

## QBQ107 – 2020s2 – Enfermagem

### Projeto da Disciplina

Professor João Carlos Setubal – [setubal@iq.usp.br](mailto:setubal@iq.usp.br)

19 de novembro de 2020

O objetivo deste projeto é familiarizar o aluno com conceitos de genômica e bioinformática, e também avaliar aprendizado de conceitos de biologia molecular apresentados nas aulas.

Para realização do projeto basta que cada aluno tenha acesso à internet e seja capaz de preparar as respostas num computador pessoal que possua um processador de textos.

### Formato dos relatórios

*Os relatórios são individuais.*

Todo relatório precisa de uma **capa ou cabeçalho**, onde devem ser incluídas as seguintes informações:

disciplina QBQ107 – Introdução à Biologia Molecular – 2020s2

Nome completo do aluno

Número USP

Título: Projeto QBQ107

*Relatórios devem estar em formato PDF. O nome do arquivo deve ser o nome do aluno. Por exemplo: joaosetubal.pdf*

Os relatórios devem ser enviados ao professor por email ([setubal@iq.usp.br](mailto:setubal@iq.usp.br)) **usando obrigatoriamente no campo assunto as palavras: projeto qbq107**. Somente se você fizer isto você receberá um recibo de entrega!

**prazo de entrega do relatório: até meio dia de 30 de novembro de 2020**

### Atividade 1: Conceitos básicos

Todo organismo tem um **genoma**, definido como sendo o conjunto de suas moléculas de DNA contendo informação genética. Mesmo vírus, que não tem status de organismos independentes, possuem informação genética armazenada em moléculas de DNA ou de RNA.

Neste projeto, trabalharemos com genomas de bactérias e de arqueias. Estes dois conjuntos de organismos são coletivamente conhecidos como **procariotos**. Embora distintos, esses organismos se caracterizam por serem unicelulares. Além dos procariotos, encontramos na biosfera os eucariotos e os vírus.

Iremos trabalhar com *representações computacionais* dos genomas de procariotos, o que é fruto de seu sequenciamento. Isto quer dizer que para fins deste projeto, um genoma é uma coleção de arquivos de computador, contendo informações diversas sobre o genoma enquanto molécula. Um desses arquivos é a *sequência do genoma*. Por exemplo, a sequência do genoma da bactéria *Xanthomonas citri* cepa A306 tem esta cara:

```
>NC_003919.1 Xanthomonas axonopodis pv. citri str. 306, complete genome
TTTTCTGCGCGCCGCAGCCAGCAACAGATGATCATACTTTGATGGATGCTTGGCCCCGCT
GTCTGGAACGTCTCGAAGCCGAATTCCCGCCCCGAAGATGTCCACACCTGGTTGAAACCCC
TGCAAGCCGAAGATCGCGGCGACAGCATCGTGCTGTACGCTCCGAACGCCTTCATCGTCTG
```

Estas são apenas as 4 primeiras linhas do arquivo, de um total de 86261 linhas. Na primeira linha está o cabeçalho, que mostra a qual organismo pertence o genoma. NC\_003919.1 é um identificador utilizado pelo *National Center for Biotechnology Information* (NCBI) para identificar esse genoma. Esse identificador recebe o nome de **Accession Number**. Nas demais linhas aparece a sequência de DNA, que pode ser entendida como uma longa cadeia de nucleotídeos, aqui representados pelas letras A, C, G, e T. A representação num arquivo de computador exige que a cadeia seja linear, mas na maioria das bactérias a molécula na verdade é circular; ou seja, para que possamos armazenar a sequência num arquivo, temos que linearizar a representação da molécula, “cortando-a” numa certa posição. A primeira base à direita desse corte é a primeira que aparece na linearização (T no exemplo acima); a base imediatamente à esquerda do ponto de corte será a última.

Outro conceito importante é o de **amplicons**. Para nós, um amplicon é um componente de um genoma que tem sua replicação independente dos demais amplicons. Os amplicons que nos interessam são os cromossomos; os demais, no caso de procariotos, se chamam plasmídeos. O exemplo acima é do cromossomo de *X. citri*. A maioria dos procariotos tem um único cromossomo.

Neste projeto iremos nos ocupar principalmente com os *genes* dos genomas, em particular os genes codificadores de proteínas (GCPs). No NCBI também é possível encontrar um arquivo com as sequências de todos os GCPs que se conhecem para um determinado genoma. Por exemplo, o primeiro gene de *X. citri*, na linearização apresentada, pode ser representado da seguinte forma:

```
>WP_011050021.1: chromosomal replication initiator protein DnaA
MDAWPRCLERLEAEFPPEDVHTWLKPLQAEDRGDSIVLYAPNAFIVDQVRERYLPRIRELLAYFVGNGDV
ALAVGSRPRAPEPAPAPVAVPSAPQAAPIVPFAGNLDShyTfANFVEGRSNQLGLAAAIQAAQKPGDRAH
```

Veja que o formato é semelhante ao do genoma: mesmo tipo de cabeçalho, e nas linhas seguintes aparece a sequência, que é de aminoácidos a não de nucleotídeos (ou seja, a porção codificadora desse gene foi traduzida de nucleotídeos para aminoácidos utilizando o código genético). Estão mostradas apenas as primeiras 3 linhas do arquivo, de um total de 8.

Nos repositórios de dados a porção codificadora de um GCP recebe o nome de *Coding Sequence*, abreviado para **CDS**.

### Roteiro para atividade 1

1. Cada aluno tem designado a si dois de genomas de procariotos, conforme a tabela tabelaGenomasQBQ107.pdf disponível no site da disciplina (<http://www.iq.usp.br/setubal/qbq107/2020/tabelaGenomas.pdf>). Os dois genomas são aqui chamados de *genoma-consulta* e *genoma-alvo*. Na tabela constam os **Accession Numbers** desses genomas. Cada “genoma” na verdade é apenas o cromossomo principal desse procarioto; estamos fazendo a simplificação de ignorar outros amplicons do organismo, caso existam.
2. Acesse a página <https://www.ncbi.nlm.nih.gov/genome/browse/#!/overview/>
3. Clique na aba “Prokaryotes”
4. Na caixa de busca, coloque o Accession Number do seu genoma-consulta.
5. Deve aparecer uma tabela com apenas 1 linha, que mostra dos dados do seu genoma-consulta.
6. Uma das colunas dessa tabela se chama **CDS**. Nesta coluna há um número. Clique nesse número. Deve aparecer a lista dos genes codificadores de proteína do seu genoma-consulta.
7. Nessa lista você deve procurar 3 genes de proteínas (enzimas) que foram mencionadas nas aulas de biologia molecular como enzimas ou proteínas importantes nos processos de replicação, transcrição, tradução, reparo de DNA, ou expressão gênica. Utilize a coluna **protein name** e/ou o recurso de busca (CTRL-F) para achar esses genes.

**Relatório da atividade 1:** apresente uma lista com 3 genes do seu genoma-consulta com as seguintes informações *para cada gene* (sua ficha):

nome do organismo / genoma utilizado; indique se o organismo é bactéria ou arqueia

locus tag (é uma coluna da lista de proteínas)

coordenadas genômicas (start e stop)

fita (fita mais ou fita menos) (strand)

nome da proteína

sequência em aminoácidos, que você pode obter clicando no link da coluna *protein product*. Na página que aparecer, haverá um pequeno link à esquerda chamado GenPept: clique nele, e vai aparecer um menu. Nesse menu escolha a opção FASTA. Desta forma a sequência vai aparecer num formato fácil de copiar e colar.

**Para cada um dos genes que você escolheu, explique a função da proteína** codificada por esse gene no processo do qual ela participa. Para esta explicação você deve usar o material das aulas, eventualmente complementado com material que você mesmo pode achar na Internet.

## Atividade 2: Conceitos básicos

Nesta atividade cada aluno irá comparar as sequências dos genes relatados na atividade anterior contra as sequências de genes do genoma-alvo, utilizando a ferramenta BLAST.

A ferramenta ou programa BLAST (*Basic Local Alignment Search Tool*) é capaz de comparar sequências de nucleotídeos e de proteínas entre si. Tal comparação tem muitas utilidades; uma delas é verificar a presença de um determinado gene *X* de organismo *A* em outro organismo *B*. Se *X* estiver presente em *B* como gene *Y*, então a ferramenta nos devolve um bom alinhamento entre *X* e *Y*. Um exemplo de alinhamento de duas sequências de proteínas é mostrado logo abaixo. A palavra ‘bom’ foi sublinhada porque este alinhamento precisa ser de *boa qualidade*, para podermos concluir que de fato *X* e *Y* representam o mesmo gene (ou seja, codificam proteínas que tem a mesma função nos respectivos organismos).

Como medir a qualidade de um alinhamento? Há várias formas, sendo que as principais são as seguintes:

O alinhamento deve incluir grande parte de *X* e de *Y*; idealmente deveria incluir *toda* a sequência *X* e *toda* a sequência *Y*. Mas é raro isto acontecer. Em geral, desejamos que ao menos **80%** de *X* e ao menos **80%** de *Y* participem do alinhamento. Vamos chamar esta medida de **cobertura**. (Veja mais abaixo uma explicação mais detalhada do que é cobertura.)

Para medir a similaridade entre *X* e *Y* utilizaremos a medida *percentual de identidade*. Esta medida nos informa quantas posições do alinhamento tem aminoácidos que são idênticos, dividido pelo tamanho do alinhamento e multiplicado por 100.

Devemos também ter uma medida da *significância estatística* do alinhamento. Para entender este conceito, considere o alinhamento da sequência LLL contra a sequência LLL (onde L representa o aminoácido leucina). O percentual de similaridade é 100%, pois as sequências são idênticas. Mas tal alinhamento não é estatisticamente significativo; ele pode ser fruto do acaso

(que é favorecido pelo fato de que leucinas são muito comuns em proteínas). Para medir significância estatística utilizaremos uma medida chamada **e-valor**, ou valor esperado (**e-value** em inglês).

E-valor nos dá o número esperado de alinhamentos ao acaso para a dada sequência de consulta e as dadas sequências do banco de sequências onde se fez a busca com BLAST. Assim, se um e-valor for igual a 3 para um dado alinhamento, podemos inferir que é alta a probabilidade de que esse alinhamento seja devido ao acaso. Um valor de 0,1 já indica que essa probabilidade é bem menor. Na prática, utilizaremos como limiar o valor  $10^{-5}$ , que é 0,00001. Ou seja, se um e-valor for igual ou menor a  $10^{-5}$ , diremos que o alinhamento é estatisticamente significativo; caso contrário, não. O fundamento teórico desta medida e deste limiar está além do escopo desta disciplina.

Note que comumente se utiliza a seguinte notação para representar e-valores menores do que 1: a letra e, seguida de um sinal negativo, seguida de um número, que é o valor do expoente. Por exemplo, o valor  $10^{-5}$  nessa notação é e-5. Se fosse e-47, isto significaria 1 sobre 1 elevado a 47, ou  $10^{-47}$ , portanto um número muito próximo de zero. Quanto mais próximo de zero for o e-valor, mais significativo é o alinhamento.

Exemplo de um alinhamento de BLAST:

```
>lcl|35099 t
Length=499

Score = 604 bits (1558), Expect = 0.0, Method: Compositional matrix adjust.
Identities = 301/499 (60%), Positives = 365/499 (73%), Gaps = 25/499 (5%)

Query 21 DAACSEAAGDKSAMMHDALFERFSARLKAQVGPEVYASWFARLKLHTVSKSVVRFTVPTT 80
          DA C E ++ LF+ S++L+ QVG +VYASWF RLK +VS ++V +VPT
Sbjct 23 DARCLETTCEE-----LFKNVSSKLEDQVGSVDVYASWFQRLKFRSVSHNIVYLSVPTN 75

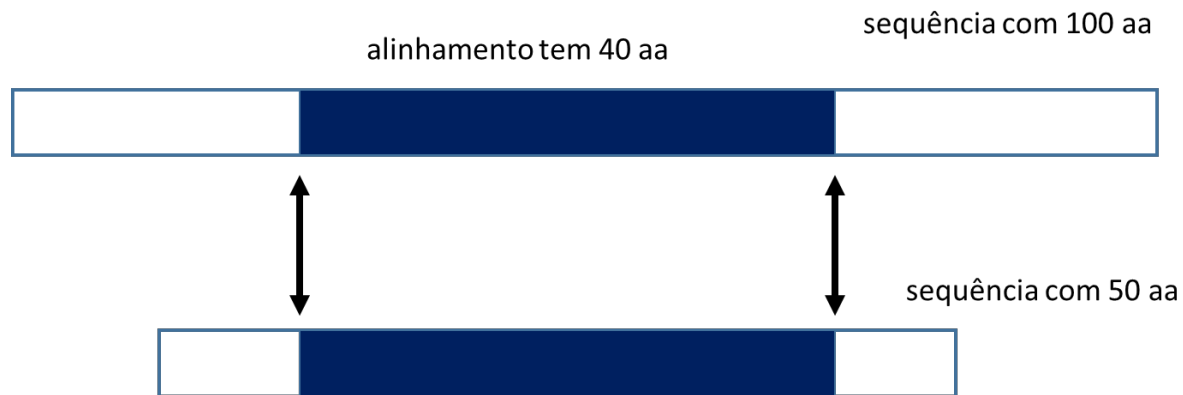
Query 81 FLKSWINNRYMDLITSLVQSEDPDVLKVEILVRSASRPVRPAQTEERAQPVQEVGAAPRN 140
          FLK+WI NRY+D IT L Q + VEI+VRS+ + P++T +
Sbjct 76 FLKAWIKNRYIDTITKLFQESISSIQGVEIIVRSAA--LMPSETS-----S 119

Query 141 KSFIPQSATAPAAQPMAAQATLRQGGSGPLFGSPLDTRFTFDIFVEGSSNRVALAAAKT 200
          S I +A P + P+FGSPLD++F F F+EG SNRVALAAA T
Sbjct 120 SSAIAHTTAKPPIINTGKISTIQKGQSINPVFGSPLDSKFVFSNFIEGPSNRVALAAHT 179
```

A figura acima mostra apenas parte de um alinhamento. A primeira linha indica o nome da sequência-sujeito (lcl|35099) (a sequência que BLAST achou ao fazer a busca com base na sequência-consulta). A segunda linha indica que essa sequência tem 499 aminoácidos. Nas duas linhas seguintes há informações sobre o alinhamento. As partes que nos interessam são: Expect (é o e-valor); Identities, que é o **percentual de identidade**; e o denominador de Identities e de Positives, que é o tamanho do alinhamento em aminoácidos. Nas linhas seguintes aparece o alinhamento propriamente dito, em seções. Cada seção tem 3 linhas, representando um trecho do alinhamento. São necessárias diferentes seções porque o alinhamento não cabe numa única linha, então é necessário quebrar o alinhamento em várias

seções. As 3 linhas de cada seção mostram a sequência-consulta na primeira linha (*query*), a sequência-sujeito na terceira linha (*subject*), e uma linha no meio que mostra as posições onde os aminoácidos são idênticos. O sinal de positivo (+) indica aminoácidos não-idênticos mas de propriedades físico-químicas semelhantes.

Neste parágrafo, explico melhor o conceito de **cobertura**. Considere o alinhamento ilustrado abaixo:



Neste alinhamento os retângulos representam sequências. A sequência de cima tem 100 aa e a sequência de baixo tem 50 aa. No entanto, o alinhamento entre elas resultou num tamanho de apenas 40 aa (a parte indicada pelas flechas). Neste caso, podemos dizer que o alinhamento cobre 40% da sequência de cima (40/100) e 80% da sequência de baixo (40/50). Note portanto que a sequência de cima tem uma cobertura diferente da sequência de baixo. Para calcular a cobertura é essencial saber o tamanho das sequências alinhadas.

Não confunda cobertura com percentual de identidade! esse percentual sempre aparece na saída do BLAST (veja *Identities*). As coberturas, você mesmo que terá que calcular com base na explicação acima.

### Resumo dos critérios do que é um bom alinhamento

*Todos os critérios abaixo devem ser satisfeitos:*

**cobertura:** pelo menos 80% da sequência-consulta e pelo menos 80% da sequência-sujeito

**percentual de identidade:** pelo menos 70%

**e-value:** no máximo  $10^{-10}$ , que é a mesma coisa que  $1e-10$ . Veja que  $1e-20$  ou  $3e-45$ , etc, são bem menores do que  $1e-10$ , e portanto bem melhores.

Os genes que você deve ter escolhido na atividade 1 são genes fundamentais para qualquer organismo celular. Por esse motivo, espera-se que estejam presentes em qualquer bactéria e em qualquer arqueia. Isto quer dizer que todas as buscas que você fará na atividade 2 devem

gerar alinhamentos (não necessariamente *bons* alinhamentos; não é uma exigência deste projeto que você necessariamente encontre *bons* alinhamentos).

Entretanto, a qualidade do alinhamento irá variar entre diferentes pares de genoma-consulta e genoma-alvo. Essa qualidade será tanto melhor quanto mais próximos esses dois genomas forem filogeneticamente entre si. A distância filogenética é uma forma de quantificar a distância evolutiva entre dois quaisquer organismos. Por exemplo, a distância filogenética que separa humanos de chimpanzés é muito menor do que a distância que separa humanos de ratos, e menor ainda do que a distância que separa humanos de moscas. Com as bactérias e arqueias se passa o mesmo. O que está escrito neste parágrafo é importante para a interpretação do alinhamento que é pedida para o relatório da atividade 2.

## Roteiro para Atividade 2

Para esta atividade, você deve ter as sequências em aminoácidos dos genes que descreveu na atividade 1. Estas serão as suas *sequências-consulta*. Além disso, na mesma tabela do website em que você pegou o genoma-consulta, você deve pegar o nome do seu *genoma-alvo*.

Para cada uma das sequências-consulta:

1. Acesse o programa BLAST em <https://blast.ncbi.nlm.nih.gov/Blast.cgi>
2. Escolha a opção Protein BLAST (também conhecida por blastp)
3. Insira a sequência na caixa **Enter Query Sequence**
4. No quadro **Choose Search Set** e no subquadro **Organism** digite o nome do *genoma-alvo* (**não** digite o accession number). O sistema tem que reconhecer o nome que vc digitou. Note que para muitos procaríotos, além de aparecer o nome da espécie também pode aparecer uma opção com designação da cepa (por exemplo: Escherichia coli K12). Você deve sempre escolher a opção genérica (sem cepa).
5. Rode BLAST apertando o botão azul BLAST.

Ao fazer isto, BLAST irá comparar a sua sequência-consulta contra *todas* as sequências de genes codificadores de proteínas do genoma-alvo (estas são as suas potenciais sequências-sujeito), e reportar todos os alinhamentos resultantes, começando pelo melhor (menor e-valor).

## Relatório da Atividade 2

Informe o nome dos organismos do seu genoma-consulta e do seu genoma-alvo.

Seu relatório deve apresentar os alinhamentos que você obteve, para cada um dos seus 3 genes. Apresente *apenas o primeiro alinhamento de cada caso*. Você pode copiar e colar o alinhamento reportado por BLAST.

Inclua comentários sobre a *qualidade e cobertura* desse alinhamento com base nos conceitos básicos explicados acima. As informações sobre a qualidade do alinhamento fazem parte do resultado do BLAST, conforme ilustrado acima.

Você deve no final interpretar o alinhamento obtido. Isto é muito importante! Se o alinhamento for bom, o que você pode concluir? e se for ruim? o que sua interpretação tem a ver com os genomas que foram comparados (o genoma-consulta e o genoma-alvo)? Compare os 3 alinhamentos que você obteve entre si. Qual foi o melhor? qual foi o pior? justifique.

**Nota bene:** o resultado esperado de sua sua busca com BLAST é que você encontrará bons “hits”. O motivo principal para isso é que os genes e suas proteínas que estudamos nas aulas, e que você deve ter escolhido para fazer suas buscas, são genes fundamentais para qualquer organismo. Assim sendo, a probabilidade é alta de que você vai encontrar bons hits. Entretanto, ocasionalmente poderá acontecer que uma particular dupla de genoma-consulta e genoma-alvo tenha como resultado “no hits”. Este resultado poderia ocorrer por exemplo se seu genoma-consulta é uma bactéria, e seu genoma-alvo é uma arqueia. Esses 2 tipos de organismos são muito distantes filogeneticamente. Caso isto tenha acontecido com sua dupla de genomas, você deverá escolher outro genoma-alvo da tabela do website (ou seja, um genoma-alvo atribuído a outro aluno) e repetir a busca. De qualquer forma, no seu relatório você sempre deve indicar que genoma-consulta e que genoma-alvo você usou, mesmo que sejam os mesmos atribuídos a você na tabela. Se mesmo assim persistir seu resultado de “no hits”, uma possível explicação é que você não fez uma boa escolha de genes!

**Nota bene 2:** uma dica para aumentar a chance de obter hits caso estes não tenham sido obtidos é alterar o parâmetro *word size* do BLASTP de 6 para 3. Esse parâmetro é encontrado na seção *Algorithm parameters* que aparece no final da página do BLAST, logo depois do botão de execução. Para ver os parâmetros disponíveis basta clicar no sinal de ‘+’.

### Atividade 3: conceitos básicos

Toda CDS de um gene de proteína pode ser representada por uma sequência de nucleotídeos e por uma sequência de aminoácidos. Para ir de uma para a outra basta fazer a tradução usando o código genético (se a tabela de código genético tem U ao invés de T, basta interpretar U como sendo T). Na atividade 2 você usou a versão em aminoácidos dos genes que escolheu. Nesta atividade, você vai repetir a atividade 2, mas agora usando a versão em nucleotídeos.

Vamos então ver que cara tem um alinhamento de nucleotídeos, através do seguinte exemplo:

	Score	Expect	Identities	Gaps	Strand	
	219 bits(242)	2e-54	646/985(66%)	22/985(2%)	Plus/Plus	
Query	358		AACCCGAAATATACATTTGATACTTTCGTCATCGGCAAGGGGAACCAAATGGCCCATGCC			417
Sbjct	5157129		AATCCAAAATATACATTTGATACTTTGTTATCGGCTCTGGTAACCGTTTGGCCCATGCA			5157188



```

Query 418      GCAGCGTTAGTTGTGTTCGGAAGAGCCCGGGACAATGTATAATCCGTTGTTTTCTACGGG 477
                ||| | |||| | || | ||| | | | | |||| | || | ||| | |||| |
Sbjct 5157189  GCTTCATTAGCTGTAGCCGAGGCGCCAGCTAAAGCGTATAACCCACTCTTTATTTACGGG 5157248

Query 478      GCGGTTGGTCTGGGCAAAACCCACCTAATGCACGCTATCGGTAACAAATGTTAGAAACC 537
                || | |||| | || | || | || | |||| | || | ||| | | | | |||
Sbjct 5157249  GGAGTTGGGCTTGGAAGACGCATTTAATGCACGCAATTGGTCATTATGTAATTGAACAT 5157308

Query 538      GATCCGACTAGTAACATCAAATACGTCACCTAGCGAATCCTTCACCAATGAATTGATCAAT 597
                |||| | | | ||| | || | ||| | || | |||| | || | |||
Sbjct 5157309  AATCCAAATGCAAAAGTTGTATATTTATCATCAGAAAAATTTACAAATGAATTTATTAAC 5157368

Query 598      GCCATCCAAACGAAAAACAGGAAGCGTTTCGGGAAGAATACCGCAACGTCGACCTGTTA 657
                | || | | ||| | | | |||| | | ||| | |||| | || | |||
Sbjct 5157369  TCTATTCGTGATAATAAGGCTGTTGATTTTCGTAATAAATATCGCAACGTAGATGTTTTA 5157428

```

Estamos vendo apenas uma parte de um alinhamento mais longo (o alinhamento completo tem 985 posições, e acima estão mostradas apenas 300 posições, ou colunas). Este alinhamento é semelhante àquele visto anteriormente. Mas desta vez, a linha do meio apenas indica as posições onde a base de cima é igual à base de baixo, por meio de traços verticais. Os indicadores do alinhamento no cabeçalho indicam que o percentual de identidade neste particular alinhamento é de 66% e o seu e-value é de  $10^{-54}$ .

### Roteiro da atividade 3:

Usaremos informações já descritas nos roteiros das atividades 1 e 2. Para cada gene que você escolheu:

1. Acesse a página que tem a lista dos genes codificadores de proteína do seu genoma-consulta. Encontre e clique no link da coluna **geneID** correspondente ao gene escolhido.
2. **Se a coluna geneID não tiver links, siga para as instruções abaixo.**
3. Na tela que aparecer, rolar para baixo até que apareça a seção *Reference assembly*. Nesta seção da página, o primeiro segmento se chama *Genomic* (depois vem outro segmento chamado mRNA and Protein(s)). No segmento *Genomic* clique no link chamado FASTA. Esse clique o levará para a sequência em nucleotídeos do gene escolhido, permitindo que você faça cópia da sequência.
4. De posse dessa sequência, repita os passos da atividade 2, mas desta vez escolha a opção **nucleotide blast**. Dentro desta opção existem sub-opções de algoritmo; você deve escolher a sub-opção **blastn** (veja a seção *Program selection*).
5. Caso o resultado do blastn seja “no hits”, repita a busca utilizando um genoma-alvo que seja diferente do seu genoma-alvo. Escolha um genoma alvo que seja do mesmo gênero que seu genoma-consulta, mas de espécie diferente. Por exemplo, se seu genoma-consulta for *Xanthomonas citri*, você poderia escolher como genoma-alvo *Xanthomonas campestris*.

**Caso a coluna geneID não tenha links:**

Neste caso será necessário baixar um arquivo com as sequências nucleotídicas de todos genes do seu genoma-consulta (infelizmente – não tem outro jeito), e depois abrir esse arquivo e achar manualmente seu gene dentro do arquivo.

Para baixar o arquivo correto, coloque no seu browser o seguinte endereço:

[https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi?db=nucleotide&mode=text&rettype=fasta\\_cds\\_na&id=accession number do genoma-consulta](https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi?db=nucleotide&mode=text&rettype=fasta_cds_na&id=accession number do genoma-consulta)

Exemplo:

[https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi?db=nucleotide&mode=text&rettype=fasta\\_cds\\_na&id=NZ\\_CP026082.1](https://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi?db=nucleotide&mode=text&rettype=fasta_cds_na&id=NZ_CP026082.1)

Ao fazer isso, deverá se abrir uma janelinha de salvamento de um arquivo chamado **sequence.txt** em seu computador. Esse arquivo contém as sequências em nucleotídeos de todos os genes do seu genoma-consulta. Para quem usa windows, sugiro abrir o arquivo com o utilitário WordPad. Para achar seu gene nesse arquivo, busque pelo accession number da proteína (esse aparece na coluna *protein product* que você viu no finzinho do roteiro da atividade 1).

**Nota bene:** esses arquivos são relativamente grandes (alguns megabytes; varia conforme o tamanho do genoma).

Se por qualquer motivo você tiver dificuldade em recuperar essa sequência, escreva ao professor, indicando o accession number do seu genoma-consulta e o accession number da proteína de interesse.

**Nota bene 2:** Uma dica para aumentar a chance de obter hits caso estes não tenham sido obtidos é alterar o parâmetro *word size* do BLASTN de 11 para 7. Esse parâmetro é encontrado na seção *Algorithm parameters* que aparece no final da página do BLAST, logo depois do botão de execução. Para ver os parâmetros disponíveis basta clicar no sinal de '+',

### Relatório da Atividade 3

Este relatório deve seguir basicamente as mesmas instruções do relatório da atividade 2. O resultado esperado é que a busca no genoma-alvo resulte em no hits ou alinhamentos ruins, mas a busca com genoma do mesmo gênero resulte em bons hits. Se seu resultado for diferente deste, você deve verificar se executou todos os passos corretamente.

**Nota bene 3:** para a atividade 3 você não precisa (e nem deve) calcular a cobertura da sequência-alvo. Calcule apenas a cobertura da sequência-consulta.

### Critério de correção dos relatórios

Serão avaliados os seguintes 10 quesitos, com 1 ponto para cada quesito

1. A1: fez atividade 1
2. A1: escolheu genes relevantes em biologia molecular mencionados nas aulas
3. A1: apresentou a ficha completa de cada gene
4. A1: explicou corretamente a função de cada gene, mesmo que tenha escolhido genes em desacordo com o estipulado em (2)
5. A2: fez atividade 2
6. A2: apresentou 3 alinhamentos
7. A2: comentou qualidade e cobertura de cada alinhamento e comparou os 3 alinhamentos entre si, apontando corretamente qual o melhor e qual o pior
8. A3: fez atividade 3
9. A3: apresentou 3 alinhamentos
10. A3: analisou os alinhamentos obtidos conforme (7)