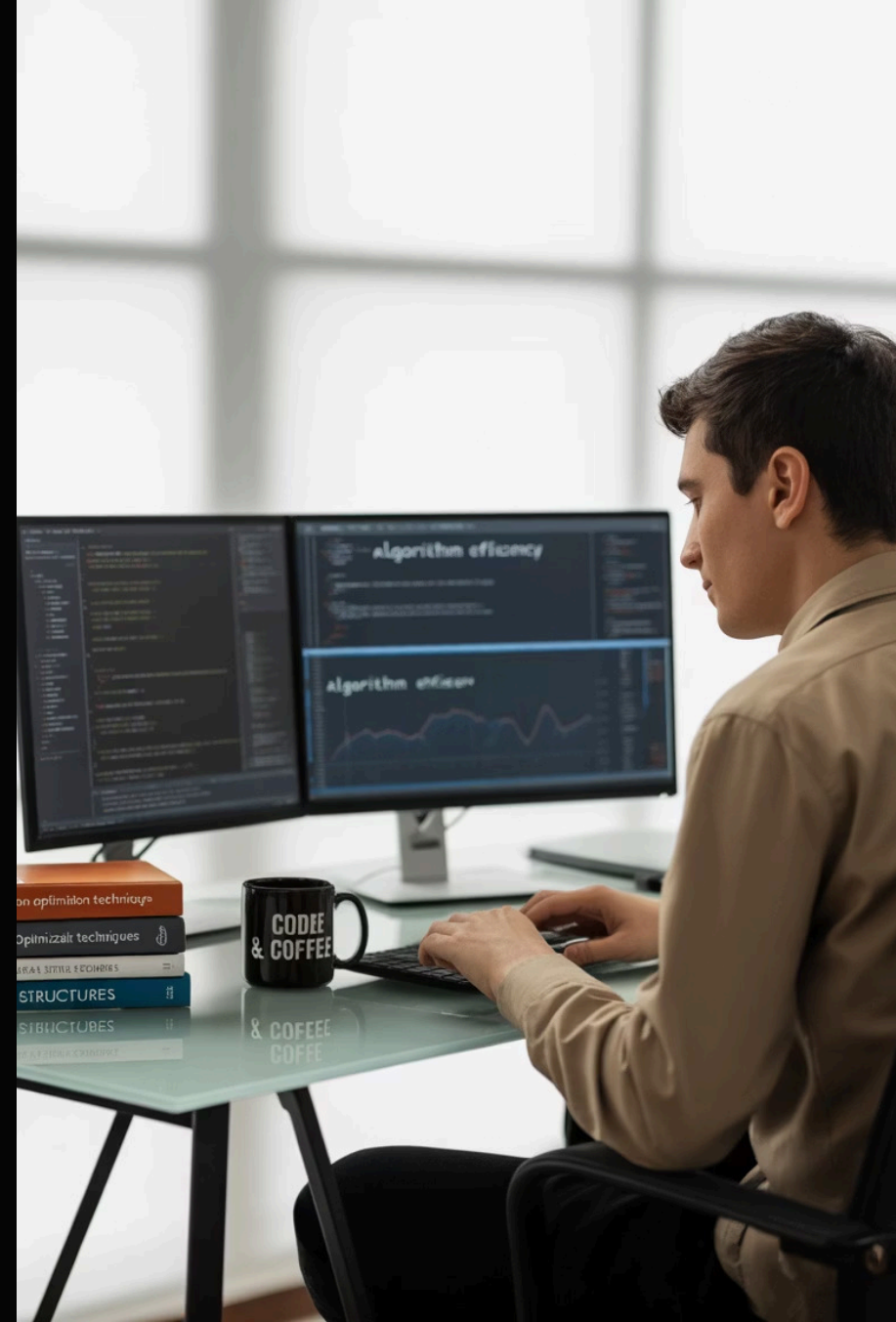


Python Advanced Concepts and Performance Optimization

A Technical Deep Dive Into Python Performance Optimization and Advanced Features

Instructor: Chandrashekar Babu <training@chandrashekar.info>

Website: <https://www.chandrashekar.info> | <https://www.slashprog.com>



Course Overview

Welcome to this intensive course on Python performance optimization and advanced concurrency. As experienced developers, you already understand Python's fundamentals, but this course will take your skills to the next level by exploring how to make your Python applications significantly faster, more efficient, and capable of handling complex concurrent workloads.

Day 1

01

Python Performance Profiling

Tools and techniques for identifying bottlenecks in your code

02

Threads, Processes, Generators and Coroutines

Understanding concurrency models and their appropriate use cases

03

Managing Concurrency with Event Loops

Exploring gevent, asyncio, uvloop, and modern concurrency patterns

04

Thread Pool vs Process Pool Executors

Practical guidance for implementing parallel execution strategies

Day 2

01

Python Performance Optimization Techniques

Exploring JIT/AOT compilers, PyPy, Numba, Cython and other performance-enhancing tools

02

Python Threads vs Process

Understanding the GIL, Python 3.12 enhancements, and strategies for effective concurrency

03

Decorators and AOP Patterns

Advanced implementation techniques and real-world applications of Aspect-Oriented Programming

Day 3

01

Gradual Typing

Optional static typing in Python, mypy, type hints, and advanced typing techniques

02

Buffer Protocol

Memory views, efficient data handling, and implementing custom buffer protocols

03

PySpark

Distributed computing, RDDs, DataFrames, and practical applicationsProgramming

About me

A FOSS Technologist and Corporate Trainer with 30+ years of experience in software development, technology consulting and corporate training.

A Pythonista since 1999

Also a polyglot with proficiency in many programming languages that include – Perl, Ruby, Tcl, PHP, Bash, Java, Clojure, C, Zig, Rust, x86 assembly language and many more...

Expertise in diverse technology domains in the Linux ecosystem (Linux Kernel / Device Driver development, MySQL / PostgreSQL Database systems, Apache web server administration and tuning, Web development using PHP, Rails, Flask, software architecture and design principles, agile / adaptive software methodologies).

To know more about me: kindly visit <https://www.chandrashekar.info/>

Your brief introductions please...

Your designation / role in your organization

Your experience

Your comfort-level in Python

Your comfort-level in OOP concepts and dynamic programming paradigms

Which other programming language you are familiar with (if any)

What You'll Learn Today



Advanced Profiling Techniques

How to systematically identify performance bottlenecks in your codebase using tools like cProfile, line_profiler, and memory_profiler with practical examples



Event-Driven Programming

Implementing high-performance event loops with gevent, asyncio and uvloop to handle thousands of concurrent connections



Concurrency Models in Python

When and how to use threads, processes, generators and coroutines effectively, including the fundamental differences between concurrency and parallelism



Parallel Execution Frameworks

Comparing Thread Pool and Process Pool executors with hands-on exercises to demonstrate real-world performance implications

Prerequisites

To get the most from today's session, please ensure you have:

Required Knowledge

- Strong Python programming fundamentals
- Experience with backend development
- Understanding of basic threading concepts
- Familiarity with Python 3.9+

Technical Setup

- Python 3.9+ installed
- Jupyter Notebook or JupyterLab
- Required packages: cProfile, line_profiler, gevent, asyncio
- Access to a multi-core system for parallel processing demos

This course assumes you're comfortable with Python syntax, decorators, context managers, and basic object-oriented programming. We won't be covering these fundamentals but will build upon them extensively.



Chapter 1: Python Performance Profiling

The Performance Mindset

Before diving into specific tools, it's crucial to develop a performance-oriented mindset when approaching Python development:

Measure First, Optimise Later

Resist the urge to optimise prematurely. Always collect data to identify actual bottlenecks rather than guessing where performance issues might be. Donald Knuth famously stated, "Premature optimization is the root of all evil."

Profile in Production-Like Environments

Development environments often mask performance issues that only appear at scale. Set up profiling in environments that closely mimic production with realistic data volumes and traffic patterns.

Optimise the Critical Path

Focus your efforts on the 20% of code that accounts for 80% of execution time. Optimising rarely-called functions will have minimal impact on overall performance.

Consider Trade-offs

Every optimization comes with trade-offs in terms of code readability, maintainability, and development time. Always evaluate whether the performance gain justifies these costs.

Remember: Python's flexibility makes it easy to write code quickly, but this same flexibility can lead to hidden performance pitfalls that are only discovered through systematic profiling.

cProfile: Python's Built-in Profiler

cProfile is Python's standard profiling tool, providing detailed statistics about function calls in your code:

Key Features

- Deterministic profiling (accounts for all function calls)
- Minimal performance overhead
- Built into the standard library
- Provides detailed call count and timing statistics
- Easy integration with existing codebases

```
import cProfile
import pstats
from pstats import SortKey

# Profile the function execution
cProfile.run('my_function()', 'profile_stats')

# Analyse the results
p = pstats.Stats('profile_stats')
p.strip_dirs().sort_stats(SortKey.CUMULATIVE).print_stats(10)
```

cProfile works by instrumenting your code and recording precise timing information about each function call. It's an excellent first step for identifying bottlenecks, especially in larger applications where the performance issues aren't immediately obvious.

The profiler provides several key metrics including **ncalls** (number of calls), **tottime** (total time spent in the function), **cumtime** (cumulative time in the function and all functions it calls), and **percall** (average time per call).

cProfile Output Analysis

Understanding how to interpret cProfile output is critical for effective optimization:

```
214748 function calls (205097 primitive calls) in 0.423 seconds
```

Ordered by: cumulative time

```
ncalls  tottime  percall  cumtime  percall filename:lineno(function)
    1    0.000    0.000    0.423    0.423 {built-in method builtins.exec}
    1    0.000    0.000    0.423    0.423 main.py:1()
    1    0.001    0.001    0.422    0.422 main.py:50(process_data)
  1000    0.010    0.000    0.410    0.000 data_processor.py:15(transform_entry)
  1000    0.350    0.000    0.400    0.000 data_processor.py:22(_complex_calculation)
  9651    0.050    0.000    0.050    0.000 {built-in method time.sleep}
```

In this example output:

- **process_data** is our main function, consuming 0.422 seconds total (cumtime)
- It calls **transform_entry** 1,000 times, which cumulatively takes 0.410 seconds
- The **_complex_calculation** function is called 1,000 times and consumes 0.400 seconds
- There's also significant time spent in **time.sleep**, indicating potential waiting periods

This analysis immediately points us to **_complex_calculation** as our primary bottleneck, which is consuming nearly 95% of our total execution time. This would be our first target for optimization.

Visualizing Profile Data

Textual profiling output can be difficult to interpret for complex applications. Visualization tools can help identify bottlenecks more effectively:

Popular Visualization Tools

- **SnakeViz**: Interactive visualizations for cProfile data
- **pyprof2calltree**: Convert profile data to KCacheGrind format
- **gprof2dot**: Generate call graphs from profiling data
- **Flame Graphs**: Hierarchical visualization showing time distribution

These tools transform raw profiling data into intuitive visualizations that help identify performance bottlenecks at a glance, particularly in complex codebases with deep call hierarchies.



Example flame graph visualization of profiling data showing call hierarchy and time distribution

```
# Generate a call graph using gprof2dot
pip install gprof2dot
python -m cProfile -o profile.pstats your_script.py
gprof2dot -f pstats profile.pstats | dot -Tpng -o profile_callgraph.png
```

timeit: Micro-benchmarking Tool

While cProfile helps identify bottlenecks in large applications, timeit is perfect for comparing specific implementations of functions or code snippets:

Key Features

- Designed for small code snippet benchmarking
- Runs code multiple times for statistical significance
- Minimizes system and Python interpreter variability
- Available as module and Jupyter magic command
- Excellent for algorithm comparison

```
import timeit

# Compare list comprehension vs. for loop
list_comp = timeit.timeit(
    '[i*2 for i in range(1000)]',
    number=10000
)

for_loop = timeit.timeit(
    """
    result = []
    for i in range(1000):
        result.append(i*2)
    """,
    number=10000
)

print(f"List comprehension: {list_comp:.5f} seconds")
print(f"For loop: {for_loop:.5f} seconds")
```

When used in Jupyter notebooks, the %%timeit magic provides an even more convenient interface:

```
%%timeit
# This cell will be timed over multiple runs
result = []
for i in range(10000):
    result.append(i*2)
```

line_profiler: Line-by-Line Analysis

When you need to go deeper than function-level profiling, line_profiler provides detailed timing information for each line of code:

Advantages

- Pinpoints exact lines causing bottlenecks
- Reveals inefficient constructs within functions
- Shows time spent on each line as percentage
- Perfect for optimizing complex functions
- Available as Jupyter magic with %lprun

```
# Install with pip install line_profiler
from line_profiler import LineProfiler

def process_data(items):
    result = []
    for item in items:
        if complex_condition(item):
            result.append(transform(item))
    return sorted(result, key=lambda x: x.priority)

# Profile the function
lp = LineProfiler()
lp_wrapper = lp(process_data)
lp_wrapper(data)
lp.print_stats()
```

The output will show timing for individual lines, helping you identify specific operations that consume disproportionate amounts of time, even within a single function. This granular view is invaluable for optimizing complex algorithms.

line_profiler Output Example

Timer unit: 1e-06 s

Total time: 5.73893 s

File: data_processor.py

Function: process_data at line 42

Line #	Hits	Time	Per Hit	% Time	Line Contents
=====					
42					def process_data(items):
43	1	18.0	18.0	0.0	result = []
44	1000	2548.0	2.5	0.0	for item in items:
45	1000	1240385.0	1240.4	21.6	if complex_condition(item):
46	534	3380921.0	6331.3	58.9	result.append(transform(item))
47	1	1115056.0	1115056.0	19.4	return sorted(result, key=lambda x: x.priority)

Analysis of this output reveals:



Identify

Line 46 (transform function) is consuming nearly 59% of execution time



Evaluate

The sorting operation (line 47) takes 19.4% despite being called only once



Optimize

The complex_condition check (line 45) uses 21.6% of time and should be improved



Validate

Re-profile after each optimization to confirm improvements

This granular view enables targeted optimizations that wouldn't be obvious from function-level profiling alone.

memory_profiler: Finding Memory Bottlenecks

While execution time is often the primary concern, memory usage can be equally important, especially for long-running services or applications processing large datasets:

Key Capabilities

- Line-by-line memory consumption tracking
- Identification of memory leaks
- Monitoring of peak memory usage
- Integration with Jupyter notebooks
- Time-based memory usage monitoring

```
# Install with pip install memory_profiler
from memory_profiler import profile

@profile
def process_large_dataset(filename):
    data = []
    with open(filename) as f:
        for line in f:
            data.append(parse_line(line))

    # Process the data
    result = [transform(item) for item in data]

    # Further operations
    return aggregate(result)
```

The output will show memory usage for each line, helping identify where large allocations occur. This is particularly valuable for optimizing applications that process large datasets or need to run in memory-constrained environments.

Common memory issues include: accumulating data in loops without bounds, inefficient data structures for the task, and keeping references to large objects longer than necessary.

Common Performance Bottlenecks in Python

Inefficient Data Structures

Using lists when sets or dictionaries would be more appropriate. For example, checking membership with `if x in large_list` is $O(n)$, while `if x in large_set` is $O(1)$.

```
# Slow (O(n) lookups)
values = [1, 2, 3, ... 10000]
if 9999 in values: # Linear search
```

```
# Fast (O(1) lookups)
values = set([1, 2, 3, ... 10000])
if 9999 in values: # Hash table
lookup
```

String Concatenation in Loops

Building strings with `+=` in loops creates new string objects repeatedly. Use `join()` or `io.StringIO` instead.

```
# Slow
result = ""
for i in range(10000):
    result += str(i) # Creates new
string each time
```

```
# Fast
result = "".join(str(i) for i in
range(10000))
```

Excessive Function Calls

Python function calls have overhead. Moving frequently called functions to C extensions or replacing with inline operations can improve performance.

```
# Function call overhead example
def square(x):
    return x * x
```

```
# Called in tight loop - expensive
result = [square(i) for i in
range(1000000)]
```

```
# Inline - faster
result = [i * i for i in range(1000000)]
```


More Performance Bottlenecks to Watch For

Global Variable Access

Accessing global variables is slower than local variables. When a variable will be used frequently in a function, consider creating a local reference.

```
# Slow - repeated global lookups
def process_items(items):
    for item in items:
        result =
        global_transformation(item)

# Faster - single global lookup
def process_items(items):
    transform = global_transformation
    for item in items:
        result = transform(item)
```

Unnecessary I/O Operations

File and network I/O are orders of magnitude slower than in-memory operations. Batch I/O operations where possible and consider caching results.

```
# Slow - separate I/O for each item
for item in items:
    with open('log.txt', 'a') as f:
        f.write(str(item) + '\n')

# Faster - batched I/O
with open('log.txt', 'a') as f:
    for item in items:
        f.write(str(item) + '\n')
```

Not Using Built-in Functions

Python's built-in functions are implemented in C and are typically much faster than equivalent Python implementations.

```
# Slow - manual sum calculation
total = 0
for num in numbers:
    total += num

# Fast - built-in sum function
total = sum(numbers)
```

Python-Specific Optimization Techniques

Using Collections Module

The collections module provides specialized container datatypes that are often more efficient than the standard ones:

- **defaultdict**: Dictionary with default factory for missing keys
- **Counter**: Dict subclass for counting hashable objects
- **deque**: Efficient append/pop from both ends
- **namedtuple**: Tuple subclass with named fields

```
from collections import defaultdict
```

```
# Instead of:
```

```
counts = {}
for word in document:
    if word not in counts:
        counts[word] = 0
    counts[word] += 1
```

```
# Use:
```

```
counts = defaultdict(int)
for word in document:
    counts[word] += 1
```

Leveraging itertools

The itertools module provides memory-efficient tools for creating iterators for efficient looping:

- **combinations**: Generate all combinations
- **permutations**: Generate all permutations
- **groupby**: Group sequential items
- **chain**: Combine multiple iterables

```
from itertools import islice, cycle
```

```
# Instead of:
```

```
def circular_slice(items, start, length):
    result = []
    size = len(items)
    for i in range(length):
        result.append(items[(start + i) % size])
    return result
```

```
# Use:
```

```
def circular_slice(items, start, length):
    return list(islice(cycle(items), start, start + length))
```

Leveraging NumPy for Numerical Operations

When dealing with numerical computations, NumPy can provide dramatic performance improvements:

NumPy Performance Benefits

- Vectorized operations instead of Python loops
- Efficient memory layout for numerical data
- Implementation in C with minimal Python overhead
- Optimized algorithms for mathematical operations
- SIMD (Single Instruction, Multiple Data) support

```
import numpy as np
import time

# Create data
size = 10_000_000
data = list(range(size))

# Python implementation
start = time.time()
python_result = [x**2 for x in data]
python_time = time.time() - start

# NumPy implementation
np_data = np.array(data)
start = time.time()
np_result = np_data**2
numpy_time = time.time() - start

print(f"Python: {python_time:.3f}s")
print(f"NumPy: {numpy_time:.3f}s")
print(f"NumPy is {python_time/numpy_time:.1f}x faster")
```

For this simple squaring operation, NumPy is typically 30-100x faster than pure Python. The speedup can be even more dramatic for complex operations like matrix multiplication, where NumPy leverages highly optimized BLAS libraries.

Practical Example: Optimizing a Data Processing Pipeline

Let's examine a real-world example of profiling and optimizing a data processing function:

```
# Original implementation
def process_dataset(filename):
    results = []
    with open(filename, 'r') as f:
        for line in f:
            if line.strip():
                data = parse_line(line)
                if is_valid(data):
                    transformed = transform_data(data)
                    results.append(transformed)

    # Sort by priority
    results.sort(key=lambda x: x['priority'])

    # Calculate statistics
    total = sum(r['value'] for r in results)
    average = total / len(results) if results else 0

    return {
        'results': results,
        'total': total,
        'average': average
    }
```

Profiling this function might reveal several opportunities for optimization, which we'll explore on the next slide.

Optimized Data Processing Pipeline

After profiling, we can apply several optimizations:

```
# Optimized implementation
def process_dataset(filename):
    # Use a generator expression to avoid loading everything into memory
    def parse_file():
        with open(filename, 'r') as f:
            for line in f:
                line = line.strip()
                if line: # Skip empty lines early
                    yield line

    # Combine operations and use generator pipeline
    valid_data = (
        transform_data(data)
        for line in parse_file()
        if (data := parse_line(line)) and is_valid(data)
    )

    # Materialize results only once
    results = sorted(valid_data, key=lambda x: x['priority'])

    # Calculate statistics in one pass
    count, total = 0, 0
    for r in results:
        count += 1
        total += r['value']

    return {
        'results': results,
        'total': total,
        'average': total / count if count else 0
    }
```

Key optimizations include: using generators to avoid loading the entire file into memory at once, combining operations to reduce loop overhead, and calculating statistics in a single pass through the data.



Hands-on Exercise: Profiling and Optimizing

For this exercise, you'll apply profiling tools to identify and fix performance issues in a real-world scenario:

1. Clone the exercise repository: `git clone https://github.com/chandrashekar-babu/python-optimization-workshop.git`
2. Navigate to the exercise directory: `cd python-optimization-workshop/exercises/profiling`
3. Run the unoptimized script: `python process_data.py`
4. Apply cProfile to identify bottlenecks: `python -m cProfile -o profile.stats process_data.py`
5. Analyze the profiling results and identify at least three performance issues
6. Implement your optimizations and measure the performance improvement

The goal is to achieve at least a 5x performance improvement while maintaining identical output. We'll review solutions together after the exercise.



Chapter 2: Threads, Processes, Generators and Coroutines

**Concurrent
Programming**

Understanding Python's Concurrency Models

Python offers several concurrency models, each with distinct characteristics:

Threads

Lightweight concurrent units that share memory space. Limited by Global Interpreter Lock (GIL) for CPU-bound tasks but effective for I/O-bound operations.

Processes

Separate Python interpreters with independent memory spaces. Not limited by GIL and ideal for CPU-bound tasks, but have higher overhead than threads.

Generators

Functions that can pause and resume execution, yielding values incrementally. Useful for memory-efficient iteration and implementing cooperative multitasking.

Coroutines

Advanced generators that can receive values and handle exceptions. Foundation of Python's async/await syntax for managing concurrent I/O operations.

The key to effective Python concurrency is understanding which model is appropriate for your specific workload. Using the wrong model can lead to suboptimal performance or even reduced throughput compared to sequential code.

Concurrency vs. Parallelism

Concurrency

Concurrency is about **structure** - it's the composition of independently executing tasks.

- Multiple tasks making progress over a period of time
- Not necessarily executing simultaneously
- Tasks may take turns using a resource
- Primarily about managing multiple things at once
- Often implemented via task switching

Best for: I/O-bound workloads where tasks spend time waiting (network requests, file operations, database queries)

Parallelism

Parallelism is about **execution** - it's the simultaneous execution of multiple tasks.

- Multiple tasks executing at literally the same moment
- Requires multiple processors/cores
- Each task has dedicated resources
- Primarily about doing multiple things at once
- Implemented via multiple execution units

Best for: CPU-bound workloads where tasks perform intensive calculations

In Python, threads provide concurrency but not necessarily parallelism (due to the GIL), while processes can provide true parallelism across multiple cores.

The Global Interpreter Lock (GIL)

The GIL is a mutex that protects access to Python objects, preventing multiple threads from executing Python bytecode simultaneously:

GIL Implications

- Only one thread can execute Python code at once
- CPU-bound threads don't run in parallel
- I/O operations release the GIL
- C extensions can release the GIL
- Present in CPython (standard Python)
- Not present in some implementations (e.g., Jython, IronPython)



The GIL ensures thread safety for Python's memory management but limits parallel execution of Python code.

The GIL simplifies Python's implementation and many C extensions rely on it for thread safety. It's crucial to understand that Python threads are still effective for I/O-bound tasks because the GIL is released during I/O operations, allowing other threads to run.

For CPU-bound tasks requiring true parallelism, multiprocessing is usually more effective than threading.

Introduction to Generators

Generators are a powerful Python feature that allows you to create iterators using a more concise syntax than defining a class with `__iter__` and `__next__` methods:

Key Characteristics

- Functions that use `yield` instead of `return`
- Maintain state between calls
- Lazy evaluation - values computed on demand
- Memory efficient for large sequences
- Can represent infinite sequences

```
def count_up_to(max):  
    count = 1  
    while count <= max:  
        yield count  
        count += 1  
  
# Using the generator  
counter = count_up_to(5)  
print(next(counter)) # 1  
print(next(counter)) # 2  
  
# Or in a loop  
for num in count_up_to(5):  
    print(num) # 1, 2, 3, 4, 5
```

Generators are especially useful when working with large datasets that would consume too much memory if loaded all at once. For example, processing a multi-gigabyte log file line by line using a generator avoids loading the entire file into memory.

They're also a building block for more advanced concurrency patterns in Python.

Generator Expressions

Generator expressions provide a concise way to create generators, similar to list comprehensions but with parentheses instead of square brackets:

```
# List comprehension - creates entire list in memory
squares_list = [x*x for x in range(1000000)]
# Uses a lot of memory immediately

# Generator expression - creates a generator
squares_gen = (x*x for x in range(1000000))
# Uses minimal memory until values are requested

# Memory usage comparison
import sys
list_size = sys.getsizeof(squares_list) # ~8MB
gen_size = sys.getsizeof(squares_gen)  # ~112 bytes
```

When to Use Generator Expressions

- Processing large datasets
- When you only need to iterate once
- Creating data transformation pipelines
- When memory efficiency is important
- When you don't need random access to elements

Generator expressions can be chained together to create memory-efficient data processing pipelines, with each stage transforming the data for the next.

Generator expressions support the same syntax as list comprehensions, including filtering and nested loops:

```
# Filter even numbers and square them
even_squares = (x*x for x in range(100) if x % 2 == 0)

# Process nested structures
flattened = (item for sublist in nested_list for item in sublist)
```

Advanced Generator Techniques

Generators have additional capabilities beyond simple iteration:

send(), throw(), and close()

```
def counter():
    count = 0
    while True:
        # yield returns a value AND receives
        # the value from .send()
        increment = yield count
        count += increment if increment else 1

c = counter()
print(next(c))    # 0 (must call next() first)
print(c.send(2))  # 2 (count += 2)
print(c.send(10)) # 12 (count += 10)
c.close()         # Exit the generator
```

yield from Expression

```
def chain_generators(*iterables):
    # Delegate to each iterable in sequence
    for it in iterables:
        yield from it

# Example usage
def range_1_to_3():
    yield 1
    yield 2
    yield 3

def range_4_to_6():
    yield 4
    yield 5
    yield 6

for num in chain_generators(range_1_to_3(),
                             range_4_to_6()):
    print(num) # Prints 1, 2, 3, 4, 5, 6
```

yield from is particularly useful for creating generator pipelines and implementing coroutines. It delegates to a subgenerator, allowing for more modular generator-based code.

From Generators to Coroutines

Coroutines evolved from generators but serve a different purpose:

Generators

- Primary purpose: lazy iteration
- Produce data
- Pull-based (consumer pulls data)
- One-way communication (yield values)
- Used for iterating over sequences

Coroutines

- Primary purpose: concurrent execution
- Consume and produce data
- Push-based (producer pushes data)
- Two-way communication (yield and send)
- Used for asynchronous operations

```
# A simple coroutine
def grep(pattern):
    print(f"Looking for {pattern}")
    while True:
        line = yield # Receive a value without sending one
        if pattern in line:
            print(f"Found: {line}")

# Using the coroutine
search = grep("python")
next(search) # Prime the coroutine
search.send("no match here") # Nothing happens
search.send("python is great") # Prints "Found: python is great"
```

Coroutines became the foundation for Python's modern `async/await` syntax, providing a more intuitive way to write asynchronous code.

Introduction to Threading

Python's threading module provides a high-level interface for working with threads:

Basic Threading Concepts

- Threads share the same memory space
- Lightweight compared to processes
- Subject to the GIL limitation
- Good for I/O-bound tasks
- Require careful synchronization

```
import threading
import time

def worker(name):
    print(f"Worker {name} starting")
    time.sleep(2) # Simulate I/O operation
    print(f"Worker {name} finished")

# Create and start threads
threads = []
for i in range(5):
    t = threading.Thread(target=worker, args=(i,))
    threads.append(t)
    t.start()

# Wait for all threads to complete
for t in threads:
    t.join()

print("All workers done")
```

When a thread performs I/O operations (like network requests, file operations, or the `time.sleep()` in our example), it releases the GIL, allowing other threads to execute. This makes threading effective for I/O-bound applications despite the GIL.

However, threads introduce complexity due to shared state and potential race conditions, which we'll address in the following slides.

Thread Synchronization

When multiple threads access shared data, synchronization mechanisms are necessary to prevent race conditions:

Common Synchronization Primitives

- **Lock:** Exclusive access to a resource
- **RLock:** Reentrant lock for nested locking
- **Semaphore:** Control access to a limited resource
- **Event:** Signal between threads
- **Condition:** Complex synchronization scenarios
- **Barrier:** Wait for multiple threads to reach a point

```
import threading

# Shared resource
counter = 0
counter_lock = threading.Lock()

def increment_counter(count):
    global counter
    for _ in range(count):
        # Critical section
        with counter_lock:
            counter += 1

# Create threads
t1 = threading.Thread(target=increment_counter,
                      args=(10000,))
t2 = threading.Thread(target=increment_counter,
                      args=(10000,))

# Start and join threads
t1.start()
t2.start()
t1.join()
t2.join()

print(f"Final counter value: {counter}")
```

Without the lock, the final counter value would likely be less than 20,000 due to race conditions where both threads read and increment the value simultaneously, causing some increments to be lost.

Thread Safety Patterns

Beyond basic locks, several patterns can help maintain thread safety:

Thread-Local Storage

Thread-local storage provides thread-specific data that isn't shared, eliminating the need for synchronization for that data.

```
import threading

# Create thread-local storage
local_data = threading.local()

def process_request(request_id):
    # Each thread has its own copy of
    # these attributes
    local_data.request_id = request_id
    local_data.results = []

    # Process the request

    local_data.results.append(f"Processed {request_id}")
    print(local_data.results)

threads = [

    threading.Thread(target=process_request, args=(i,))
    for i in range(5)
]
for t in threads:
    t.start()
```

Immutable Data Structures

Immutable objects can't be modified after creation, making them inherently thread-safe.

```
# Using immutable data structures
import collections

# Mutable (not thread-safe)
results = [] # List is mutable

# Immutable (thread-safe)
Point = collections.namedtuple('Point', ['x', 'y'])
p = Point(1, 2) # Cannot be modified after creation

# Functional style with immutable data
def add_point(points, x, y):
    return points + (Point(x, y),) # Creates new tuple
```

Queue-Based Communication

Using thread-safe queues to communicate between threads avoids shared state issues.

```
import queue
import threading

# Create a thread-safe queue
work_queue = queue.Queue()

def worker():
    while True:
        try:
            # Get work from the queue
            item = work_queue.get(timeout=1)
            print(f"Processing {item}")
            work_queue.task_done()
        except queue.Empty:
            break

# Start worker threads
threads = [threading.Thread(target=worker) for _ in range(3)]
for t in threads:
    t.daemon = True
    t.start()

# Add work to the queue
for i in range(10):
    work_queue.put(i)

# Wait for all work to be processed
work_queue.join()
```

Introduction to Multiprocessing

Python's multiprocessing module allows you to bypass the GIL limitation by using separate Python processes:

Multiprocessing Characteristics

- Each process has its own Python interpreter and memory space
- True parallelism across multiple CPU cores
- Higher overhead than threading
- Data must be explicitly shared or passed between processes
- Ideal for CPU-bound tasks

```
import multiprocessing
import time

def cpu_bound_task(n):
    """A CPU-intensive function."""
    count = 0
    for i in range(n):
        count += i
    return count

if __name__ == '__main__':
    # Create a pool of worker processes
    with multiprocessing.Pool() as pool:
        # Map the function to multiple inputs
        results = pool.map(
            cpu_bound_task,
            [10000000, 20000000, 30000000, 40000000]
        )

    print(f"Results: {results}")
```

The multiprocessing module mirrors the threading API in many ways, making it relatively easy to switch between the two. However, there are important differences in how data is shared and passed between execution units.

Inter-Process Communication

Since processes don't share memory, Python provides several mechanisms for communication between them:

Queues

Thread-safe FIFO queues for exchanging data between processes.

```
from multiprocessing import Process, Queue

def worker(q):
    q.put('Result from worker')

if __name__ == '__main__':
    q = Queue()
    p = Process(target=worker, args=(q,))
    p.start()
    result = q.get() # Blocks until data is available
    p.join()
    print(result)
```

Pipes

A two-way communication channel between processes.

```
from multiprocessing import Process, Pipe

def worker(conn):
    conn.send('Hello from worker')
    response = conn.recv()
    print(f"Worker received: {response}")
    conn.close()

if __name__ == '__main__':
    parent_conn, child_conn = Pipe()
    p = Process(target=worker, args=(child_conn,))
    p.start()
    print(f"Parent received: {parent_conn.recv()}")
    parent_conn.send('Hello from parent')
    p.join()
```

Shared Memory

Direct access to shared memory segments for high-performance data sharing.

```
from multiprocessing import Process, shared_memory
import numpy as np

def worker(shm_name, shape, dtype):
    # Attach to existing shared memory
    existing_shm = shared_memory.SharedMemory(name=shm_name)
    # Create array backed by shared memory
    np_array = np.ndarray(shape, dtype=dtype, buffer=existing_shm.buf)
    # Modify the array
    np_array[0] = 88
    existing_shm.close()

if __name__ == '__main__':
    # Create shared memory and array
    shm = shared_memory.SharedMemory(create=True, size=10)
    np_array = np.ndarray((2,), dtype=np.int64, buffer=shm.buf)
    np_array[0] = 42

    # Create and run process
    p = Process(target=worker, args=(shm.name, np_array.shape, np_array.dtype))
    p.start()
    p.join()

    print(f"After modification: {np_array[0]}")

    # Clean up
    shm.close()
    shm.unlink()
```

Process Pools and Worker Patterns

For common parallel processing tasks, the Pool class provides a higher-level interface:

Pool Methods

- **map()**: Apply function to every item, return results in order
- **apply()**: Apply function to single argument, blocking
- **apply_async()**: Non-blocking version of apply()
- **map_async()**: Non-blocking version of map()
- **starmap()**: Like map() but unpacks arguments from a sequence
- **imap()**: Lazy version of map()
- **imap_unordered()**: Like imap() but results can be out of order

```
from multiprocessing import Pool
import time

def process_chunk(chunk):
    """Process a chunk of data."""
    result = 0
    for i in range(chunk[0], chunk[1]):
        result += i
    return result

if __name__ == '__main__':
    # Define chunks of work
    chunks = [(1, 1000000), (1000000, 2000000),
              (2000000, 3000000), (3000000, 4000000)]

    # Sequential processing
    start = time.time()
    sequential_result = sum(process_chunk(chunk)
                             for chunk in chunks)
    sequential_time = time.time() - start

    # Parallel processing
    start = time.time()
    with Pool(processes=4) as pool:
        parallel_result = sum(pool.map(process_chunk,
                                       chunks))
    parallel_time = time.time() - start

    print(f"Sequential: {sequential_time:.2f}s")
    print(f"Parallel: {parallel_time:.2f}s")
    print(f"Speedup: {sequential_time/parallel_time:.1f}x")
```

Choosing Between Concurrency Models

I/O-Bound Tasks

- Tasks spend time waiting for external resources
- Examples: network requests, file operations, database queries
- Best approach: **Threads** or **Asyncio**
- GIL doesn't limit performance since Python releases GIL during I/O
- Asyncio has less overhead than threads but requires async-compatible libraries

High-Concurrency Servers

- Handling many simultaneous connections
- Examples: web servers, API services, chat servers
- Best approach: **Asyncio** or **Gevent**
- Can handle thousands of connections with minimal resources
- Non-blocking I/O maximizes throughput



CPU-Bound Tasks

- Tasks perform intensive calculations
- Examples: data processing, scientific computing, image processing
- Best approach: **Multiprocessing**
- Bypasses GIL limitation by using multiple processes
- Higher overhead but enables true parallelism

High-Volume Data Processing

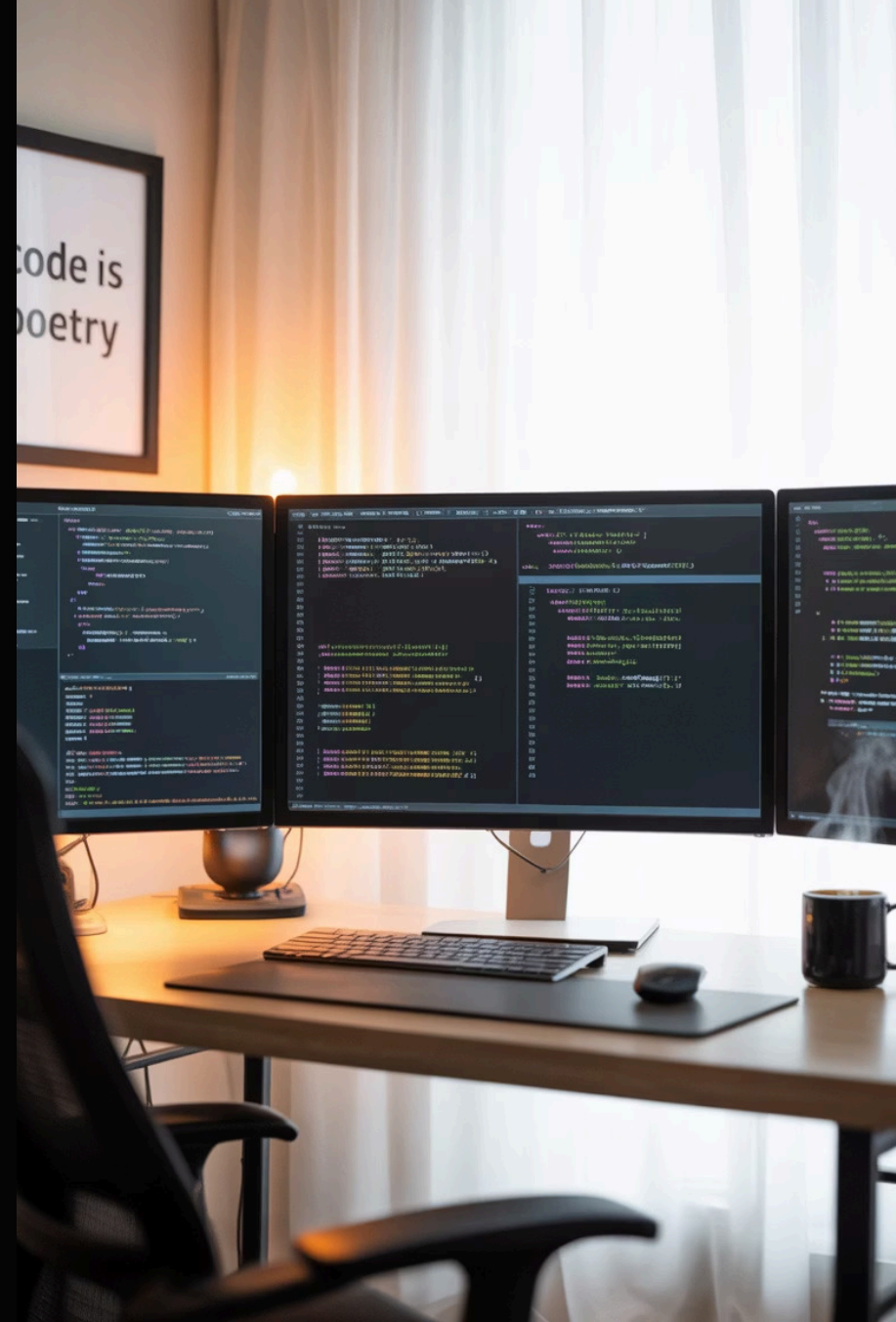
- Processing large datasets incrementally
- Examples: ETL pipelines, log processing
- Best approach: **Generators** combined with parallelism
- Memory-efficient processing of large datasets
- Can be combined with multiprocessing for parallel chunks

Hands-on Exercise: Concurrency Models

For this exercise, you'll implement the same task using different concurrency models and compare their performance:

1. Navigate to the exercise directory: `cd python-optimization-workshop/exercises/concurrency`
2. The task involves downloading and processing multiple files:
 - Sequential implementation is provided as a baseline
 - Implement using threading
 - Implement using multiprocessing
 - Implement using asyncio (if time permits)
3. Run the benchmark script to compare performance: `python benchmark.py`
4. Analyze the results and explain why certain models perform better for this specific task

This exercise will demonstrate how the nature of the workload (I/O-bound vs. CPU-bound) affects the relative performance of different concurrency models.



Chapter 3: Managing Concurrency with Event Loops

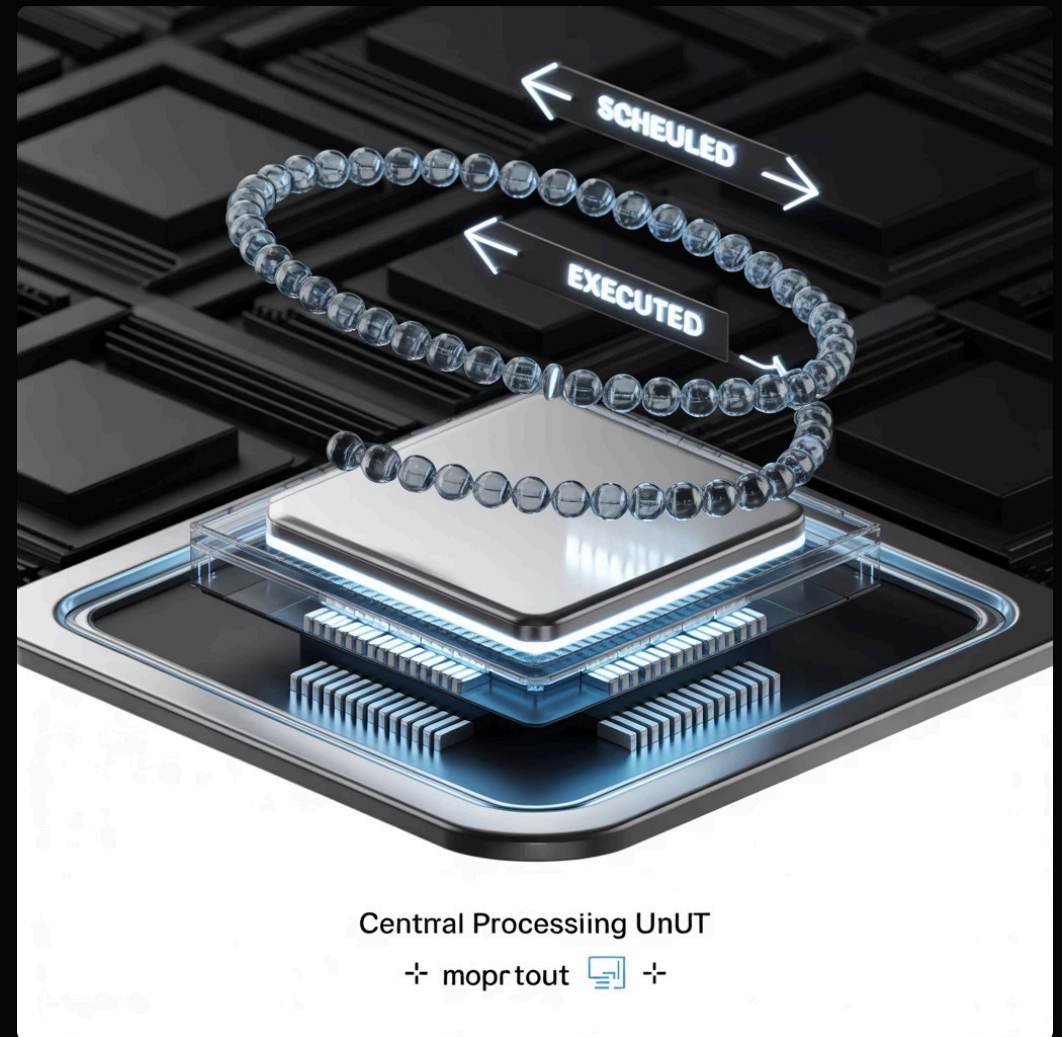


Understanding Event Loops

An event loop is the core of asynchronous programming, managing the execution of different tasks:

Event Loop Concepts

- Single-threaded, non-blocking execution model
- Tasks yield control when waiting for I/O
- Loop decides which task to run next
- Ideal for I/O-bound workloads
- Can handle thousands of concurrent operations



The event loop continuously checks for completed I/O operations and runs the appropriate callbacks.

Event loops differ from threading in that they provide concurrency within a single thread through cooperative multitasking. Tasks voluntarily yield control rather than being preemptively scheduled by the operating system.

Python has several event loop implementations, including:

- The built-in `asyncio` module's event loop
- Third-party libraries like `gevent` and `Tornado`
- High-performance alternatives like `uvloop`

Introduction to gevent

gevent is a coroutine-based Python networking library that uses greenlet to provide a high-level synchronous API on top of an asynchronous event loop:

Key Features

- Lightweight "greenlets" (microthreads)
- Implicit cooperative multitasking
- Monkey patching of standard library
- Familiar synchronous API
- High performance for networking

```
import gevent
from gevent import import socket

# Create a function to fetch a URL
def fetch(url):
    print(f'Fetching {url}')
    sock = socket.socket()
    sock.connect((url, 80))
    sock.send(f'GET / HTTP/1.0\r\nHost: {url}\r\n\r\n'.encode())
    data = sock.recv(1024)
    sock.close()
    print(f'Finished {url}: {len(data)} bytes')
    return data

# Run multiple fetches concurrently
urls = ['www.google.com', 'www.example.com', 'www.python.org']
jobs = [gevent.spawn(fetch, url) for url in urls]
gevent.joinall(jobs)
results = [job.value for job in jobs]
```

One of gevent's key advantages is its ability to patch standard library functions to be cooperative using monkey patching, allowing existing synchronous code to run concurrently without major changes.

Monkey Patching with gevent

Monkey patching is a technique where standard library modules are replaced at runtime with gevent-compatible versions:

How Monkey Patching Works

- Replaces blocking I/O operations with non-blocking versions
- Enables automatic yielding during I/O operations
- Makes existing code work with gevent without changes
- Should be done as early as possible in your application
- Can be selective (patch only specific modules)

```
# Import gevent and monkey patch BEFORE other imports
from gevent import monkey
monkey.patch_all()
```

```
# Now standard library modules will use gevent
import requests
import time
```

```
def fetch_url(url):
    print(f'Fetching {url}')
    start = time.time()
    response = requests.get(url)
    elapsed = time.time() - start
    print(f'Got {len(response.content)} bytes from "
          f"{url} in {elapsed:.2f} seconds")
    return response.content
```

```
# These will run concurrently thanks to monkey patching
urls = [
    "http://www.google.com",
    "http://www.github.com",
    "http://www.stackoverflow.com"
]
```

```
import gevent
jobs = [gevent.spawn(fetch_url, url) for url in urls]
gevent.joinall(jobs)
```

Monkey patching affects the entire Python process, so it should be used with care, especially in library code. It's generally safer to use it in application code where you control the entire environment.

Practical gevent Patterns

Pools and Groups

Managing collections of greenlets for concurrent execution with controlled parallelism.

```
from gevent.pool import Pool

def task(n):
    import time
    time.sleep(0.1) # Simulated I/O
    return n * n

# Limit concurrency to 5
# simultaneous tasks
pool = Pool(5)
results = pool.map(task, range(100))
print(f"Results: {len(results)} tasks completed")
```

Timeouts and Exception Handling

Adding timeouts to prevent blocked operations and handling exceptions in concurrent code.

```
import gevent
from gevent import Timeout

def potentially_blocking():
    import time
    time.sleep(10) # Long operation
    return "Completed"

# Set a timeout for the operation
timeout = Timeout(2)
timeout.start()

try:
    result = potentially_blocking()
    print(result)
except Timeout:
    print("Operation timed out")
finally:
    timeout.cancel()
```

Event-Based Communication

Coordinating between greenlets using events for signaling.

```
import gevent
from gevent.event import Event

# Create an event object
ready = Event()

def producer():
    print("Producer: Working...")
    gevent.sleep(2)
    print("Producer: Ready")
    ready.set() # Signal that data is ready

def consumer():
    print("Consumer: Waiting for data...")
    ready.wait() # Block until the event is set
    print("Consumer: Got the data!")

gevent.joinall([
    gevent.spawn(producer),
    gevent.spawn(consumer)
])
```

Introduction to asyncio

asyncio is Python's built-in library for asynchronous programming using the `async/await` syntax:

Key Components

- Coroutines defined with `async/await` syntax
- Event loop managing execution
- Tasks representing concurrent operations
- Futures representing eventual results
- Streams for network I/O

```
import asyncio
import time

async def say_after(delay, what):
    await asyncio.sleep(delay)
    print(what)

async def main():
    print(f'Started at {time.strftime('%X')}')

    # Run these concurrently
    task1 = asyncio.create_task(
        say_after(1, 'hello'))
    task2 = asyncio.create_task(
        say_after(2, 'world'))

    # Wait for both tasks to complete
    await task1
    await task2

    print(f'Finished at {time.strftime('%X')}')

asyncio.run(main())
```

Unlike `gevent`, `asyncio` requires explicit use of `async/await` syntax. Code must be written specifically to work with `asyncio`, but this explicitness can make the concurrency model clearer and easier to reason about.

asyncio Coroutines and Tasks

Coroutines and tasks are the basic building blocks of asyncio-based applications:

Coroutines

- Defined using `async def` syntax
- Can pause execution with `await`
- Return values using `return`
- Can be awaited by other coroutines
- Don't do anything until scheduled

```
async def fetch_data(url):  
    # Coroutine function  
    print(f"Fetching {url}")  
    await asyncio.sleep(1) # Simulate I/O  
    return f"Data from {url}"
```

Tasks

- Wraps a coroutine to schedule execution
- Runs concurrently with other tasks
- Can be cancelled or have timeouts
- Keeps track of state and result
- Created with `asyncio.create_task()`

```
async def main():  
    # Create tasks to run concurrently  
    task1 = asyncio.create_task(fetch_data("url1"))  
    task2 = asyncio.create_task(fetch_data("url2"))  
  
    # Wait for both to complete  
    results = await asyncio.gather(task1, task2)  
    print(results)
```

The key difference between a coroutine and a task is that a coroutine is just a function definition, while a task is a scheduled execution of a coroutine in the event loop.

asyncio Practical Patterns

Gathering Results

Running multiple coroutines concurrently and collecting their results.

```
import asyncio

async def fetch(url):
    await asyncio.sleep(1) # Simulate
    network delay
    return f"Result from {url}"

async def main():
    urls = ["url1", "url2", "url3", "url4"]

    # Run all coroutines concurrently
    results = await asyncio.gather(
        *[fetch(url) for url in urls]
    )

    # Process results
    for url, result in zip(urls, results):
        print(f"{url}: {result}")

asyncio.run(main())
```

Timeouts

Adding timeouts to prevent operations from blocking indefinitely.

```
import asyncio

async def slow_operation():
    await asyncio.sleep(10) # Very
    slow operation
    return "Completed"

async def main():
    try:
        # Set a 2-second timeout
        result = await asyncio.wait_for(
            slow_operation(),
            timeout=2
        )
        print(result)
    except asyncio.TimeoutError:
        print("Operation timed out")

asyncio.run(main())
```

Cancellation

Cancelling tasks that are no longer needed.

```
import asyncio

async def cancellable_operation():
    try:
        while True:
            print("Working...")
            await asyncio.sleep(0.5)
    except asyncio.CancelledError:
        print("Operation was cancelled")
        # Cleanup code here
        raise # Re-raise to propagate
        cancellation

async def main():
    task =
    asyncio.create_task(cancellable_oper
        ation())

    # Let it run for a bit
    await asyncio.sleep(2)

    # Cancel the task
    task.cancel()

    try:
        await task
    except asyncio.CancelledError:
        print("Main: task was cancelled")

asyncio.run(main())
```

Asynchronous Context Managers and Iterators

asyncio extends Python's context manager and iterator protocols with asynchronous versions:

Async Context Managers

Context managers that can perform asynchronous operations during entry and exit.

```
import asyncio

class AsyncResource:
    async def __aenter__(self):
        # Asynchronous setup
        await asyncio.sleep(1)
        print("Resource acquired")
        return self

    async def __aexit__(self, exc_type, exc_val, exc_tb):
        # Asynchronous cleanup
        await asyncio.sleep(1)
        print("Resource released")

async def main():
    async with AsyncResource() as resource:
        print("Using resource")
        await asyncio.sleep(1)

asyncio.run(main())
```

Async Iterators

Iterators that can perform asynchronous operations to produce the next value.

```
import asyncio

class AsyncCounter:
    def __init__(self, stop):
        self.current = 0
        self.stop = stop

    def __aiter__(self):
        return self

    async def __anext__(self):
        if self.current < self.stop:
            await asyncio.sleep(0.5) # Simulate I/O
            self.current += 1
            return self.current - 1
        raise StopAsyncIteration

async def main():
    # Iterate asynchronously
    async for i in AsyncCounter(5):
        print(f"Got {i}")

asyncio.run(main())
```

These features allow for more expressive and cleaner asynchronous code by integrating with Python's existing patterns.

Introduction to uvloop

uvloop is a drop-in replacement for asyncio's event loop that significantly improves performance:

Key Features

- Built on top of libuv (same library used by Node.js)
- 2-4x faster than asyncio's default loop
- Compatible with asyncio API
- Optimized for I/O operations
- Simple to integrate into existing asyncio code

```
import asyncio
import uvloop

# Install uvloop as the default event loop
uvloop.install()

async def hello_world():
    print("Hello World")
    await asyncio.sleep(0.1)
    print("Goodbye World")

# The asyncio.run() will now use uvloop
asyncio.run(hello_world())
```

Benchmarks typically show uvloop performing 2-4 times faster than the standard asyncio event loop for I/O-bound operations, making it an excellent choice for high-performance asyncio applications.

uvloop is particularly beneficial for applications that handle many concurrent connections, such as web servers, API services, and other network-intensive applications.

However, note that uvloop is not compatible with Windows directly (though it works in WSL), and it may have some minor compatibility issues with certain asyncio features.

Introduction to trio

trio is an alternative async/await-based concurrency library that emphasizes correctness and usability:

trio's Design Philosophy

- Focus on usability and correctness
- Structured concurrency model
- Built-in cancellation and timeout handling
- Simplified API compared to asyncio
- Strong emphasis on avoiding common pitfalls

```
import trio

async def child1():
    print("child1 started")
    await trio.sleep(1)
    print("child1 finished")

async def child2():
    print("child2 started")
    await trio.sleep(2)
    print("child2 finished")

async def parent():
    print("parent started")
    async with trio.open_nursery() as nursery:
        # Both children run concurrently
        nursery.start_soon(child1)
        nursery.start_soon(child2)
        # Nursery blocks until all children complete
    print("parent finished")

trio.run(parent)
```

trio's "nursery" concept is central to its structured concurrency model. A nursery is a context manager that ensures all child tasks complete (or are properly cancelled) before the parent continues. This helps prevent common issues like "forgotten" tasks or improper error handling.

Choosing an Event Loop Implementation

With multiple event loop implementations available, it's important to choose the right one for your specific needs:

asyncio

Best for: General-purpose async code with standard library integration

- **Pros:** Built-in, well-documented, widely supported
- **Cons:** Not the fastest implementation, more complex API
- **Use when:** Writing portable code, integrating with other asyncio libraries

uvloop

Best for: High-performance asyncio applications

- **Pros:** 2-4x performance boost, drop-in replacement
- **Cons:** Limited Windows support, occasional compatibility issues
- **Use when:** Performance is critical, especially for network servers

gevent

Best for: Making existing synchronous code concurrent

- **Pros:** Works with standard libraries via monkey patching
- **Cons:** Can be hard to reason about, less explicit
- **Use when:** Retrofitting concurrency into existing codebases

trio

Best for: Emphasis on correctness and clarity

- **Pros:** Structured concurrency, intuitive API
- **Cons:** Smaller ecosystem, not compatible with asyncio
- **Use when:** Building new projects with emphasis on robustness

Many libraries like FastAPI, Starlette, and aiohttp work with multiple event loop implementations, allowing you to choose the one that best fits your needs while maintaining compatibility with the ecosystem.

Best Practices for Event-Loop Based Concurrency

1

Never Block the Event Loop

The single most important rule for event loop-based concurrency is to never block the loop with CPU-intensive operations or blocking I/O.

When you need to perform such operations, use:

- `asyncio.to_thread()` for blocking I/O
- `run_in_executor()` for CPU-bound functions
- Consider moving CPU-intensive operations to separate processes

2

Handle Cancellation Properly

Ensure your coroutines handle cancellation gracefully, performing necessary cleanup when cancelled:

```
async def task():
    try:
        # Do work
        await asyncio.sleep(1)
    except asyncio.CancelledError:
        # Clean up resources
        print("Task cancelled, cleaning up")
        raise # Re-raise to propagate cancellation
```

3

Use Timeouts Liberally

Add timeouts to prevent operations from blocking indefinitely:

```
try:
    result = await asyncio.wait_for(
        potentially_slow_operation(),
        timeout=5.0
    )
except asyncio.TimeoutError:
    # Handle timeout case
    pass
```

4

Avoid Mixing Concurrency Models

Choose one concurrency model and stick with it. Mixing threading, multiprocessing, and asyncio in the same application can lead to complex interactions and bugs. If you need to use multiple models, keep them clearly separated with well-defined interfaces.



Hands-on Exercise: Event Loop Concurrency

For this exercise, you'll implement a simple web crawler using different event loop-based concurrency approaches:

1. Navigate to the exercise directory: `cd python-optimization-workshop/exercises/event_loops`
2. The task involves crawling multiple web pages and extracting links:
 - Implement using standard `asyncio`
 - Implement using `gevent`
 - Optionally, implement using `uvloop` or `trio`
3. Run the benchmark script to compare performance: `python benchmark.py`
4. Pay attention to:
 - Code clarity and error handling
 - Performance characteristics
 - Memory usage during execution

This exercise will demonstrate how different event loop implementations handle I/O-bound concurrency and help you understand their practical differences.



Chapter 4: Thread Pool Executor vs Process Pool Executor

The concurrent.futures Module

concurrent.futures provides a high-level interface for asynchronously executing callables using threads or processes:

Key Components

- **Executor:** Abstract base class that provides methods to execute calls
- **ThreadPoolExecutor:** Uses a pool of threads to execute calls
- **ProcessPoolExecutor:** Uses a pool of processes to execute calls
- **Future:** Encapsulates the asynchronous execution of a callable

```
from concurrent.futures import ThreadPoolExecutor
import time

def task(n):
    print(f'Processing {n}')
    time.sleep(1) # Simulate I/O or computation
    return n * n

# Create a thread pool
with ThreadPoolExecutor(max_workers=4) as executor:
    # Submit tasks and get Future objects
    futures = [executor.submit(task, i) for i in range(10)]

    # Wait for all tasks to complete and get results
    for future in futures:
        result = future.result()
        print(f'Result: {result}')
```

The concurrent.futures module simplifies parallel execution by providing a common interface for both thread-based and process-based parallelism. This makes it easy to switch between the two models based on the nature of your workload.

ThreadPoolExecutor in Depth

ThreadPoolExecutor manages a pool of worker threads for executing tasks asynchronously:

Key Characteristics

- Shares memory with the main process
- Limited by the GIL for CPU-bound tasks
- Lower overhead than processes
- Ideal for I/O-bound tasks
- Thread safety considerations apply

```
from concurrent.futures import ThreadPoolExecutor
import requests
import time

def fetch_url(url):
    """Fetch a URL and return the response text."""
    start = time.time()
    response = requests.get(url)
    return {
        'url': url,
        'status': response.status_code,
        'time': time.time() - start
    }

urls = [
    'https://www.python.org',
    'https://www.github.com',
    'https://www.stackoverflow.com',
    'https://www.google.com',
    'https://www.bbc.co.uk'
]

# Number of worker threads
with ThreadPoolExecutor(max_workers=5) as executor:
    # map() applies the function to each item in parallel
    results = list(executor.map(fetch_url, urls))

for result in results:
    print(f'{result["url"]}: {result["status"]} in '
          f'{result["time"]:.2f}s')
```

ThreadPoolExecutor is particularly well-suited for network I/O, file I/O, and other operations where threads spend time waiting rather than actively using the CPU.

ProcessPoolExecutor in Depth

ProcessPoolExecutor manages a pool of worker processes for executing tasks in parallel:

Key Characteristics

- Separate memory space for each process
- Not limited by the GIL
- Higher overhead than threads
- Ideal for CPU-bound tasks
- Data must be serialized between processes

```
from concurrent.futures import ProcessPoolExecutor
import time

def cpu_intensive_task(n):
    """A CPU-bound task that calculates the sum of squares."""
    start = time.time()
    result = sum(i*i for i in range(n))
    return {
        'input': n,
        'result': result,
        'time': time.time() - start
    }

inputs = [10000000, 20000000, 30000000, 40000000]

# Sequential execution for comparison
start = time.time()
sequential_results = [cpu_intensive_task(n) for n in inputs]
sequential_time = time.time() - start

# Parallel execution with processes
start = time.time()
with ProcessPoolExecutor() as executor:
    parallel_results = list(executor.map(
        cpu_intensive_task, inputs))
parallel_time = time.time() - start

print(f"Sequential: {sequential_time:.2f}s")
print(f"Parallel: {parallel_time:.2f}s")
print(f"Speedup: {sequential_time/parallel_time:.2f}x")
```

ProcessPoolExecutor is well-suited for CPU-intensive tasks that can be parallelized, such as numerical computations, data processing, and simulations.

Advanced Executor Patterns

as_completed() for Dynamic Results

Process results as they become available rather than waiting for all tasks to finish:

```
from concurrent.futures import
ThreadPoolExecutor, as_completed
import time
import random

def task(id):
    # Simulate variable completion
    times
    sleep_time = random.uniform(0.5,
3.0)
    time.sleep(sleep_time)
    return f'Task {id} completed in
{sleep_time:.2f}s'

with
ThreadPoolExecutor(max_workers=1
0) as executor:
    # Submit 20 tasks
    futures = [executor.submit(task, i)
for i in range(20)]

    # Process results as they complete
    for future in
as_completed(futures):
        try:
            result = future.result()
            print(result)
        except Exception as e:
            print(f'Task generated an
exception: {e}')
```

wait() with Timeouts

Wait for tasks with timeout control and handle different completion states:

```
from concurrent.futures import
ThreadPoolExecutor, wait,
FIRST_COMPLETED
import time

def task(id, duration):
    time.sleep(duration)
    return f'Task {id} completed'

with
ThreadPoolExecutor(max_workers=5)
as executor:
    # Submit tasks with different
    durations
    futures = [
        executor.submit(task, 1, 1), #
Quick
        executor.submit(task, 2, 2), #
Medium
        executor.submit(task, 3, 10) #
Very slow
    ]

    # Wait for first task to complete or
    timeout
    done, not_done = wait(
        futures,
        timeout=3,
        return_when=FIRST_COMPLETED
    )

    # Process completed tasks
    for future in done:
        print(future.result())

    # Cancel remaining tasks
    for future in not_done:
        future.cancel()
        print(f'Cancelled a task')
```

Callbacks on Completion

Execute code when a task completes without blocking:

```
from concurrent.futures import
ThreadPoolExecutor
import time

def task(id):
    time.sleep(1)
    return f'Result from task {id}'

def callback(future):
    # This is called when the task
    completes
    try:
        result = future.result()
        print(f'Callback received: {result}')
    except Exception as e:
        print(f'Task failed with: {e}')

with
ThreadPoolExecutor(max_workers=5)
as executor:
    for i in range(5):
        future = executor.submit(task, i)
        # Add a callback to be executed
        when done
        future.add_done_callback(callback)

    print("All tasks submitted, main
thread continues...")
    # Main thread can do other work
    here
```

Chaining Tasks with Executors

Complex workflows often require chaining tasks, where the output of one task becomes the input for another:

```
from concurrent.futures import ThreadPoolExecutor
import time

def fetch_data(url):
    """Simulates fetching data from a URL."""
    print(f"Fetching {url}")
    time.sleep(1) # Simulate network delay
    return f"Data from {url}"

def process_data(data):
    """Simulates processing the fetched data."""
    print(f"Processing {data}")
    time.sleep(0.5) # Simulate processing time
    return f"Processed {data}"

def save_result(processed_data):
    """Simulates saving the processed data."""
    print(f"Saving {processed_data}")
    time.sleep(0.5) # Simulate I/O operation
    return f"Saved {processed_data}"

urls = [f"url{i}" for i in range(5)]

with ThreadPoolExecutor(max_workers=5) as executor:
    # First stage: fetch data
    future_data = {executor.submit(fetch_data, url): url for url in urls}

    # Second stage: process data as it becomes available
    future_processed = {}
    for future in as_completed(future_data):
        url = future_data[future]
        try:
            data = future.result()
            # Submit the processing task for this result
            future_processed[executor.submit(process_data, data)] = data
        except Exception as e:
            print(f"Error fetching {url}: {e}")

    # Third stage: save processed data
    future_saved = {}
    for future in as_completed(future_processed):
        data = future_processed[future]
        try:
            processed = future.result()
            # Submit the save task for this result
            future_saved[executor.submit(save_result, processed)] = processed
        except Exception as e:
            print(f"Error processing {data}: {e}")

    # Wait for all save operations to complete
    for future in as_completed(future_saved):
        processed = future_saved[future]
        try:
            result = future.result()
            print(f"Completed workflow: {result}")
        except Exception as e:
            print(f"Error saving {processed}: {e}")
```

This pattern allows for efficient pipeline processing, where each stage can be parallelized and data flows through the pipeline as soon as it's ready for the next stage.

Choosing Between Thread and Process Pools

ThreadPoolExecutor

Best for: I/O-bound tasks

Advantages:

- Lower memory overhead
- Faster startup times
- Shared memory access
- More efficient for many small tasks

Disadvantages:

- Limited by GIL for CPU work
- Requires thread safety consideration
- Potential for race conditions

Examples: Web scraping, file I/O, API calls

ProcessPoolExecutor

Best for: CPU-bound tasks

Advantages:

- True parallelism across CPU cores
- Not limited by GIL
- Process isolation prevents shared state issues

Disadvantages:

- Higher memory overhead
- Slower startup time
- Data serialization overhead
- Limited by number of CPU cores

Examples: Data processing, numerical computation, image processing

32x

Memory overhead

A process typically uses about 32 times more memory than a thread in Python

10ms

Thread creation

Average time to create a new thread

100ms

Process creation

Average time to create a new process

4x

Process speedup

Typical speedup for CPU-bound tasks on a quad-core processor

Hybrid Approaches

For complex applications, combining thread and process pools can provide the best of both worlds:

```
from concurrent.futures import ThreadPoolExecutor, ProcessPoolExecutor
import time
import requests
import hashlib

def fetch_url(url):
    """Fetch content from a URL (I/O-bound)."""
    start = time.time()
    response = requests.get(url)
    return {
        'url': url,
        'content': response.content,
        'time': time.time() - start
    }

def process_content(content_dict):
    """Process the content with CPU-intensive operations."""
    url = content_dict['url']
    content = content_dict['content']
    start = time.time()

    # CPU-intensive operations
    result = hashlib.sha256(content).hexdigest()
    for i in range(1000000): # Simulate more CPU work
        result = hashlib.sha256(result.encode()).hexdigest()

    return {
        'url': url,
        'result': result[:10] + '...',
        'time': time.time() - start
    }

urls = [
    'https://www.python.org',
    'https://www.github.com',
    'https://www.stackoverflow.com',
    'https://www.wikipedia.org',
    'https://www.bbc.co.uk'
]

start = time.time()

# First stage: Use threads for I/O-bound URL fetching
with ThreadPoolExecutor(max_workers=5) as thread_executor:
    # Fetch all URLs concurrently
    content_dicts = list(thread_executor.map(fetch_url, urls))

# Second stage: Use processes for CPU-bound content processing
with ProcessPoolExecutor() as process_executor:
    # Process all content concurrently
    results = list(process_executor.map(process_content, content_dicts))

total_time = time.time() - start

# Print results
for result in results:
    print(f'{result["url"]}: {result["result"]} in {result["time"]:.2f}s')
print(f"Total time: {total_time:.2f}s")
```

This hybrid approach uses threads for I/O-bound tasks (URL fetching) and processes for CPU-bound tasks (content processing), maximizing efficiency by using the right tool for each part of the workload.



Hands-on Exercise: Comparing Thread and Process Pools

For this exercise, you'll implement the same image processing task using both `ThreadPoolExecutor` and `ProcessPoolExecutor` to compare their performance:

1. Navigate to the exercise directory: `cd python-optimization-workshop/exercises/executors`
2. The task involves processing a set of images with various transformations:
 - Load the image (I/O-bound)
 - Apply transformations (CPU-bound)
 - Save the result (I/O-bound)
3. Implement using `ThreadPoolExecutor`
4. Implement using `ProcessPoolExecutor`
5. Implement a hybrid approach using both
6. Run the benchmark script to compare performance: `python benchmark.py`

This exercise will demonstrate the practical differences between thread and process pools for mixed workloads that include both I/O-bound and CPU-bound components.

Course Summary and Next Steps

Today we've covered a comprehensive range of Python performance optimization and concurrency topics:

1

Performance Profiling

- Measuring code performance with cProfile, timeit, and line_profiler
- Identifying bottlenecks in Python applications
- Optimizing memory usage with memory_profiler
- Implementing performance improvements based on profiling data

2

Concurrency Models

- Understanding the differences between threads, processes, generators, and coroutines
- Working with Python's Global Interpreter Lock (GIL)
- Choosing the appropriate concurrency model for different workloads
- Implementing thread-safe code patterns

3

Event Loop Concurrency

- Building asynchronous applications with asyncio
- Using gevent for implicit cooperative multitasking
- Exploring high-performance alternatives like uvloop
- Implementing best practices for event-loop based concurrency

4

Parallel Execution

- Working with ThreadPoolExecutor for I/O-bound tasks
- Using ProcessPoolExecutor for CPU-bound tasks
- Implementing advanced patterns for task chaining and management
- Combining approaches for optimal performance in complex applications

For further learning, explore: NumPy and SciPy for numerical computing, Dask for scalable parallel computing, and Cython or Numba for performance-critical code.

Thank you for participating! Any questions before we conclude?