

Atividade Prática em Sala de Aula

Gabriel Rodrigues

2025-11-05

Atividade Prática realizada no dia 05/11/2025, na disciplina de “Filogenia Molecular”, ministrada pelo Prof. Marco A. Costa, na Universidade Estadual de Santa Cruz (UESC) pelo Programa de Pós-Graduação em Genética e Biologia Molecular.

Atividade

Editar as sequências contidas nos eletroferogramas anexos, gerando um arquivo word com as sequências Forward e Reverse, e as sequências consenso. Verificar qual espécie no banco de dados do NCBI apresenta maior identidade com cada sequência consenso obtida.

Metodologia

A atividade deveria ser realizada a partir do software [BioEdit](#), entretanto não foi possível encontrar binários do software que fossem suportados pelo sistema operacional Linux (Archinux), que atualmente está sendo utilizado pelo discente.

Para contornar a situação, serão adotadas estratégias de processamento para processar os dados resultantes do sequenciamento Sanger (eletroferograma), que sigam os mesmos passos do tutorial apresentado em sala de aula.

Análise dos arquivos .ab1

Arquivos *.ab1* são dados brutos vindos de sequenciadores tipo Sanger. Esses arquivos possuem a sequência do DNA analisada, um eletroferograma (um gráfico mostrando os picos de fluorescência para cada base) e uma pontuação de qualidade associada à cada base da sequência.

Nesse caso vamos utilizar o software [CutePeaks](#) para realizar o *basecalling*, a chamada das bases individuais da sequência analisada, para conferir seu pico de fluorescência e sua qualidade.

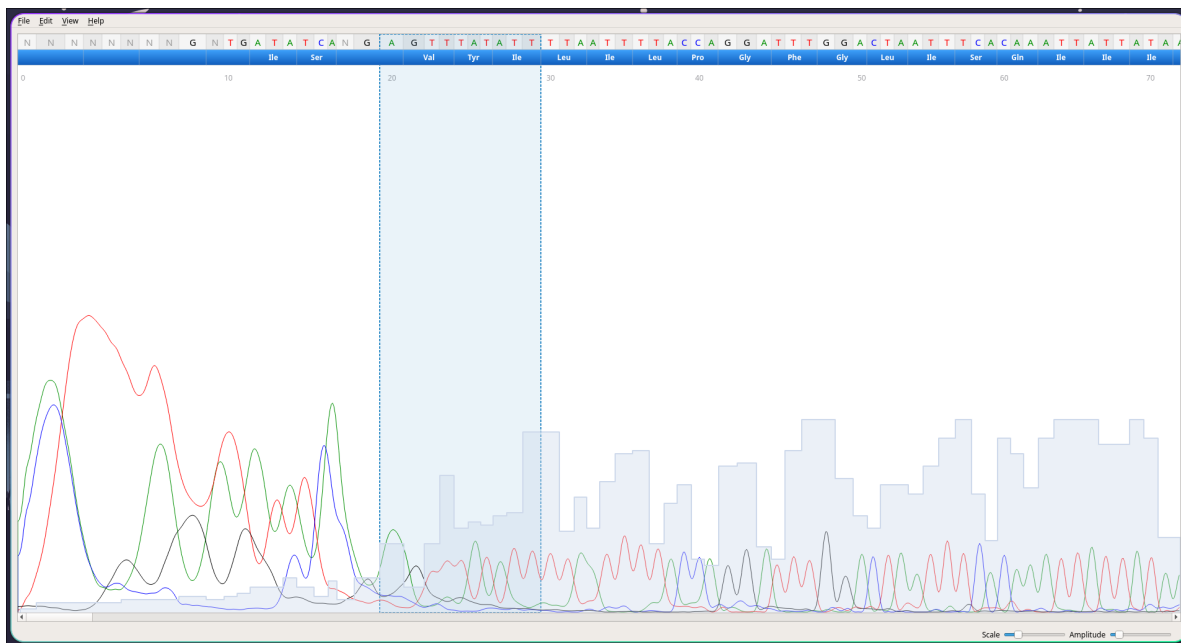


Figure 1: Exemplo de arquivo .ab1 lido pelo software CutePeaks

Conversão de .ab1 para .fastq

Como não será possível trabalhar diretamente com o software *BioEdit*, será necessário realizar a conversão dos arquivos arquivos .ab1 para o formato .fastq, formato padrão dos dados brutos de sequenciadores de nova geração, como os da plataforma *Illumina*.

Para realizar essa conversão, foi utilizada a biblioteca BioPython (suportada pela linguagem de programação Python) para criar o arquivo `abi2fastq.py` com o seguinte código:

```
# Importando a biblioteca BioPython
from Bio import SeqIO
import sys

# Definindo a entrada e a saída dos arquivos
entrada = sys.argv[1]
saida = sys.argv[2]

# Realizando a conversão dos dados
records = SeqIO.parse(entrada, "abi")
count = SeqIO.write(records, saida, "fastq")
```

```
python abi2fastq.py COC_2F.ab1 COC_2F.fastq
```

[illegible]

Checagem da qualidade dos arquivos fastq

3

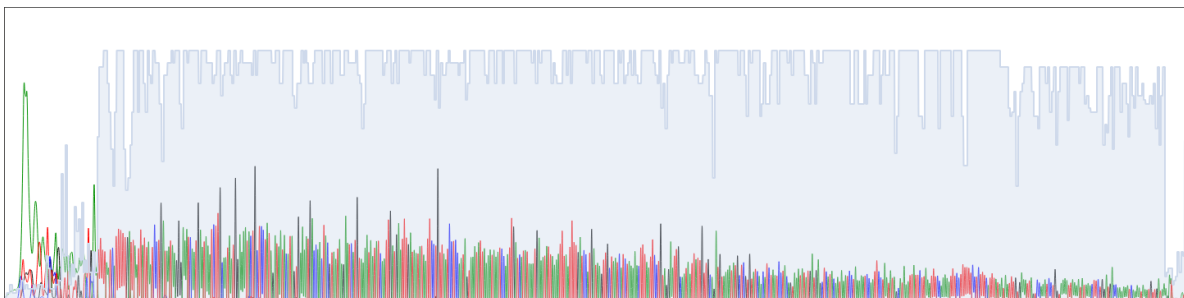


Figure 2: Qualidade do eletroferograma do arquivo RBA_2R no software CutePeaks, a qualidade é representada pelas barras cinzas

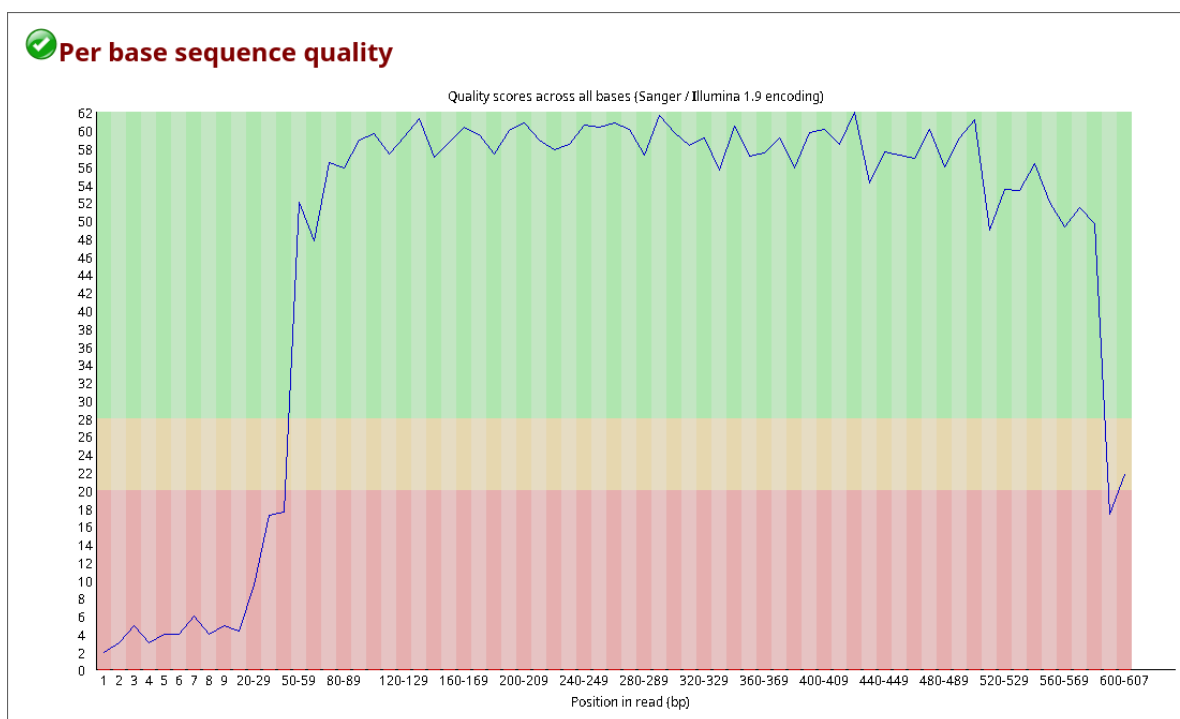


Figure 3: Checagem de qualidade do arquivo fastq (RBA_2R), através do software FastQC

Remoção das bases de baixa qualidade

Para a remoção das bases de baixa qualidade foi utilizado o software [Seqtk](#), em sua função `trimfq`, que remove automaticamente as sequências de baixa qualidade ($Phred < 15$) e bases não definidas.

As análises de qualidade também foram realizadas nas sequências após a remoção das bases de baixa qualidade.

Conversão de FASTQ para FASTA

Após a trimagem das sequências e remoção de bases de não definidas, os arquivos *FASTQ* foram convertidos para o formato *FASTA* através da biblioteca **BioPython**, seguindo a mesma lógica da conversão dos arquivos *.ab1*:

```
# Importando a biblioteca BioPython
from Bio import SeqIO
import sys

# Definindo a entrada e a saída dos arquivos
entrada = sys.argv[1]
saida = sys.argv[2]

# Realizando a conversão dos dados
records = SeqIO.parse(entrada, "fastq")
count = SeqIO.write(records, saida, "fasta")
```

O arquivo final `fastq2fasta.py` então é chamado no terminal linux:

```
python fastq2fasta.py COC_2F.fastq COC_2F.fasta
```

Um arquivo em formato FASTA se assemelha à um arquivo FASTQ, porém sem os índices de qualidade por base (*Phred Score*) anexados ao arquivo das sequências. Outra diferença são os símbolos adotados, o *header* de uma sequência (seu cabeçalho) agora é marcado pelo símbolo `>` antes do nome da sequência.

Obtenção da Sequência Consenso

A sequência consenso é o resultado da união via sobreposição das sequências *forward* e *reverse* do sequenciamento. Para fazer a montagem dessas sequências em uma única final. Para isso pode-se utilizar o software [CAP3](#), que alinha a sobreposição das sequências para estabelecer as contigs finais, uma técnica denominada de *Overlap Layout Consensus*.

A entrada para o software CAP3 foi um arquivo FASTA com a junção das duas sequências (*forward* e *reverse*) de cada amostra, dessa maneira (Arquivo **COC_merge.fasta**):

>13F

```
TTTATATTTTAATTTTACCAGGATTTGGACTAATTTACAAAATTATTATAAATGAAAGAG
GAAAAAAGGAAATTTTGGAAATTTAAGAATAATTTATGCTATATTAGGAATTGGATTTT
TAGGATTTATTGTATGAGCTCATCATATTTACTGTAGGATTAGATGTTGATACACGAG
CATATTTTACATCTGCAACAATAATTATTGCAATTCCTACAGGAATTAAAGTTTTAGAT
GATTAGCAACTTATCATGGATCAAAATTTAAATTTTAATATTTTCAATTTATATGATCAATTG
GATTTATTTTAATATTTTACTATTGGAGGATTAACAGGAATTATATTATCAAATTCATCAA
TTGATATTATTTTACATGATTCTTATTACGTAGTTGGTCATTTTCACTATGTATTATCTA
TAGGAGCAGTATTTTCCATTATTGCAAGATTATTTCATTGATTTCCCTTATTATCAGGAT
TAATAATTAATCAAAAATGATTAATAATTTCAATTTTTTTTTTATATTTCATTGGAATTAATT
TCACTTTTTTTCCTCAACATTTTTTAGGATTAATATC
```

>13R

```
AAAAAAAAAATTGAAATTTTAATCATTTTTTGATTAATTATTAATCCTGATAATAAGGGAAA
TCAATGAATAAATCTTGCAATAATGGAAAATACTGCTCCTATAGATAATACATAGTGAAA
ATGACCAACTACGTAATAAGAATCATGTAAAATAATATCAATTGATGAATTTGATAATAT
AATTCCTGTTAATCCTCCAATAGTAAATATTTAAATAAATCCAATTGATCATATAAATGA
AATATTTAAATTTAATTTTGATCCATGATAAGTTGCTAATCATCTAAAACTTTAATTCC
TG TAGGAATTGCAATAATTATTGTTGCAGATGTAAAATATGCTCGTGTATCAACATCTAA
TCCTACAGTAAATATATGATGAGCTCATACAATAAATCCTAAAAATCCAATTCCTAATAT
AGCATAAATTATTCTTAAATTTCCAAAAATTCCTTTTTTCTCTTTTCATTTATAATAAT
TTGTGAAATTAGTCCAAATCCTGGTAAAATTTAAATATAAACTTCTGGATGACCAAAAAA
TCAAAATAAATGTTGATAA
```

Alinhamento das sequências

As sequências finais foram Alinhadas via [BLASTn online](#) contra o banco de dados `core_nt`, uma versão reduzida e clusterizada do banco total de nucleotídeos `nt`.

Resultados

Qualidade das Sequências

A qualidade sequências dos arquivos *FASTQ* foram realizadas com o software **FastQC** e depois agragadas com o software **MultiQC**. As qualidades foram medidas antes e depois da remoção de sequências de baixa qualidade.

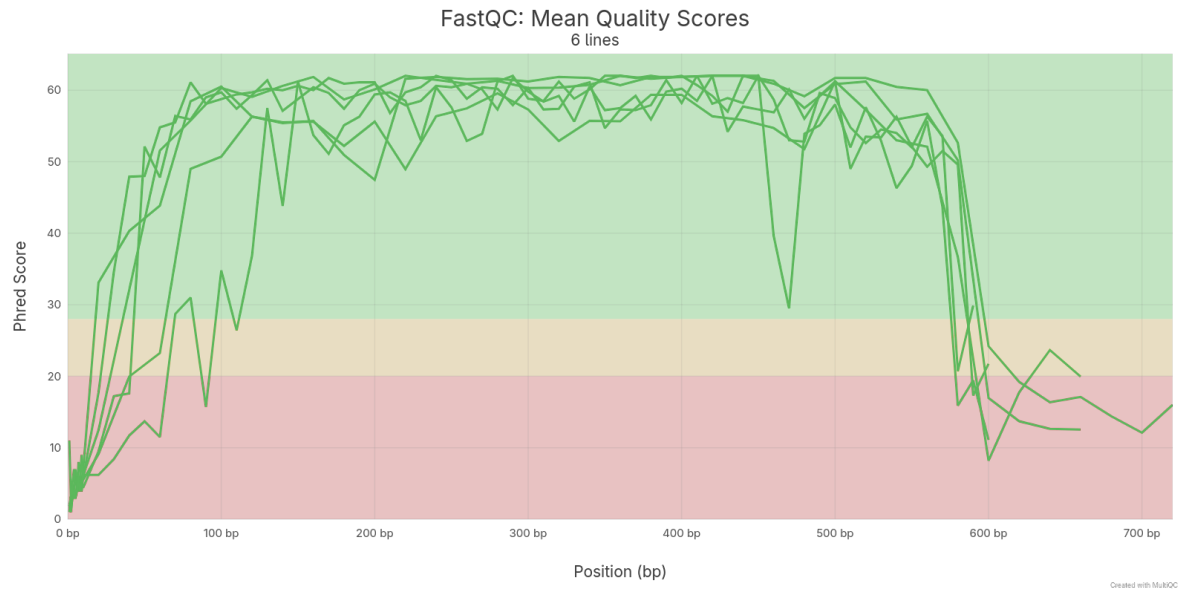


Figure 4: Qualidade das sequências SEM a remoção das bases de baixa qualidade

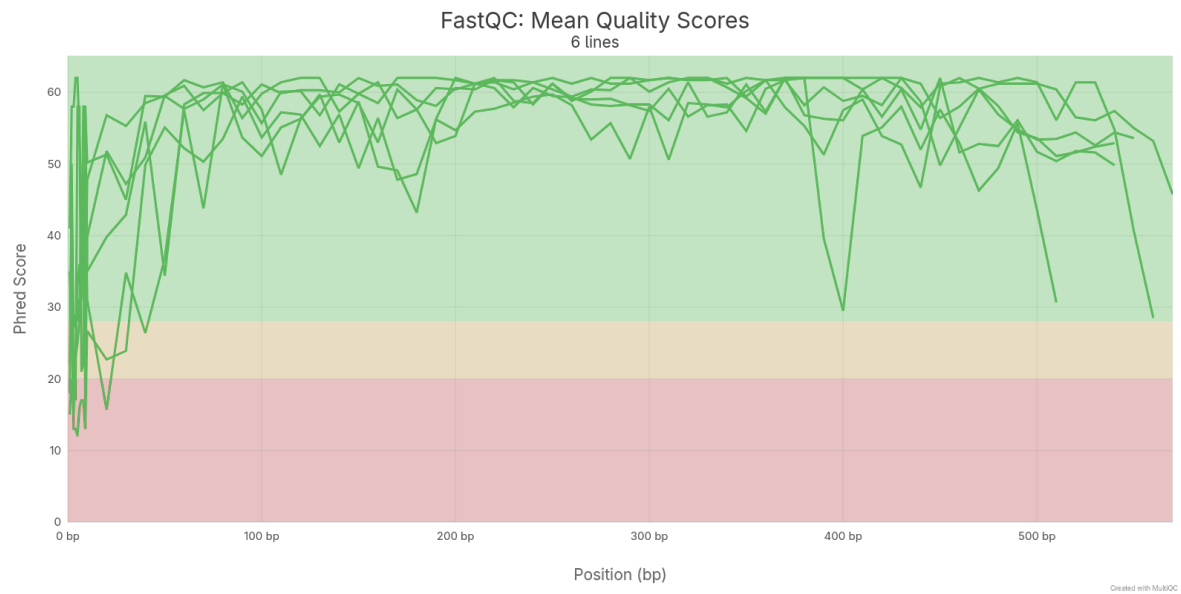


Figure 5: Qualidade das sequências COM a remoção das bases de baixa qualidade

Sequências montadas

As sequências fasta (*forward* e *reverse*) foram unidas em arquivos únicos para cada código e amostra, **COC_merged.fasta**, **MAC_merged.fasta** e **RBA_merged.fasta**. Esses arquivos foram montados utilizando o software CAP3 e tiveram como resultado final os seguintes tamanhos em bases:

Table 1: Tamanho das sequências em nucleotídeos

Amostra	Tamanho (nt)
COC	615
MAC	607
RBA	617

Como resultado da montagem, temos o arquivo FASTA final para cada sequência:

```
>COC
TTATCAACATTTATTTTGATTTTTGGTCATCCAGAAGTTTATATTTTAATTTTACCAGG
ATTTGGACTAATTTACAAAATTATTATAAATGAAAGAGGAAAAAAGGAAATTTTGGAAA
TTTAAGAATAATTTATGCTATATTAGGAATTGGATTTTTAGGATTTATTGTATGAGCTCA
TCATATATTTACTGTAGGATTAGATGTTGATACACGAGCATATTTTACATCTGCAACAAT
AATTATTGCAATTCCTACAGGAATTAAAGTTTTTAGATGATTAGCAACTTATCATGGATC
AAAATTTAAATTTAATATTTTCATTTATATGATCAATTGGATTTATTTTAATATTTACTAT
TGGAGGATTAACAGGAATTATATTATCAAATTCATCAATTGATATTATTTTACATGATTC
TTATTACGTAGTTGGTCATTTTCACTATGTATTATCTATAGGAGCAGTATTTTCCATTAT
TGCAAGATTTATTCATTGATTTCCCTTATTATCAGGATTAATAATTAATCAAAAATGATT
AAAATTTCAATTTTTTTTTTATATTCATTGGAATTAATTTCACTTTTTTTCCTCAACATTT
TTTAGGATTAATATC

>RBA
TTATCAACATTTATTTTGATTTTTGGTCATCCAGAAGTTTATATTTTAATTTTACCAGG
ATTTGGACTAATTTACAAAATTATTATAAATGAAAGAGGAAAAAAGGAAATTTTGGAAA
TTTAAGAATAATTTACGCTATATTAGGAATTGGATTTTTAGGATTTATTGTATGAGCTCA
TCATATATTTACTGTAGGATTAGATGTTGATACACGAGCATATTTTACATCTGCAACAAT
AATTATTGCAATTCCTACAGGAATTAAAGTTTTTAGATGATTAGCAACTTATCATGGATC
AAAATTTAAATTTAATATTTTCATTTATATGATCAATTGGATTTATTTTAATATTTACTAT
TGGAGGATTAACAGGAATTATATTATCAAATTCATCAATTGATATTATTCTACATGATTC
TTATTACGTAGTTGGTCATTTTCACTATGTATTATCTATAGGAGCAGTATTTTCCATTAT
TGCAAGATTTATTCATTGATTTCCCTTATTATCAGTGATTAATAATTAATCAAAAATGAT
TAAAATTTCAATTTTTTTTTTATATTCATTGGAATTAATTTAACTTTTTTTCCTCAACATT
TTTTAGG
```


>MAC

```
TTATCAACATTTATTTTGATTTTTTGGTCATCCAGAAGTTTATATTTTAATTTTACCAGG
ATTTGGACTAATTTACAAAATTATTATAAATGAAAGAGGAAAAAAGGAAATTTTGGAAA
TTTAAGAATAATTTATGCTATATTAGGAATTGGATTTTATAGGATTTATTGTATGAGCTCA
TCATATATTTACTGTAGGATTAGATGTTGATACACGAGCATATTTTACATCTGCAACAAT
AATTATTGCAATTCCTACAGGAATTAAGTTTTAGATGATTAGCAACTTATCATGGATC
AAAATTAATTTTAATATTTTCATTTATATGATCAATTGGATTTATTTTAATATTTACTAT
TGGAGGATTAACAGGAATTATATTATCAAATTCATCAATTGATATTATTTTACATGATTC
TTATTACGTAGTTGGTCATTTTCACTATGTATTATCTATAGGAGCAGTATTTTCCATTAT
TGCAAGATTTATTCATTGATTTCCCTTATTATCAGGATTAATAATTAATCAAAAATGATT
AAAATTTCAATTTTTTTTTTATATTCATTGGAATTAATTTCACTTTTTTTCCTCAACATTT
TTTAGGATTTAATATCA
```

Alinhamento no BLAST

Foi realizado o alinhamento global das sequências em um banco de dados de nucleotídeos (*NCBI Core NT*), os resultados foram os seguintes:

Table 2: Melhores resultados do alinhamento no BLAST para cada amostra

Amostra	Melhor Hit	Organismo
COC	H2 cytochrome c oxidase subunit I (COI) gene	<i>Partamona rustica</i>
MAC	H2 cytochrome c oxidase subunit I (COI) gene	<i>Partamona rustica</i>
RBA	H5 cytochrome c oxidase subunit I (COI) gene	<i>Partamona rustica</i>

Os alinhamentos tiveram índices de identidade e cobertura maiores que 99% para os melhores alinhamentos.

As sequências aparentam ser derivadas da espécie *Partamona rustica*, especificamente vindas de subunidades H do gene **citocromo c oxidase**, sendo uma enzima na cadeia respiratória de transporte de elétrons das células.