

Las secuencias de ADN con las que se trabajaron para el desarrollo de las soluciones secuencial y paralelas fueron E-Coli y la Salmonella. Que son básicamente archivos FASTA (extensión .fna) de 4,7 MB de tamaño en promedio.

B. Lenguaje y Librerías

Python es un lenguaje de programación poderoso y versátil, el módulo multiprocessing permite la programación paralela y concurrente en Python, y MPI4py es una biblioteca que implementa el estándar MPI para desarrollar aplicaciones paralelas y distribuidas en Python. Estas herramientas son fundamentales para aprovechar el poder de procesamiento de sistemas multiprocesador y clústeres de computadoras, mejorando el rendimiento y la eficiencia.

1) *Python*: Es un lenguaje de programación de alto nivel y fácil de aprender. Es ampliamente utilizado en diversos campos, incluyendo desarrollo web, análisis de datos, inteligencia artificial y computación científica. Python es conocido por su sintaxis legible y su amplia gama de bibliotecas y frameworks.

2) *Multiprocessing*: Es un módulo incorporado en Python que proporciona funcionalidades para la programación paralela y concurrente. Permite la ejecución de tareas en paralelo utilizando múltiples procesos en lugar de un solo proceso. Esto es especialmente útil para tareas intensivas en CPU que se pueden dividir en subprocesos independientes.

3) *MPI4py*: Es una biblioteca de Python que implementa el estándar Message Passing Interface (MPI). MPI es un modelo de programación utilizado para desarrollar aplicaciones paralelas y distribuidas en sistemas de memoria distribuida. Permite la comunicación y la coordinación entre múltiples procesos que se ejecutan en diferentes nodos de un clúster de computadoras.

Algunas de las otras librerías que utilizamos son las siguientes:

- NumPy: Biblioteca para computación numérica en Python, con estructuras eficientes y funciones matemáticas de alto nivel.
- Matplotlib.pyplot: Biblioteca de visualización en Python para crear gráficos y visualizaciones.
- Time: Módulo para trabajar con el tiempo, medir el tiempo de ejecución y manipular fechas y horas.
- Bio import SeqIO: Parte de Biopython, permite leer, escribir y manipular datos de secuencias biológicas.
- CV2: Biblioteca de visión por computadora (OpenCV) para procesamiento de imágenes y análisis de video.
- Argparse: Módulo para analizar argumentos de línea de comandos en programas Python.
- Tqdm import tqdm: Biblioteca para mostrar una barra de progreso en bucles iterativos en Python.
- Numba: Biblioteca para compilar funciones y bucles numéricos en Python para obtener un rendimiento más rápido.

C. Práctica

Para la ejecución de la aplicación se solicitaba que fuera a través de una línea de comandos, en el que se pudiera seleccionar la cantidad de procesadores que se soliciten.

Después de empezar a implementar las cadenas de ADN de prueba, nos enfrentamos al problema de que la memoria de las máquinas utilizadas para el procesamiento se llenaba y colapsaba por la cantidad de información procesada. Por lo

que también se hizo necesario poder seleccionar el porcentaje total de información que se va a procesar.

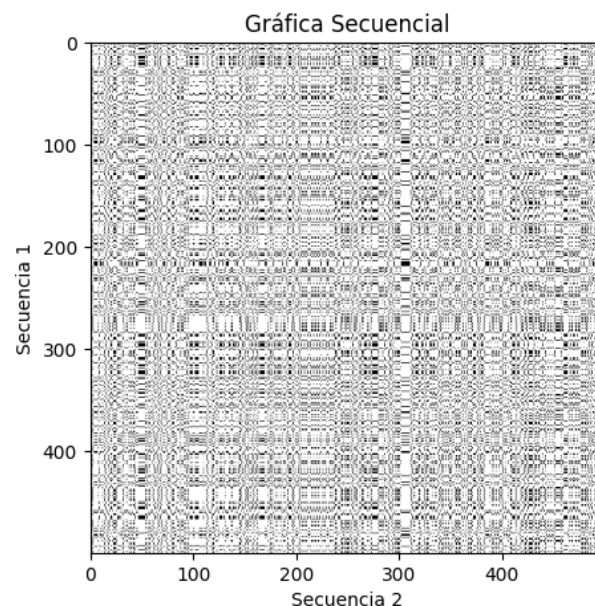
Una vez realizadas todas las configuraciones en el código necesarias para cumplir con estos requerimientos, la instrucción que se debe ejecutar luce así:

```
lpython appDotplot.py --file1=Salmonella.fna --file2=E_coli.fna --ps 0.2 --threads 2 4 8 --mpi
```

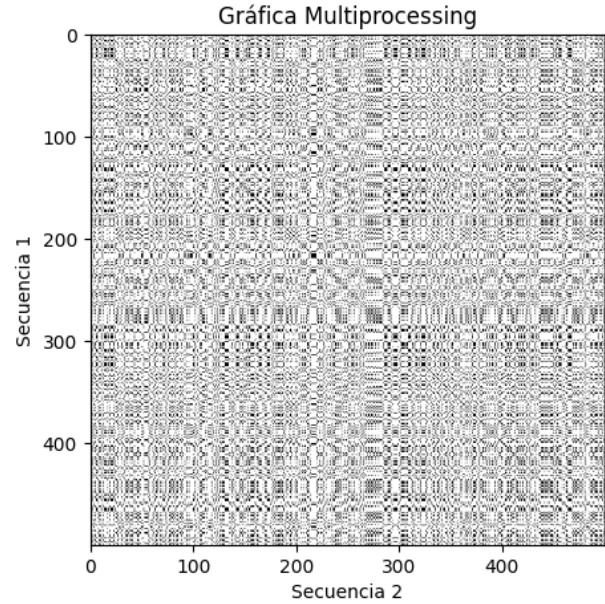
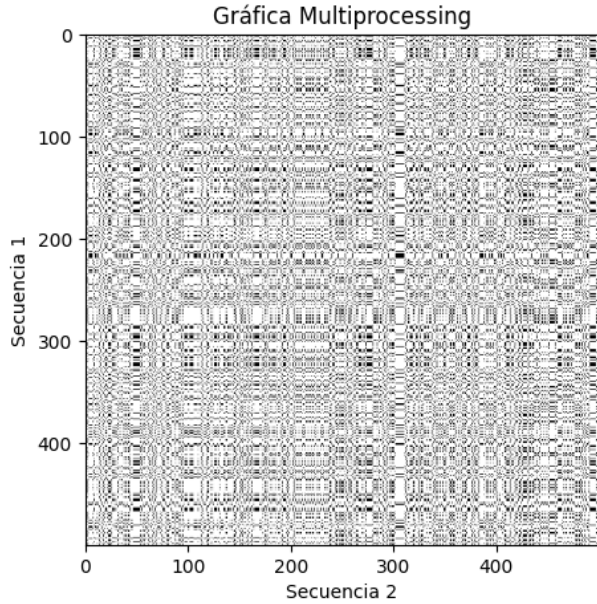
Por efectos prácticos, en nuestro caso decidimos usar el 20% de todos los datos, y hacer un zoom en la matriz resultante con un tamaño de 500x500, esto con el fin de agilizar los procesos y poder visualizar los cambios.

Encontramos entonces, que efectivamente las gráficas que generan cada una de las librerías no cambia o no presentan alteraciones como lo vamos a ver a continuación, sin embargo, la diferencia fundamental radica en la cantidad de tiempo que se toma cada proceso en realizar las operaciones.

1) *Secuencial*: En este caso, corresponde a realizar los procesamientos a pura fuerza bruta y sin paralelización.

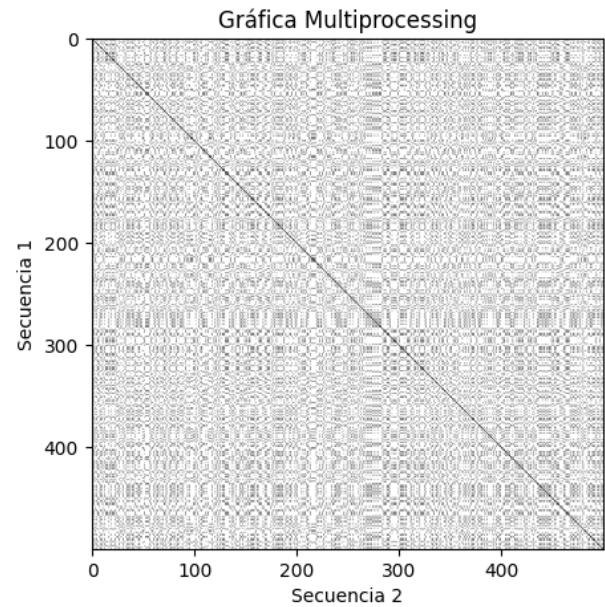
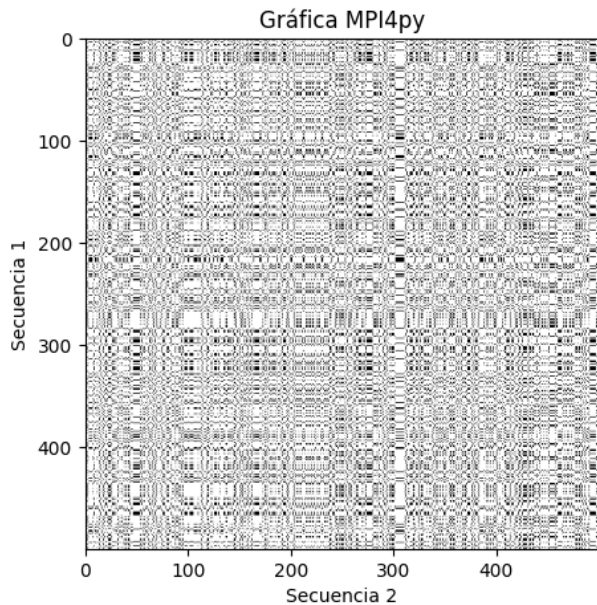


2) *Multiprocessing*: En el caso de multiprocessing, corresponde a una de las librerías comentadas anteriormente, que básicamente divide las tareas por la cantidad de procesadores ingresados.



Y este sería el resultado filtrado como lo especificado:

3) *MPI4py*: En el caso de *MPI4py*, esta librería permite la comunicación y la coordinación entre múltiples procesos que se ejecutan en diferentes nodos.



III. TIEMPOS Y ACELERACIÓN

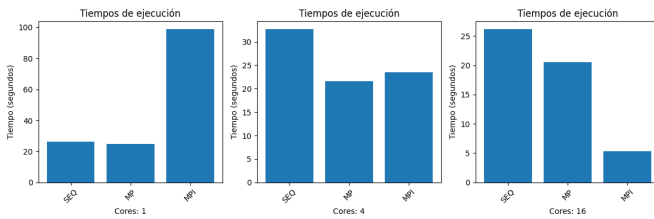
A continuación analizaremos los tiempos, la aceleración y eficiencia total de los tiempos de ejecución dependiendo la cantidad de procesadores ingresados.

4) *Filtro*: Dado el caso se encontrara dos secuencias parecidas en ciertos puntos, se requería implementar un filtro, que cambiara los valores en la matriz, si el valor en la matriz era 0 dejarlo en blanco, si había coincidencia y era diagonal se pintará de negro, si no que se pintará de gris.

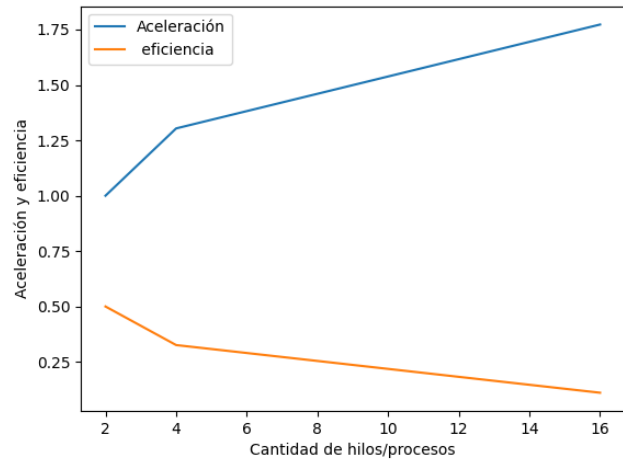
Este sería el resultado sin filtrar:

A. Tiempos

En esta sección revisaremos la gráfica de tiempos que está dividida por la cantidad de procesadores en las cuales está representado la librería utilizada y el tiempo total utilizado. Para la muestra hemos utilizado 2, 4 y 16 cores.



Como podemos observar el claro ganador a mayor número de procesadores es MPI, pero vemos que presenta una falencia grande al momento de presentar pocos procesadores.



B. Rendimiento

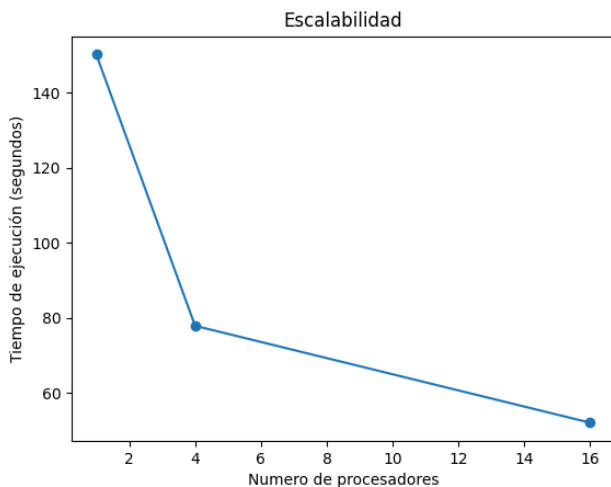
Python puede presentar limitaciones en aceleración y escalabilidad debido a su naturaleza de alto nivel. Multiprocessing puede proporcionar cierta aceleración al aprovechar múltiples núcleos, pero la escalabilidad puede verse limitada. MPI4py ofrece una mejor aceleración y escalabilidad al aprovechar sistemas de memoria distribuida y comunicación eficiente entre procesos.

IV. CONCLUSIÓN

En resumen, si lo que se busca es aceleración y escalabilidad y se tiene la capacidad para aumentar el número de núcleos, las bibliotecas especializadas como multiprocessing y MPI4py suelen ser más efectivas que el uso directo de Python. Multiprocessing es útil para aprovechar múltiples núcleos en un solo sistema, mientras que MPI4py es más adecuado para la programación paralela y distribuida en sistemas de memoria distribuida. La elección de la biblioteca depende de la naturaleza y los requisitos de tu proyecto, así como de la arquitectura del sistema en el que se ejecutará.

AGRADECIMIENTOS

Gracias al profesor Reinel Tabares Soto por la materia impartida. A Google por prestarnos sus equipos. Y a OpenAI por despejarnos dudas.



Como vemos al aumentar la cantidad de cores el tiempo tiende a disminuir, hasta llegar a un punto en el que el tiempo por más que se quisiera no podría disminuir.

La elección de la biblioteca de paralelización dependerá de los requisitos y características específicas del proyecto, así como de la arquitectura del sistema en el que se ejecutará y la cantidad de código que sea posible paralelizar.