# Boyer–Moore string-search algorithm

Anders Søberg <andersob@ifi.uio.no>
Martin Mihle Nygaard <martimn@ifi.uio.no>

IN4120 2022 Science Fair

# Boyer–Moore string-search [1]

- Proposed by Boyer and Moore in 1977.
- Standard benchmark for string-search
- E.g used by several `grep` implementations

# Proposition: Naïve string-search

Brute-force all alignments of a pattern $P$ in text $T$.

# Proposition: Naïve string-search

Brute-force all alignments of a pattern $P$ in text $T$.

## Example

$T$ = "EKSEMPEL" and $P$ = "EMP" gives
$n - m + 1 = 6$ alignments to try (where
$n$ and $m$ is the length of $T$ and $P$)

# Proposition: Naïve string-search

Brute-force all alignments of a pattern $P$ in text $T$.

## Example

$T =$ "EKSEMPEL" and $P =$ "EMP" gives $n - m + 1 = 6$ alignments to try (where $n$ and $m$ is the length of $T$ and $P$)

E K S E M P E L

# Proposition: Naïve string-search

Brute-force all alignments of a pattern $P$ in text $T$.

## Example

$T =$ "EKSEMPEL" and $P =$ "EMP" gives $n - m + 1 = 6$ alignments to try (where $n$ and $m$ is the length of $T$ and $P$)

```
E K S E M P E L
E M P - - - - -
```

# Proposition: Naïve string-search

Brute-force all alignments of a pattern $P$ in text $T$.

## Example

$T =$ "EKSEMPEL" and $P =$ "EMP" gives
$n - m + 1 = 6$ alignments to try (where
$n$ and $m$ is the length of $T$ and $P$)

```
E K S E M P E L
E M P - - - - -
- E M P - - - -
```

# Proposition: Naïve string-search

Brute-force all alignments of a pattern $P$ in text $T$.

## Example

$T = $ "EKSEMPEL" and $P = $ "EMP" gives
$n - m + 1 = 6$ alignments to try (where
$n$ and $m$ is the length of $T$ and $P$)

```
E K S E M P E L
E M P - - - - -
- E M P - - - -
- - E M P - - -
```

# Proposition: Naïve string-search

Brute-force all alignments of a pattern $P$ in text $T$.

## Example

$T =$ "EKSEMPEL" and $P =$ "EMP" gives
$n - m + 1 = 6$ alignments to try (where
$n$ and $m$ is the length of $T$ and $P$)

```
E K S E M P E L
E M P - - - - -
- E M P - - - -
- - E M P - - -
- - - E M P - -
```

# Proposition: Naïve string-search

Brute-force all alignments of a pattern $P$ in text $T$.

### Example

$T =$ "EKSEMPEL" and $P =$ "EMP" gives $n - m + 1 = 6$ alignments to try (where $n$ and $m$ is the length of $T$ and $P$)

```
E K S E M P E L
E M P - - - - -
- E M P - - - -
- - E M P - - -
- - - E M P - -
- - - - E M P -
```

# Proposition: Naïve string-search

Brute-force all alignments of a pattern $P$ in text $T$.

## Example

$T =$ "EKSEMPEL" and $P =$ "EMP" gives $n - m + 1 = 6$ alignments to try (where $n$ and $m$ is the length of $T$ and $P$)

```
E K S E M P E L
E M P - - - - -
- E M P - - - -
- - E M P - - -
- - - E M P - -
- - - - E M P -
- - - - - E M P
```

# Improvement

Skip as many alignments as possible.

# Improvement

Skip as many alignments as possible.

### Idea
- ▶ Compare the last aligned characters in $P$ to $T$, and shift along $T$ based on result.

# Improvement

Skip as many alignments as possible.

## Idea

- Compare the last aligned characters in $P$ to $T$, and shift along $T$ based on result.
- Length of shift based on characters and suffixes of $P$.

# Improvement

Skip as many alignments as possible.

## Idea
- ▶ Compare the last aligned characters in $P$ to $T$, and shift along $T$ based on result.
- ▶ Length of shift based on characters and suffixes of $P$.
- ▶ Do some preprocessing of $P$.

# Shift rules

The bad character rule ($delta_1$)

# Shift rules

## The bad character rule ($delta_1$)

If aligned character $\in P$:

- ▶ Shift $P$ so that next occurrence of the character is aligned.

```
E K S E M P E L
S E M P - - - -
- - S E M P - -
```

# Shift rules

## The bad character rule ($delta_1$)

If aligned character $\in P$:
- ▶ Shift $P$ so that next occurrence of the character is aligned.

```
E K S E M P E L
S E M P - - - -
- - S E M P - -
```

Else, aligned character $\notin P$:
- ▶ Shift $P$ by it's own length, $m$.

```
E K S E M P E L
E M P - - - - -
- - - E M P - -
```

# Shift rules

## The good suffix rule ($delta_2$)

Suppose for an alignment, $S$ and $P$ shares a suffix $t$.

▶ If a copy of $t$ exists in $P$, shift to it. (Case 1)

# Shift rules

## The good suffix rule ($delta_2$)

Suppose for an alignment, $S$ and $P$ shares a suffix $t$.
- If a copy of $t$ exists in $P$, shift to it. (Case 1)

## Example

```
        ... FLÅKLYPA-GRAND-PRIX-SUFFIX-SUFFIX ...
Case 1:              FIX-SUFFIX
```

# Shift rules

## The good suffix rule ($delta_2$)

Suppose for an alignment, $S$ and $P$ shares a suffix $t$.
- ▶ If a copy of $t$ exists in $P$, shift to it. (Case 1)

## Example

```
        ... FLÅKLYPA-GRAND-PRIX-SUFFIX-SUFFIX ...
Case 1:             FIX-SUFFIX
                         FIX-SUFFIX
```

# Shift rules

## The good suffix rule ($delta_2$)

Suppose for an alignment, $S$ and $P$ shares a suffix $t$.
- If a copy of $t$ exists in $P$, shift to it. (Case 1)
- Else if, a prefix of $P$ matches a suffix of $t$, shift to it. (Case 2)

## Example

```
        ... FLÅKLYPA-GRAND-PRIX-SUFFIX-SUFFIX ...
Case 1:             FIX-SUFFIX
Case 2:                    FIX-SUFFIX
```

# Shift rules

## The good suffix rule ($delta_2$)

Suppose for an alignment, $S$ and $P$ shares a suffix $t$.
- If a copy of $t$ exists in $P$, shift to it. (Case 1)
- Else if, a prefix of $P$ matches a suffix of $t$, shift to it. (Case 2)

## Example

```
        ... FLÅKLYPA-GRAND-PRIX-SUFFIX-SUFFIX ...
Case 1:             FIX-SUFFIX
Case 2:                     FIX-SUFFIX
                                    FIX-SUFFIX
```

# Shift rules

## The good suffix rule ($delta_2$)

Suppose for an alignment, $S$ and $P$ shares a suffix $t$.

- If a copy of $t$ exists in $P$, shift to it. (Case 1)
- Else if, a prefix of $P$ matches a suffix of $t$, shift to it. (Case 2)
- Else, shift $P$ past $t$.

## Example

```
        ... FLÅKLYPA-GRAND-PRIX-SUFFIX-SUFFIX ...
Case 1:            FIX-SUFFIX
Case 2:                    FIX-SUFFIX
                                   FIX-SUFFIX
```

# Shift rules

## The Galil rule [2]

- ▶ Proposed by Galil in 1979.
- ▶ Improvement by making the comparisons faster.
- ▶ Speeds up multi-matching, with *periodic* patterns.
- ▶ Improves worst cases significantly.

# Shift rules

## The Galil rule [2]

- ▶ Proposed by Galil in 1979.
- ▶ Improvement by making the comparisons faster.
- ▶ Speeds up multi-matching, with *periodic* patterns.
- ▶ Improves worst cases significantly.

## Example

$T = $ "KEBABABA" and $P = $ "BABA", without and with Galil rule.

# Shift rules

## The Galil rule [2]

▶ Proposed by Galil in 1979.
▶ Improvement by making the comparisons faster.
▶ Speeds up multi-matching, with *periodic* patterns.
▶ Improves worst cases significantly.

## Example

$T$ = "KEBABABA" and $P$ = "BABA", without and with Galil rule.

```
          ↓
    K E B A B A B A
    B A B A - - - -
```

# Shift rules

## The Galil rule [2]

▶ Proposed by Galil in 1979.
▶ Improvement by making the comparisons faster.
▶ Speeds up multi-matching, with *periodic* patterns.
▶ Improves worst cases significantly.

## Example

$T = $ "KEBABABA" and $P = $ "BABA", without and with Galil rule.

```
        ↓ ↓
  K E B A B A B A
  B A B A - - - -
```

# Shift rules

## The Galil rule [2]

- ▶ Proposed by Galil in 1979.
- ▶ Improvement by making the comparisons faster.
- ▶ Speeds up multi-matching, with *periodic* patterns.
- ▶ Improves worst cases significantly.

## Example

$T$ = "KEBABABA" and $P$ = "BABA", without and with Galil rule.

```
    ↓ ↓ ↓
K E B A B A B A
B A B A - - - -
```

# Shift rules

## The Galil rule [2]

- ▶ Proposed by Galil in 1979.
- ▶ Improvement by making the comparisons faster.
- ▶ Speeds up multi-matching, with *periodic* patterns.
- ▶ Improves worst cases significantly.

## Example

$T = $ "KEBABABA" and $P = $ "BABA", without and with Galil rule.

```
    ↓ ↓ ↓   ↓
  K E B A B A B A
  B A B A - - - -
  - - B A B A - -
```

# Shift rules

## The Galil rule [2]

▶ Proposed by Galil in 1979.
▶ Improvement by making the comparisons faster.
▶ Speeds up multi-matching, with *periodic* patterns.
▶ Improves worst cases significantly.

## Example

$T =$ "KEBABABA" and $P =$ "BABA", without and with Galil rule.

```
     ↓ ↓ ↓ ↓ ↓
 K E B A B A B A
 B A B A - - - -
 - - B A B A - -
```

# Shift rules

## The Galil rule [2]

▶ Proposed by Galil in 1979.
▶ Improvement by making the comparisons faster.
▶ Speeds up multi-matching, with *periodic* patterns.
▶ Improves worst cases significantly.

## Example

$T =$ "KEBABABA" and $P =$ "BABA", without and with Galil rule.

```
            ↓
    ↓  ↓  ↓  ↓  ↓
    K  E  B  A  B  A  B  A
    B  A  B  A  -  -  -  -
    -  -  B  A  B  A  -  -
```

# Shift rules

## The Galil rule [2]

▶ Proposed by Galil in 1979.
▶ Improvement by making the comparisons faster.
▶ Speeds up multi-matching, with *periodic* patterns.
▶ Improves worst cases significantly.

## Example

$T =$ "KEBABABA" and $P =$ "BABA", without and with Galil rule.

```
        ↓ ↓
      ↓ ↓ ↓ ↓ ↓
    K E B A B A B A
    B A B A - - - -
    - - B A B A - -
```

# Shift rules

## The Galil rule [2]
- ▶ Proposed by Galil in 1979.
- ▶ Improvement by making the comparisons faster.
- ▶ Speeds up multi-matching, with *periodic* patterns.
- ▶ Improves worst cases significantly.

## Example

$T =$ "KEBABABA" and $P =$ "BABA", without and with Galil rule.

```
         ↓ ↓
     ↓ ↓ ↓ ↓ ↓     ↓
   K E B A B A B A
   B A B A - - - -
   - - B A B A - -
   - - - - B A B A
```

# Shift rules

## The Galil rule [2]

- ▶ Proposed by Galil in 1979.
- ▶ Improvement by making the comparisons faster.
- ▶ Speeds up multi-matching, with *periodic* patterns.
- ▶ Improves worst cases significantly.

## Example

$T$ = "KEBABABA" and $P$ = "BABA", without and with Galil rule.

```
        ↓ ↓
     ↓ ↓ ↓ ↓ ↓ ↓ ↓
    K E B A B A B A
    B A B A - - - -
    - - B A B A - -
    - - - - B A B A
```

# Shift rules

## The Galil rule [2]

► Proposed by Galil in 1979.
► Improvement by making the comparisons faster.
► Speeds up multi-matching, with *periodic* patterns.
► Improves worst cases significantly.

## Example

$T$ = "KEBABABA" and $P$ = "BABA", without and with Galil rule.

```
        ↓ ↓     ↓
      ↓ ↓ ↓ ↓ ↓ ↓ ↓
    K E B A B A B A
    B A B A - - - -
    - - B A B A - -
    - - - - B A B A
```

# Shift rules

## The Galil rule [2]

- ▶ Proposed by Galil in 1979.
- ▶ Improvement by making the comparisons faster.
- ▶ Speeds up multi-matching, with *periodic* patterns.
- ▶ Improves worst cases significantly.

## Example

$T =$ "KEBABABA" and $P =$ "BABA", without and with Galil rule.

```
      ↓ ↓ ↓ ↓
    ↓ ↓ ↓ ↓ ↓ ↓ ↓
  K E B A B A B A
  B A B A - - - -
  - - B A B A - -
  - - - - B A B A
      11 comps.
```

# Shift rules

## The Galil rule [2]

- ▶ Proposed by Galil in 1979.
- ▶ Improvement by making the comparisons faster.
- ▶ Speeds up multi-matching, with *periodic* patterns.
- ▶ Improves worst cases significantly.

## Example

$T =$ "KEBABABA" and $P =$ "BABA", without and with Galil rule.

```
      ↓ ↓ ↓ ↓
    ↓ ↓ ↓ ↓ ↓ ↓ ↓                      ↓
    K E B A B A B A          K E B A B A B A
    B A B A - - - -          B A B A - - - -
    - - B A B A - -
    - - - - B A B A
        11 comps.
```

# Shift rules

## The Galil rule [2]

- ▶ Proposed by Galil in 1979.
- ▶ Improvement by making the comparisons faster.
- ▶ Speeds up multi-matching, with *periodic* patterns.
- ▶ Improves worst cases significantly.

## Example

$T$ = "KEBABABA" and $P$ = "BABA", without and with Galil rule.

```
        ↓ ↓ ↓ ↓
      ↓ ↓ ↓ ↓ ↓ ↓ ↓                      ↓ ↓
      K E B A B A B A              K E B A B A B A
      B A B A - - - -              B A B A - - - -
      - - B A B A - -
      - - - - B A B A
         11 comps.
```

# Shift rules

## The Galil rule [2]

▶ Proposed by Galil in 1979.
▶ Improvement by making the comparisons faster.
▶ Speeds up multi-matching, with *periodic* patterns.
▶ Improves worst cases significantly.

## Example

$T$ = "KEBABABA" and $P$ = "BABA", without and with Galil rule.

```
      ↓ ↓ ↓ ↓
    ↓ ↓ ↓ ↓ ↓ ↓ ↓                    ↓ ↓ ↓
    K E B A B A B A          K E B A B A B A
    B A B A - - - -          B A B A - - - -
    - - B A B A - -
    - - - - B A B A
       11 comps.
```

# Shift rules

## The Galil rule [2]

► Proposed by Galil in 1979.
► Improvement by making the comparisons faster.
► Speeds up multi-matching, with *periodic* patterns.
► Improves worst cases significantly.

## Example

$T$ = "KEBABABA" and $P$ = "BABA", without and with Galil rule.

```
      ↓ ↓ ↓ ↓
    ↓ ↓ ↓ ↓ ↓ ↓ ↓              ↓ ↓ ↓     ↓
    K E B A B A B A            K E B A B A B A
    B A B A - - - -            B A B A - - - -
    - - B A B A - -            - - B A B A - -
    - - - - B A B A
         11 comps.
```

# Shift rules

## The Galil rule [2]
- ▶ Proposed by Galil in 1979.
- ▶ Improvement by making the comparisons faster.
- ▶ Speeds up multi-matching, with *periodic* patterns.
- ▶ Improves worst cases significantly.

## Example
$T$ = "KEBABABA" and $P$ = "BABA", without and with Galil rule.

```
        ↓ ↓ ↓ ↓
      ↓ ↓ ↓ ↓ ↓ ↓ ↓                    ↓ ↓ ↓ ↓ ↓
      K E B A B A B A              K E B A B A B A
      B A B A - - - -              B A B A - - - -
      - - B A B A - -              - - B A B A - -
      - - - - B A B A
           11 comps.
```

# Shift rules

## The Galil rule [2]
- ▶ Proposed by Galil in 1979.
- ▶ Improvement by making the comparisons faster.
- ▶ Speeds up multi-matching, with *periodic* patterns.
- ▶ Improves worst cases significantly.

## Example

$T =$ "KEBABABA" and $P =$ "BABA", without and with Galil rule.

```
      ↓ ↓ ↓ ↓                        ↓
    ↓ ↓ ↓ ↓ ↓ ↓ ↓              ↓ ↓ ↓ ↓ ↓
    K E B A B A B A            K E B A B A B A
    B A B A - - - -            B A B A - - - -
    - - B A B A - -            - - B A B A - -
    - - - - B A B A
         11 comps.
```

# Shift rules

## The Galil rule [2]

- ▶ Proposed by Galil in 1979.
- ▶ Improvement by making the comparisons faster.
- ▶ Speeds up multi-matching, with *periodic* patterns.
- ▶ Improves worst cases significantly.

## Example

$T =$ "KEBABABA" and $P =$ "BABA", without and with Galil rule.

```
        ↓ ↓ ↓ ↓                    ↓ ↓
      ↓ ↓ ↓ ↓ ↓ ↓                ↓ ↓ ↓ ↓
      K E B A B A B A            K E B A B A B A
      B A B A - - - -            B A B A - - - -
      - - B A B A - -            - - B A B A - -
      - - - - B A B A
          11 comps.
```

# Shift rules

## The Galil rule [2]

▶ Proposed by Galil in 1979.
▶ Improvement by making the comparisons faster.
▶ Speeds up multi-matching, with *periodic* patterns.
▶ Improves worst cases significantly.

## Example

$T$ = "KEBABABA" and $P$ = "BABA", without and with Galil rule.

```
      ↓ ↓ ↓ ↓                        ↓ ↓
    ↓ ↓ ↓ ↓ ↓ ↓ ↓              ↓ ↓ ↓ ↓ ↓    ↓
    K E B A B A B A            K E B A B A B A
    B A B A - - - -            B A B A - - - -
    - - B A B A - -            - - B A B A - -
    - - - - B A B A            - - - - B A B A
         11 comps.
```

# Shift rules

## The Galil rule [2]

▶ Proposed by Galil in 1979.
▶ Improvement by making the comparisons faster.
▶ Speeds up multi-matching, with *periodic* patterns.
▶ Improves worst cases significantly.

## Example

$T$ = "KEBABABA" and $P$ = "BABA", without and with Galil rule.

```
        ↓ ↓ ↓ ↓                    ↓ ↓
      ↓ ↓ ↓ ↓ ↓ ↓ ↓            ↓ ↓ ↓ ↓ ↓ ↓ ↓
      K E B A B A B A          K E B A B A B A
      B A B A - - - -          B A B A - - - -
      - - B A B A - -          - - B A B A - -
      - - - - B A B A          - - - - B A B A
          11 comps.                  9 comps.
```

# Performance

- Best case performance is $O(n/m)$

# Performance

- Best case performance is $O(n/m)$
  - When you can skip $m$ elements at every step

# Performance

- Best case performance is $O(n/m)$
  - When you can skip $m$ elements at every step
- Worst case – two cases

# Performance

- Best case performance is $O(n/m)$
  - When you can skip $m$ elements at every step
- Worst case – two cases
  - If the pattern does not appear in the text, the worst case is $O(m + n)$

# Performance

- Best case performance is $O(n/m)$
  - When you can skip $m$ elements at every step
- Worst case – two cases
  - If the pattern does not appear in the text, the worst case is $O(m + n)$
  - Proven in "Fast Pattern Matching in Strings"

# Performance

▶ Worst case for the original algorithm when the pattern does appear in the text is $O(mn)$.

# Performance

- Worst case for the original algorithm when the pattern does appear in the text is $O(mn)$.

- This can be seen by considering the case where the pattern and text both contain the same letter repeatedly:

# Performance

- Worst case for the original algorithm when the pattern does appear in the text is $O(mn)$.

- This can be seen by considering the case where the pattern and text both contain the same letter repeatedly:

```
A A A A A A A A
A A A - - - - -
```

# Performance

▶ Worst case for the original algorithm when the pattern does appear in the text is $O(mn)$.

▶ This can be seen by considering the case where the pattern and text both contain the same letter repeatedly:

<pre>
A A A A A A A A
A A A - - - - -
</pre>

# Performance

- Worst case for the original algorithm when the pattern does appear in the text is $O(mn)$.

- This can be seen by considering the case where the pattern and text both contain the same letter repeatedly:

```
A A A A A A A A
A A A - - - - -
```

# Performance

- Worst case for the original algorithm when the pattern does appear in the text is $O(mn)$.

- This can be seen by considering the case where the pattern and text both contain the same letter repeatedly:

```
A A A A A A A
A A A - - - -
```

# Performance

- Worst case for the original algorithm when the pattern does appear in the text is $O(mn)$.

- This can be seen by considering the case where the pattern and text both contain the same letter repeatedly:

```
A A A A A A A A
A A A – – – – –
– A A A – – – –
```

# Performance

▶ Worst case for the original algorithm when the pattern does appear in the text is $O(mn)$.

▶ This can be seen by considering the case where the pattern and text both contain the same letter repeatedly:

```
A A A A A A A A
A A A - - - - -
- A A A - - - -
- - A A A - - -
```

# Performance

▶ Worst case for the original algorithm when the pattern does appear in the text is $O(mn)$.

▶ This can be seen by considering the case where the pattern and text both contain the same letter repeatedly:

```
A A A A A A A A
A A A - - - - -
- A A A - - - -
- - A A A - - -
- - - A A A - -
```

# Performance

- Worst case for the original algorithm when the pattern does appear in the text is $O(mn)$.

- This can be seen by considering the case where the pattern and text both contain the same letter repeatedly:

```
A A A A A A A A
A A A - - - - -
- A A A - - - -
- - A A A - - -
- - - A A A - -
- - - - A A A -
```

# Performance

- Worst case for the original algorithm when the pattern does appear in the text is $O(mn)$.

- This can be seen by considering the case where the pattern and text both contain the same letter repeatedly:

```
A A A A A A A A
A A A - - - - -
- A A A - - - -
- - A A A - - -
- - - A A A - -
- - - - A A A -
- - - - - A A A
```

# Performance

▶ Worst case for the original algorithm when the pattern does appear in the text is $O(mn)$.

▶ This can be seen by considering the case where the pattern and text both contain the same letter repeatedly:

▶ With Galil rule, linear in general

```
A A A A A A A A
A A A - - - - -
- A A A - - - -
- - A A A - - -
- - - A A A - -
- - - - A A A -
- - - - - A A A
```

# Performance

- Preprocessing:
  - Varies depending on implementation.
  - Usually $O(m)$ or $O(km)$, where $k$ is the size of the alphabet
  - Space complexity $O(k)$

A practical simplification: just use the 'bad character rule'.

- Worst case is $O(nm)$, but average case is $O(n)$
- Trivial to implement.

A practical simplification: just use the 'bad character rule'.

- Worst case is $O(nm)$, but average case is $O(n)$
- Trivial to implement.
  - See implementation

# References

[1] Robert S. Boyer and J. Strother Moore. "A Fast String Searching Algorithm". In: *Commun. ACM* 20.10 (Oct. 1977), pp. 762–772. ISSN: 0001-0782. DOI: 10.1145/359842.359859. URL: https://doi.org/10.1145/359842.359859.

[2] Zvi Galil. "On Improving the Worst Case Running Time of the Boyer-Moore String Matching Algorithm". In: *Commun. ACM* 22.9 (Sept. 1979), pp. 505–508. ISSN: 0001-0782. DOI: 10.1145/359146.359148. URL: https://doi.org/10.1145/359146.359148.

[3] R. Nigel Horspool. "Practical fast searching in strings". In: *Software: Practice and Experience* 10.6 (1980), pp. 501–506. DOI: 10.1002/spe.4380100608. URL: https://onlinelibrary.wiley.com/doi/abs/10.1002/spe.4380100608.

[4] Donald E. Knuth, James H. Morris Jr., and Vaughan R. Pratt. "Fast Pattern Matching in Strings". In: *SIAM Journal on Computing* 6.2 (1977), pp. 323–350. DOI: 10.1137/0206024. URL: https://doi.org/10.1137/0206024.

[5] Wikipedia. *Boyer–Moore–Horspool algorithm — Wikipedia, The Free Encyclopedia*. Nov. 7, 2022. URL: https://en.wikipedia.org/w/index.php?title=Boyer%E2%80%93Moore%E2%80%93Horspool_algorithm.