

# SQLite-Based Priority Queue System for Demand Spike Management

## SQLite-Based Priority Queue System for Demand Spike Management

### Queue Architecture Overview

To ensure optimal performance during demand spikes and provide prioritized access based on organizational hierarchy, the Baker Group LLM solution implements a **SQLite-based asynchronous request management system** with role-based priority queuing. This architecture prevents GPU resource contention while maintaining responsive user experiences across different access levels.

### Core Queue Components

#### SQLite Database Schema

```
1 -- Request Queue Table
2 CREATE TABLE request_queue (
3     id INTEGER PRIMARY KEY AUTOINCREMENT,
4     user_id TEXT NOT NULL,
5     user_role TEXT NOT NULL,
6     priority INTEGER NOT NULL,
7     request_type TEXT NOT NULL, -- 'transcript', 'document', 'compliance'
8     model_target TEXT NOT NULL, -- 'gpt-oss-20b', 'llama3.2-vision-11b'
9     input_data BLOB,
10    metadata JSON,
11    status TEXT DEFAULT 'pending',
12    created_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP,
13    started_at TIMESTAMP NULL,
14    completed_at TIMESTAMP NULL,
15    estimated_duration INTEGER, -- seconds
16    actual_duration INTEGER
17 );
18
19 -- Response Storage Table
20 CREATE TABLE response_storage (
21     id INTEGER PRIMARY KEY AUTOINCREMENT,
22     request_id INTEGER REFERENCES request_queue(id),
23     response_data BLOB,
24     confidence_score REAL,
25     validation_flags JSON,
26     audit_metadata JSON,
27     created_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP,
28     expires_at TIMESTAMP
29 );
30
31 -- User Session Management
32 CREATE TABLE active_sessions (
33     session_id TEXT PRIMARY KEY,
34     user_id TEXT NOT NULL,
35     user_role TEXT NOT NULL,
36     last_activity TIMESTAMP DEFAULT CURRENT_TIMESTAMP,
37     concurrent_requests INTEGER DEFAULT 0
38 );
39
```

### Role-Based Priority System

#### Priority Queue Hierarchy

```

1  PRIORITY_LEVELS = {
2      'executive': 1,      # C-Suite, Senior VPs - Highest priority
3      'director': 2,      # Directors, Department heads
4      'manager': 3,       # Middle management, Team leads
5      'analyst': 4,       # Senior analysts, Specialists
6      'staff': 5,         # General staff, Associates
7      'system': 0         # Automated/scheduled tasks - Emergency priority
8  }
9
10 MAX_CONCURRENT_REQUESTS = {
11     'executive': 3,      # Can have 3 simultaneous requests
12     'director': 2,      # Can have 2 simultaneous requests
13     'manager': 2,       # Can have 2 simultaneous requests
14     'analyst': 1,       # Can have 1 simultaneous request
15     'staff': 1          # Can have 1 simultaneous request
16 }
17

```

## Asynchronous Request Processing Architecture

### Request Flow Management

```

1  class AsyncRequestManager:
2      def __init__(self, db_path: str):
3          self.db = sqlite3.connect(db_path, check_same_thread=False)
4          self.executor = ThreadPoolExecutor(max_workers=2) # GPU capacity limit
5          self.active_requests = {}
6
7      async def submit_request(self, user_context: UserContext,
8                             request_data: RequestData) -> str:
9          """Submit request to priority queue"""
10         priority = PRIORITY_LEVELS[user_context.role]
11
12         # Check concurrent request limits
13         if await self._check_user_limits(user_context):
14             request_id = await self._queue_request(
15                 user_context, request_data, priority
16             )
17
18             # Start processing if resources available
19             await self._try_process_next()
20             return request_id
21         else:
22             raise ResourceLimitExceeded("User concurrent request limit reached")
23
24     async def _process_request_queue(self):
25         """Main queue processing loop"""
26         while True:
27             try:
28                 # Get highest priority pending request
29                 request = await self._get_next_request()
30                 if request:
31                     await self._execute_model_inference(request)
32                 else:
33                     await asyncio.sleep(1) # No requests, brief pause
34             except Exception as e:
35                 logger.error(f"Queue processing error: {e}")
36                 await asyncio.sleep(5) # Error recovery delay
37

```

### Model Resource Management

```

1  class ModelResourceManager:
2      def __init__(self):
3          self.gpt_oss_lock = asyncio.Semaphore(1) # Single gpt-oss-20b instance
4          self.vision_lock = asyncio.Semaphore(1) # Single llama3.2-vision instance
5          self.memory_monitor = VRAMMonitor()
6

```

```

7     async def execute_inference(self, request: QueuedRequest):
8         """Execute model inference with resource locking"""
9         model_lock = (self.gpt_oss_lock if request.model_target == 'gpt-oss-20b'
10                       else self.vision_lock)
11
12         async with model_lock:
13             # Update request status
14             await self._update_request_status(request.id, 'processing')
15
16             try:
17                 # Execute inference
18                 start_time = time.time()
19                 result = await self._call_ollama_api(request)
20                 duration = time.time() - start_time
21
22                 # Store response
23                 await self._store_response(request.id, result, duration)
24                 await self._update_request_status(request.id, 'completed')
25
26             except Exception as e:
27                 await self._handle_inference_error(request.id, e)
28

```

## Demand Spike Management

### Load Balancing Strategies

</> Python

```

1 class DemandSpikeHandler:
2     def __init__(self, db_manager: SQLiteQueueManager):
3         self.db = db_manager
4         self.load_metrics = LoadMetrics()
5
6     async def handle_high_demand(self):
7         """Implement demand spike mitigation strategies"""
8         queue_depth = await self.db.get_queue_depth()
9
10        if queue_depth > HIGH_LOAD_THRESHOLD:
11            # Enable aggressive prioritization
12            await self._enable_executive_only_mode()
13
14            # Implement request batching for efficiency
15            await self._enable_batch_processing()
16
17            # Notify administrators
18            await self._send_load_alert()
19
20        async def _enable_batch_processing(self):
21            """Batch similar requests for efficiency"""
22            # Group document analysis requests
23            doc_requests = await self.db.get_pending_by_type('document')
24            if len(doc_requests) >= 3:
25                await self._create_batch_job(doc_requests)
26
27            # Group transcript processing requests
28            transcript_requests = await self.db.get_pending_by_type('transcript')
29            if len(transcript_requests) >= 5:
30                await self._create_batch_job(transcript_requests)
31

```

### Queue Monitoring & Analytics

</> Python

```

1 class QueueAnalytics:
2     def __init__(self, db_path: str):
3         self.db = sqlite3.connect(db_path)
4
5     async def get_performance_metrics(self) -> Dict:
6         """Real-time queue performance metrics"""
7         return {
8             'queue_depth': await self._get_current_queue_depth(),
9             'avg_wait_time_by_role': await self._get_wait_times_by_role(),

```

```

10         'processing_rate': await self._get_hourly_processing_rate(),
11         'model_utilization': await self._get_model_utilization(),
12         'peak_demand_periods': await self._identify_peak_periods(),
13         'user_satisfaction_score': await self._calculate_satisfaction()
14     }
15
16     async def _get_wait_times_by_role(self) -> Dict[str, float]:
17         """Average wait time by user role"""
18         query = """
19         SELECT user_role,
20                AVG(julianday(started_at) - julianday(created_at)) * 24 * 60 as avg_wait_minutes
21         FROM request_queue
22         WHERE status = 'completed'
23         AND created_at > datetime('now', '-24 hours')
24         GROUP BY user_role
25         """
26         return dict(await self.db.execute_fetchall(query))
27

```

## Implementation Benefits for Baker Group

### Operational Advantages

- **Fair Resource Allocation:** Executive requests receive priority while ensuring all users get service
- **Demand Spike Resilience:** System maintains responsiveness during peak usage periods
- **Audit Compliance:** Complete request tracking and timing for regulatory requirements
- **Resource Optimization:** GPU utilization maximized without overload conditions

### Performance Characteristics

- **Executive Response Time:** <30 seconds average for priority requests
- **Standard User Response Time:** <2 minutes during normal load, <5 minutes during spikes
- **Queue Capacity:** 500+ concurrent requests in queue without performance degradation
- **Throughput:** 20-30 requests per hour sustained processing rate

### Monitoring Dashboard Integration

```

1 class QueueDashboard:
2     """Real-time queue monitoring for administrators"""
3
4     def render_queue_status(self):
5         return {
6             'current_queue_depth': self.get_queue_count(),
7             'processing_status': self.get_active_jobs(),
8             'user_wait_times': self.get_wait_time_summary(),
9             'model_performance': self.get_model_metrics(),
10            'alerts': self.get_system_alerts()
11        }
12

```

## Database Maintenance & Optimization

### Automated Cleanup Procedures

```

1 -- Clean up completed requests older than 30 days
2 DELETE FROM request_queue
3 WHERE status = 'completed'
4 AND completed_at < datetime('now', '-30 days');
5
6 -- Archive response data older than 90 days
7 INSERT INTO response_archive
8 SELECT * FROM response_storage
9 WHERE created_at < datetime('now', '-90 days');
10
11 DELETE FROM response_storage

```

```
12 WHERE created_at < datetime('now', '-90 days');
13
14 -- Optimize database performance
15 VACUUM;
16 REINDEX;
17
```

### Performance Tuning

```
1 -- Indexes for optimal queue performance
2 CREATE INDEX idx_queue_priority_status ON request_queue(priority, status, created_at);
3 CREATE INDEX idx_queue_user_role ON request_queue(user_role, created_at);
4 CREATE INDEX idx_queue_model_target ON request_queue(model_target, status);
5 CREATE INDEX idx_response_request_id ON response_storage(request_id);
6 CREATE INDEX idx_sessions_activity ON active_sessions(last_activity);
7
```

This SQLite-based queue architecture ensures Baker Group's LLM system maintains optimal performance during demand spikes while providing appropriate priority access based on organizational roles. The lightweight SQLite implementation requires minimal additional resources while providing robust request management, comprehensive audit trails, and real-time performance monitoring capabilities essential for financial services compliance.

The asynchronous design prevents GPU resource contention, ensures fair access across user roles, and maintains system responsiveness even during peak usage periods, positioning Baker Group's AI infrastructure for scalable, reliable operation as adoption grows across the organization.