

Rabin-Karp for string matching

...

Endre Valen, Martin Toft, Sergey Jakobsen
endrelv@ifi.uio.no, matoft@ifi.uio.no, sergeyj@ifi.uio.no

History and application

Michael O. Rabin and Richard M. Karp (1987)

Example application: Detecting plagiarism

String matching

$T \leftarrow \text{Text}$

$P \leftarrow \text{Pattern}$

Problem: $P \in T?$

Naive algorithm

$i, j \leftarrow 0, 0$

1. Compare $T[i]$ with $P[j]$
2. If match: increase i and j
 - 2.1. If $j = \text{len}(P)$: $P \in T$ ✓
3. If not match: increase i , $j \leftarrow 0$
4. If $i = \text{len}(T)$: $P \notin T$ ✗
5. Repeat step 1

Rabin-Karp

What makes Rabin Karp special? Rolling hash!

1. Hash P
2. Hash $|P|$ first characters of T
3. After this, we can find the next hash values in T in constant time by applying Horner's method!

Importance of hash function

We want to limit the amount of collisions

The hash function should map every unique input to a unique value

Easier said than done in practice

Demonstration

Assume we have -

a rolling hash function: $\text{hash}(\text{string}) \rightarrow \text{Int}$

a pattern: “ABA”

a text: “AABABCABACABBA”

A	A	B	A	B	C	A	B	A	C	A	B	B	A
---	---	---	---	---	---	---	---	---	---	---	---	---	---

$\text{hash}(\text{"ABA"}) \rightarrow 33$

A	A	B	A	B	C	A	B	A	C	A	B	B	A
---	---	---	---	---	---	---	---	---	---	---	---	---	---

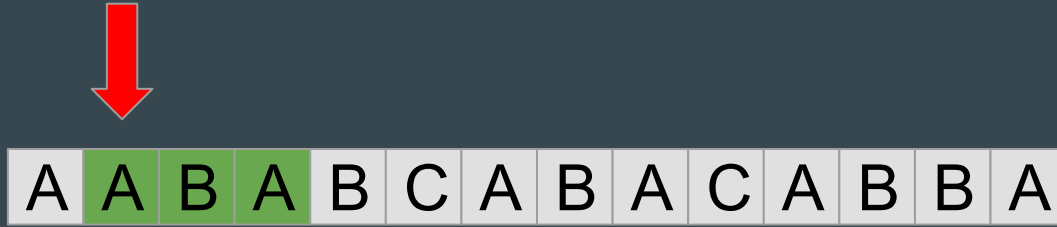
$\text{hash}(\text{"AAB"}) \rightarrow 17$

$\text{hash}(\text{"ABA"}) \rightarrow 33$

A	A	B	A	B	C	A	B	A	C	A	B	B	A
---	---	---	---	---	---	---	---	---	---	---	---	---	---

$\text{hash}(\text{"ABA"}) \rightarrow 33$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



$\text{hash}(\text{"ABA"}) \rightarrow 33$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



A	A	B	A	B	C	A	B	A	C	A	B	B	A
---	---	---	---	---	---	---	---	---	---	---	---	---	---

$\text{hash}(\text{"ABA"}) \rightarrow 33$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



A	A	B	A	B	C	A	B	A	C	A	B	B	A
---	---	---	---	---	---	---	---	---	---	---	---	---	---

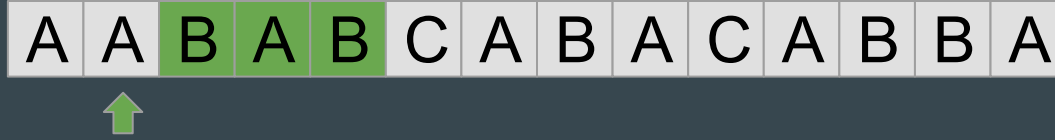
$\text{hash}(\text{"ABA"}) \rightarrow 33$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



$\text{hash}(\text{"ABA"}) \rightarrow 33$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



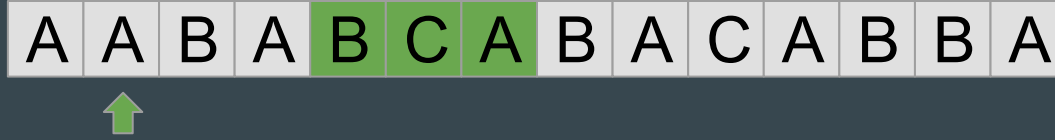
$\text{hash}(\text{"BAB"}) \rightarrow 45$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



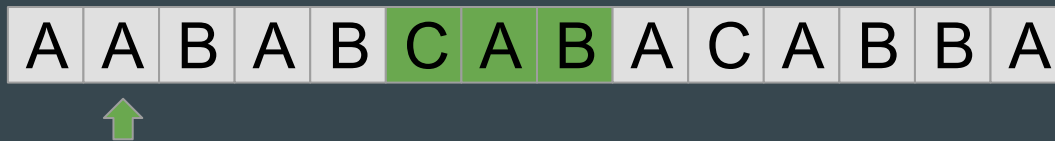
$\text{hash}(\text{"ABC"}) \rightarrow 28$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



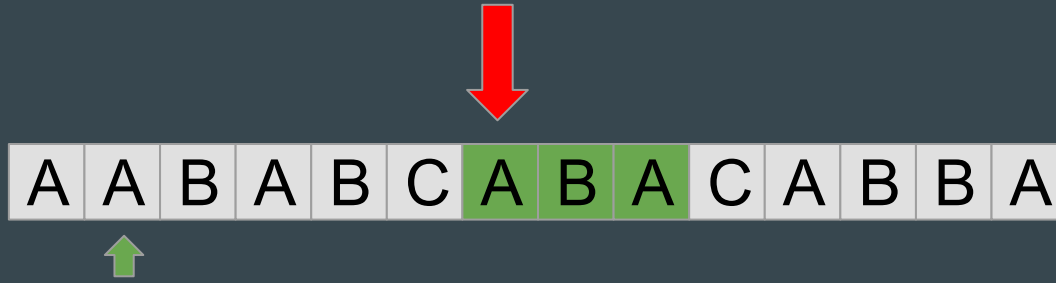
$\text{hash}(\text{"BCA"}) \rightarrow 13$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



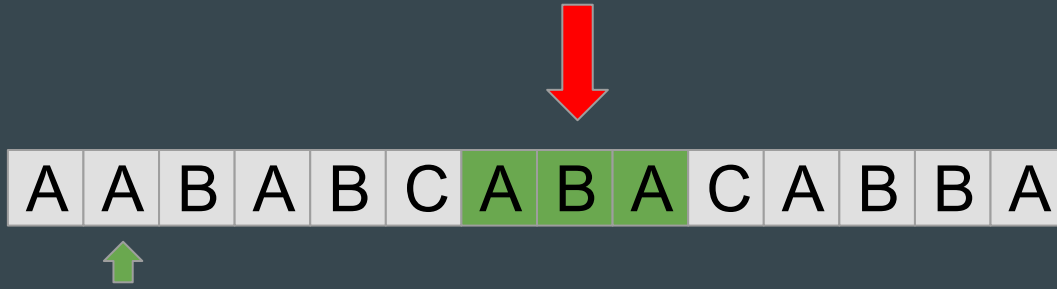
$\text{hash}(\text{"CAB"}) \rightarrow 17$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



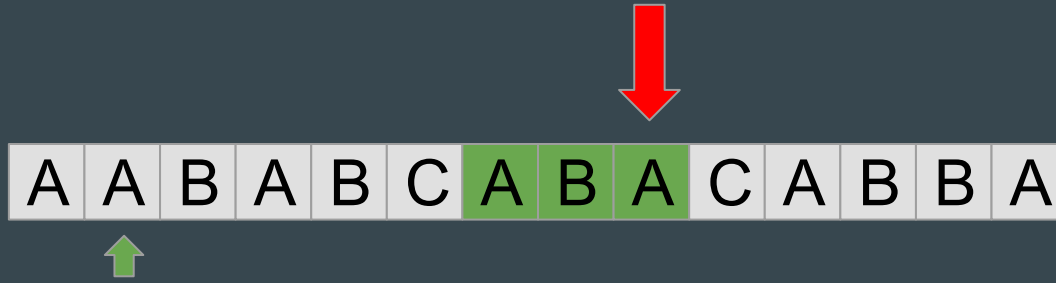
$\text{hash}(\text{"ABA"}) \rightarrow 33$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



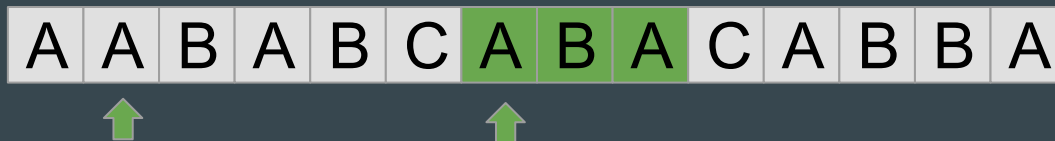
$\text{hash}(\text{"ABA"}) \rightarrow 33$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



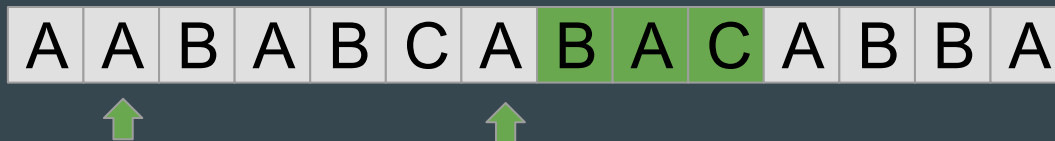
$\text{hash}(\text{"ABA"}) \rightarrow 33$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



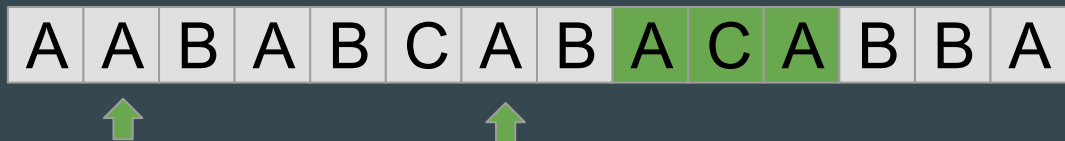
$\text{hash}(\text{"ABA"}) \rightarrow 33$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



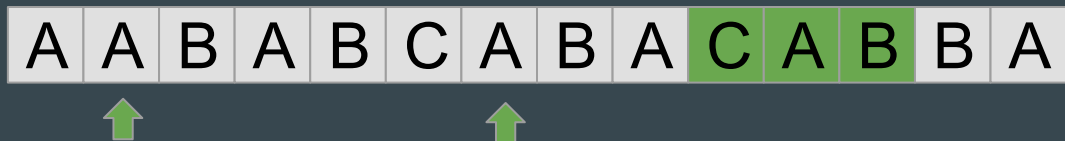
$\text{hash}(\text{"BAC"}) \rightarrow 51$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



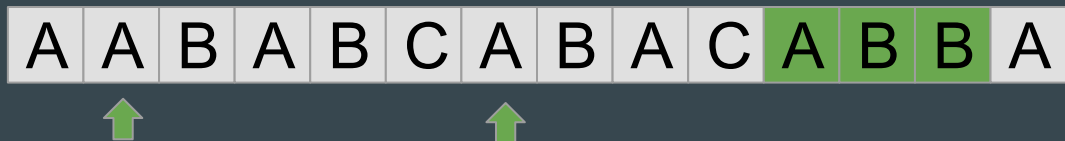
$\text{hash}(\text{"ACA"}) \rightarrow 43$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



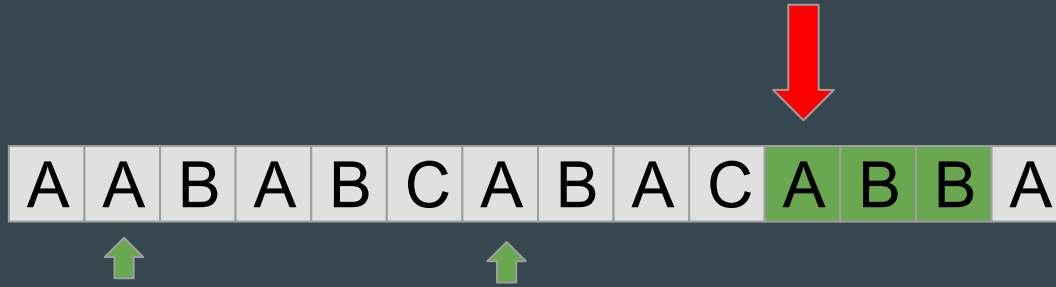
$\text{hash}(\text{"CAB"}) \rightarrow 17$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



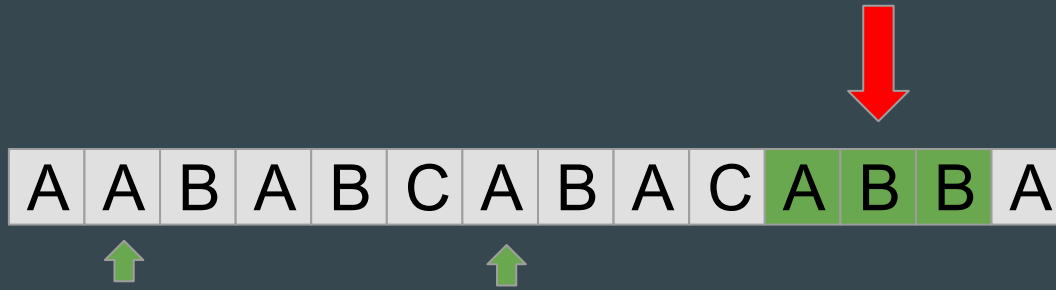
$\text{hash}(\text{"ABB"}) \rightarrow 33$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



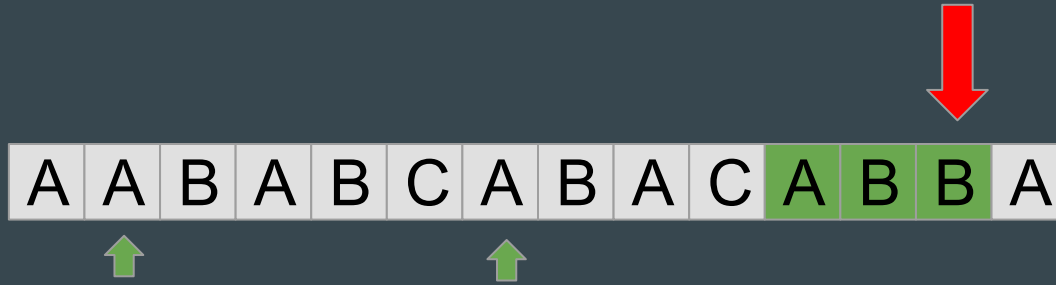
$\text{hash}(\text{"ABB"}) \rightarrow 33$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



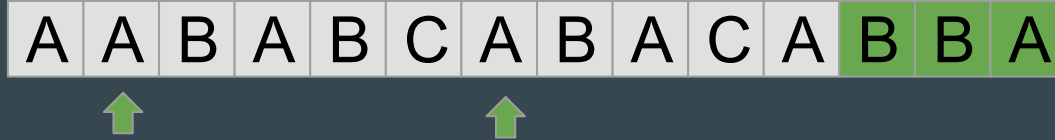
$\text{hash}(\text{"ABB"}) \rightarrow 33$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



$\text{hash}(\text{"ABB"}) \rightarrow 33$

$\text{hash}(\text{"ABA"}) \rightarrow 33$



$\text{hash}(\text{"BBA"}) \rightarrow 47$

$\text{hash}(\text{"ABA"}) \rightarrow 33$

Complexity and advantages/disadvantages

Average/best case: $O(n)$

Worst case: $O(mn) = O(m) + O(m) * O(n)$

Hash collisions and false positives!

Thanks for watching!

Sources:

- https://en.wikipedia.org/wiki/Rabin%E2%80%93Karp_algorithm
- <https://www.uio.no/studier/emner/matnat/ifi/IN3130/h21/slides/forelesning-3---string-search-pk.pdf>