAM
  - BLOSUM
  - Aminoácidos
  - Nucleotídeos

Profª Drª Silvana Giulitti 22

**Matrizes de Substituição**

### Matrizes de Substituição PAM

- PAM – Percent Accept Mutation
- Desenvolvida por Margaret Dayhoff *et al* (1978)
- Considerou seqüências de aminoácidos com pelo menos 85% de similaridade
- As substituições de aminoácidos foram estimados
  - 1572 mudanças em 71 grupos de seqüências de proteínas.
- Matriz mais utilizada PAM 250



Profª Drª Silvana Giulitti 23

**Matrizes de Substituição**

### Matriz PAM 250

	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W
C	12																			
S	2	13																		
T	2	3	11																	
P	3	6	6	11																
A	2	1	1	2	11															
G	3	1	0	1	1	15														
N	4	1	0	1	0	0	2													
D	5	0	0	1	0	1	2	4												
E	5	0	0	1	0	1	5	4												
Q	5	1	1	0	0	1	2	2	4											
H	3	1	1	0	1	2	2	1	1	3	6									
R	4	0	1	0	2	3	0	1	1	1	2	6								
K	5	0	0	1	1	2	1	0	0	1	0	3	5							
M	3	2	1	2	1	3	2	3	2	1	2	0	0	6						
I	2	1	0	2	1	3	2	3	2	2	2	2	2	5						
L	6	3	2	3	2	4	5	4	3	2	5	3	3	4	2	6				
V	2	1	0	1	0	1	2	2	2	2	2	2	2	2	4	2	4			
F	4	3	3	5	4	5	4	6	5	5	2	4	5	0	1	3	1	9		
Y	0	3	3	5	3	5	2	4	4	4	0	4	4	2	1	1	2	7	10	
W	2	3	3	5	6	7	4	7	7	5	3	2	3	4	3	7	6	0	0	17



Profª Drª Silvana Giulitti 24


**Matrizes de Substituição**


**Matriz PAM**

- **Valor Zero** → frequência de substituição entre dois aminoácidos é esperada ao acaso
- **Valor menor que zero** → frequência é menor que a esperada. Substituição de dois aminoácidos ao acaso.
- **Valor maior que Zero** → frequência maior que a esperada. Substituição não é ao acaso. Indica maior probabilidade de relação com ancestral



Grupo de Bioinformática  
Depto. de Genética - FMUSP/USP
Profª Drª Silvana Giulitti


**Matrizes de Substituição**


**Matriz de Substituição BLOSUM**

- **BLOSUM – Blocks Substitution Matrix**
- Desenvolvidas por Henikoff e Henikoff, (1992)
- Aminoácidos são organizados em blocos
- Utilizou-se 500 famílias de proteínas
- Matriz mais utilizada BLOSUM62



Grupo de Bioinformática  
Depto. de Genética - FMUSP/USP
Profª Drª Silvana Giulitti


**Matrizes de Substituição**


**Matriz BLOSUM62**

	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W
C	9																			
S	-1	4																		
T	-1	1	5																	
P	-3	-1	-1	7																
A	0	1	0	-1	4															
G	-2	0	-2	-2	0	6														
N	-3	1	0	-2	-2	0	6													
D	-3	0	-1	-1	-2	-1	1	6												
E	-4	0	-1	-1	-1	-2	0	2	5											
Q	-3	0	-1	-1	-1	-2	0	0	2	5										
H	-3	-1	-2	-2	-2	-2	-1	-1	0	0	8									
R	-3	-1	-1	-3	-1	-2	0	2	0	3	0	5								
Y	-3	0	-1	-1	-1	-2	0	-1	1	1	-1	2	5							
M	-1	-1	-2	-1	-3	-2	-3	-2	0	-2	-1	-1	5							
I	-1	-2	-1	-3	-1	-4	-3	-3	-3	-3	-3	-3	1	4						
L	-1	-2	-1	-3	-1	-4	-3	-4	-3	-2	-3	-2	2	2	4					
V	-1	-2	0	-2	0	-3	-3	-3	-2	-3	-3	-2	1	3	1	4				
F	-2	-2	-4	-2	-3	-3	-3	-3	-3	-3	-3	-3	0	0	0	-1	6			
Y	-3	-2	-3	-2	-3	-3	-3	-3	-1	2	-3	-3	-1	-1	-1	-3	3	7		
W	-3	-2	-4	-3	-2	-4	-3	-3	-3	-3	-3	-3	-1	-1	-2	-3	1	2	11	



Grupo de Bioinformática  
Depto. de Genética - FMUSP/USP
Profª Drª Silvana Giulitti


**Matrizes de Substituição**


**Matriz BLOSUM**

- **Valor Zero** → probabilidade de substituição entre dois aminoácidos iguais
- **Valor menor que zero** → maior probabilidade de substituição de dois aminoácidos ser ao acaso
- **Valor maior que Zero** → maior probabilidade de substituição entre dois aminoácidos não ser por acaso. Indica maior probabilidade de relação com ancestral

Grupo de Bioinformática  
Depto. de Genética - FMUSP/USP
Profª Drª Silvana Giulitti



 **Matrizes de Substituição** 

**PAM** **X** **BLOSUM**

- Calculadas de alinhamentos globais
- Sequências utilizadas com pelo menos 85% de similaridade
- As matrizes são extrapolações da PAM 1
- Usada para traçar origens da Evolução das proteínas



- Calculadas de alinhamentos locais
- Pode-se selecionar a similaridade entre as sequências
- Cada matriz é gerada do resultado de uma análise
- Usadas para encontrar domínios conservados

Profª Drª Silvana Giuliatti 29

 **Métodos de Palavras** 



- Alinham sequências mais rapidamente.
- Procuram por partes curtas idênticas (palavras ou k-tuplas).
- Pesquisas em bancos de dados: FASTA e BLAST
- Seguem um método heurístico.

Profª Drª Silvana Giuliatti 30

 **FASTA** 

- Desenvolvido por Pearson e Lipman (1988).
- Uma sequência de proteína ou DNA com todas as sequências num **banco de dados**.
- Apresenta os alinhamentos locais da sequência analisada com as sequências do banco.



Profª Drª Silvana Giuliatti 31

 **FASTA** 

- Algoritmo escrito em linguagem C.
- Mais lento que BLAST.
- Procura por um número k de consecutivas letras (aminoácidos ou nucleotídeos): palavras ou k-tuplas.

Profª Drª Silvana Giuliatti 32



 **FASTA** 



Grupo de Bioinformática  
Depto. de Genética - FMUSP

A sequência de entrada deve estar no formato FASTA.

O algoritmo pode ser dividido em 4 etapas:

- a) **seleção das 10 melhores regiões.**
- b) **re-classificação das 10 melhores regiões.**
- c) **seleção das seqüências mais semelhantes.**
- d) **alinhamento das seqüências selecionadas.**

Profª Drª Silvana Giuliatti 33



 **Métodos de Palavras**  
**BLAST** 

Grupo de Bioinformática  
Depto. de Genética - FMUSP

Basic Local Alignment Sequence Tool

- Alinhamento de uma seqüência de proteína ou DNA com todas as seqüências num banco de dados.
- Apresenta os alinhamentos locais da seqüência analisada com as seqüências do banco.
- Mais rápido que FASTA.
- Algoritmo escrito em linguagem C.
- Procura por um número k de consecutivas letras (aminoácidos ou nucleotídeos): palavras ou k-tuplas.



Profª Drª Silvana Giuliatti 34

 **BLAST** 

Grupo de Bioinformática  
Depto. de Genética - FMUSP

- Procura por palavras que são mais significantes
- A significância é incorporada ao algoritmo através de matrizes de pontuação.
- Buscar por identidades de comprimento k:
  - **11 para nucleotídeos**
  - **3 para aminoácidos**

Profª Drª Silvana Giuliatti 35

 **BLAST** 

Grupo de Bioinformática  
Depto. de Genética - FMUSP

- O algoritmo pode ser dividido em 4 etapas:
  - a) **montagem da lista de palavras.**
  - b) **procura pelas palavras em cada seqüência do banco.**
  - c) **extensão.**
  - d) **alinhamento das seqüências.**

Profª Drª Silvana Giuliatti 36

**BLAST**

**d) Alinhamento das seqüências**

Determina se cada HSP é estatisticamente significativa.

**Depois de determinar se HSP é estatisticamente significativa, faz o alinhamento dos melhores segmentos.**

Profª Drª Silvana Giuliatti 37

**EVALUE**

O mais usado score é o **Evalue**: proporciona uma estimativa do número de falsos positivos esperados.

**Interpretação do Valor Esperado: Evalue**

- $E < 10^{-100}$  ⇒ valor muito baixo. Genes homólogos ou idênticos.
- $E < 10^{-3}$  ⇒ valor moderado. Genes podem estar relacionados.
- $E > 1$  ⇒ valor alto. Prováveis genes sem relação.
- $0,5 < E < 1$  ⇒ Região duvidosa - "Twilight zone"

**Twilight zone:** nessa região, nada é garantido sobre o significado das similaridades observadas. Homologia ou não, nunca é garantida nessa área.

Profª Drª Silvana Giuliatti 38

**Alinhamento Múltiplo**

- O alinhamento múltiplo de seqüências identifica resíduos ou regiões conservadas ou equivalentes em estruturas.

```

Bradi3ves  |  GDGDFEETKPSRQVYQVAGTSLGQ@HTFYRPLRSTLQDELPLPPLRGLG  |  366
DnaB13com  |  DLGSRFETKPSRQVYQVAGTSLGQ@HTFYRPLRSTLQDELPLPPLRGLG  |  363
Egub13E22  |  GDHDELLTQVPSRQVYQVAGTSLGQ@HTFYRPLRSTLQDELPLPPLRGLG  |  365
PanaR3seq  |  GSDKALQVPSRQVYQVAGTSLGQ@HTFYRPLRSTLQDELPLPPLRGLG  |  369
DnaB13888  |  GRDGLKPLRPSRQVYQVAGTSLGQ@HTFYRPLRSTLQDELPLPPLRGLG  |  377
ActA1Arg1  |  GFRDQVPSRQVYQVAGTSLGQ@HTFYRPLRSTLQDELPLPPLRGLG  |  406
  
```

Profª Drª Silvana Giuliatti 39

**CLUSTALW**

- Método heurístico
- Rápido e eficiente.
- Faz alinhamento progressivo dos perfis e seqüências mais distantes
- O mais usado: ClustalW (Thompson et al, 1994)

Profª Drª Silvana Giuliatti 40

**CLUSTALW**

- **Algoritmo de 3 etapas:**
  - Alinhamento em pares de todas as sequências para determinar similaridade entre elas.
  - Definir a ordem do alinhamento progressivo baseado na similaridade.
  - Construir o alinhamento múltiplo baseado na ordem definida.

Qual a ordem do alinhamento?

Profª Drª Silvana Giuliatti 41

**CLUSTALW**

**Etapa 1:** Alinhamento em pares de todas sequências para determinar similaridade entre elas.

- Usa método de alinhamento (global) de pares de sequências
- Usa matriz de substituição e penalidade por gaps.

Profª Drª Silvana Giuliatti 42

**CLUSTALW**

- Usa os alinhamentos em pares para calcular uma “distância genética” entre todos os pares de sequências.
- Constrói uma matriz de valores de distâncias.



Profª Drª Silvana Giuliatti 43

**CLUSTALW**

**Etapa 2:** Definir a ordem do alinhamento progressivo baseado na similaridade.

- Definir as sequências mais próximas: árvore de similaridade.
  - Usa matriz de distâncias para calcular a árvore.
  - Método de junção por vizinhos (neighbor-joining)

Profª Drª Silvana Giuliatti 44

 **CLUSTALW** 



Grupo de Bioinformática  
Instituto de Genética - UNESP/USP

Faculdade de Medicina  
UNESP - Ribeirão Preto

**Etapa 3:** Construir o alinhamento múltiplo baseado na ordem definida.

- Combinar os alinhamentos começando com os grupos mais próximos para os mais distantes

Profª Drª Silvana Giuliatti 45

 **Referências** 

Grupo de Bioinformática  
Instituto de Genética - UNESP/USP

Faculdade de Medicina  
UNESP - Ribeirão Preto

- ~ http://www.sxc.hu
- ~ http://creative.gettyimages.com/source/home/home.aspx
- 📖 Mount, D. W., Bioinformatics - Sequence and Genome Analysis, ed. CSHL, 2ª edição.
- 📖 Gibas, C e Jambeck, P., Desenvolvendo a Bioinformática. Ed. Campus.

Profª Drª Silvana Giuliatti 46