



UNIVERSIDAD
DE ANTIOQUIA

HILOS Y PROCESOS PARA ACCELERAR PROCESOS DE BÚSQUEDA

Daniel Lujan Agudelo
Luis Guillermo Sánchez Cubides
Jose David Gómez Muñeton
Juan Pablo Arango Gaviria



HILOS



V/S



**PRO-
CESOS**



Introducción

El manejo de grandes volúmenes de datos es crucial para sectores como ciencia de datos, comercio electrónico e inteligencia artificial. Las búsquedas secuenciales se vuelven ineficaces a medida que los datos crecen, pero el paralelismo, usando hilos y procesos, permite mejorar drásticamente el rendimiento y la escalabilidad, ofreciendo soluciones rápidas y eficientes.

Hilos



- ✓ Espacio de memoria compartido entre hilos, lo que facilita la comunicación.
- ✓ El cambio de contexto es menos costoso.
- ✓ Menos costoso en terminos de memoria y CPU.
- ✓ I/O bound
- ✗ Es necesario manejar problemas de sincronización como condiciones de carrera.
- ✗ Los recursos compartidos pueden llevar a bloqueos entre hilos.

Procesos

- ✓ Espacio de memoria separado para cada proceso.
- ✓ Como tienen memoria independiente se eliminan los problemas de condiciones de carrera
- ✓ CPU bound
- ✗ Mayor uso de la CPU.
- ✗ Como cada proceso tiene su propio espacio de memoria incrementa el consumo de la misma.



Metodología

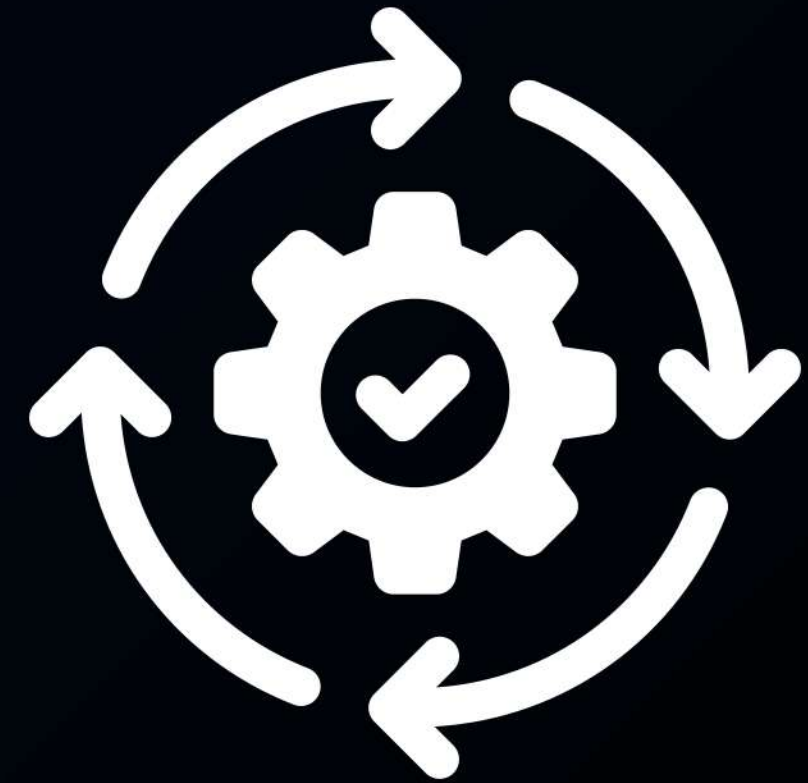
Implementación de las técnicas de
búsqueda



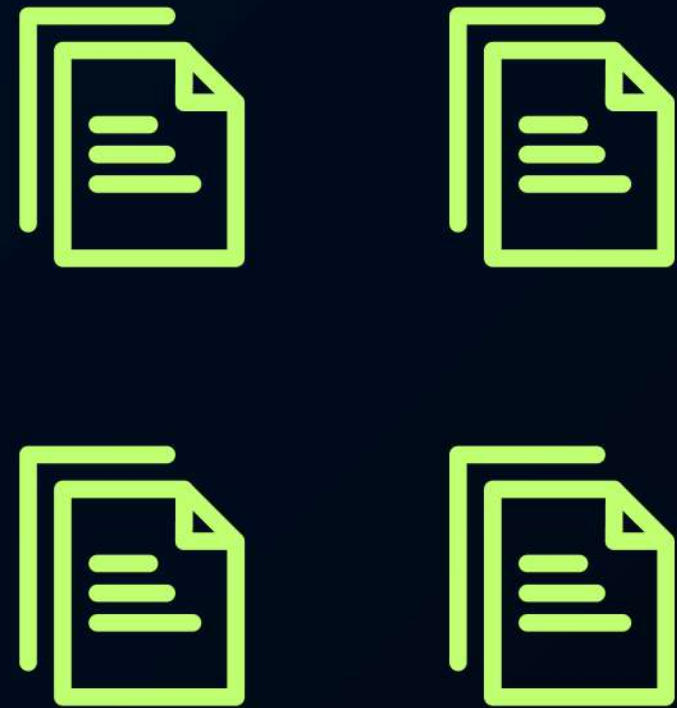
Configuración del entorno
experimental



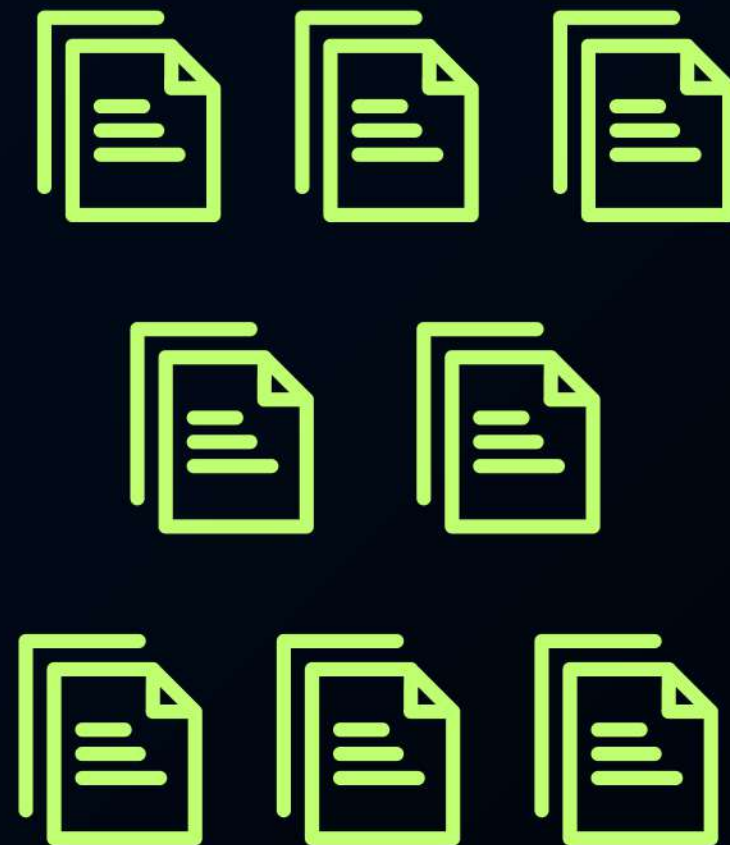
Ejecución de pruebas en datasets de
diferentes tamaños



El Enfrentamiento...



4 Archivos
497,664 Registros C/U



8 Archivos
248,832 Registros C/U



12 Archivos
165,888 Registros C/U

Multithreading

Se implementa **multithreading** para realizar la búsqueda de un valor en múltiples archivos. Para cada archivo se crea un hilo independiente y este recorre línea por línea para comparar el contenido con el valor buscado.

Si un hilo encuentra el valor, todos los demás detienen su ejecución, evitando trabajo innecesario y mejorando el tiempo total de la búsqueda.



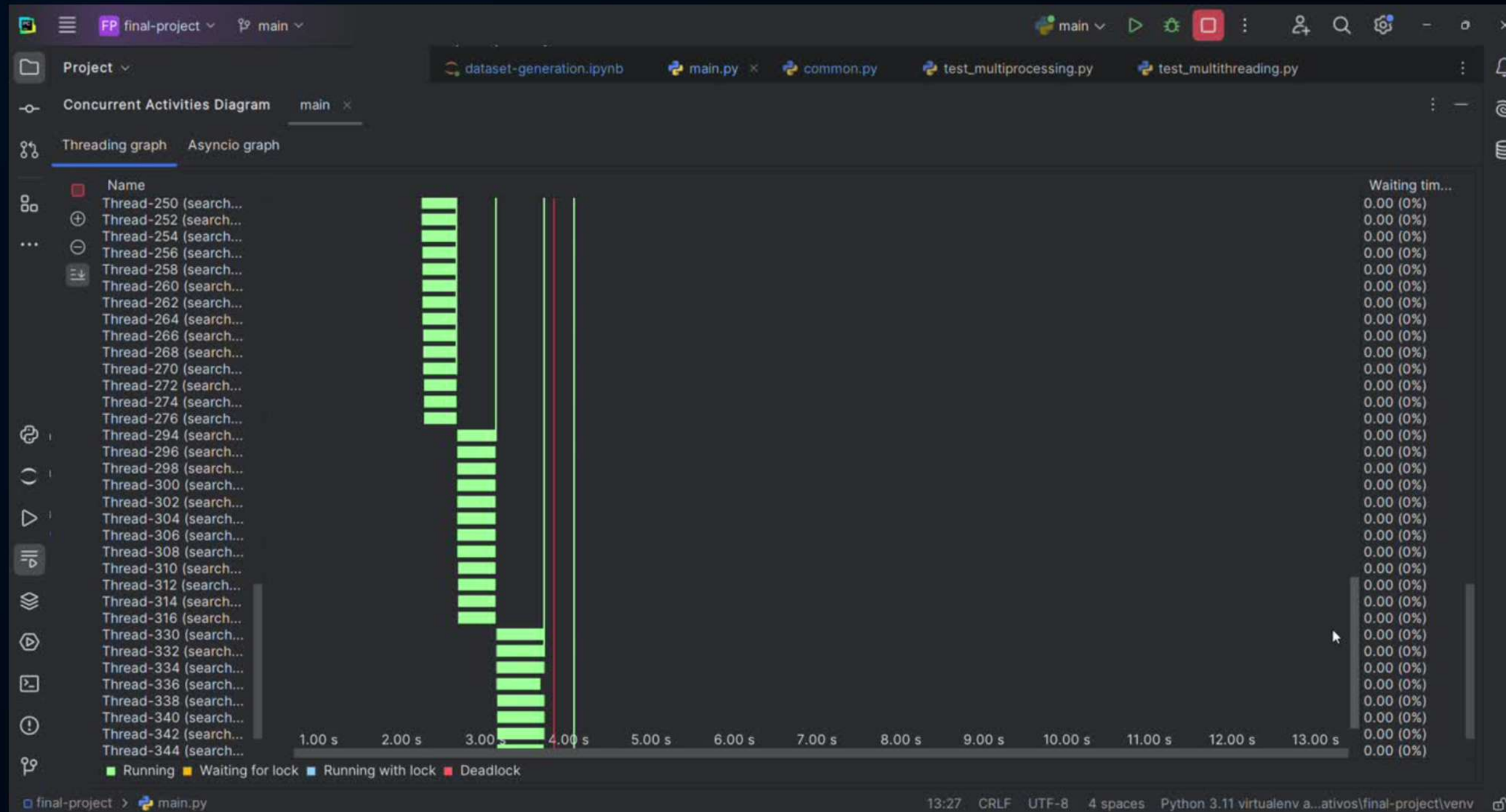
Multiprocessing

En este metodo en lugar de usar hilos, se crean proceos independientes, donde cada proceso es el encargado de buscar en un archivo en especifico. Cada proceso se ejecuta de forma aislada, trabajando en su propio espacio de memoria y evitando cualquier interferencia entre los mismos.

La implementación utiliza un pool de procesos y a través del método map, el pool asigna a cada proceso un archivo y un valor a buscar, ejecutando la búsqueda en paralelo y retornando una lista de resultados.



Demo



Demo

The image shows a Windows environment with a terminal window and the Task Manager application.

Terminal Window:

```
→ final-project ⚡ (P main)
▶ python main.py
Target values values generated and saved in ./temp/target_values/
Ejecutando test con Multithreading...
Test 1/5 - COMPLETED
Test 2/5 - COMPLETED
Test 3/5 - COMPLETED
Test 4/5 - COMPLETED
```

Task Manager - Procesos:

Nombre	Estado	CPU	Memoria
Pantalla de bloqueo predeterminada de Windows (2)		0%	
pg_ctl - starts/stops/restarts the PostgreSQL server		0%	
pgAgent - PostgreSQL Scheduling Agent		0%	
PostgreSQL Server (8)		0%	
Print driver host for applications		0%	
Procesador de comandos de Windows		0%	
Proceso de host para tareas de Windows		0%	
Proceso de host para tareas de Windows		0%	
Python (5)		0%	
Python (2)		10.0%	
Realtek HD Audio Universal Service		0%	
scdaemon.exe		0%	
Servicio Administración de máquinas virtuales		0%	
Servicio de proceso de host de Hyper-V		0%	
ShellHost		0%	
SpotifyWidgetProvider.exe		0%	
SQL Server VSS Writer - 64 Bit		0%	
updatechecker.exe		0%	
User OOBE Broker		0%	
Usermode Font Driver Host		0%	
Usermode Font Driver Host		0%	

Uso de recursos

Starting Multithreading

Ejecutando test con Multithreading...

Test 1/3 - COMPLETED

Test 2/3 - COMPLETED

Test 3/3 - COMPLETED

Test finalizado. Resultados guardados en 'results-multithreading.csv'

Multithreading completed!

Peak CPU usage: 62.50%

Peak Memory usage: 66.10 MB

Elapsed time: 4.51 seconds

Starting Multiprocessing

Ejecutando test con Multiprocessing...

Test 1/3 - COMPLETED

Test 2/3 - COMPLETED

Test 3/3 - COMPLETED

Test finalizado. Resultados guardados en 'results-multiprocessing.csv'

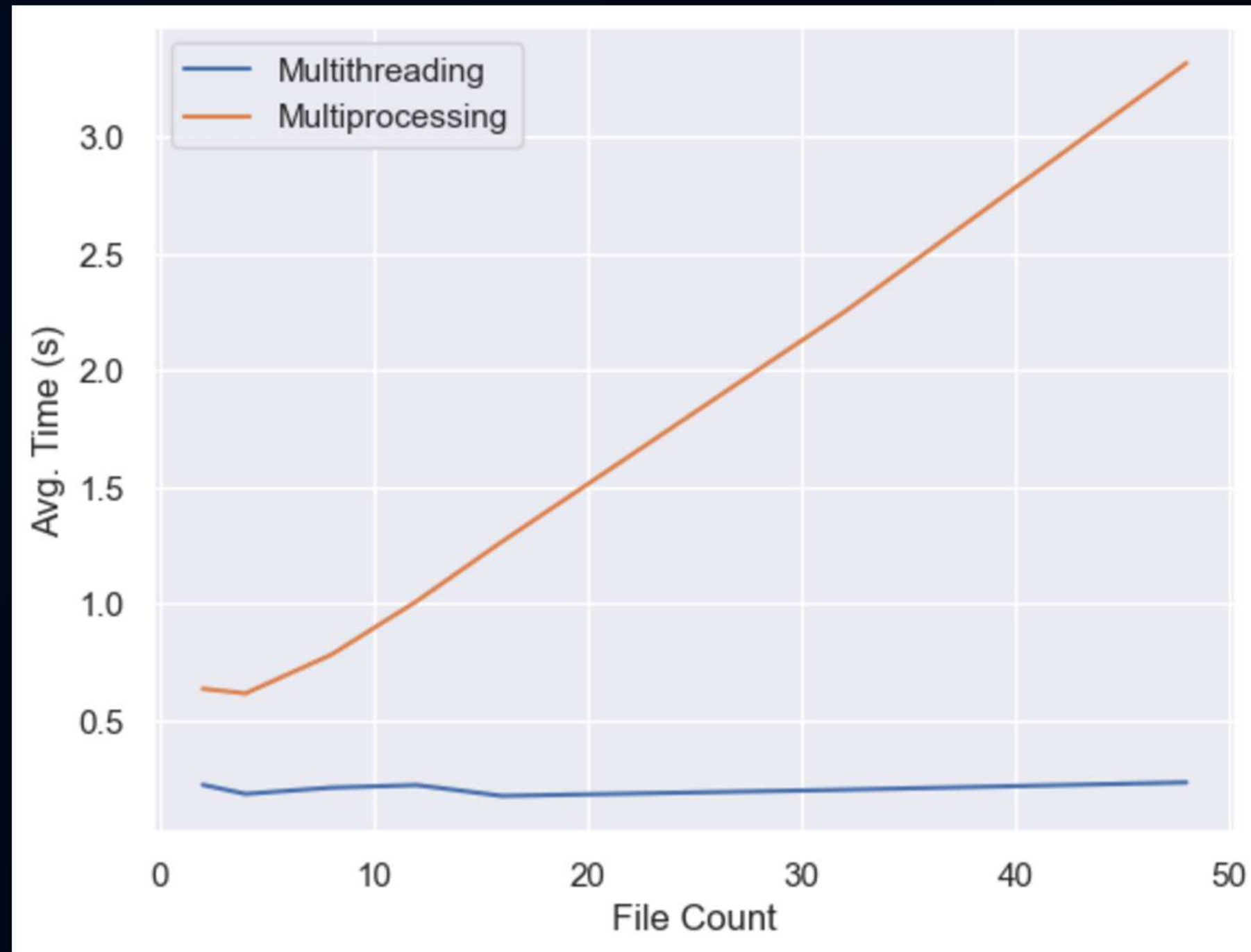
Multiprocessing completed!

Peak CPU usage: 71.90%

Peak Memory usage: 65.00 MB

Elapsed time: 30.54 seconds

Análisis Estadístico



Análisis Estadístico

H0: No hay diferencias significativas en los tiempos de ejecución entre las técnicas de búsqueda.

H1: Existen diferencias significativas en los tiempos de ejecución entre las técnicas.

Análisis Estadístico

	SCE	DF	F	PR(>F)
Técnica de paralelismo	202.358518	1	447.732571	2.096616e-73
Residual	252.195307	558	----	----

Análisis Estadístico

Se rechaza la hipótesis nula:

H1: Existen diferencias significativas en los tiempos de ejecución entre las técnicas.

... Y el ganador es 🥁



Conclusiones

- El rendimiento de la técnica de procesamiento paralelo mediante multithreading se ve altamente afectado con el número de archivos, por la cantidad de cambios de contexto.
- Es mejor utilizar procesamiento paralelo mediante multiprocessing cuando las tareas que se requieren son de alto consumo de CPU.
- El multiprocessing es menos adecuado para este tipo de tareas debido a la sobrecarga asociada con la creación y manejo de procesos en comparación con los hilos.



Referencias

1 What is multiprocessing?

<https://www.techtarget.com/searchdatacenter/definition/multiprocessing>

2 Entendiendo los procesos, hilos y multihilos.

<https://medium.com/@diego.coder/entendiendo-los-procesos-hilos-y-multihilos-9423f6e40ca7>

3 Python time.time() vs time.perf_counter().

https://superfastpython.com/time-time-vs-time-perf_counter/



Gracias



Repositorio
<https://github.com/daniel-lujan/multiprocessing-multithreading>



Notebook
<https://github.com/daniel-lujan/multiprocessing-multithreading/blob/main/notebooks/EDA.ipynb>