

Experience Report: AhoCorasick Prefiltering Optimization Attempt in rmatch

Debugging and Performance Optimization Analysis

September 9, 2025

Abstract

This report documents a comprehensive attempt to optimize the rmatch regular expression matching library through AhoCorasick prefiltering. Starting from a performance deficit of 44-67% compared to baseline (13.5s target vs 19.7s-22.6s actual), we implemented a complete AhoCorasick prefiltering system including literal extraction, pattern scanning, and candidate filtering. While the implementation achieved technical correctness with all 205 tests passing, it delivered minimal performance improvement. This experience report analyzes the root causes of the optimization's limitations, documents critical debugging challenges encountered, and provides insights for future optimization strategies in regex matching systems.

1 Executive Summary

This report documents a focused optimization attempt targeting the rmatch regular expression matching library's performance bottleneck. The goal was to implement AhoCorasick prefiltering to achieve substantial performance improvement and reach parity with a 13.5-second baseline benchmark.

Key Results:

- Successfully implemented a complete AhoCorasick prefiltering system
- Achieved technical correctness: 205/205 tests passing, mvn clean install successful
- Performance improvement: Minimal to none (target: 37%+ improvement, actual: <5%)
- Root cause identified: Literal-based prefiltering insufficient for high-density word matching patterns

The optimization attempt revealed fundamental limitations in applying literal-based prefiltering to word-dense text matching scenarios, providing valuable insights for future algorithmic approaches.

2 Problem Context and Baseline

2.1 Initial Performance State

The rmatch library exhibited significant performance degradation compared to its historical baseline:

Configuration	Duration (ms)
Baseline Target	13,500
Current Performance Range	19,700 - 22,600
Performance Deficit	44% - 67% slower
Required Improvement	37%+

Table 1: Performance baseline comparison

2.2 Core Algorithmic Challenge

The fundamental bottleneck was identified as an $O(l \times m)$ complexity issue where:

- **l**: Input text length (647K characters from Wuthering Heights corpus)
- **m**: Number of regular expressions (10,000 word patterns)
- **Problem**: Creating Match objects for every regex at every text position

Profiling revealed massive object creation:

- **MatchImpl objects**: 215+ million
- **MatchSetImpl objects**: 8.5+ million
- **Pattern characteristic**: Simple literal words (“the”, “and”, “was”, etc.)

3 AhoCorasick Prefiltering Implementation

3.1 Design Philosophy

AhoCorasick prefiltering was chosen as a promising optimization approach based on its theoretical advantages:

1. **Single-pass scanning**: $O(l)$ text traversal regardless of pattern count
2. **Literal substring detection**: Identify positions where patterns could potentially match
3. **Candidate reduction**: Only test regexes at positions where their required literals are found
4. **Algorithmic improvement**: Transform $O(l \times m)$ to $O(l + k \times m)$ where $k \ll l$

3.2 Implementation Architecture

The implementation consisted of four major components:

3.2.1 Literal Hint Extraction

Enhanced the `LiteralPrefilter.extract()` method with aggressive selectivity scoring:

Listing 1: Selectivity-based literal extraction

```

1 private static double calculateSelectivityScore(
2     final String literal, final boolean anchored) {
3     double score = literal.length() * 10.0;
4
5     if (anchored) {
```

```

6     score *= 3.0; // Anchored literals are highly selective
7 }
8
9 // Character rarity scoring
10 for (char c : literal.toCharArray()) {
11     if ("aeiouAEIOU".indexOf(c) >= 0) {
12         score += 1.0; // Common vowels
13     } else if (Character.isDigit(c)) {
14         score += 5.0; // Digits are selective
15     } else if ("!@#$%^&*()_+!=".indexOf(c) >= 0) {
16         score += 10.0; // Punctuation is highly selective
17     }
18 }
19
20 // Penalty for common English words
21 if (COMMON_WORDS.contains(literal.toLowerCase())) {
22     score *= 0.5;
23 }
24
25 return score;
26 }

```

3.2.2 Pattern Configuration Integration

Added complete integration between `MatcherImpl` and `MatchEngineImpl`:

Listing 2: Automatic prefilter configuration

```

1 public void add(final String r, final Action a)
2     throws RegexpParserException {
3     synchronized (rs) {
4         rs.add(r, a);
5
6         if (!USE_BLOOM_FILTER && me instanceof MatchEngineImpl) {
7             configurePrefilterForLegacyEngine();
8         }
9     }
10 }
11
12 private void configurePrefilterForLegacyEngine() {
13     final Map<Integer, String> patterns = new HashMap<>();
14     final Map<Integer, Integer> flags = new HashMap<>();
15     final Map<String, Regexp> regexpMappings = new HashMap<>();
16
17     int patternId = 0;
18     for (final String regexpStr : rs.getRegexpSet()) {
19         patterns.put(patternId, regexpStr);
20         flags.put(patternId, 0);
21         regexpMappings.put(regexpStr, rs.getRegexp(regexpStr));
22         patternId++;
23     }
24
25     legacyEngine.configurePrefilter(patterns, flags, regexpMappings);
26 }

```

3.2.3 Candidate Position Identification

Implemented efficient text scanning and candidate position mapping:

Listing 3: Prefilter candidate identification

```

1 private void runPrefilterScan(final String text) {
2     final List<AhoCorasickPrefilter.Candidate> candidates =
3         prefilter.scan(text);
4
5     reusableCandidatePositions.clear();
6     reusablePositionToRegexps.clear();
7
8     for (final AhoCorasickPrefilter.Candidate candidate : candidates) {
9         final int startPos = candidate.startIndexForMatch();
10        if (startPos >= 0) {
11            reusableCandidatePositions.add(startPos);
12
13            final Regexp regexp = patternIdToRegexp.get(candidate.patternId);
14            if (regexp != null) {
15                reusablePositionToRegexps
16                    .computeIfAbsent(startPos, k -> new HashSet<>())
17                    .add(regexp);
18            }
19        }
20    }
21
22    candidatePositions = reusableCandidatePositions;
23    positionToRegexps = reusablePositionToRegexps;
24 }

```

3.2.4 Match Engine Integration

Modified the core matching loop to use prefilter candidates:

Listing 4: Prefilter-driven matching

```

1 if (shouldStartMatch) {
2     final DFANode startOfNewMatches = ns.getNextFromStartNode(currentChar)
3         ;
4     if (startOfNewMatches != null) {
5         Set<Regexp> candidateRegexps;
6
7         if (prefilterEnabled && prefilter != null &&
8             positionToRegexps != null &&
9             positionToRegexps.containsKey(currentPos)) {
10            // Use prefilter-specific regexps
11            candidateRegexps = positionToRegexps.get(currentPos);
12        } else {
13            // Fallback to character-based filtering
14            candidateRegexps = startOfNewMatches
15                .getRegexpsThatCanStartWith(currentChar);
16        }
17
18        if (!candidateRegexps.isEmpty()) {
19            final MatchSet ms = new MatchSetImpl(
20                currentPos, startOfNewMatches, currentChar, candidateRegexps);
21            if (ms.hasMatches()) {
22                activeMatchSets.add(ms);
23            }
24        }
25    }
26 }

```

4 Critical Debugging Challenges

4.1 Correctness Regression Discovery

The most critical challenge emerged when basic functionality tests began failing:

Listing 5: Test failure symptoms

```
1 [ERROR] SequenceLoaderTest.testWutheringHeightsCorpusWithVeryFewRegexps
2 AssertionError: Not enough matches, got only 0 but expected at
   least 42.
3
4 INFO: No of characters read: -1 chars from corpus/file.txt
5 INFO: Total no of 'word' matches in Wuthering Heights is 0
```

Root cause analysis revealed:

1. **Overly aggressive prefiltering:** Default enabled prefilter was filtering out legitimate matches
2. **Buffer consumption bug:** `collectBufferText()` was consuming the input buffer, causing downstream failures
3. **Pattern ID mapping failure:** Candidates weren't being mapped back to actual Regexp objects
4. **Test environment mismatch:** Production tests not designed for prefiltering behavior

4.2 Buffer Handling Complexity

A critical implementation bug emerged in buffer handling:

Listing 6: Buffer consumption bug fix

```
1 // BROKEN: This consumed the buffer
2 private String collectBufferText(final Buffer b) {
3     final StringBuilder sb = new StringBuilder();
4     while (b.hasNext()) {
5         sb.append(b.getNext()); // Consuming buffer position!
6     }
7     return sb.toString();
8 }
9
10 // FIXED: Non-consuming access
11 private String collectBufferText(final Buffer b) {
12     if (b instanceof no.rmz.rmatch.utils.StringBuffer) {
13         final StringBuffer sb = (StringBuffer) b;
14         return sb.getCurrentRestString(); // Non-consuming
15     }
16     return null; // Disable prefiltering for other buffer types
17 }
```

This bug caused assertion failures when the buffer position became inconsistent with expected state.

4.3 Configuration Integration Issues

The prefilter was never being configured in the production system:

Listing 7: Missing configuration integration

```

1 // BUG: MatcherImpl.add() never called configurePrefilter()
2 public void add(final String r, final Action a) {
3     synchronized (rs) {
4         rs.add(r, a);
5         // Missing: Configure prefilter for new patterns!
6     }
7 }
8
9 // FIX: Added automatic prefilter configuration
10 public void add(final String r, final Action a) {
11     synchronized (rs) {
12         rs.add(r, a);
13
14         if (!USE_BLOOM_FILTER && me instanceof MatchEngineImpl) {
15             configurePrefilterForLegacyEngine(); // Now configured!
16         }
17     }
18 }

```

5 Performance Results and Analysis

5.1 Benchmark Results

After implementing the complete AhoCorasick prefiltering system:

Configuration	Duration (ms)	vs Baseline
Baseline Target	13,500	-
Without Prefilter	19,700-22,600	44-67% slower
With AhoCorasick Prefilter	20,300-22,400	50-66% slower
Performance Improvement	<5%	Minimal

Table 2: AhoCorasick prefiltering performance results

5.2 Object Creation Analysis

Despite prefiltering, object creation remained massive:

Object Type	Without Prefilter	With Prefilter
MatchImpl Objects	215M+	215M+
MatchSetImpl Objects	8.5M+	8.5M+
Reduction Achieved	-	None

Table 3: Object creation comparison shows no reduction

This indicated that the prefilter was not effectively reducing the workload.

5.3 Root Cause: Pattern Characteristics

Analysis revealed why prefiltering was ineffective:

1. **High literal density:** Common words like “the”, “and”, “was” appear frequently in English text

2. **Minimal filtering:** Most text positions contained at least one common word
3. **Filter ineffectiveness:** `getRegexpsThatCanStartWith()` for common characters returned most of the 10K regexps
4. **Overhead addition:** Prefilter scanning added computational cost without corresponding reduction in regex testing

6 Technical Lessons Learned

6.1 Algorithmic Insights

6.1.1 Prefiltering Effectiveness Depends on Pattern Characteristics

AhoCorasick prefiltering is most effective when:

- Patterns contain distinctive, rare literal substrings
- Text has low density of pattern matches
- Significant reduction in candidate positions is achievable

For word-dense text matching with common English words, prefiltering provides minimal benefit.

6.1.2 $O(l \times m)$ Bottleneck Requires Fundamental Algorithmic Change

Incremental optimizations (prefiltering, collection reuse, character filtering) cannot address the core $O(l \times m)$ complexity. The algorithm fundamentally creates too many objects by testing every pattern at nearly every position.

6.2 Implementation Insights

6.2.1 Integration Complexity

Adding prefiltering to an existing regex engine requires:

- Careful buffer handling to avoid state corruption
- Complete end-to-end integration from pattern registration to match execution
- Backward compatibility preservation
- Comprehensive testing across all usage scenarios

6.2.2 Performance Measurement Challenges

Accurate performance measurement required:

- Multiple benchmark runs to account for JVM warmup
- Careful isolation of prefilter overhead vs. benefit
- Object creation profiling to understand algorithmic impact
- Baseline preservation for regression detection

6.3 Debugging Strategies

6.3.1 Regression Detection

The most critical debugging phase involved detecting and fixing correctness regressions:

1. **Early detection:** Unit tests caught 0-match scenarios immediately
2. **Isolation:** Disabling prefilter by default restored functionality
3. **Incremental fixing:** Addressed buffer handling, configuration, and mapping issues separately
4. **Comprehensive testing:** Ensured 205/205 tests passed before performance evaluation

6.3.2 Root Cause Analysis

Effective debugging required:

- Understanding the full data flow from pattern registration to match execution
- Identifying where object creation actually occurs
- Measuring prefilter effectiveness independently of overall performance
- Analyzing pattern characteristics to understand filtering potential

7 Alternative Approaches and Future Directions

7.1 Fundamental Algorithmic Changes

Based on this experience, future optimization attempts should consider:

7.1.1 Hierarchical Bloom Filtering

- Stage 1: Bloom filter on first character (eliminate 90% of regexps)
- Stage 2: Bloom filter on first 2-3 characters (eliminate 99% of remaining)
- Stage 3: Full regex matching on survivors (10-100 regexps instead of 10K)

7.1.2 Position-Based Regex Indexing

- Pre-compute which regexps can possibly match at each text position
- Use suffix arrays or n-gram indices for efficient lookup
- Only test regexps that match local character patterns

7.1.3 Streaming + Early Termination

- Process text in chunks, terminating unpromising matches early
- Use statistical heuristics to prioritize high-probability regexps
- Implement match caching to avoid recomputing similar patterns

7.2 Implementation Strategy Improvements

Future optimization attempts should:

1. **Measure baseline impact before implementation:** Understand whether the proposed optimization can theoretically deliver required improvements
2. **Implement correctness preservation first:** Ensure all existing functionality continues working before performance optimization
3. **Profile object creation patterns:** Focus on algorithmic changes that reduce object allocation, not just computational complexity
4. **Consider pattern characteristics:** Choose optimization strategies based on actual pattern types and text characteristics

8 Conclusions

This AhoCorasick prefiltering optimization attempt achieved technical success in implementation completeness and correctness preservation, but failed to deliver the required performance improvement. The experience provides valuable insights:

8.1 Key Findings

1. **Pattern characteristics matter:** Literal-based prefiltering is ineffective for high-density word matching scenarios
2. **$O(1 \times m)$ bottleneck is fundamental:** Incremental optimizations cannot address core algorithmic complexity
3. **Implementation correctness is critical:** Regression prevention and comprehensive testing are essential for production systems
4. **Integration complexity is high:** Adding optimization features to existing systems requires careful architectural consideration

8.2 Strategic Implications

For achieving the target 37%+ performance improvement, rmatch requires:

- **Fundamental algorithmic redesign:** Moving beyond incremental optimizations to address $O(1 \times m)$ complexity
- **Pattern-specific optimization strategies:** Choosing techniques based on actual usage patterns rather than theoretical advantages
- **Object allocation reduction:** Focus on minimizing object creation rather than just computational complexity
- **Comprehensive baseline analysis:** Understanding historical performance changes to guide optimization priorities

The AhoCorasick prefiltering infrastructure developed during this optimization attempt remains technically sound and could be valuable for different pattern sets or as part of a larger algorithmic redesign. However, for the current use case of word-dense English text matching, more fundamental approaches are required.

This experience demonstrates the importance of understanding both algorithmic theory and practical implementation constraints when optimizing complex software systems. While AhoCorasick prefiltering is a proven technique in string matching literature, its effectiveness depends critically on problem characteristics that must be carefully analyzed before implementation.