

Algoritmo BLAST

Aluno: Marllus de Melo Lustosa

Sumário:

1. Alinhamento de sequências
 - 1.1. Alinhamento global
 - 1.2. Alinha local
 - 1.3. Global x Local
 - 1.4. Alinhamento ótimo x heurístico
 - 1.5. Ferramentas de alinhamento
2. FASTA vs BLAST
3. BLAST
 - 3.1. Funcionamento do algoritmo
 - 3.1.1. Semeadura
 - 3.1.2. Extensão
 - 3.1.3. Avaliação
 - 3.2. Família de programas BLAST–NCBI

1. Alinhamento de sequências

- Encontrar um grau de similaridade entre sequências de nucleotídeos ou proteínas.
 - Definir sequências homólogas;
 - Definir fragmentos similares entre sequências;
 - Determinar características entre sequências;

1.1. Alinhamento global

- Sequências são alinhadas de ponta a ponta;
- Pode incluir grandes pedaços com baixa similaridade;
- Útil para comparar sequências cujas semelhanças sejam esperadas em toda a sua extensão;

LGPSSKQTGKGS-SRIWDN
LN-ITKSAGKGAIMRLGDA

Fig1. Exemplo de alinhamento global entre duas sequências. [1]

1.2. Alinhamento local

- São alinhados um ou mais segmentos com alta similaridade entre as sequências;
- Útil quando não se tem nenhum conhecimento sobre a semelhança entre as sequências a comparar;

```
                tccCAGTTATGTCAGgggacacgagcatgcagagac
                |||||
aattgccgccgtcgttttcagCAGTTATGTCAGatc
```

Fig2. Exemplo de alinhamento local entre duas sequências. [2]

1.3. Global x Local

Global: As sequências são alinhadas de ponta a ponta;

Local: Pedacos das sequências é que são comparados;

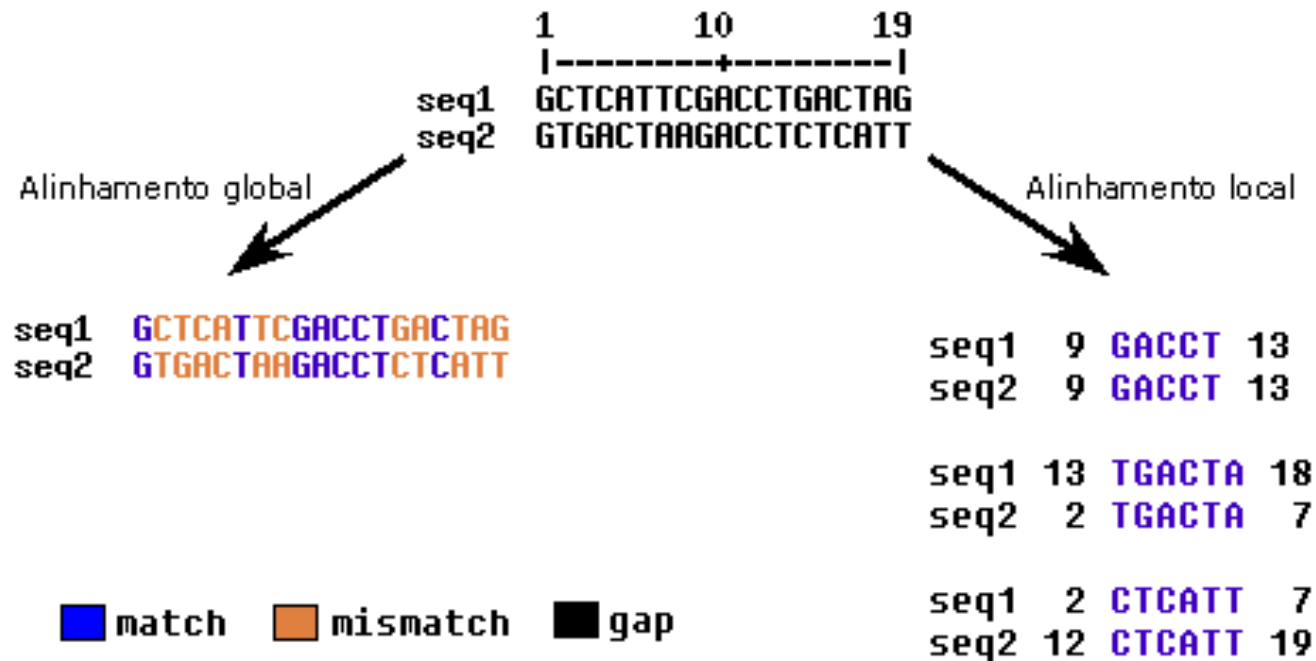


Fig3. Exemplo de alinhamento global e local entre duas sequências. [3]

1.4. Alinhamento ótimo x heurístico

heurística -- do dicionário *Houaiss*

Acepções

| substantivo feminino

1 arte de inventar, de fazer descobertas; ciência que tem por objeto a descoberta dos fatos

1.1 Rubrica: história.

ramo da História voltado à pesquisa de fontes e documentos

1.2 Rubrica: informática.

método de investigação baseado na aproximação progressiva de um dado problema

1.3 Rubrica: pedagogia.

método educacional que consiste em fazer descobrir pelo aluno o que se lhe quer ensinar

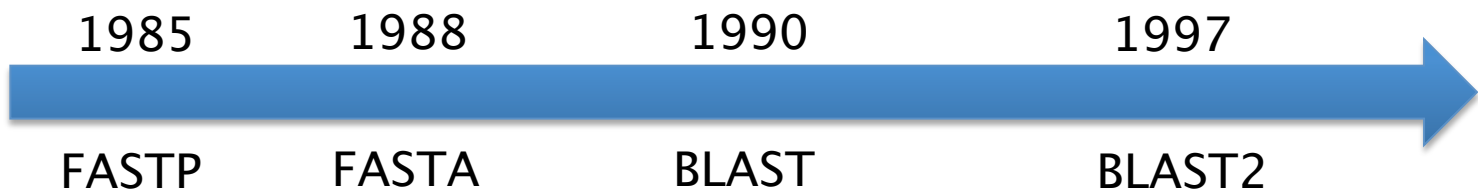
- Alinhamento **ótimo**: produz o melhor resultado computacionalmente possível;
- Alinhamento **heurístico**: produz um resultado o mais próximo possível do resultado ótimo, mas, principalmente, produz um resultado de maneira muito veloz;

1.5. Ferramentas de alinhamento

Programa	Tipo de Alinhamento	Precisão do Alinhamento	Número de seqüências a serem alinhadas
BLAST2Sequences	Local	Heurístico	2
SWAT (Smith–Waterman)	Local	Ótimo	2
ClustalW	Global	Heurístico	N
Multalin	Global	Heurístico	N
Needleman–Wunsch	Global	Ótimo	2

Tab1. Exemplo de alinhamento global e local entre duas seqüências. [3]

2. FASTA vs BLAST



- BLAST e FASTA são algoritmos de alinhamento local;
- BLAST é mais rápido que o FASTA;
- BLAST é mais preciso que o FASTA;
- BLAST é mais versátil e mais amplamente utilizado que o o FASTA;
- Partem da ideia básica: Um bom alinhamento contém subsequências de identidade absoluta (pequenas palavras de similaridade exata) [5].

3. BLAST

- Basic Local Alignment Search Tool;
- Ferramenta de alinhamento mais utilizada no mundo;
- O artigo onde a ferramenta foi publicada é o mais citado da história das ciências biológicas;
- É um algoritmo de alinhamento simples, heurístico e local;
- Alinha um sequência de entrada contra uma base de dados desejada;

3. BLAST

Programa	Seqüência de entrada	Tipo de seqüência alvo
blastp	proteína	proteína
blastn	nucleotídeo	nucleotídeo
blastx	nucleotídeo traduzido	proteína
Tblastn	proteína	nucleotídeo traduzido
Tblastx	nucleotídeo traduzido	nucleotídeo traduzido

Tab2. Família de programas BLAST. [4 – Adaptada]

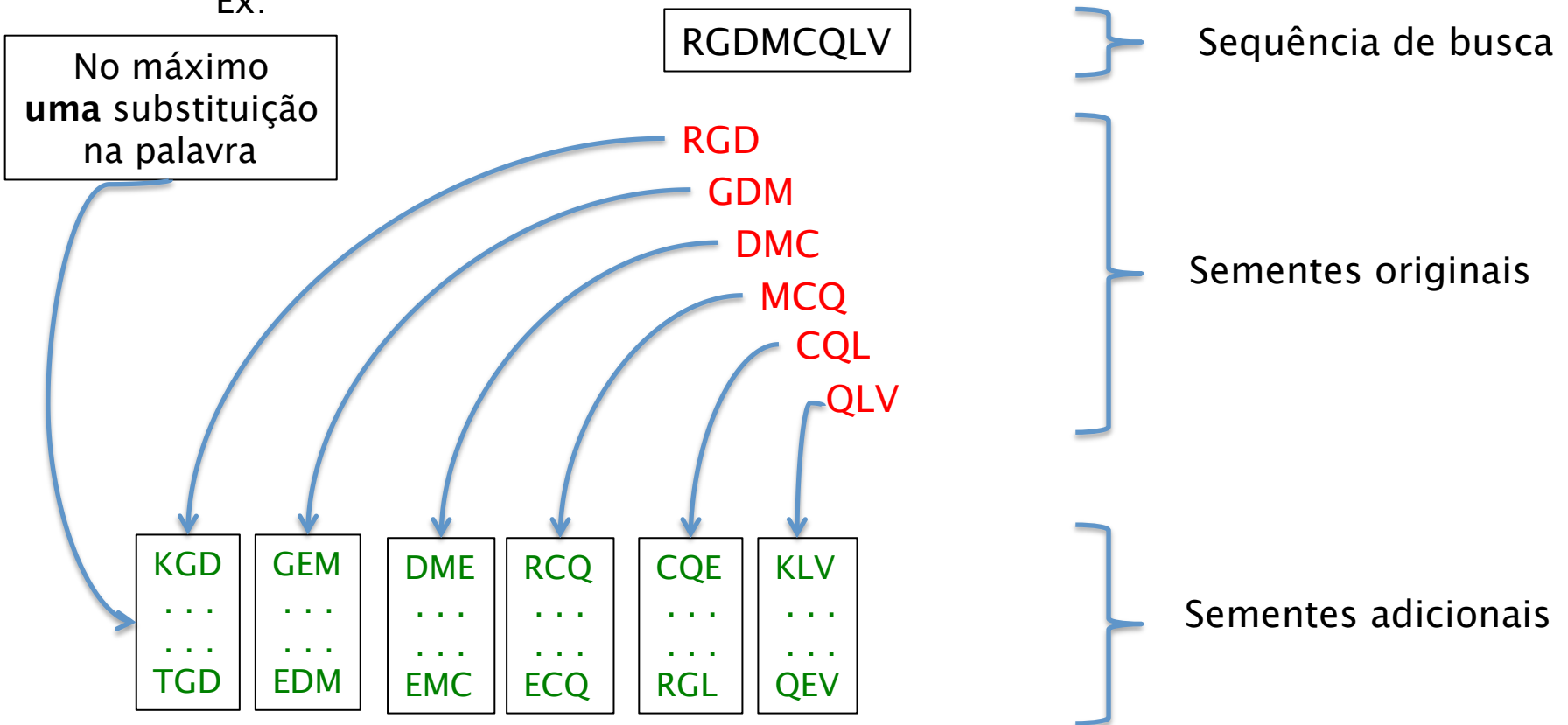
3.1. Funcionamento do algoritmo

- Consiste em 3 etapas heurísticas:
 - **Semeadura;**
 - Separa a sequência de busca em palavras;
 - Identifica onde começa o alinhamento;
 - **Extensão;**
 - Extende o alinhamento das sementes;
 - **Avaliação;**
 - Determina quais alinhamentos são significantes;

3.1.1. Semeadura

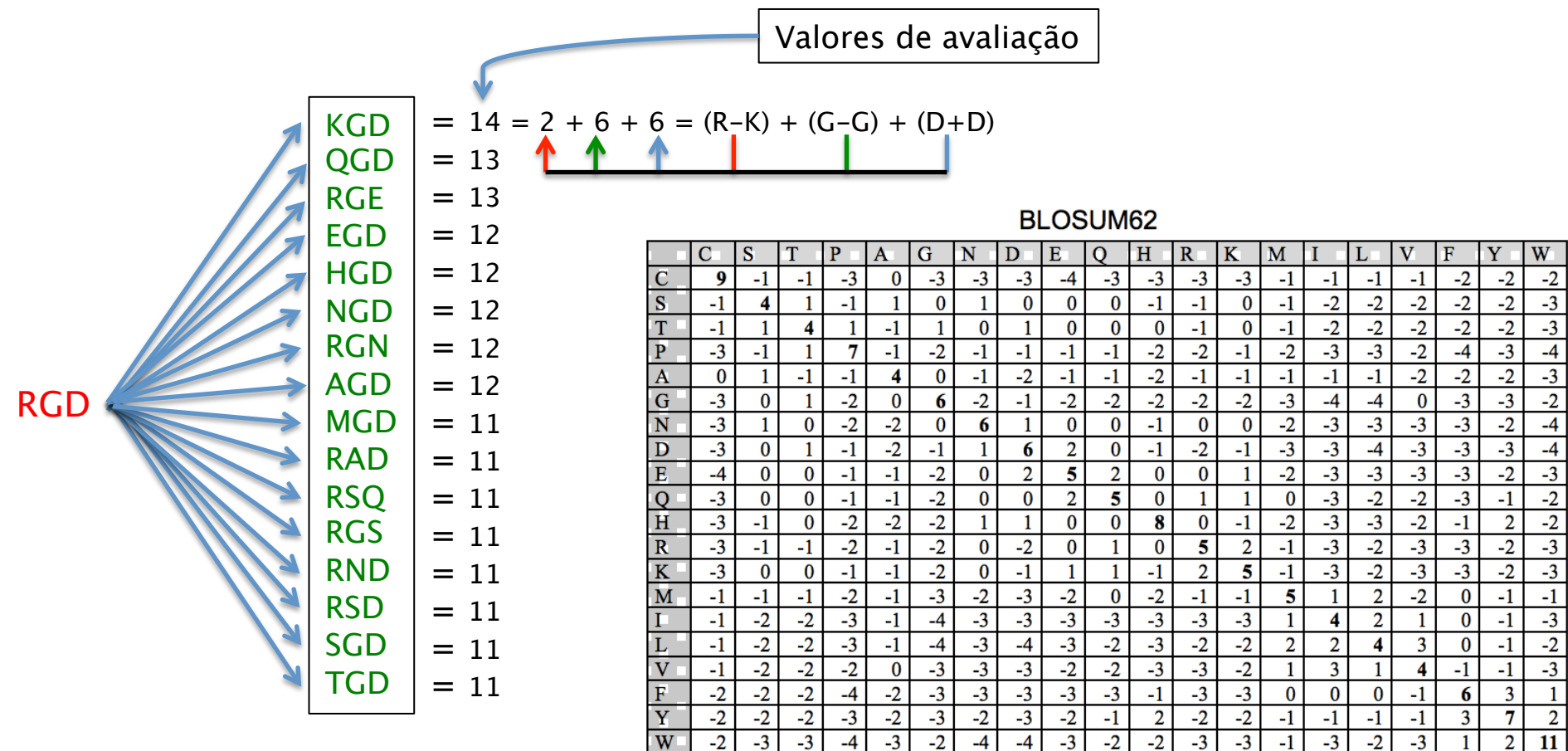
- Dada uma sequência de entrada, identifique todas as sequências de um tamanho específico (sementes).

Ex:



3.1.1. Semeadura

- Comparar as sementes adicionais com a semente original correspondente, utilizando uma matriz de substituição (recomenda-se a matriz BLOSUM62 [6]).



3.1.1. Semeadura

- Definir um valor mínimo de avaliação na seleção das sementes adicionais;
 - Padrão BLAST, em geral, é utilizado o valor 12 para palavras de tamanho 3;

KGD	= 14	}	Palavras válidas
QGD	= 13		
RGE	= 13		
EGD	= 12		
HGD	= 12		
NGD	= 12		
RGN	= 12		
AGD	= 12	}	Palavras que serão excluídas do conjunto das sementes válidas
MGD	= 11		
RAD	= 11		
RSQ	= 11		
RGS	= 11		
RND	= 11		
RSD	= 11		
SGD	= 11		
TGD	= 11		

3.1.1. Semeadura

- Conjunto de sementes válidas para a busca no banco de dados:

Sementes originais + sementes adicionais

RGD



Semente original

KGD
QGD
RGE
EGD
HGD
NGD
RGN
AGD



Sementes adicionais

3.1.1. Semeadura

- Realizar a busca pelas sementes no banco de dados (prioridade para sementes originais);

Ex:

Sementes de busca:	GDM CQL
Sequência encontrada no BD:	E G D M K C Q L W

3.1.2. Extensão

- Extender o alinhamento das sementes;

Ex:

Sementes de busca:	← RGDM – CQLV →
Sequência encontrada no BD:	EGDMKCQLW

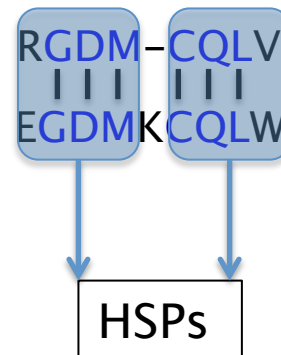
- Extende cada semente para a direita e para esquerda, considerando os seguintes critérios [7]:
 - A pontuação (*score*) da semente for maior que um valor T ;
 - Possuir outra semente a uma certa distância máxima entre elas;
 - O *score* da extensão com *gaps* tem pontuação normalizada de pelo menos S_g *bits*;
- Muitas vezes é necessário adicionar *gaps* (buracos) para corrigir o alinhamento;
 - *Gaps* são vistos pelo BLAST como penalidades;
 - Quanto menos *gaps*, melhor o alinhamento;

3.1.2. Extensão

- Extender o alinhamento das sementes;

Ex:

Sementes de busca:
Sequência encontrada no BD:



- HSPs (*High-scoring Segment Pair*): São alinhamentos locais que atingem os *scores* mais altos em uma busca;
- MSPs (*Maximal-scoring Segment Pair*): São os maiores HSPs encontrados na busca;

3.1.3. Avaliação

– **Score bruto:**

$$S = \text{soma}(\text{matches}) - \text{soma}(\text{mismatches}) - \text{soma}(\text{penalidades de gap})$$

– **Score normalizado (Bit score):**

$$S' = (\lambda S - \ln K) / \ln 2$$

– **E-value** (probabilidade de alinhamentos terem ocorrido ao acaso [2]):

$$E(S) = Kmne^{-\lambda S}$$

ou

$$E(S') = mn2^{-S'}$$

legenda

m = Tamanho do banco de dados
n = Tamanho da sequência de entrada
 λ = Escala da matriz de *scores*
K = Escala do tamanho do espaço de busca

Penalidades de *gap*:

(*gap open + gap extension*)

Gap open: é definido um valor

Gap extension: é definido um valor

3.1.3. Avaliação

Sementes de busca: R**GDM**–**CQL**V
 | | | | |
Sequência encontrada no BD: E**GDM**K**CQL**W

Definir *Gap Costs*:

Gap open: 5

Gap extension: 2

De acordo com [8], valores de:

Match/mismatches = 1/-3 => Sequências 99% conservadas

***Match/mismatches* = 1/-2 => Sequências 95% conservadas**

Match/mismatches = 1/-1 => Sequências 75% conservadas

Matches: $6 * 1 = 6$

Mismatches: $2 * 2 = 4$

Gap open: $1 * 5 = 5$

Gap extension: $1 * 2 = 2$

Resultado parcial

Similaridade: 6/9 (67%)

$S = 6 - 4 - 5 - 2 = -5$

- ## Resultado do BLAST

Valores calculados na “mão”

Referências

- [1] http://www.bdttd.ucb.br/tede/tde_arquivos/13/TDE-2004-11-12T134348Z-143/Publico/Dissert_Felipe%20Lieberman.pdf
- [2] <http://minerva.ufpel.edu.br/~lmaia.faem/aula4.pdf>
- [3] <http://www.biotechnologia.com.br/revista/bio29/bioinf.pdf>
- [4] http://www.bdttd.ucb.br/tede/tde_arquivos/13/TDE-2004-11-12T134348Z-143/Publico/Dissert_Felipe%20Lieberman.pdf
- [5] http://csb.stanford.edu/class/public/readings/Bioinformatics_I_Lecture6/Altschul_JMB_90_BLAST_Sequence_alignment.pdf
- [6] Henikoff, S. & Henikoff, J.G. (1992) "Amino acid substitution matrices from protein blocks." Proc. Natl. Acad. Sci. USA 89:10915–10919.
- [7] <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC146917/pdf/253389.pdf>
- [8] <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC55814/>

OBRIGADO

Contatos

www.marllus.com

marlluslustosa@gmail.com