

# TRABAJO DE PRÁCTICAS

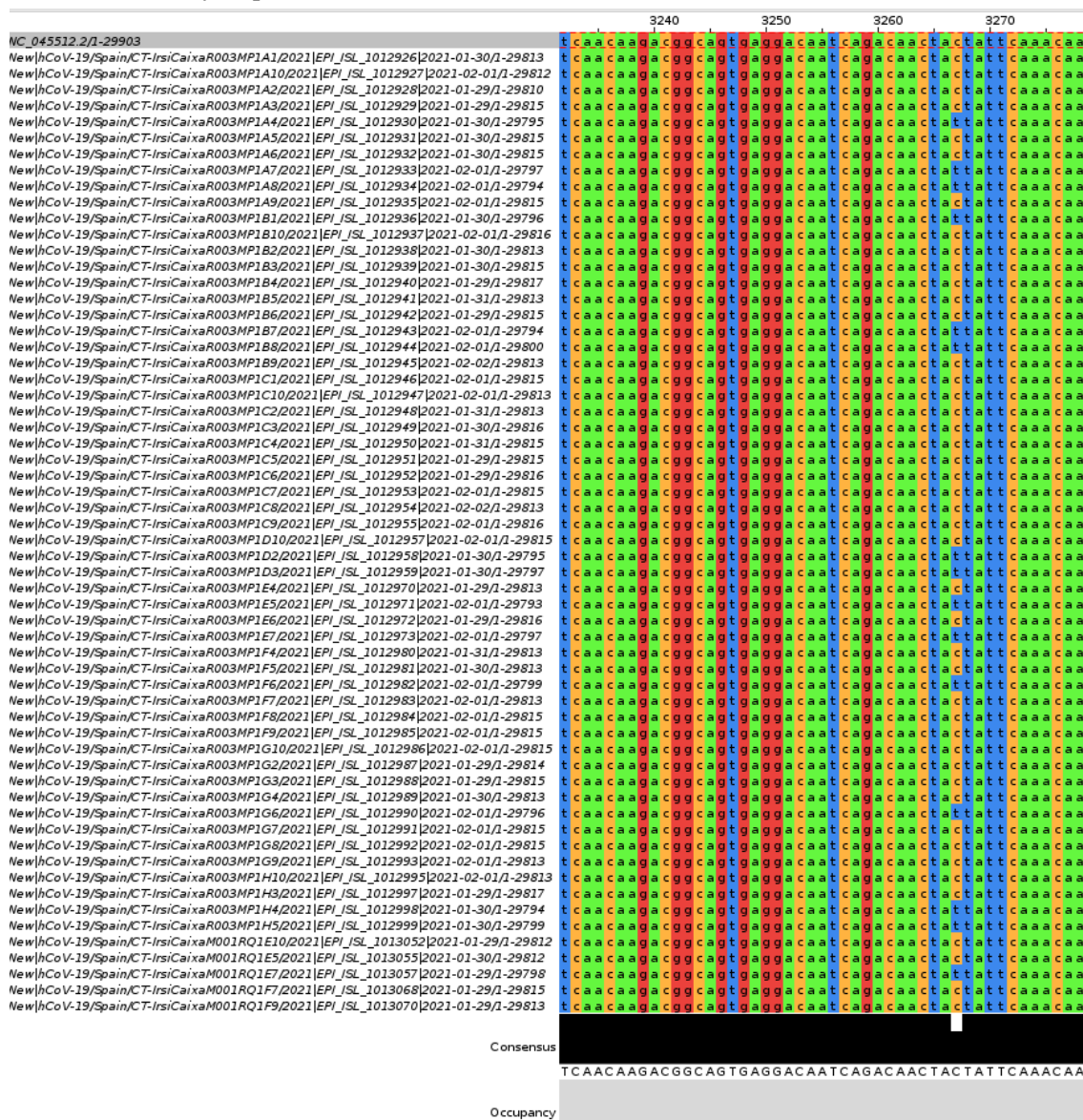
El objetivo de este trabajo de prácticas es estudiar con las herramientas tratadas en las anteriores prácticas la variedad y relación entre las secuencias proporcionadas (secuencias de nucleótidos de SARS-CoV-2, en concreto el grupo catHC4).

## Tarea 1: Multialineamiento

Utilizando [mafft](#) se va a crear un multialineamiento entre las secuencias del fichero 'gisaid\_hcov-19\_2021\_03\_23\_catHC4.fasta', la secuencia de la variante británica ('gisaid\_hcov-19\_2021\_03\_23\_B117.fasta') y la secuencia de referencia ('RefSeqWuhan.fasta'). Se ha realizado el alineamiento en la versión online de la herramienta ya que es mucho más rápida (22 segundos), el comando equivalente en la terminal es:

```
% mafft --reorder --maxambiguous 0.05 --addfragments fragments --auto input
```

El método utilizado ha sido FFT-NS-fragment que está por defecto y una estrategia elegida automáticamente que ha sido multipair ya que son 411 secuencias de 29903 bases, un tamaño razonable para mafft. El alineamiento final tiene una longitud de 29906 y está guardado en el fichero 'alineamiento.fasta' y es posible visualizarlo con Jalview.



## Tarea 2: Detección de mutaciones y errores de secuenciación

Para la detección de mutaciones se ha creado un programa en python que lee el alineamiento del fichero 'alineamiento.fasta' utilizando el módulo Biopython y posteriormente calcula la medida de conservación para cada posición en el alineamiento mediante la frecuencia de cada nucleótido.

$$f_v(i) = \frac{n_v(i)}{n(i)} \Rightarrow C(i) = - \sum_{v=1}^4 f_v(i) \ln(f_v(i))$$

Tras calcular la medida de conservación se busca secuencias que contengan mutaciones en posiciones cuya medida de conservación sea menor a 0.05 (posiciones con pocas variaciones). En concreto se van a buscar estas mutaciones en las siguientes posiciones: 266 a 13484 (Gen ORF1A), 21563 a 25384 (Proteína Spike), 25393 a 26220 (Gen ORF3A) y 28274 a 29533 (Gen N).

Para la ejecución del programa es necesario tener instalada una versión de python3 así como el módulo Biopython.

```
$ sudo apt-get update  
$ sudo apt-get install python3.6  
$ pip install biopython
```

Con el siguiente comando se pueden obtener tanto un fichero con las medidas de conservación 'conservaciones.txt' como las variaciones con más mutaciones respecto a la secuencia de referencia en 'mutacionesImportantes.txt':

```
$ python3 mutaciones.py
```

Los resultados de este programa muestran algunas secuencias con muchas mutaciones en estas posiciones (en las casillas se indica el número de mutaciones de cada grupo de cada secuencia):

Identificador de secuencia <sup>1</sup>	Posiciones con mutaciones			
	ORF1A	Spike	ORF3A	N
EPI_ISL_1020573	8	6	0	2
EPI_ISL_1012978	5	1	1	8
EPI_ISL_1013002	3	0	0	12
EPI_ISL_1208626	15	2	0	1
EPI_ISL_1208631	15	1	0	1
EPI_ISL_1208653	8	8	2	1

Estas mutaciones pueden afectar de forma muy negativa al individuo al estar en posiciones con un índice de conservación alto.

---

<sup>1</sup> El identificador de las secuencias se ha acortado por legibilidad, tiene el siguiente formato: Secuencia New[hCoV-19/Spain/CT-HUVH-09286/2021|EPI\_ISL\_xxxxxxx|2021-MM-DD

### Tarea 3: Construcción de una filogenia

Para la construcción de un árbol filogenético se han probado dos herramientas: RAxML y FastTree. Se ha clonado el repositorio de GitHub <https://github.com/stamatak/standard-RAxML> y se han seguido las instrucciones de instalación de la versión Pthreads SSE3 (fichero README). También se ha descargado el programa precompilado para Linux de la página inicial de FastTree en el apartado de Instalación(<http://www.microbesonline.org/fasttree/#Install>) .

Utilizando varios modelos distintos se ha buscado el árbol filogenético de máxima verosimilitud comparando la puntuación Log-likelihood en un modelo Gamma(20). Los resultados obtenidos son del mejor árbol para cada modelo y aparecen en la siguiente tabla. En tiempo de ejecución los resultados han sido muy similares.

	Modelo	Puntuación Log-likelihood
RAMxML	GTR+CAT	-54510.497986
	GTR+GAMMA	-54509.235990
	GTR+CATI	-54443.181115
	GTR+GAMMAI	-54443.669419
FastTree (Gamma20-based likelihood)	GTR+CAT	-56494.199000
	WAG+CAT	-57594.642000
	LG+CAT	-57594.642000
	JTT+CAT	-57594.642000

Los modelos que se han probado para la matriz de sustituciones y optimizaciones son:

- GTR(Generalized time-reversible)
- WAG(Whelan and Goldman)
- JTT(Jones-Taylor-Thornton)
- LG(Le and Gascuel)

Los modelos evolutivos que se han probado son:

- CAT: Más rápido al realizar aproximadamente  $\frac{1}{4}$  de las operaciones que el modelo GAMMA precisa. Solución menos precisa pero bastante aproximada.  
\*CATI: después de la búsqueda con CAT, se evalúan los árboles finales bajo una estimación de GAMMA en vez de la GAMMA predeterminada.
- GAMMA: Más lento, pero más preciso. Puede no funcionar para árboles muy grandes con más de 10.000 taxones.  
\*GAMMAI: igual que GAMMA pero con una estimación de la proporción de los sitios invariables.

El mejor árbol filogenético (GTR+CATI con RAxML) se encuentra en el fichero 'RAxML\_bestTree.gtrcati.txt'. El comando utilizado para obtener este árbol es:

```
$ ./raxmlHPC-PTHREADS-SSE3 -s alineamiento.fasta -n gtrcati.txt -m GTRCATI -p 111
```

El formato de salida de ambas herramientas son árboles Newick por lo que es posible visualizar dichos árboles en MEGA-X o cualquier otro visualizador. En la imagen posterior se muestra una parte del árbol donde se puede ver la secuencia de referencia marcada con un rombo rojo y la secuencia correspondiente a la variante británica en color verde.

