



EXPERIMENTAL COMPARISON STRING MATCHING ALGORITHMS

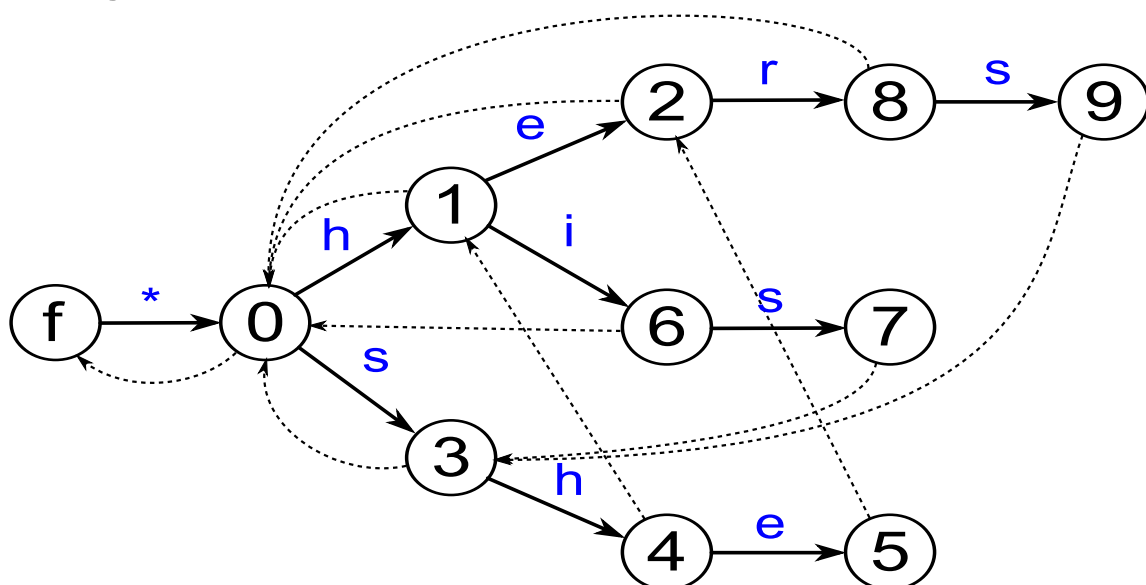
The *multiple exact string matching* problem asks us to search a text T for each pattern in \mathcal{P} using some exact string matching algorithm. In this paper, we give an overview of three such algorithms: Aho-Corasick, Shift-And, and Karabin. We compare the theoretical performance of these algorithms to experimental results. We assume a constant alphabet size of σ and equal-length patterns, $|P_j| = l$ for every j .

THE ALGORITHMS

AHO-CORASICK

The algorithm. We can construct a trie from the patterns \mathcal{P} , then add a link to the node that represents the longest proper suffix of the node and make note of pattern occurrences for each node. The automaton can be then simulated with the text T as input to find occurrences.

Aho-Corasick automaton. For the patterns $\mathcal{P} = \{\text{he, she, his, hers}\}$ we can produce the following automaton:



SHIFT-AND

The algorithm. Shift-And maintains a bitvector D with the longest prefixes of patterns.

```
# preprocessing
for (i = 0; i < m; i++)
    B[P[i/l][i%l]] += 1
for (i = 0; i < m; i++)
    pBegin += 2i
for (i = l - 1; i < l; i++)
    pEnd += 2i
# search
for (i = 0; i < n; i++)
    D = ((D << 1) | B[P[i/l][i%l]])
    if D & pEnd ≠ 0:
```

Bitparallelism. If w is a word, the bitvectors D , p are represented in $\lceil m/w \rceil$ words stored in $\sigma \lceil m/w \rceil$ words.

ON OF MULTIPLE EXACT THMS

Paul Saikko

of length n for a set of patterns \mathcal{P} with total length $\|\mathcal{P}\| = m$. We could
algorithm, but much better algorithms exist for this task. Here we give
p-Rabin. Each provides a different approach to solving the problem. We
l results and identify some of their strenghts and weaknesses, assuming
ry $P_j \in \mathcal{P}$.

KARP-RABIN

and scans the text, it main-
which keeps track of the
rns it has encountered.

Karp-Rabin hash function. For some fixed pos-
itive integers r and q , the karp-rabin hash of a
string $S = s_0s_1 \dots s_{m-1}$ is

$$H(S) = \sum_{i=0}^{m-1} (s_i r^{m-1-i}) \mod q.$$

```
i++) :  
= 2i  
i += 1) :  
< m; i += 1) :
```

This is an example of a rolling hash function – if
we know the hash $H(T_{[i \dots i+m]})$, we can compute
 $H(T_{[i+1 \dots i+1+m]})$ in constant time.

```
i++) :  
oBegin) & B[T[i]]  
yield i
```

the length of a machine
oBegin, and pEnd can be
machine words. B can be
s and D can be updated

Multiple string matching. We can precompute
the hash value for every pattern and store them
in a data structure that supports constant-time
lookups. Potential occurrences in the text can
be found by computing $H(T_{[i \dots i+l]})$ for each $i \in$
 $[0 \dots n - l)$, and comparing them to the precom-
puted pattern hashes. Every potential occurrence
must be checked, which leads to poor worst-case

Time complexity.

- Preprocessing: $O(m)$
- Search: $O(n)$

EXPERIMENTAL RESULTS

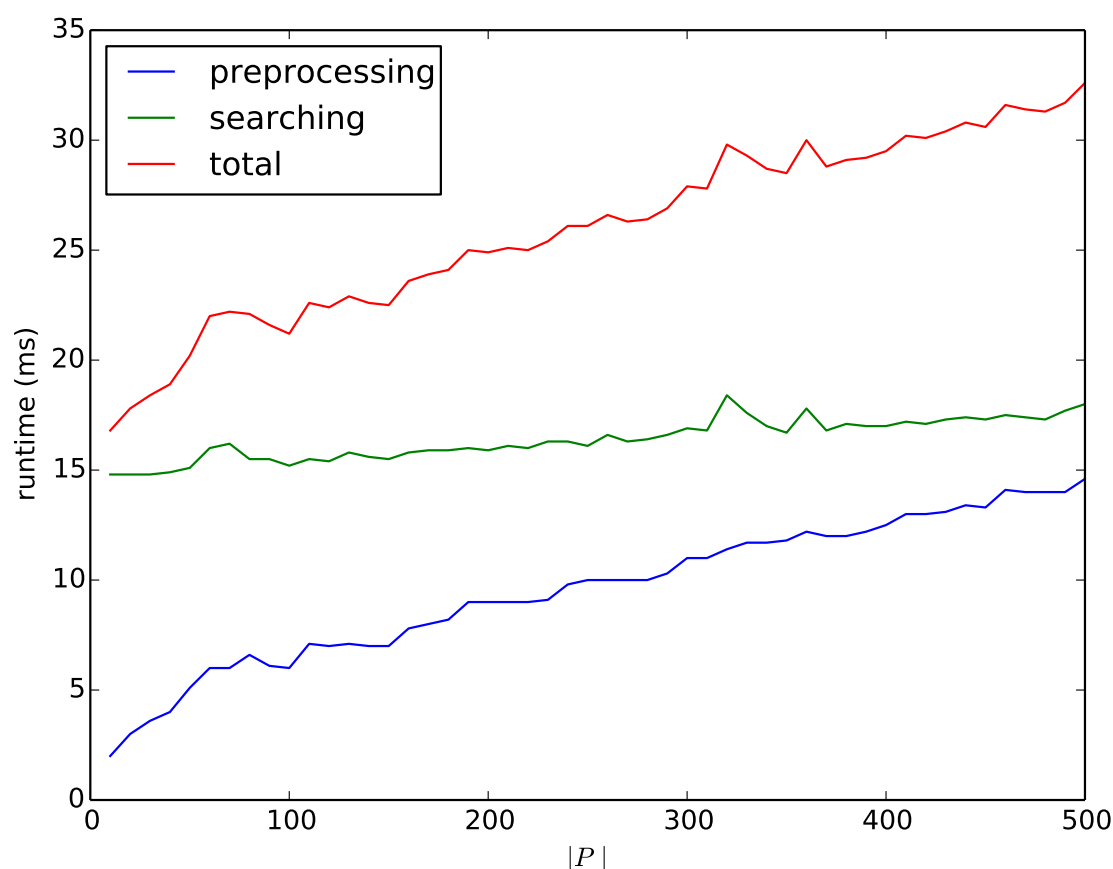
METHODOLOGY

Implementations. The algorithms shown above were implemented with java version 1.7. The source code for the implementations is available on github at <https://github.com/psaikko/string-algorithms-project>.

Experiments. The experiment data was gathered using python scripts to automate running the algorithms. Every plotted data point is an average of 10 runs of the algorithm.

COMPARISON WITH THEORETICAL PERFORMANCE

With constant text length. We plot the preprocessing and search times for each algorithm as the number of patterns to search for is increased.

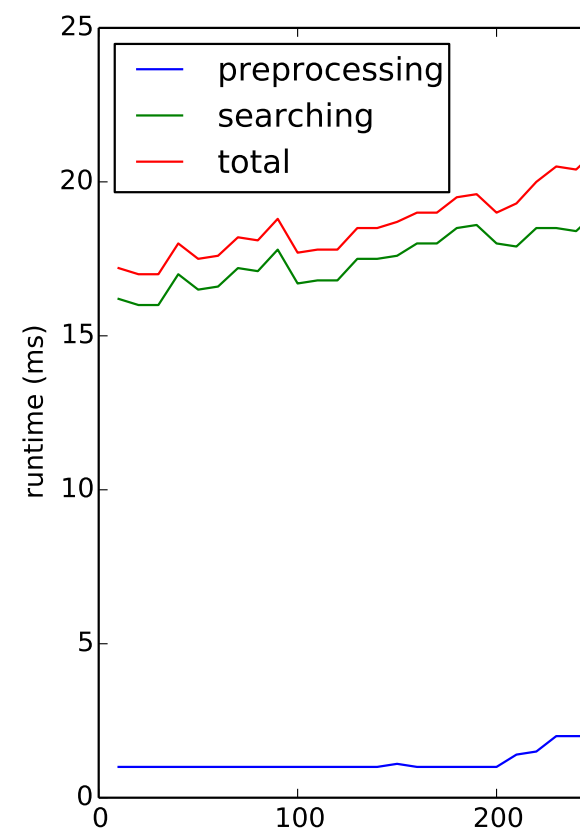


Aho-Corasick. We see search time staying roughly constant as expected, while preprocessing time increases linearly as we add more patterns.

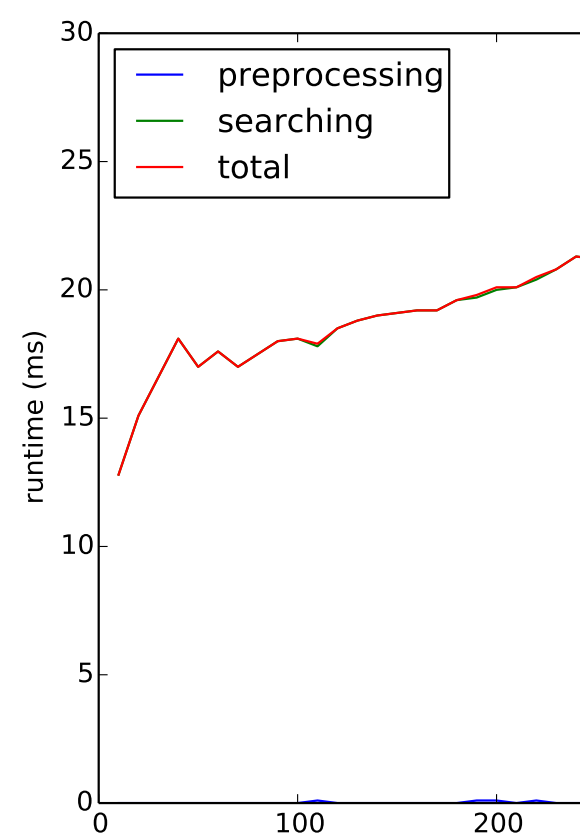
in $O(\lceil m/w \rceil)$ bitwise op

Time complexity.

- Preprocessing: $O(m)$
- Search: $O(n \lceil m/w \rceil)$



Karp-Rabin. The algorithm performs well despite poor worst-case preprocessing is quite fast



Shift-And. The preprocessing time is much shorter than the search time despite the same asymptotic

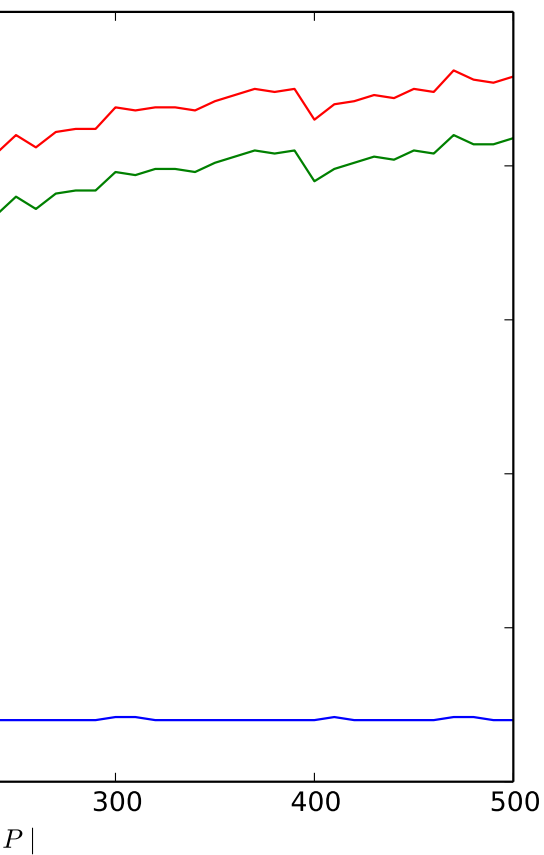
operations at each step.

behavior.

Average time complexity.

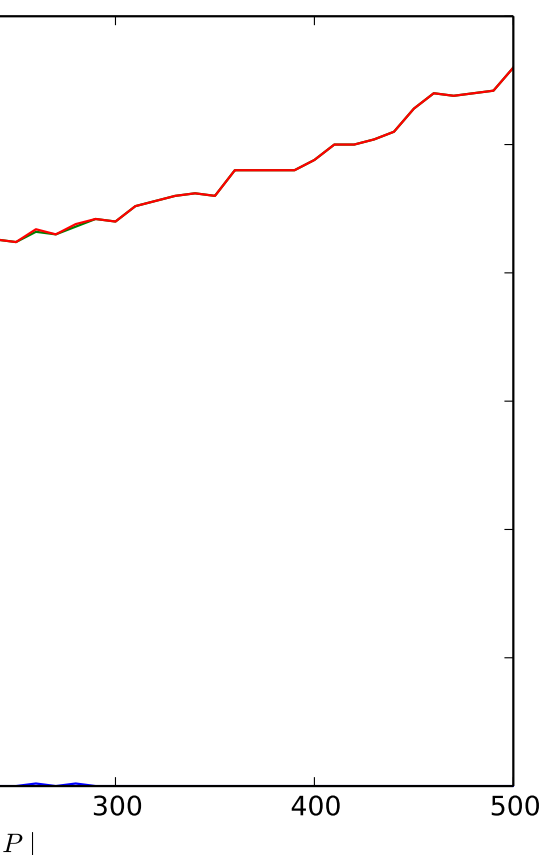
- Preprocessing: $O(m)$
- Search: $O(n + m)$

EDGE CASES

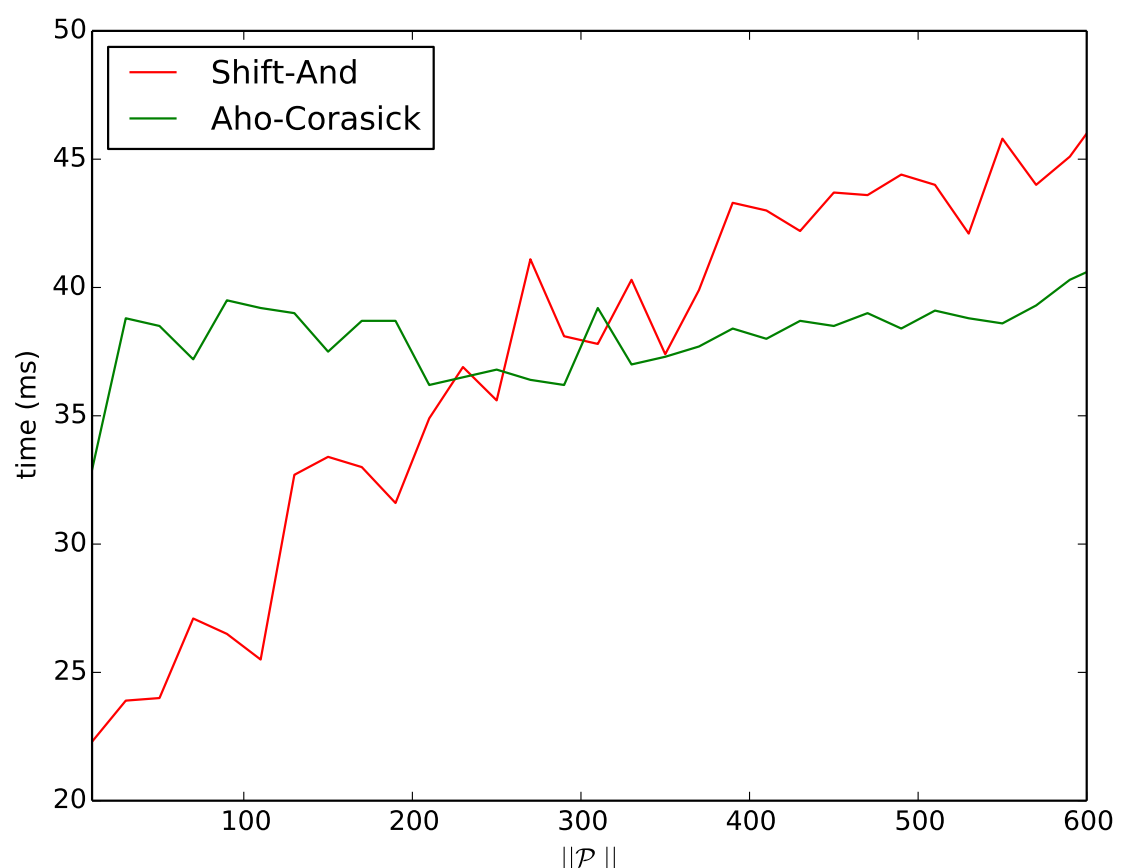
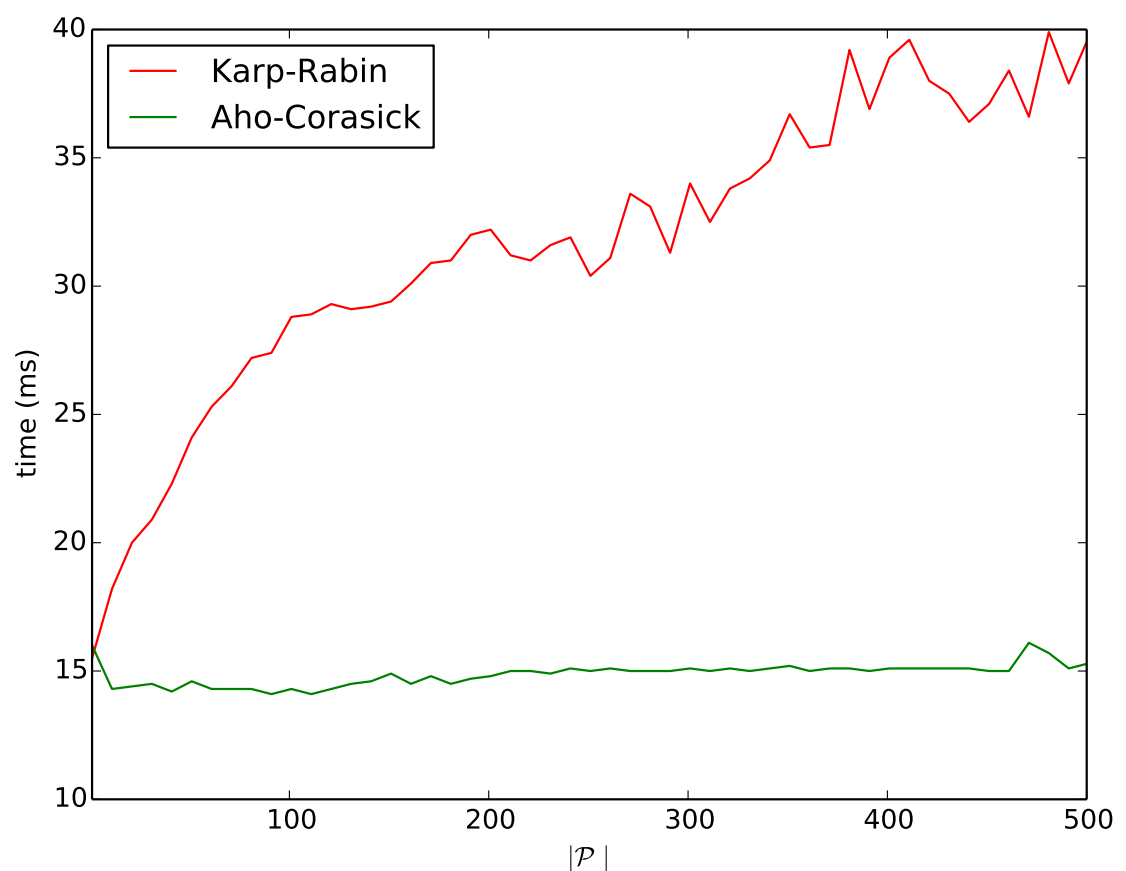


Algorithm generally performs well in the worst-case behavior for search, but not for preprocessing.

Many short patterns. A combination of the Karp-Rabin algorithm's $O(n + m)$ time complexity and hash collisions leads it to perform poorly.



Preprocessing time for Shift-And is much faster than other algorithms, despite its worst-case time complexity.



Small total pattern length. Although Shift-And scales poorly with total pattern length, the bit-parallel algorithm is fast when $\|P\|$ is not large.