

Design And Analysis Of Algorithms

Bonus Homework Report

Enes ELDES - 22050111017

1 Boyer–Moore Implementation

The Boyer–Moore algorithm was implemented using the *bad character heuristic*. The general logic is familiar with the GeeksForGeeks implementation.

1.1 Unicode-Aware Modification

Standard Boyer–Moore implementations typically assume ASCII characters with a fixed-size table (usually 256). However, during testing, cases involving Unicode characters were encountered.

To solve this, the implementation dynamically determines the maximum character value appearing in both the text and the pattern, and constructs the bad character table accordingly.

This ensures algorithm is reliable, but causes overhead as seen in the results table. **Some test cases that are best case for Boyer Moore are not faster with this implementation.**

1.2 Design Summary

- Uses bad character heuristic only
- Dynamically sized bad character table
- Correct handling of empty patterns and edge cases

2 Pre-Analysis Strategy

The selection strategy is based on the following features:

1. **Pattern Length:** Very short patterns use the Naive algorithm.
2. **Repetition Ratio:** Detected via the LPS array; high repetition favors KMP.
3. **Character Diversity:** High diversity favors Boyer–Moore due to better skipping behavior.
4. **Text Entropy:** Shannon entropy is used to estimate randomness. High entropy favors Rabin–Karp, while low entropy favors KMP.

3 Experimental Results

See the output of the `./test.sh` command in the Figures 1 and 2.

Test Case	RabinKarp (μs)	GoCrazy (μs)	BoyerMoore (μs)	KMP (μs)	Native (μs)	Winner
Simple Match	3.473	N/A	4.438	3.134	3.786	🏆 KMP
No Match	1.157	N/A	1.742	1.431	0.741	🏆 Naive
Single Character	6.863	N/A	6.777	6.186	6.758	🏆 KMP
Pattern at End	5.808	N/A	3.077	1.543	1.002	🏆 Naive
Pattern at Beginning	2.421	N/A	2.140	1.423	1.074	🏆 Naive
Overlapping Patterns	2.802	N/A	3.055	1.825	2.102	🏆 KMP
Long Text Multiple Matches	5.738	N/A	6.229	7.279	4.736	🏆 Naive
Pattern Longer Than Text	0.160	N/A	0.154	0.677	0.190	🏆 BoyerMoore
Entire Text Match	1.270	N/A	2.054	1.374	1.063	🏆 Naive
Repeating Pattern	3.483	N/A	3.794	2.307	2.045	🏆 Naive
Case Sensitive	2.561	N/A	2.878	2.400	1.768	🏆 Naive
Numbers and Special Characters	3.810	N/A	3.565	3.540	2.287	🏆 Naive
Unicode Characters	7.004	N/A	130.285	2.668	1.651	🏆 Naive
Very Long Text	39.241	N/A	30.738	41.983	28.691	🏆 Naive
Pattern with Spaces	4.026	N/A	3.796	3.570	2.625	🏆 Naive
All Same Character	14.056	N/A	14.122	6.874	13.418	🏆 KMP
Alternating Pattern	21.103	N/A	20.233	12.492	17.430	🏆 KMP
Long Pattern	11.901	N/A	10.021	12.272	7.852	🏆 Naive
Pattern at Boundaries	2.535	N/A	3.058	3.374	1.934	🏆 Naive
Near Matches	1.778	N/A	3.664	1.978	1.652	🏆 Naive
Empty Pattern	0.868	N/A	0.822	1.074	1.131	🏆 BoyerMoore
Empty Text	0.160	N/A	0.163	0.403	0.179	🏆 RabinKarp
Both Empty	0.527	N/A	0.316	0.542	0.326	🏆 BoyerMoore
Single Character Pattern	7.551	N/A	7.363	6.598	6.191	🏆 Naive
Complex Overlap	3.156	N/A	3.855	2.357	2.713	🏆 KMP
DNA Sequence	10.977	N/A	9.100	8.131	9.933	🏆 KMP
Palindrome Pattern	6.389	N/A	5.141	6.409	4.171	🏆 Naive
Worst Case for Naive	6.859	N/A	12.809	8.849	6.190	🏆 Naive
Best Case for Boyer-Moore	3.119	N/A	5.957	6.594	1.477	🏆 Naive
KMP Advantage Case	1.535	N/A	2.500	8.588	4.395	🏆 RabinKarp

Figure 1: Execution time comparison of string matching algorithms across test cases

Test Case	Choice	PreA+Choice (μs)	vs Naive	vs RabinKarp	vs BoyerMoore	vs KMP
Simple Match	Naive		1.78 N/A	+0.72 μs	-1.76 μs	+0.53 μs
No Match	Naive	0.41	N/A	+0.04 μs	-0.93 μs	-0.28 μs
Single Character	Naive		1.59 N/A	-0.02 μs	-1.59 μs	+0.15 μs
Pattern at End	BoyerMoore		3.62 +2.98 μs	+2.86 μs	N/A	+2.80 μs
Pattern at Beginning	BoyerMoore		2.79 +2.45 μs	+2.32 μs	N/A	+1.73 μs
Overlapping Patterns	Naive		1.03 N/A	+0.02 μs	-0.67 μs	-0.11 μs
Long Text Multiple Matches	Naive		1.47 N/A	-0.35 μs	-1.66 μs	-1.15 μs
Pattern Longer Than Text	BoyerMoore		1.40 +1.28 μs	+1.29 μs	N/A	+0.89 μs
Entire Text Match	BoyerMoore		2.89 +2.59 μs	+2.45 μs	N/A	+1.89 μs
Repeating Pattern	Naive		1.08 N/A	-0.12 μs	-0.75 μs	-0.07 μs
Case Sensitive	Naive		0.75 N/A	-0.25 μs	-0.97 μs	-0.45 μs
Numbers and Special Characters	BoyerMoore		3.53 +2.79 μs	+2.17 μs	N/A	+2.26 μs
Unicode Characters	Naive		1.10 N/A	+0.33 μs	-49.09 μs	-0.04 μs
Very Long Text	Naive		9.42 N/A	-1.27 μs	-2.26 μs	-4.19 μs
Pattern with Spaces	BoyerMoore		3.65 +2.54 μs	+2.64 μs	N/A	+2.27 μs
All Same Character	KMP		4.77 -0.22 μs	-0.42 μs	-1.96 μs	N/A
Alternating Pattern	Naive		9.64 N/A	-1.78 μs	-0.58 μs	+1.97 μs
Long Pattern	BoyerMoore		7.99 +4.34 μs	+4.49 μs	N/A	+3.25 μs
Pattern at Boundaries	Naive		1.13 N/A	-0.46 μs	-1.38 μs	-0.65 μs
Near Matches	Naive		1.08 N/A	+0.25 μs	-0.62 μs	-1.75 μs
Empty Pattern	Naive		1.21 N/A	+0.50 μs	+0.53 μs	+0.51 μs
Empty Text	Naive		0.44 N/A	+0.26 μs	+0.27 μs	-0.05 μs
Both Empty	Naive		0.51 N/A	-0.22 μs	+0.19 μs	-0.18 μs
Single Character Pattern	Naive		4.44 N/A	-0.10 μs	-1.05 μs	-0.47 μs
Complex Overlap	KMP		2.70 +0.95 μs	+0.92 μs	-0.06 μs	N/A
DNA Sequence	Naive		4.87 N/A	-0.38 μs	-1.34 μs	-0.56 μs
Palindrome Pattern	RabinKarp		17.25 +15.78 μs	N/A	+14.34 μs	+14.17 μs
Worst Case for Naive	KMP		17.64 +8.22 μs	+15.30 μs	+12.63 μs	N/A
Best Case for Boyer-Moore	Naive		2.45 N/A	+0.11 μs	-1.01 μs	-0.58 μs
KMP Advantage Case	KMP		19.69 +13.63 μs	+17.31 μs	+16.25 μs	N/A

Figure 2: Execution time comparison with pre-analysis strategy (green: pre-analysis was faster, red: pre-analysis was slower).

3.1 Key Observations

- KMP dominates in highly repetitive patterns and structured data (e.g., DNA sequences).
- Boyer–Moore excels when character diversity is high and mismatches occur early.
- Rabin–Karp performs best in high-entropy texts but is sensitive to hashing overhead.

- Pre-analysis algorithm is not perfectly designed since overhead (red labels) is seen in most of the test cases in the Figure 2. It sometimes selects the inappropriate algorithm, too.

References

- GeeksForGeeks, *Boyer Moore Algorithm for Pattern Searching*. <https://www.geeksforgeeks.org/dsa/boyer-moore-algorithm-for-pattern-searching/>