

## O programa *BLAST*: guia prático de utilização

ISSN 0102 0110  
Dezembro, 2007

**Empresa Brasileira de Pesquisa Agropecuária  
Embrapa Recursos Genéticos e Biotecnologia  
Ministério da Agricultura, Pecuária e Abastecimento**

## *Documentos* 224

### **O programa *BLAST*: guia prático de utilização**

***Embrapa Recursos Genéticos e Biotecnologia*  
Brasília, DF  
2007**

Exemplares desta edição podem ser adquiridos na

Embrapa Recursos Genéticos e Biotecnologia

Serviço de Atendimento ao Cidadão

Parque Estação Biológica, Av. W/5 Norte (Final) –

Brasília, DF CEP 70770-900 – Caixa Postal 02372 PABX: (61) 448-4600 Fax: (61) 340-3624

<http://www.cenargen.embrapa.br>

e.mail:sac@cenargen.embrapa.br

Comitê de Publicações

Presidente: *Sergio Mauro Folle*

Secretário-Executivo: *Maria da Graça Simões Pires Negrão*

Membros: *Arthur da Silva Mariante*

*Maria de Fátima Batista*

*Maurício Machain Franco*

*Regina Maria Dechechi Carneiro*

*Sueli Correa Marques de Mello*

*Vera Tavares de Campos Carneiro*

Supervisor editorial: *Maria da Graça S. P. Negrão*

Normalização bibliográfica: *Maria Iara Pereira Machado*

Editoração eletrônica: *Daniele Alves de Loiola*

1ª edição

1ª impressão (2007):

**Todos os direitos reservados**

A reprodução não autorizada desta publicação, no todo ou em parte, constitui violação dos direitos autorais (Lei nº 9.610).

**Dados Internacionais de Catalogação na Publicação (CIP)  
Embrapa Recursos Genéticos e Biotecnologia**

P 965 O programa *BLAST*: guia prático de utilização Embrapa Recursos Genéticos e Biotecnologia / Alexandre Morais do Amaral, Marcelo da Silva Reis, Felipe Rodrigues da Silva (editores). -- Brasília, DF: Embrapa Recursos Genéticos e Biotecnologia, 2007. 24 p. -- (Documentos / Embrapa Recursos Genéticos e Biotecnologia, 1676 - 1340; 224).

1. BLAST - ferramenta de busca - similaridade - seqüências biológicas. 2. DNA. 3. Aminoácidos. 4. Embrapa Recursos Genéticos e Biotecnologia. I. Amaral, Alexandre Morais do. II. Reis, Marcelo da Silva. III. Silva, Felipe Rodrigues da. IV. Série.

631.5233 - CDD 21.

## **Editores**

Alexandre Morais do Amaral

Pesquisador da EMBRAPA Recursos Genéticos e Biotecnologia,  
Brasília, DF

Marcelo da Silva Reis

Laboratório de Biotecnologia do Centro APTA Citros "Sylvio Moreira" - IAC,  
Cordeirópolis, SP.

Felipe Rodrigues da Silva

Pesquisador da EMBRAPA Recursos Genéticos e Biotecnologia,  
Brasília, DF

**Programa  
e  
Resumos Expandidos**

14-15 de Junho de 2007

BRASÍLIA – DF

## Sumário

<b>Introdução</b> .....	<b>7</b>
<b>Conceitos</b> .....	<b>8</b>
Formato FASTA .....	10
Programas do BLAST .....	10
"Blast local" .....	11
Banco de Dados .....	12
CDS .....	14
Filtragem ( <i>Filtering</i> ) .....	14
Matriz de substituição .....	14
BLOSUM .....	15
Relatório do BLAST ( <i>report</i> ) .....	15
Identificadores .....	18
Identificadores de seqüências (gi) .....	18
Identificador de proteína (Protein ID) .....	18
Número de acesso .....	18
Alinhamento .....	19
Nucleotídeos .....	19
Aminoácido .....	19
Score .....	20
E-value .....	20
Identidade .....	21
<i>Gap</i> .....	21
Similaridade .....	21
Homologia .....	21
<b>Exemplos de alinhamentos</b> .....	<b>22</b>
Proteína .....	22
Nucleotídeo .....	23
<b>Glossário de termos computacionais</b> .....	<b>23</b>
<b>REFERÊNCIAS</b> .....	<b>24</b>

## Introdução

Com a disponibilidade de grande volume de dados genéticos, gerados principalmente em projetos de seqüenciamento de genomas, e o envolvimento cada vez maior de novos pesquisadores a estes projetos, ferramentas da bioinformática que permitam a interpretação adequada de tais informações têm apresentado grande demanda por parte daqueles que se interessam pelo exame comparativo das seqüências. Sem dúvida, o procedimento mais comum e difundido para a análise de seqüências é a utilização do programa BLAST para identificar rapidamente a existência de seqüências semelhantes entre si. Estas semelhanças podem indicar homologia e permitir a inferência de função. Ou seja, através da similaridade identificada pode-se atribuir - ainda que provisoriamente - uma função à seqüência sem a necessidade de realização imediata de experimento biológico.

O termo “BLAST” é uma sigla/trocadilho em inglês. A palavra *blast* significa “explosão”. O algoritmo, entretanto, recebe este nome por ser uma *Basic Local Alignment Search Tool* (Ferramenta Básica de Busca de Alinhamentos Locais). Como o nome indica, trata-se de uma ferramenta de busca de similaridade entre seqüências biológicas (DNA ou aminoácidos). O algoritmo BLAST e o programa computacional que o implementa foram desenvolvidos, na década de 1980, nos Estados Unidos (ALTSCHUL et al., 1990). Este algoritmo prioriza alinhamentos de locais específicos da seqüência, com respaldo estatístico, em lugar de realizar alinhamentos “globais”, como os realizados por programas de alinhamento de múltiplas seqüências, como o Clustal (CHENNA et al., 2003) (apresentam uma boa revisão de toda a família de programas). Com isto, o programa é capaz de identificar relação entre seqüências que compartilham similaridade, mesmo que esta ocorra somente em algumas regiões isoladas (ALTSCHUL et al., 1990). Tal fato indica que se o recurso utilizado pelo programa fosse o alinhamentos global, ou seja, com o uso imediato do comprimento total da seqüência, menor número de similaridades seriam encontradas, sobretudo no caso de “domínios” e “motivos”. Devido ao seu uso constante, a sigla BLAST é hoje empregada como substantivo e até mesmo verbo (“blaster uma seqüência”).

---

### Saiba Mais (o algoritmo BLAST):

O BLAST é uma derivação do algoritmo Smith e Waterman (1981), que se caracteriza por apresentar a pontuação máxima do alinhamento local de duas seqüências. O algoritmo Smith-Waterman emprega um *método exaustivo* conhecido como *programação dinâmica*, e assim garante encontrar a pontuação máxima. Por sua vez, o algoritmo BLAST emprega uma *heurística* baseada em *palavras*, o que o torna 50 vezes mais rápido. Quando se considera o tamanho dos bancos de dados contra os quais são realizadas buscas atualmente (mais de 83 bilhões de bases de DNA no *GenBank* em dezembro de 2007), é simples entender porque velocidade é importante, mesmo que às custas da perda de *sensibilidade* (PERTSEMLIDIS e FONDON, 2001).

Uma busca utilizando BLAST compara a seqüência de interesse (*query* ou seqüência de entrada), contra um banco de dados (*subject*). A origem da seqüência de interesse assim como a pergunta que se tenta responder irão determinar o banco de dados a ser utilizado. Desta forma, um pesquisador tentando clonar um determinado gene (digamos, a Rubisco de uma planta recentemente descoberta) pode comparar a seqüência

que acabou de ser gerada em um seqüenciador automático contra todas as seqüências de proteínas já descritas em qualquer parte do mundo. Já um outro pesquisador, diante do resultado do seqüenciamento aleatório de uma biblioteca de cDNA, pode optar por “blastar” as seqüências recém-geradas em seu projeto contra um banco de dados formado pelas seqüências anteriormente geradas no projeto para avaliar o grau de redundância do mesmo.

Nos exemplos do parágrafo anterior, o primeiro pesquisador poderia acessar o BLAST em um *site* de acesso público, ou seja, um local que disponibiliza este recurso para qualquer pessoa com acesso à internet. Ainda que exista um número considerável de servidores públicos de BLAST, a imensa maioria dos usuários no mundo realiza suas buscas no NCBI (*National Center for Biotechnology Information*, ou Centro Nacional para Informação Biotecnológica (<http://www.ncbi.nlm.nih.gov>)), pertencente ao governo dos Estados Unidos. O segundo pesquisador, por outro lado, talvez esteja realizando suas buscas em um servidor com acesso restrito aos participantes do projeto (denominado “BLAST local”).

Nas duas situações, o programa irá identificar um ou mais trechos da seqüência de entrada que são semelhantes a porções de seqüências presentes no banco de dados em questão e fornecer um relatório a respeito desta comparação, com várias informações, tais como o nome e descrição das seqüências do banco de referência e a demonstração numérica e gráfica dos pareamentos significativos.

O objetivo desta publicação é permitir para o usuário da ferramenta BLAST o melhor entendimento e compreensão das informações obtidas após a sua utilização, na interpretação dos resultados da pesquisa. Enfim, a finalidade é aumentar a exploração e o entendimento de dados gerados em estudos que utilizam informação genética a partir da análise de seqüências.

## Conceitos

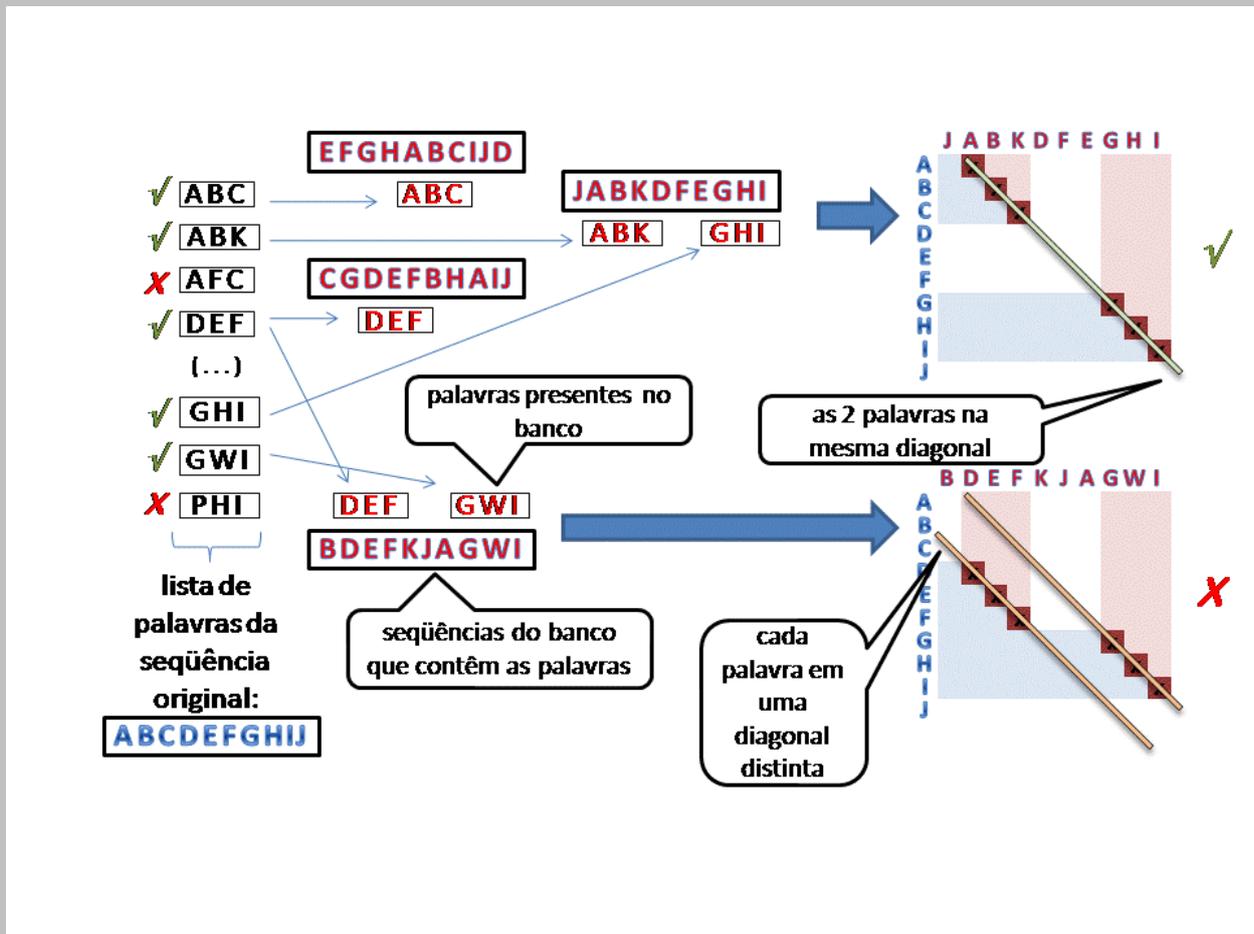
De forma bastante resumida, o algoritmo BLAST pode ser dividido em três estágios básicos (ALTSCHUL et al., 1990):

- (1) Compilação de uma lista de “palavras” de alta pontuação;
- (2) Procura destas palavras no banco de dados;
- (3) Extensão de alinhamentos a partir das palavras encontradas.

### Saiba mais (os três estágios básicos do algoritmo):

1. No primeiro estágio é gerada uma lista de todas as “palavras” encontradas na seqüência que está sendo buscada. “Palavra”, no caso do BLAST, é um trecho de comprimento definido da seqüência. Os tamanhos-padrão de palavras são 3 aminoácidos para seqüências protéicas e 11 nucleotídeos para seqüências nucléicas. No caso de seqüências protéicas, a lista inicial inclui todas as combinações possíveis de 3 aminoácidos consecutivos, como esquematizado à esquerda na Figura 1. Esta lista é refinada através do cálculo da pontuação do alinhamento, sem *gaps* (*i.e.*, apenas comparando palavras inteiras, ou seja, contínuas, sem interrupção), destas palavras contra todas as 8.000 ( $20^3$ ) palavras possíveis. Apenas “palavras” com pontuação acima de um determinado limiar são selecionadas. Pode-se notar que “palavras”

que não ocorrem na seqüência original, mas que pontuam acima do limiar com palavras presentes na seqüência buscada, são incluídas na lista de trechos selecionados. Em média, 50 palavras serão incluídas na lista para cada aminoácido da seqüência original. Em suma, no primeiro estágio, compila-se uma lista de “palavras” relevantes à seqüência buscada.



**Figura 1.** Exemplo de procedimento do programa BLAST para busca de seqüências protéicas, com tentativas iniciais de todas as combinações possíveis de 3 aminoácidos consecutivos.

2. Palavras idênticas às presentes na lista são identificadas nas seqüências do banco de dados. Como o banco de dados teve sua lista de palavras pré-compilada e indexada no momento em que foi criado, não no instante em que se executa a busca, este passo é extremamente rápido.

3.a. Caso sejam encontradas duas palavras na mesma diagonal, o algoritmo inicia tentativas de estender o alinhamento, em ambas as direções e sem *gaps*, a partir da região idêntica representada pela 2ª palavra. Este mesmo algoritmo atribui nota (*score*) para o alinhamento durante a sua formação: enquanto esta nota permanecer alta a extensão é mantida, mas é interrompida quando são detectadas variações que levem à diminuição consistente deste valor. Caso o *score* deste alinhamento sem *gaps* atinja um determinado patamar, um alinhamento que permite *gaps* é disparado a partir da região já alinhada.

3.b. Se o programa identifica alinhamento com alto grau de continuidade (ou seja, com poucas regiões de inserção ou deleção), é produzido um alinhamento da comparação, realizada uma extensão deste mesmo alinhamento, a partir de adaptação do algoritmo de Smith-Waterman, e os resultados estatisticamente significativos são demonstrados graficamente para o usuário.

A utilização do programa com suas várias possibilidades também leva ao uso de uma série de recursos com expressões bastante peculiares e que são de uso muito comum durante a sua manipulação, tanto na preparação e submissão dos dados assim como na análise dos resultados. Dentre tais, estão as expressões FASTA, banco de dados, identidade, alinhamento, e-value etc..., como demonstrado abaixo.

## Formato FASTA

A submissão de dados para uma análise através do algoritmo BLAST requer a formatação desta informação de tal maneira que possa ser reconhecida e interpretada pelo programa. O formato FASTA é a forma de apresentação da seqüência, representada pelo código que utiliza uma letra para cada nucleotídeo ou aminoácido, segundo as normas IUB/IUPAC, antecedida por uma linha que pode conter qualquer registro dado pelo usuário, mas que normalmente é utilizada para descrever o que representa, a origem da seqüência ou mesmo qualquer comentário de interesse de quem está manuseando a seqüência (Fig. 2). Esta “identificação” não é obrigatória para a utilização do programa BLAST se o objetivo é submeter apenas uma seqüência (*query*), ou seja, individualmente, para comparação com um banco de dados em que os seus registros apresentem tal identificação. Preferencialmente, cada linha da seqüência deve conter no máximo 80 caracteres.

Estes dados (seqüências) são opcionalmente identificados, em seu cabeçalho, ou seja, no texto que antecede os nucleotídeos ou aminoácidos, pelo símbolo “>” (“maior que”) seguido de curta descrição da seqüência em uma linha (também idealmente, com menos de 80 caracteres), onde são informados pelo usuário a denominação, base de dados, comentário etc...

O símbolo “>” deve vir na primeira coluna e sem espaço entre ele e a primeira letra da identificação. Este formato permite também que o programa identifique, ao reconhecer novamente o símbolo, o início de uma próxima seqüência que ocorra em seguida, se o usuário submeter simultaneamente, para análise, mais de uma seqüência (formato “multi-FASTA”).

```
>gi|37221963|gb|AAN78258.1| resistance protein [Arachis simpsonii]  
LAKALYNSICDRFECACFLFNVRTISDQEEGLVRLQQTLLSKLLGEWEIKVRSVEEGISMIKEKLSKKRA  
LIVLDDVNKIEQLKALAGECDWFSYGTRIVITRDKYLLTAHKVEIKYKMKLLSDPESLELFCWNAFKISR  
PKENYEDLSNQAIIHYAQ
```

**Figura 2.** Exemplo de seqüência de aminoácidos em formato FASTA, com identificação (cabeçalho) destacado em negrito.

## Programas do BLAST

De forma geral, o programa permite a comparação de seqüências de cinco formas (“sabores”) diferentes: Blastp, Blastn, Blastx, tBlastn e tBlastx. Cada comparação a ser solicitada ao algoritmo tem um objetivo e um programa do BLAST a ser utilizado, que é baseado na necessidade do usuário e natureza dos dados e do banco de referência.

Basicamente, serão descritos aqui aqueles programas mais freqüentemente utilizados e que permitem a comparação da seqüência de entrada contra um banco de dados, ambos podendo ser apresentados, simultaneamente (MULTI-FASTA) ou não, como conjunto de nucleotídeos ou aminoácidos (Tabelas 1 e 2).

Embora menos usuais, análises mais específicas são oferecidas por outros programas também derivados do BLAST, tais como o PSI-BLAST (*Position specific iterative BLAST*), que é um refinamento do BLAST e que gera resultados considerando-se uma tabela (matriz) construída automaticamente a partir da frequência de cada aminoácido em cada posição da seqüência da proteína.

**Tabela 1.** Descrição dos programas presentes no algoritmo BLAST para comparação de seqüências.

Programa	Descrição
Blastp	Compara a seqüência de aminoácido de entrada ( <i>query</i> ) contra um banco de dados de seqüências de proteínas ( <i>subject</i> ).
Blastn	Compara a seqüência de nucleotídeos de entrada contra um banco de dados de seqüências de nucleotídeos.
Blastx	Compara a seqüência de nucleotídeos de entrada traduzida para todas as seqüências de leitura possíveis contra um banco de dados de seqüências de proteínas. É o programa mais utilizado em grandes projetos de seqüenciamento, pois permite identificar possíveis proteínas a partir de uma seqüência de nucleotídeos desconhecida.
tBlastn	Compara a seqüência de aminoácido de entrada contra um banco de dados de seqüências de nucleotídeos traduzidas para todas as seqüências de leitura possíveis.
tBlastx	Compara as seis seqüências de leitura possíveis de um nucleotídeo contra um banco de dados de nucleotídeos traduzidos para todas as seqüências de leitura possíveis.

Fonte: [www.ncbi.nlm.nih.org](http://www.ncbi.nlm.nih.org)

**Tabela 2.** Simplificação das diferentes comparações possíveis com o uso do BLAST.

Seqüência de entrada ("query")	Banco de dados ("subject")	
a.a.	a.a.	<b>Blastp</b>
nt	nt	<b>Blastn</b>
nt*	a.a.	<b>Blastx</b>
a.a.	nt*	<b>tBlastn</b>
nt*	nt*	<b>tBlastx</b>

a.a. = aminoácido; nt = nucleotídeo; \* = traduzidos para todas as seqüências (fases, "frames") de leitura possíveis

#### "Blast local"

O chamado "Blast local", como o próprio nome indica, trata-se de uma versão do programa BLAST disponível para *download* no NCBI, tanto em versão *web* quanto versão terminal (*blastall*), que o pesquisador pode utilizar para analisar, em um computador local de seu laboratório, seqüências confidenciais (*i.e.*, seqüências sigilosas, que ainda não tornaram-se públicas). O BLAST local, disponível em diversas plataformas (Windows, Solaris, Linux, etc), pode ser obtido através no endereço <ftp://ftp.ncbi.nih.gov/blast/> .

Para a versão *web* funcionar, é necessário algum conhecimento de informática, pois é preciso configurar algum servidor *web* (*i.e.*, *software* que exhibe páginas *web* para outras máquinas, p.e. Apache) na máquina hospedeira.

banco de dados para a utilização do programa BLAST (ver seção 2.3). Já o blastall é o programa BLAST propriamente dito. Embora o blastall funcione em plataforma Windows, por questões de compatibilidade e de eficiência é altamente recomendável que seja utilizado um sistema operacional do tipo UNIX (Linux, Solaris, etc). Uma distribuição Linux bem fácil de instalar e de utilizar é a Ubuntu (<http://www.ubuntu.com>).

### Saiba mais (Blast local em terminal):

Para disparar o blastall (BLAST) em terminal, digite o seguinte comando:

```
nohup blastall -p blastx -e 1e-5 -F F -d nr -i meu_multifasta.fasta -o meu_multifasta.blastx 1> arq.out 2> arq.err &
```

Onde:

*nohup* : "no hang up"; evita que o fechamento do terminal encerre a execução do programa (o comando *nohup* serve para qualquer comando terminal UNIX, não somente para o blastall);

*blastall* : arquivo binário (executável) da ferramenta BLAST;

*-p blastx* : indica que o blast é de uma seqüência nucléica contra um banco de proteínas;

*-e 1e-5* : *cutoff* (valor de corte) do *e-value* (significância estatística), geralmente assumido como o valor  $10^{-5}$ ;

*-F F* : desabilita o filtro de baixa complexidade (*i.e.*, não esconde regiões de baixa complexidade da seqüência em questão);

*-d nr* : o banco de dados (database) utilizado; neste exemplo, o nr;

*-i meu\_multifasta.fasta* : especifica arquivo multifasta de entrada;

*-o meu\_multifasta.blastx* : especifica arquivo do relatório (*report*) do blast de saída;

*1> arq.out* : redireciona mensagens de avisos (*warnings*) para esse arquivo especificado (isto também serve pra qualquer comando UNIX);

*2> arq.err* : redireciona mensagens de erros para esse arquivo especificado (isto também serve pra qualquer comando UNIX);

*&* : força o programa chamado à rodar em *background* (isto também serve pra qualquer comando UNIX).

### Banco de Dados

O banco de dados nada mais é que o local onde são armazenadas as seqüências (nucleotídeos ou aminoácidos) que são utilizadas como referência para a comparação. Como dito anteriormente, este banco pode ser de acesso restrito ou irrestrito. O banco de dados de acesso restrito se caracteriza por não compartilhar publicamente seu conteúdo e o programa BLAST é instalado para funcionamento em sigilo ("Blast local"). Os bancos públicos permitem o acesso irrestrito de dados por qualquer usuário, sem prévia autorização. Eventualmente, estes bancos públicos podem restringir (ou limitar), temporariamente, o acesso de algum usuário se detectarem excesso de demanda, o que também leva muitos usuários a preferirem o "Blast local". A descrição do mecanismo de limitação por excesso de uso do Blast do NCBI pode ser encontrada em [http://www.ncbi.nlm.nih.gov/blast/Blast.cgi?CMD=Web&PAGE\\_TYPE=BlastFAQs#QueueTime](http://www.ncbi.nlm.nih.gov/blast/Blast.cgi?CMD=Web&PAGE_TYPE=BlastFAQs#QueueTime).

Há 3 grandes bancos públicos mundiais de seqüências e que trocam diariamente dados entre si: EMBL (*European Molecular Biology Laboratory*, <http://www.embl-heidelberg.de/>), GenBank (*National Center for Biotechnology Information*, <http://www.ncbi.nlm.nih.gov/Genbank/index.html>) e DDBJ (*DNA Databank of Japan*, <http://www.ddbj.nig.ac.jp/>).

Conforme o banco a ser utilizado e o objetivo do usuário, há também a possibilidade de restringir a busca dentro do conjunto de seqüências. Ou seja, definir qual coleção de dados será utilizada dentro deste banco. O *GenBank* - o banco público de seqüências mais utilizado no mundo, com dados de mais de 100.000 organismos e que compartilha vários recursos de utilização com o EMBL e DDBJ - possui várias alternativas de busca, seja pela natureza dos dados [DNA, proteína, EST (sítio com seqüência expressa), STS (sítio com seqüência única marcada), GSS (basicamente, EST e DNA genômico) etc...], seja pelo tipo de organismo (humano, plantas, camundongo, fungos, procariotos etc...). De fato, a grande maioria dos usuários do programa BLAST utiliza o *GenBank* com acesso público para a comparação de suas seqüências. Nas Tabelas 3 e 4, são descritas informações de bancos de dados contidos no *GenBank*.

**Tabela 3.** Descrição dos bancos de dados disponíveis no algoritmo BLAST para comparação de seqüências de nucleotídeos.

Banco de dados	Descrição
nr	Todas as seqüências do <i>GenBank</i> + EMBL + DDBJ + PDB (sem seqüências de EST, STS, GSS ou HTGS). Atualmente, não exclui redundância.
refseq_rna	Seqüências de mRNA do NCBI.
refseq_genomic	Seqüências genômicas do NCBI.
est	Todos os ESTs dos bancos de dados de <i>GenBank</i> + EMBL + DDBJ.
est_human	Sub-divisão do banco "est" com coleção apenas de seqüências do ser humano.
est_mouse	Sub-divisão do banco "est" com coleção apenas de seqüências do camundongo.
est_others	Coleção de todas as seqüências originadas de ESTs, com exceção de seres humanos e camundongo.
gss ( <i>genome survey sequence</i> )	Seqüências originadas de genomas, como EST, mas também inclui dados de PCR originadas do banco "Alu" e pontas de cosmídeos, BAC e YAC, por exemplo.
htgs ( <i>high throughput genomic sequences</i> )	Seqüências depositadas em massa no banco, através de centros de seqüenciamento. Ou seja, sem "refinamento" dos dados.
Pat	Seqüências de nucleotídeos derivadas do banco de patentes do <i>GenBank</i> .
pdb	Seqüências originadas de banco de dados para estrutura tridimensional de proteínas.
month	Todas as seqüências recentes, revisadas ou não, dos bancos <i>GenBank</i> + EMBL + DDBJ + PDB depositadas nos últimos 30 dias.
alu_repeats	Banco com a tradução de todas as repetições de Alu selecionadas do REPBASE, ideal para eliminar efeitos destas repetições presentes nas seqüências de entrada ( <a href="ftp://ncbi.nlm.nih.gov/pub/jmc/alu">ftp://ncbi.nlm.nih.gov/pub/jmc/alu</a> ). Refere-se a dados de genoma de humanos e outros primatas.
dbsts	Dados de STS ( <i>Sequence Tag Site</i> ) dos bancos <i>GenBank</i> + EMBL + DDBJ.
chromosome	Genomas e cromossomos completos do NCBI, presentes também no banco Refseq_genomic
wgs	Montagem de seqüências de projetos de genoma completo por shotgun ( <i>whole genome shotgun</i> )
env_nt	Seqüências geradas de amostras com DNA de vários organismos simultaneamente (metagenoma). Não estão presentes no banco nr.

**Tabela 4.** Descrição dos bancos de dados disponíveis no algoritmo BLAST para comparação de seqüências de proteínas.

Banco de dados	Descrição
nr	Todos os arquivos não redundantes ( <i>i.e.</i> , ausência de arquivos “repetidos”) dos bancos de dados: seqüências codantes (“CDS”) traduzidas do <i>GenBank</i> +PDB+SwissProt+PIR+PRF.
refseq	Coleção (não redundante) de todas as seqüências de proteínas do NCBI.
swissprot	Banco que armazena, a partir do banco EMBL, seqüências de proteínas com alto grau de anotação (atribuição de função, estrutura de domínios, modificações pós-tradução, entre outros), baixo nível de redundância e alta integração com outros bancos de dados.
pat	Seqüências de proteínas derivadas da divisão de patentes do <i>GenBank</i> .
pdb	Seqüências derivadas do banco de dados “Brookhaven”, para estrutura tridimensional de proteínas.
env_nr	Tradução dos CDs não redundantes dos depósitos env_nt.
month	Todos os CDS novos ou revisados, traduzidos do <i>GenBank</i> +PDB+SwissProt+PIR+PRF, e liberados nos últimos 30 dias.

Fonte: [www.ncbi.nlm.nih.org](http://www.ncbi.nlm.nih.org)

## CDS

Se parte da seqüência do nucleotídeo codifica uma proteína, uma tradução conceitual, chamada CDS (*coding sequence*), é anotada. Em outras palavras, os CDS são as seqüências codantes ou a porção de uma seqüência que permite a formação dos códons (unidade de três nucleotídeos) responsáveis pela codificação dos aminoácidos. No arquivo da seqüência, o código traz a localização, se houver, do códon de iniciação e de terminação.

## Filtragem (*Filtering*)

É o processo de retirada (omissão) de regiões da seqüência de entrada que apresentam baixa “complexidade”, ou seja, que trazem pouca informação do ponto de vista biológico, tais como regiões poliA, e que, presentes, podem levar a *scores* que não correspondem necessariamente à realidade do alinhamento. Embora o programa apresente esta filtragem como já automaticamente acionada, eventualmente ela pode ser desconsiderada pelo usuário se o seu objetivo é obter informações adicionais sobre a seqüência.

## Matriz de substituição

A avaliação da qualidade do alinhamento de proteínas, um procedimento muito importante durante a análise, é realizada com o uso da matriz de substituição, que atribui uma pontuação (*score*) para cada resíduo substituído dentro de um alinhamento: alterações drásticas são penalizadas com valores baixos, e até negativos, pois a presença de aminoácidos não idênticos no alinhamento pode representar repercussões biológicas relevantes, enquanto alterações “tênuas” (também sob o ponto de vista biológico) são premiadas com valores altos. Além disso, conservações significativas são melhor pontuadas que conservações irrelevantes.

Se o alinhamento dos resíduos resulta em similaridade das propriedades físico-químicas e a substituição é freqüente entre proteínas da mesma família, a substituição é considerada conservadora e indica possível fator evolutivo envolvido. Do ponto de vista biológico, a substituição tem diferentes efeitos, as quais podem ser ponderados de acordo com as propriedades do aminoácido. Como se pode perceber, para aminoácidos, a identidade entre as seqüências pondera a repercussão estrutural da proteína a partir da presença de diferentes

aminoácidos na cadeia. As matrizes de substituição, que fornecem valores referentes à probabilidade que um determinado aminoácido x tem de sofrer uma mutação e originar o aminoácido y, consideram todos os pareamentos possíveis de aminoácidos. Percebe-se que esta abordagem também permite gerar a probabilidade de uma determinada mutação ocorrer ao longo da evolução.

Há matrizes de substituição que podem ser escolhidas no programa BLAST conforme a necessidade do usuário, tais como PAM (“Point Accepted Mutation”) e BLOSUM (“Blocks Substitution Matrix”). Atualmente verifica-se a preferência pela BLOSUM, que é uma matriz mais recente e baseada nos dados genômicos gerados a partir dos anos 90, ao contrário da matriz PAM, estabelecida na década de 70.

## BLOSUM

A denominação BLOSUM significa matriz de substituição em blocos. Esta matriz dá a pontuação para o alinhamento a partir da avaliação da frequência de substituições em blocos de alinhamentos locais em proteínas relacionadas. A matriz mais comumente utilizada é a BLOSUM62.

A construção da matriz específica é baseada em determinada distância evolucionária, que é indicada na sua denominação. Como exemplo, na matriz BLOSUM62 (a mais utilizada atualmente) o alinhamento cujos *scores* foram derivados foi construída a partir de blocos com, no máximo, 62% de identidade (*i.e.* se duas seqüências são mais de 62% idênticas o alinhamento entre elas não será empregado no cálculo da matriz). Ou seja, quanto maior o número da matriz BLOSUM, menos divergentes devem ser as seqüências relacionadas.

### **Relatório do BLAST (*report*)**

Toda a submissão correta de dados de seqüências ao programa que aplica o algoritmo BLAST irá gerar um relatório da análise realizada (Figs. 3.a. e 3.b.). Este relatório, o “*report* do BLAST”, traz uma série de informações referentes à análise e, dentre elas, a demonstração gráfica dos alinhamentos estatisticamente significativos, a lista das seqüências (com respectivas identificações) que apresentaram similaridade no banco de dados e o grau desta similaridade. Cabe salientar que as versões gráficas do relatório sofrem modificações freqüentes, ao contrário do texto em si.

NCBI Blast:Fulanina - Mozilla Firefox

http://www.ncbi.nlm.nih.gov/BLAST/blast.cgi

BLAST Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help My NCBI (Sign In) (Register)

NCBI BLAST/blast/ Formatting Results - 8AWP3B7E014 [Reformat these Results] [Edit and Resubmit] [Sign in above to save your search strategy]

Job Title: Fulanina

BLASTN 2.2.17 (Jun-24-2007)

Reference: Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schäffer, Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1990). "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs", *Nucleic Acids Res.* 25:3389-3402.

RID: 8AWP3B7E014

Database: All GenBank+EMBL+DDBJ+PDB sequences (but no EST, STS, GSS, environmental\_samples or phase 0, 1 or 2 HTGS sequences) 5,435,119 sequences; 20,979,113,122 total letters

If you have any problems or questions with the results of this search please refer to the [BLAST FAQ](#) [Taxonomy reports](#)

Query = Fulanina  
Length=600

Distribution of 25 Blast Hits on the Query Sequence

Mouse-over to show details and scores, click to show alignments

Distance tree of results **NEW**

Legend for links to other resources: [U](#) UniGene [E](#) GEO [G](#) Gene [S](#) Structure [M](#) Map Viewer

Sequences producing significant alignments: (Click headers to sort columns)

Accession	Description	Max score	Total score	Query coverage	E value	Max ident	Links
<a href="#">CT834773.1</a>	Oryza sativa (indica cultivar-group) cDNA clone:OSIGCEA029H06, full inse	592	592	99%	4e-166	82%	<a href="#">U</a>
<a href="#">NM_001050059.1</a>	Oryza sativa (japonica cultivar-group) Os01g0606100 (Os01g0606100) r	592	592	99%	4e-166	82%	<a href="#">UG</a>
<a href="#">AB028207.1</a>	Oryza sativa (japonica cultivar-group) genomic DNA, chromosome 1	592	985	99%	4e-166	100%	
<a href="#">AB028209.4</a>	Oryza sativa (japonica cultivar-group) genomic DNA, chromosome 1, PAC	592	592	99%	4e-166	82%	
<a href="#">AB028210.1</a>	Oryza sativa (japonica cultivar-group) genomic DNA, chromosome 4	93.3	398	36%	1e-15	100%	
<a href="#">AL653012.3</a>	Oryza sativa genomic DNA, chromosome 4, BAC clone: OSJNBb0069N01.	93.3	93.3	26%	1e-15	72%	
<a href="#">AB023410.4</a>	Oryza sativa (japonica cultivar-group) genomic DNA, chromosome 1, BAC	75.2	75.2	18%	3e-10	75%	
<a href="#">AB023011.2</a>	Oryza sativa (japonica cultivar-group) genomic DNA, chromosome 1, PAC	75.2	75.2	18%	3e-10	75%	
<a href="#">AB023046.2</a>	Oryza sativa (japonica cultivar-group) genomic DNA, chromosome 1, PAC	68.0	68.0	10%	4e-08	83%	
<a href="#">AB028215.1</a>	Oryza sativa (japonica cultivar-group) genomic DNA, chromosome 9	59.0	187	28%	2e-05	80%	

Alignments

Get selected sequences Select all Deselect all Distance tree of results

```
>[emb|CT834773.1] Oryza sativa (indica cultivar-group) cDNA clone:OSIGCEA029H06,
full insert sequence
Length=600
Score = 592 bits (656), Expect = 4e-166
Identities = 499/608 (82%), Gaps = 11/608 (1%)
Strand=Plus/Plus
Query 1  ACATCAATCAACCTTCCGCCCTTCCGCGGGTCAATCCCGCTCCCTCGCTGGGACCCCA 60
Sbjct 112  ACGTCATCACTTCAACCCAGCCCGCGGGTCTCCCGCCCTCCCTCGCTGGGACCCCA 171
Query 61  GCGCGAGCACTTCTCGCTCCGCGGGTCAATCCCGCTCCCTCGCTGGGACCCCA 120
Sbjct 172  GCGCGAGCACTTCTCGCTCCGCGGGTCAATCCCGCTCCCTCGCTGGGACCCCA 231
Query 121  GGGTACCGAGTGCATCGCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCT 177
Sbjct 232  GGGTACCGAGTGCATCGCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCT 291
Query 178  GACGCGCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCT 237
Sbjct 292  GACGCGCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCT 351
Query 238  GCGCGAGCACTTCTCGCTCCGCGGGTCAATCCCGCTCCCTCGCTGGGACCCCA 293
Sbjct 352  GCGCGAGCACTTCTCGCTCCGCGGGTCAATCCCGCTCCCTCGCTGGGACCCCA 410
Query 294  ATCCGCTGGACCGCCACACCGGACCGCTTCTCTCTCTCTCTCTCTCTCTCTCTCTCT 353
Sbjct 411  ATCCGCTGGACCGCCACACCGGACCGCTTCTCTCTCTCTCTCTCTCTCTCTCTCTCT 470
Query 354  CTCTCCGTAAGGGGACCTCCGCGGGTCAATCCCGCTCCCTCGCTGGGACCCCA 410
Sbjct 471  CTCTCCGTAAGGGGACCTCCGCGGGTCAATCCCGCTCCCTCGCTGGGACCCCA 530
Query 411  CCGTACGCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCT 470
Sbjct 531  CCGTACGCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCT 590
Query 471  ACGGACCGCACTTCTCGCTCCGCGGGTCAATCCCGCTCCCTCGCTGGGACCCCA 530
Sbjct 591  ACGGACCGCACTTCTCGCTCCGCGGGTCAATCCCGCTCCCTCGCTGGGACCCCA 650
Query 531  CTCAGCGCCCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCT 590
Sbjct 651  CTCAGCGCCCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCT 710
Query 591  GCGCGAGCACTTCTCGCTCCGCGGGTCAATCCCGCTCCCTCGCTGGGACCCCA 599
Sbjct 711  GCGCGAGCACTTCTCGCTCCGCGGGTCAATCCCGCTCCCTCGCTGGGACCCCA 600
```

```
>[ref|NM_001050059.1] Oryza sativa (japonica cultivar-group) Os01g0606100 (Os01g0606100)
mRNA, complete cds
Length=1518
Score = 592 bits (656), Expect = 4e-166
Identities = 499/608 (82%), Gaps = 11/608 (1%)
Strand=Plus/Plus
CDS: Putative 1 1 I I T F A P S R A V N P A S L A N D P 60
Query 1 1 ACATCAATCAACCTTCCGCCCTTCCGCGGGTCAATCCCGCTCCCTCGCTGGGACCCCA 60
```

Figura 3.a. Exemplo de relatório gráfico da versão web de BLAST, disponível na página do NCBI.



## Identificadores

Todas as seqüências depositadas nos bancos de dados do NCBI recebem automaticamente uma identificação e que, dependendo do código empregado, fornece uma série de informações relevantes quanto à origem e natureza destes dados.

### Identificadores de seqüências (gi)

O *gi* é um código único, representado por números e atribuído a toda seqüência de nucleotídeos ou proteína traduzida depositada no banco de dados do *GenBank*, não importando a origem. É uma espécie de “RG” da seqüência: individual, intransferível e não modificável. Mesmo que haja dois registros de seqüências dentro do mesmo arquivo, a versão em nucleotídeo apresentará um código e a proteína, outro. Com alterações realizadas na seqüência, será registrado um novo *gi* que pode não guardar qualquer relação numérica com o arquivo original – que não será extinto nem modificado.

### Identificador de proteína (Protein ID)

O identificador *Protein ID* é atribuído a uma seqüência de aminoácidos que pertence a algum arquivo. Como padrão, adota-se um formato composto de três letras seguidas de cinco dígitos, um ponto e o número da versão. Por exemplo, o código CAC40990.1 se refere a uma proteína com *gi* 14331118 e número de acesso AJ404328.1.

Além do *gi*, cada seqüência apresenta, concomitantemente, uma identificação que representa o código do local onde foi realizado o primeiro depósito daquela seqüência presente no *GenBank*. As origens podem ser representadas pelos seguintes códigos: gb (*GenBank*), emb (EMBL), dbj (DDBJ), ref (RefSeq), sp (Swiss Prot), pdb (*Protein Databank*), pir (PIR), prf (PRF), tpg (TPA *GenBank*), tpe (TPA EMBL) e tpj (TPA DDBJ). A expressão TPA (*Third party annotation*) significa que são seqüências que, mesmo já presentes no banco de dados, sofreram alteração de anotação posterior por terceiros, ou seja, indivíduos que não os autores ou a própria equipe do banco de dados.

### Número de acesso

O número de acesso ou *accession number* é o identificador do registro da seqüência depositada no *GenBank*, que combina letras e números, e que pertence então à coleção de seqüências do banco de dados.

Normalmente, este identificador compreende a combinação de uma letra seguida de cinco dígitos ou duas letras e seis dígitos. Ele representa o relatório completo da seqüência e não somente a seqüência em si, ao contrário do “Sequence Identifier”, que é representado pelo código “gi” ou “ProteinID”.

Exemplos:

gi|14331118|emb|CAC40990.1|

gi|41052472|dbj|BAD07483.1|

gi|15218218|ref|NP\_173005.1|

gi|92894031|gb|ABE92044.1|



idênticos, recebem pontuação positiva na Matriz de Substituição escolhida para a busca dizemos que eles são similares. A similaridade, representado por sinal positivo (“+”) depende, portanto, da Matriz de Substituição (Fig. 5).

```
Query 13 SASENK---KKGMVLPFDPHSITFDEVVYSVD 41
      S +ENK K+GM+L F+PH ITFDEV YSVD
Sbjct 532 SVTENKHYGKRGMILSFEPHCITFDEVTYSVD 563
```

**Figura 5.** Exemplo de representação da similaridade em alinhamentos entre aminoácidos. O sinal positivo (+) indica que a par de aminoácidos alinhados pontua positivamente na Matriz de Substituição empregada no alinhamento.

### Score

Nota atribuída pelo algoritmo e baseada no número de pareamentos perfeitos (*match*) e imperfeitos (*mismatch*) entre a seqüência de entrada e alguma seqüência do banco de dados. É consequência do número de inserções, deleções e substituições neste pareamento. Contudo, mesmo quando o *score* é alto e, em princípio, melhor o pareamento, este valor não pode ser analisado individualmente, mas acompanhado do respaldo estatístico (*e-value*). Este *score* de um alinhamento (S) é dado pela soma de *scores* de *matches*, *gaps* (inserções e deleções) e, no caso de aminoácidos, substituições, sendo estas últimas obtidas através de uma tabela (normalmente a BLOSUM) que reflete o grau de relevância desta “troca” de aminoácidos. Enfim, o valor do *score* dá uma indicação se o alinhamento é bom ou não, sendo o seu valor positivamente correlacionado com a qualidade deste alinhamento (ou seja, quanto maior, melhor).

### E-value

É uma das expressões mais utilizadas na análise de seqüências e representa o valor estatístico (probabilidade) que indica se o alinhamento é real ou foi obtido meramente pelo acaso naquele banco de dados (“falso positivo”). Em outras palavras, é o número esperado de falsos positivos que obteriam *score* igual ou maior que o reportado em um determinado alinhamento entre a seqüência de entrada e uma do banco.

Fundamentalmente, quanto menor o *e-value*, menores as chances daquele resultado ser consequência do acaso.

Como exemplo, um *e-value* = 0.02, dentre outras interpretações, pode significar que em um determinado alinhamento há 2 chances em 100 (ou 1 em 50) de que as similaridades sejam resultantes meramente do acaso. Do ponto de vista estatístico, tal número pode representar um valor significativo, mas do ponto de vista biológico representa, em princípio, valor bastante aquém do ideal – embora alguns pesquisadores considerem que ainda assim, tal resultado pode representar alguma informação biológica (Pertsemliadis e Fondon, 2001). Pode-se perceber então que quanto menor o *e-value* (quanto mais próximo de zero), maior a confiabilidade da predição (mais significativo é o alinhamento). Adicionalmente, em geral estabelece-se uma “nota de corte”, que determina que sejam considerados para uma primeira análise apenas *e-values* menores que *e-5* ou *e-10*. Geralmente, o valor do *e-value* é apresentado de três formas possíveis, por exemplo: 0.0,  $3e-25$  e 0.12. O valor  $3e-25$  é uma abreviação de  $3 \times 10^{-25}$ , o que indica que este número está mais próximo de 0.0 (máximo de confiabilidade) do que de 0.12.

Outro fator importante a ser considerado é que buscas dentro de um banco de dados utilizando-se seqüência de entrada muito curta podem resultar em homologia total, porém com valor de *e-value* pior (alto), pois o

cálculo também pondera o comprimento da seqüência, uma vez que seqüências menores tendem naturalmente a ter maiores chances (probabilidade) de encontrarem trechos semelhantes em um banco de dados e, por isso, podem representar apenas um evento do acaso.

$$E = K m n e^{-\lambda S}$$

Onde E é o *e-value*, m e n são os comprimentos das seqüências, K e  $\lambda$  (lambda) são parâmetros e S é o *score*. Através da fórmula, percebe-se então que o *e-value* decresce (i.e., melhora) à medida que aumenta o *score* encontrado para o alinhamento das duas seqüências.

### Identidade

É o número de resíduos (letras) similares (*matches*) identificados no alinhamento e expresso, em porcentagem, a partir da comparação com o comprimento deste alinhamento. Nos resultados do BLAST, os “positivos” (que já não representam a identidade, mas a similaridade) indicam a conservação evolutiva, ou seja, são a soma do número de aminoácidos idênticos e aqueles que são diferentes na comparação mas que apresentam *score* positivo na tabela empregada.

### Gap

Espaço introduzido em um alinhamento para compensar regiões de inserção ou deleção em uma das duas seqüências (entrada ou banco). À medida que aumenta o número de *gaps* no alinhamento, para poder acomodar o pareamento entre as seqüências, há o aumento da penalização aplicada ao *score* deste mesmo alinhamento. Ou seja, embora a introdução de *gaps* resulte em um melhor pareamento, há como consequência uma diminuição ponderada do *score*, o que, em algum grau, irá afetar desfavoravelmente o *e-value*.

### Similaridade

Similaridade é o grau da semelhança entre as seqüências. Este valor é baseado na identidade e/ou conservação da seqüência. No programa BLAST, este termo é baseado nos valores das tabelas de matrizes. Logicamente, o julgamento da similaridade vai ser limitado à disponibilidade de seqüências do banco de dados para a comparação.

Uma prática comum na sua utilização e interpretação é o uso do *hit* (seqüência do banco) mais similar nos relatórios que analisam vários genes ou proteínas.

### Homologia

A expressão se refere à similaridade atribuída a descendentes originados de ancestral comum. Ou seja, seqüências homólogas são seqüências que, supostamente, têm a mesma origem. A homologia pode ser classificada como paralogia ou ortologia. Se a homologia é resultante de especiação, ou seja, a seqüência ocorre em uma espécie que originou outra espécie que, por sua vez, apresenta uma cópia desta mesma seqüência, elas são ortólogas, embora possam não ser responsáveis por função semelhante. Por outro lado,

se há duas cópias da mesma seqüência no mesmo organismo, ou seja, uma duplicação da seqüência, estas são parálogas (FITCH, 2000).

Freqüentemente, na comparação entre seqüências, de forma errônea utiliza-se a expressão, por exemplo, “93% de homologia” enquanto o correto é dizer-se que há “93% de identidade” entre elas, uma vez que a terminologia “homólogo” se refere basicamente à origem e não ao grau de similaridade. Ou, colocado de outra forma, a porcentagem de identidade é uma medida concreta, enquanto homologia é uma hipótese sustentada por esta evidência. Em suma, a homologia se descreve apenas binominalmente (homólogo ou não-homólogo).

### Exemplos de alinhamentos

Proteína

[gi|41052472|dbj|BAD07483.1|](#) PDR-type ABC transporter 1 [Nicotiana tabacum]

[gi|75326590|sp|Q76CU2|PDR1 TOBAC](#) Pleiotropic drug resistance protein 1 (NtPDR1)

Length=1434

Score = 64.7 bits (156), Expect = 1e-09  
Identities = 33/42 (78%), Positives = 35/42 (83%), Gaps = 2/42 (4%)

```
Query 2 SPQITSTQEGDSASE--NKKKGMVLPFDPHSITFDEVVYSVD 41
      S QITST GDS SE N KKGMLVPF+PHSITFD+VVYSVD
Sbjct 806 SSQITSTDGGDSISESQNNKGMVLPFEPHSITFDDVVYSVD 847
```

Algumas informações que podem ser obtidas do alinhamento acima:

1. A seqüência de referência (banco de dados), embora presente no *GenBank* (como verificado pelo código *gi*), foi originalmente depositada no banco de dados DDBJ (*dbj*) e também no Swiss Prot (*sp*) e apresenta 1434 aminoácidos;
2. O *score* indica a pontuação atribuída pelo programa, que é baseada no número de pareamentos perfeitos (*match*), inserções e deleções (*gaps*) e substituições, entretanto tem pouca utilidade analisado isoladamente;
3. Considerando o *e-value* resultante (*Expect*), há uma probabilidade de  $1 \times 10^{-9}$  do alinhamento ter sido apenas casual, o que não se pode considerar como um ótimo valor, sobretudo por se tratar de seqüência curta (*i.e.*, poucos resíduos);
4. O alinhamento se iniciou no 2º aminoácido da seqüência de entrada e foi interrompido em seu 41º aminoácido, enquanto que na seqüência de referência o alinhamento se iniciou no 806º aminoácido e foi finalizado no 847º aminoácido;
5. Neste trecho da seqüência de referência, de 42 aminoácidos da seqüência de entrada, 33 são idênticos aos da seqüência de referência, enquanto 35 apresentaram “similaridade”, ou seja, além dos 33 idênticos entre as duas seqüências, dois aminoácidos não são idênticos (foram substituídos) mas pontuam positivamente na Matriz de Substituição empregada (*i.e.*, são conservados – representados pelo sinal “+” no alinhamento).
6. Há duas deleções presentes na seqüência de entrada de 42 aminoácidos [*Gaps* = 2/42 (4%)].



*Programação Dinâmica* – Tipo de programação (construção de algoritmos) onde a solução ótima de um estado inicial é atingida a partir do aproveitamento da solução ótima de um sub-problema do mesmo estado inicial. Este tipo de programação visa principalmente evitar o recálculo, tornando o algoritmo mais rápido.

## REFERÊNCIAS

- ALTSCHUL, S. F.; GISH, W.; MILLER, W.; MYERS, E. W.; LIPMAN, D. J. Basic local alignment search tool. **Journal of Molecular Biology**, Amsterdam, v. 215, n. 3, p. 403-410, 1990.
- ALTSCHUL, S. F.; MADDEN, T. L.; SCHAFFER, A. A.; ZHANG, J.; ZHANG, Z.; MILLER, W.; LIPMAN, D. J. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. **Nucleic Acids Research**, London, v. 25, n. 17, p. 3389-33402, 1997.
- CHENNA, R.; SUGAWARA, H.; KOIKE, T.; LOPEZ, R.; GIBSON, T. J.; HIGGINS, D. G.; THOMPSON, J. D. Multiple sequence alignment with the Clustal series of programs. **Nucleic Acids Research**, London, v. 31, p. 3497–3500, 2003.
- FITCH, W. M. Homology a personal view on some of the problems. **Trends in Genetics: DNA Differentiation & Development**, Amsterdam, v. 16, n. 5, p. 227-231, 2000.
- HENIKOFF, S.; HENIKOFF, J. G. Amino acids substitution matrices from protein blocks. **Proceedings of the National Academy of Sciences of the United States of America**, Washington, v. 89, n. 22, p. 10915-10919, 1992.
- KOSKI, L. B.; GOLDING, G. B. The closest BLAST hit is often not the nearest neighbor. **Journal of Molecular Evolution**, New York, v. 52, n. 6, p. 540-542, 2001.
- KRAUTHAMMER, M.; RZHETSKY, A.; MOROZOV, P.; FRIEDMAN, C. Using BLAST for identifying gene and protein names in journal articles. **Gene**, Amsterdam, v. 259, n. 1-2, p. 245-252, 2000.
- PERTSEMLIDIS, A.; FONDON, J. W. 3<sup>rd</sup>. Having a BLAST with bioinformatics (and avoiding BLASTphemy). **Genome Biology**, v. 2, n. 10, p. reviews2002.1-2002.10, 2001.
- SMITH, T. F.; WATERMAN, M. S. Identification of common molecular subsequences. **Journal of Molecular Biology**, Amsterdam, v. 147, p. 195-197, 1981.