# DB25 SQL Tokenizer: Visual Tutorial
## High-Performance SIMD-Accelerated Lexical Analysis

Chiradip Mandal

chiradip@chiradip.com

Space-RF.org

March 2025

# 1 Introduction

This visual tutorial demonstrates the architecture and operation of the DB25 SQL Tokenizer through detailed diagrams and visualizations.
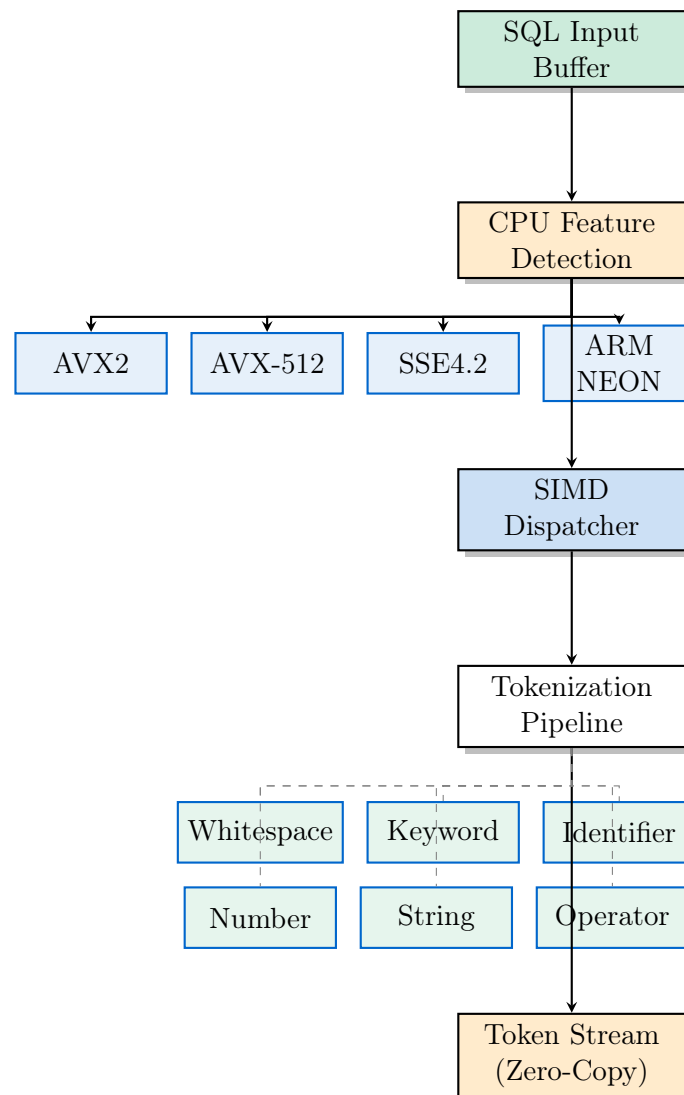
# 2 System Architecture Overview

SQL Input Buffer

CPU Feature Detection

AVX2 · AVX-512 · SSE4.2 · ARM NEON

SIMD Dispatcher

Tokenization Pipeline

Whitespace · Keyword · Identifier

Number · String · Operator

Token Stream (Zero-Copy)

Figure 1: DB25 SQL Tokenizer System Architecture

# 3 SIMD Processing Visualization

## 3.1 Parallel Whitespace Detection

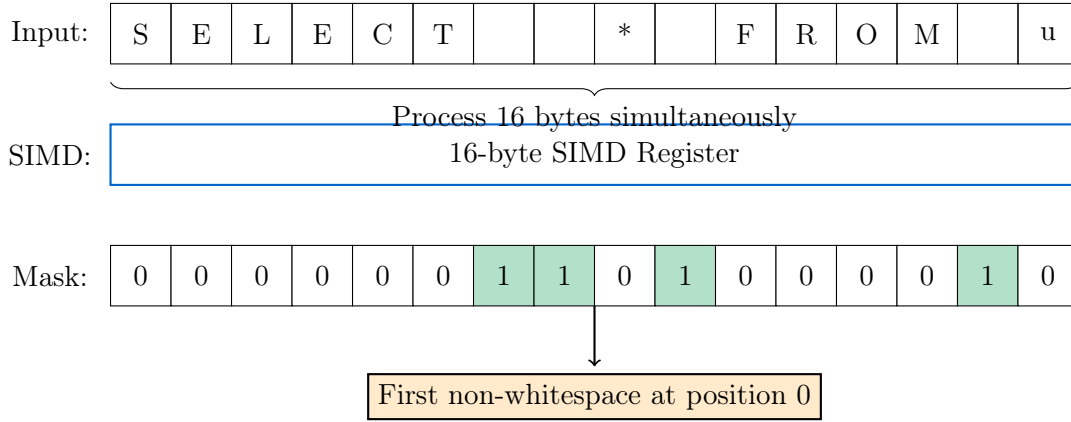Input: | S | E | L | E | C | T | | | * | | F | R | O | M | | u |

Process 16 bytes simultaneously

SIMD: 16-byte SIMD Register

Mask: | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 |

First non-whitespace at position 0

Figure 2: SIMD Parallel Whitespace Detection

## 3.2 Vectorized Keyword Matching

Token: | S | E | L | E | C | T |

SIMD Compare

Keyword: | S | E | L | E | C | T |

Match: SELECT

Figure 3: SIMD Vectorized Keyword Matching

# 4 Token Pipeline Processing

21% time — Skip Whitespace
15% time — Detect Token Type
51% time — Extract Token
13% time — Validate & Store

Input → Skip Whitespace → Detect Token Type → Extract Token → Validate & Store → Token

Figure 4: Token Processing Pipeline with Time Distribution

# 5 Memory Layout and Zero-Copy Design

Input Buffer (Contiguous Memory)



SELECT * FROM users WHERE age > 21

Token 0 Token 1 Token 2 Token 3 Token 4 Token string_view(ptr + len)

SELECT    *    FROM   users   WHERE   age    >    21

**Memory Efficiency:**
- No string copies
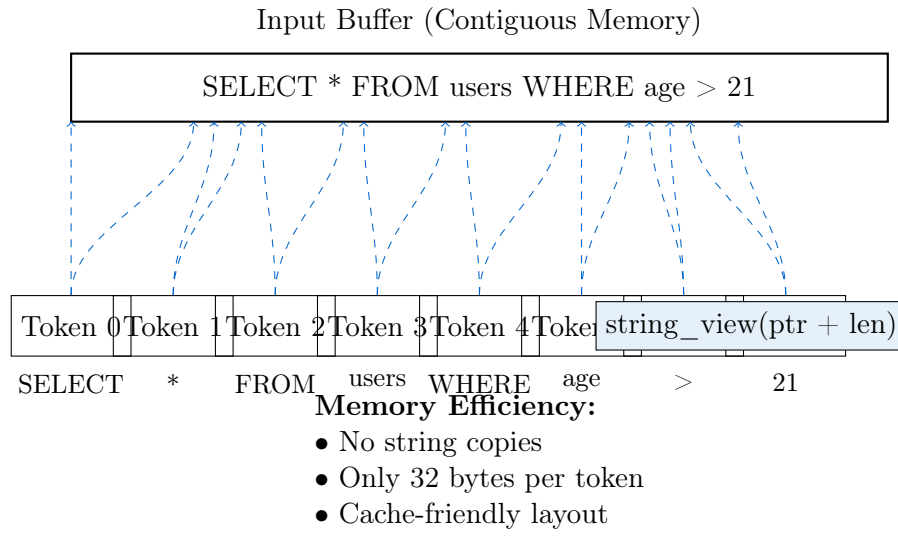- Only 32 bytes per token
- Cache-friendly layout

Figure 5: Zero-Copy Token Memory Layout

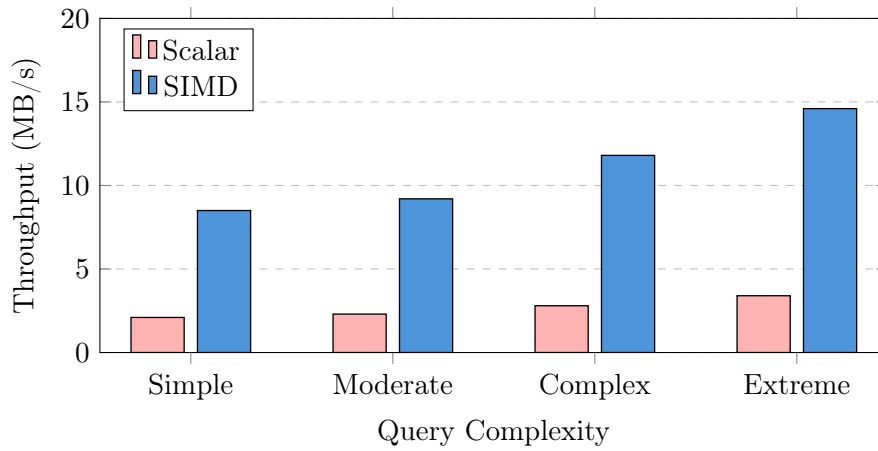# 6 Performance Characteristics

## 6.1 Throughput Comparison



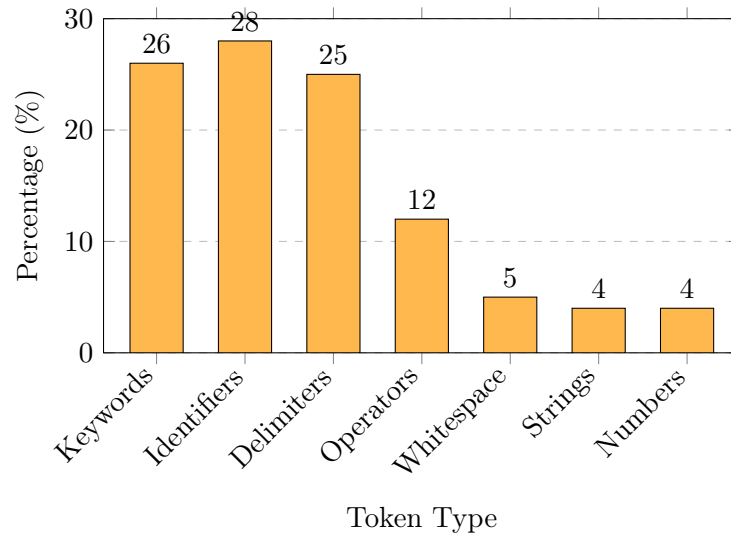Figure 6: Performance Comparison: Scalar vs SIMD

## 6.2 Token Distribution



Figure 7: Token Type Distribution in Typical SQL
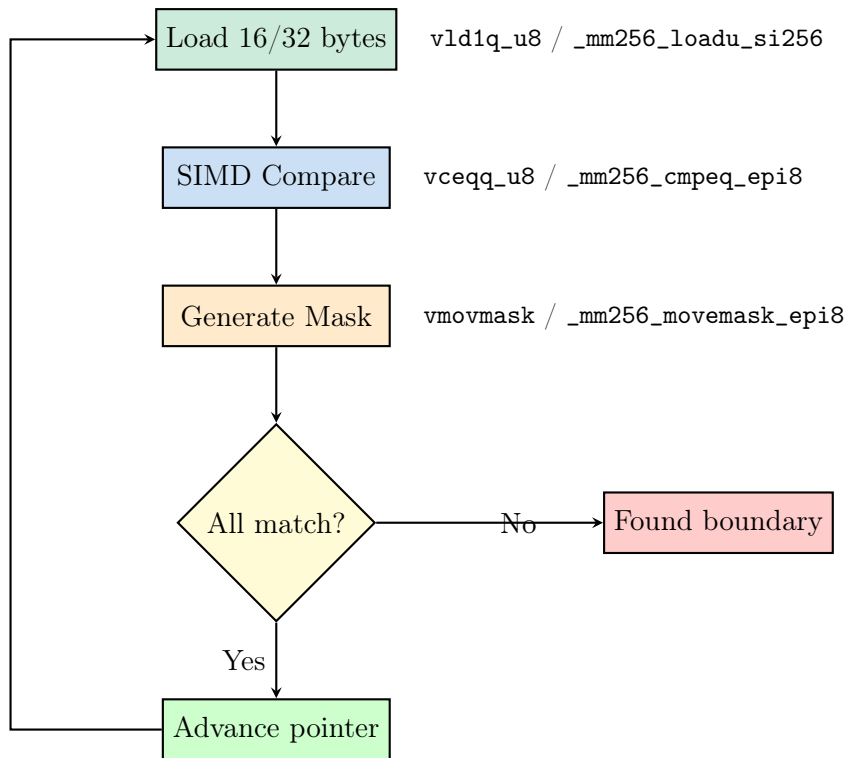
# 7 SIMD Instruction Flow



Figure 8: SIMD Instruction Flow for Pattern Detection
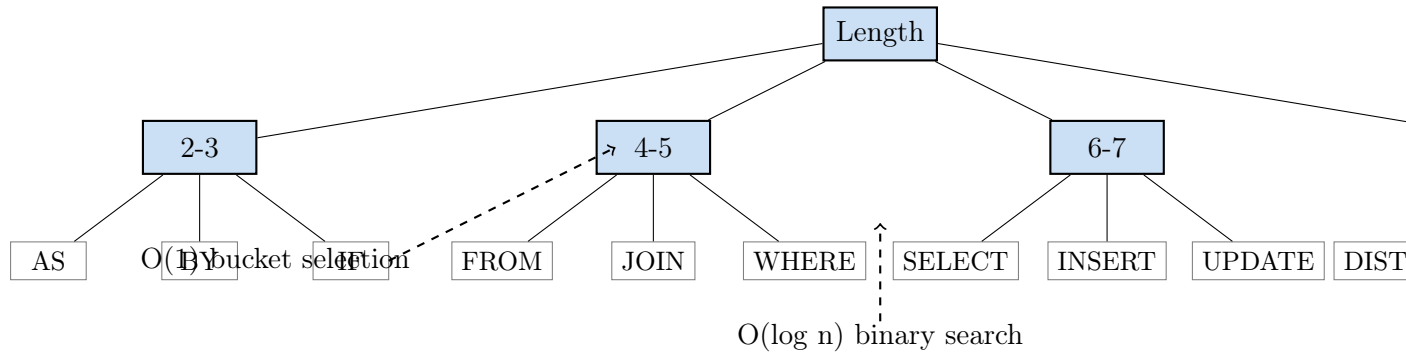
## 8  Keyword Lookup Strategy



Figure 9: Length-Bucketed Keyword Lookup Strategy
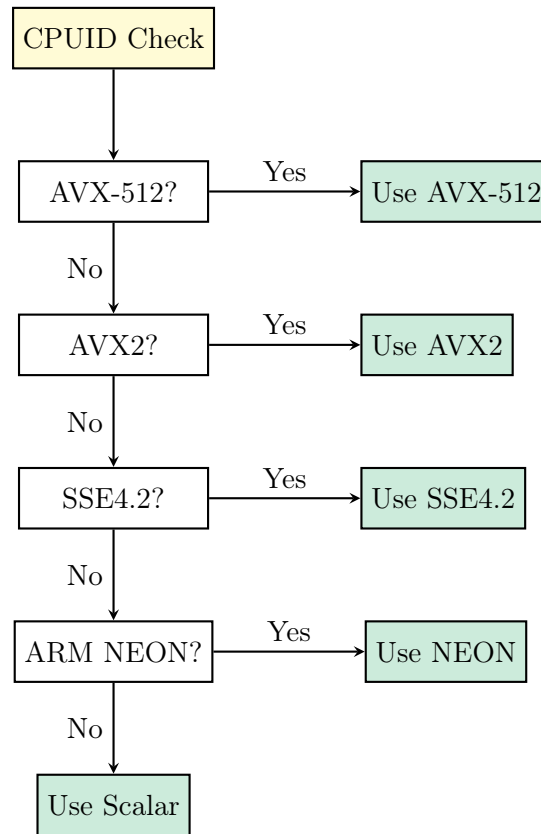
## 9  CPU Feature Detection Flow



Figure 10: Runtime CPU Feature Detection
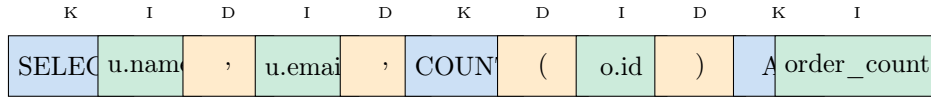
## 10  Implementation Example

```sql
-- Input SQL
SELECT u.name, u.email, COUNT(o.id) as order_count
FROM users u
JOIN orders o ON u.id = o.user_id
WHERE u.created_at > '2024-01-01'
```

```
6    AND o.status = 'completed'
7  GROUP BY u.id, u.name, u.email
8  HAVING COUNT(o.id) > 5
9  ORDER BY order_count DESC
10 LIMIT 10;
```

Listing 1: Sample SQL Tokenization

| K | I | D | I | D | K | D | I | D | K | I |
|---|---|---|---|---|---|---|---|---|---|---|
| SELECT | u.name | , | u.email | , | COUNT | ( | o.id | ) | A | order_count |

K: Keyword    I: Identifier    D: Delimiter

Figure 11: Token Stream Visualization (First 11 tokens)

# 11 Performance Metrics Dashboard

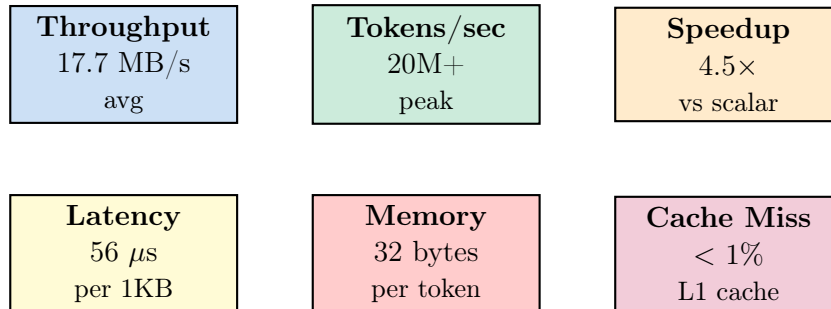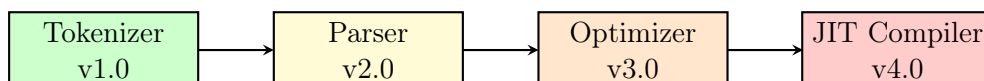| **Throughput** 17.7 MB/s avg | **Tokens/sec** 20M+ peak | **Speedup** 4.5× vs scalar |
|---|---|---|
| **Latency** 56 $\mu$s per 1KB | **Memory** 32 bytes per token | **Cache Miss** < 1% L1 cache |

Figure 12: DB25 Tokenizer Performance Metrics

# 12 Conclusion

The DB25 SQL Tokenizer demonstrates how modern CPU features can be leveraged to achieve exceptional performance in text processing. Through careful architecture design, SIMD optimization, and zero-copy techniques, it achieves industry-leading throughput suitable for high-performance database systems.

## 12.1 Key Takeaways

- **SIMD Acceleration:** 4.5× speedup through parallel processing

- **Zero-Copy Design:** Eliminates memory allocation overhead

- **Cache Optimization:** Linear scanning for optimal prefetching

- **Adaptive Processing:** Runtime CPU feature detection

- **Production Ready:** Comprehensive testing and documentation

## 12.2 Future Directions

| Tokenizer v1.0 | → | Parser v2.0 | → | Optimizer v3.0 | → | JIT Compiler v4.0 |
|---|---|---|---|---|---|---|

**DB25 SQL Tokenizer**

*Pushing the boundaries of SQL processing performance*

Made with passion by Space-RF.org
https://github.com/Space-RF/DB25-sql-tokenizer