

CRITICAL INTERVIEW CONCEPTS (Memorize These)

The GIL (Global Interpreter Lock) - PDF 3.1

MOST IMPORTANT: Only ONE thread executes Python bytecode at a time!

- **Threading HELPS:** I/O-bound work (network, disk, sleep)
 - GIL is released during I/O operations
 - Multiple threads can wait for I/O concurrently
- **Threading DOESN'T HELP:** CPU-bound work (math, parsing)
 - Use `multiprocessing` instead (separate processes)
- **For interview:** Assume I/O-bound request processing

Thread Count Decision (PDF 2.2.1)

Rule of Thumb:

- I/O-bound: `os.cpu_count() * 2` (workers wait on I/O)
- CPU-bound: `os.cpu_count()` (avoid context switching)
- **Trade-offs:**
 - More threads = better throughput BUT more memory/overhead
 - Each thread has stack space (~8MB on Linux)
 - Context switching overhead increases
- **Interview answer:** "Start with 4-8 for I/O-bound, then profile and tune based on latency/throughput metrics"

Single Queue vs Partitioned Queues (PDF 2.4)

Aspect	Single Queue	Partitioned
Ordering	No guarantee	Per-partition
Load balancing	Excellent	Poor (hot partition)
Complexity	Simple	More complex
Use case	Unordered tasks	User requests

Common Pitfalls (PDF 3.3) - MEMORIZE!

1. **No `join()`:** Threads keep running after main exits
2. **Shared state without locks:** Race conditions, data corruption
3. **Ignoring thread exceptions:** Silent failures (CRITICAL!)
4. **Infinite blocking:** Always use timeouts on `get()`
5. **Deadlocks:** Multiple locks acquired in different orders
6. **Increment race:** `counter += 1` is NOT atomic!

Key Queue Methods (PDF 3.2.3)

```

1 q = queue.Queue(maxsize=10) # Bounded queue
2 q.put(item, block=True, timeout=None) # Add item
3 item = q.get(block=True, timeout=None) # Remove item
4 q.task_done() # Signal one task complete
5 q.join() # Wait until all tasks done
6 q.qsize() # Approximate size (for monitoring)
7 q.empty() # Check if empty (not reliable!)
8 q.full() # Check if full (not reliable!)

```

Thread Safety Patterns

```

# Pattern 1: Lock for shared state (CRITICAL!)
lock = threading.Lock()
with lock: # Acquires and releases automatically
    shared_counter += 1 # This operation is NOT atomic!

# Pattern 2: Event for signaling
shutdown_event = threading.Event()
shutdown_event.set() # Signal all threads
shutdown_event.clear() # Reset flag
if shutdown_event.is_set(): # Check flag
    break # Worker exits

# Pattern 3: Queue for communication (already thread-safe!)
queue.Queue() # No lock needed - built-in synchronization!

# Pattern 4: RLock for reentrant locking
rlock = threading.RLock()
with rlock:
    with rlock: # Can acquire same lock again (Lock would deadlock!)
        shared_data.update()

# Pattern 5: Semaphore for resource pooling
sem = threading.Semaphore(3) # Allow 3 concurrent accesses
with sem:
    limited_resource.use() # Only 3 threads can be here

# Pattern 6: Condition for complex coordination
condition = threading.Condition()
# Producer:
with condition:
    produce_item()
    condition.notify() # Wake one waiting consumer
# Consumer:
with condition:
    while not item_available():
        condition.wait() # Release lock and wait for notification
    consume_item()

```

Why Operations Aren't Atomic

```

1 # counter += 1 is actually THREE operations:
2 # 1. Read counter value (e.g., 5)
3 # 2. Add 1 (5 + 1 = 6)
4 # 3. Write back (counter = 6)

```

COMPLETE PRODUCTION IMPLEMENTATION

```

1      """
2      INTERVIEW-READY REQUEST EXECUTOR
3      =====
4      Time to implement: 30-40 minutes
5      Lines of code: ~200-250
6
7      PDF References:
8      - Section 2.4: Partitioning strategy
9      - Section 3.2: Threading fundamentals
10     - Section 4.1: Complete template
11
12     ARCHITECTURE DECISION TREE:
13     -----
14
15     Q: Do requests need ordering?
16     YES -> Use partitioned queues (this implementation)
17     NO  -> Use single queue (simpler, better load balancing)
18
19     Q: How to handle backpressure?
20     - Fast-fail (block=False): Producer gets immediate feedback
21     - Block with timeout: Producer waits briefly
22     - Block forever: BAD! Cascading failures
23
24     Q: How many threads?
25     - I/O-bound: cpu_count() * 2
26     - CPU-bound: use multiprocessing instead!
27
28     import queue
29     import threading
30     import time
31     import os
32     from typing import Any, Dict, List, Optional
33     from dataclasses import dataclass
34     from collections import defaultdict
35
36     # =====
37     # DOMAIN MODEL
38     # =====
39
40     @dataclass
41     class Request:
42         """
43             Request with partition key for ordering.
44
45             INTERVIEW Q: Why partition_key?
46             A: Routes related requests (same user) to same worker
47                 for sequential processing. Hash determines worker.
48         """
49         request_id: str
50         partition_key: str # e.g., user_id
51         payload: Any
52         timestamp: float = None
53
54         def __post_init__(self):
55             if self.timestamp is None:
56                 self.timestamp = time.time()
57
58     # =====
59     # MAIN EXECUTOR
60     # =====
61
62     class PartitionedRequestExecutor:
63         """
64             Production-ready executor with:

```

```

1. Request partitioning (ordering per key)
2. Backpressure detection and fast-fail
3. Graceful shutdown
4. Comprehensive error handling
5. Thread-safe metrics
6. Worker failure detection

PDF Section 4.1, 2.4
"""

def __init__(self,
             num_workers: int = None,
             max_queue_size_per_worker: int = 100,
             process_func: Optional[callable] = None):
    """
INTERVIEW Q: Why default num_workers to None?
A: Calculate optimal thread count based on workload.
    For I/O-bound: cpu_count() * 2

PDF Section 2.2.1
"""
    if num_workers is None:
        num_workers = os.cpu_count() * 2
        print(f"Auto-detected {num_workers} workers "
              f"(CPU count: {os.cpu_count()})")

    self.num_workers = num_workers
    self.max_queue_size_per_worker = max_queue_size_per_worker
    self.process_func = process_func or self._default_process

    # One queue per worker (partition)
    # CRITICAL: Bounded queues provide backpressure
    self.queues = [
        queue.Queue(maxsize=max_queue_size_per_worker)
        for _ in range(num_workers)
    ]

    # Worker thread management
    self.workers: List[threading.Thread] = []
    self.shutdown_event = threading.Event()

    # Thread-safe metrics
    # INTERVIEW Q: Why need lock for metrics?
    # A: Dict updates not atomic. Even reads need sync.
    self.metrics_lock = threading.Lock()
    self.metrics = {
        'submitted': 0,
        'processed': 0,
        'failed': 0,
        'rejected': 0,
        'per_worker': defaultdict(lambda: {
            'processed': 0,
            'failed': 0,
            'last_active': None
        })
    }

    # Tracking for correctness testing (PDF 2.2.2)
    self.submitted_ids = set()
    self.processed_ids = set()
    self.failed_ids = set()

    self._start_workers()

def _start_workers(self):
    """
Start worker threads.

```

```

134
135     INTERVIEW Q: Daemon vs non-daemon?
136     A: daemon=False allows graceful shutdown with join().
137     Daemon threads killed abruptly on main exit.
138
139     PDF Section 3.2.1
140     """
141     for worker_id in range(self.num_workers):
142         t = threading.Thread(
143             target=self._worker_loop,
144             args=(worker_id,),
145             name=f"Worker-{worker_id}",
146             daemon=False # For graceful shutdown
147         )
148         t.start()
149         self.workers.append(t)
150     print(f"Started {self.num_workers} workers")
151
152     def _worker_loop(self, worker_id: int):
153         """
154             Main worker loop - CRITICAL PATTERN
155
156             MUST HAVE:
157             1. Check shutdown flag in loop
158             2. Timeout on get() to check shutdown periodically
159             3. Try/except around ALL processing
160             4. task_done() in finally block
161
162             INTERVIEW Q: What if exception in worker?
163             A: Doesn't propagate to main! Must catch here.
164
165             PDF Section 2.2.4, 3.2.3
166             """
167             my_queue = self.queues[worker_id]
168             print(f"Worker-{worker_id} started")
169
170             while not self.shutdown_event.is_set():
171                 try:
172                     # CRITICAL: Use timeout to check shutdown
173                     # NEVER: my_queue.get() (blocks forever!)
174                     request = my_queue.get(timeout=0.5)
175
176                     try:
177                         # Process request
178                         self.process_func(worker_id, request)
179
180                         # Success metrics (thread-safe)
181                         with self.metrics_lock:
182                             self.metrics['processed'] += 1
183                             self.metrics['per_worker'][worker_id]['processed'] += 1
184                             self.metrics['per_worker'][worker_id]['last_active'] = time.time()
185                             self.processed_ids.add(request.request_id)
186
187                         except Exception as e:
188                             # CRITICAL: Catch ALL exceptions
189                             # Don't let one bad request kill worker!
190                             print(f"[Worker-{worker_id}] ERROR: "
191                                   f"request {request.request_id}: {e}")
192
193                             with self.metrics_lock:
194                                 self.metrics['failed'] += 1
195                                 self.metrics['per_worker'][worker_id]['failed'] += 1
196                                 self.failed_ids.add(request.request_id)
197
198                         finally:
199                             # ALWAYS mark task done (even on error!)
200                             # Required for queue.join() to work
201                             my_queue.task_done()

```

```

202
203     except queue.Empty:
204         # No work, loop back to check shutdown
205         continue
206
207     print(f"Worker-{worker_id} shutdown complete")
208
209     def _default_process(self, worker_id: int, request: Request):
210         """Default processing (override with process_func)"""
211         print(f"[Worker-{worker_id}] Processing "
212               f"{request.request_id} (key: {request.partition_key})")
213         time.sleep(0.1) # Simulate I/O
214
215     def submit(self, request: Request, block: bool = False) -> bool:
216         """
217             Submit request with fast-fail backpressure.
218
219             INTERVIEW Q: Why block=False?
220             A: Fast-fail gives producer control to handle rejection
221                 (retry, return 503, etc.). Blocking causes pile-up.
222
223             INTERVIEW Q: How does partitioning work?
224             A: hash(partition_key) % num_workers
225                 Same key -> same hash -> same worker (ordering!)
226
227             PDF Section 2.4.3 (Consistent Hashing)
228             """
229             if self.shutdown_event.is_set():
230                 raise RuntimeError("Executor shutting down")
231
232             # PARTITIONING: Hash-based routing
233             # CRITICAL: Same key always goes to same worker
234             worker_id = hash(request.partition_key) % self.num_workers
235             target_queue = self.queues[worker_id]
236
237             try:
238                 target_queue.put(request, block=block)
239
240                 with self.metrics_lock:
241                     self.metrics['submitted'] += 1
242                     self.submitted_ids.add(request.request_id)
243
244                 return True
245
246             except queue.Full:
247                 # BACKPRESSURE: Reject when queue full
248                 print(f"[BACKPRESSURE] Rejected {request.request_id} "
249                       f"(Worker-{worker_id} queue full)")
250
251                 with self.metrics_lock:
252                     self.metrics['rejected'] += 1
253
254                 return False
255
256             def get_queue_depths(self) -> Dict[int, int]:
257                 """Monitor queue depth for backpressure detection"""
258                 return {i: self.queues[i].qsize()
259                        for i in range(self.num_workers)}
260
261             def is_backpressure(self, threshold: float = 0.8) -> bool:
262                 """
263                     Detect backpressure condition.
264
265                     INTERVIEW Q: How detect backpressure?
266                     A: Monitor queue depth. If > 80% full, system overwhelmed.
267
268                     INTERVIEW Q: Why check ANY worker, not average?
269                     A: Hot partition can bottleneck entire system.
270

```

```

271
272 PDF Section 2.2.3
273 """
274     for worker_id, depth in self.get_queue_depths().items():
275         if depth > (self.max_queue_size_per_worker * threshold):
276             return True
277     return False
278
279 def detect_dead_workers(self, timeout: float = 5.0) -> List[int]:
280     """
281     Detect workers that haven't processed in timeout seconds.
282
283     INTERVIEW Q: How detect dead workers?
284     A: Track last_active timestamp. If stale, worker may be stuck.
285
286     PDF Section 2.2.4
287     """
288     dead = []
289     now = time.time()
290
291     with self.metrics_lock:
292         for worker_id in range(self.num_workers):
293             last = self.metrics['per_worker'][worker_id]['last_active']
294             if last and (now - last) > timeout:
295                 dead.append(worker_id)
296
297     return dead
298
299 def verify_correctness(self) -> Dict[str, Any]:
300     """
301     Verify no dropped/duplicate requests.
302
303     INTERVIEW Q: How test correctness?
304     A: Invariant: submitted = processed + failed + rejected
305
306     PDF Section 2.2.2
307     """
308     with self.metrics_lock:
309         submitted = self.metrics['submitted']
310         processed = self.metrics['processed']
311         failed = self.metrics['failed']
312         rejected = self.metrics['rejected']
313
314         accounted_for = processed + failed + rejected
315
316         return {
317             'submitted': submitted,
318             'accounted_for': accounted_for,
319             'match': submitted == accounted_for,
320             'dropped': submitted - accounted_for,
321             'duplicate_processed': len(self.processed_ids) != processed,
322             'duplicate_failed': len(self.failed_ids) != failed
323         }
324
325 def get_metrics(self) -> Dict[str, Any]:
326     """Get all metrics (thread-safe)"""
327     with self.metrics_lock:
328         return {
329             **self.metrics,
330             'queue_depths': self.get_queue_depths(),
331             'backpressure': self.is_backpressure(),
332             'correctness': self.verify_correctness()
333         }
334
335 def shutdown(self, wait: bool = True, timeout: float = None):
336     """
337     Graceful shutdown sequence.
338
339     SHUTDOWN STEPS:
340     1. Stop accepting new requests
341     2. Wait for queues to drain (if wait=True)
342     3. Signal workers to stop (Event.set())
343     4. Wait for workers to exit (join with timeout)
344
345     INTERVIEW Q: Why Event instead of boolean?
346     A: Event is thread-safe. Boolean would need lock.
347     Event.set() and is_set() are atomic.
348
349     PDF Section 3.2.4
350     """
351     print("\n" + "="*50)
352     print("SHUTDOWN INITIATED")
353     print("="*50)
354
355     if wait:
356         print("Draining queues...")
357         for i, q in enumerate(self.queues):
358             q.join() # Wait for this queue to empty
359             print("All queues drained")
360
361     # Signal all workers
362     self.shutdown_event.set()
363
364     # Wait for workers
365     print("Waiting for workers...")
366     for worker in self.workers:
367         worker.join(timeout=timeout)
368
369     print("="*50)
370     print("SHUTDOWN COMPLETE")
371     print("="*50)
372     print("\nFinal Metrics:")
373     for key, val in self.get_metrics().items():
374         print(f" {key}: {val}")
375
376     # =====#
377     # ALTERNATIVE: SINGLE QUEUE IMPLEMENTATION
378     # =====#
379
380     class SingleQueueExecutor:
381         """
382         Simpler alternative: Single shared queue.
383
384         TRADE-OFFS vs Partitioned:
385         + Better load balancing (no hot partition problem)
386         + Simpler implementation
387         - No ordering guarantee
388         - All workers contend on one lock
389
390         USE WHEN: Tasks are independent, no ordering needed
391
392         PDF Section 2.4
393         """
394
395         def __init__(self, num_workers: int = None, max_queue_size: int = 1000):
396             if num_workers is None:
397                 num_workers = os.cpu_count() * 2
398
399             self.num_workers = num_workers
400             self.queue = queue.Queue(maxsize=max_queue_size) # Single queue!
401             self.workers = []
402             self.shutdown_event = threading.Event()
403
404             self._start_workers()
405
406             def _start_workers(self):
407                 for i in range(self.num_workers):
408                     t = threading.Thread(target=self._worker_loop, args=(i,))
409                     t.start()

```

```

409     self.workers.append(t)
410
411     def _worker_loop(self, worker_id: int):
412         while not self.shutdown_event.is_set():
413             try:
414                 # All workers pull from SAME queue
415                 # Whoever is free gets next task (good load balancing!)
416                 request = self.queue.get(timeout=0.5)
417                 self._process(worker_id, request)
418                 self.queue.task_done()
419             except queue.Empty:
420                 continue
421
422             def _process(self, worker_id: int, request: Request):
423                 print(f"Worker-{worker_id}: {request.request_id}")
424                 time.sleep(0.1)
425
426             def submit(self, request: Request, block: bool = False) -> bool:
427                 try:
428                     self.queue.put(request, block=block)
429                     return True
430                 except queue.Full:
431                     return False
432
433             def shutdown(self):
434                 self.queue.join()
435                 self.shutdown_event.set()
436                 for w in self.workers:
437                     w.join()
438
439 # =====
440 # TESTING & USAGE EXAMPLES
441 # =====
442
443     def test_correctness():
444         """
445             Test that no requests are dropped or duplicated.
446
447             INTERVIEW Q: How do we test this?
448             A: Submit N requests, verify N processed/failed/rejected.
449
450             PDF Section 2.2.2
451         """
452         print("\n" + "="*50)
453         print("CORRECTNESS TEST")
454         print("=".*50)
455
456         executor = PartitionedRequestExecutor(
457             num_workers=3,
458             max_queue_size_per_worker=5
459         )
460
461         # Submit 20 requests
462         num_requests = 20
463         for i in range(num_requests):
464             req = Request(
465                 request_id=f"req-{i}",
466                 partition_key=f"user-{i % 3}",
467                 payload={}
468             )
469             executor.submit(req, block=False)
470
471         # Wait and check
472         executor.shutdown(wait=True)
473
474         correctness = executor.verify_correctness()
475         print("\nCorrectness Check:")
476         print(f"  Submitted: {correctness['submitted']}")
477         print(f"  Accounted: {correctness['accounted_for']}")
478
479         print(f"  Match: {correctness['match']}")
480
481         if correctness['match']:
482             print("  PASS: No dropped requests!")
483         else:
484             print(f"  FAIL: {correctness['dropped']} requests dropped!")
485
486     def test_backpressure():
487         """
488             Test backpressure handling.
489
490             INTERVIEW Q: What happens when overwhelmed?
491             A: Bounded queue rejects, producer gets immediate feedback.
492         """
493         print("\n" + "="*50)
494         print("BACKPRESSURE TEST")
495         print("=".*50)
496
497         executor = PartitionedRequestExecutor(
498             num_workers=2,
499             max_queue_size_per_worker=3 # Small to trigger backpressure
500         )
501
502         # Submit many requests quickly
503         rejected = 0
504         for i in range(20):
505             req = Request(
506                 request_id=f"req-{i}",
507                 partition_key="user-1", # All to same worker!
508                 payload={}
509             )
510
511             if not executor.submit(req, block=False):
512                 rejected += 1
513
514         print(f"\nRejected due to backpressure: {rejected}")
515         executor.shutdown(wait=True)
516
517
518     def test_partitioning():
519         """
520             Test that same partition_key goes to same worker.
521
522             INTERVIEW Q: How verify partitioning works?
523             A: Track which worker processes each key. Same key should
524                 always go to same worker.
525         """
526         print("\n" + "="*50)
527         print("PARTITIONING TEST")
528         print("=".*50)
529
530         key_to_worker = {}
531
532         def track_worker(worker_id: int, request: Request):
533             key = request.partition_key
534             if key in key_to_worker:
535                 if key_to_worker[key] != worker_id:
536                     print(f"ERROR: {key} went to multiple workers!")
537             else:
538                 key_to_worker[key] = worker_id
539
540             print(f"Worker-{worker_id}: {request.request_id} "
541                  f"(key: {key})")
542             time.sleep(0.05)
543
544         executor = PartitionedRequestExecutor(
545             num_workers=3,
546             process_func=track_worker

```

```

547
548
549 # Submit requests for 3 users
550 users = ["alice", "bob", "charlie"]
551 for i in range(15):
552     req = Request(
553         request_id=f"req-{i}",
554         partition_key=users[i % 3],
555         payload={}
556     )
557     executor.submit(req)
558
559 executor.shutdown(wait=True)
560
561 print("\nPartitioning Results:")
562 for key, worker_id in key_to_worker.items():
563     print(f" {key} -> Worker-{worker_id}")
564
565
566 if __name__ == "__main__":
567     # Run all tests
568     test_correctness()
569     test_backpressure()
570     test_partitioning()
571
572     print("\n" + "="*50)
573     print("ALL TESTS COMPLETE")
574     print("="*50)

```

INTERVIEW WALKTHROUGH (30 MIN CODING)

Step 1: Basic Structure (5 min)

```

1 import queue, threading, time
2
3 class RequestExecutor:
4     def __init__(self, num_workers):
5         self.queues = [queue.Queue()
6                         for _ in range(num_workers)]
6         self.workers = []
7         self.shutdown = threading.Event()
8         self._start_workers()
9

```

Step 2: Worker Loop (10 min)

```

1 def _worker_loop(self, worker_id):
2     my_queue = self.queues[worker_id]
3
4     while not self.shutdown.is_set():
5         try:
6             req = my_queue.get(timeout=0.5)
7                 # Process request
8                 print(f"Worker-{worker_id}: {req}")
9                 my_queue.task_done()
10            except queue.Empty:
11                continue

```

Step 3: Submit with Partitioning (5 min)

```

1 def submit(self, request, block=False):
2     worker_id = hash(request.key) % len(self.queues)
3     try:
4         self.queues[worker_id].put(request, block=block)
5         return True
6     except queue.Full:
7         return False    # Backpressure

```

Step 4: Shutdown (5 min)

```

1 def shutdown(self, wait=True):
2     if wait:
3         for q in self.queues:
4             q.join()
5
6         self.shutdown.set()
7         for w in self.workers:
8             w.join()

```

Step 5: Add Metrics (5 min)

```

1 def __init__(self, ...):
2     # ... existing code ...
3     self.lock = threading.Lock()
4     self.processed = 0
5
6     def _worker_loop(self, worker_id):
7         # ... after processing ...

```

```
8     with self.lock:  
9         self.processed += 1
```

FOLLOW-UP QUESTIONS - PREPARE ANSWERS

Q1: Thread Count (PDF 2.2.1)

Q: How many threads should we use?

A: Start with `cpu_count() * 2` for I/O-bound. Trade-offs:

- More threads = higher throughput but more memory
- Each thread 8MB stack space
- Too many = context switching overhead
- Profile latency p50/p99 and tune

Q2: Testing Correctness (PDF 2.2.2)

Q: How ensure no dropped/duplicate requests?

A: Track invariants:

```
1 submitted_count == processed + failed + rejected  
2 len(processed_ids) == processed_count # No duplicates
```

Q3: Worker Failures (PDF 2.2.4)

Q: What if a worker thread dies?

A:

- Try/except in worker loop catches exceptions
- Track last_active timestamp per worker
- Monitor: if stale > 5s, worker may be stuck
- Recovery: restart worker thread or requeue task

Q4: Backpressure (PDF 2.2.3, 2.3.2)

Q: When signal backpressure to upstream?

A: Signal when:

- Queue depth > 80% capacity
- Request latency p99 > threshold (e.g., 5s)
- Worker saturation > 90%
- Monitor over time window (avoid false alarms)

Q5: Hot Partition (PDF 2.4.4)

Q: What about hot partition problem?

A:

- Trade-off: Ordering vs load balancing
- One user sends 1000 req = 1 worker overloaded
- Solutions:
 - Sub-partitioning (split hot keys)
 - Dynamic rebalancing (complex!)
 - Monitor per-worker queue depth
- Accept limitation for most use cases

Q6: Single Queue Alternative

Q: Why not use single queue?

A: Trade-offs:

- Single queue: Better load balancing, no ordering
- Partitioned: Ordering guarantee, hot partition issue
- Choice depends on requirements

KEY TALKING POINTS (MEMORIZIZE)

Why Threading for This Problem?

"This is an I/O-bound workload - processing requests involves network calls, database queries, or external API calls. During I/O operations, Python releases the GIL, allowing other threads to execute. This means we can handle multiple requests concurrently even with the GIL."

Why Partitioned Queues?

"We need to guarantee that requests from the same user are processed in order - for example, login before update profile. By hashing the `user_id` and routing to the same worker, we ensure sequential processing per user while still processing different users in parallel."

Why Fast-Fail Backpressure?

"When a queue is full, blocking would cause requests to pile up at the producer, potentially exhausting resources. Fast-fail with `block=False` gives the producer immediate feedback to handle it appropriately - retry with exponential backoff, return HTTP 503, or route to another instance."

Why Bounded Queues?

"Unbounded queues can lead to out-of-memory errors if requests arrive faster than we can process them. Bounded queues provide a natural backpressure mechanism - when full, we reject and let upstream handle it."

QUICK REFERENCE CARD

Essential Threading APIs

```

1 # Thread creation
2 t = threading.Thread(target=func, args=(a,))
3 t.start() # Begin execution
4 t.join() # Wait for completion
5
6 # Synchronization
7 lock = threading.Lock()
8 with lock: # Acquire/release
9     shared_state += 1
10
11 event = threading.Event()
12 event.set() # Signal
13 event.is_set() # Check
14
15 # Queue operations
16 q = queue.Queue(maxsize=10)
17 q.put(item, block=False) # Add
18 item = q.get(timeout=1) # Remove
19 q.task_done() # Mark complete
20 q.join() # Wait all done

```

Common Mistakes to Avoid

1. Forgetting `task_done()` after `get()`

2. Using `get()` without timeout (blocks forever!)

3. Shared state without locks (race conditions)

4. Assuming exceptions propagate (they don't!)

5. Using `daemon=True` when need graceful shutdown

6. Not checking `shutdown_event` in worker loop

Interview Time Management

- 0-5 min: Clarify requirements, draw architecture
- 5-20 min: Core implementation (basic executor)
- 20-30 min: Add partitioning + backpressure
- 30-40 min: Add error handling + metrics
- 40-50 min: Testing + discussion
- 50-60 min: Follow-up questions