# Building a High-Performance Mini SQLite-Based Database System

## From Concept to Implementation

Database Systems Architecture

July 8, 2025

Project Overview & Objectives
Core Architecture Overview
Core Components Design
Indexing Strategy
Query Processing Pipeline
Concurrency & Parallelization
Performance Optimizations
Implementation Roadmap
Testing & Benchmarking
Conclusion & Next Steps

# Table of Contents

## Project Vision

### Goal

Build a SQLite-inspired database system optimized for modern hardware constraints
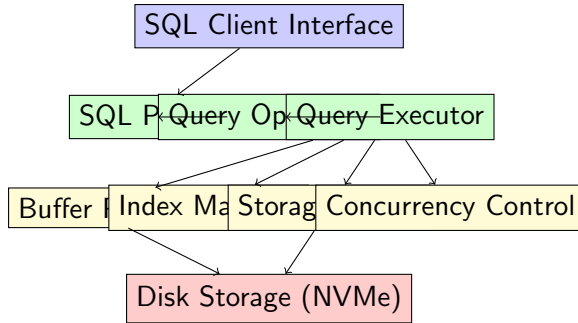
**Traditional SQLite Constraints (2000s):**

- Limited RAM (64MB-1GB)
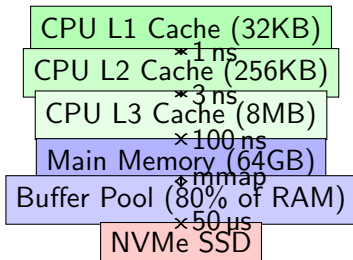- Slow disk I/O
- Single-core CPUs
- Conservative memory usage

**Modern Reality (2025):**

- Abundant RAM (16-128GB)
- Fast NVMe SSDs
- Multi-core CPUs (8-32 cores)
- Speed is everything

## System Architecture - High Level

SQL Client Interface

SQL P Query Op Query Executor

Buffer F Index Ma Storag Concurrency Control

Disk Storage (NVMe)

## Memory-First Architecture



CPU L1 Cache (32KB)

×1 ns

CPU L2 Cache (256KB)

×3 ns

CPU L3 Cache (8MB)

×100 ns

Main Memory (64GB)

mmap

Buffer Pool (80% of RAM)

×50 µs

NVMe SSD

**Key Principles:**

- Keep hot data in RAM
- Minimize disk I/O
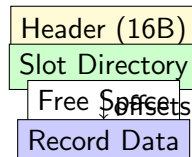- Optimize for cache locality
- Use memory-mapped files

**Implementation:**

- 80-90% RAM for buffer pool
- Async disk writes
- Prefetching algorithms

# Storage Engine - Data Structures

## Listing 1: Page Structure

```
1    class Page {
2        static const size_t PAGE_SIZE = 4096;
3
4        struct Header {
5            uint32_t page_id;
6            uint16_t free_space;
7            uint16_t slot_count;
8            uint32_t checksum;
9        };
10
11        Header header;
12        std::vector<uint16_t> slot_directory;
13        std::vector<uint8_t> data;
14
15   public:
16        bool insert_record(const Record&
                record);
17        Record get_record(uint16_t slot_id);
```

Header (16B)

Slot Directory

Free Space ↓ offsets

Record Data
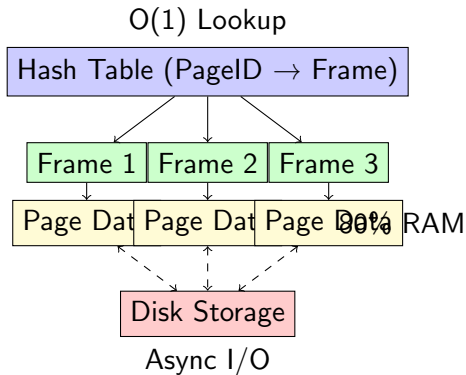
**Advantages:**

- Cache-friendly layout
- Variable-length records

## Buffer Pool Manager

### Listing 2: Buffer Pool Implementation

```
1    class BufferPool {
2        struct BufferFrame {
3            Page* page;
4            uint32_t page_id;
5            bool is_dirty;
6            uint32_t pin_count;
7            std::chrono::steady_clock::time_point last_access;
8        };
9
10       std::unordered_map<uint32_t, BufferFrame*> page_table;
11       std::vector<BufferFrame> buffer_frames;
12       std::mutex buffer_mutex;
13
14   public:
15       Page* get_page(uint32_t page_id);
16       void unpin_page(uint32_t page_id, bool is_dirty);
17       void flush_all_pages();
18
19   private:
```
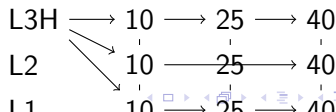
# Buffer Pool - Memory Management Strategy

O(1) Lookup

Hash Table (PageID → Frame)

Frame 1 | Frame 2 | Frame 3

Page Dat | Page Dat | Page Dat 80% RAM

Disk Storage

Async I/O

## Skip List Index Implementation

Listing 3: Skip List Node

```
template<typename Key,
    typename Value>
class SkipListNode {
    Key key;
    Value value;
    std::vector<SkipListNode
        *> forward;

public:
    SkipListNode(Key k,
        Value v, int level
        )
        : key(k), value(v),
            forward(level
            + 1) {}

    Key get_key() const {
        return key; }
```

# Hash Index for Equality Queries

## Listing 4: Robin Hood Hash Table

```cpp
template<typename Key, typename Value>
class RobinHoodHashTable {
    struct Entry {
        Key key;
        Value value;
        uint32_t hash;
        uint32_t distance;   // Distance from ideal position
        bool is_occupied;
    };

    std::vector<Entry> table;
    size_t capacity;
    size_t size;

public:
    bool insert(const Key& key, const Value& value);
    bool find(const Key& key, Value& value);
    bool remove(const Key& key);
```

# Query Processing Flow

SQL Query

Lexer/Tokenizer    st...           ...tom tokenizer

Parser (AST)       Re...           ...nt parser

Query Optimizer    Co...           ...mization

Execution Engine   Vo...           ... model

Result Set

SELECT
——
columns
—— FROM
 +——
WHERE

## SQL Parser - AST Construction

### Listing 5: Abstract Syntax Tree Nodes

```cpp
class ASTNode {
public:
    virtual ~ASTNode() = default;
    virtual void accept(ASTVisitor& visitor) = 0;
};

class SelectStatement : public ASTNode {
    std::vector<std::unique_ptr<Expression>> select_list;
    std::unique_ptr<FromClause> from_clause;
    std::unique_ptr<WhereClause> where_clause;
    std::unique_ptr<OrderByClause> order_by_clause;

public:
    void accept(ASTVisitor& visitor) override {
        visitor.visit(*this);
    }
};

class TableReference : public ASTNode {
```

# Query Optimizer - Cost-Based Decisions

## Listing 6: Query Optimization Framework

```cpp
1    class QueryOptimizer {
2        struct Statistics {
3            size_t table_size;
4            size_t distinct_values;
5            double selectivity;
6            std::unordered_map<std::string, size_t> column_stats;
7        };
8
9        std::unordered_map<std::string, Statistics> table_stats;
10
11   public:
12       std::unique_ptr<ExecutionPlan> optimize(const ASTNode& query);
13
14   private:
15       double estimate_cost(const ExecutionPlan& plan);
16       std::vector<ExecutionPlan> generate_plans(const ASTNode& query);
17       ExecutionPlan select_best_plan(const std::vector<ExecutionPlan>& plans);
18
19       // Optimization rules
```

# CPU Parallelization with OpenMP

## Listing 7: Parallel Table Scan

```
 1          class TableScan {
 2              std::vector<Page*> pages;
 3              Predicate where_condition;
 4
 5          public:
 6              std::vector<Record> scan_parallel() {
 7                  std::vector<Record> results;
 8                  std::mutex results_mutex;
 9
10                  #pragma omp parallel for
11                  for (size_t i = 0; i < pages.size(); ++i) {
12                      std::vector<Record> local_results;
13
14                      // Scan page i
15                      for (const auto& record : pages[i]->get_records()) {
16                          if (where_condition.evaluate(record)) {
17                              local_results.push_back(record);
18                          }
19                      }
```
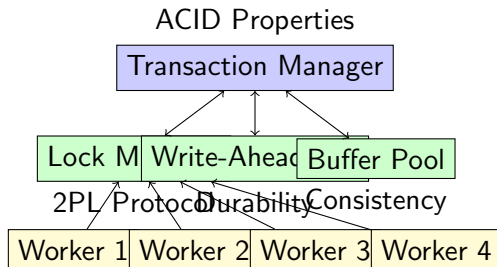
# SIMD Vectorization for Batch Operations

## Listing 8: Vectorized Comparison Operations

```
1    #include <immintrin.h>
2
3    class VectorizedOperations {
4    public:
5        // Compare 8 integers simultaneously
6        static std::vector<bool> compare_greater_than(
7            const std::vector<int32_t>& values,
8            int32_t threshold) {
9
10           std::vector<bool> results(values.size());
11           const __m256i threshold_vec = _mm256_set1_epi32(threshold);
12
13           size_t i = 0;
14           for (; i + 8 <= values.size(); i += 8) {
15               __m256i values_vec = _mm256_loadu_si256(
16                   reinterpret_cast<const __m256i*>(&values[i]));
17
18               __m256i cmp_result = _mm256_cmpgt_epi32(values_vec, threshold_vec);
19
```

## Concurrency Control Architecture



**Concurrency Strategy:**

- **Two-Phase Locking (2PL)**: Strict locking protocol
- **Write-Ahead Logging**: Ensure durability and recovery

# Cache-Conscious Data Structures
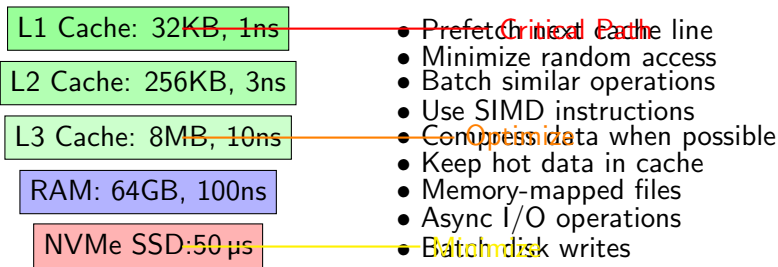
## Listing 9: Cache-Friendly Record Layout

```cpp
// Traditional approach - poor cache locality
struct Record_Bad {
    std::string* name;      // Pointer to heap
    int* age;               // Pointer to heap
    double* salary;         // Pointer to heap
};

// Cache-friendly approach
struct Record_Good {
    char name[64];          // Fixed-size, inline storage
    int age;                // Value stored inline
    double salary;          // Value stored inline

    // Padding to align to cache line boundary
    char padding[64 - sizeof(name) - sizeof(age) - sizeof(salary)];
} __attribute__((aligned(64)));

// Columnar storage for analytical queries
class ColumnStore {
```

## Async I/O and Prefetching

### Listing 10: Asynchronous Disk Operations

```cpp
#include <future>
#include <thread>

class AsyncDiskManager {
    std::thread io_thread;
    std::queue<std::packaged_task<void()>> io_queue;
    std::mutex queue_mutex;
    std::condition_variable cv;

public:
    AsyncDiskManager() : io_thread(&AsyncDiskManager::io_worker, this) {}

    std::future<Page*> read_page_async(uint32_t page_id) {
        auto task = std::packaged_task<Page*()>(
            [this, page_id]() { return read_page_from_disk(page_id); });

        auto future = task.get_future();

        {
```

# Memory Hierarchy Optimization

L1 Cache: 32KB, 1ns

L2 Cache: 256KB, 3ns

L3 Cache: 8MB, 10ns

RAM: 64GB, 100ns

NVMe SSD: 50 μs

Critical Path

Optimize

Minimize

- Prefetch next cache line
- Minimize random access
- Batch similar operations
- Use SIMD instructions
- Compress data when possible
- Keep hot data in cache
- Memory-mapped files
- Async I/O operations
- Batch disk writes

## Development Phases

Phase 1: Foundation
- Basic storage engine
- Simple B+ tree index
- Basic SQL parser

Phase 2: Core Features
- Transaction support
- Concurrency control

Phase 3: Optimization
- Query optimization
- Parallel execution
- SIMD operations
- Cache optimization

Phase 4: Advanced
- Advanced indexing
- Compression
- Distributed features

## Phase 1: Foundation Implementation

### Listing 11: Basic Storage Engine Skeleton

```
1    // Core classes to implement first
2    class StorageEngine {
3        std::unique_ptr<BufferPool> buffer_pool;
4        std::unique_ptr<DiskManager> disk_manager;
5
6    public:
7        bool create_table(const std::string& table_name,
8                          const TableSchema& schema);
9        bool insert_record(const std::string& table_name,
10                           const Record& record);
11        std::vector<Record> scan_table(const std::string& table_name);
12    };
13
14    class DiskManager {
15        std::fstream db_file;
16        std::mutex file_mutex;
17
18    public:
19        Page* read_page(uint32_t page_id);
```

# Phase 2: Adding SQL and Transactions

## Listing 12: SQL Parser and Transaction Manager

```
1    // SQL Parser using recursive descent
2    class SQLParser {
3        std::regex select_pattern;
4        std::regex insert_pattern;
5        std::regex create_table_pattern;
6
7    public:
8        SQLParser() {
9            select_pattern = std::regex(R"(SELECT\s+(.+)\s+FROM\s+(\w+)(?:\s+WHERE\s+(.+))?
                )");
10           insert_pattern = std::regex(R"(INSERT\s+INTO\s+(\w+)\s+VALUES\s*\((.+)\))");
11       }
12
13       std::unique_ptr<Statement> parse(const std::string& sql);
14   };
15
16   // Transaction Manager
17   class TransactionManager {
18       std::atomic<uint64_t> next_trn_id{1};
```

# Phase 3: Performance Optimization

## Listing 13: Parallel Query Execution

```cpp
// Thread pool for parallel execution
class ThreadPool {
    std::vector<std::thread> workers;
    std::queue<std::function<void()>> tasks;
    std::mutex queue_mutex;
    std::condition_variable condition;
    bool stop = false;

public:
    ThreadPool(size_t num_threads = std::thread::hardware_concurrency());

    template<typename F>
    auto enqueue(F&& f) -> std::future<decltype(f())> {
        auto task = std::make_shared<std::packaged_task<decltype(f())()>>(
            std::forward<F>(f));

        auto future = task->get_future();
        {
            std::unique_lock<std::mutex> lock(queue_mutex);
```

# Phase 4: Advanced Features

## Listing 14: Compression and Advanced Indexing

```
1      // Data compression for storage efficiency
2      class CompressionManager {
3      public:
4          // Dictionary compression for strings
5          std::vector<uint8_t> compress_string_column(
6              const std::vector<std::string>& strings) {
7              std::unordered_map<std::string, uint32_t> dictionary;
8              std::vector<uint32_t> compressed_data;
9
10             uint32_t next_id = 0;
11             for (const auto& str : strings) {
12                 if (dictionary.find(str) == dictionary.end()) {
13                     dictionary[str] = next_id++;
14                 }
15                 compressed_data.push_back(dictionary[str]);
16             }
17
18             return serialize_compressed_data(dictionary, compressed_data);
19         }
```

## Performance Testing Framework

### Listing 15: Benchmark Suite Implementation

```
1    #include <chrono>
2    #include <random>
3
4    class DatabaseBenchmark {
5        std::unique_ptr<StorageEngine> engine;
6        std::mt19937 rng;
7
8    public:
9        // TPC-C style benchmark
10       void benchmark_oltp_workload() {
11           const size_t num_transactions = 10000;
12           const size_t num_threads = std::thread::hardware_concurrency();
13
14           auto start = std::chrono::high_resolution_clock::now();
15
16           std::vector<std::thread> workers;
17           for (size_t i = 0; i < num_threads; ++i) {
18               workers.emplace_back([this, num_transactions, num_threads, i]() {
19                   for (size_t i = i; i < num_transactions; i += num_threads) {
```

## Performance Metrics & Targets

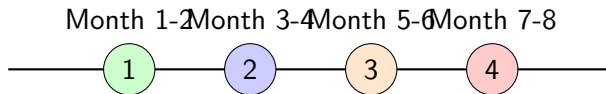| Operation | Target | Measurement | Comparison |
|-----------|--------|-------------|------------|
| Point Query | ¡ 0.1ms | std::chrono | vs SQLite |
| Range Scan | 1M rows/sec | Throughput | vs PostgreSQL |
| Insert | 10K TPS | Transactions/sec | vs MySQL |
| Parallel Scan | 4x speedup | vs serial | CPU cores |
| Memory Usage | 80% buffer hit | Cache ratio | vs traditional |

Query Time (ms)

SQLite

## Key Takeaways

### Architecture Principles

- **Memory-first design**: Leverage abundant RAM for speed
- **CPU parallelization**: Use all available cores with OpenMP
- **Cache-conscious data structures**: Optimize for modern CPU hierarchy
- **Async I/O**: Overlap computation and disk operations

### Implementation Strategy

- **Start simple**: Basic storage engine with file I/O
- **Iterate quickly**: Add features incrementally
- **Measure everything**: Benchmark each optimization

# Development Timeline

Month 1-2  Month 3-4  Month 5-6  Month 7-8

(1)      (2)      (3)      (4)

Basic SQL  Storage  Transactions  Parallelization  Advanced Features

Deliverables:
• Working prototype
• Performance benchmarks
• Documentation
• Test suite

## Questions & Discussion

# Thank You!

### Contact & Resources

- **Source code**: Will be available on GitHub
- **Documentation**: Comprehensive API documentation
- **Benchmarks**: Performance comparison results
- **Community**: Open source contribution guidelines