# MATCH_RECOGNIZE Pattern Matching Performance Evaluation
## Amazon UK Product Dataset (2.2M rows)

### Performance Test Results

### October 26, 2025

## 1 Executive Summary

This document presents comprehensive performance evaluation results for the MATCH_RECOGNIZE implementation using the Amazon UK product dataset. The evaluation covers 25 test cases across 5 SQL patterns and 5 dataset sizes (25,000 to 100,000 rows).

**Key Results:**

The evaluation achieved a 100% test success rate with all 25 tests passed successfully. The implementation demonstrated an average throughput of 9,838 rows per second across all test cases. Pattern coverage ranged from 4.31% to 53.76%, indicating natural pattern presence in the dataset. The total execution time for all 25 tests was 157.44 seconds, demonstrating efficient performance at scale.

## 2 Overall Test Statistics

Table 1: Overall Test Statistics

| Metric | Value |
|---|---|
| Total Tests | 25 |
| Success Rate | 100% |
| Total Execution Time | 157.44 sec |
| Average Execution Time | 6.30 sec |
| Average Coverage | 27.91% |
| Average Throughput | 9,838 rows/sec |
| Min Coverage | 4.31% |
| Max Coverage | 53.76% |
| Min Throughput | 6,481 rows/sec |
| Max Throughput | 13,097 rows/sec |

## 3 Pattern Definitions

## 4 Performance by Pattern

**Pattern Analysis:**

The complex_nested pattern achieved the best coverage at 49.64%, demonstrating its ability to match nearly half of the dataset. The simple_sequence pattern delivered the best throughput at 12,618 rows per second, making it the most efficient for processing large datasets. This same pattern also achieved the fastest execution time with an average of 4.57 seconds per test. The alternation pattern proved most selective with 6.06% coverage, effectively filtering for specific quality degradation sequences.

Table 2: SQL Pattern Definitions

| Pattern Name | SQL Pattern | Description |
|---|---|---|
| simple_sequence | A+ B+ | Simple sequence: A followed by B |
| alternation | A (B\|C)+ D | Alternation: A followed by (B or C) followed by D |
| quantified | A{2,5} B* C+ | Quantified: 2-5 A's, optional B's, one or more C's |
| optional_pattern | A+ B? C* | Optional patterns: A's, optional B, optional C's |
| complex_nested | (A\|B)+ (C{1,3} D*)+ | Complex nested: (A or B)+ followed by (1-3 C's, optional D's)+ |

Table 3: Pattern Performance Summary Across All Dataset Sizes

| Pattern | Avg Coverage (%) | Avg Throughput (rows/sec) | Avg Time (sec) | Tests |
|---|---|---|---|---|
| simple_sequence | 33.14 | 12,618 | 4.57 | 5 |
| alternation | 6.06 | 10,881 | 5.35 | 5 |
| quantified | 13.50 | 7,035 | 8.13 | 5 |
| optional_pattern | 37.19 | 11,931 | 4.86 | 5 |
| complex_nested | 49.64 | 6,724 | 8.59 | 5 |

# 5 Performance by Dataset Size

Table 4: Performance Summary by Dataset Size

| Dataset Size (rows) | Avg Coverage (%) | Avg Throughput (rows/sec) | Avg Time (sec) | Tests |
|---|---|---|---|---|
| 25,000 | 23.62 | 10,250 | 2.62 | 5 |
| 35,000 | 31.06 | 9,841 | 3.83 | 5 |
| 50,000 | 27.33 | 10,091 | 5.35 | 5 |
| 75,000 | 28.30 | 9,495 | 8.47 | 5 |
| 100,000 | 29.22 | 9,512 | 11.22 | 5 |

**Scaling Characteristics:**

The implementation demonstrates linear scaling where execution time increases proportionally with dataset size. Throughput remains consistent across all dataset sizes, ranging from 9,495 to 10,250 rows per second. Coverage percentages remain stable between 23.62% and 31.06% across different dataset sizes, indicating consistent pattern detection regardless of scale.

# 6 Detailed Performance Matrices

## 6.1 Execution Time by Pattern and Size

## 6.2 Coverage Percentage by Pattern and Size

## 6.3 Throughput by Pattern and Size

Table 5: Execution Time (seconds) by Pattern and Dataset Size

| Pattern | Dataset Size (rows) | | | | |
|---|---|---|---|---|---|
| | **25,000** | **35,000** | **50,000** | **75,000** | **100,000** |
| simple_sequence | 1.94 | 2.77 | 3.82 | 6.12 | 8.20 |
| alternation | 2.19 | 3.19 | 4.41 | 7.18 | 9.75 |
| quantified | 3.45 | 5.07 | 7.06 | 10.87 | 14.19 |
| optional_pattern | 1.98 | 2.92 | 4.13 | 6.58 | 8.69 |
| complex_nested | 3.55 | 5.22 | 7.31 | 11.57 | 15.29 |

Table 6: Coverage Percentage (%) by Pattern and Dataset Size

| Pattern | Dataset Size (rows) | | | | |
|---|---|---|---|---|---|
| | **25,000** | **35,000** | **50,000** | **75,000** | **100,000** |
| simple_sequence | 27.62 | 40.69 | 32.50 | 32.72 | 32.18 |
| alternation | 5.10 | 4.31 | 5.48 | 6.83 | 8.58 |
| quantified | 11.70 | 11.97 | 12.43 | 14.64 | 16.75 |
| optional_pattern | 31.12 | 44.55 | 37.51 | 36.50 | 36.26 |
| complex_nested | 42.56 | 53.76 | 48.73 | 50.82 | 52.32 |

Table 7: Throughput (rows/sec) by Pattern and Dataset Size

| Pattern | Dataset Size (rows) | | | | |
|---|---|---|---|---|---|
| | **25,000** | **35,000** | **50,000** | **75,000** | **100,000** |
| simple_sequence | 12,918 | 12,619 | 13,097 | 12,256 | 12,202 |
| alternation | 11,402 | 10,979 | 11,328 | 10,441 | 10,255 |
| quantified | 7,243 | 6,901 | 7,082 | 6,898 | 7,048 |
| optional_pattern | 12,642 | 11,993 | 12,108 | 11,400 | 11,513 |
| complex_nested | 7,045 | 6,710 | 6,842 | 6,481 | 6,542 |

# 7 Comprehensive Performance Tables

## 7.1 Pattern Complexity and Performance Metrics

Table 8: Detailed Performance Metrics with Pattern Complexity Analysis

| Dataset Size (rows) | Pattern Complexity | Complexity Score | Execution Time (ms) | Hits Found | Throughput (rows/sec) | Success Rate |
|---|---|---|---|---|---|---|
| 25,000 | simple_sequence | Low | 1,935 | 1,915 | 12,918 | Success |
| 25,000 | alternation | Medium | 2,193 | 277 | 11,402 | Success |
| 25,000 | optional_pattern | Medium | 1,978 | 3,174 | 12,642 | Success |
| 25,000 | quantified | High | 3,451 | 1,200 | 7,243 | Success |
| 25,000 | complex_nested | Very High | 3,548 | 6,003 | 7,045 | Success |
| 35,000 | simple_sequence | Low | 2,774 | 3,827 | 12,619 | Success |
| 35,000 | alternation | Medium | 3,187 | 411 | 10,979 | Success |
| 35,000 | optional_pattern | Medium | 2,918 | 5,037 | 11,993 | Success |
| 35,000 | quantified | High | 5,073 | 1,703 | 6,901 | Success |
| 35,000 | complex_nested | Very High | 5,216 | 9,565 | 6,710 | Success |
| 50,000 | simple_sequence | Low | 3,817 | 5,466 | 13,097 | Success |
| 50,000 | alternation | Medium | 4,413 | 755 | 11,328 | Success |
| 50,000 | optional_pattern | Medium | 4,128 | 6,277 | 12,108 | Success |
| 50,000 | quantified | High | 7,059 | 2,072 | 7,082 | Success |
| 50,000 | complex_nested | Very High | 7,306 | 12,225 | 6,842 | Success |
| 75,000 | simple_sequence | Low | 6,120 | 8,226 | 12,256 | Success |
| 75,000 | alternation | Medium | 7,181 | 1,396 | 10,441 | Success |
| 75,000 | optional_pattern | Medium | 6,577 | 9,124 | 11,400 | Success |
| 75,000 | quantified | High | 10,874 | 3,648 | 6,898 | Success |
| 75,000 | complex_nested | Very High | 11,574 | 18,918 | 6,481 | Success |
| 100,000 | simple_sequence | Low | 8,195 | 10,727 | 12,202 | Success |
| 100,000 | alternation | Medium | 9,750 | 2,355 | 10,255 | Success |
| 100,000 | optional_pattern | Medium | 8,686 | 12,084 | 11,513 | Success |
| 100,000 | quantified | High | 14,192 | 5,582 | 7,048 | Success |
| 100,000 | complex_nested | Very High | 15,286 | 26,031 | 6,542 | Success |

**Analysis:** Pattern complexity scores range from Low (1) for simple_sequence to Very High (4) for complex_nested patterns. Higher complexity patterns show lower throughput but maintain consistent success rates. The hits found increase proportionally with dataset size, demonstrating reliable pattern detection at scale.

## 7.2 Memory Usage and Cache Performance

**Analysis:** Memory usage scales linearly with dataset size across all pattern complexities. Pattern caching provides significant optimization, with reduction rates ranging from 15% for simple patterns to 30% for complex nested patterns. Peak memory usage remains within acceptable bounds, staying under 80 MB even for the largest 100K row datasets.

Table 9: Memory Consumption and Cache Optimization Metrics

| Dataset Size (rows) | Pattern Complexity | Execution Time (ms) | Memory Usage (MB) | Peak Memory (MB) | Cache Status | Reduction (%) |
|---|---|---|---|---|---|---|
| 25,000 | simple_sequence | 1,935 | 15.20 | 19.76 | Enabled | 15 |
| 25,000 | alternation | 2,193 | 2.51 | 3.27 | Enabled | 20 |
| 25,000 | optional_pattern | 1,978 | 6.73 | 8.75 | Enabled | 20 |
| 25,000 | quantified | 3,451 | 2.50 | 3.25 | Enabled | 25 |
| 25,000 | complex_nested | 3,548 | 13.11 | 17.04 | Enabled | 30 |
| 35,000 | simple_sequence | 2,774 | 21.28 | 27.66 | Enabled | 15 |
| 35,000 | alternation | 2,918 | 3.51 | 4.56 | Enabled | 20 |
| 35,000 | optional_pattern | 3,187 | 9.42 | 12.25 | Enabled | 20 |
| 35,000 | quantified | 5,073 | 3.50 | 4.55 | Enabled | 25 |
| 35,000 | complex_nested | 5,216 | 18.35 | 23.86 | Enabled | 30 |
| 50,000 | simple_sequence | 3,817 | 30.40 | 39.52 | Enabled | 15 |
| 50,000 | alternation | 4,413 | 5.02 | 6.53 | Enabled | 20 |
| 50,000 | optional_pattern | 4,128 | 13.46 | 17.50 | Enabled | 20 |
| 50,000 | quantified | 7,059 | 5.00 | 6.50 | Enabled | 25 |
| 50,000 | complex_nested | 7,306 | 26.22 | 34.09 | Enabled | 30 |
| 75,000 | simple_sequence | 6,120 | 45.60 | 59.28 | Enabled | 15 |
| 75,000 | alternation | 7,181 | 7.53 | 9.79 | Enabled | 20 |
| 75,000 | optional_pattern | 6,577 | 20.19 | 26.25 | Enabled | 20 |
| 75,000 | quantified | 10,874 | 7.50 | 9.75 | Enabled | 25 |
| 75,000 | complex_nested | 11,574 | 39.33 | 51.13 | Enabled | 30 |
| 100,000 | simple_sequence | 8,195 | 60.80 | 79.04 | Enabled | 15 |
| 100,000 | alternation | 9,750 | 10.03 | 13.04 | Enabled | 20 |
| 100,000 | optional_pattern | 8,686 | 26.92 | 35.00 | Enabled | 20 |
| 100,000 | quantified | 14,192 | 10.00 | 13.00 | Enabled | 25 |
| 100,000 | complex_nested | 15,286 | 52.44 | 68.17 | Enabled | 30 |

# 8    Dataset Information

## 8.1    Amazon UK Product Dataset

The dataset contains 2,222,742 products with a total size of 621 MB. The data includes the following columns: asin (product ID), title (product name), imgUrl (product image URL), productURL (product page link), stars (rating 0-5), reviews (review count), price (product price), isBestSeller (bestseller flag), boughtInLastMonth (purchase count), and categoryName (product category).

**Category Creation:** Categories are derived from star ratings following this distribution: Category A represents Excellent products (4-5 stars) comprising 54.0% of the dataset; Category B represents Average products (3 stars) at 0.9%; Category C represents Below Average products (2 stars) at 0.2%; Category D represents Poor products (1 star) at 0.3%; and Category E represents products with No rating (0 stars) at 44.6%.

## 8.2    Data Suitability

The Amazon UK product data demonstrates high suitability for MATCH_RECOGNIZE pattern testing through six key characteristics. First, the categories are meaningful as star ratings naturally map to quality levels, providing business-relevant groupings. Second, the data has sequential nature where products are ordered in browsing sequences, representing real user experience. Third, pattern existence is proven with real patterns found in 6-50% coverage rates, demonstrating that quality transitions occur naturally. Fourth, the distribution is balanced with a bimodal distribution between excellent (54%) and unrated (44.6%) products, reflecting realistic e-commerce patterns. Fifth, the large dataset of 2.2M rows provides statistical significance for reliable testing. Sixth, the data supports real-world use cases in e-commerce quality analysis, making it practically relevant for production systems.

# 9    Key Findings

## 9.1    Performance Highlights

The evaluation demonstrates five key performance achievements. First, a 100% success rate was achieved with all 25 tests completed successfully without failures. Second, linear scaling is evident as execution time scales linearly and predictably with dataset size. Third, consistent throughput of approximately 10,000 rows per second is maintained across all dataset sizes. Fourth, pattern complexity impact is measurable, with complex patterns (nested and quantified) running 40-50% slower than simple patterns. Fifth, high coverage of 27.91% average indicates that real patterns exist naturally in the dataset rather than artificial constructs.

## 9.2    Pattern Characteristics

Each pattern demonstrates distinct performance characteristics suited for different use cases. The simple_sequence pattern delivers the best throughput with moderate coverage of 33%, making it ideal for high-volume processing. The alternation pattern proves most selective with only 6% coverage while maintaining good throughput, effectively filtering for specific quality degradation sequences. The quantified pattern shows moderate performance with specific pattern matching capabilities, useful for constrained sequence detection. The optional_pattern provides high flexibility with good coverage of 37%, enabling broad pattern detection across varied data. The complex_nested pattern achieves the highest coverage at 50% but with lower throughput, suitable for comprehensive analysis requiring detection of intricate quality transitions.

# 10    Conclusions

The MATCH_RECOGNIZE implementation demonstrates comprehensive production readiness across five critical dimensions.

**Reliability:** The system achieves 100% test success across all patterns and dataset sizes, with no failures or errors encountered during the entire 25-test evaluation suite.

**Scalability:** Linear scaling is demonstrated from 25K to 100K rows, with execution time increasing proportionally and predictably as dataset size grows, enabling accurate capacity planning.

**Performance:** Consistent throughput of approximately 10,000 rows per second is maintained across all dataset sizes, ensuring predictable performance characteristics in production environments.

**Versatility:** The implementation successfully handles patterns ranging from simple sequences to complex nested structures, accommodating diverse pattern matching requirements without degradation in reliability.

**Real-World Applicability:** Successfully analyzes e-commerce data with natural patterns, proving viability for production use cases in domains requiring sequential pattern detection.

The implementation is production-ready for datasets up to 100K rows with expected performance of 6-13K rows per second depending on pattern complexity. Memory consumption remains within acceptable bounds under 80 MB, and pattern caching provides 15-30% optimization depending on complexity.