# Connection Pool Exhaustion

## 1. The Real Problem This Exists to Solve

Database connections are expensive resources (TCP connection, authentication, memory). Connection pools reuse connections across requests to avoid overhead. When all connections are in use, new requests must wait or fail. Connection pool exhaustion causes cascading failures, timeout errors, and complete service outages.

Real production scenario:

- API with 20-connection database pool
- Average request: 50ms database query
- Sudden traffic spike: 1000 requests/second
- **Without proper pool management:**
    - 20 connections immediately occupied
    - Request 21 waits for available connection
    - Request 22-1000 queue up (wait timeout: 30s)
    - Slow query appears (10 seconds)
    - 1 connection stuck for 10s
    - Only 19 connections available
    - Queue backs up further
    - Requests start timing out (30s)
    - Users see: "Database connection timeout"
    - 980 requests fail
    - Service effectively down

- **With proper pool management:**
    - Monitor pool utilization
    - Alert when >80% utilized
    - Implement connection timeout (5s)
    - Fail fast instead of queueing forever
    - Auto-scale pool size based on load
    - Release connections in finally blocks
    - Identify and kill slow queries
    - Circuit breaker prevents cascade
    - Service stays healthy

**The fundamental problem**: Database connections are a finite, shared resource. Without careful management, a single slow query or traffic spike can consume all connections, blocking all other requests and causing a complete outage. Connection pool exhaustion is a system-wide failure that cascades across all services.

Without pool management:

- Complete service outage
- Cascading failures
- Long timeout delays (30s+)
- Memory exhaustion
- Unrecoverable state

With proper management:

- Bounded wait times
- Fail fast
- Observable metrics
- Auto-recovery
- Graceful degradation

## 2. The Naive / Incorrect Approaches (IMPORTANT)

### ❌ Incorrect Approach #1: No Connection Pool (New Connection Per Request)

```javascript
// Incorrect: Create new connection for every request
app.get('/api/users', async (req, res) => {
  // Create new connection
  const client = new Client({
    host: 'localhost',
    database: 'myapp',
  });

  await client.connect();

  const result = await client.query('SELECT * FROM users');

  await client.end();

  res.json(result.rows);
});
```

**Why it seems reasonable:**

- Simple, no pool configuration
- Clean connection per request
- No connection leaks (closed every time)

**How it breaks:**

```
Connection overhead:
- TCP handshake: 50ms
- SSL handshake: 100ms
- Authentication: 50ms
- Total per request: 200ms

Load test results:
- 100 requests/second
- Total time: 200ms × 100 = 20 seconds of overhead
- Database: "Too many connections" error
- Max connections: 100 (PostgreSQL default)
- Request 101: CONNECTION REFUSED

Production:
- P50 latency: 200ms (just connection overhead)
```

```
- P99 latency: 5000ms (connection failures + retries)
- Throughput: 100 req/s max (connection limit)
```

**Production symptoms:**

```
Database logs:
ERROR: sorry, too many clients already

App logs:
Error: connect ECONNREFUSED
Error: Connection terminated unexpectedly

Monitoring:
- Open connections: Sawtooth pattern (spike to 100, drop to 0, repeat)
- Query time: 5ms
- Total request time: 250ms (50× overhead)
```

## ❌ Incorrect Approach #2: Infinite Pool Size (No Limit)

```javascript
// Incorrect: Pool with no max size
const pool = new Pool({
  host: 'localhost',
  database: 'myapp',
  // No max property! Will grow infinitely
});

app.get('/api/users', async (req, res) => {
  const client = await pool.connect();
  const result = await client.query('SELECT * FROM users');
  client.release();
  res.json(result.rows);
});
```

**Why it seems reasonable:**

- No connection limit
- Never blocks on pool.connect()
- Scales with demand

**How it breaks:**

```
Traffic spike:
- 10,000 concurrent requests
- Pool creates 10,000 connections
- Database max_connections: 100
- 9,900 connections rejected
- Database overloaded
- Memory exhaustion (each connection: 10MB)
- 10,000 × 10MB = 100GB memory
- OOM kill

Database side:
```

```
- 100 connections accepted
- 9,900 connections in SYN_RECV (pending)
- Database CPU: 100% (context switching)
- All queries slow (resource contention)
```

**Production symptoms:**

```
Database monitoring:
- Active connections: 100 (max)
- Waiting connections: 9,900
- Memory: 98% (OOM imminent)
- CPU: 100% (thrashing)

App monitoring:
- Connection errors: 9,900/10,000
- Memory: 50GB (connection overhead)
- ERROR: "too many clients already"
```

## ✖ Incorrect Approach #3: Never Releasing Connections

```javascript
// Incorrect: Forget to release connection
const pool = new Pool({ max: 20 });

app.get('/api/users', async (req, res) => {
  const client = await pool.connect();

  const result = await client.query('SELECT * FROM users');

  // Forgot to call client.release()!

  res.json(result.rows);
});
```

**Why it seems reasonable:**

- Simpler code (no finally block)
- Query works fine
- Pool should handle cleanup (wrong!)

**How it breaks:**

```
Request timeline:
Request 1: Acquires connection 1, completes, never releases
Request 2: Acquires connection 2, completes, never releases
...
Request 20: Acquires connection 20, completes, never releases
Request 21: pool.connect() → WAITS FOREVER

Pool state:
- Total connections: 20
- Available: 0
- In use: 20 (all leaked)
```

```
- Waiting: 1000+

Result:
- First 20 requests: Work fine
- All subsequent requests: Timeout or hang
- Service completely broken
- Requires restart to recover
```

**Production symptoms:**

```
Logs:
TimeoutError: ResourceRequest timed out

Monitoring:
- Pool size: 20 (max)
- Idle connections: 0
- Active queries: 0 (queries finished, connections not released!)
- Waiting clients: 5,000+

Manual inspection:
SELECT count(*) FROM pg_stat_activity WHERE state = 'idle';
Result: 20 (all connections idle but not returned to pool)
```

## ❌ Incorrect Approach #4: Long-Running Queries Hold Connections

```javascript
// Incorrect: Slow query blocks connection for minutes
const pool = new Pool({ max: 20 });

app.get('/api/analytics', async (req, res) => {
  const client = await pool.connect();

  try {
    // Slow query: 5 minutes
    const result = await client.query(`
      SELECT
        date_trunc('day', created_at) as day,
        count(*) as total,
        sum(amount) as revenue
      FROM orders
      WHERE created_at > NOW() - INTERVAL '2 years'
      GROUP BY day
      ORDER BY day
    `);

    res.json(result.rows);
  } finally {
    client.release();
  }
});
```

**Why it seems reasonable:**

- Connection properly released
- Analytics query needs all data
- finally block ensures cleanup

**How it breaks:**

```
Timeline:
T0: Analytics request arrives, acquires connection 1
T1: Slow query starts (5 minutes)
T5: Connection 1 still in use
T10: 10 requests arrive, use connections 2-11
T60: Connection 1 still in use
T120: Connection 1 still in use
T300: Query finally completes, connection released

Impact:
- 1 connection occupied for 5 minutes
- Only 19 connections available for fast queries
- If 5 analytics requests: only 15 connections left
- Pool effectively reduced by 25%
- Fast queries start queuing
- Users experience slowness
```

**Production symptoms:**

```
Monitoring:
- Pool utilization: 95%
- Active queries: 5 long-running, 15 normal
- Query time: P50 = 50ms, P99 = 300,000ms (5 min)

Database (pg_stat_activity):
SELECT pid, state, query_start, query
FROM pg_stat_activity
WHERE state = 'active'
AND query_start < NOW() - INTERVAL '10 seconds';

Result:
5 queries running for >1 minute
Query: SELECT * FROM orders WHERE created_at > ...
```

## ❌ Incorrect Approach #5: No Connection Timeout

```javascript
// Incorrect: Wait forever for connection
const pool = new Pool({
  max: 20,
  // No connectionTimeoutMillis!
});

app.get('/api/users', async (req, res) => {
  // Will wait FOREVER if pool exhausted
  const client = await pool.connect();
```

```
  try {
    const result = await client.query('SELECT * FROM users');
    res.json(result.rows);
  } finally {
    client.release();
  }
});
```

**Why it seems reasonable:**

- Patient waiting for available connection
- Eventually a connection will be free
- No errors thrown

**How it breaks:**

```
Scenario:
- Pool: 20 connections (all in use)
- 1000 requests arrive
- All 1000 call pool.connect()
- All 1000 wait indefinitely
- Node.js event loop blocked
- No new requests processed
- Health check fails (timeout)
- Load balancer removes instance
- Traffic shifts to other instances
- Other instances also exhaust pools
- Cascading failure

User experience:
- Browser: Request pending... (spinning forever)
- After 2 minutes: Browser timeout
- User: "Site is down"
```

**Production symptoms:**

```
Application:
- Event loop lag: 30,000ms
- Heap memory: Growing (queued promises)
- Active handles: 1,000+ (all waiting)

Load balancer:
- Health check timeout: Failed
- Instance marked unhealthy
- Removing from pool

Result: Complete outage
```

## 3. Correct Mental Model (How It Actually Works)

Connection pools maintain a fixed set of reusable database connections. Requests acquire connections from the pool, use them, then return them. When the pool is exhausted, requests wait (with timeout) or fail fast.

**The Pool Lifecycle**

```
1. Initialize pool (create min connections)
2. Request arrives
3. Acquire connection from pool
   - If available: Return immediately
   - If none available: Wait (up to timeout)
   - If timeout exceeded: Throw error
4. Execute query
5. Release connection back to pool
6. Connection available for next request
```

**Pool States**

```
Pool size: 20 (max)
├─ Available: 15 (idle, ready to use)
├─ In use: 5 (active queries)
└─ Waiting: 0 (requests waiting for connection)

Healthy: Available > 0, Waiting = 0
Warning: Available < 20%, Waiting > 0
Critical: Available = 0, Waiting > 10
```

**Resource Management**

```javascript
// Correct pattern
const client = await pool.connect(); // Acquire
try {
  await client.query(...); // Use
} finally {
  client.release(); // ALWAYS release
}
```

## 4. Correct Design & Algorithm

### Strategy 1: Fixed Pool with Timeout

```javascript
const pool = new Pool({
  max: 20, // Max connections
  min: 5,  // Min idle connections
  idleTimeoutMillis: 30000, // Close idle after 30s
  connectionTimeoutMillis: 5000, // Wait max 5s for connection
});
```

### Strategy 2: Monitor and Alert

```
setInterval(() => {
  const total = pool.totalCount;
  const idle = pool.idleCount;
  const waiting = pool.waitingCount;

  if (idle === 0 && waiting > 0) {
    alert('Pool exhausted!');
  }
}, 10000);
```

**Strategy 3: Circuit Breaker**

```
if (pool.waitingCount > 50) {
  // Pool overloaded, fail fast
  throw new Error('Service temporarily unavailable');
}
```

## 5. Full Production-Grade Implementation

```
import { Pool, PoolClient } from 'pg';
import { EventEmitter } from 'events';

/**
 * Configuration
 */
interface PoolConfig {
  max: number;
  min: number;
  idleTimeoutMillis: number;
  connectionTimeoutMillis: number;
  statementTimeout: number;
}

const config: PoolConfig = {
  max: 20, // Max connections in pool
  min: 5,  // Min idle connections
  idleTimeoutMillis: 30000, // Close idle connections after 30s
  connectionTimeoutMillis: 5000, // Wait max 5s for connection
  statementTimeout: 30000, // Kill queries after 30s
};

/**
 * Monitored connection pool
 */
class MonitoredPool extends EventEmitter {
  private pool: Pool;
  private metrics = {
    acquireCount: 0,
```

```typescript
    releaseCount: 0,
    errorCount: 0,
    timeoutCount: 0,
  };

  constructor(config: PoolConfig) {
    super();

    this.pool = new Pool({
      host: process.env.DB_HOST || 'localhost',
      database: process.env.DB_NAME || 'myapp',
      user: process.env.DB_USER || 'postgres',
      password: process.env.DB_PASSWORD,
      max: config.max,
      min: config.min,
      idleTimeoutMillis: config.idleTimeoutMillis,
      connectionTimeoutMillis: config.connectionTimeoutMillis,
      statement_timeout: config.statementTimeout,
    });

    this.setupMonitoring();
  }

  /**
   * Setup pool event listeners
   */
  private setupMonitoring() {
    this.pool.on('connect', (client) => {
      console.log('[Pool] New connection established');
    });

    this.pool.on('acquire', (client) => {
      this.metrics.acquireCount++;
      this.emit('acquire');
    });

    this.pool.on('remove', (client) => {
      console.log('[Pool] Connection removed');
    });

    this.pool.on('error', (err, client) => {
      this.metrics.errorCount++;
      console.error('[Pool] Unexpected error:', err);
      this.emit('error', err);
    });
  }

  /**
   * Acquire connection with monitoring
   */
  async connect(): Promise<PoolClient> {
    const startTime = Date.now();
```

```typescript
    try {
      const client = await this.pool.connect();

      const waitTime = Date.now() - startTime;

      if (waitTime > 1000) {
        console.warn(`[Pool] Slow connection acquire: ${waitTime}ms`);
      }

      // Wrap release to track metrics
      const originalRelease = client.release.bind(client);
      client.release = (err?: Error | boolean) => {
        this.metrics.releaseCount++;
        return originalRelease(err);
      };

      return client;
    } catch (error: any) {
      if (error.message.includes('timeout')) {
        this.metrics.timeoutCount++;
      }
      throw error;
    }
  }

  /**
   * Get pool statistics
   */
  getStats() {
    return {
      total: this.pool.totalCount,
      idle: this.pool.idleCount,
      waiting: this.pool.waitingCount,
      acquireCount: this.metrics.acquireCount,
      releaseCount: this.metrics.releaseCount,
      errorCount: this.metrics.errorCount,
      timeoutCount: this.metrics.timeoutCount,
      leaks: this.metrics.acquireCount - this.metrics.releaseCount,
    };
  }

  /**
   * Check pool health
   */
  isHealthy(): boolean {
    const stats = this.getStats();

    // Unhealthy if no idle connections and requests waiting
    if (stats.idle === 0 && stats.waiting > 0) {
      return false;
    }
```

```
      // Unhealthy if too many leaks
      if (stats.leaks > 5) {
        return false;
      }

      return true;
    }

    /**
     * Close pool
     */
    async close(): Promise<void> {
      await this.pool.end();
    }
  }

  // Initialize pool
  const pool = new MonitoredPool(config);

  /**
   * Express middleware: Check pool health
   */
  app.use(async (req, res, next) => {
    const stats = pool.getStats();

    // Circuit breaker: If pool exhausted, fail fast
    if (stats.idle === 0 && stats.waiting > 50) {
      return res.status(503).json({
        error: 'Service temporarily unavailable (database overload)',
        retryAfter: 10,
      });
    }

    next();
  });

  /**
   * Middleware: Track connection usage per request
   */
  app.use((req, res, next) => {
    const connectionAcquired = Date.now();

    res.on('finish', () => {
      const duration = Date.now() - connectionAcquired;

      if (duration > 5000) {
        console.warn(`[Request] Slow request: ${req.path} (${duration}ms)`);
      }
    });

    next();
```

```typescript
});

/**
 * Correct usage pattern
 */
app.get('/api/users', async (req, res) => {
  let client: PoolClient | null = null;

  try {
    // Acquire connection
    client = await pool.connect();

    // Execute query
    const result = await client.query('SELECT id, name, email FROM users LIMIT
100');

    res.json(result.rows);
  } catch (error: any) {
    if (error.message.includes('timeout')) {
      res.status(503).json({
        error: 'Database connection timeout',
        retryAfter: 5,
      });
    } else {
      res.status(500).json({ error: 'Internal error' });
    }
  } finally {
    // ALWAYS release connection
    if (client) {
      client.release();
    }
  }
});

/**
 * Transaction pattern
 */
async function transferMoney(fromId: number, toId: number, amount: number) {
  let client: PoolClient | null = null;

  try {
    client = await pool.connect();

    await client.query('BEGIN');

    await client.query(
      'UPDATE accounts SET balance = balance - $1 WHERE id = $2',
      [amount, fromId]
    );

    await client.query(
      'UPDATE accounts SET balance = balance + $1 WHERE id = $2',
```

```javascript
      [amount, toId]
    );

    await client.query('COMMIT');
  } catch (error) {
    if (client) {
      await client.query('ROLLBACK');
    }
    throw error;
  } finally {
    if (client) {
      client.release();
    }
  }
}

/**
 * Connection pool metrics endpoint
 */
app.get('/metrics/pool', (req, res) => {
  const stats = pool.getStats();

  res.json({
    ...stats,
    health: pool.isHealthy() ? 'healthy' : 'unhealthy',
    utilizationPercent: ((stats.total - stats.idle) / stats.total) * 100,
  });
});

/**
 * Health check with pool status
 */
app.get('/health', (req, res) => {
  const stats = pool.getStats();

  if (!pool.isHealthy()) {
    return res.status(503).json({
      status: 'unhealthy',
      reason: 'Connection pool exhausted',
      stats,
    });
  }

  res.status(200).json({
    status: 'healthy',
    stats,
  });
});

/**
 * Monitoring: Alert on pool exhaustion
 */
```

```typescript
setInterval(() => {
  const stats = pool.getStats();

  const utilization = ((stats.total - stats.idle) / stats.total) * 100;

  if (utilization > 80) {
    console.warn(`[Alert] High pool utilization: ${utilization.toFixed(1)}%`);
  }

  if (stats.waiting > 0) {
    console.warn(`[Alert] ${stats.waiting} requests waiting for connection`);
  }

  if (stats.leaks > 5) {
    console.error(`[Alert] Potential connection leak: ${stats.leaks} unreleased`);
  }
}, 10000); // Every 10 seconds

/**
 * Query timeout handling
 */
async function queryWithTimeout(
  client: PoolClient,
  query: string,
  params: any[],
  timeout: number = 5000
): Promise<any> {
  // Set statement timeout for this query
  await client.query(`SET LOCAL statement_timeout = ${timeout}`);

  try {
    return await client.query(query, params);
  } catch (error: any) {
    if (error.message.includes('statement timeout')) {
      throw new Error(`Query timeout after ${timeout}ms`);
    }
    throw error;
  }
}

/**
 * Separate pool for analytics (slow queries)
 */
const analyticsPool = new MonitoredPool({
  max: 5, // Smaller pool for analytics
  min: 1,
  idleTimeoutMillis: 60000,
  connectionTimeoutMillis: 10000,
  statementTimeout: 300000, // 5 minutes
});

app.get('/api/analytics/revenue', async (req, res) => {
```

```typescript
  let client: PoolClient | null = null;

  try {
    // Use separate pool for slow queries
    client = await analyticsPool.connect();

    const result = await client.query(`
      SELECT
        date_trunc('day', created_at) as day,
        sum(amount) as revenue
      FROM orders
      WHERE created_at > NOW() - INTERVAL '90 days'
      GROUP BY day
      ORDER BY day
    `);

    res.json(result.rows);
  } catch (error: any) {
    res.status(500).json({ error: error.message });
  } finally {
    if (client) {
      client.release();
    }
  }
});

/**
 * Kill long-running queries
 */
async function killLongQueries(thresholdMs: number = 60000) {
  let client: PoolClient | null = null;

  try {
    client = await pool.connect();

    const result = await client.query(`
      SELECT
        pid,
        now() - query_start as duration,
        query
      FROM pg_stat_activity
      WHERE state = 'active'
        AND now() - query_start > interval '${thresholdMs} milliseconds'
        AND query NOT ILIKE '%pg_stat_activity%'
    `);

    for (const row of result.rows) {
      console.warn(`[Kill] Terminating query (PID ${row.pid}): ${row.query}`);
      await client.query(`SELECT pg_terminate_backend(${row.pid})`);
    }

    return result.rows.length;
```

```typescript
    } finally {
      if (client) {
        client.release();
      }
    }
  }
}

// Run every minute
setInterval(() => {
  killLongQueries(60000).then((count) => {
    if (count > 0) {
      console.log(`[Kill] Terminated ${count} long-running queries`);
    }
  });
}, 60000);

/**
 * Graceful shutdown with pool cleanup
 */
process.on('SIGTERM', async () => {
  console.log('[Shutdown] Closing connection pools...');

  await pool.close();
  await analyticsPool.close();

  console.log('[Shutdown] Pools closed');
  process.exit(0);
});

/**
 * Connection leak detector
 */
class ConnectionLeakDetector {
  private acquisitions = new Map<PoolClient, { stack: string; timestamp: number }>();

  trackAcquire(client: PoolClient) {
    const stack = new Error().stack || '';
    this.acquisitions.set(client, {
      stack,
      timestamp: Date.now(),
    });
  }

  trackRelease(client: PoolClient) {
    this.acquisitions.delete(client);
  }

  detectLeaks(thresholdMs: number = 60000) {
    const now = Date.now();
    const leaks: Array<{ duration: number; stack: string }> = [];
```

```
    for (const [client, info] of this.acquisitions.entries()) {
      const duration = now - info.timestamp;

      if (duration > thresholdMs) {
        leaks.push({ duration, stack: info.stack });
      }
    }

    return leaks;
  }
}

const leakDetector = new ConnectionLeakDetector();

// Check for leaks every minute
setInterval(() => {
  const leaks = leakDetector.detectLeaks(60000);

  if (leaks.length > 0) {
    console.error(`[Leak] Detected ${leaks.length} potential connection leaks:`);
    leaks.forEach((leak) => {
      console.error(`  - Held for ${leak.duration}ms`);
      console.error(`  - Stack trace: ${leak.stack.split('\n').slice(0,
5).join('\n')}`);
    });
  }
}, 60000);
```

## 6. Correct Usage Patterns (Where This Shines)

### Pattern 1: Always Release in Finally

```
const client = await pool.connect();
try {
  await client.query(...);
} finally {
  client.release(); // ALWAYS executes
}
```

### Pattern 2: Separate Pools for Different Workloads

```
const fastPool = new Pool({ max: 20 }); // Fast queries
const slowPool = new Pool({ max: 5 });  // Analytics
```

### Pattern 3: Circuit Breaker on Pool Exhaustion

```
if (pool.waitingCount > 50) {
  throw new ServiceUnavailableError();
```

```
}
```

## 7. Failure Modes & Edge Cases

### Pool Exhaustion

**Problem:** All connections in use, requests timeout.

**Mitigation:** Increase pool size, kill slow queries, circuit breaker.

### Connection Leaks

**Problem:** Connections never released, pool depletes over time.

**Mitigation:** Always use finally blocks, leak detection.

### Slow Query Starvation

**Problem:** One slow query blocks connection for minutes.

**Mitigation:** Statement timeout, separate pools, kill queries.

## 8. Performance Characteristics & Tradeoffs

### Pool Size

- **Too small:** Frequent exhaustion, high wait times
- **Too large:** Database overload, memory waste
- **Sweet spot:** 2-4× number of CPU cores

### Connection Overhead

- **No pool:** 200ms per request (connection setup)
- **With pool:** 0ms (reuse existing)

## 9. Foot-Guns & Common Mistakes (DO NOT SKIP)

### Mistake 1: Forgetting to Release

**Fix:** ALWAYS use finally block.

### Mistake 2: No Connection Timeout

**Fix:** Set connectionTimeoutMillis (5-10s).

### Mistake 3: No Statement Timeout

**Fix:** Set statement_timeout (30s for API, 5min for analytics).

### Mistake 4: Single Pool for All Workloads

**Fix:** Separate pools for fast/slow queries.

### Mistake 5: Infinite Pool Size

**Fix:** Set max connections (20 for API, 5 for analytics).

## 10. When NOT to Use This (Anti-Patterns)

### Serverless Functions

Use single connection or connection proxy (RDS Proxy, PgBouncer).

### Microservices with Low Traffic

Smaller pool (5-10 connections) sufficient.

## 11. Related Concepts (With Contrast)

### Connection Pooler (PgBouncer)

**Difference:** Pool at application level, PgBouncer pools at database level.

### Read Replicas

**Related:** Separate pools for primary (writes) vs replicas (reads).

### Circuit Breaker

**Related:** Prevents cascading failure when pool exhausted.

## 12. Production Readiness Checklist

### Pool Configuration

- ☐ Set max connections (20 for API)
- ☐ Set min idle connections (5)
- ☐ Set connection timeout (5-10s)
- ☐ Set statement timeout (30s for API, 5min for analytics)
- ☐ Set idle timeout (30s)

### Resource Management

- ☐ Always release connections in finally block
- ☐ Use transactions properly (BEGIN/COMMIT/ROLLBACK)
- ☐ Separate pools for fast/slow queries
- ☐ Connection leak detection

### Monitoring

- ☐ Track pool utilization (% used)
- ☐ Alert when utilization >80%
- ☐ Alert when waiting count >0
- ☐ Dashboard: idle, active, waiting connections
- ☐ Track connection leaks (acquire - release)

### Error Handling

- ☐ Handle connection timeout gracefully
- ☐ Return 503 when pool exhausted

- [ ] Circuit breaker for overload protection
- [ ] Kill long-running queries (>60s)