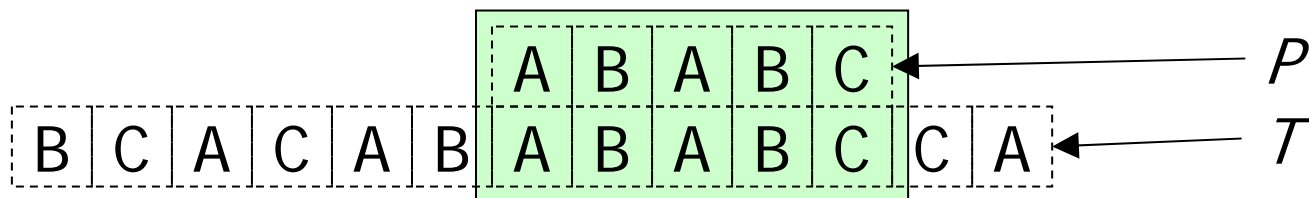# String Matching

Alexandre David

B2-206

# The Problem

- Given a text $T$ and a patter $P$, find an occurrence of $P$ inside $T$ or return *no match*.

- $T$ is of size $t$, $P$ is of size $p$.

- Example:

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

⟵ $P$

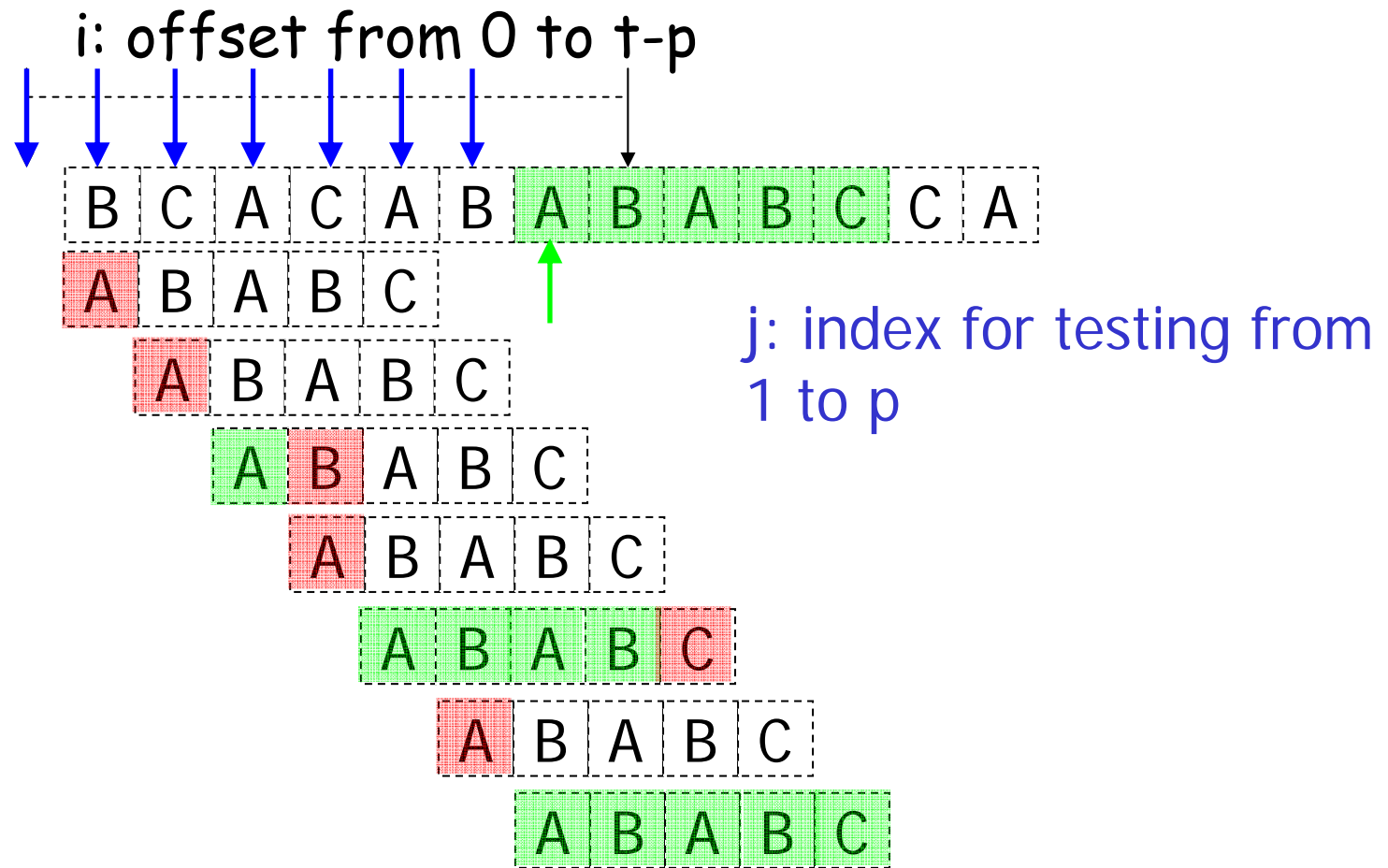| B | C | A | C | A | B | A | B | A | B | C | C | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

⟵ $T$

# Naïve Solution

- Compare P to T starting at position 1.
  - If mismatch, move P to the right and try again.
  - If match, return current position.

- Worst-case: $(t-p+1)*p$ comparisons, that is $O((t-p)*p)$. If $p=O(t)$ then we have $O(t*p)$.

```
naïve_find(T,P):
p = length(P)
t = length(T)
for i = 0 to t-p do
  ok = true
  for j = 1 to p do
    if P[j] != T[i+j] then
      ok = false
      break
    fi
  done
  if ok then return i+1
done
return -1
```
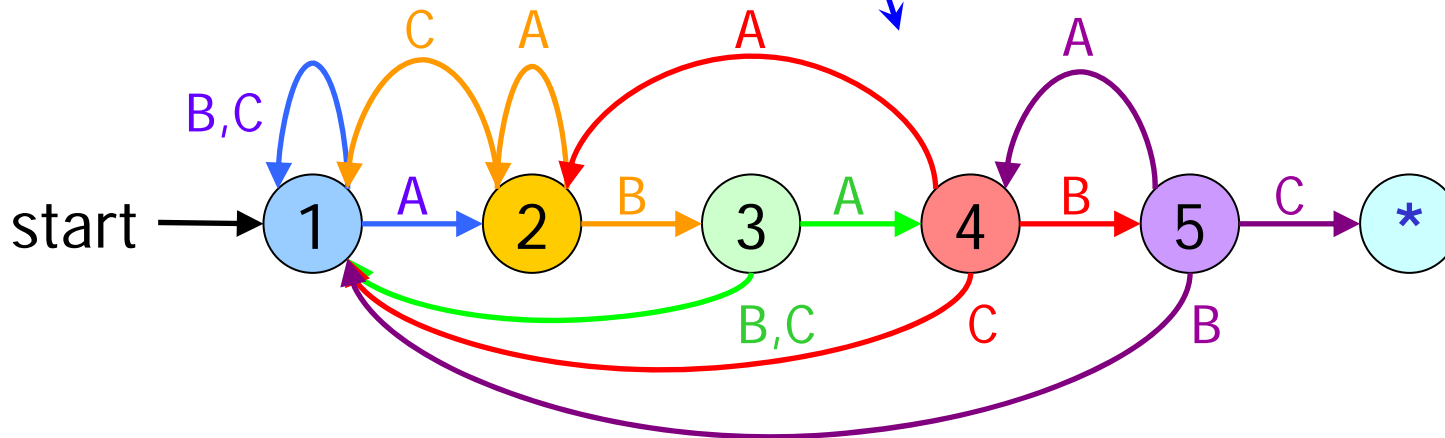
AA1

# Example

i: offset from 0 to t-p

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

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

j: index for testing from 1 to p

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

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

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

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

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

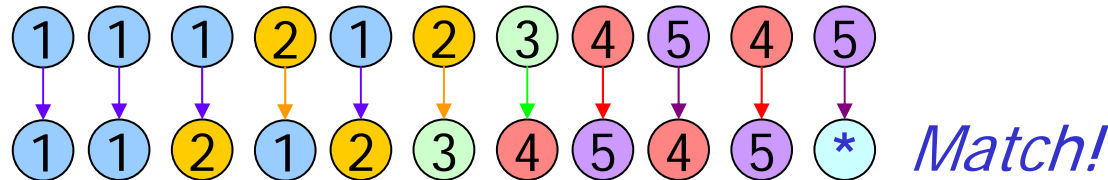| A | B | A | B | C |
|---|---|---|---|---|

# Solution With Finite Automata

- Given *P*, it is possible to construct a finite automaton that is used to scan *T* in *O(t)*.

- Idea: Remember the last matched sub-string and re-use the information.

- Matching = reach the state **\***.
  No match = get stuck in the automaton.

- Pre-processing required: Construct the automaton in *O(p\*|alphabet|)*.
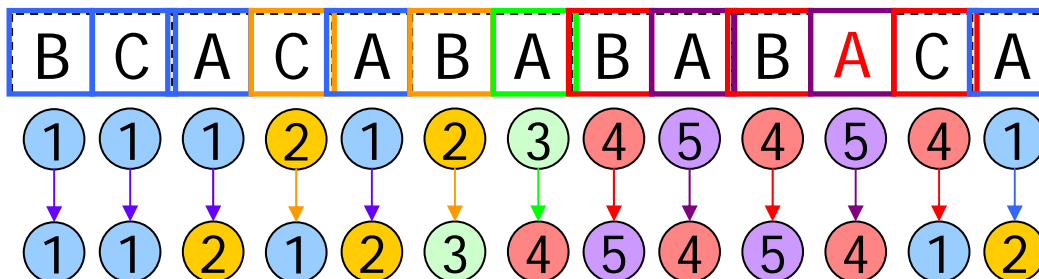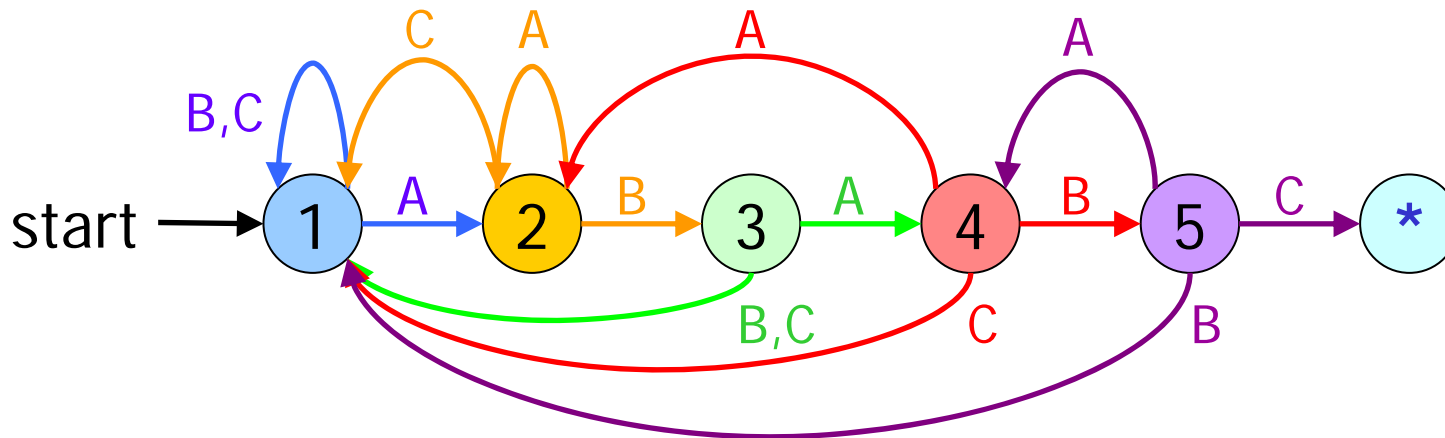
# Example With Automaton

A B A B C



start → 1

B C A C A B A B A B C C A

1 1 1 2 1 2 3 4 5 4 5
↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓
1 1 2 1 2 3 4 5 4 5 *    *Match!*

# Example With Automaton

A B A B C



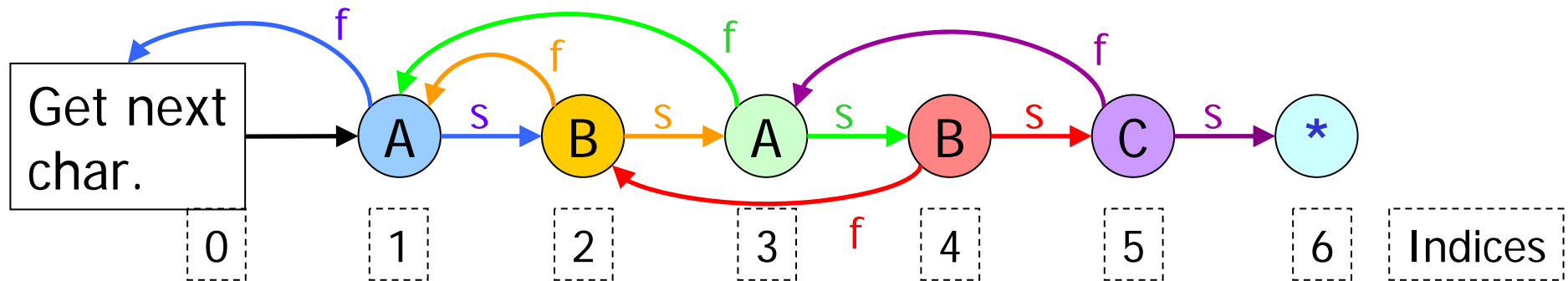*End of string, stuck in the automaton ⇒ No match!*

# Knuth-Morris-Pratt Flowchart

- Given $P$, it is possible to construct a finite flowchart that is used to scan $T$ in $O(t+p)$.

- Idea is to remember the maximum of matchable characters before the $i^{th}$ position.

- Matching = reach the state $*$.
No match = get stuck in the automaton.

- Pre-processing: Construct the flowchart in $O(p^2)$.

# Example With Flowchart

A B A B C → Read and test character. → s: success, read, test
f: fail, test



Get next char.

```
f          f          f          f
A → B → A → B → C → *
   s    s    s    s    s
0    1    2    3    4    5    6    Indices
                     f
```

Construct *next* table (f):

| A | B | A | B | C | - |
|---|---|---|---|---|---|
| 0 | 1 | 1 | 2 | 3 | |

# Example With Flowchart

| 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| $n$ | A | B | A | B | C | * |
|   | 0 | 1 | 1 | 2 | 3 |   |

| B | C | A | C | A | B | A | B | A | B | C | C | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | A | A | A | B | A | B | A | B | C | B | C |   |
| $n$ | 0 | 0 | 2 | 1 | 2 | 3 | 4 | 5 | 3 | 5 | 6 |   |
| 1 | $n$ | $n$ |   | A |   |   |   |   | A |   | * |   |
|   | 1 | 1 |   | 0 |   |   |   |   | 4 |   |   |   |
|   |   |   |   | $n$ |   |   |   |   |   |   |   |   |
|   |   |   |   | 1 |   |   |   |   |   |   |   |   |

*Match!*

# Boyer-Moore Algorithm

- Idea is to skip text without checking it. Scan from right to left, use heuristics to decide how far to jump.

- Average running time $O(t/p)$, worst $O(t*p)$.

| B | C | A | C | **D** | A | B | **E** | B | C | A | B | A | B | C | C | A |

| A | B | A | B | C |

D not in P ⇒    | A | B | A | B | C |

E not in P ⇒    | A | B | A | B | C |

| A | B | A | B | C |

| A | B | A | B | C |    *Match!*

# Rabin-Karp Algorithm

- Use a hash function to identify equal strings! Very powerful for multi-pattern matching.

- Trick: Use a special hash function. Treat the character as a number in some base (usually a big prime) and compute the next hash iteratively.
  Hopefully, we have few collisions.

- Average running time $O(t)$, worst $O(t*p)$.

# Rabin-Karp Algorithm

- Hash update = "shift" in the corresponding base.

- In practice, useful to use base 256 for characters and a prime as the hash table size.

  - Very fast and hash performs reasonably well.

# Example With Rabin-Karp

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

$\longrightarrow$ $hash_p$ $O(p)$.

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

Initial hash $O(p)$. Test $O(1) \rightarrow$ no.

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

Update hash $O(1)$. Test $O(1) \rightarrow$ no.

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

Updates of hash $O(1)$ + tests $O(1)$... Test $O(1) \rightarrow$ yes.

Test P $O(p) \rightarrow$ yes.