

# UNLOCK THE LOCK: MASTERING PYTHON MULTITHREADING DEADLOCKS BEFORE THEY MASTER YOU



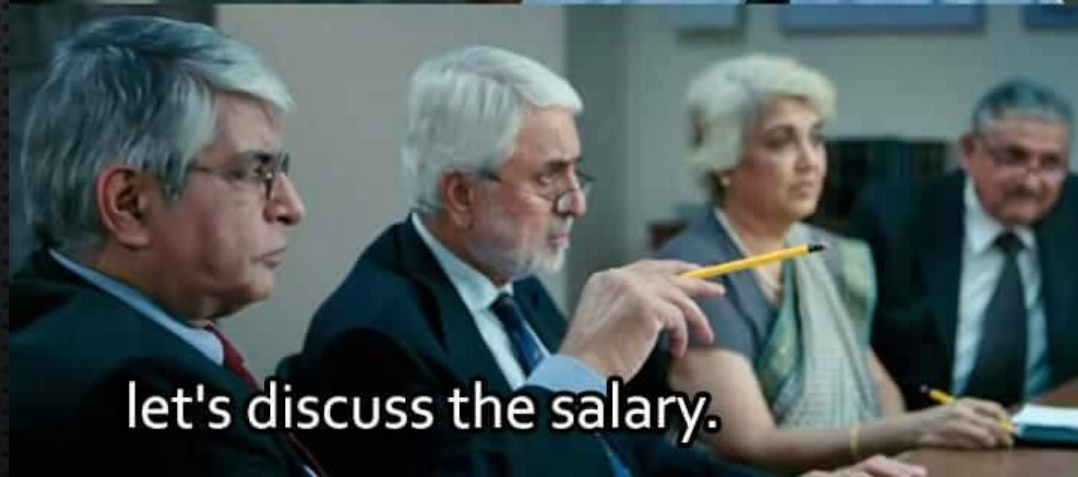


## About Me

- Principal Data Scientist @ Monarch Tractor Asia Pacific Pte Ltd
- Loves to learn through Communities (Slack, Discord)
- Teaches Bootcamps (Le Wagon)
- Writes
  - Medium: <https://hanqi01.medium.com/>
  - freeCodeCamp: <https://www.freecodecamp.org/news/the-logic-philosophy-and-science-of-software-testing-handbook-for-developers/>
- Plays Badminton, Piano, Drums, Chess, Sings
- Linkedin: <https://www.linkedin.com/in/hanqi91> (QR code)
  - Please give your 1<sup>st</sup> time speaker some feedback!
- Code: [https://github.com/gitgithan/pyconsg2025\\_multithreading](https://github.com/gitgithan/pyconsg2025_multithreading)
- Slides: Organizer will collect and share (also on github later)







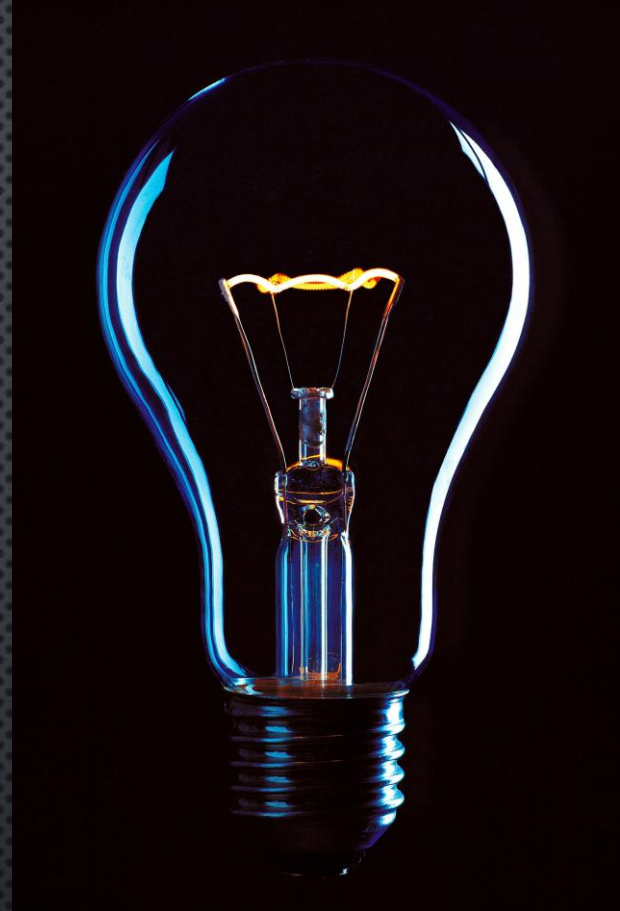


# Learning Objectives

- How to use VSCode Debugger
- Understand stepping behaviour
  - F10 (Step Over), F11 (Step In), Shift+F11 (Step Out), F5 (Continue)
- Understand Deadlock Spectrum
- See a Heisenbug
- Why use nested parallelism
- General Computer Science Concepts
  - Blocking
  - Decoupling with Queues
  - Synchronization with Locks
  - Deadlocks

## Pre-requisites

- MultiThreading and Pre-emptive multitasking
- ThreadPoolExecutor
- Blocking vs Non-blocking code
- Locks
- Attention to Detail
- Curiosity
- Bias to action
- Comfortable switching between certainty and uncertainty
- Your worst case: You're completely lost, but gain transferrable VSCode debugging skills
- My worst case: Code deadlocks when it shouldn't, wing it for the rest of the talk





# Talk Map

- PEP3148
- Deadlock Spectrum
- Script overview and Demo
- Deadlock Mechanics
- Debugging Techniques
- Deviation in Control Flow
- Influencing Script Termination
- Unswallow RuntimeError: The HeisenBug
- Overview of Solutions
- Nested Concurrency Patterns
- What else can you learn
- Bonus: 2 types of shutdown and Context Manager
- FAQ
- Resources



# PEP3148

peps.python.org/pep-3148/#threadpoolexecutor

- internal future methods
  - Module Functions
    - Check Prime Example
    - Web Crawl Example
  - Rationale
  - Reference Implementation
  - References
  - Copyright

[Page Source \(GitHub\)](#)

## [ThreadPoolExecutor](#)

The `ThreadPoolExecutor` class is an `Executor` subclass that uses a pool of threads to execute calls asynchronously.

Deadlock can occur when the callable associated with a `Future` waits on the results of another `Future`. For example:

```
import time
def wait_on_b():
    time.sleep(5)
    print(b.result()) # b will never complete because it is waiting on a.
    return 5

def wait_on_a():
    time.sleep(5)
    print(a.result()) # a will never complete because it is waiting on b.
    return 6

executor = ThreadPoolExecutor(max_workers=2)
a = executor.submit(wait_on_b)
b = executor.submit(wait_on_a)
```

And:

```
def wait_on_future():
    f = executor.submit(pow, 5, 2)
    # This will never complete because there is only one worker thread and
    # it is executing this function.
    print(f.result())

executor = ThreadPoolExecutor(max_workers=1)
executor.submit(wait_on_future)
```

Why this example?

- 5 lines, deceptively simple
- Uses `ThreadPoolExecutor` instead of lower level `threading.Thread`

Why would you ever do this?

- No reason, usually
- Awareness of debugger behaviours
- Improved thinking to write better code, debug future problems

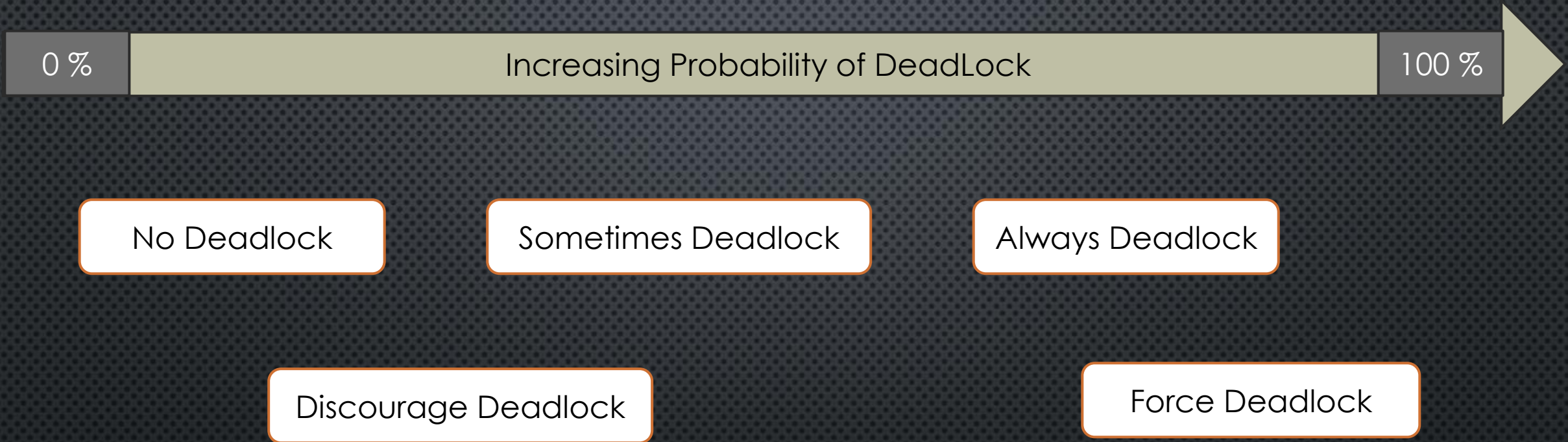
Pitfalls

- Possibly different behaviour when run in Jupyter Notebook

Whole talk using Python 3.10.15



# Deadlock Spectrum





## Script Overview and Demo

Module name	Behaviour
no_deadlock.py	No deadlock
sometimes_deadlock.py	More often <b>deadlock</b> than <b>no deadlock</b>
always_deadlock.py	Deadlock*
discourage_deadlock.py	More often <b>no deadlock</b> than <b>deadlock</b>
force_deadlock.py	Deadlock
unswallow_runtimeerror.py	No deadlock**

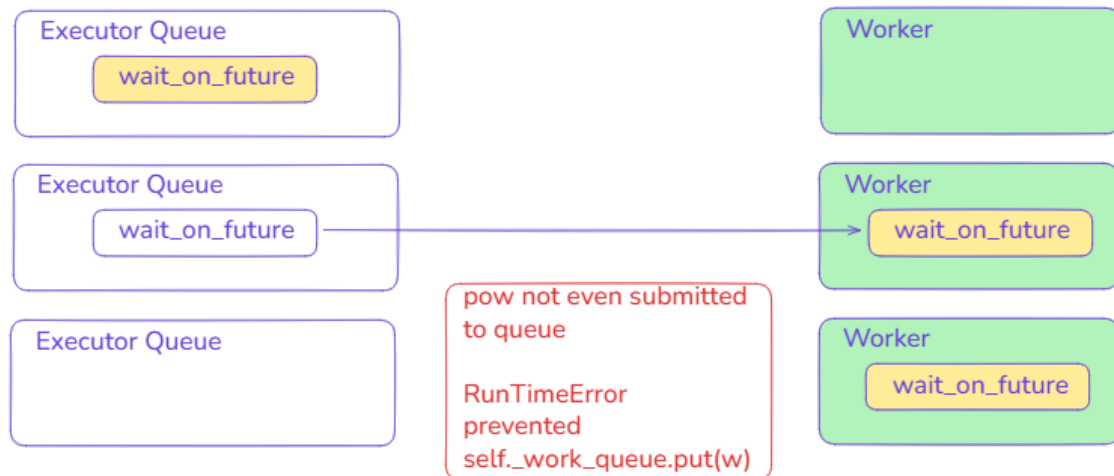
\* During debugging, ensure enough `time.sleep(seconds)`

\*\* RuntimeError visible on terminal

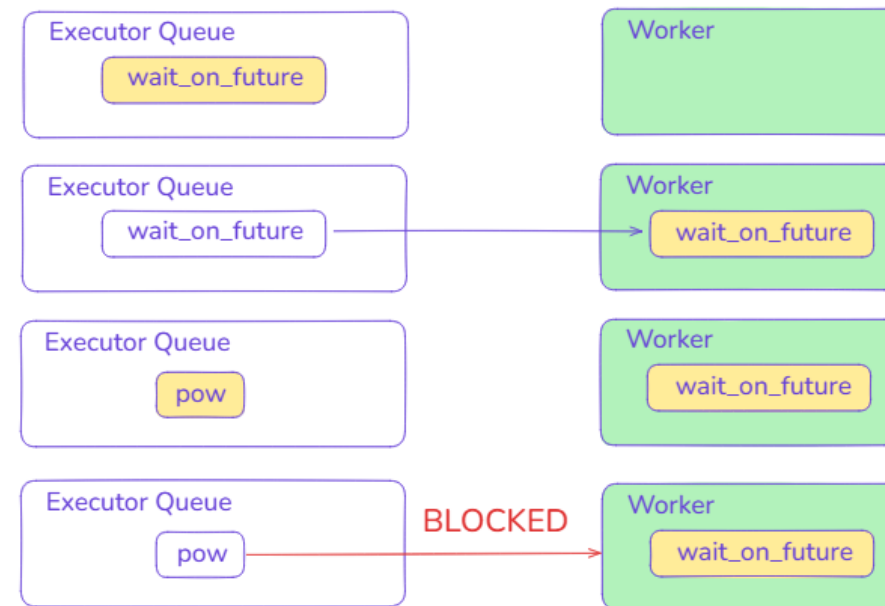


# Deadlock Mechanics

## No Deadlock



## Deadlock



```
no_deadlock.py > wait_on_future
1 from concurrent.futures import ThreadPoolExecutor
2 import concurrent.futures.thread as thread_module
3 import time
4
5
6 def wait_on_future():
7     f = executor.submit(pow, 5, 2)
8     # This will never complete because there is only one worker thread and
9     # it is executing this function.
10    print(f.result())
11
12
13 executor = ThreadPoolExecutor(max_workers=1)
14 executor.submit(wait_on_future)
```



# Debugging Techniques

- Watch window
  - Single Expression
  - Auto updates (Saves time from typing repeatedly in Debug Console)
  - Defined or not depends on current scope
- Call Stack
  - Helps track what has been executed, and where to put earlier breakpoints
  - Shows when and how many threads/processes started
  - Shows which thread is breakpoint paused at (PAUSED ON BREAKPOINT), and which thread is paused because other threads are paused (PAUSED)
- Debug Console
  - REPL (as many expression as you want)
- Variables
  - Investigate complex objects by clicking instead of typing in Debug Console
  - Choose important items to move to Watch window



# Deviation in Control Flow

## No Deadlock (no\_deadlock.py)

- `executor.submit(wait_on_future)` in `main scope`
- `_python_exit` in `thread.py`
  - Sets global variable `_shutdown = True`
  - Note that global is only global to specific module (`thread._shutdown`)
- `executor.submit(pow, 5, 2)` in `wait_on_future`
  - Reads `_shutdown` as `True`, **raise `RuntimeError`**



## Deadlock (always\_deadlock.py)

- `executor.submit(wait_on_future)` in `main scope`
- `time.sleep` in `main scope`
  - Global variable `_shutdown` is still `False`
- `executor.submit(pow, 5, 2)` in `wait_on_future`
  - Successfully submits `self._work_queue.put(w)`



# Influencing Script Termination

## **sometimes\_deadlock.py**

- Undeterministic behaviour

## **discourage\_deadlock.py**

- **executor.shutdown(wait=False)** increases chance of main script exiting first

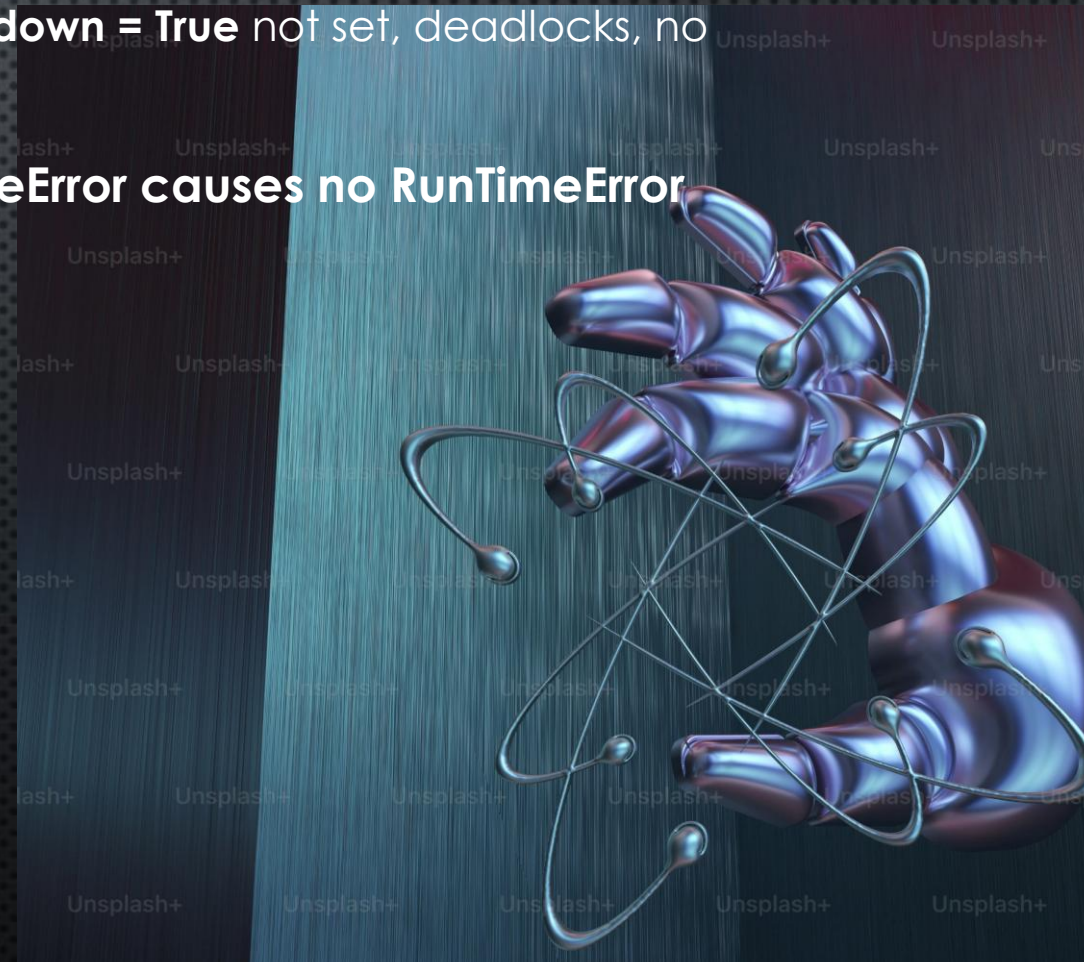
## **force\_deadlock.py**

- **f = executor.submit(wait\_on\_future)**
- **f.result()**
  - Blocking using condition variables, stronger because no time limit like `time.sleep`
    - `threading.Condition()` which use `RLock()` by default



# Unswallow RuntimeError: The HeisenBug

- Need no deadlock scenario to see RuntimeError (`_shutdown = True`)
- Need **`f.result()`** or **`f.exception()`** to see exceptions propagated through futures to the submitting scope
- **`f.result()`** blocks main from exit, `_python_exit` not run, **`shutdown = True`** not set, deadlocks, no RuntimeError
- **Heisenbug: Adding `f.result()` to observe the RuntimeError causes no RuntimeError**



# Overview of Solutions

```
solution1_max_workers.py > ...
1 from concurrent.futures import ThreadPoolExecutor
2 import concurrent.futures.thread as thread_module
3 import time
4
5
6 def wait_on_future():
7     print("Running wait_on_future")
8     f = executor.submit(pow, 5, 2)
9
10    print(f.result())
11
12
13 executor = ThreadPoolExecutor(max_workers=2) # with 2 workers, this will not deadlock
14 f = executor.submit(wait_on_future)
15 f.result()
```

```
solution2_processpool.py > wait_on_future
1 from concurrent.futures import ThreadPoolExecutor, ProcessPoolExecutor
2 import concurrent.futures.thread as thread_module
3 import time
4
5
6 def wait_on_future():
7     print("Running wait_on_future")
8     executor = ProcessPoolExecutor(max_workers=1)
9     f = executor.submit(pow, 5, 2)
10
11    print(f.result())
12
13
14 executor = ThreadPoolExecutor(max_workers=1)
15 f = executor.submit(wait_on_future)
16 f.result()
```

```
solution3_externalqueue.py > wait_on_future
1 from concurrent.futures import ThreadPoolExecutor
2 import queue
3 import time
4
5
6 def wait_on_future(result_queue):
7     f = executor.submit(pow, 5, 2)
8     result_queue.put(f) # Put the Future object in the queue
9
10
11 result_queue = queue.Queue()
12 with ThreadPoolExecutor(max_workers=1) as executor:
13     executor.submit(wait_on_future, result_queue)
14     # Retrieve result in the main thread
15     future = result_queue.get()
16     print(future.result())
```

```
solution4_timeout.py > ...
1 from concurrent.futures import ThreadPoolExecutor
2 import concurrent.futures.thread as thread_module
3 import time
4
5
6 def wait_on_future():
7     print("Running wait_on_future")
8     f = executor.submit(pow, 5, 2)
9
10    print(f.result())
11
12
13 executor = ThreadPoolExecutor(max_workers=1)
14 f = executor.submit(wait_on_future)
15 f.result(timeout=1)
```



# Overview of Solutions

Solution Filename	Key idea	Pros	Cons
solution1_max_workers.py	<p><b>Increase Worker Count:</b> Ensure enough threads are available for nested tasks.</p> <p><code>ThreadPoolExecutor(max_workers=2)</code></p>	<ul style="list-style-type: none"> <li>- Simple and effective</li> <li>- Allows concurrent execution of inner and outer tasks within the same pool.</li> </ul>	<ul style="list-style-type: none"> <li>- Increases resource consumption (more threads)</li> <li>- May not solve deadlocks if deeper nesting of tasks exceeds max_workers.</li> </ul>
solution2_processpool.py	<p><b>Separate Process Pool:</b> Offload the inner task to a <code>ProcessPoolExecutor</code> to avoid <code>ThreadPoolExecutor</code> deadlock</p> <p><code>ProcessPoolExecutor(max_workers=1)</code></p>	<ul style="list-style-type: none"> <li>- Bypasses <code>ThreadPoolExecutor</code> deadlock by using separate processes</li> <li>- Suitable for CPU-bound tasks (no GIL contention).</li> </ul>	<ul style="list-style-type: none"> <li>- Higher overhead due to process creation</li> <li>- More complex inter-process communication for data</li> <li>- Less suitable for I/O-bound tasks</li> <li>- Not ideal if shared memory access is required.</li> </ul>
solution3_externalqueue.py	<p><b>External Queue for Futures:</b> Put the Future object into a <code>queue.Queue</code> and retrieve it from an independent thread (e.g., main thread) for <code>result()</code></p> <p><code>result_queue.put(f)</code></p>	<ul style="list-style-type: none"> <li>- Breaks the direct blocking dependency on <code>f.result()</code> within the worker thread</li> <li>- Provides more control over result retrieval.</li> </ul>	<ul style="list-style-type: none"> <li>- Adds complexity due to the external queue</li> <li>- Requires explicit management of putting Future objects into and getting them from the queue</li> <li>- The main thread (or another thread) now bears the responsibility of awaiting the result.</li> </ul>
solution4_timeout.py	<p><b>Timeout:</b> Prevent indefinite blocking by raising a <code>TimeoutError</code> if the result is not available within a specified time</p> <p><code>print(f.result(timeout=1))</code></p>	<ul style="list-style-type: none"> <li>- Prevents indefinite blocking; raises <code>TimeoutError</code> instead</li> <li>- Useful for diagnosing deadlocks or unresponsive tasks.</li> </ul>	<ul style="list-style-type: none"> <li>- Does not solve the underlying deadlock; merely fails gracefully (or noisily)</li> <li>- The inner task still won't execute</li> <li>- Choosing an appropriate timeout value can be tricky and may lead to false positives or negatives</li> <li>- More of an error-handling mechanism than a true concurrency solution.</li> </ul>



# Nested Concurrency Patterns

Pattern	Outer Executor	Inner Executor	Use Case	Why Nested?
<b>Web Scraping</b>	ThreadPoolExecutor	ThreadPoolExecutor	Scrape multiple sites, then their subpages	Don't know subpage URLs upfront; need to group results per site
<b>Batch Processing</b>	ThreadPoolExecutor	ProcessPoolExecutor	Fetch data batches, then compute on each batch	Don't know compute jobs until parsing; need aggregation per batch
<b>Compute + Upload</b>	ProcessPoolExecutor	ThreadPoolExecutor	Heavy computation, then upload resulting files	Don't know file count upfront; need grouped upload handling
<b>ML Ensemble</b>	ProcessPoolExecutor	ProcessPoolExecutor	Generate bootstrap samples, then train models	Dynamic sample generation; need model aggregation

Reason for not flattening	Description	Critical Need
<b>Dynamic Task Discovery</b>	Tasks are only known after previous tasks complete	Can't pre-plan all work
<b>Result Grouping</b>	Need to aggregate/reduce results within logical groups	Maintains data relationships
<b>Resource Optimization</b>	Different task types need different executor types	Matches workload to resources
<b>Error Handling</b>	Group-level retry/failure policies	Logical error boundaries



# What else can you learn

## 1. Signaling — The Sentinel Pattern

- The use of None to tell threads to exit is a classic sentinel.
- Lesson: Simple values can be powerful signals in concurrent systems.
- Related patterns: poison pill, end-of-stream markers.

## 2. Queues — Decoupled, Asynchronous System Design

- Workers and producers communicate via a queue.
- Lesson: Queues enable asynchronous, decoupled, and safe task submission between threads.
- Foundation of many event-driven and pipeline architectures.

## 3. Resource Control — Bounded Concurrency

- max\_workers limits the number of concurrent threads.
- Lesson: Concurrency isn't about "as much as possible", but as much as needed.
- Teaches backpressure and resource safety.

## 4. Lifecycle Management — Ownership and Cleanup

- Explicit shutdown() vs. implicit \_python\_exit().
- Lesson: Resource lifecycles (like threads) require clear ownership and cleanup responsibility.
- Think: RAII, context managers, try/finally.

## 5. Task Abstraction — Futures and the submit() Interface

- Tasks submitted return Future objects, hiding execution details.
- Lesson: Future-based APIs allow for composability, polling, and callback chaining.
- Common in async frameworks, distributed systems, cloud functions.

## 6. Thread Safety via Synchronization — Locks and Atomicity

- Internals use threading.Lock, WeakKeyDictionary, etc.
- Lesson: Even with high-level abstractions, low-level race condition prevention is still required.
- Reinforces importance of safe data sharing.





# What else can you learn

## 7. Graceful Degradation — Ignore Failures at Shutdown

- `_python_exit()` doesn't retry, check state, or raise — it assumes shutdown is best-effort.
- Lesson: Graceful failure and best-effort cleanup are often more practical than perfection, especially at system exit.

## 8. Daemon vs. Non-Daemon Threads

- By default, threads are non-daemon, so they keep the process alive unless shut down.
- Lesson: Thread lifecycle directly affects process lifetime — daemon threads vs. non-daemon can make or break shutdown behavior.

## 9. Finalizers and `atexit` Hooks

- Uses `weakref.finalize` and `atexit.register` to manage global cleanup.
- Lesson: Be careful with global state and cleanup — it must not assume full interpreter health during teardown.

## 10. Blocking vs. Non-blocking APIs

- `submit()` is non-blocking; `shutdown(wait=True)` is blocking.
- Lesson: Understand and choose between immediate return and wait-for-completion based on context and backpressure needs.

## 11. Thread Leaks and Zombie Threads

- Failing to `shutdown()` or `join()` threads can lead to leaks or stalled interpreter exit.
- Lesson: Proper thread lifecycle hygiene is essential for long-running or daemonized applications.



## BONUS: 2 types of shutdown

- `_python_exit` setting `_shutdown`
- Heisenbug solution setting `executor.shutdown(wait=False)`
- Both handled by `executor.submit`
- Difference in scope (**python interpreter level** vs **executor level**)

```
118 class ThreadPoolExecutor(_base.Executor):
161     ... def submit(self, fn, /, *args, **kwargs):
● 162         ... with self._shutdown_lock, _global_shutdown_lock:
163             ... if self._broken:
164                 ... raise BrokenThreadPool(self._broken)
165         ↓
166             ... if self._shutdown:
● 167                 ... raise RuntimeError('cannot schedule new futures after shutdown')
168             ... if _shutdown:
● 169                 ... raise RuntimeError('cannot schedule new futures after '
170                 ... 'interpreter shutdown')
171         ↓
```

## BONUS: Context Manager

- concurrent.futures are usually used with context managers to automatically close the executor
- Similarly to the Heisenbug solution, both shutdown using **executor.shutdown**, albeit with differences in wait = True/False

context\_manager.py > ...

```
1 from concurrent.futures import ThreadPoolExecutor
2 import concurrent.futures.thread as thread_module
3 import time
4
5
6 def wait_on_future():
7     f = executor.submit(pow, 5, 2)
8     # This will never complete because there is only one worker thread and
9     # it is executing this function.
10    print(f.result())
11
12
13 with ThreadPoolExecutor(max_workers=1) as executor:
14     executor.submit(wait_on_future)
```

usr > lib > python3.10 > concurrent > futures > \_base.py > Executor > \_\_exit\_\_

```
571 class Executor(object):
572
573     def shutdown(self, wait=True, *, cancel_futures=False):
574         """Clean-up the resources associated with the Executor.
575
576         It is safe to call this method several times. Otherwise, no other
577         methods can be called after this one.
578
579         Args:
580             wait: If True then shutdown will not return until all running
581                   futures have finished executing and the resources used by the
582                   executor have been reclaimed.
583             cancel_futures: If True then shutdown will cancel all pending
584                             futures. Futures that are completed or running will not be
585                             cancelled.
586         """
587         pass
588
589     def __enter__(self):
590         return self
591
592     def __exit__(self, exc_type, exc_val, exc_tb):
593         self.shutdown(wait=True)
594         return False
```



# Frequently Asked Questions

## 1. How to know where to put the breakpoints?

- Put breakpoint early in execution flow, Step Over (F10) line by line on first few runs until familiar
- Later put breakpoints at critical sections you want to revisit and use Continue (F5) to make larger jumps

## 2. How to know whether there is deadlock?

- Simply observing that your script finishes execution doesn't guarantee the absence of deadlocks; you might just be "lucky" during that specific run.
- Conversely, experiencing a deadlock once doesn't mean it will *always* deadlock.
- **Halting Problem:** no tool can definitively guarantee the detection of all possible deadlocks in all possible programs *before* they actually occur or get stuck.
- If a program is still running or already stuck, we can investigate though



# Frequently Asked Questions

## 3. Are there tools that can help detect deadlocks?

- **faulthandler** (Python): This built-in Python module can dump tracebacks of all active threads when a fatal error occurs or after a specified timeout. This can reveal which threads are blocked and what they are waiting on if a deadlock causes the program to hang.
- OS-level tools (e.g., **strace** on Linux, **BPF Compiler Collection** tools): These system-level utilities can trace system calls, including lock acquisitions and releases. While powerful, interpreting their output for deadlock analysis can require a deeper understanding of operating system internals. They show you the resources threads are contending for at a given moment.
- **pydd** (Python Deadlock Detector): This library actively monitors lock operations in a running Python program. It builds a real-time dependency graph of lock acquisitions and can detect cycles in this graph, indicating an actual or imminent deadlock. It effectively tells you "a deadlock has occurred" or "threads are currently stuck."



# Resources

1. PEP3148: <https://peps.python.org/pep-3148/#threadpoolexecutor>
2. Introduction to PEP: <https://peps.python.org/pep-0001/>
3. How to use debugger: <https://www.youtube.com/watch?v=7qZBwhSlfOo>
4. VSCode python debugging: <https://code.visualstudio.com/docs/python/debugging>
5. Deterministic Logic vs Undeterministic Science: <https://www.freecodecamp.org/news/the-logic-philosophy-and-science-of-software-testing-handbook-for-developers/>
6. Synchronization with threading.Barrier: <https://www.linkedin.com/feed/update/urn:li:activity:7332378628250288128>
7. Avoiding multiprocessing race conditions: <https://www.linkedin.com/feed/update/urn:li:activity:7324384085206863872>
8. Multiprocessing and Pipes: <https://www.linkedin.com/feed/update/urn:li:activity:7323283169821061120>
9. LLMs to visualize relationship between synchronization primitives (Event, Condition, Lock, Semaphore, Barrier, Rlock, Timer): [https://www.linkedin.com/posts/hanqi91\\_is-claude-ready-for-education-yes-i-saw-activity-7313614310440587267--J33](https://www.linkedin.com/posts/hanqi91_is-claude-ready-for-education-yes-i-saw-activity-7313614310440587267--J33)
10. Problems in Concurrent algorithm design: [https://en.wikipedia.org/wiki/Dining\\_philosophers\\_problem](https://en.wikipedia.org/wiki/Dining_philosophers_problem)
11. Unintuitive blocking behaviour in nested context managers: <https://stackoverflow.com/questions/79538469/is-with-concurrent-futures-executor-blocking/79538630>
12. State of No GIL python: <https://pyfound.blogspot.com/2025/06/python-language-summit-2025-state-of-free-threaded-python.html>
13. BPF Compiler Collection (BCC): <https://github.com/iovisor/bcc>



## About Me

- Principal Data Scientist @ Monarch Tractor Asia Pacific Pte Ltd
- Loves to learn through Communities (Slack, Discord)
- Teaches Bootcamps (Le Wagon)
- Writes
  - Medium: <https://hanqi01.medium.com/>
  - freeCodeCamp: <https://www.freecodecamp.org/news/the-logic-philosophy-and-science-of-software-testing-handbook-for-developers/>
- Plays Badminton, Piano, Drums, Chess, Sings
- Linkedin: <https://www.linkedin.com/in/hanqi91> (QR code)
  - Please give your 1<sup>st</sup> time speaker some feedback!
- Code: [https://github.com/gitgithan/pyconsg2025\\_multithreading](https://github.com/gitgithan/pyconsg2025_multithreading)
- Slides: Organizer will collect and share (also on github later)

