



Electrical & Computer
Engineering Department

**UNIVERSITY OF
PELOPONNESE**

Parallel Processing in Python 3.13

Parallel Systems and Programming ECE_INF920

Students

Mantvydas Deltuva ece24601

Justinas Teselis ece24600

Professor

Παναγιώτης Αλεφραγκής

Patras, 2025

Contents

1. Introduction	4
1.1. Definition of Parallel Processing	4
1.2. Importance in Modern Computing	5
1.3. Purpose of the Document	5
2. Background	6
2.1. The Global Interpreter Lock (GIL)	6
2.2. Threads vs. Processes	7
2.3. Evolution of Parallel Processing in Python	9
3. Parallel Processing Basics in Python	10
3.1. Module threading	10
3.2. Module multiprocessing	12
3.3. Module concurrent.futures	13
4. Free-threaded mode in Python 3.13	14
5. Comparing GIL and GIL-free Performance	16
5.1. Count	16
5.2. Download Images	20
5.3. Merge Sort	24
5.4. Matrix Multiplication	28
6. Conclusion	32

Illustrations

1 pic. Comparison of Sequential and Parallel Computing Approaches	4
2 pic. Threads and Processes Feature Comparisons	8
3 pic. Module threading Example Code Snippet (part 1)	11
4 pic. Module threading Example Code Snippet (part 2)	11
5 pic. Module multiprocessing Example Code Snippet	12
6 pic. Module concurrent.futures Example Code Snippet	13
7 pic. Count Single-Threaded Code Snippet	16
8 pic. Count Multi-Threaded Code Snippet	17
9 pic. Count Multi-Process Code Snippet	17
10 pic. Count Performance Graphs	19
11 pic. Download Images Task Code Snippet	20
12 pic. Download Images Single-Threaded Code Snippet	20
13 pic. Download Images Multi-Threaded Code Snippet	21
14 pic. Download Images Multi-Process Code Snippet	21
15 pic. Download Images Performance Graphs	23
16 pic. Merge Sort Task Code Snippet	24
17 pic. Merge Sort Single-Threaded Code Snippet	24
18 pic. Merge Sort Multi-Threaded Code Snippet	25
19 pic. Merge Sort Multi-Process Code Snippet	26
20 pic. Merge Sort Performance Graphs	27
21 pic. Matrix Multiplication Task Code Snippet	28
22 pic. Matrix Multiplication Single-Threaded Code Snippet	28
23 pic. Matrix Multiplication Multi-Threaded Code Snippet	29
24 pic. Matrix Multiplication Multi-Process Code Snippet	30
25 pic. Matrix Multiplication Performance Graphs	31

1. Introduction

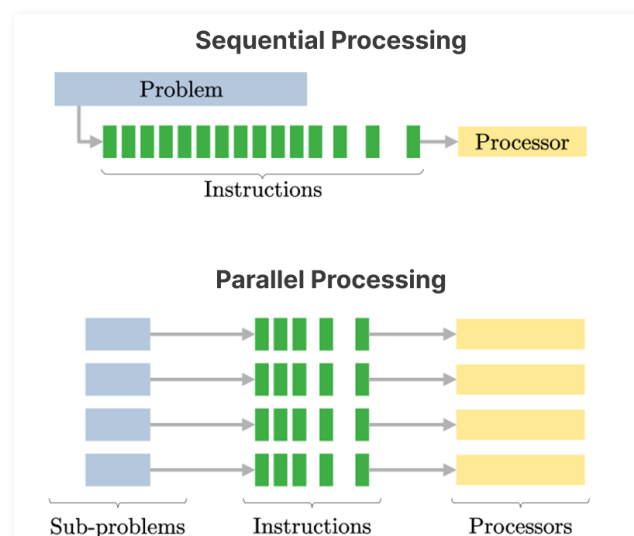
In today's era where computational demands continue to rise, parallel processing has become a cornerstone of modern computing. By leveraging multi-core processors and distributed systems, parallel processing allows for the efficient execution of tasks, enabling faster computations and improved scalability. This document explores the principles and applications of parallel processing, particularly in Python, and highlights the advancements introduced in Python 3.13.

1.1. Definition of Parallel Processing

Parallel processing provides the ability to execute multiple (usually repetitive) tasks simultaneously by distributing them across multiple threads, CPU cores or systems. It leverages today's standard multi-core processor architectures by dividing large computational workloads into smaller tasks. By doing so, it reduces overall execution time and enhances the performance of computations, particularly those involving intensive calculations or data processing.

Key aspects of parallel processing include:

- **Concurrency** - the ability to manage multiple tasks at the same time.
- **Scalability** - computing resources are efficiently utilized by managing increasing workloads.
- **Synchronization** - ensuring tasks interact seamlessly without conflicts, memory corruption or errors.



1 pic. Comparison of Sequential and Parallel Computing Approaches

1.2. Importance in Modern Computing

Parallel processing is foundational to the efficiency and scalability of modern computing systems. Its presence in various fields such as machine learning, scientific simulations, data analysis and web services is crucial for growing demands of computational resources. The primary benefits include:

- **Enhanced scalability** - modern systems often manage massive datasets and complex computations. Parallel processing allows applications to scale efficiently, accommodating these demands without bottlenecks.
- **Reduced computation time** - tasks that would take tremendous amounts of time to complete can now be completed in significantly less time, enabling quicker insights and faster decision-making.
- **Improved computational efficiency** - by dividing tasks among multiple processors or cores, parallel processing ensures optimal utilization of hardware resources, minimizing idle time and maximizing throughput.

Examples of areas that benefit from parallel processing include:

- **Real-time applications** - enhancing the responsiveness of systems such as autonomous vehicles and online gaming.
- **Scientific research** - simulating physical systems or modeling complex phenomena such as galaxies or other bodies.

1.3. Purpose of the Document

This document serves as a guide for understanding and implementing Python's parallel processing capabilities, with a focus on newly added features which were introduced in Python 3.13. Specifically, it aims to:

- **Highlight Python's progress in parallelism** – showcase Python's tools and libraries for parallel computing, including the introduction of the GIL-free mode in Python 3.13.
- **Demonstrate real-world applications** – define practical use cases where parallel processing can drive significant performance improvements and differences in implementations.

With the **experimental removal** of the **Global Interpreter Lock (GIL)** in **Python 3.13**, this guide is mostly focused on this feature for multi-threaded programming and the opportunities it unlocks for Python developers in achieving true parallelism.

2. Background

The efficiency of Python's concurrency model is a widely discussed topic, with much of the conversation revolving around its **Global Interpreter Lock (GIL)**. Python's threading model, often misunderstood due to the constraints of the GIL, has both strengths and weaknesses. Understanding the differences between threads and processes in Python is key to designing efficient, scalable applications, especially in multi-core or distributed computing environments.

2.1. The Global Interpreter Lock (GIL)

The most widely used Python implementation uses **Global Interpreter Lock (GIL)** which is a crucial feature that has its own advantages and disadvantages. It is a mutex (mutual exclusion lock) that enforces rules for native threads for executing Python code simultaneously. It disallows concurrent execution. This way it simplifies memory management, garbage collection and ensures thread safety. But the main drawback is that it significantly impacts the performance of multi-threaded programs.

The GIL was introduced in the early days of Python to simplify the implementation of it by:

- **Reducing complexity** - the lock eliminates the need for complex synchronization mechanisms across threads and thus makes it easier for development. Internal state remains consistent by preventing multiple threads from executing at the same time.
- **Managing memory** - the GIL ensures that memory allocation reference counts are updated safely (one by other) without race conditions.
- **Supporting extensions** - many extensions for Python depend on the GIL to avoid implementing their own thread-safety mechanisms. For extensions developers this reduces the complexity of thread-safe code.

Despite its advantages, the GIL provides certain drawbacks:

- **Latency in I/O-bound scenarios** - context switching among threads introduces latency, which can degrade performance for applications with many concurrent and repeating I/O tasks.
- **Underutilization of multi-core systems** – today's hardware with multiple cores is not fully utilized due to the single-threaded code execution bottleneck by the GIL.

- **"True" multi-threading** – as GIL only allows one thread to execute Python code at a time, the performance of multi-threaded applications designed for CPU-bound tasks has internal limitations by the language.
- **Workarounds** - developers often use multiprocessing or rewriting critical sections in other languages to overcome GIL-related performance issues. Development overhead and complexity increases by the size of project.

2.2. Threads vs. Processes

Threads are lightweight units of execution which live in a single specific process. This way it is using already allocated memory and resources of parent process. However, they also share the constraints of the Global Interpreter Lock in Python, which limits their effectiveness for certain tasks.

Key features of threads:

- **Low overhead** – threads create minimal memory overhead compared to processes because they operate within a single, already initiated with resources and memory, process.
- **Shared memory space** – does not require inter-process communication (IPC) because it operates within the same process. This way it can access and modify the same data structure without IPC.
- **I/O-bound tasks** - threads excel at tasks, such as reading and writing to files, network communication, or database interactions.

Limitations of threads:

- **Thread safety** – synchronization locks or semaphores are essential to ensure that race condition does not occur in shared memory space.
- **GIL dependency** – threads are ineffective for CPU-bound tasks, as GIL prevents them from executing Python code concurrently. Only one thread can do that at a time.

Processes are independent execution units, each with its own memory space and system resources. They bypass the GIL entirely, making them suitable for tasks requiring true concurrent calculation.

Key features of processes:

- **True parallelism** – since each process operates independently, they can fully utilize multi-core CPUs.

- **Independent memory space** – race condition is not present because processes do not share memory by default. This isolation enhances stability but also requires IPC mechanisms for communication between processes.
- **CPU-bound tasks** – processes are ideal for computationally intensive tasks, such as data processing, machine learning or numerical simulations.

Limitations of processes:

- **Overhead** – each process has its own CPU resources and memory allocation thus creating overhead for managing processes.
- **Complex communication** - sharing data between processes requires explicit mechanisms like queues, pipes, or shared memory.

<i>Feature</i>	<i>Threads</i>	<i>Processes</i>
Memory Sharing	Shared memory space	Independent memory space
Parallelism	Limited by the GIL in Python	True parallelism across cores
Overhead	Minimal	High
Best Use Case	I/O-bound tasks	CPU-bound tasks
Communication	Easier (shared memory)	Requires IPC mechanisms
Safety	Risk of race conditions	No shared memory race conditions

2 pic. Threads and Processes Feature Comparisons

2.3. Evolution of Parallel Processing in Python

Now, Python has multiple modules that leap towards efficient parallel processing, each with its own advantages and disadvantages. These tools have laid the groundwork for parallel computing in Python.

threading:

- Introduced lightweight threading for I/O-bound tasks.
- Threads share memory within a single process, enabling fast communication but making them dependent on the GIL.

multiprocessing:

- Enabled true parallelism by creating separate processes instead of threads.
- Each process has its own memory space, allowing them to operate independently and bypass the GIL.
- Suited for CPU-bound tasks but has higher overhead due to the need for inter-process communication (IPC) and separate memory allocation.

concurrent.futures:

- Provided a high-level interface for managing pools of threads and processes. Simplified the implementation of concurrent tasks.
- Designed for developers seeking to implement parallelism with minimal boilerplate, while offering flexibility in choosing between threads, processes, or both.

3. Parallel Processing Basics in Python

Python has several modules to implement concurrency and parallelism, each suited for specific use cases. Whether dealing with I/O-bound or CPU-bound tasks, these modules provide the tools to manage threads, processes, and tasks effectively. This section introduces the **threading**, **multiprocessing**, and **concurrent.futures** modules, explaining their key features, limitations, and ideal use cases, along with practical examples.

3.1. Module threading

Python threads are most effective for I/O-bound tasks therefore threading module can achieve concurrency in Python for tasks such as reading files, handling network requests or other tasks, where program must wait for external resources. Memory space is shared across threads, making data sharing simple but prone to race conditions. To ensure thread-safe execution, module provides synchronization primitives such as locks, semaphores. This module is constrained by the GIL for CPU-bound tasks which limits true parallelism and its effectiveness. For more information regarding this module, please refer to the official documentation on threading [1].

Below is an example of how to download and save multiple images concurrently using threading module. Code explanation:

- Threads are created to manage multiple file downloads simultaneously with universal image download function.
- `thread.start()` begins the execution of each thread.
- `thread.join()` ensures the main program waits for all threads to finish before proceeding.

Parallel Processing in Python 3.13

```
1 import threading
2 import requests
3 import time
4 import os
5
6 # Create a start event
7 start_event = threading.Event()
8
9 # Universal function for downloading an image in .jpg format from specified URL
10 def download_image(id, url, output_folder):
11     try:
12         # Construct image path based on provided output folder and ID
13         image_path = os.path.join(output_folder, f"image_{id}.jpg")
14
15         # Provide initialization feedback
16         print(f"Thread for image [{id}] is initialized, waiting for start signal...")
17
18         # Wait for start signal
19         start_event.wait()
20
21         # Provide start feedback
22         print(f"Starting download for image [{id}] from [{url}]")
23
24         # Send a GET request to the URL and stream the response
25         response = requests.get(url, stream = True)
26         response.raise_for_status()
27
28         # Download and save the image in 8KB chunks
29         with open(image_path, "wb") as image:
30             for chunk in response.iter_content(chunk_size = 8192):
31                 image.write(chunk)
32
33         # Provide completion feedback
34         print(f"Completed download for image [{id}], saved to [{image_path}]")
35     except Exception as e:
36         print(f"Failed to download image [{id}]: {e}")
```

3 pic. Module threading Example Code Snippet (part 1)

```
56 # Create threads for downloading images
57 threads = [
58     threading.Thread(
59         target = download_image,
60         args = (i, url, output)
61     )
62     for i, url in enumerate(urls)
63 ]
64
65 # Initialize all threads
66 for thread in threads:
67     thread.start()
68
69 # Signal all threads to start
70 time.sleep(3)
71 print("Releasing threads to start downloads...")
72 start_event.set()
73
74 # Wait for all threads to complete
75 for thread in threads:
76     thread.join()
77
78 # Provide finish feedback
79 print("All downloads completed.")
```

4 pic. Module threading Example Code Snippet (part 2)

3.2. Module multiprocessing

Python processes are most effective for CPU-bound tasks therefore multiprocessing module can achieve parallelism in Python for tasks such as numerical computations, data processing or other tasks, where program can benefit from multi-core CPU architectures by creating separate processes, each with its own memory space and CPU resources. To have clear communication with processes there are supported tools for it like pools, queues, or pipes. This module overcomes the GIL limitation. For more information regarding this module, please refer to the official documentation on multiprocessing [2].

Below is an example of how to compute number squares in parallel using multiprocessing module. Code explanation:

- A Pool of workers is created to divide the computation of squares across multiple processes.
- Function `pool.map()` distributes tasks to the worker processes and collects the results.
- Each process runs independently, leveraging multiple CPU cores for true parallel execution.

```

1  from multiprocessing import Pool, current_process
2  import random
3  import time
4
5  # Universal function to compute the square of a number
6  def compute_square(number):
7      # Provide start feedback
8      process_name = current_process().name
9      sleep_time = random.uniform(0.5, 2.0)
10     print(f"Process {process_name} is calculating square " +
11           f"for {number} with sleep time {sleep_time}")
12
13     # Simulate a heavy computation with a random sleep time
14     time.sleep(sleep_time)
15
16     # Return the square of the number
17     return number * number
18
19 if __name__ == "__main__":
20     # Numbers to calculate the square of
21     numbers = [1, 512, 4, 128, 16, 32, 64, 8, 256, 2]
22
23     # Create a pool of processes to calculate the squares
24     with Pool(processes = 4) as pool:
25         results = pool.map(compute_square, numbers)
26
27     # Provide finish feedback
28     print(f"Squares: {results}")

```

5 pic. Module multiprocessing Example Code Snippet

3.3. Module concurrent.futures

For unified usage of threads and processes there is concurrent.futures module. It simplifies the implementation of parallelism regarding management of threads, processes, and tasks. It also includes pool executors with helpful utility functions such as submit() or map(). Threads and processes have the same characteristics as mentioned before, where threads are dependent on GIL, and processes have overhead regarding memory. For more information regarding this module, please refer to the official documentation on concurrent.futures [3].

Below is an example of how to compute number squares in parallel using multiprocessing module. Code explanation:

- ProcessPoolExecutor is used to manage a pool of processes for computing data concurrently.
- executor.map() applies the provided function to each number in the list. This way the load is distributed evenly.
- The module simplifies thread management, reducing boilerplate code and providing more control over pools.

```

1  from concurrent.futures import ProcessPoolExecutor
2  from multiprocessing import current_process
3  import random
4  import time
5  import math
6
7  # Universal function to compute factorial of a number
8  def compute_factorial(number):
9      # Provide start feedback
10     process_name = current_process().name
11     sleep_time = random.uniform(0.5, 2.0)
12     print(f"Process {process_name} is calculating factorial " +
13           f"for {number} with sleep time {sleep_time}")
14
15     # Simulate a heavy computation with a random sleep time
16     time.sleep(sleep_time)
17
18     # Return the factorial of the number
19     return math.factorial(number)
20
21 if __name__ == "__main__":
22     # Numbers to calculate the factorial of
23     numbers = [1, 16, 4, 5, 8, 6, 2]
24
25     # Create ProcessPoolExecutor for parallel execution
26     with ProcessPoolExecutor(max_workers=4) as executor:
27         results = list(executor.map(compute_factorial, numbers))
28
29     # Provide finish feedback
30     print(f"Results: {results}")

```

6 pic. Module concurrent.futures Example Code Snippet

4. Free-threaded mode in Python 3.13

Python 3.13 was released on October 7, 2024 and it introduced significant new features and optimizations. The most notable changes include a revamped **interactive interpreter**, an experimental **Just-In-Time (JIT) compiler**, and an experimental **free-threaded mode**, which we will explore further.

The revolutionary **free-threaded mode** significantly improves Python's ability to manage multi-threaded applications. This mode disables the **Global Interpreter Lock (GIL)**, a long-standing feature introduced in Python 1.5 that prevented multiple threads from executing Python bytecode simultaneously, making it possible that only one thread can execute Python bytecode at a time. While the GIL ensured thread safety, it severely limited Python's performance in multi-threaded tasks.

Free-threaded mode addresses this limitation by enabling threads to execute in parallel across multiple CPU cores, fully utilizing multi-core processors. This feature is particularly advantageous for workloads that require multi-threading, such as simulations, data analysis, and other compute-intensive operations. With the enablement of true parallelism, significant performance improvements can be achieved, particularly in scenarios designed with threading in mind.

However, the feature currently remains **experimental**, and its impact may vary depending on the application:

- **Performance gains** - multi-threaded programs designed for concurrency can benefit from increased speed and efficiency on multi-core hardware.
- **Compatibility concerns** - legacy libraries or applications designed with the GIL assumption may encounter issues. Some libraries may not work or require updates to function correctly in free-threaded mode.
- **Single-threaded applications** - these applications might see performance decline due to the added complexity of thread management without the GIL.

To use free-threaded mode, developers must run their applications with a specialized executable, called *python3.13t* (or *python3.13t.exe*). Pre-built binaries for free-threaded mode are available in official *Windows* and

macOS installers. Alternatively, developers can also build CPython from source with the *--disable-gil* option for customized configurations [4].

At the time of writing this document, free-threaded mode is still experimental and under active development. Once major challenges are solved, it is expected to become a standard part of Python in a future release. For now, developers can try it out to see its benefits and limitations.

5. Comparing GIL and GIL-free Performance

In this section, we will explore the performance differences between traditional GIL-based execution and the new GIL-free mode by analyzing the benchmarks of their behavior across a range of scenarios, including CPU-bound tasks and I/O-bound tasks. The discussion will highlight the benefits and limitations of each approach. This analysis will provide basic understanding of how execution types influence Python application performance. **It is important to note that specific results may vary depending on each system.** Configurations, installed hardware, operating system, and Python implementation can have positive or negative effects of performance. However, the overall trends - such as the percentage of the performance gains in single-threaded, multi-threaded and multi-process of GIL-free execution versus GIL-enabled - are expected to remain consistent across most environments.

5.1. Count

We will begin our experimentation with a simple counting algorithm. The goal of this algorithm is to count to a specific number as fast as possible. Below is the code snippet for the single threaded version of this algorithm.

```
1  # A CPU-bound task: counting to a number
2  # Single Threaded (consecutive)
3  def count(n: int, start: int = 0) -> int:
4      result = 0
5      for _ in range(start, n, 1):
6          result += 1
7      return result
```

7 pic. Count Single-Threaded Code Snippet

Next, we extend this task to evaluate concurrent executions, utilizing Python's *threading* and *multiprocessing* modules. These approaches will allow us to assess the impact of threading and multiprocessing on performance.


```
1  # Multi Threaded (concurrent)
2  def threaded_count(n: int, num_threads: int) -> int:
3      threads = []
4      chunk_size = n // num_threads
5      results = [0] * num_threads
6
7      def __worker(index: int, start: int, end: int) -> None:
8          results[index] = count(end, start)
9
10     for i in range(num_threads):
11         start = i * chunk_size
12         end = (i + 1) * chunk_size if i < num_threads - 1 else n
13
14         thread = threading.Thread(target=__worker, args=(i, start, end))
15         threads.append(thread)
16
17     for thread in threads:
18         thread.start()
19
20     for thread in threads:
21         thread.join()
22
23     return sum(results)
```

8 pic. Count Multi-Threaded Code Snippet

```
1  # Multi Process (concurrent)
2  def process_count(n: int, num_processes: int) -> int:
3      chunk_size = n // num_processes
4
5      tasks = [
6          (
7              # end index
8              ((i + 1) * chunk_size if i < num_processes - 1 else n),
9              # start index
10             i * chunk_size,
11         )
12         for i in range(num_processes)
13     ]
14
15     with multiprocessing.Pool(processes=num_processes) as pool:
16         results = pool.starmap(count, tasks)
17
18     return sum(results)
```

9 pic. Count Multi-Process Code Snippet

Let us try evaluating a single target value for the counting algorithm. This initial test will provide a baseline comparison of single threaded, multi-threaded and multi-process executions with and without Global Interpreter Lock. The test was done with four threads, four processes and counting to one hundred million. The results highlight significant differences in

performance between GIL-enabled and GIL-free execution models across single-threaded, multi-threaded, and multi-process scenarios.

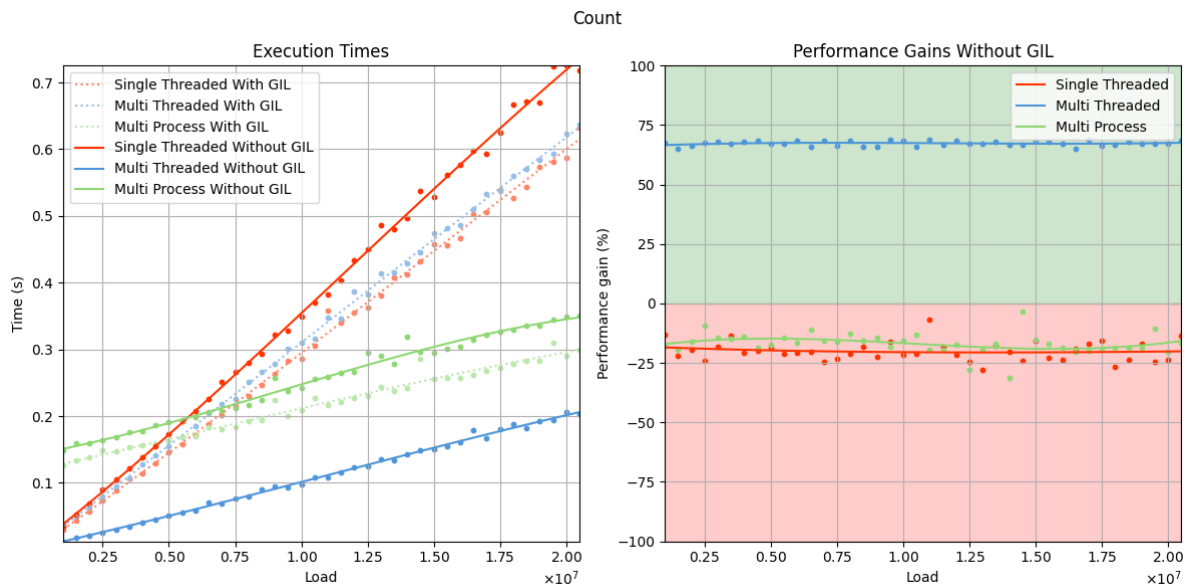
With GIL:

- **Single-threaded (2.89 seconds)** - execution performs as expected, providing a baseline for comparison. Since no concurrency is involved, the GIL does not impact this method directly.
- **Multi-threaded (3.06 seconds)** - execution is slightly slower than single-threaded due to GIL. Threads must wait for exclusive access to the GIL, which negates the potential performance gains of parallelism in CPU-bound tasks.
- **Multi-process (0.97 seconds)** - execution performs significantly faster because each process operates independently, bypassing the GIL. This demonstrates the advantage of multiprocessing for CPU-bound tasks in a GIL-enabled environment.

Without GIL:

- **Single-threaded (3.33 seconds)** - execution is slightly slower than its GIL-enabled counterpart. This may be due to implementation overheads introduced by removing the GIL.
- **Multi-threaded (0.96 seconds)** - execution sees a dramatic improvement without the GIL, completing faster than any other method. Threads can now execute concurrently and fully utilizing available CPU cores.
- **Multi-process (1.13 seconds)** - the performance of execution is slightly slower without the GIL compared to the GIL-enabled case. This may be attributed to the increased overhead of inter-process communication or other system-level variations when the GIL is not managing resource allocation.

Now, let us extend the analysis to an array of target values. For each value in the range, we will measure execution times for all methods under both GIL-enabled and GIL-free scenarios. These results will allow us to evaluate how the presence or absence of GIL influences performance as the computational workload increases. To present our findings, we will use graphical visualization to illustrate execution times across the target values. The graphs will highlight overall trends of different executions. Below is the visualization of gathered results from the target value array [1000000:20500000]. The load represents the number to which counting is happening.



10 pic. Count Performance Graphs

After evaluating the program, we can see clear differences between having the GIL enabled and disabled. With the GIL enabled, the multi-process algorithm performs the best, while the multi-threaded algorithm performs the worst, and the single-threaded algorithm performs slightly better than the multi-threaded algorithm. With the GIL disabled, the results shift entirely. The multi-threaded algorithm performs marginally better than the others, while the multi-process algorithm and single-threaded algorithm perform significantly worse, with the single-threaded being the least efficient.

From this point onward, testing of other algorithms will be performed in an equivalent manner.

5.2. Download Images

Next, we will be assessing a program used for downloading images from the internet. The goal of this program is to download images from the provided list of URLs as fast as possible. Below is the code snippet for the single threaded version of this algorithm.

```

1  # An I/O-bound task: downloading image
2  def __download_image(url: str, output_path: str) -> int:
3      size = 0
4      try:
5          response = requests.get(url, timeout=2, stream=True)
6          response.raise_for_status()
7          with open(output_path, "wb") as f:
8              for chunk in response.iter_content(chunk_size=8192):
9                  f.write(chunk)
10                 size += len(chunk)
11     except Exception as e:
12         logger.error(f"Failed to download {url}: {e}")
13     return size

```

11 pic. Download Images Task Code Snippet

```

1  # Single Threaded (consecutive)
2  def download_images(urls: list[str], output_folder_path: str) -> int:
3      size = 0
4      output_folder_path = os.path.join(
5          output_folder_path, DOWNLOAD_SINGLE_THREADED_RESULT_FOLDER_NAME
6      )
7      os.makedirs(output_folder_path, exist_ok=True)
8      for url in urls:
9          output_path = os.path.join(output_folder_path, url.split("/")[-1])
10         size += __download_image(url, output_path)
11     return size

```

12 pic. Download Images Single-Threaded Code Snippet

Next, we extend this task to evaluate concurrent executions, utilizing Python's *threading* and *multiprocessing* modules.

Parallel Processing in Python 3.13

```
1 # Multi Threaded (concurrent)
2 def threaded_download_images(
3     urls: list[str], output_folder_path: str, num_threads: int
4 ) -> None:
5     threads = []
6     chunk_size = len(urls) // num_threads
7
8     output_folder_path = os.path.join(
9         output_folder_path, DOWNLOAD_MULTI_THREADED_RESULT_FOLDER_NAME
10    )
11    os.makedirs(output_folder_path, exist_ok=True)
12
13    def __worker(start: int, end: int):
14        for url in urls[start:end]:
15            output_path = os.path.join(output_folder_path, url.split("/")[-1])
16            __download_image(url, output_path)
17
18    for i in range(num_threads):
19        start = i * chunk_size
20        end = (i + 1) * chunk_size if i < num_threads - 1 else len(urls)
21
22        thread = threading.Thread(target=__worker, args=(start, end))
23        threads.append(thread)
24
25    for thread in threads:
26        thread.start()
27
28    for thread in threads:
29        thread.join()
```

13 pic. Download Images Multi-Threaded Code Snippet

```
1 # Multi Process (concurrent)
2 def process_download_images(
3     urls: list[str], output_folder_path: str, num_processes: int
4 ) -> None:
5     output_folder_path = os.path.join(
6         output_folder_path, DOWNLOAD_MULTI_PROCESS_RESULT_FOLDER_NAME
7     )
8     os.makedirs(output_folder_path, exist_ok=True)
9
10    tasks = [
11        (
12            # url string
13            url,
14            # output path string
15            os.path.join(
16                output_folder_path,
17                url.split("/")[-1],
18            ),
19        )
20        for url in urls
21    ]
22
23    with multiprocessing.Pool(processes=num_processes) as pool:
24        pool.starmap(__download_image, tasks)
```

14 pic. Download Images Multi-Process Code Snippet

The single target value test was done with four threads, four processes and two hundred URLs. The results highlight minor differences in performance between GIL-enabled and GIL-free execution models across single-threaded, multi-threaded, and multi-process scenarios.

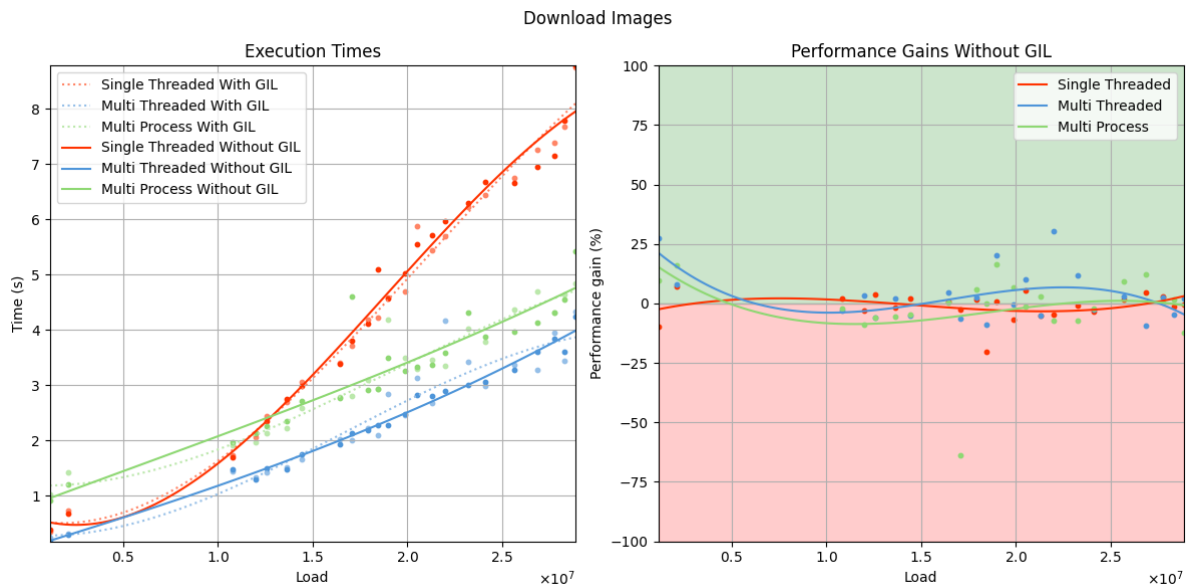
With GIL:

- **Single-threaded (20.74 seconds)** - this approach serves as the baseline for this comparison. It performs sequential processing of the URLs and is significantly slower due to the lack of concurrency.
- **Multi-threaded (9.68 seconds)** - execution shows a substantial improvement compared to single-threaded execution. In this I/O-bound scenario, the GIL does not heavily restrict threading performance because threads release the GIL during I/O operations, allowing other threads to proceed.
- **Multi-process (11.84 seconds)** - execution performs better than single-threaded but is slightly slower than multi-threaded. This is due to the higher overhead of creating and managing processes compared to threads, particularly for I/O-bound workloads.

Without GIL:

- **Single-threaded (19.94 seconds)** - execution is marginally faster without the GIL, due to minor optimizations in the underlying implementation.
- **Multi-threaded (8.19 seconds)** - execution benefits are minor from the absence of the GIL. The lack of GIL contention allows smoother thread execution even during brief periods of CPU-bound processing.
- **Multi-process (8.97 seconds)** - execution also improves without the GIL, reducing its execution time. While still slightly slower than multi-threaded execution, the gap is narrower compared to the GIL-enabled scenario.

Below is the visualization of gathered results from the target value array [1:80]. The load represents overall downloaded bytes.



15 pic. Download Images Performance Graphs

In this scenario, we observe no significant differences between running the program with the GIL enabled or disabled, although the application appears to run slightly more stably with the GIL disabled. These results could depend on numerous factors – for instance, the server might have limited the image downloads, creating a performance bottleneck. It would be better to evaluate this application in a more controlled environment.

5.3. Merge Sort

We continue our testing now with Merge Sort algorithm. Merge Sort is a divide-and-conquer sorting algorithm that divides an array into two halves, recursively sorts each half, and then merges the sorted halves back together. The process continues until the entire array is divided into individual elements, which are then combined in sorted order. Below is the code snippet for the single threaded version of this algorithm.

```

1  # A CPU-bound task: merging two sorted lists
2  def merge(left: list[int], right: list[int]) -> list[int]:
3      result = []
4      i = j = 0
5
6      while i < len(left) and j < len(right):
7          if left[i] <= right[j]:
8              result.append(left[i])
9              i += 1
10         else:
11             result.append(right[j])
12             j += 1
13
14     result.extend(left[i:])
15     result.extend(right[j:])
16
17     return result

```

16 pic. Merge Sort Task Code Snippet

```

1  # Single Threaded (consecutive)
2  def merge_sort(arr: list[int]) -> list[int]:
3      if len(arr) <= 1:
4          return arr
5
6      # Recursively split the array into halves
7      mid = len(arr) // 2
8      left = arr[:mid]
9      right = arr[mid:]
10
11     left = merge_sort(left)
12     right = merge_sort(right)
13
14     # Merge the sorted halves
15     return merge(left, right)

```

17 pic. Merge Sort Single-Threaded Code Snippet

Parallel Processing in Python 3.13

Next, we extend this task to evaluate concurrent executions, utilizing Python's *threading* and *multiprocessing* modules.

```
1  # Multi Threaded (concurrent)
2  def threaded_merge_sort(arr: list[int], num_threads: int) -> list[int]:
3      if len(arr) <= 1:
4          return arr
5
6      # Worker function
7      def worker(i: int, start: int, end: int) -> None:
8          results[i] = merge_sort(arr[start:end])
9
10     # Step 1: Divide the array into chunks
11     chunk_size = len(arr) // num_threads
12     threads = []
13     results = [None] * num_threads
14
15     # Step 2: Sort chunks in parallel
16     for i in range(num_threads):
17         start = i * chunk_size
18         end = (i + 1) * chunk_size if i != num_threads - 1 else len(arr)
19         thread = threading.Thread(target=worker, args=(i, start, end))
20         threads.append(thread)
21         thread.start()
22
23     for thread in threads:
24         thread.join()
25
26     # Step 3: Merge sorted chunks sequentially
27     while len(results) > 1:
28         next_round = []
29         for i in range(0, len(results), 2):
30             if i + 1 < len(results):
31                 next_round.append(merge(results[i], results[i + 1]))
32             else:
33                 next_round.append(results[i])
34         results = next_round
35
36     return results[0]
```

18 pic. Merge Sort Multi-Threaded Code Snippet

```
1 # Multi Process (concurrent)
2 def process_merge_sort(arr: list[int], num_processes: int) -> list[int]:
3     if len(arr) <= 1:
4         return arr
5
6     # Step 1: Divide the array into chunks
7     chunk_size = len(arr) // num_processes
8     chunks = [arr[i * chunk_size:(i + 1) * chunk_size] for i in range(num_processes - 1)]
9     chunks.append(arr[(num_processes - 1) * chunk_size:]) # Add remaining elements
10
11    # Step 2: Sort chunks in parallel
12    with multiprocessing.Pool(processes=num_processes) as pool:
13        results = pool.map(merge_sort, chunks)
14
15    # Step 3: Merge sorted chunks sequentially
16    while len(results) > 1:
17        next_round = []
18        for i in range(0, len(results), 2):
19            if i + 1 < len(results):
20                next_round.append(merge(results[i], results[i + 1]))
21            else:
22                next_round.append(results[i])
23        results = next_round
24
25    return results[0]
```

19 pic. Merge Sort Multi-Process Code Snippet

The single target value test was done with four threads, four processes and one million numbers in array. The results highlight minor change or performance decrease between GIL-enabled and GIL-free execution models across single-threaded, multi-threaded, and multi-process scenarios.

With GIL:

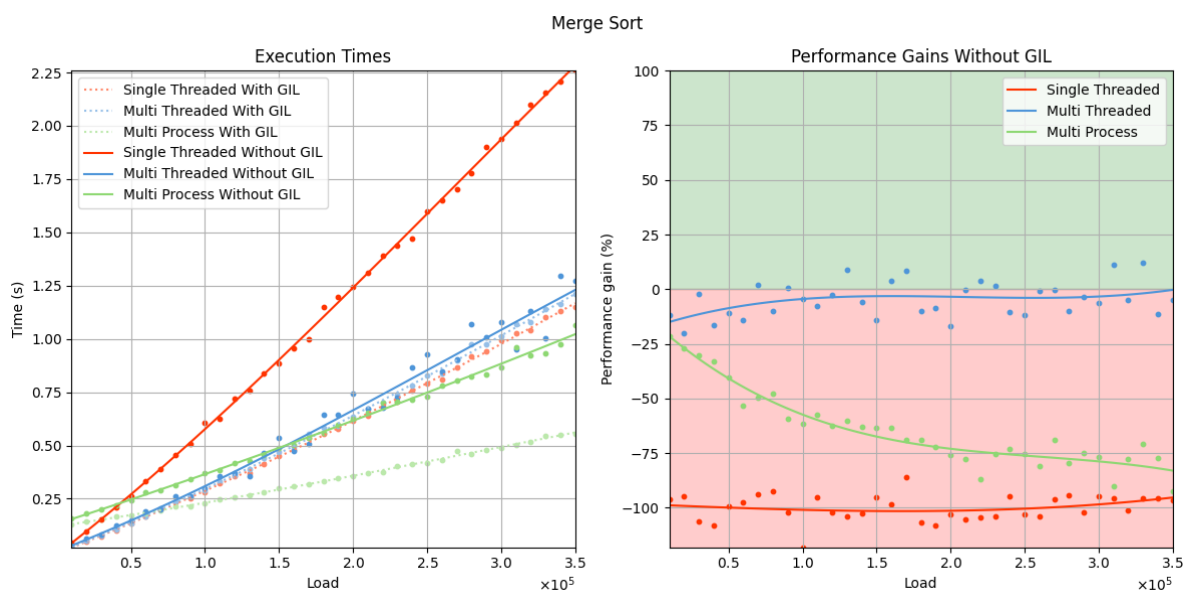
- **Single-threaded (3.83 seconds)** - execution performs well in a GIL-enabled environment as there is no concurrency. All computation occurs sequentially, avoiding the GIL overhead associated with thread context switching.
- **Multi-threaded (3.92 seconds)** - execution shows negligible improvement over single-threaded execution. In fact, it is slightly slower due to GIL contention, which serializes CPU-bound threads and introduces additional overhead.
- **Multi-process (1.64 seconds)** - multiprocessing performs significantly better, as each process runs in its own Python interpreter and bypasses the GIL entirely. This allows true parallel execution across CPU cores.

Without GIL:

- **Single-threaded (7.20 seconds)** - execution without the GIL is unexpectedly slower than its GIL-enabled counterpart. This is due to additional implementation overheads or less optimized handling in a GIL-free environment.

- **Multi-threaded (3.84 seconds)** - execution performs similarly to the GIL-enabled case. Without the GIL, threads can execute concurrently on multiple cores, but the lack of improvement suggests that other factors, such as thread management overhead, limit performance gains.
- **Multi-process (3.07 seconds)** - execution remains faster than single-threaded and multi-threaded approaches in the GIL-free scenario. However, the removal of the GIL results in slightly slower performance compared to the GIL-enabled case, due to system-level overhead in process communication.

Below is the visualization of gathered results from the target value array [10000:350000]. The load represents the count of numbers in an array.



20 pic. Merge Sort Performance Graphs

This time, we can see that with the GIL disabled, the algorithm performs worse across all tests. This could be because currently performance with GIL disabled is significantly degraded in single-threaded algorithms. Since the merge sort algorithm relies on a single thread during its final stage to merge the sorted subarrays into the final array, this bottleneck could be heavily impacting the overall performance.

5.4. Matrix Multiplication

Finally, we evaluate the matrix multiplication algorithm. Matrix multiplication involves multiplying two matrices, where each element of the result is the dot product of a row from the first matrix and a column from the second. In this case, we first transpose matrix B to optimize the multiplication. By transposing, we align the rows of matrix A with the columns of matrix B, which can improve efficiency. Below is the code snippet for the single threaded version of this algorithm.

```
1 # A CPU-bound task: matrix multiplication of row and column
2 def __row_col_multiplication(row: list[int], col: list[int]) -> int:
3     # Inner product of two vectors
4     return sum(a * b for a, b in zip(row, col))
```

21 pic. Matrix Multiplication Task Code Snippet

```
1 # Single Threaded (consecutive)
2 def matrix_multiplication(
3     matrix_a: list[list[int]], matrix_b: list[list[int]]
4 ) -> list[list[int]]:
5     if len(matrix_a[0]) != len(matrix_b):
6         raise ValueError("Matrix dimensions do not match.")
7
8     matrix_b_transposed = list(zip(*matrix_b))
9
10    # Matrix multiplication
11    return [
12        __row_col_multiplication(row, col) for col in matrix_b_transposed
13        for row in matrix_a
14    ]
```

22 pic. Matrix Multiplication Single-Threaded Code Snippet

Next, we extend this task to evaluate concurrent executions, utilizing Python's *concurrent.futures* and *multiprocessing* modules.

```
1 # Multi Threaded (concurrent)
2 def threaded_matrix_multiplication(
3     matrix_a: list[list[int]], matrix_b: list[list[int]], num_threads: int
4 ) -> list[list[int]]:
5     if len(matrix_a[0]) != len(matrix_b):
6         raise ValueError("Matrix dimensions do not match.")
7
8     # Worker function
9     def __worker(i: int) -> list[int]:
10         intermediate_row = []
11         for j in range(len(matrix_b_transposed)):
12             intermediate_row.append(
13                 __row_col_multiplication(matrix_a[i], matrix_b_transposed[j])
14             )
15         return intermediate_row
16
17     matrix_b_transposed = list(zip(*matrix_b))
18     result = []
19
20     with concurrent.futures.ThreadPoolExecutor(
21         max_workers=num_threads
22     ) as executor:
23         for i in range(len(matrix_a)):
24             result.append(executor.submit(__worker, i))
25
26     concurrent.futures.wait(result)
27     return [future.result() for future in result]
```

23 pic. Matrix Multiplication Multi-Threaded Code Snippet

```
1  # Multi Process (concurrent) worker function
2  def __process_worker(
3      i: int,
4      matrix_a: list[list[int]],
5      matrix_b_transposed: list[list[int]],
6  ) -> list[int]:
7      intermediate_row = []
8      for j in range(len(matrix_b_transposed)):
9          intermediate_row.append(
10             __row_col_multiplication(matrix_a[i], matrix_b_transposed[j])
11         )
12     return intermediate_row
13
14
15 # Multi Process (concurrent)
16 def process_matrix_multiplication(
17     matrix_a: list[list[int]], matrix_b: list[list[int]], num_processes: int
18 ) -> list[list[int]]:
19     if len(matrix_a[0]) != len(matrix_b):
20         raise ValueError("Matrix dimensions do not match.")
21
22     matrix_b_transposed = list(zip(*matrix_b))
23     tasks = [
24         (
25             # rows index
26             i,
27             # first matrix
28             matrix_a,
29             # second matrix transposed
30             matrix_b_transposed,
31         )
32         for i in range(len(matrix_a))
33     ]
34
35     with multiprocessing.Pool(processes=num_processes) as pool:
36         result = pool.starmap(__process_worker, tasks)
37
38     return result
```

24 pic. Matrix Multiplication Multi-Process Code Snippet

The single target value test was done with four threads, four processes, matrix a dimensions 543x543 and matrix b dimensions 543x777. The results highlight performance increase between GIL-enabled and GIL-free execution models across single-threaded, multi-threaded, and multi-process scenarios.

With GIL:

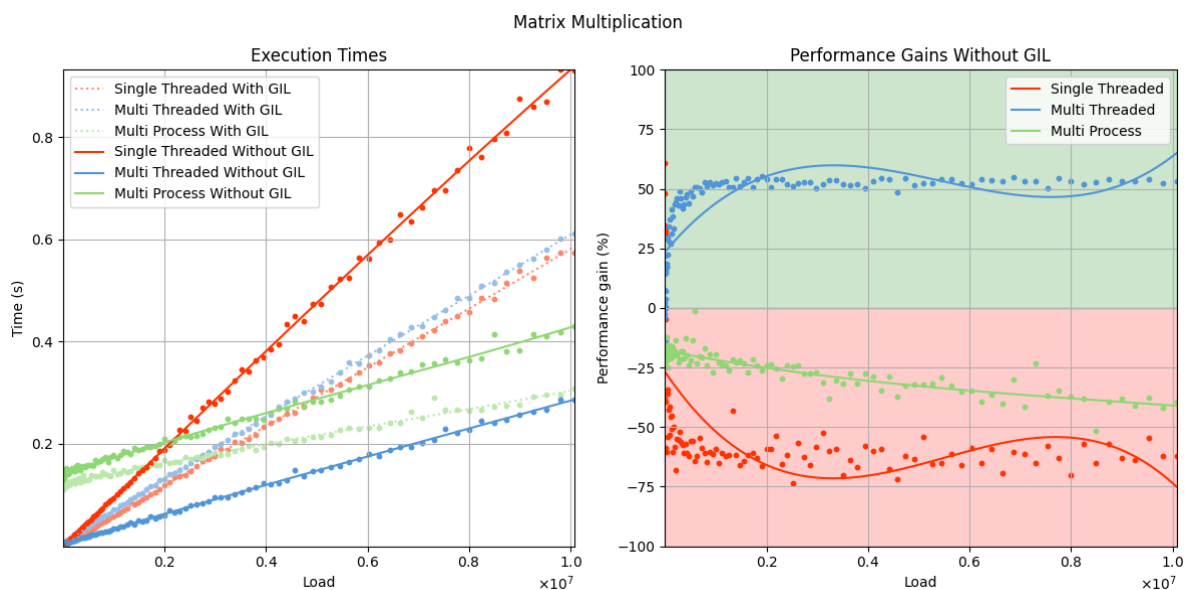
- **Single-threaded (12.98 seconds)** - execution serves as the baseline for comparison. It performs all computations sequentially, unaffected by the GIL, but cannot utilize multiple CPU cores.
- **Multi-threaded (13.52 seconds)** - execution is slightly slower than single-threaded execution. The GIL prevents true parallelism in this CPU-bound task, and thread management overhead adds to the execution time.

- **Multi-process (3.89 seconds)** - execution significantly outperforms single-threaded and multi-threaded approaches. By bypassing the GIL, it utilizes multiple CPU cores effectively for parallel computations.

Without GIL:

- **Single-threaded (19.26 seconds)** - execution without the GIL is unexpectedly slower. This performance drop may be attributed to additional implementation overheads or less efficient handling of single-threaded operations in a GIL-free environment.
- **Multi-threaded (5.71 seconds)** - execution improves dramatically without the GIL, as threads can now execute concurrently across CPU cores. This highlights the benefits of GIL removal for CPU-bound tasks.
- **Multi-process (5.87 seconds)** - execution performs similarly to multi-threaded execution in the GIL-free scenario, with a slight increase in overhead compared to the GIL-enabled case.

Below is the visualization of gathered results from the target value array. The load represents the total amount of multiplication done.



25 pic. Matrix Multiplication Performance Graphs

In the case of matrix multiplication, we see a substantial performance improvement when the GIL is disabled in multi-threaded test compared to tests with the GIL enabled. With the GIL enabled, the best algorithm was multi-process with single-threaded and multi-threaded algorithms performing similarly worse. After disabling the GIL, the situation changed, and the multi-threaded algorithm outperformed all the other tests.

6. Conclusion

The experimental removal of the Global Interpreter Lock (GIL) in Python 3.13 marks a significant step forward in addressing one of the most longstanding limitations of the language, particularly in the context of multi-threaded programming. With the increasing demand for concurrent processing, driven by the rise of artificial intelligence and machine learning applications, this move is seen as an essential shift in making Python more competitive in modern development environments. The removal of the GIL holds the potential to unlock substantial performance improvements in multi-threaded workloads, which have been limited for years due to GIL's enforcement of single-threaded bytecode execution.

However, the transition away from the GIL is not without its challenges. At the time of writing, many libraries still depend on GIL being supported and could not run or do not run properly, for example, the popular math library *NumPy*. Furthermore, our testing revealed mixed results when running multi-threaded applications in the GIL-free environment. In some cases, we saw significant performance improvements, while in others, there were no noticeable benefits or even a decline in performance compared to when the GIL was enabled.

The removal of the GIL is undoubtedly a positive and necessary direction for Python's future, but further optimization and library support are needed to fully realize the benefits. With continued development and fine-tuning, Python could finally overcome its multi-threading limitations, enabling faster and more scalable applications in a variety of use cases.

Sources

- [1] "threading - Thread-based parallelism - Python 3.13.1 documentation," Python Software Foundation, 27 December 2024. [Online]. Available: <https://docs.python.org/3/library/threading.html>. [Accessed 28 December 2024].
- [2] "multiprocessing - Process-based parallelism - Python 3.13.1 documentation," Python Software Foundation, 27 December 2024. [Online]. Available: <https://docs.python.org/3/library/multiprocessing.html>. [Accessed 28 December 2024].
- [3] "concurrent.futures - Launching parallel tasks - Python 3.13.1 documentation," Python Software Foundation, 27 December 2024. [Online]. Available: <https://docs.python.org/3/library/concurrent.futures.html>. [Accessed 28 December 2024].
- [4] A. Turner and T. Wouters, "What's New In Python 3.13 - Python 3.13.1 documentation," Python Software Foundation, 27 December 2024. [Online]. Available: <https://docs.python.org/3.13/whatsnew/3.13.html>. [Accessed 28 December 2024].