



O que é possível fazer sabendo o básico de Python?

V Workshop de Python para Dados Biológicos

Beatriz Rodrigues Estevam

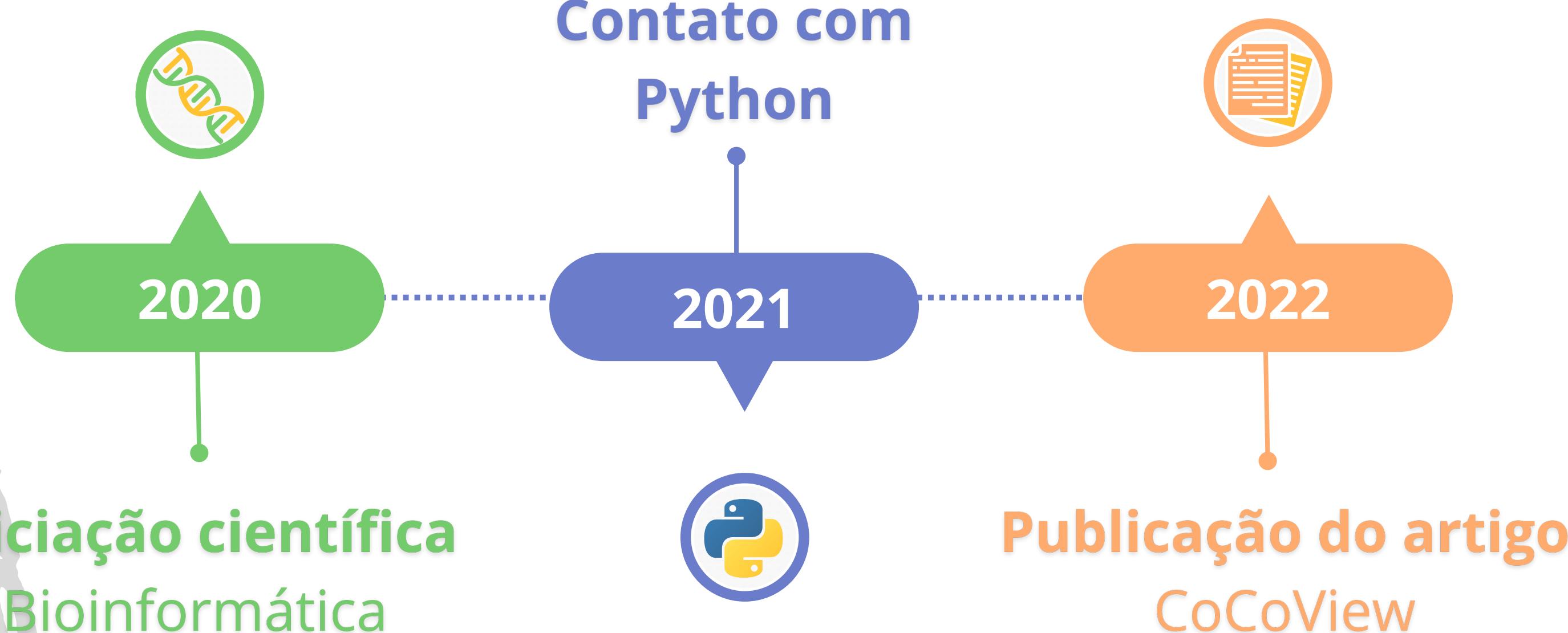
bia.estevam.25@usp.br

Sumário



-  Introdução
-  Sequências conservadas
-  Sequence Logos
-  CoCoView: A codon conservation viewer
-  Exemplos e aplicações

Por que "O que é possível fazer sabendo o básico de Python?"?



Por que "O que é possível fazer sabendo o básico de Python?"?

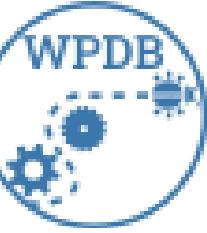


**Ferramenta para visualização da conservação
de sequências de códons.**

**Publicação do artigo
CoCoView**



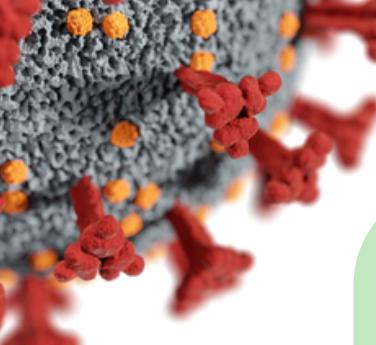
Por que "O que é possível fazer sabendo o básico de Python?"?



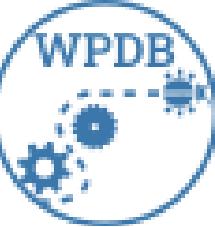
**Ferramenta para visualização da conservação
de sequências de códons.**

**Publicação do artigo
CoCoView**

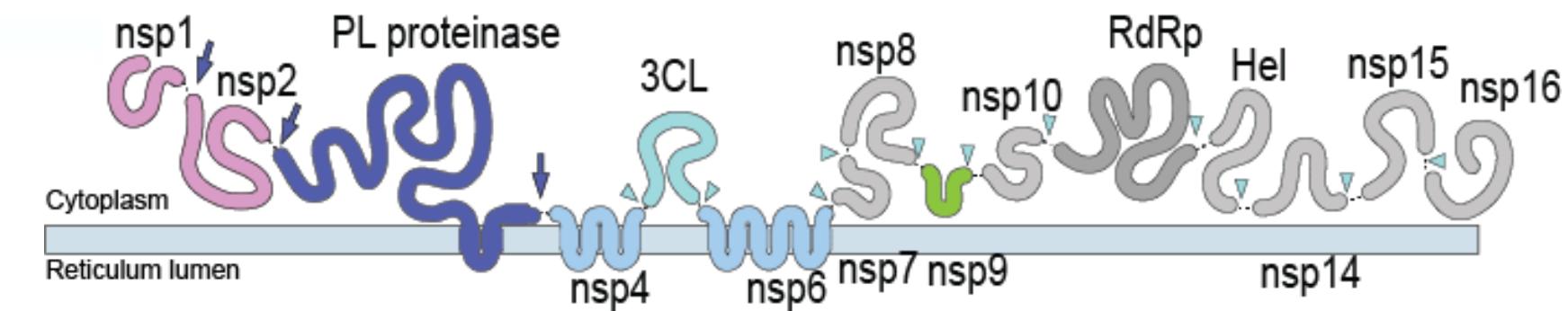
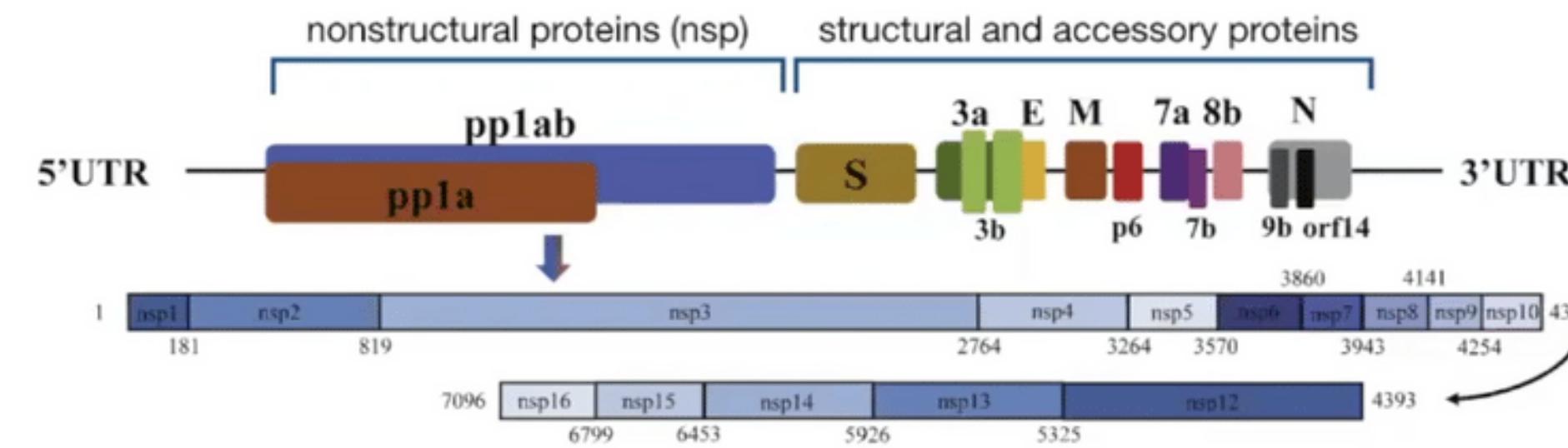




Como separar as proteínas?

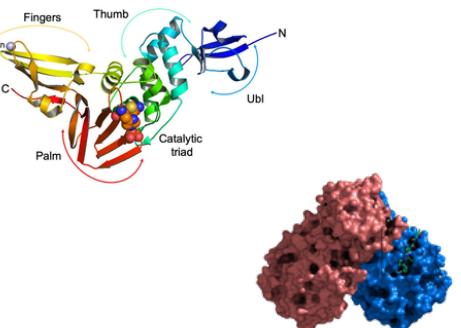


Transcrição,
Tradução e pós-
processamento

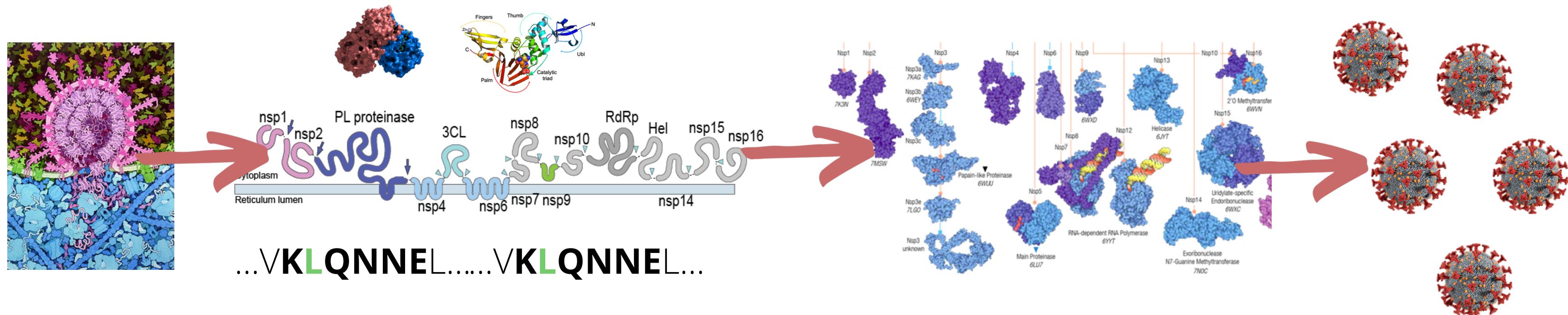


...VKLQNNE...

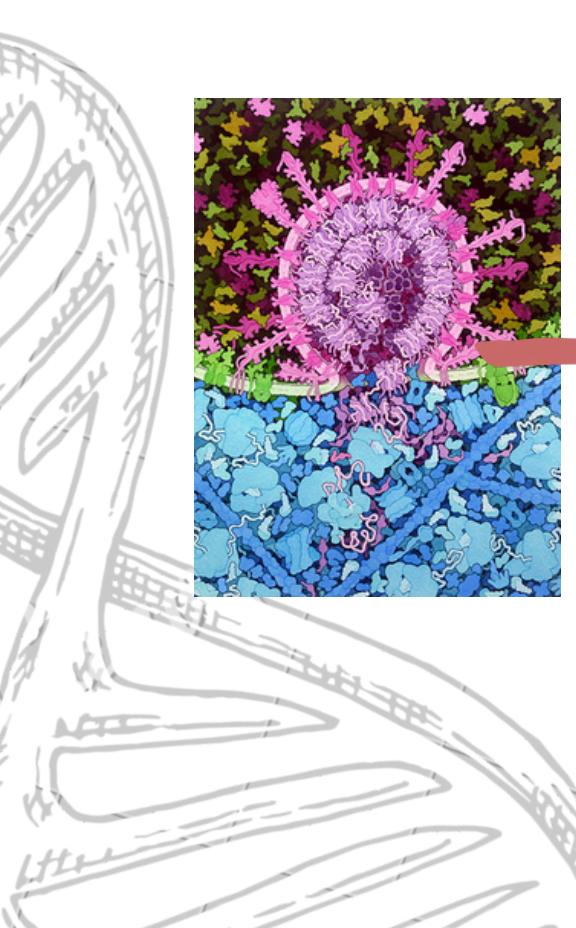
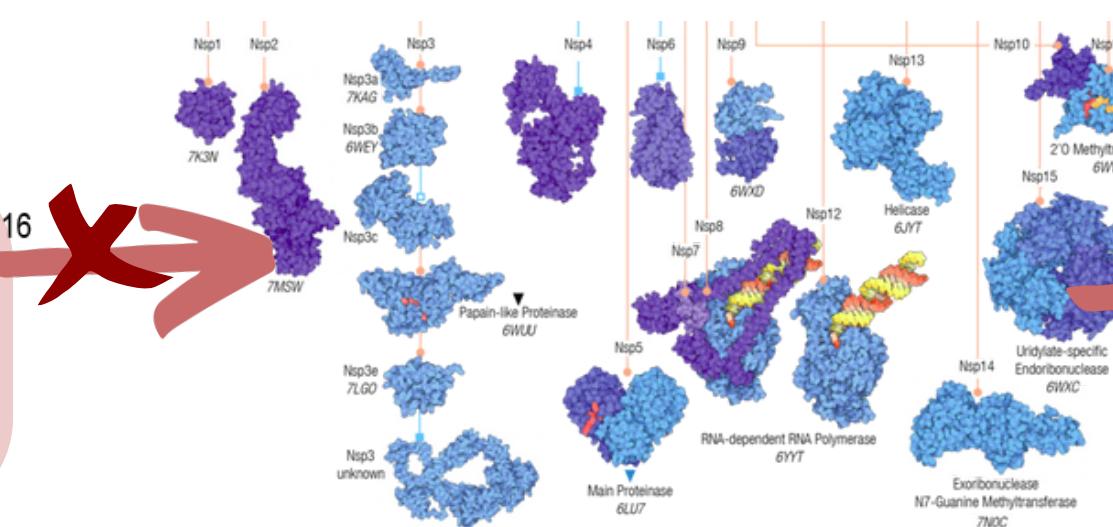
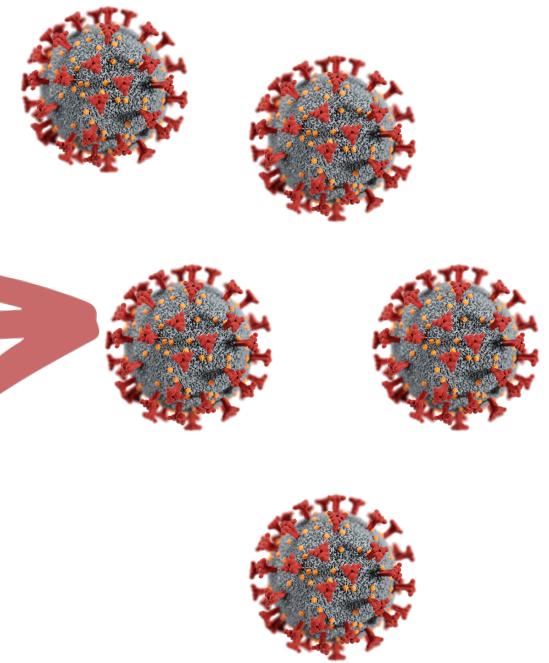
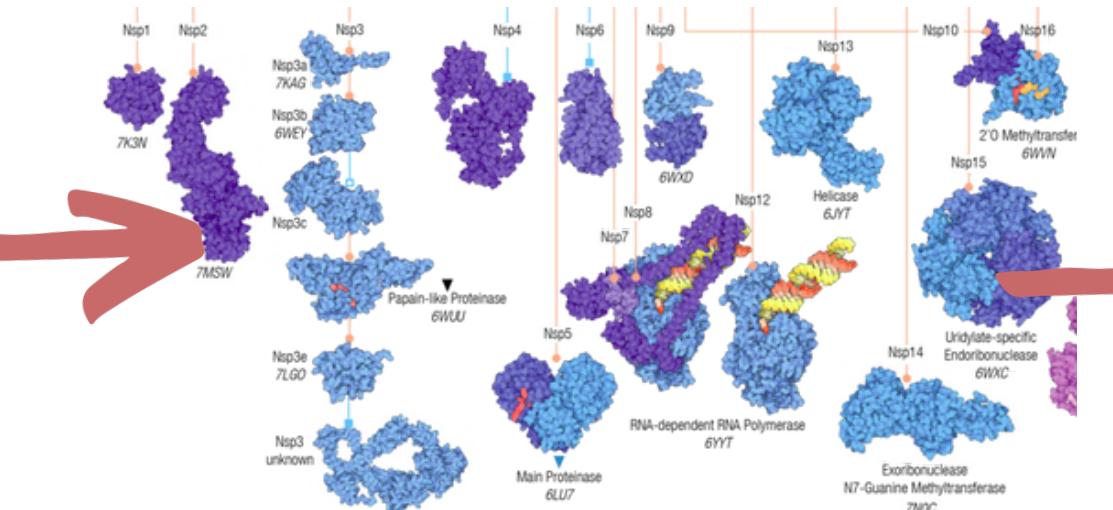
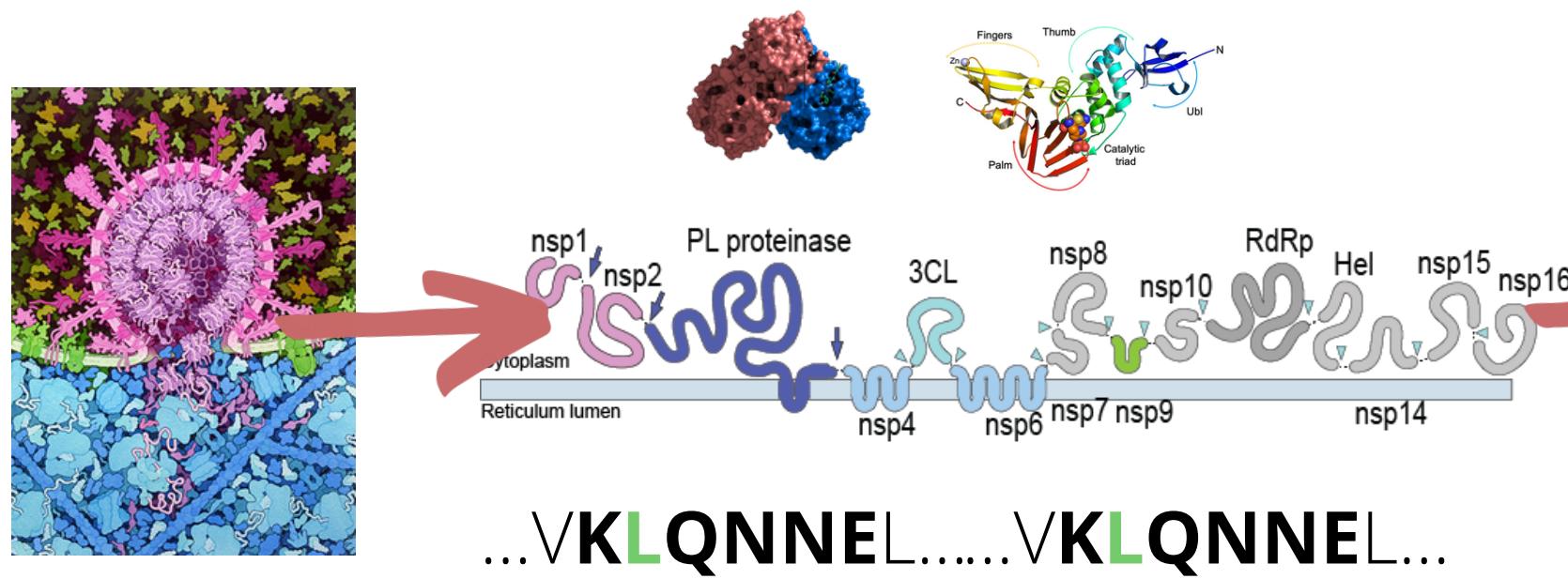
nível de
aminoácido



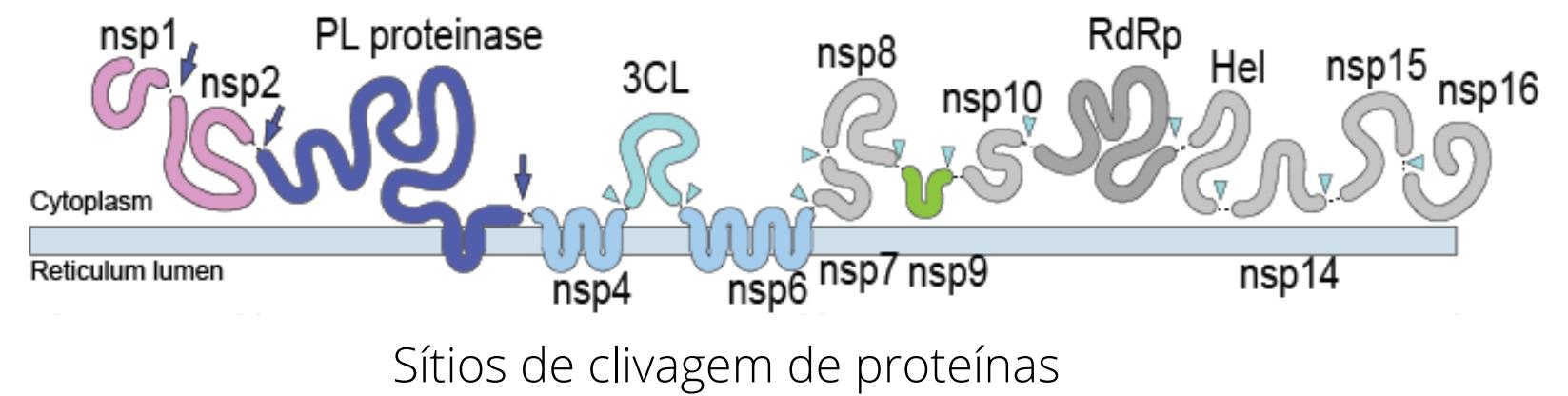
Sequências conservadas



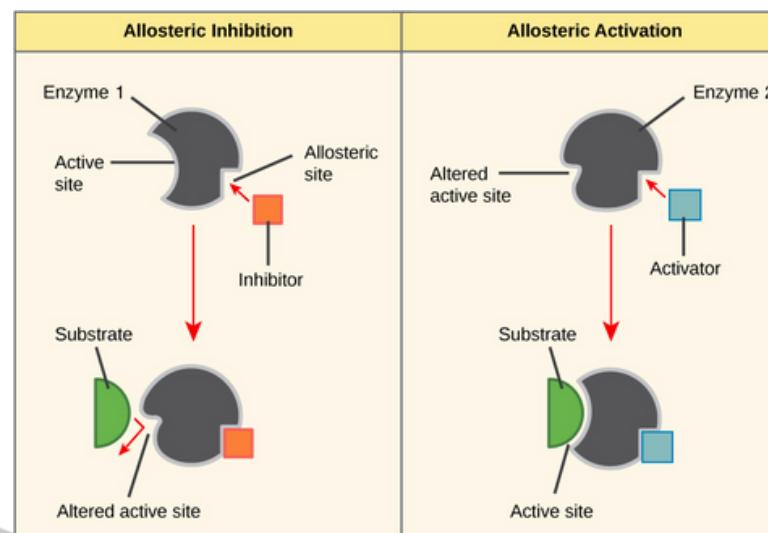
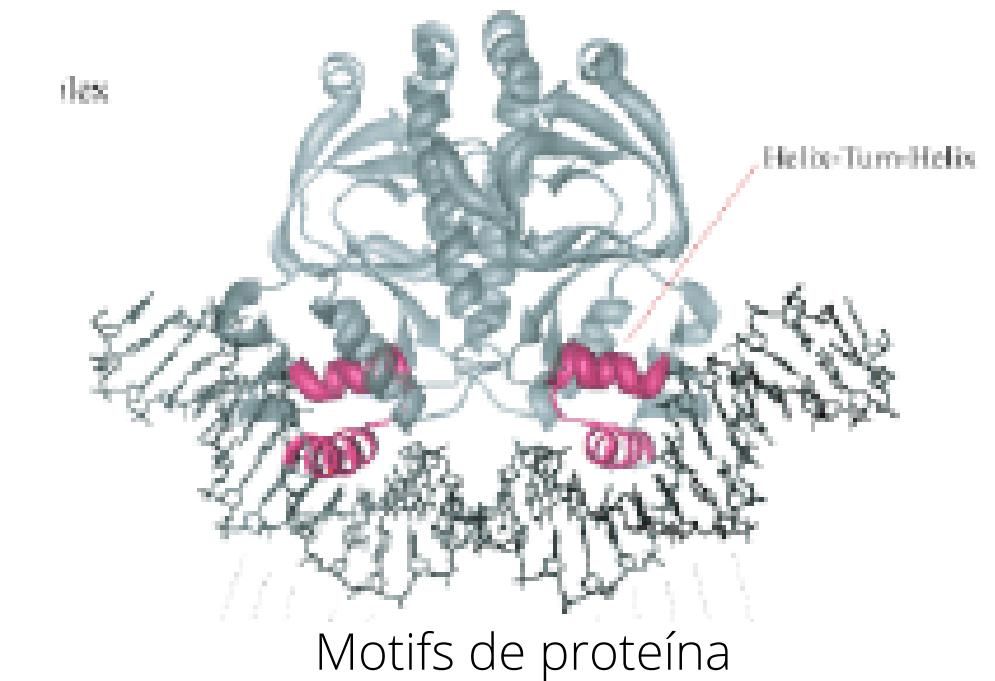
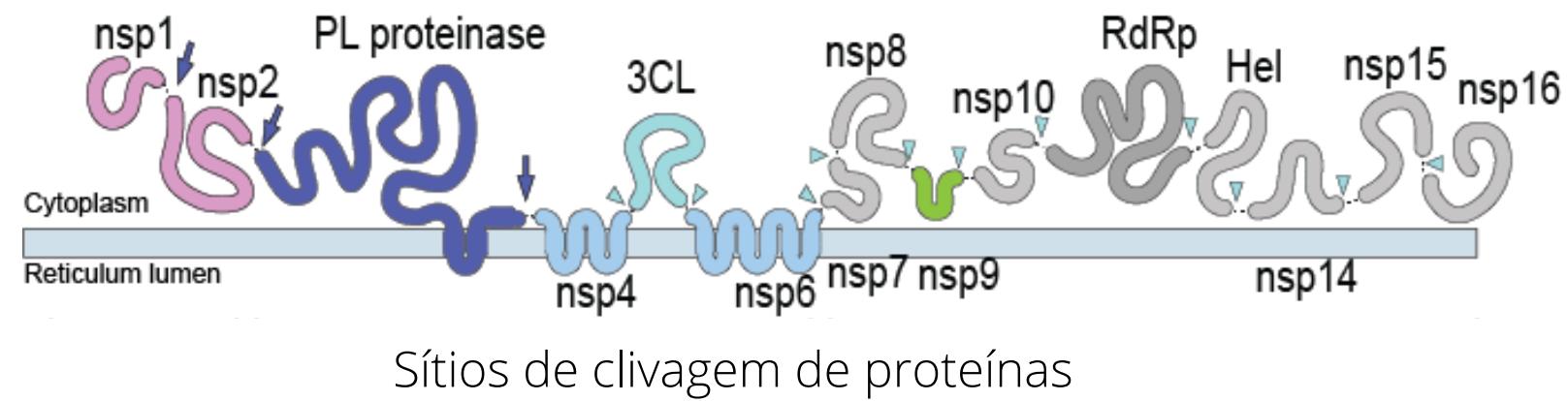
Sequências conservadas



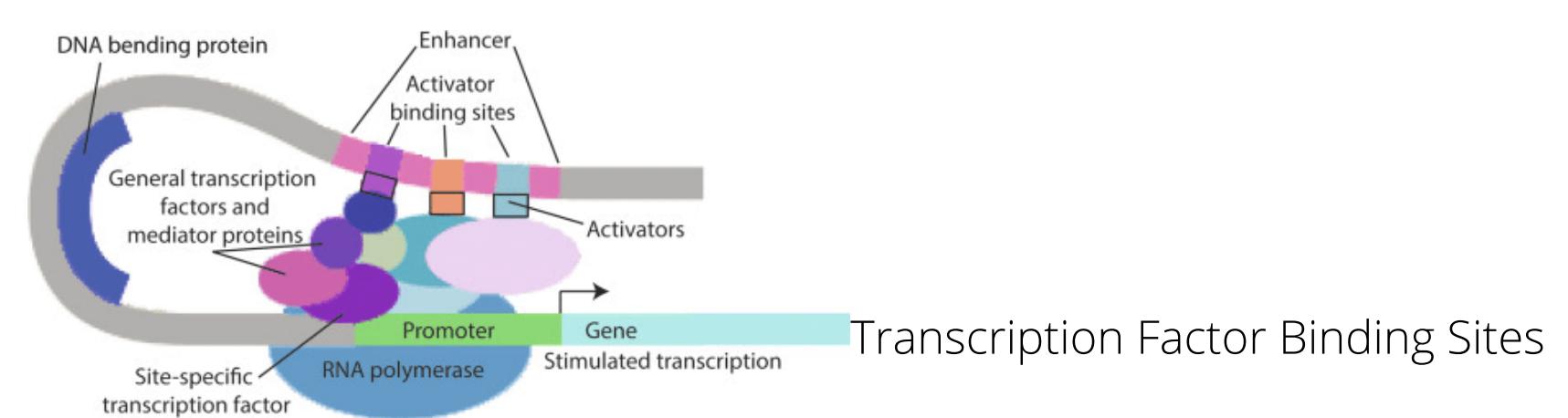
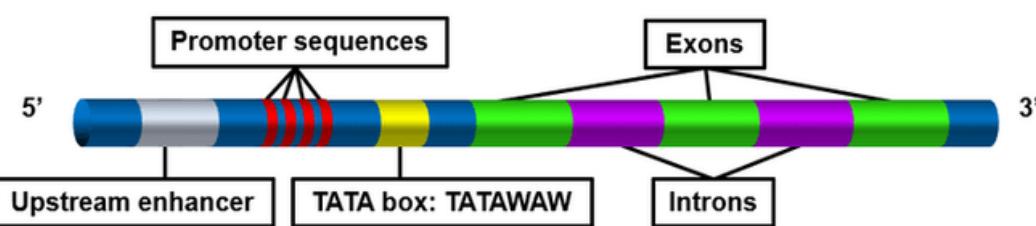
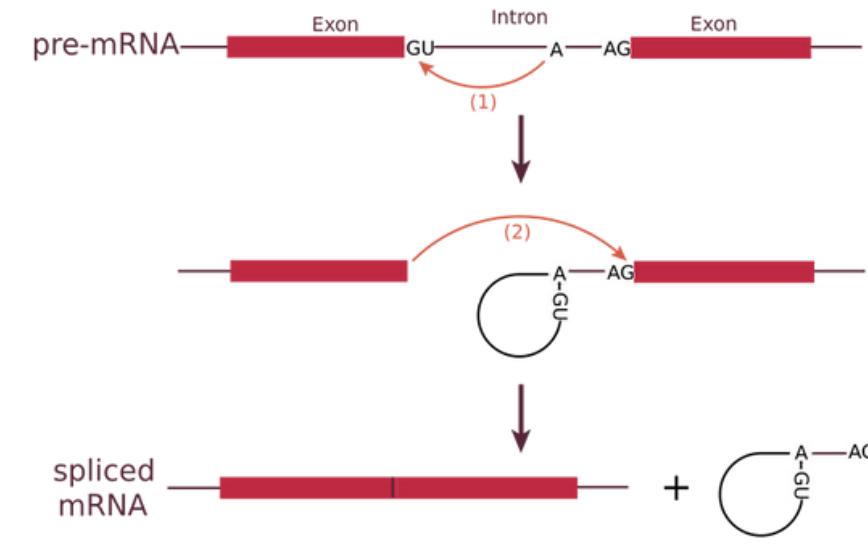
Sequências conservadas



Sequências conservadas



Sítios allostéricos e sítios ativos de enzimas



Sequências conservadas



A alta conservação de sequências indica possíveis funções relevantes (Se mudam drasticamente, são desfavorecidas)



Existe interesse em identificar e representar essas sequências

Sequence Logos



Forma de representação de sequências conservadas introduzida por Schneider and Stephens (1990)

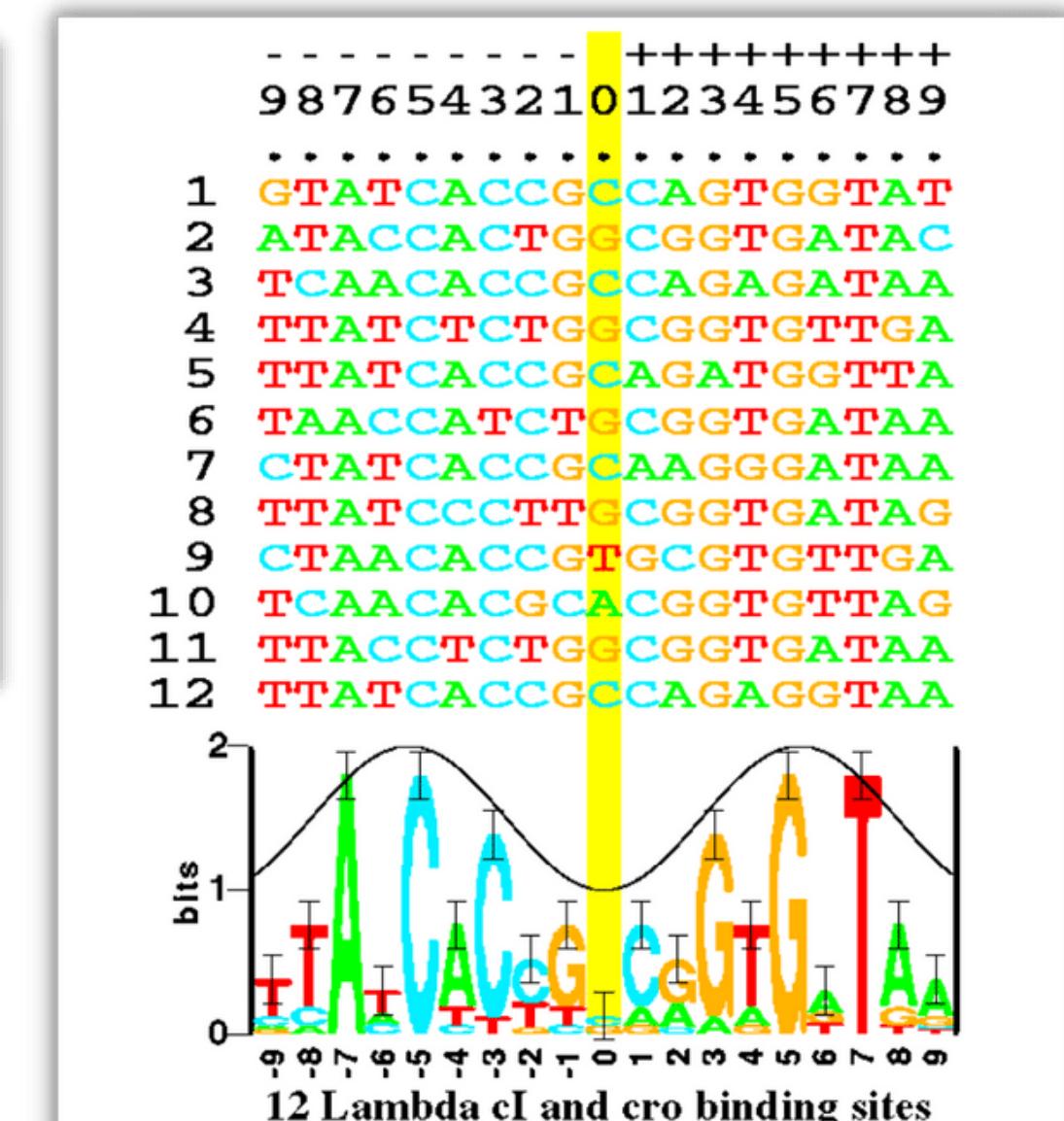
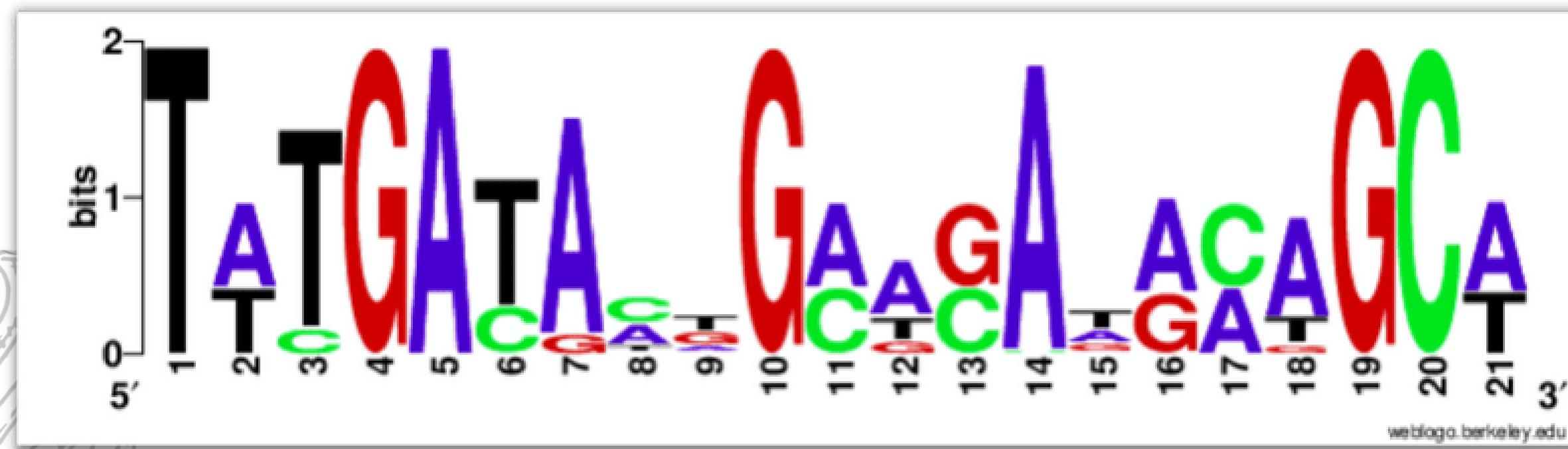


Fig. 1. Some aligned sequences and their sequence logo. At the top of the figure are listed the 12 DNA sequences from the P_L and P_R control regions in bacteriophage lambda. These are bound by both the cI and cro proteins [16]. Each even numbered sequence is the complement of the preceding odd numbered sequence. The sequence logo, described in detail in the text, is at the bottom of the figure. The cosine wave is positioned to indicate that a minor groove faces the center of each symmetrical protein. Data which support this assignment are given in reference [17].

Entendendo os Sequence Logos



Conservação

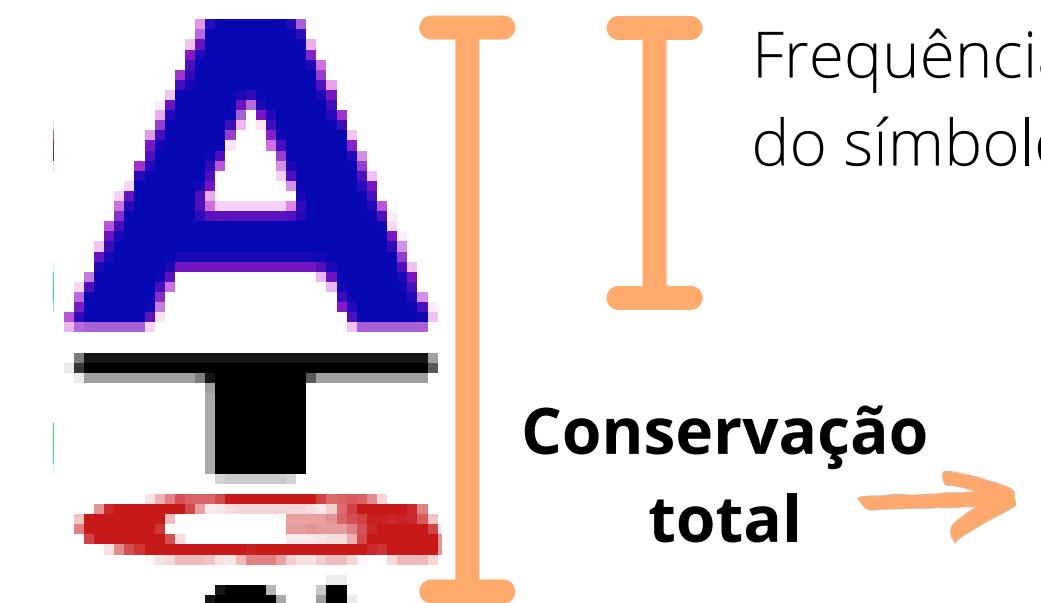
Quanto maior

=

Mais conservado



Posição na sequência



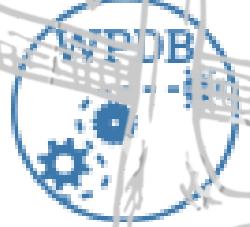
Frequência do símbolo

Quanto mais símbolos na posição

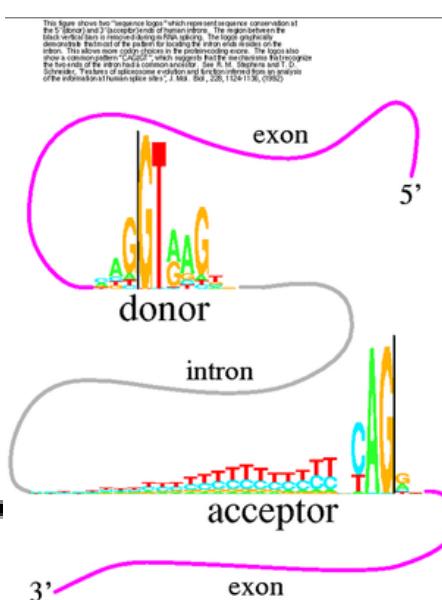
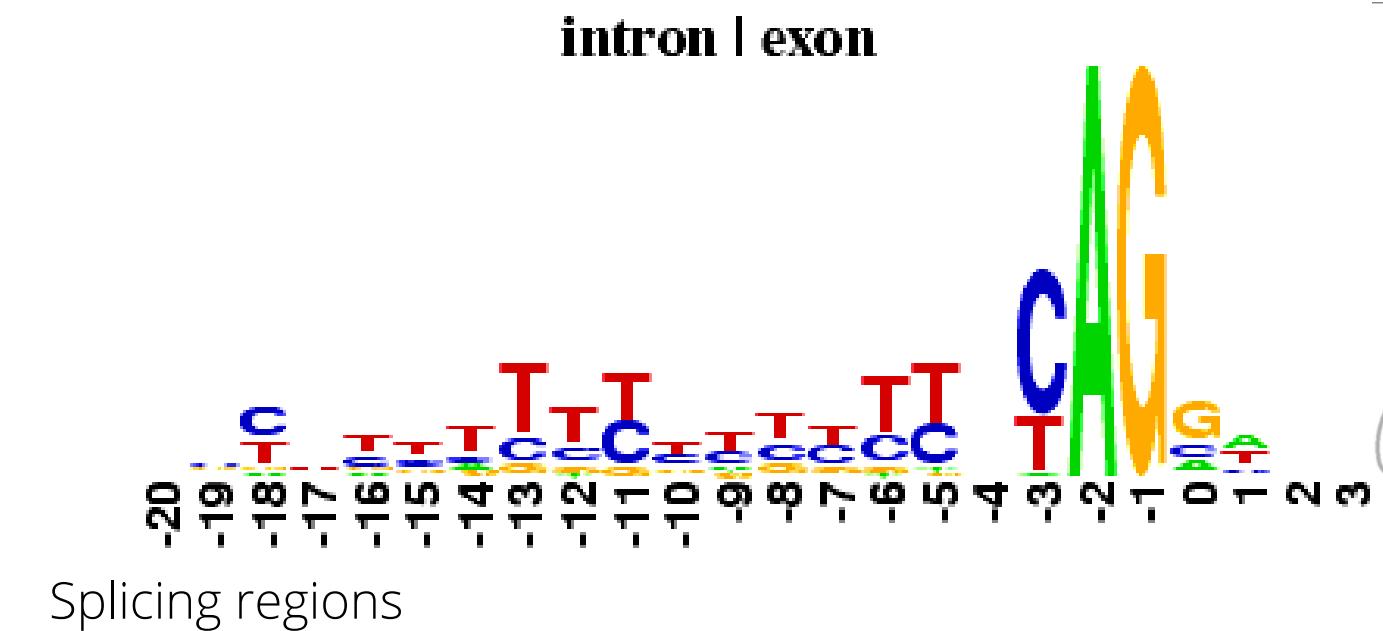
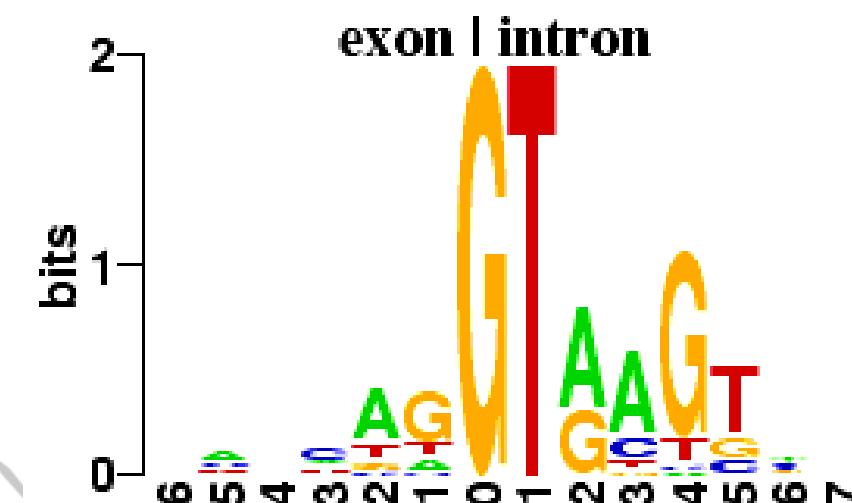
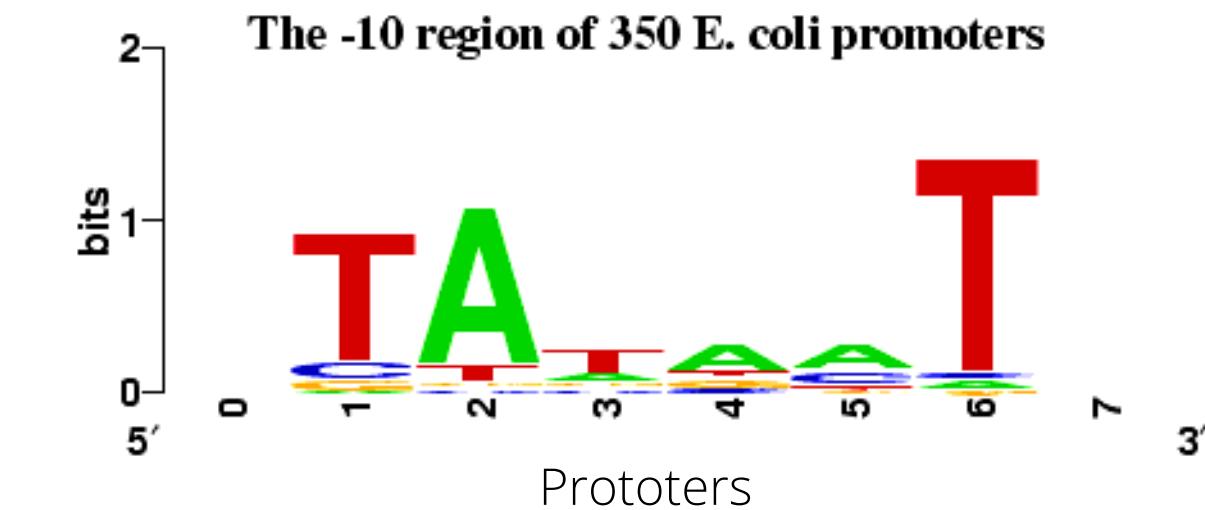
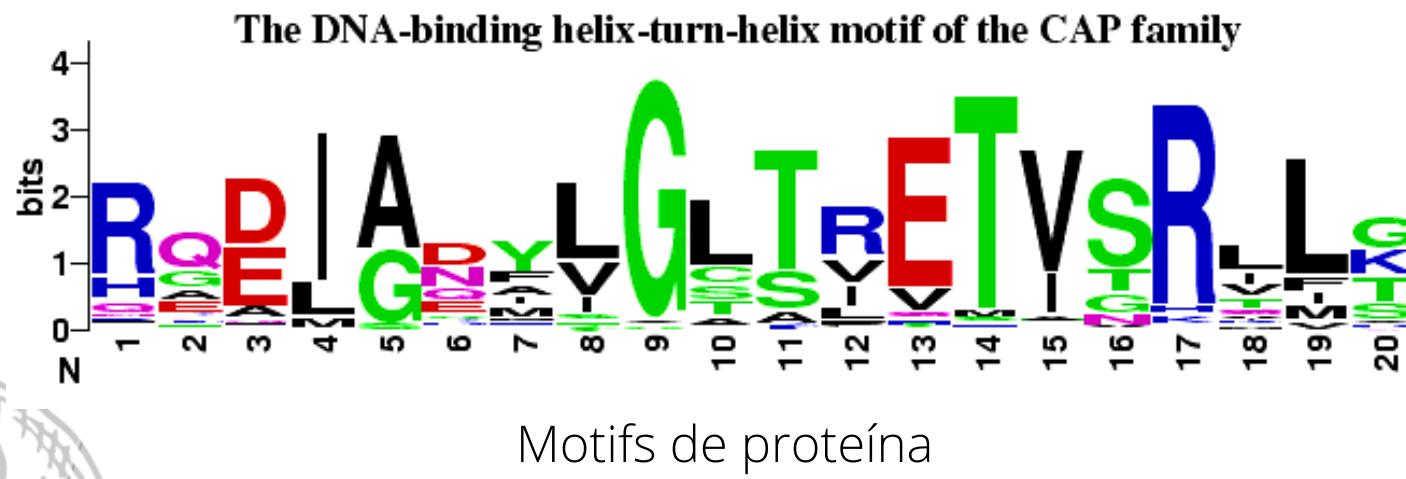
=

Menor é a conservação total

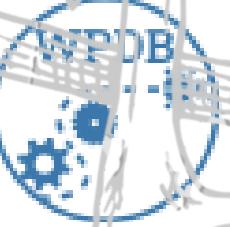
Sequence Logos



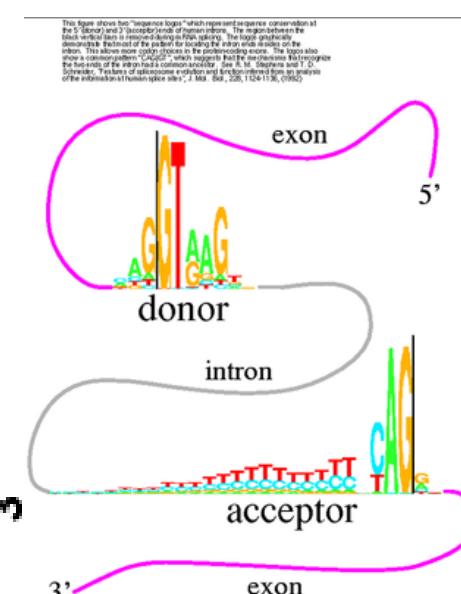
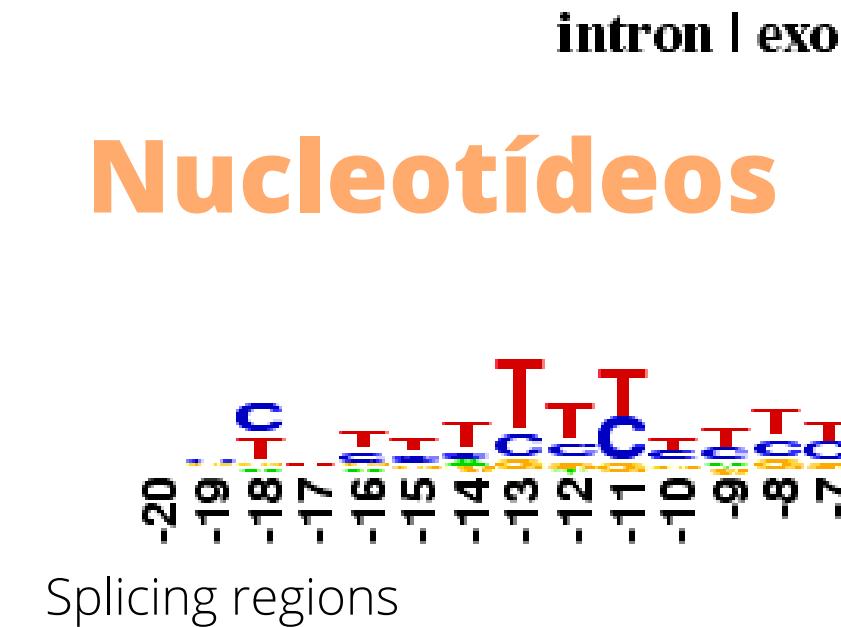
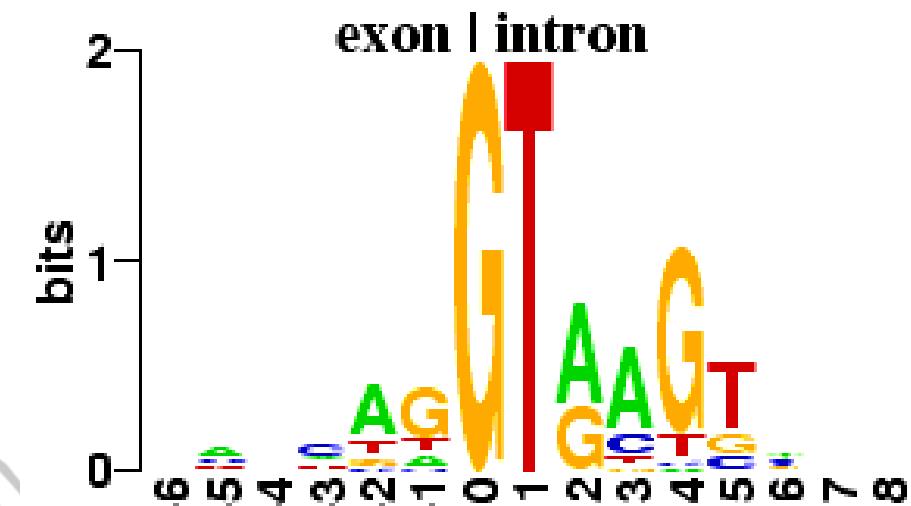
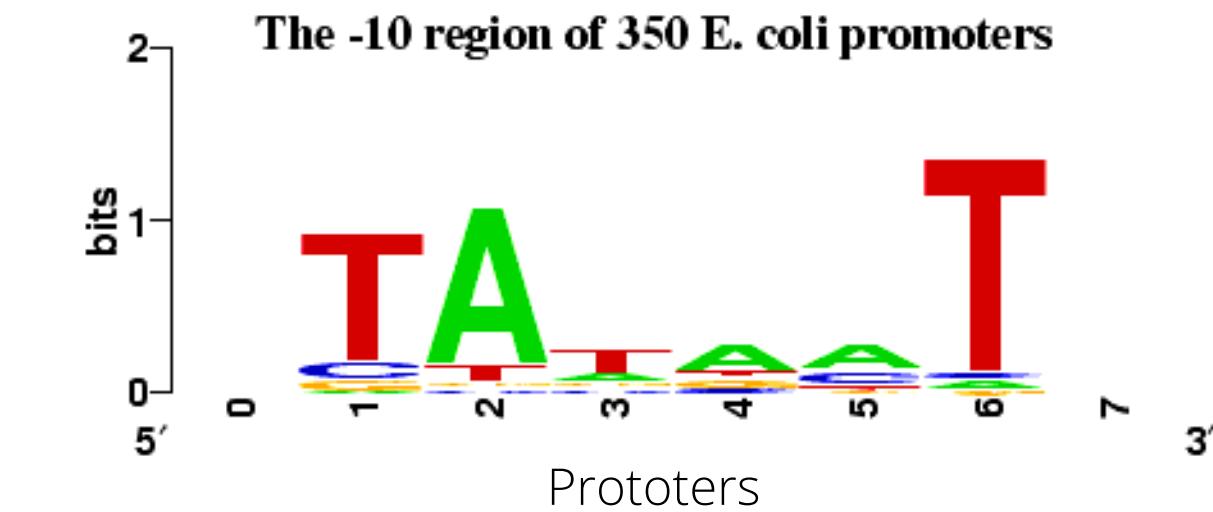
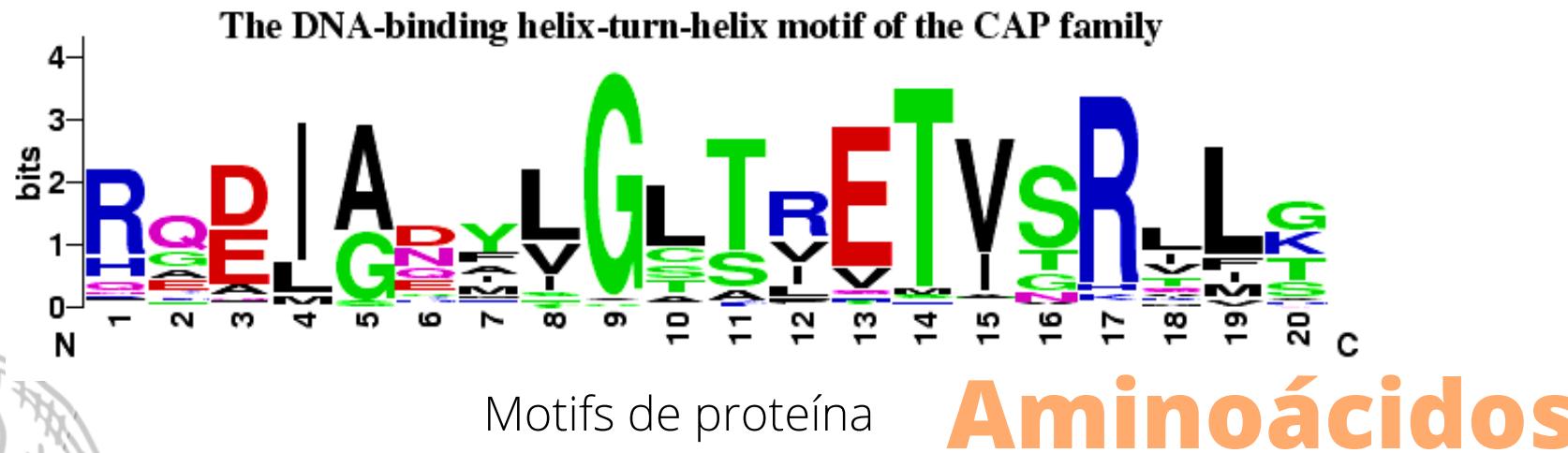
Forma de representação de sequências conservadas introduzida por Schneider and Stephens (1990)



Sequence Logos

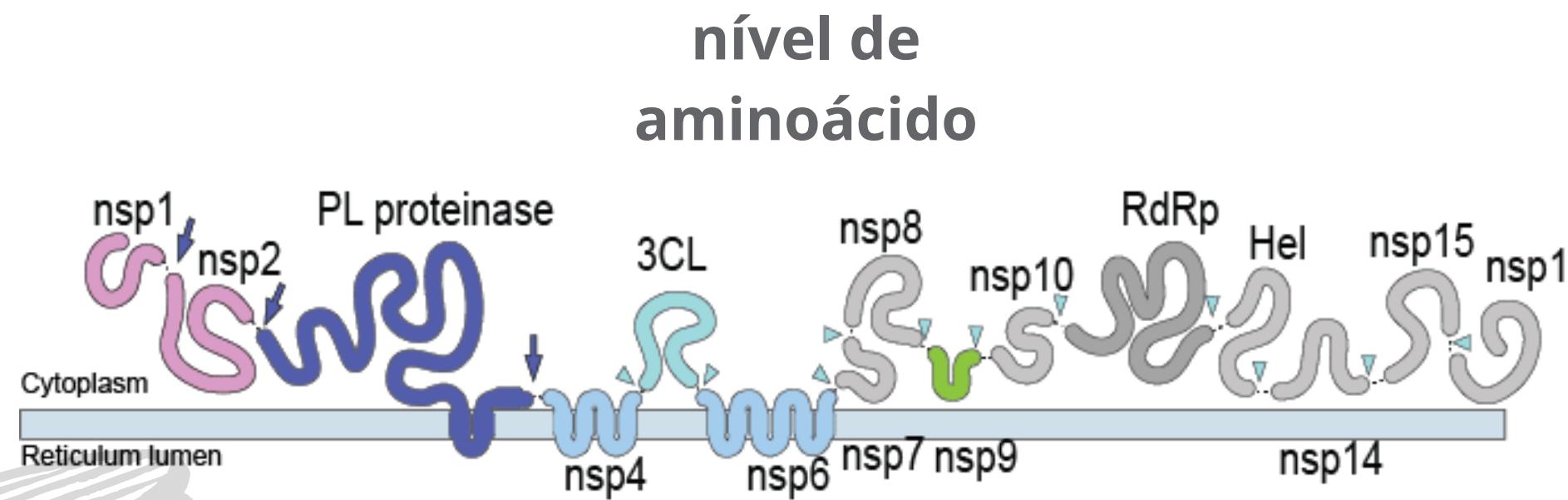
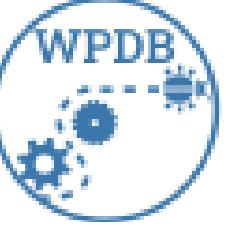


Forma de representação de sequências conservadas introduzida por Schneider and Stephens (1990)



Um caso especial

Vírus recente

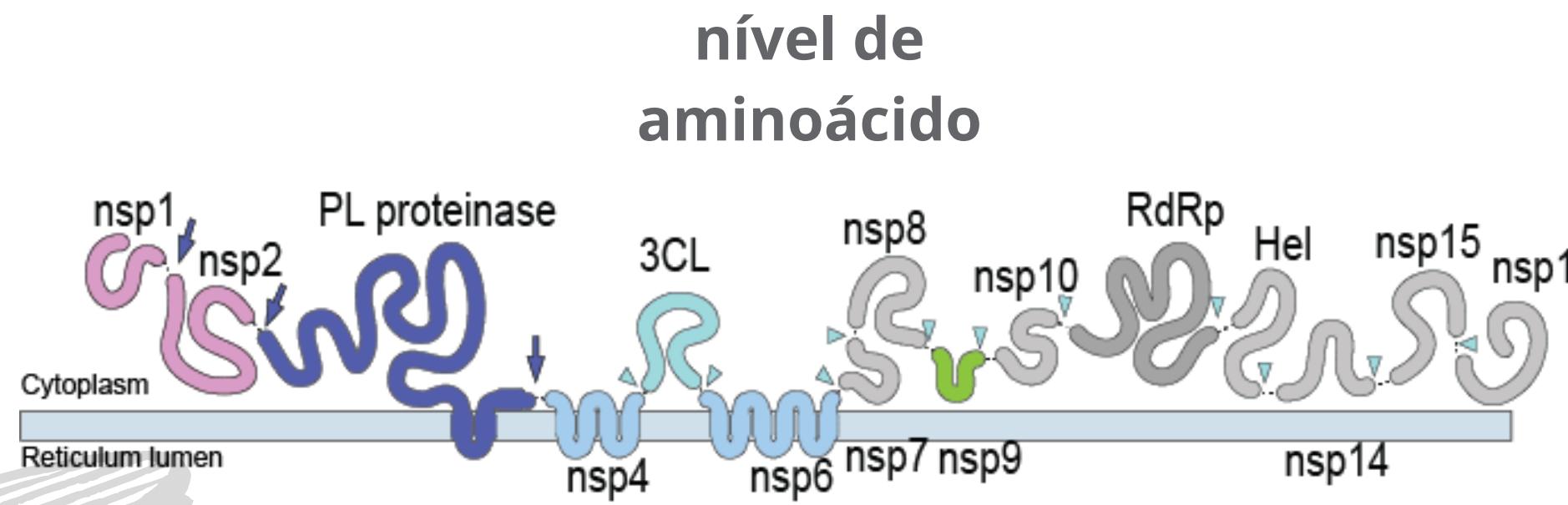
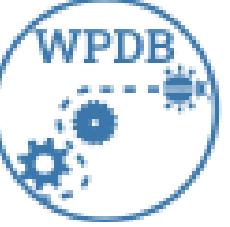


Mutações sinônimas
AGT (Ser) → AGC (Ser)

Mutações não sinônimas
AGT (Ser) → AGA (Arg)

Um caso especial

Vírus recente



Mutações sinônimas

AGT (Ser) → AGC (Ser)

Mutações não sinônimas

AGT (Ser) → AGA (Arg)

Unir o nível de aminoácidos ao nível de nuclétideos = Códons

Por que "O que é possível fazer sabendo o básico de Python?"?



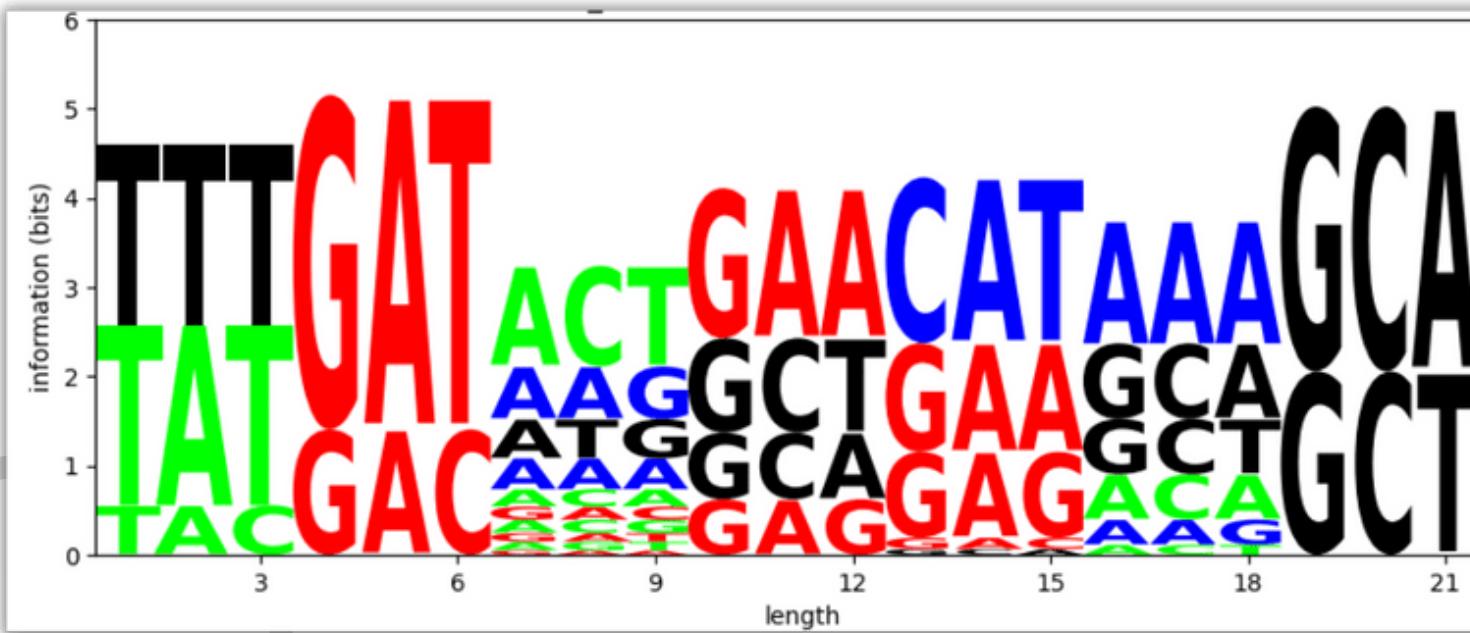
Ferramenta para visualização da conservação
de sequências de códons.

Publicação do artigo
CoCoView



Codon Sequence logo

CoCoView: A Codon Conservation Viewer



labbces/CoCoView

Contributors: 2 Issues: 0 Stars: 1 Forks: 0

labbces/CoCoView

Contribute to labbces/CoCoView development by creating an account on GitHub.

MethodsX

METHOD ARTICLE | VOLUME 9, 101803, JANUARY 01, 2022

CoCoView - A codon conservation viewer via sequence logos

Beatriz Rodrigues Estevam * Diego Mauricio Riaño-Pachón

Open Access • Published: July 30, 2022 • DOI: <https://doi.org/10.1016/j.mex.2022.101803> •

Check for updates

PlumX Metrics

Abstract

Sequence logos are a simple way to display a set of aligned sequences, and they are useful to identify conserved patterns. Since their introduction, several tools have been developed for generating these representations at the single residue level (amino acids or nucleotides). We have developed a tool to build sequence logos of protein-coding sequences at the codon level, allowing more accurate analysis of coding-sequences as they represent synonymous and non-synonymous changes instead of showing only changes that imply acid substitutions. We built CoCoView on top of the Logomaker Python library to generate codon sequence logos from a multiple sequence alignment of protein-coding genes.



Visualização geral

External Libraries

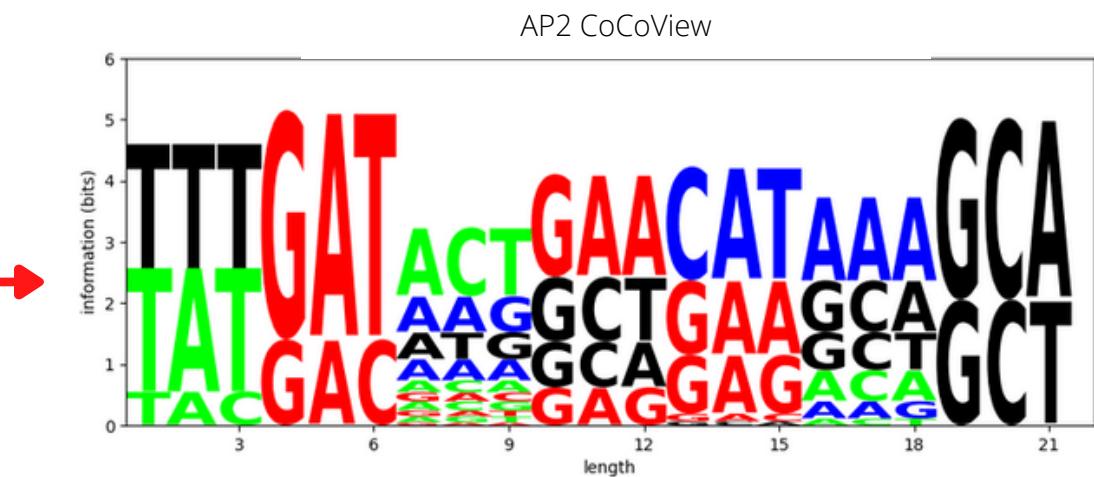
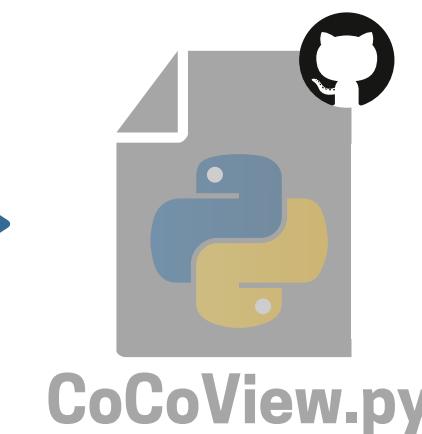
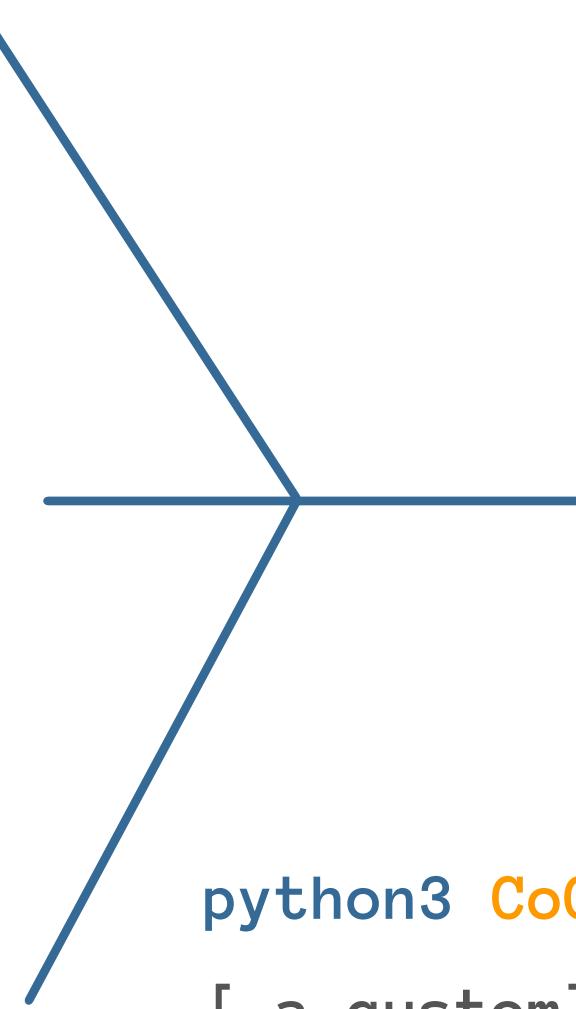
Mandatory



```
> sequence ID 1
AATGGCTGCATGCATGCTGATCGA
> sequence ID 2
AATGGCTGCATGCATGCTGATCGA
> sequence ID 3
AATGGCTGCATGCATGCTGATCGA
> sequence ID 4
AATGGCTGCATGCATGCTGATCGA
> sequence ID 5
AATGGCTGCATGCATGCTGATCGA
```

Graphic features

	Optional
prefixFileName (-p)	degreeOfUncertainty (-d)
imageTitle (-i)	matrixLogoType (-m)
alphaColor (-a)	datasetType (-t)
customPalletFile (-c)	logoFormat (-l)



```
python3 CoCoView.py file.fasta [-p FilesPrefixName] [-i AP2]
[-a custom] [-c customPaletteFile.json] [-d 0.0] [-m bit]
[-t nonredundant] [-l png]
```



Visualização geral



External Libraries

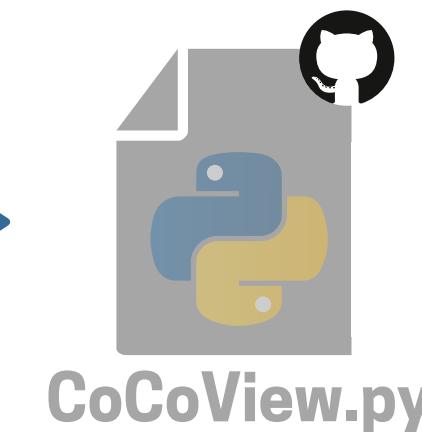
Mandatory



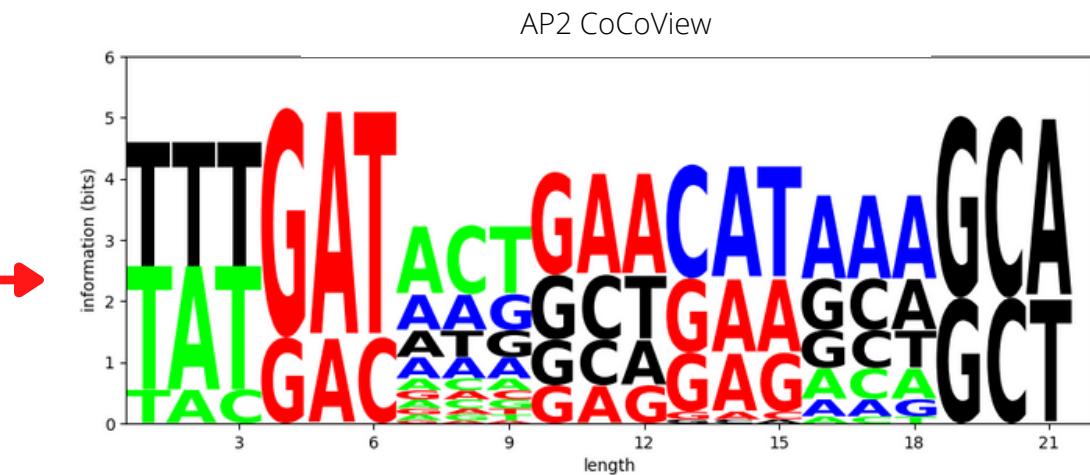
```
> sequence ID 1  
AATGGCTGCATGCATGCTGATCGA  
> sequence ID 2  
AATGGCTGCATGCATGCTGATCGA  
> sequence ID 3  
AATGGCTGCATGCATGCTGATCGA  
> sequence ID 4  
AATGGCTGCATGCATGCTGATCGA  
> sequence ID 5  
AATGGCTGCATGCATGCTGATCGA
```

Graphic features

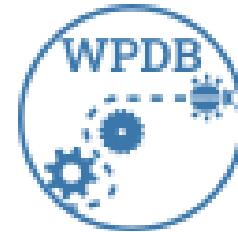
Optional	
<code>prefixFileName (-p)</code>	<code>degreeOfUncertainty (-d)</code>
<code>imageTitle (-i)</code>	<code>matrixLogoType (-m)</code>
<code>alphaColor (-a)</code>	<code>datasetType (-t)</code>
<code>customPalletFile (-c)</code>	<code>logoFormat (-l)</code>



```
python3 CoCoView.py file.fasta [-p FilesPrefixName] [-i AP2]  
[-a custom] [-c customPaletteFile.json] [-d 0.0] [-m bit]  
[-t nonredundant] [-l png]
```



Funcionalidades



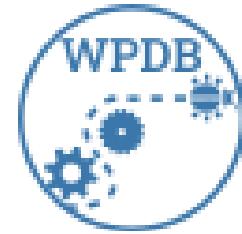
--alphaColor



--customPaletteFile

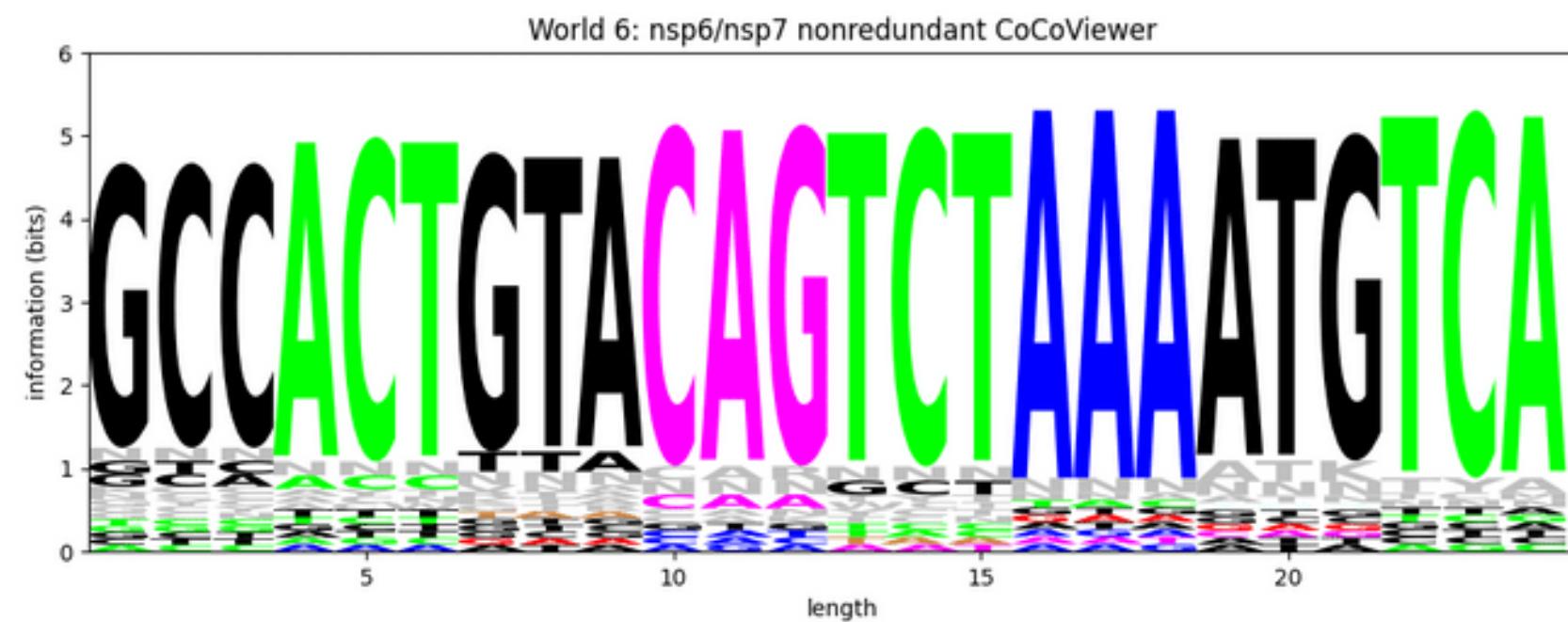
-a custom -c file.json

```
1 { "custom":{  
2     "GCT": "dimgray",  
3     "GCC": "black",  
4     "GCA": "lightcoral",  
5     "GCG": "brown",  
6     "TTT": "red",  
7     "TTC": "salmon",  
8     "ATT": "tomato",  
9     "ATC": "bisque",  
10    "ATA": "burlywood",  
11    "TTA": "tan",  
12    "TTG": "orange",  
13    "CTT": "wheat",  
14    "CTC": "gold",  
15    "CTA": "khaki",  
16    "CTG": "olive",  
17    "ATG": "yellow",  
18    "CCT": "olivedrab",  
19    "CCC": "forestgreen".
```

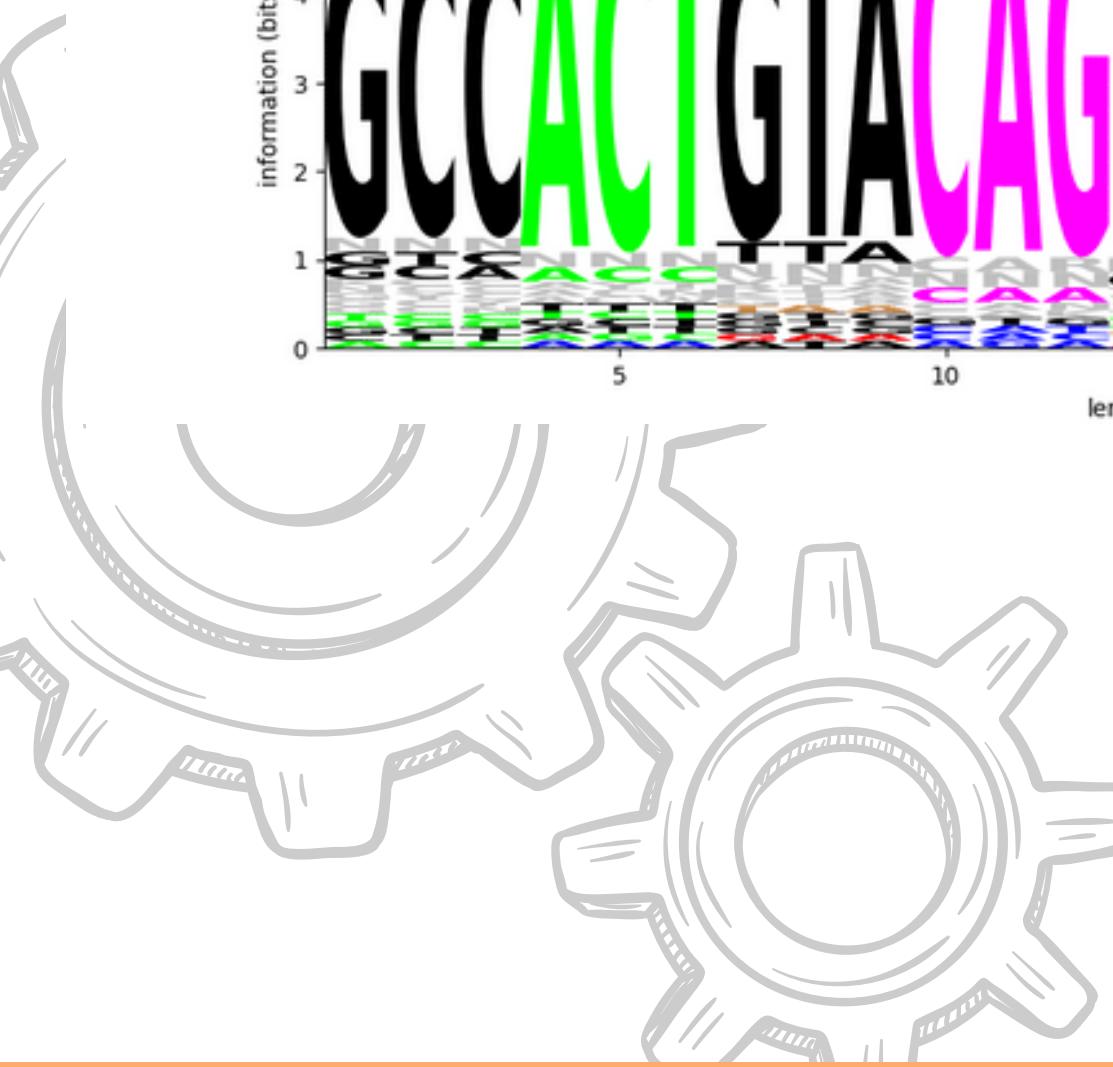
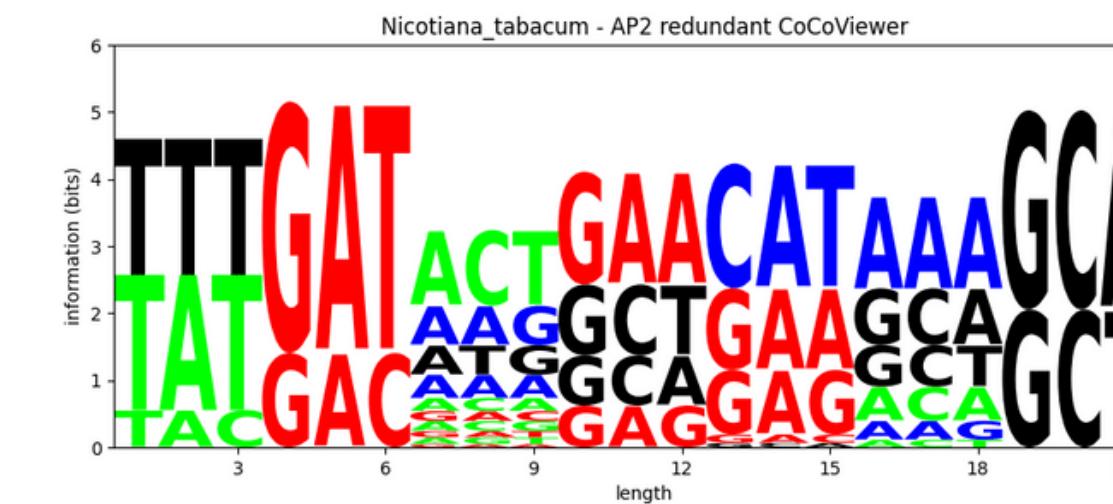
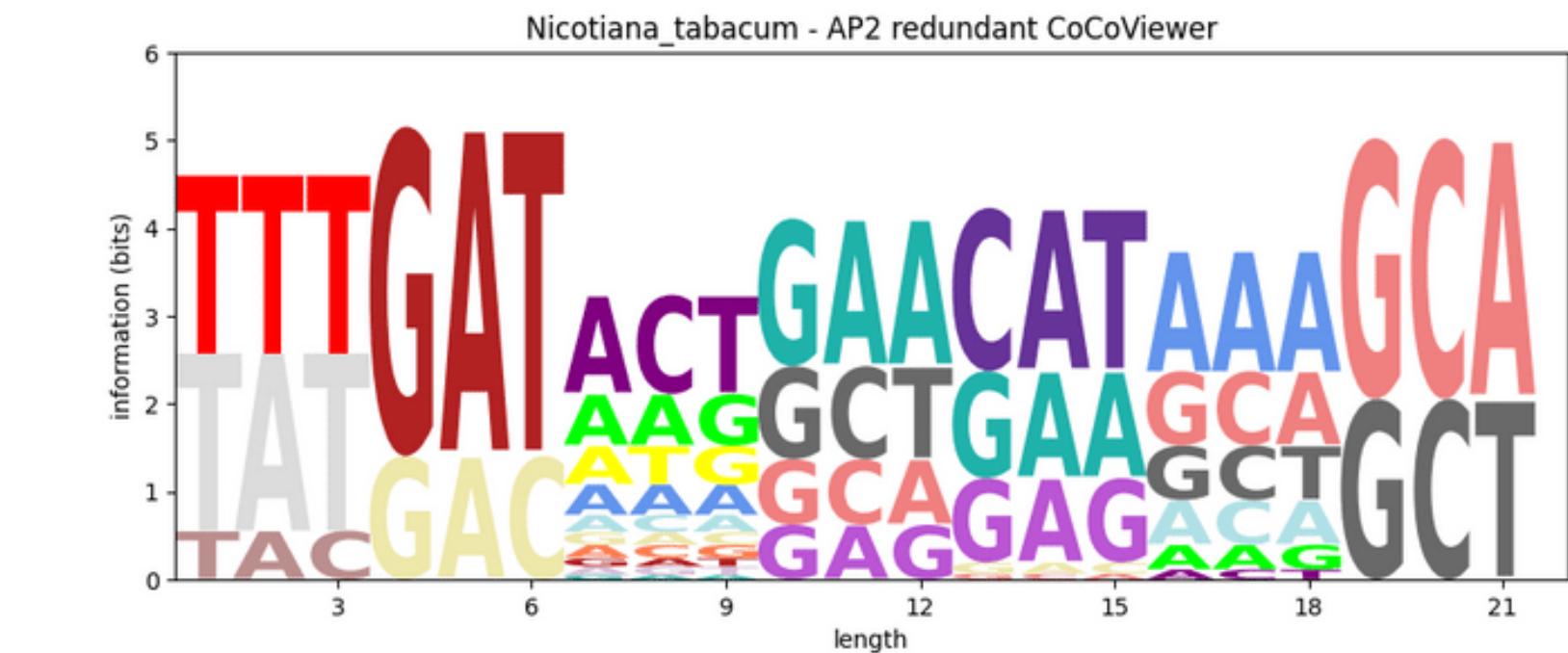


Funcionalidades

--alphaColor



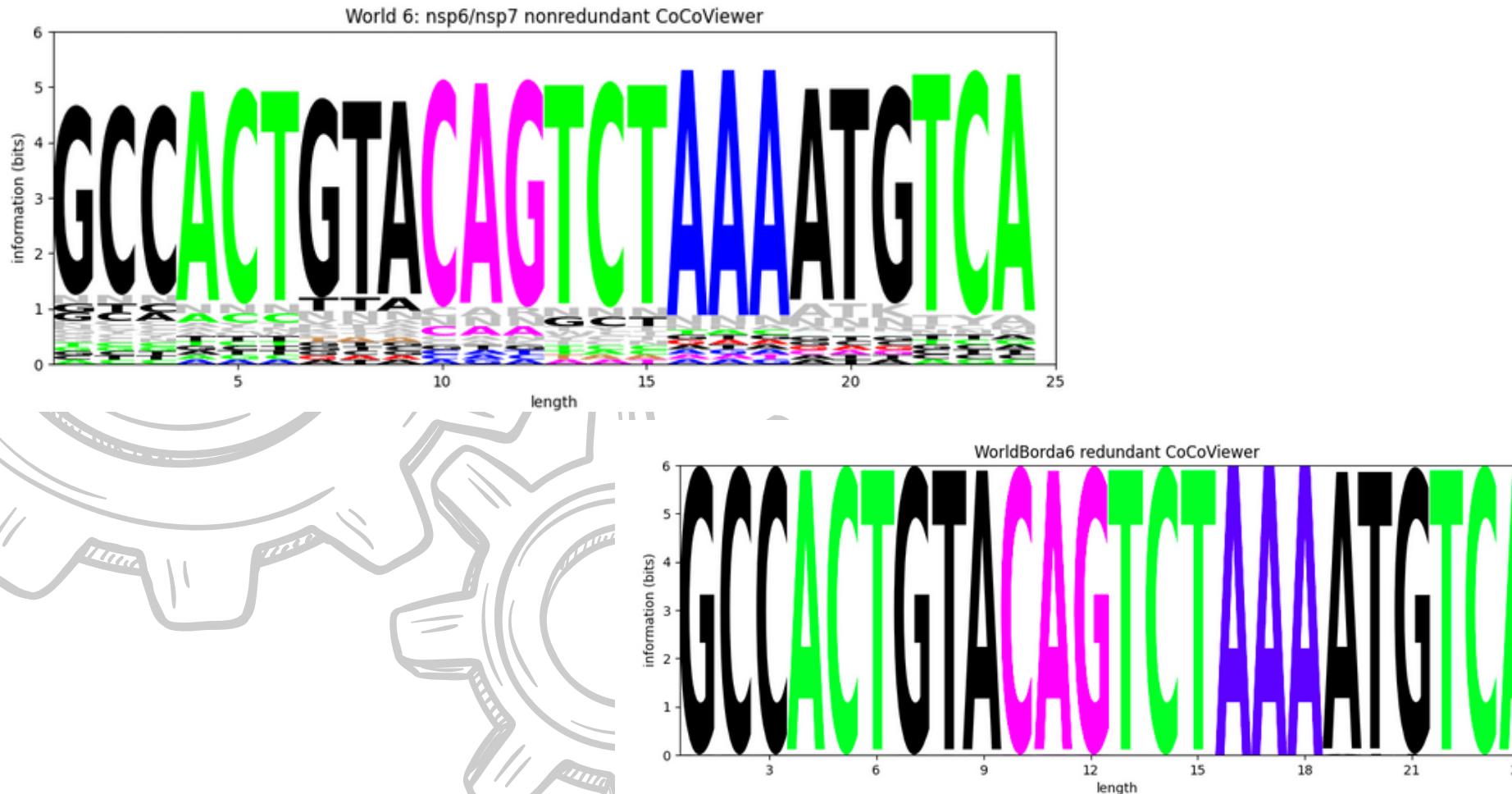
--customPaletteFile



Funcionalidades

--datasetType

- Datasets pequenos
- Redundante x Não redundante
- Remover sequências duplicadas



--degreeOfUncertainty

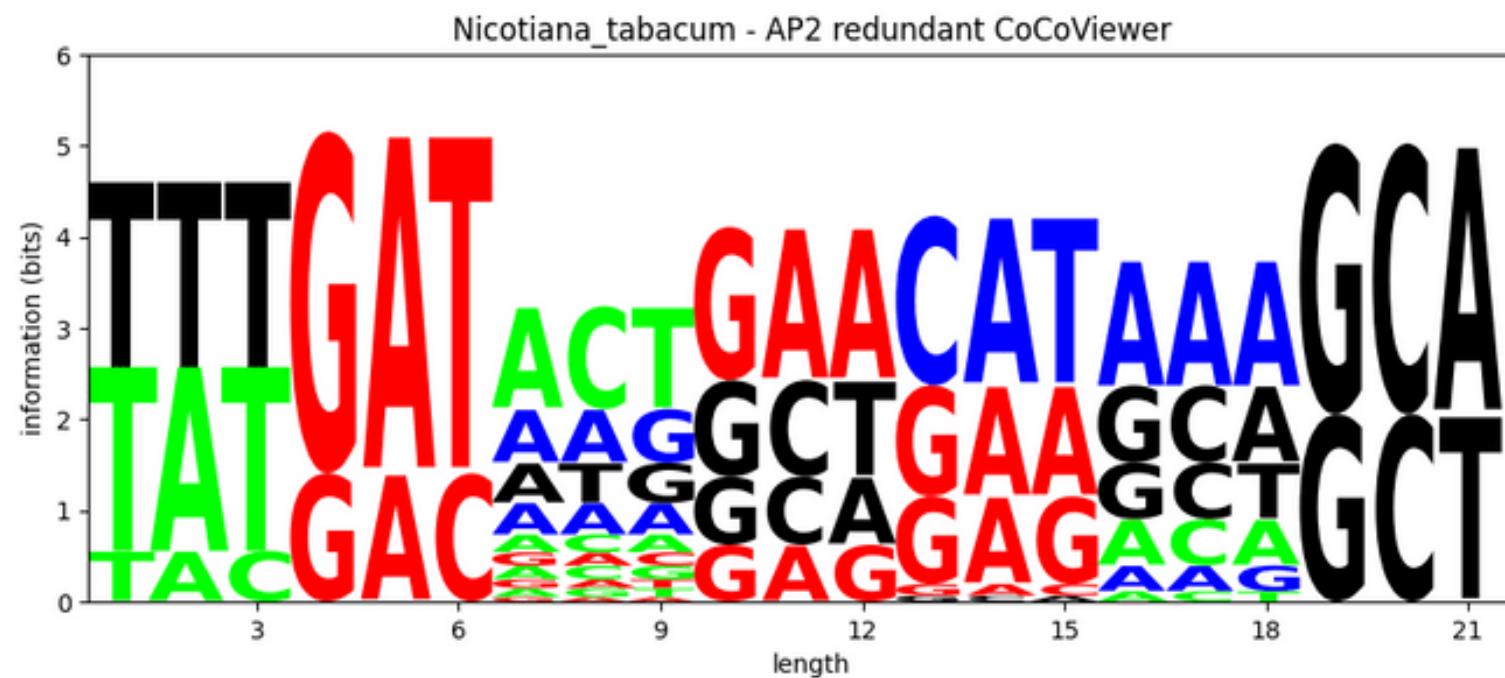
- Remover sequências (parcialmente) desconhecidas



Funcionalidades



--imageTitle



Nicotiana_tabacum - AP2 redundant CoCoViewer

--prefixFileName



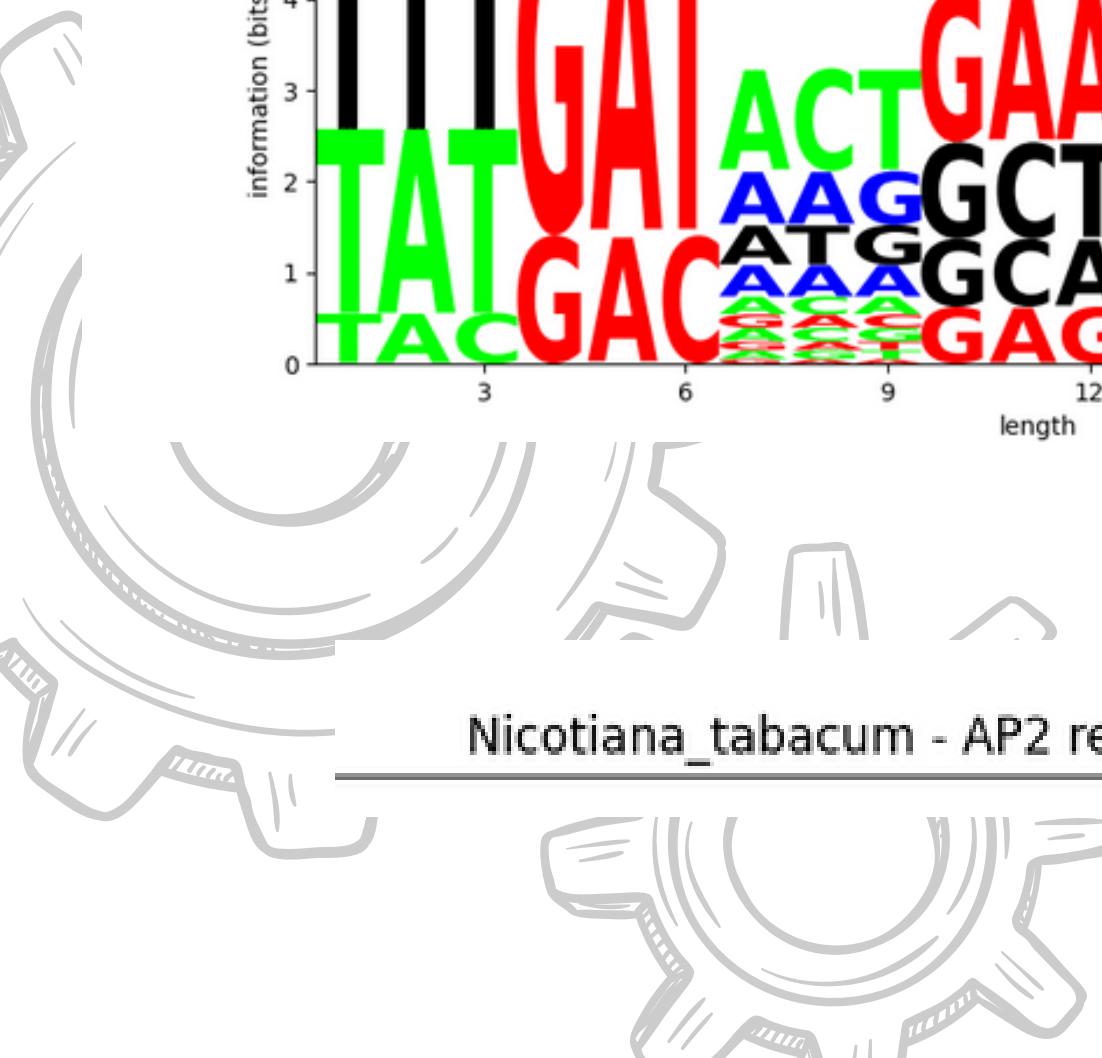
prefixFileName + datasetType +
"bits.matrix"



prefixFileName + datasetType +
"probability.matrix"



prefixFileName + datasetType +
matrixLogoType + "CoCoView" +
logoFormat



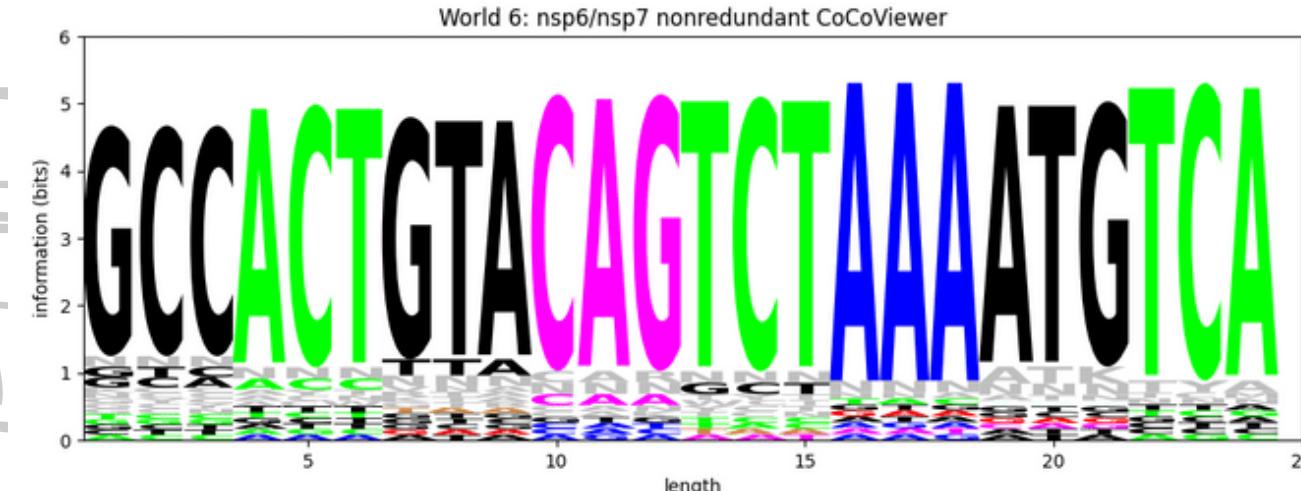
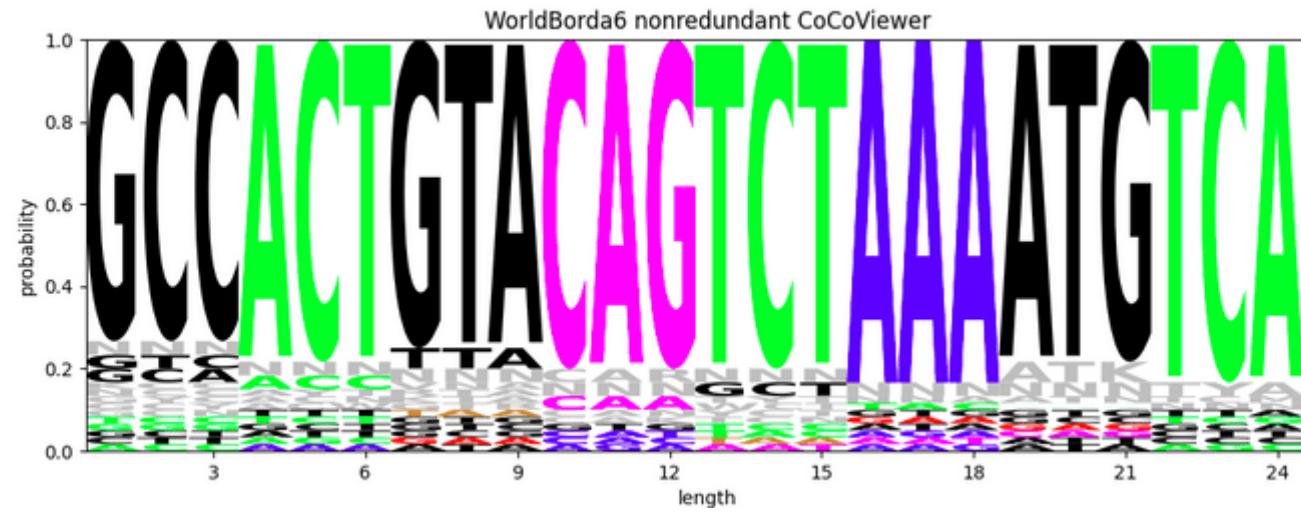
CoCoView: A Codon Conservation Viewer

Funcionalidades



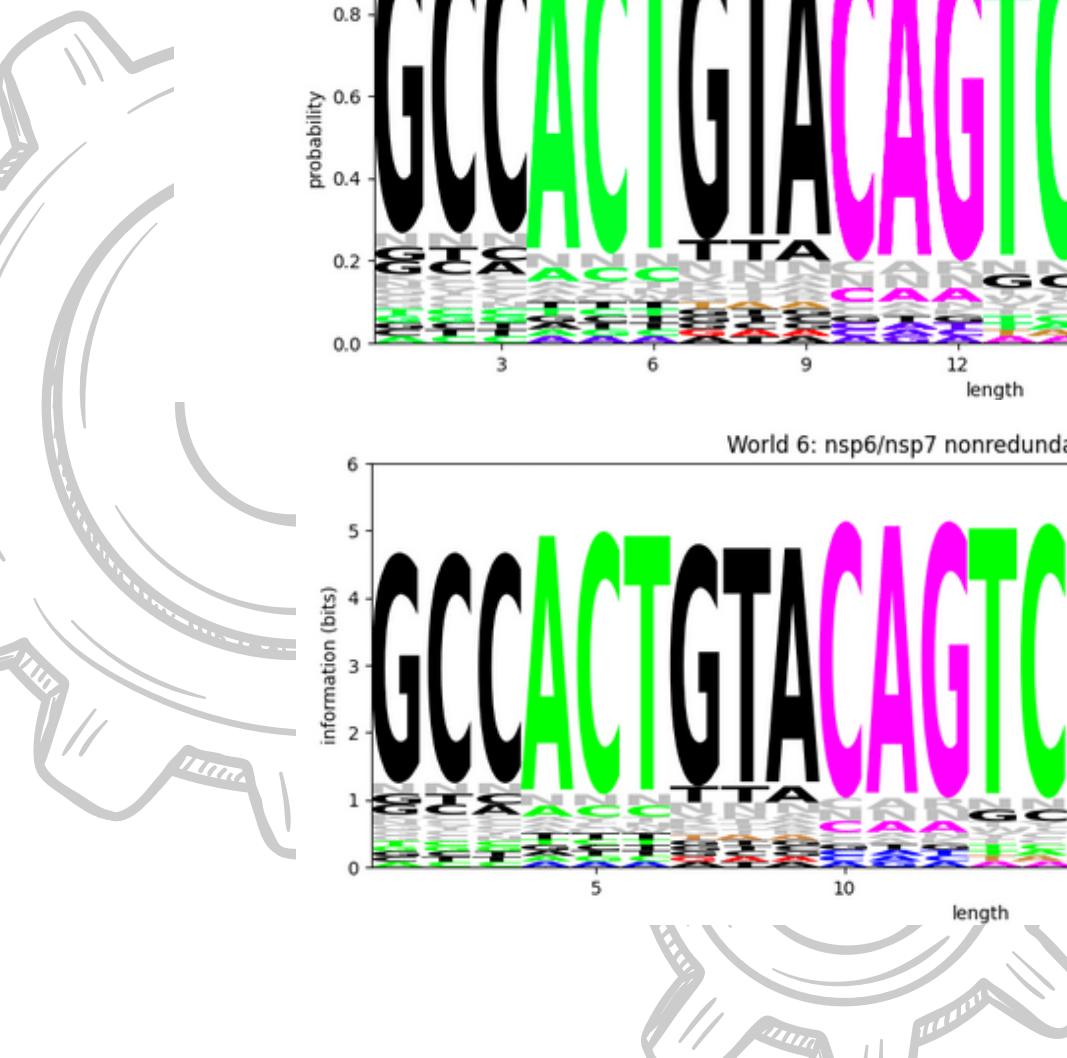
--matrixLogoType

- Frequência x Bits



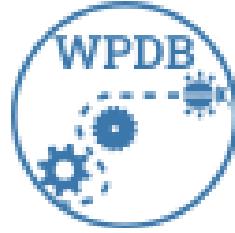
--logoFormat

- png e pdf

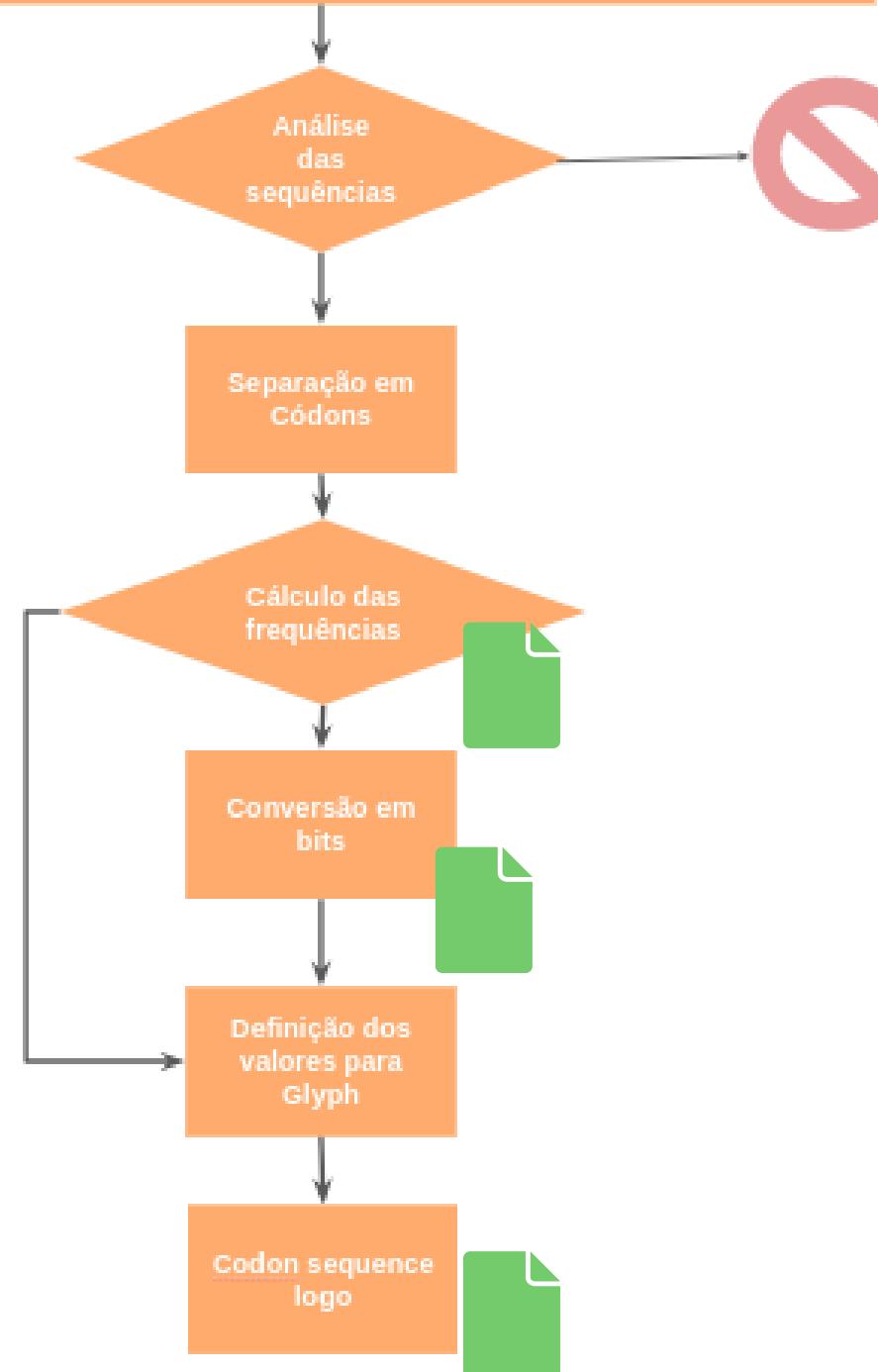


Codon Sequence logo

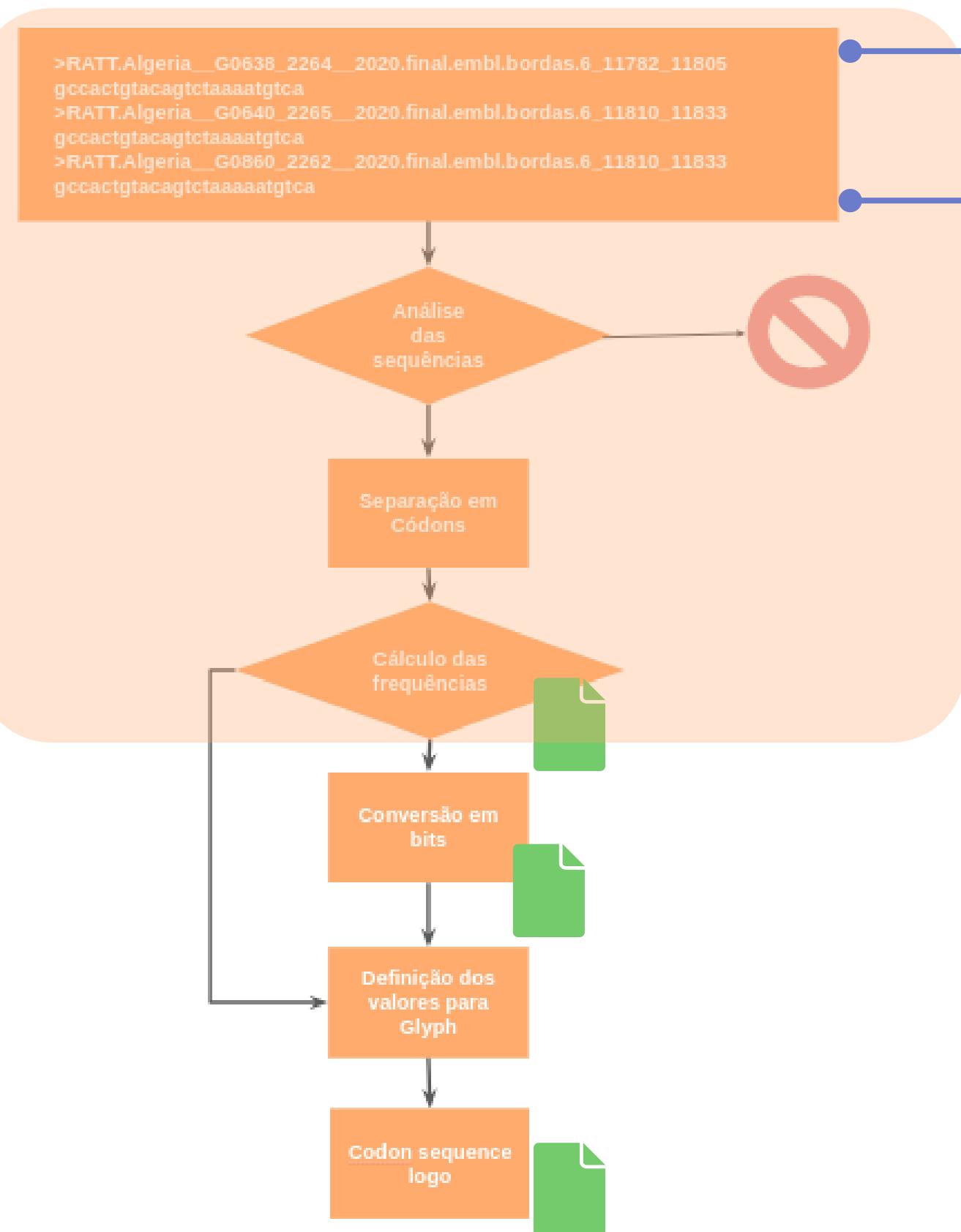
Script overview



```
>RATT.Algeria_G0638_2264_2020.final.embl.bordas.6_11782_11805  
gcactgtacagtctaaaaatgtca  
>RATT.Algeria_G0640_2265_2020.final.embl.bordas.6_11810_11833  
gcactgtacagtctaaaaatgtca  
>RATT.Algeria_G0860_2262_2020.final.embl.bordas.6_11810_11833  
gcactgtacagtctaaaaatgtca
```



Codon Sequence logo Script overview



argparse



Comunicação do script com o terminal

biopython



Operações com as seqüências

MÓDULO AlignIO

- Lê o FASTA
- Pega o tamanho da seqüência
- Seleciona as seqüências dentro do FASTA

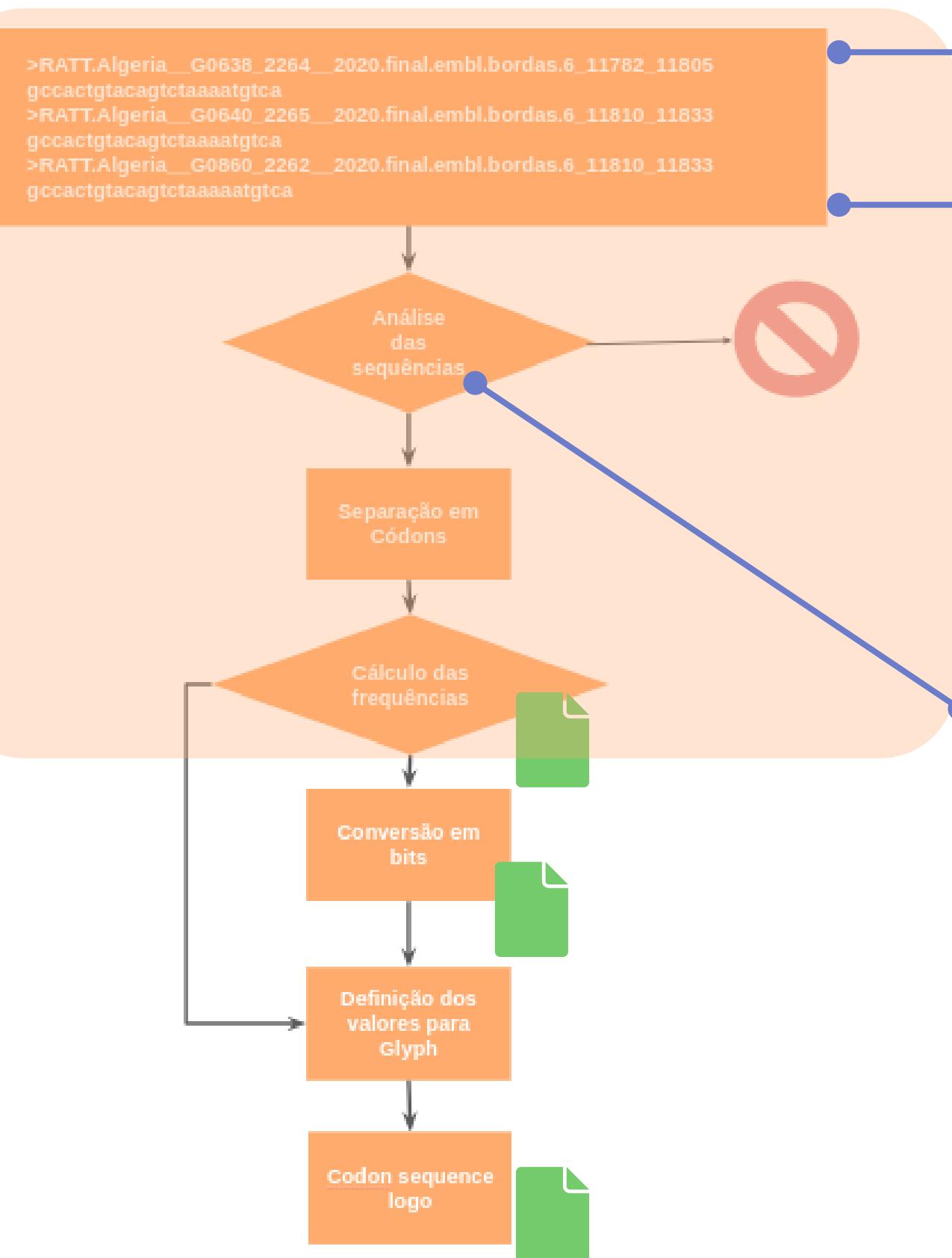
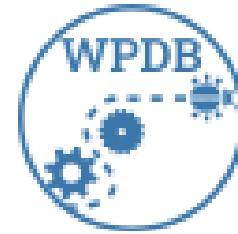
```
alignment = AlignIO.read(fastaFile, "fasta")
SeqLength = alignment.get_alignment_length()

if SeqLength % 3 != 0:
    print(f'Your sequence length {args.SeqLength} is not a multiple of 3')
    exit()
```

```
for record in alignment:
    notKnownNucleotide = 0
```



Codon Sequence logo Script overview



argparse



Comunicação do script com o terminal

biopython



Operações com as seqências

MÓDULO AlignIO

- Lê o FASTA
- Pega o tamanho da seqêncie
- Seleciona as seqências dentro do FASTA

AAAAAAA
AAAAAAAXXX
AAAXXXXXXXX

--degreeOfUncertainty
--datasetType



```
alignment = AlignIO.read(fastaFile, "fasta")
SeqLength = alignment.get_alignment_length()
```

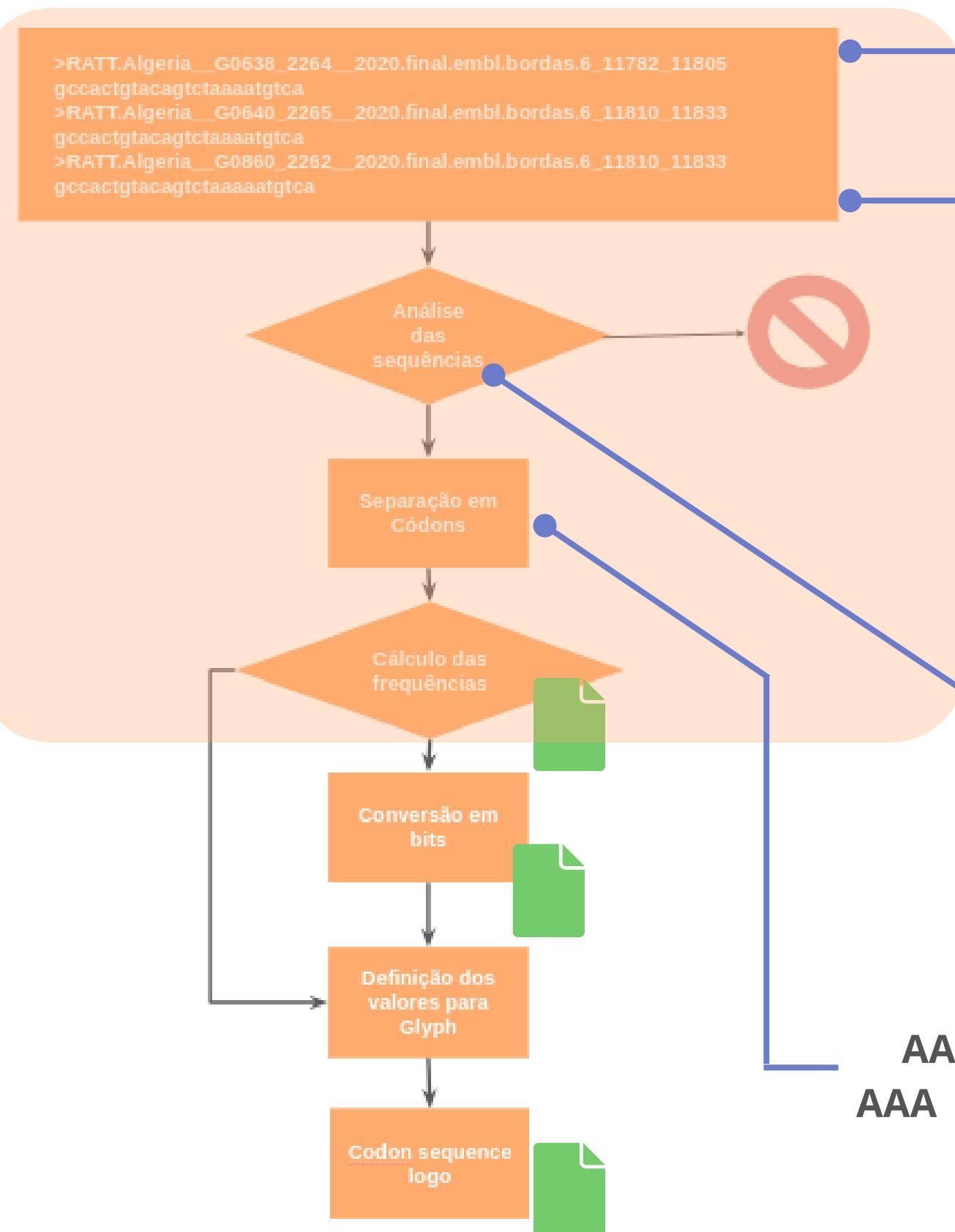
```
if SeqLength % 3 != 0:
    print(f'Your sequence length {args.SeqLength} is not a multiple of 3')
    exit()
```

```
for record in alignment:
    notKnownNucleotide = 0
```



Codon Sequence logo

Script overview



argparse



Comunicação do script com o terminal

biopython



Operações com as seqüências

MÓDULO AlignIO

- Lê o FASTA
- Pegar o tamanho da seqüência
- Seleciona as seqüências dentro do FASTA

```
alignment = AlignIO.read(fastaFile, "fasta")
SeqLength = alignment.get_alignment_length()
```

```
if SeqLength % 3 != 0:
    print(f'Your sequence length {args.SeqLength} is not a multiple of 3')
    exit()
```

```
for record in alignment:
    notKnownNucleotide = 0
```

--degreeOfUncertainty
--datasetType



AAAAAAAAXXXX
AAAAAAAAXXXX
AAAXXXXXXXX

AAAAAAAAXXXX
AAA AAA AAX XXX

seq 1	AAG	AAT	AAX	XXX
seq 2	AAG	AAA	AAX	XXX
seq 3	AAA	AAA	AAX	XXX
seq 4	AAA	AAA	AAX	XXX

Dicionário= { posição: {códon : frequência} }

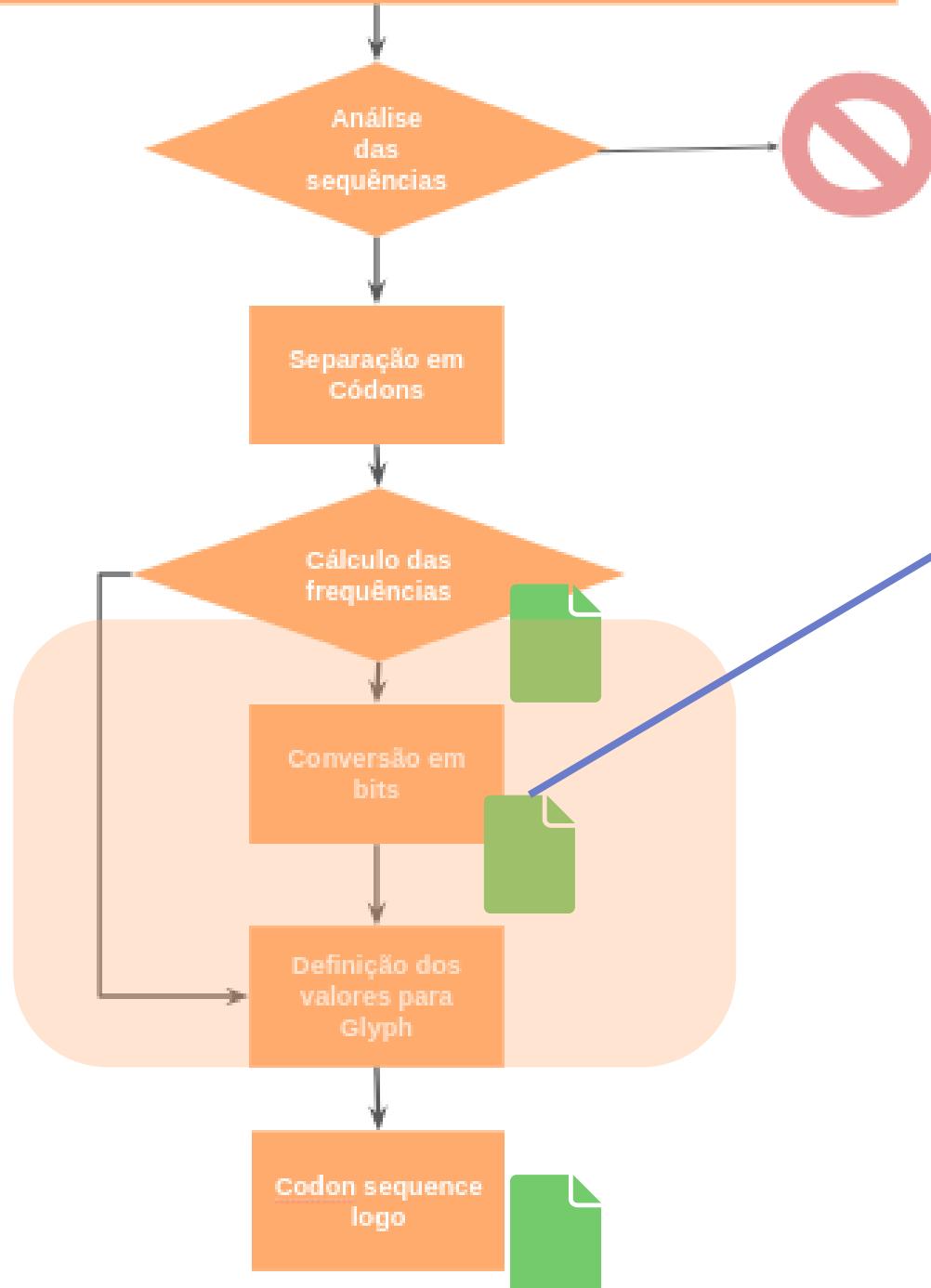


Codon Sequence logo

Script overview



```
>RATT.Algeria_G0638_2264_2020.final.embl.bordas.6_11782_11805  
gccactgtacagtctaaaaatgtca  
>RATT.Algeria_G0640_2265_2020.final.embl.bordas.6_11810_11833  
gccactgtacagtctaaaaatgtca  
>RATT.Algeria_G0860_2262_2020.final.embl.bordas.6_11810_11833  
gccactgtacagtctaaaaatgtca
```



4.2.3 matrix functions

`logomaker.transform_matrix(*args, **kwargs)`

matriz de probabilidade --> Matriz de Bits



4.2.2 Glyph class

`class logomaker.Glyph(**kwargs)`

A Glyph represents a character, drawn on a specified axes at a specified position, rendered using specified styling such as color and font_name.

Attributes

p: (number) x-coordinate value on which to center the Glyph.

c: (str) The character represented by the Glyph.

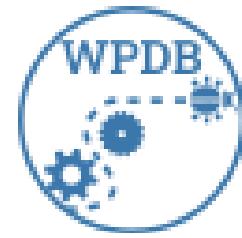
floor: (number) y-coordinate value where the bottom of the Glyph extends to. Must be < ceiling.

ceiling: (number) y-coordinate value where the top of the Glyph extends to. Must be > floor.

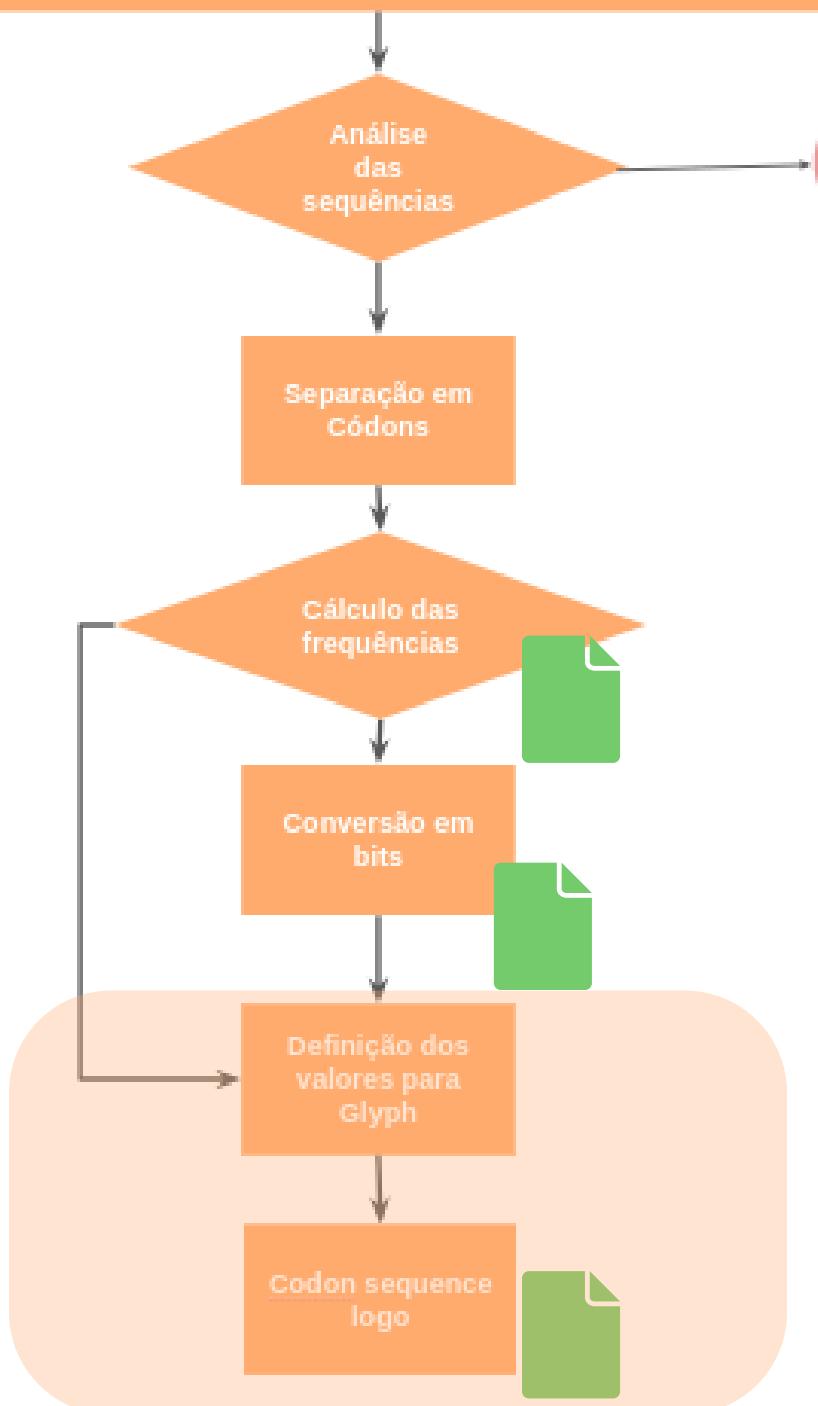
ax: (matplotlib Axes object) The axes object on which to draw the Glyph.

Codon Sequence logo

Script overview



```
>RATT.Algeria_G0638_2264_2020.final.embl.bordas.6_11782_11805  
gcacactgtacagtctaaaaatgtca  
>RATT.Algeria_G0640_2265_2020.final.embl.bordas.6_11810_11833  
gcacactgtacagtctaaaaatgtca  
>RATT.Algeria_G0860_2262_2020.final.embl.bordas.6_11810_11833  
gcacactgtacagtctaaaaatgtca
```



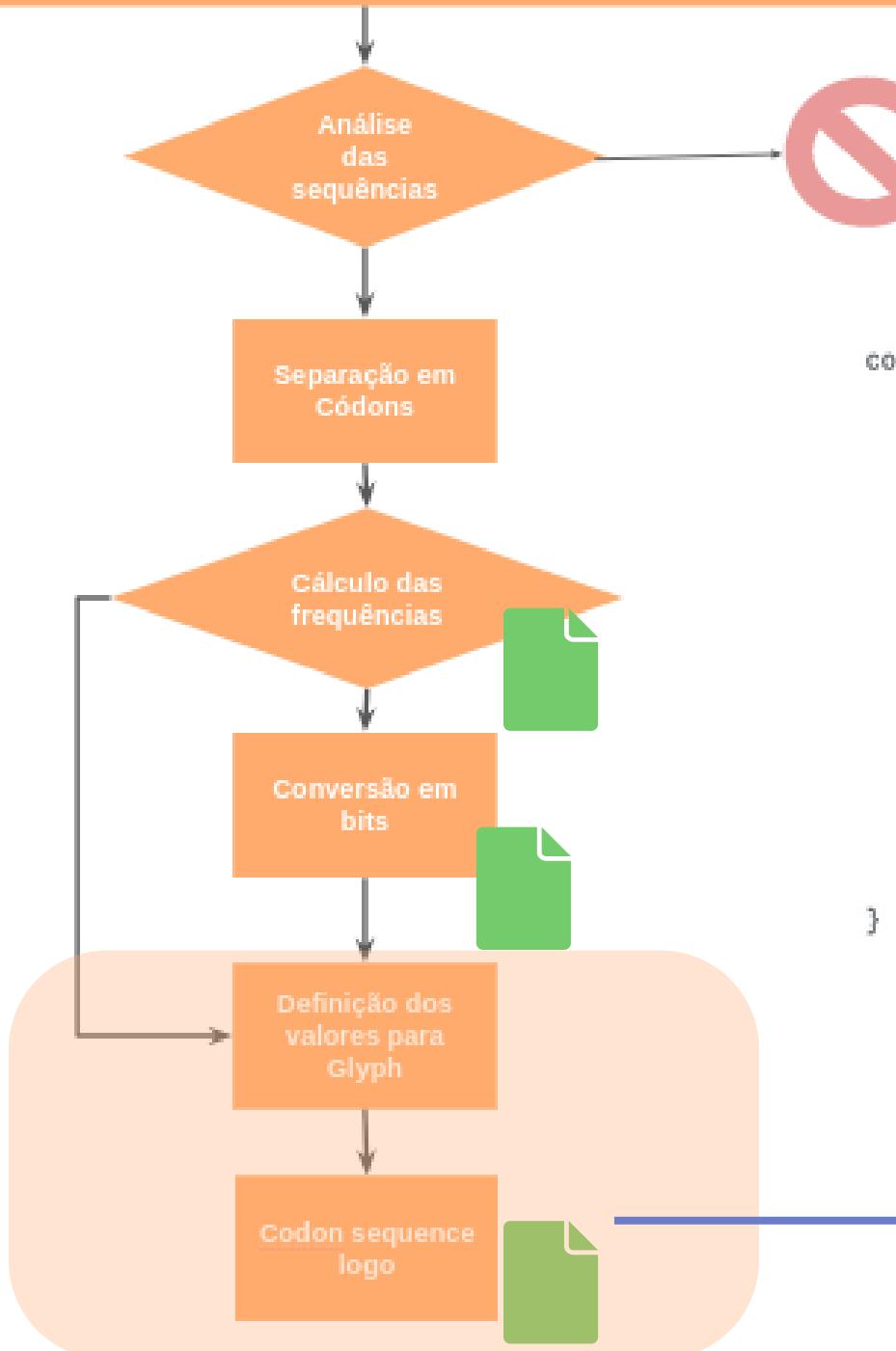
- Onde cada símbolo entra no gráfico;
- Cor, nome, formatos (.pdf, .png).
- --alphaColor, --customPaletteFile, --imageTitle, --prefixFileName, --logoFormat

Codon Sequence logo

Script overview



```
>RATT.Algeria_G0638_2264_2020.final.embl.bordas.6_11782_11805  
gcacactgtacagtctaaaaatgtca  
>RATT.Algeria_G0640_2265_2020.final.embl.bordas.6_11810_11833  
gcacactgtacagtctaaaaatgtca  
>RATT.Algeria_G0860_2262_2020.final.embl.bordas.6_11810_11833  
gcacactgtacagtctaaaaatgtca
```

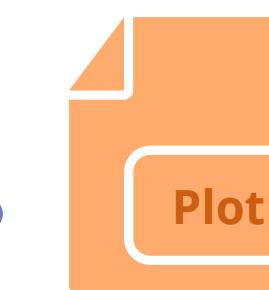


- Onde cada símbolo entra no gráfico;
- Cor, nome, formatos (.pdf, .png).
- --alphaColor, --customPaletteFile, --imageTitle, --prefixFileName, --logoFormat

```
,  
color_palett = {  
    'weblogo_protein':  
        ('GCT': 'black', 'GCC': 'black', 'GCA': 'black', 'GCG': 'black', 'TTT': 'black', 'TTC': 'black', 'ATT': 'black', 'ATC': 'black', 'ATA': 'black', 'TTA': 'bla  
        'TCC': 'lime', 'TCA': 'lime', 'TCG': 'lime', 'AGT': 'lime', 'AGC': 'lime', 'ACT': 'lime', 'ACC': 'lime', 'ACA': 'lime', 'ACG': 'lime', 'TAT': 'lime', 'T  
        'charge':  
            ('GCT': 'dimgrey', 'GCC': 'dimgrey', 'GCA': 'dimgrey', 'GCG': 'dimgrey', 'TTT': 'dimgrey', 'TTC': 'dimgrey', 'ATT': 'dimgrey', 'ATC': 'dimgrey', 'ATA': 'dim  
            'TCC': 'dimgrey', 'TCA': 'dimgrey', 'TCG': 'dimgrey', 'AGT': 'dimgrey', 'AGC': 'dimgrey', 'ACT': 'dimgrey', 'ACC': 'dimgrey', 'ACA': 'dimgrey', 'ACG': 'd  
        'hydrophobicity':  
            ('GCT': 'green', 'GCC': 'green', 'GCA': 'green', 'GCG': 'green', 'TTT': 'black', 'TTC': 'black', 'ATT': 'black', 'ATC': 'black', 'ATA': 'black', 'TTA': 'bla  
            'AGT': 'green', 'AGC': 'green', 'ACT': 'green', 'ACC': 'green', 'ACA': 'green', 'ACG': 'green', 'TAT': 'black', 'TAC': 'black', 'GAT': 'mediumblue', 'GA  
        'chemistry':  
            ('GCT': 'black', 'GCC': 'black', 'GCA': 'black', 'GCG': 'black', 'TTT': 'black', 'TTC': 'black', 'ATT': 'black', 'ATC': 'black', 'ATA': 'black', 'TTA': 'bla  
            'TCG': 'lime', 'AGT': 'lime', 'AGC': 'lime', 'ACT': 'lime', 'ACC': 'lime', 'ACA': 'lime', 'ACG': 'lime', 'TAT': 'lime', 'TAC': 'lime', 'GAT': 'red', 'GA
```



Matplotlib

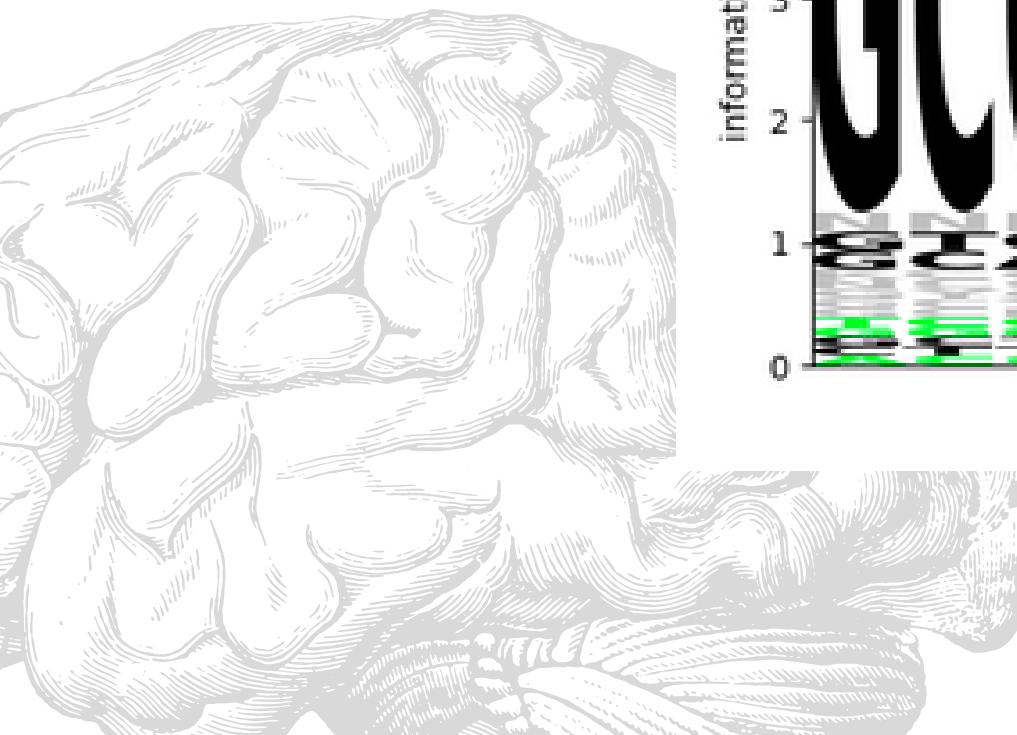
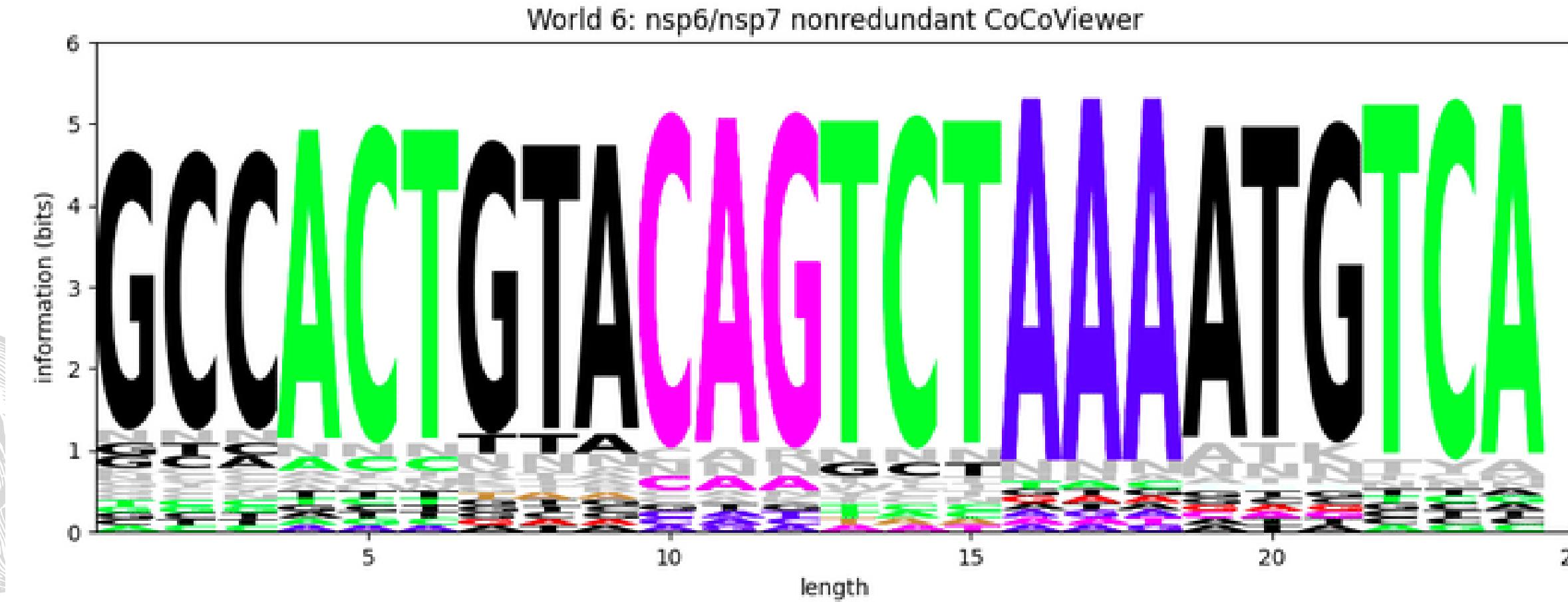


CoCoView: A Codon Conservation Viewer

Exemplo



```
python3 CoCoView.py ./CoCoView/examples/example1/Bordas_World_6.nt.fasta -p "World6" -i  
"World 6 : nsp6/nsp7 nonredundant" -a "weblogo_protein" -d 100 -m "bit" -t "nonredundant"  
-l "png"
```



CoCoView: A Codon Conservation Viewer

Exemplo



The screenshot shows the GitHub repository page for `labbces/CoCoView`. The repository has 1 branch and 0 tags. The main file listed is `README.md`, which has been updated by `Beatriz-Estevam` on 29 Jul. The repository is public and has 1 star, 1 watching, and 0 forks. The `About` section indicates no description, website, or topics provided. The `Packages` section shows no packages published, with a link to "Publish your first". The right sidebar contains the `About`, `Releases`, and `Packages` sections. A red hand cursor is pointing at the `Code` button in the top right corner of the main repository area.

`labbces / CoCoView` Public

Code Issues Pull requests Actions Projects Wiki Security Insights Settings

main · 1 branch · 0 tags

Beatriz-Estevam Update README.md zfeesse on 29 Jul 47 commits

- examples Adding new example files and images 2 months ago
- images Adding new example files and images 2 months ago
- test add example1 and update principal script 5 months ago
- CoCoView.py correcting minor error 2 months ago
- LICENSE Initial commit 7 months ago
- README.md Update README.md 2 months ago
- customPalette.json Adding option for a custom palette 3 months ago

`README.md`

CoCoView

<https://github.com/labbces/CoCoView/fork>

About
No description, website, or topics provided.
Readme
CC0-1.0 license
1 star
1 watching
0 forks

Releases
No releases published
Create a new release

Packages
No packages published
Publish your first

`labbces / CoCoView` Public

Code Issues Pull requests Actions Projects Wiki Security Insights Settings

main · 1 branch · 0 tags

Beatriz-Estevam Update README.md zfeesse on 29 Jul 47 commits

- examples Adding new example files and images 2 months ago
- images Adding new example files and images 2 months ago
- test add example1 and update principal script 5 months ago
- CoCoView.py correcting minor error 2 months ago
- LICENSE Initial commit 7 months ago
- README.md Update README.md 2 months ago
- customPalette.json Adding option for a custom palette 3 months ago

`README.md`

CoCoView

`git@github.com:labbces/CoCoView.git`

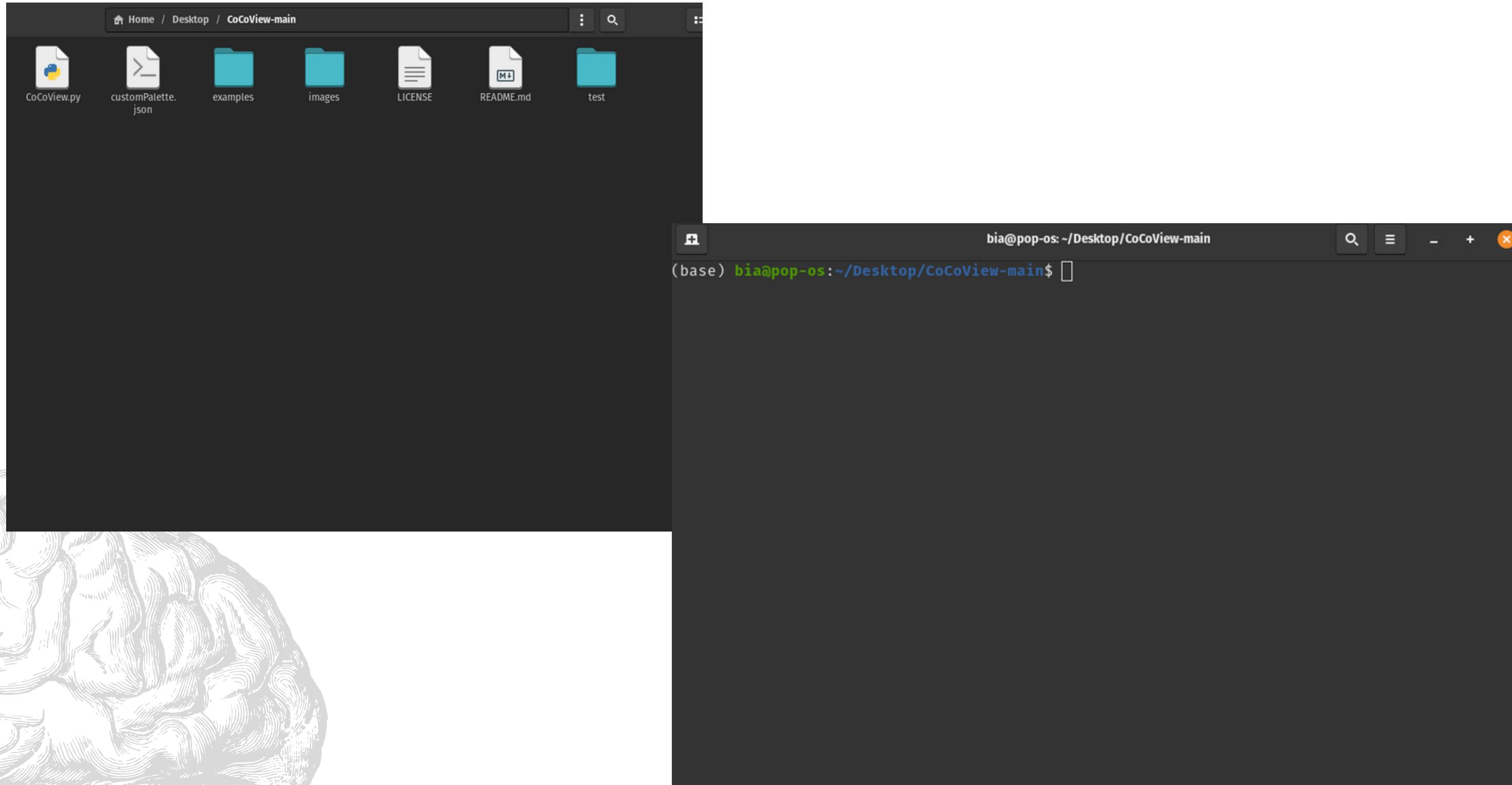
About
No description, website, or topics provided.
Readme
CC0-1.0 license
1 star
1 watching
0 forks

Releases
No releases published
Create a new release

Packages
No packages published
Publish your first package

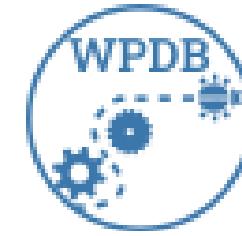
CoCoView: A Codon Conservation Viewer

Exemplo



CoCoView: A Codon Conservation Viewer

Exemplo

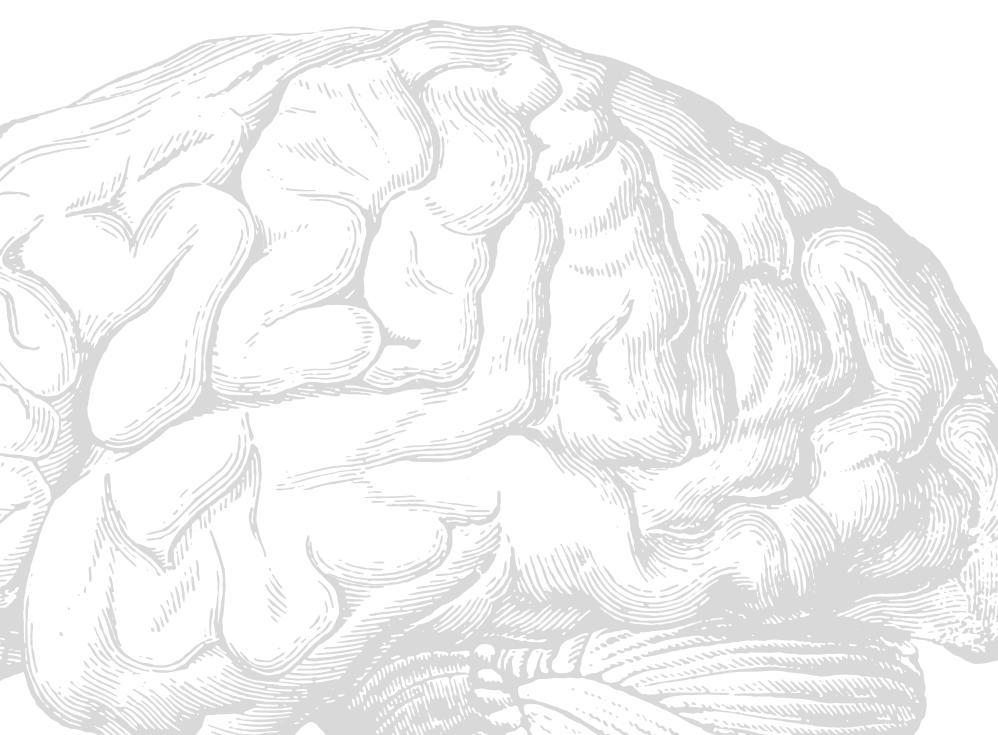


The image shows two screenshots of the CoCoView application interface. Both screenshots feature a dark-themed interface with a top navigation bar containing the CoCoView logo, a file icon, the file name "CoCoView_run.ipynb", a star icon, and menu items: Arquivo, Editar, Ver, Inserir, Ambiente de execução, Ferramentas, and Ajuda. Below the menu is a sidebar titled "Arquivos" with icons for search, upload, folder, and refresh. A red hand cursor is shown pointing at the "sample_data" folder. The main workspace is divided into two sections: "+ Código" and "+ Texto". A play button icon is located between these sections. The bottom portion of the interface shows a large, detailed anatomical illustration of a brain.

CoCoView.py
file.fasta

CoCoView: A Codon Conservation Viewer

Exemplo



```
+ Código + Texto
```

```
[ ] pip install logomaker
[ ] pip install pandas
[ ] pip install argparse
[ ] pip install matplotlib
[ ] pip install biopython
[ ] !python CoCoView.py AP2.Nta.veryShort.nt.aln.fa -p "Nicotiana_tabacumAP2" -i
```

The screenshot shows a terminal window with several pip installation commands for Python packages like logomaker, pandas, argparse, matplotlib, and biopython. Below these, a command is being run to execute the CoCoView script with specific parameters. The terminal interface includes tabs for 'Código' and 'Texto', and a status bar showing file paths and execution times (e.g., [11] 3s, [12] 5s). To the right of the terminal is a file browser window titled 'Arquivos' (Files) showing a directory structure with files related to the analysis, such as 'sample_data', 'AP2.Nta.veryShort.nt.aln.fa', 'CoCoView.py', and multiple versions of 'Nicotiana_tabacumAP2.redundant....'. A red circle highlights the terminal's command line area.

Validação do método

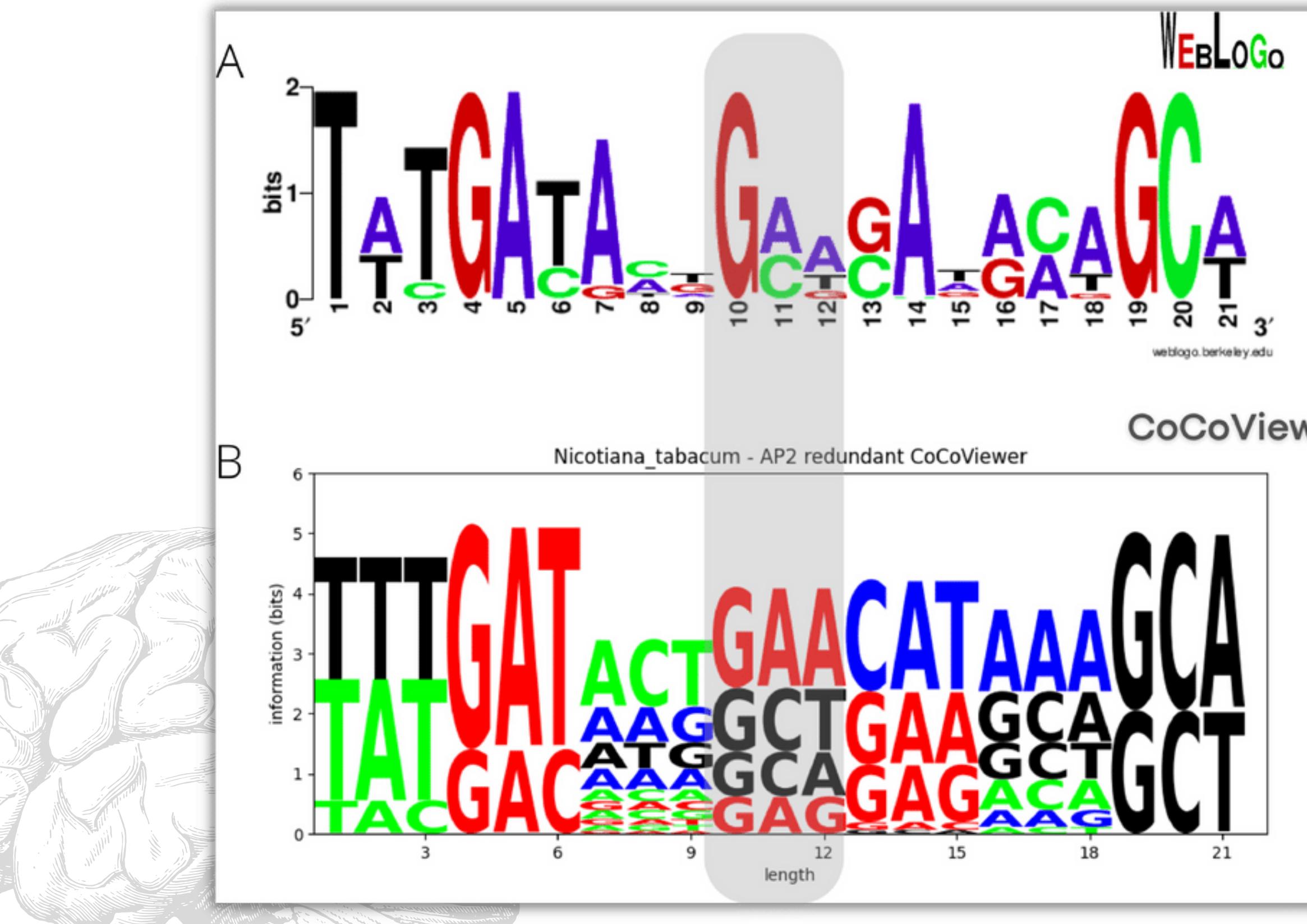
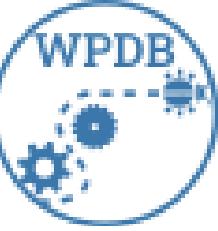


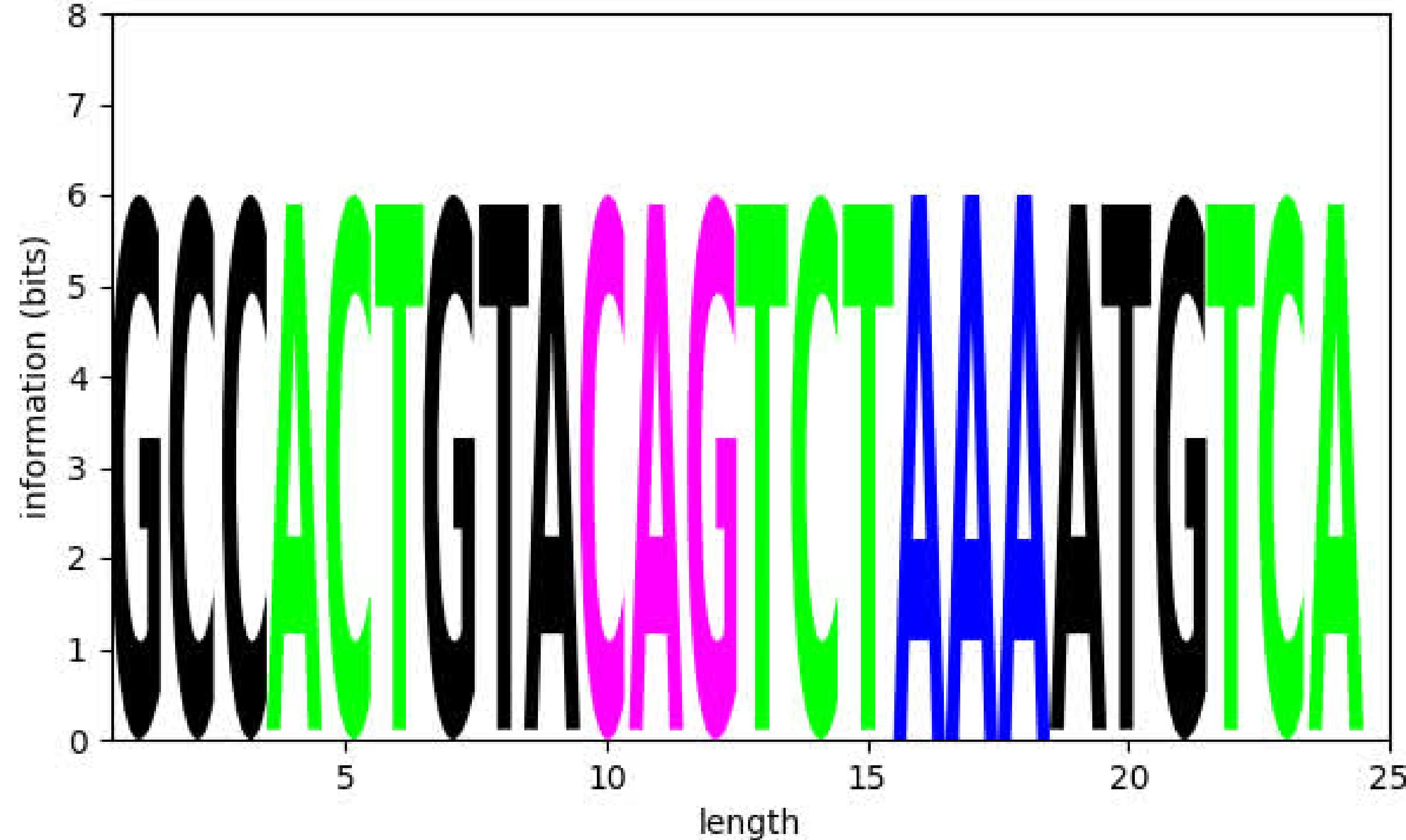
Fig. 1. CoCoView logo based on a multiple sequence alignment of a region of AP2 transcription factor coding sequences from *Nicotiana tabacum*. (A) Sequence logo generated using WebLogo [4], representing a per-nucleotide analysis. (B) Sequence logo generated using CoCoView (per-codon analysis). A per-nucleotide analysis could erroneously suggest that some codons are common, which can be ruled out on a per-codon visualization. Exemplified by the codon "GAT", at the position highlighted in gray on both sequence logos, which can be interpreted as a common codon in the per-nucleotide analysis. However, in the per-codon analysis, this codon does not occur at this position.

CoCoView: A Codon Conservation Viewer

Exemplo de aplicação



2020-01-10 - 2020-01-16 - nt NonRedundant (6)



Muito obrigada!

Agradecimentos

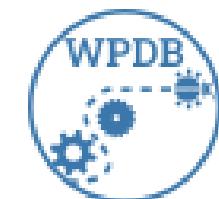
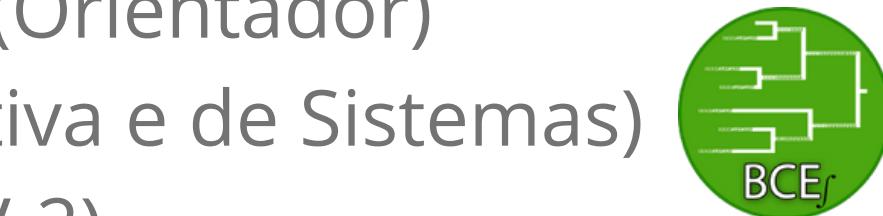
Prof. Dr. rer. nat. Diego Mauricio Riaño Pachón (Orientador)

LabBCES (Laboratório de Biologia Computacional, Evolutiva e de Sistemas)

Danilo Brito Rocha (Projeto do SARS-CoV-2)

USP/CENA

**Equipe Workshop Python para
Dados Biológicos**



Prática!

bia.estevam.25@usp.br 

@BiaREstevam 

<https://github.com/labbc ces/CoCoView>

Estevam BR, Riaño-Pachón, DM. **CoCoView - A codon conservation viewer via sequence logos.**

<https://doi.org/10.1016/j.mex.2022.101803>, 2022.

<https://labbc es.netlify.app/>

Muito obrigada!

Agradecimentos

Prof. Dr. rer. nat. Diego Mauricio Riaño Pachón (Orientador)
LabBCES (Laboratório de Biologia Computacional, Evolutiva e de Sistemas)
Danilo Brito Rocha (Projeto do SARS-CoV-2)



USP/CENA
**Equipe Workshop Python para
Dados Biológicos**



Perguntas?

bia.estevam.25@usp.br 
@BiaREstevam 

<https://github.com/labbces/CoCoView>

Estevam BR, Riaño-Pachón, DM. **CoCoView - A codon conservation viewer via sequence logos.**

<https://doi.org/10.1016/j.mex.2022.101803>, 2022.

<https://labbces.netlify.app/>