# The GIL and API Performance

The Past, Present, and Free-Threaded Future

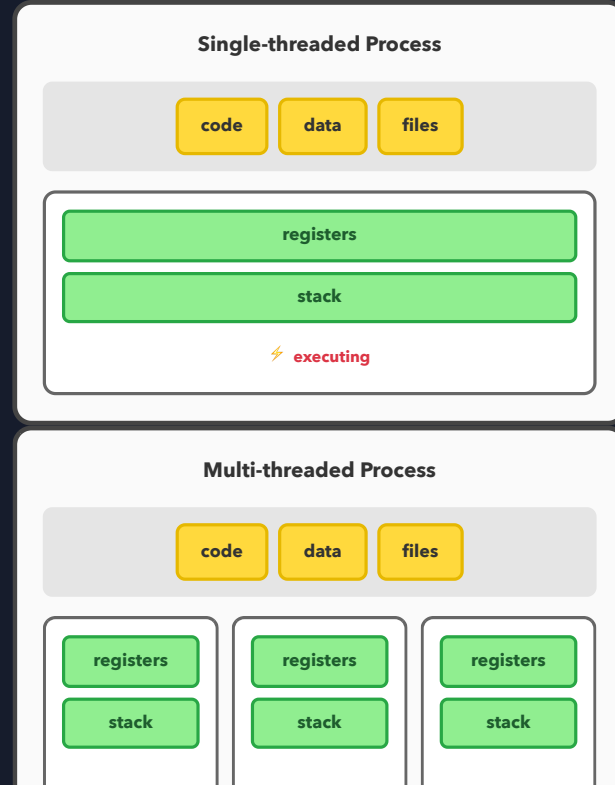# whoami

**Ruben Hias**

Software Engineer at TechWolf

**floww.dev**: Workflow automation for developers
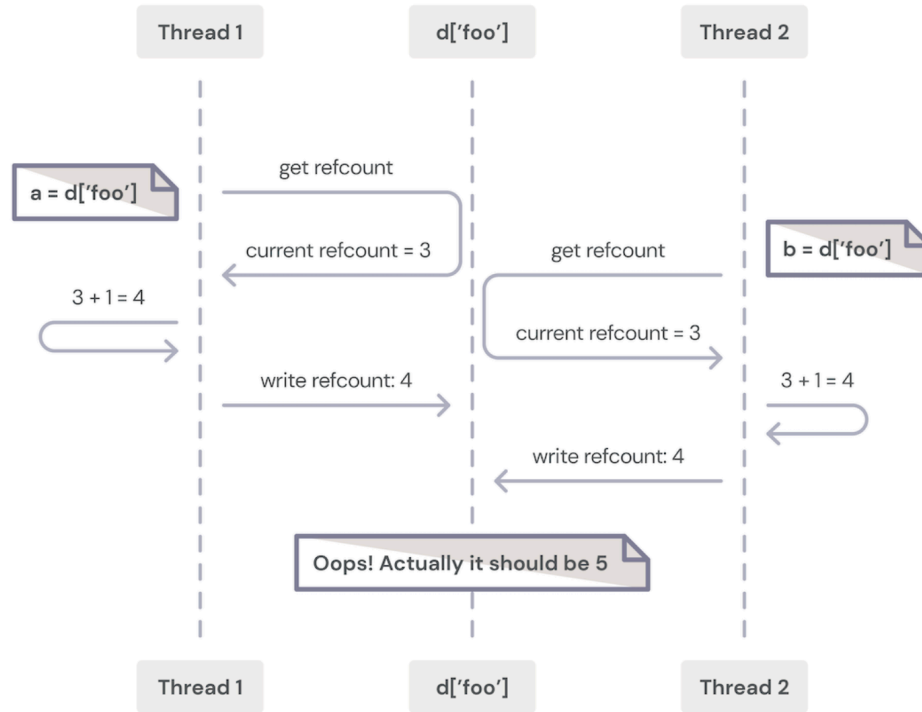
# Refresher: Threads vs Processes

- **Process**: Independent program with its own memory space
  - Isolated memory → safer but higher overhead
  - Communication via IPC
- **Thread**: Lightweight execution unit within a process
  - Shared memory space → efficient but requires synchronization
  - Direct communication via shared data
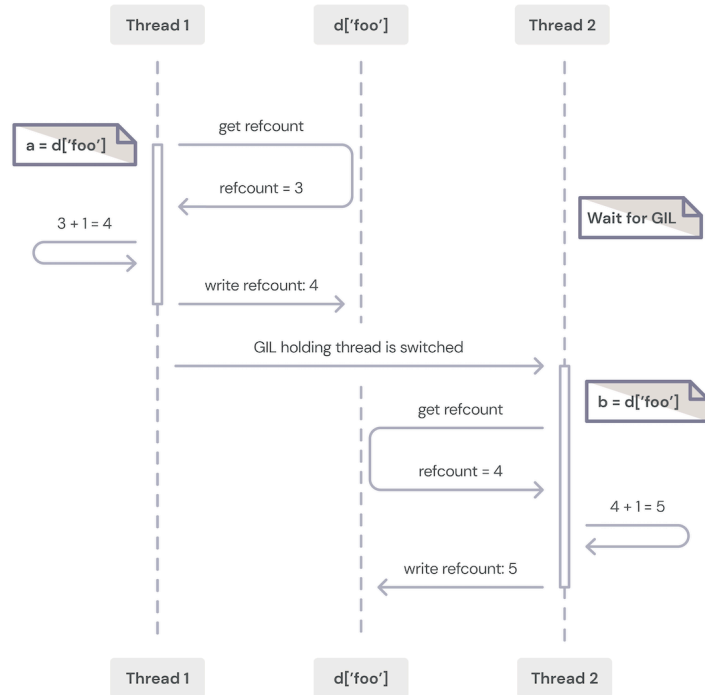
# GIL: Global Interpreter Lock

- *A lock that allows only **one thread** to execute Python bytecode at a time*
- CPython Implementation detail
- Protects internal data structures

GIL Serializing Reference Count Increment

# Commit `1984f1e`

gvanrossum committed on Aug 4, 1992

```
* Makefile adapted to changes below.

* split pythonmain.c in two: most stuff goes to pythonrun.c, in the library.
* new optional built-in threadmodule.c, build upon Sjoerd's thread.{c,h}.
* new module from Sjoerd: mmmodule.c (dynamically loaded).
* new module from Sjoerd: sv (svgen.py, svmodule.c.proto).
* new files thread.{c,h} (from Sjoerd).
* new xxmodule.c (example only).
* myselect.h: bzero -> memset
* select.c: bzero -> memset; removed global variable
```

main · v3.15.0a1 ··· 2.0

```c
Python/ceval.c

127    +
128    + void *
129    + save_thread()
130    + {
131    + #ifdef USE_THREAD
132    +     if (interpreter_lock) {
133    +         void *res;
134    +         res = (void *)current_frame;
135    +         current_frame = NULL;
136    +         release_lock(interpreter_lock);
137    +         return res;
138    +     }
139    +     else
140    +         return NULL;
141    + #endif
142    + }
143    +
144    + void
145    + restore_thread(x)
146    +     void *x;
147    + {
148    + #ifdef USE_THREAD
149    +     if (interpreter_lock) {
150    +         int err;
151    +         err = errno;
152    +         acquire_lock(interpreter_lock, 1);
153    +         errno = err;
154    +         current_frame = (frameobject *)x;
155    +     }
156    + #endif
157    + }
```
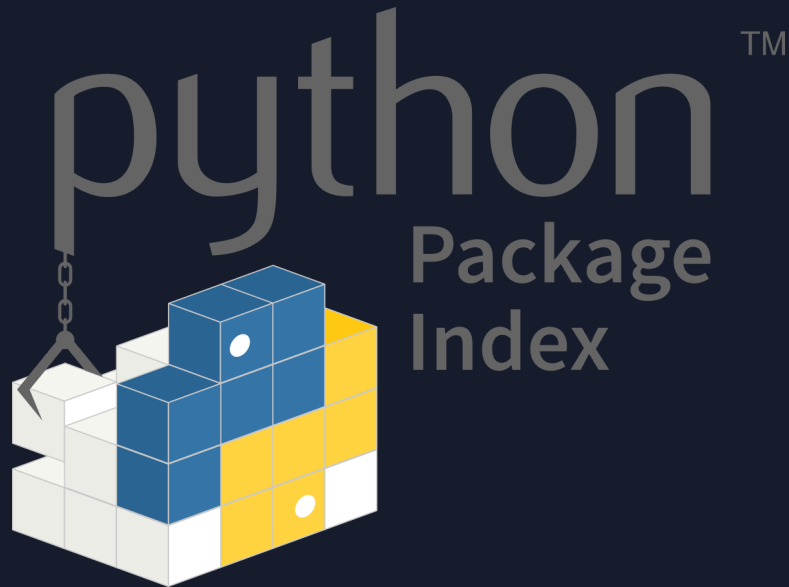
Why even have threads then?

# Why even have threads then?

## 1. I/O

- HDD Read speed: 0.5-2 MB/s
- Network Latency: 100's of ms
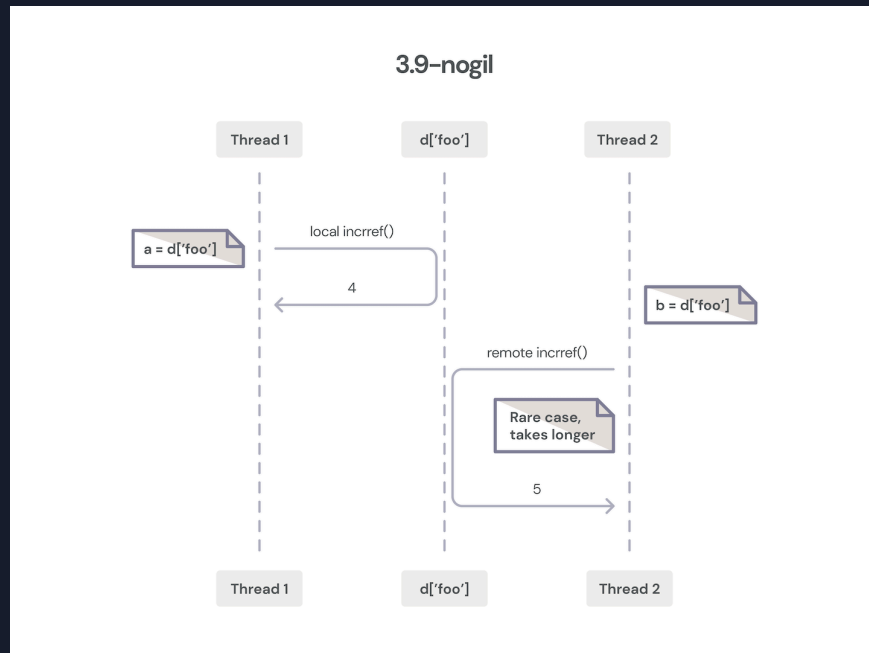- Network Bandwidth: 14.4 KB/s

# Why even have threads then?

# GIL removal attempts

- 1996
- 2007
- 2015 (x2)
- 2021

# Core innovations

- `nogil` ~ Sam Gross
- Biased reference counting & immortal objects
- Thread-safe memory allocator
- Many more small tweaks
- **Officially accepted as PEP**

# Recap: How we got here

- GIL is a CPython implementation detail to protect internal state
- Multiple removal attempts failed (1996, 2007, 2015) due to performance penalties
- Python 3.13 finally achieves free-threading via biased reference counting
- Trade-off: small single-threaded slowdown for true parallelism
- Now the question is: How do we use this?

```
uv run -p 3.14t main.py
```

# Basic multi-threading

```python
from concurrent.futures import ThreadPoolExecutor

N_THREADS = 8


def do_work(n: int = 10_000_000) -> int:
    return sum(i * i for i in range(n))


with ThreadPoolExecutor(max_workers=N_THREADS) as ex:
    futures = [
        ex.submit(do_work)
        for _ in range(N_THREADS)
    ]
```

**3.14**: 2.352 s ± 0.014 s

**3.14t**: 768.9 ms ± 11.7 ms

**~3x speedup**

# Single-threaded

```python
def do_work(n: int = 10_000_000) -> int:
    return sum(i * i for i in range(n))

for i in range(8):
    do_work()
```

**3.14**: 2.341 s ± 0.022 s

**3.14t**: 2.370 s ± 0.043 s

## 1.2% overhead

# Freethreading for CPU-bound applications

- True parallelism possible: 3x+ speedup for parallel CPU-bound tasks

- Best for embarrassingly parallel workloads (independent calculations)

- Single-threaded overhead now minimal (~3% vs Python 3.14)

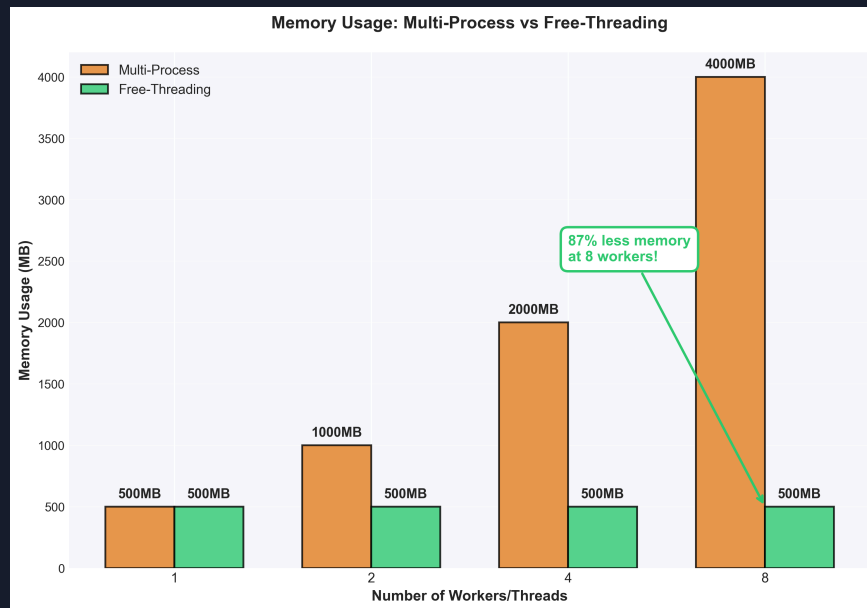- Thread overhead matters: use thread pools, avoid creating too many threads

For the API developers

# How Python APIs Scale Today

- Concurrent requests require parallel execution

- GIL blocks thread-based parallelism

- Solution: spawn multiple worker processes (Gunicorn, uWSGI)

- But processes have trade-offs…

# Issues with processes

- Processes don't share memory space

- Data sharing is hard
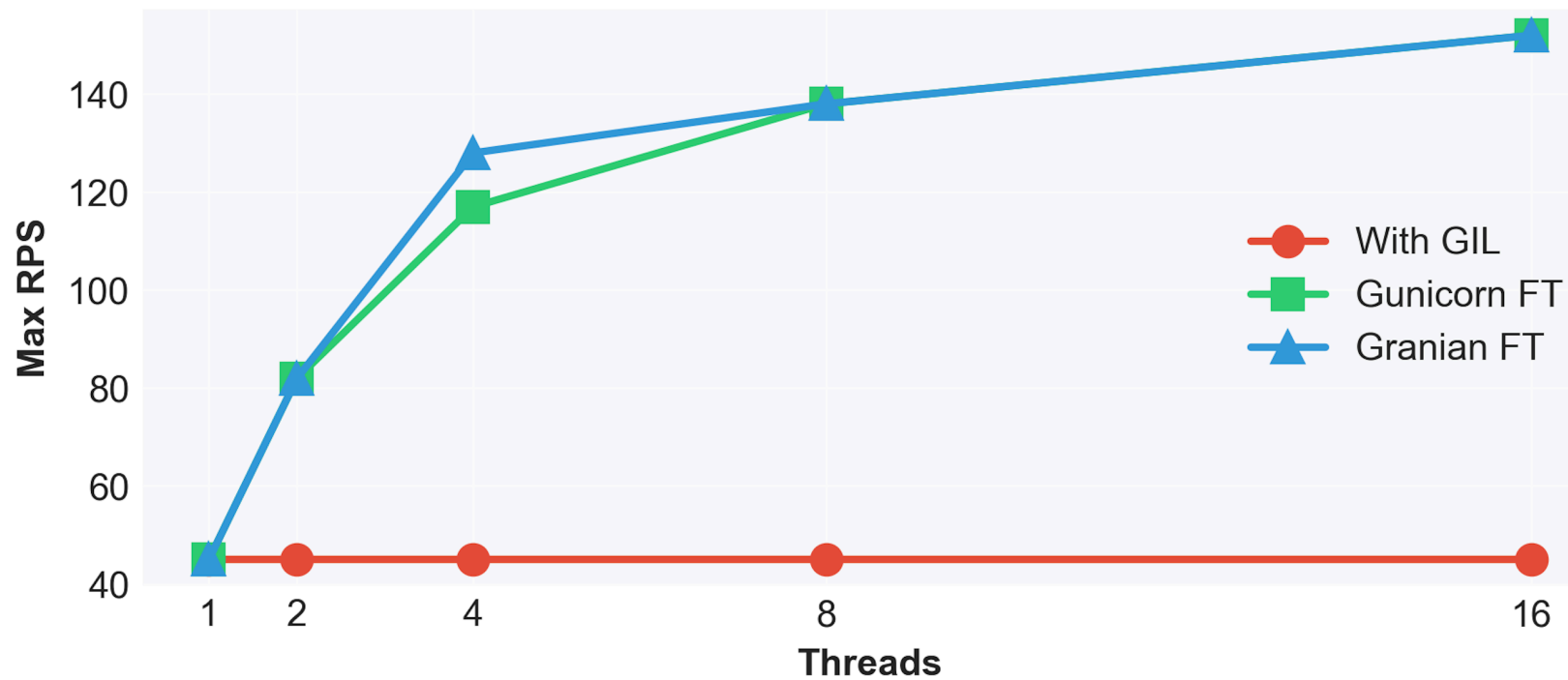
- Communication between processes is slow



Memory Usage: Multi-Process vs Free-Threading

BUT doesn't have GIL contention

# Benchmark Time

**Configurations tested:**

- Traditional: Gunicorn threads, Granian threads
- Async: Uvicorn/Granian ASGI, Gevent
- Free-threading: Python 3.14t with threads
- Baseline: Multi-process (current best practice)

**Question:** Can free-threading match multi-process performance?
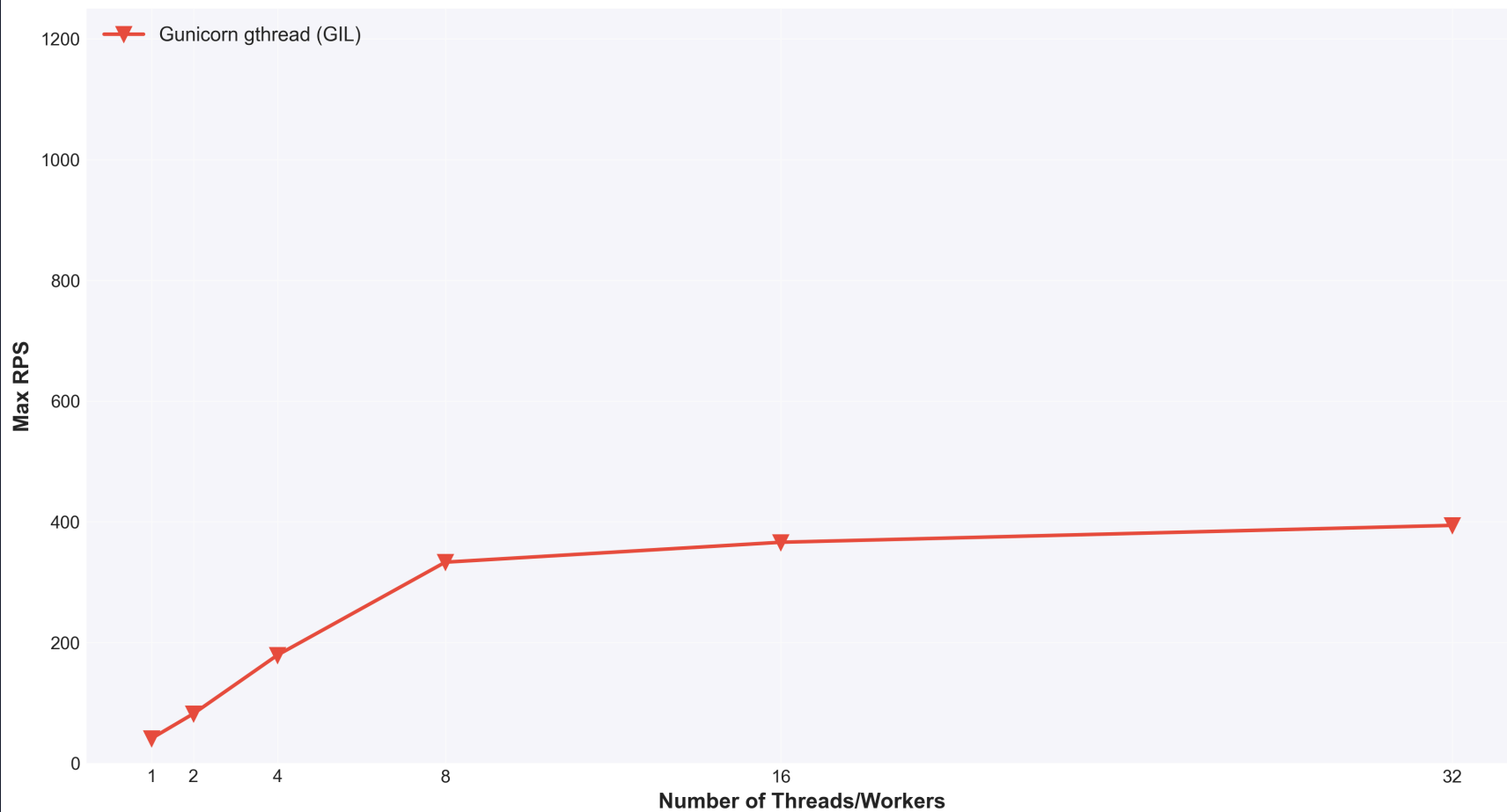
CPU-bound API Throughput Comparison
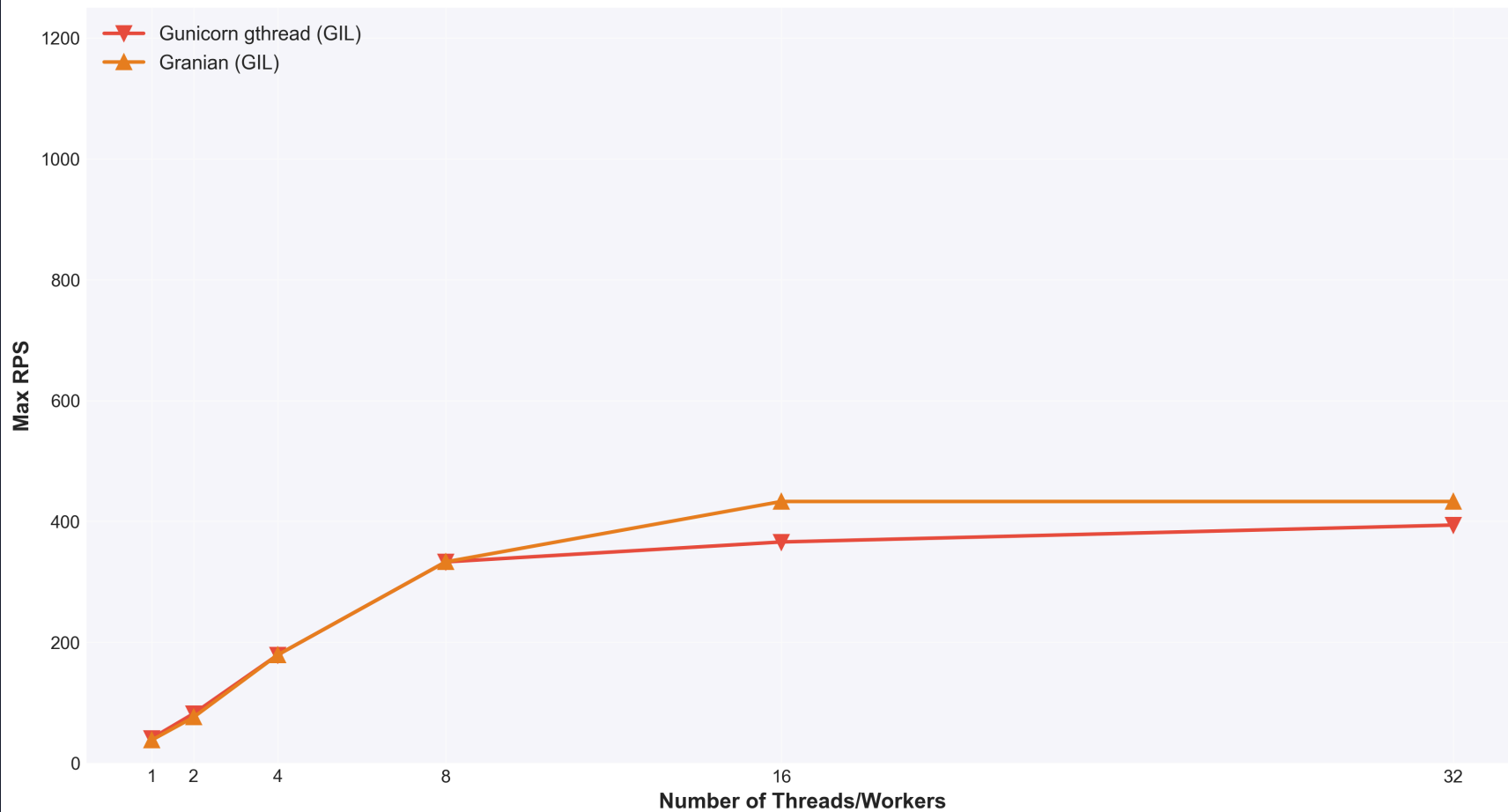
# Benchmark: Mixed Workload

Testing a realistic API: 80% I/O operations, 20% CPU work

- Comparing different scaling strategies
- Varies thread/worker count from 1 to 32
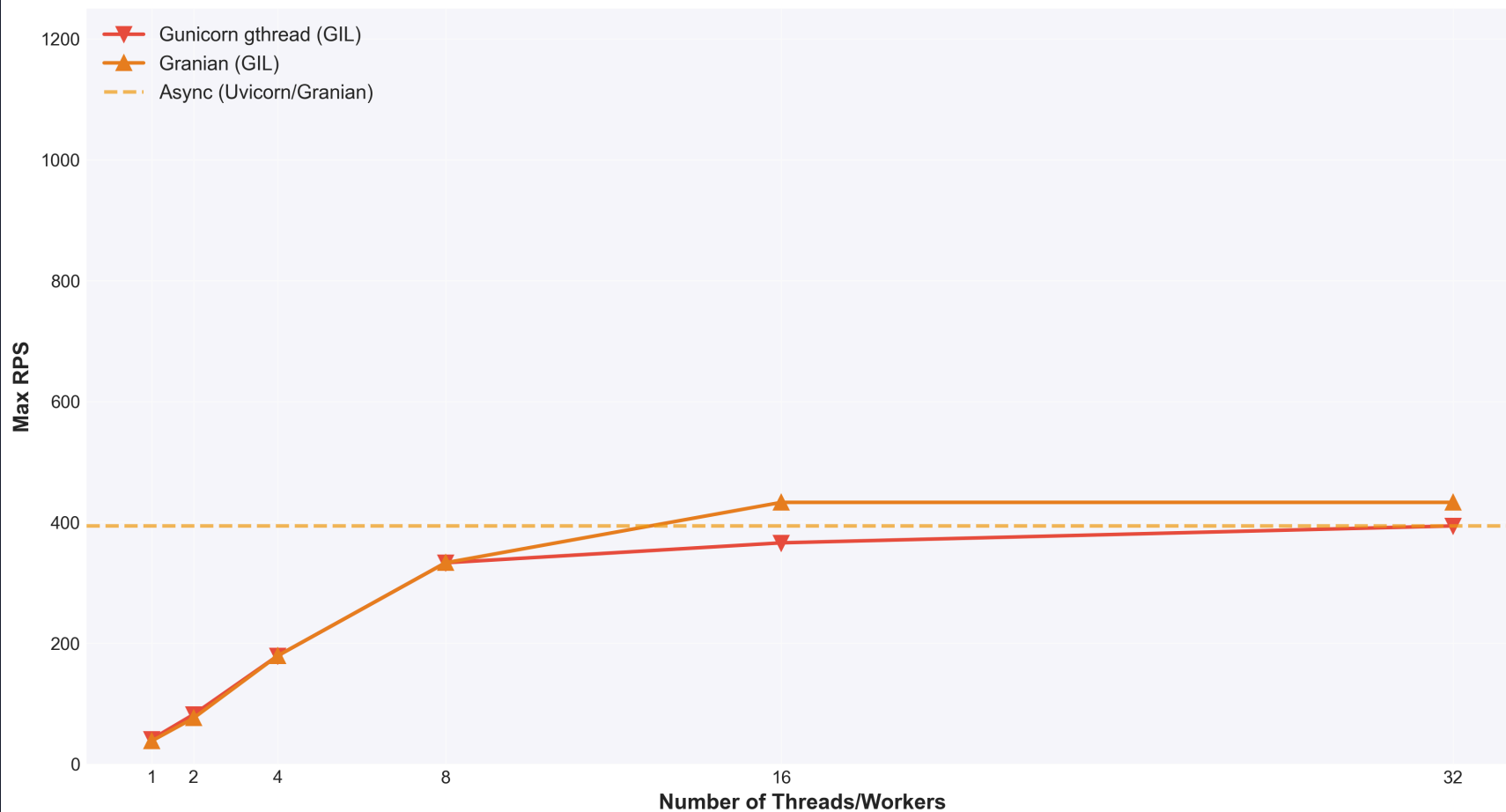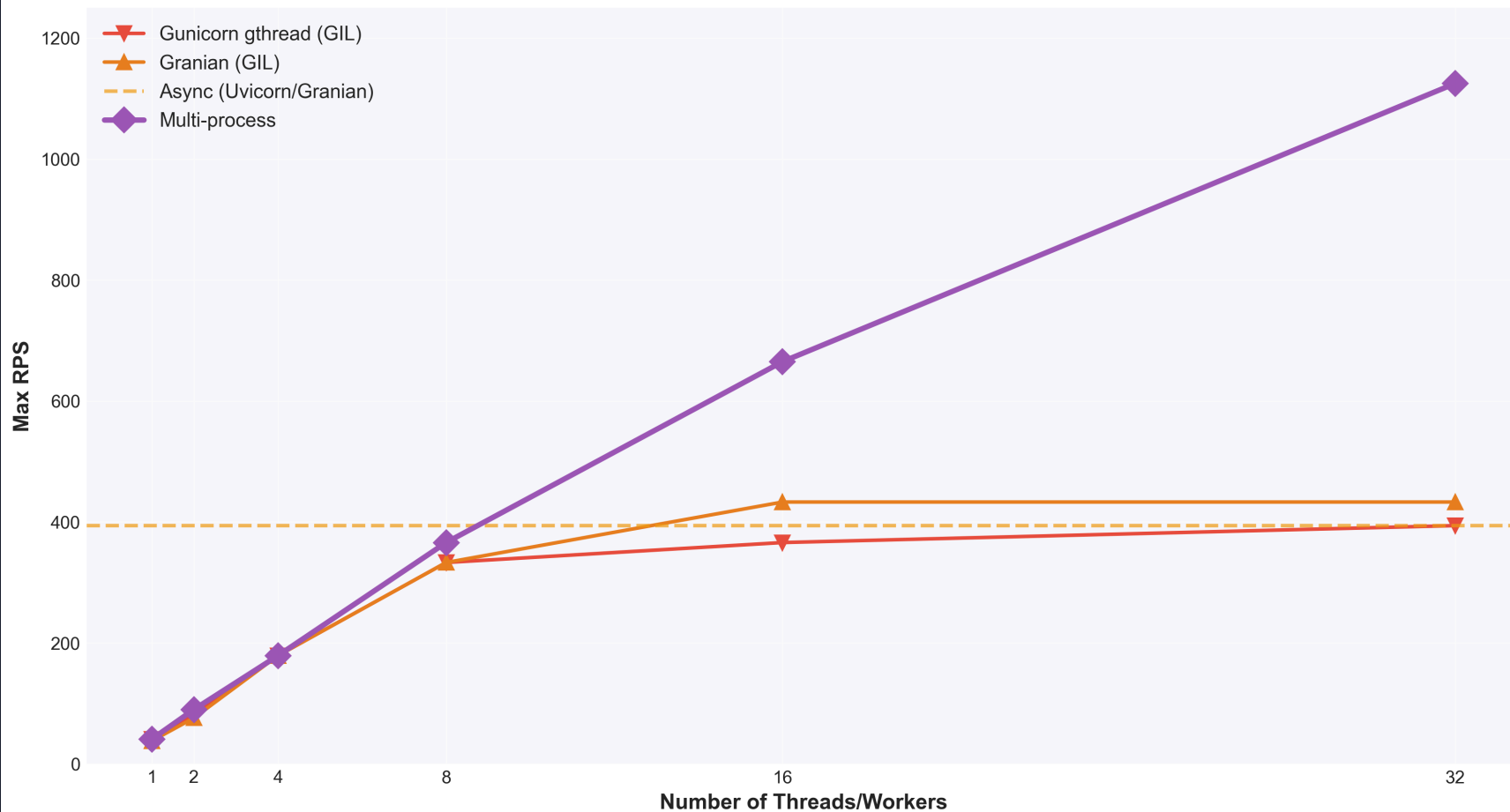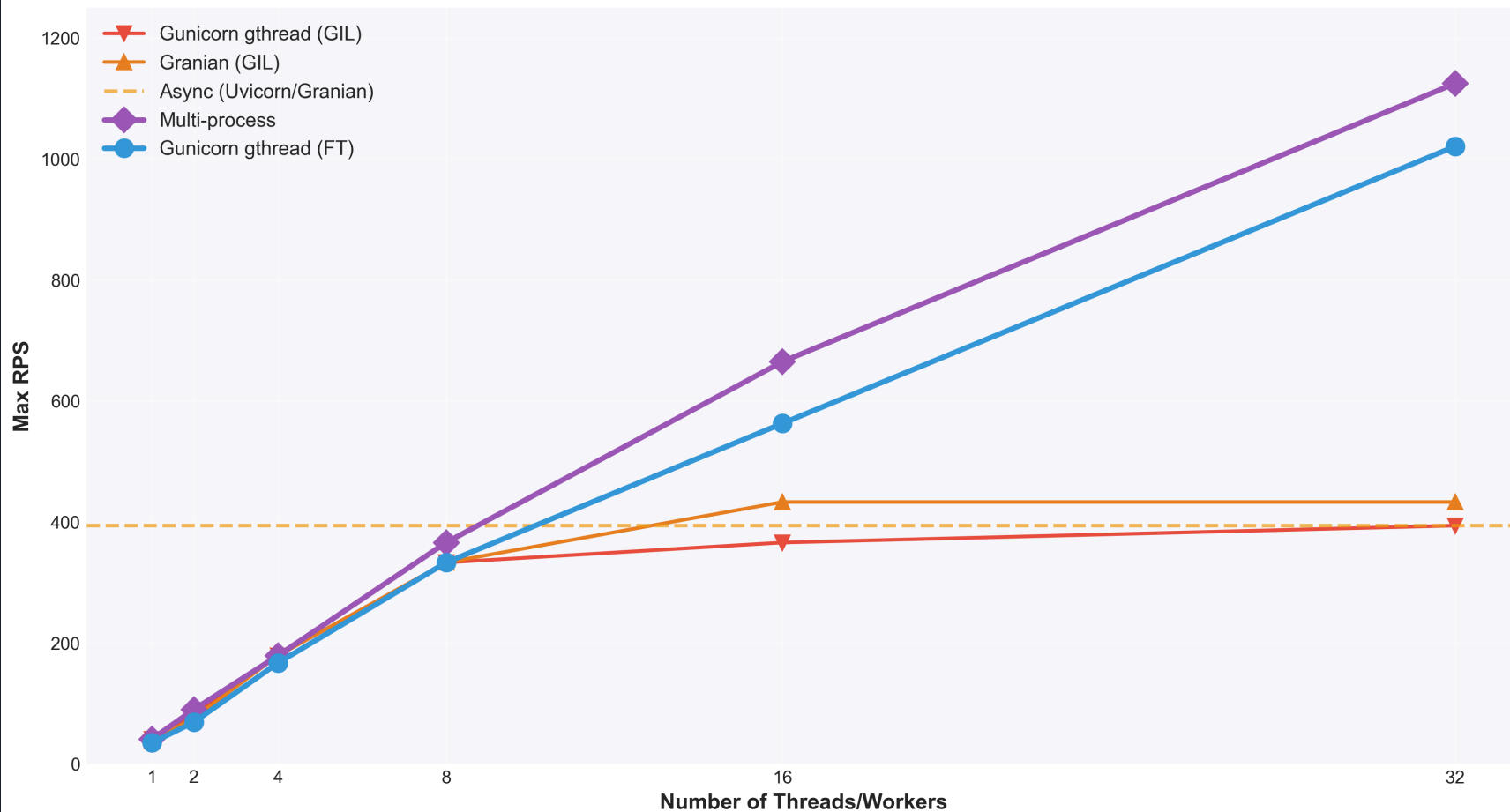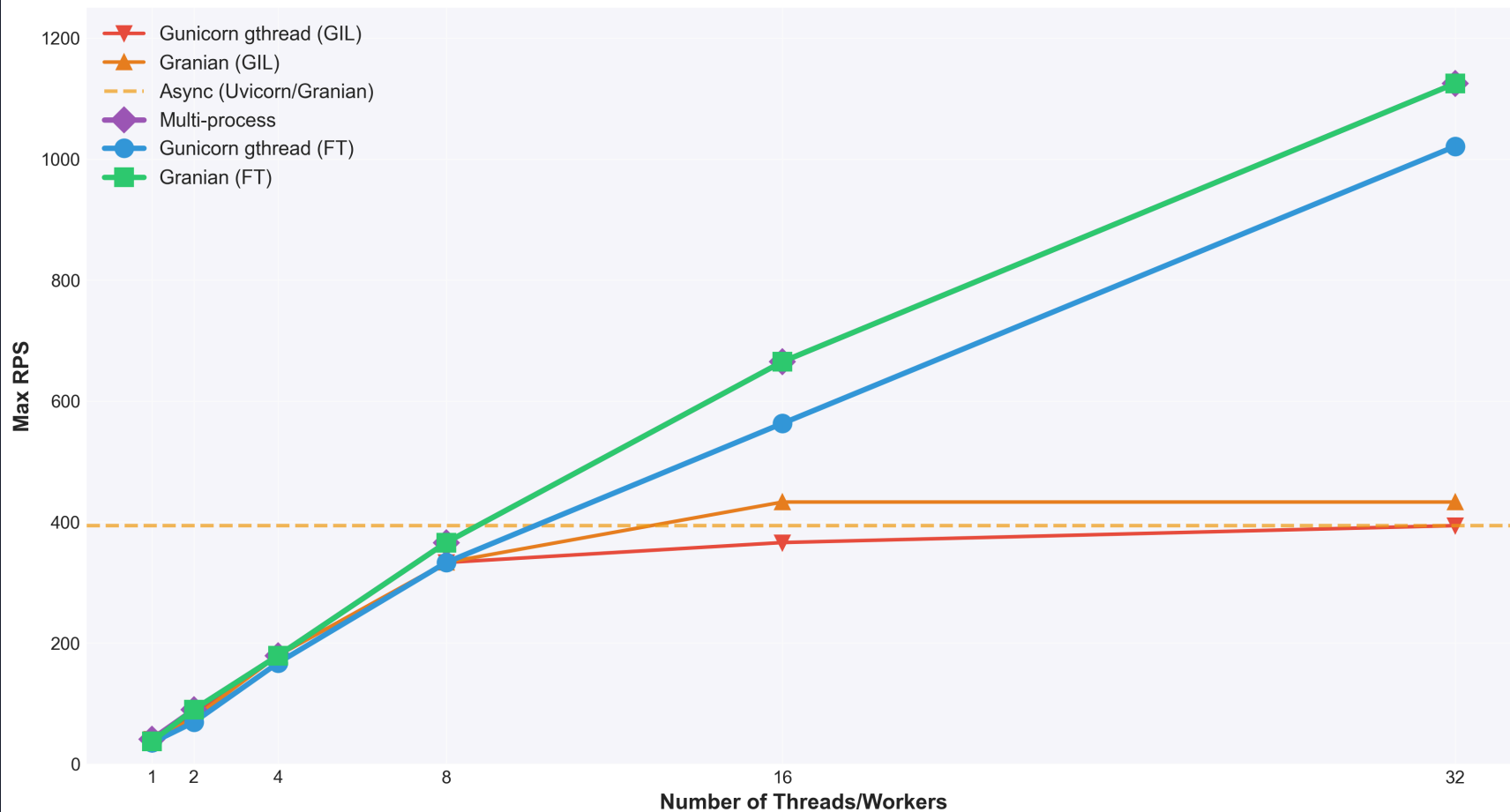- Measures maximum requests per second (RPS)

Mixed Workload Scaling (80% I/O, 20% CPU)

Mixed Workload Scaling (80% I/O, 20% CPU)

Mixed Workload Scaling (80% I/O, 20% CPU)

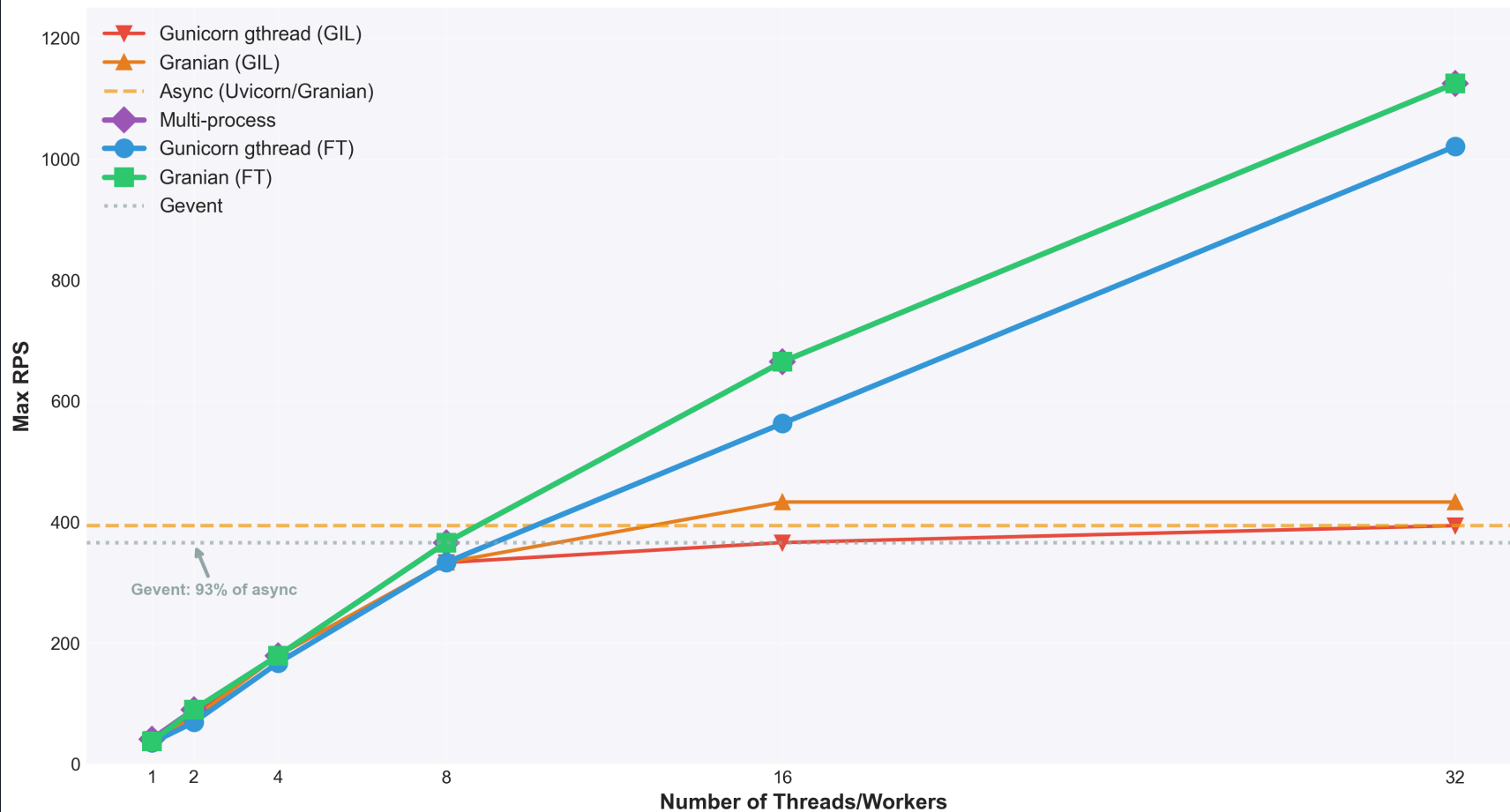Mixed Workload Scaling (80% I/O, 20% CPU)

Mixed Workload Scaling (80% I/O, 20% CPU)

Mixed Workload Scaling (80% I/O, 20% CPU)

Mixed Workload Scaling (80% I/O, 20% CPU)

# Key Takeaways

- GIL creates a scaling ceiling regardless of server choice

- Async helps with I/O concurrency but can't escape CPU-bound GIL limits

- Free-threading enables linear scaling: 2.5x improvement at 32 threads

- For I/O-heavy APIs: async/Gevent remain excellent choices today

- Modern webservers like Granian can better leverage threads

# Conclusions

- The GIL was the right trade-off in 1997, but hardware has changed
- Free-threading delivers true parallelism without memory overhead
- 2.5x scaling improvements for mixed workloads in production-ready Python
- For new projects: consider free-threading. For existing: async/multi-process work well

The future of Python performance is multi-threaded