



UNIVERSITÄT
DES
SAARLANDES



ZBI ZENTRUM FÜR
BIOINFORMATIK

Exact Pattern Matching with Automata

Algorithms for Sequence Analysis

Sven Rahmann

Summer 2021

Recall the Pattern Matching Problem

alice was beginning to get very tired of sitting by her sister on the bank and of having nothing to do once or twice she had peeped into the book her sister was reading but it had no pictures or conversations in it and what is the use of a book thought alice without pictures or conversation

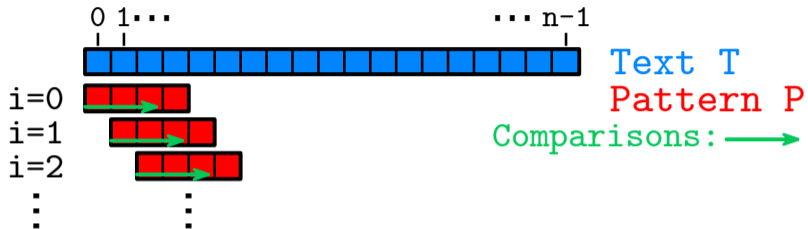
Task

Find all occurrences of a given string in another (longer) string.

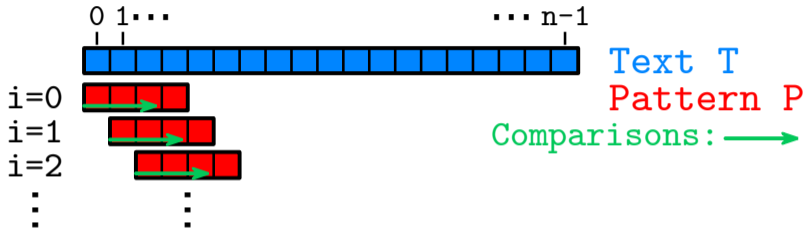
Goals

- As **fast** as possible (running time)
- As **easily** as possible (algorithm/implementation)

Recall the naïve algorithm



Recall the naïve algorithm



Ideas

- Can we shift window by more than one character?
→ Horspool algorithm (and others)
- We “touch” the same characters in T multiple times.
Can we “reuse” information from preceding comparisons?
→ Automata based algorithms (now)

Finite Automata Revisited

Deterministic Finite Automata (DFA)

Definition (DFA)

A **DFA** is a tuple $(Q, q_0, F, \Sigma, \delta)$, where

- Q is a finite set of **states**,
- $q_0 \in Q$ is a **start state**,
- $F \subset Q$ is a set of **accepting states**,
- Σ is an input **alphabet**, and
- $\delta : Q \times \Sigma \rightarrow Q$ is a **transition function**.

DFA – Example

- Q is a finite set of **states**,
- $q_0 \in Q$ is a **start state**,
- $F \subset Q$ is a set of **accepting states**,
- Σ is an input **alphabet**, and
- $\delta : Q \times \Sigma \rightarrow Q$ is a **transition function**.

Accept the strings over $\{a, b, c\}$, where 4 divides the sum of the number of as and bs:

Non-Deterministic Finite Automata (NFA)

Definition (NFA)

An **NFA** is a tuple $(Q, Q_0, F, \Sigma, \Delta)$, where

- Q is a finite set of **states**,
- $Q_0 \subset Q$ is a set of **start states**,
- $F \subset Q$ is a set of **accepting states**,
- Σ is an input **alphabet**, and
- $\Delta: Q \times \Sigma \rightarrow 2^Q$ is a (non-deterministic) **transition function**.

NFA – Example

- Q is a finite set of **states**,
- $Q_0 \subset Q$ is a set of **start states**,
- $F \subset Q$ is a set of **accepting states**,
- Σ is an input **alphabet**, and
- $\Delta: Q \times \Sigma \rightarrow 2^Q$ is a (non-deterministic) **transition function**.

Accept the strings over $\{a, b, c\}$, where 3 or 4 divides the sum of the number of as and bs:

Extending the Transition Function

Original NFA transition function: $\Delta: Q \times \Sigma \rightarrow 2^Q$

For notational convenience, we make the following definitions.

Extension to sets of states

- $\Delta(A, c) := \bigcup_{q \in A} \Delta(q, c)$ for a **set** of states A and $c \in \Sigma$.

Extension to strings

- $\Delta(A, \epsilon) := A$, where ϵ is the empty string, and
- $\Delta(A, xc) := \Delta(\Delta(A, x), c)$, where $x \in \Sigma^*$ and $c \in \Sigma$.

NFAs for Pattern Matching

NFA to Solve the Pattern Matching Problem

Goal

For given pattern $P \in \Sigma^*$, construct NFA that recognizes all strings Σ^*P .

NFA to Solve the Pattern Matching Problem

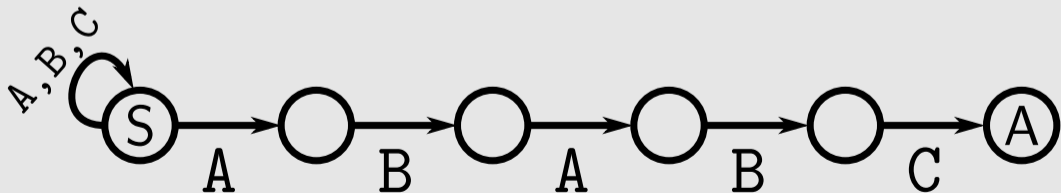
Goal

For given pattern $P \in \Sigma^*$, construct NFA that recognizes all strings Σ^*P .

Approach

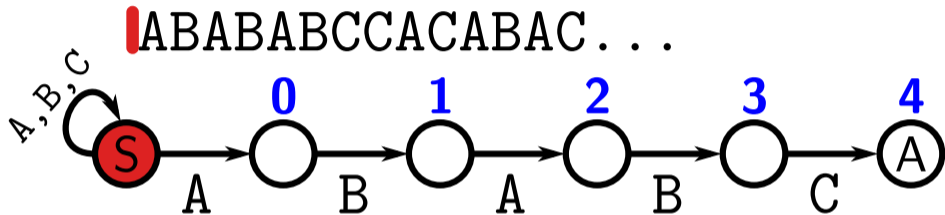
- “Linear chain” of states
- Start state remains always active

Example: $\Sigma = \{A, B, C\}$ and $P = ABABC$



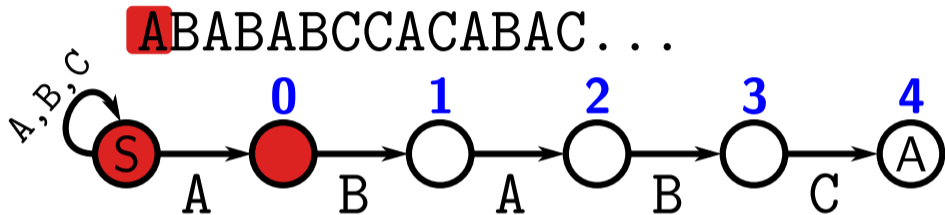
Using our NFA

$\Sigma = \{A, B, C\}$ and $P = ABABC$



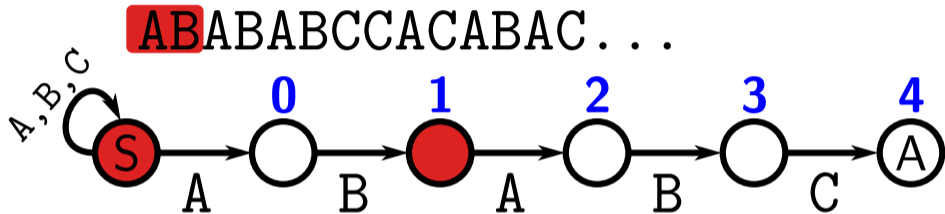
Using our NFA

$\Sigma = \{A, B, C\}$ and $P = ABABC$



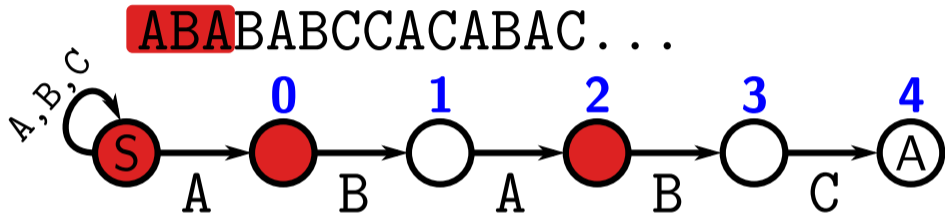
Using our NFA

$\Sigma = \{A, B, C\}$ and $P = ABABC$



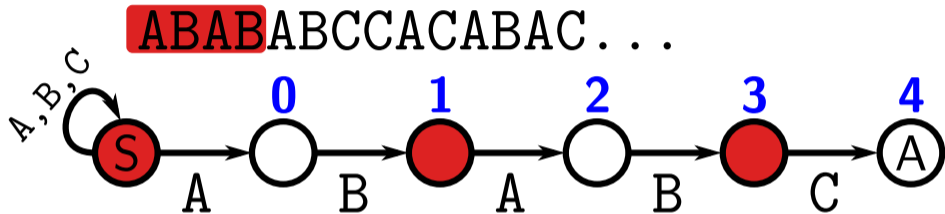
Using our NFA

$\Sigma = \{A, B, C\}$ and $P = ABABC$



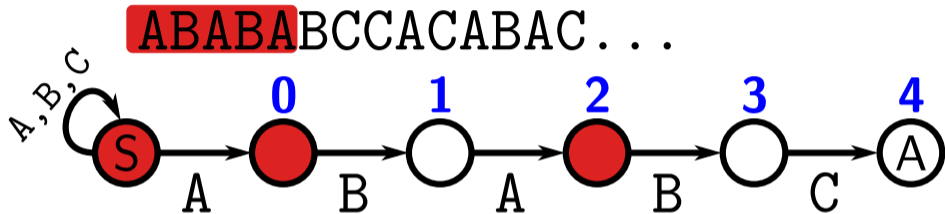
Using our NFA

$\Sigma = \{A, B, C\}$ and $P = ABABC$



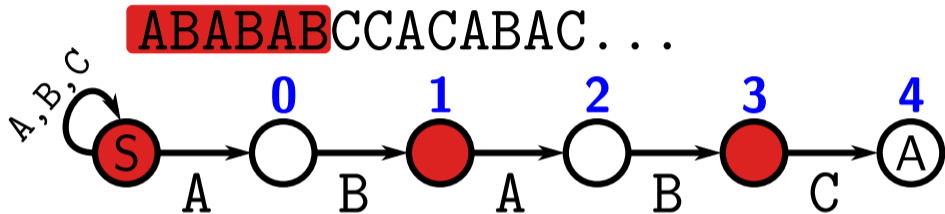
Using our NFA

$\Sigma = \{A, B, C\}$ and $P = ABABC$



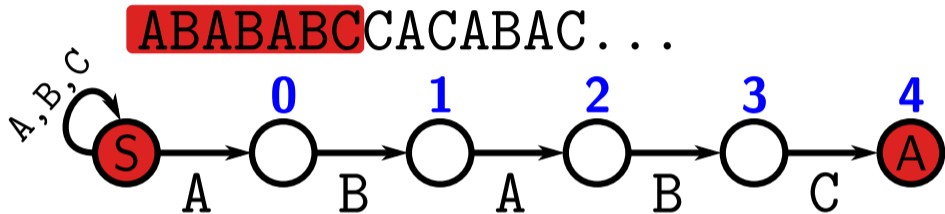
Using our NFA

$\Sigma = \{A, B, C\}$ and $P = ABABC$



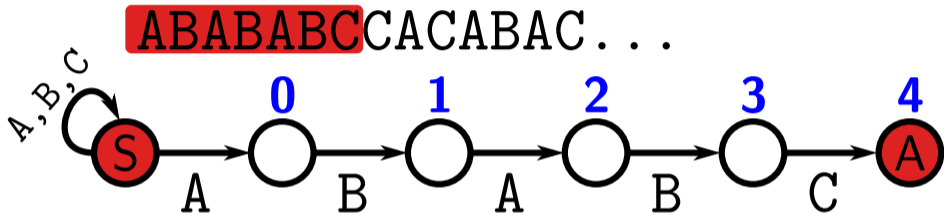
Using our NFA

$\Sigma = \{A, B, C\}$ and $P = ABABC$



Using our NFA

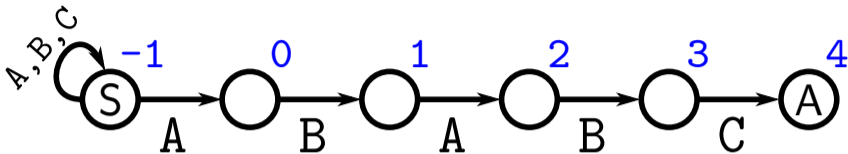
$\Sigma = \{A, B, C\}$ and $P = ABABC$



Things left to do

- Formally define this automaton,
- Give efficient implementation.

Pattern Matching NFA (Formal)



- state set $Q = \{-1, 0, \dots, m-1\}$, where $m = |P|$

- start states $Q_0 = \{-1\}$

- accepting states $F = \{m-1\}$

- transition function Δ :

For $q = -1$:
$$\Delta(-1, c) = \begin{cases} \{-1, 0\} & \text{if } c = P[0], \\ \{-1\} & \text{otherwise.} \end{cases}$$

For $q \in \{0, \dots, m-2\}$:
$$\Delta(q, c) = \begin{cases} \{q+1\} & \text{if } c = P[q+1], \\ \emptyset & \text{otherwise.} \end{cases}$$

For $q = m-1$:
$$\Delta(m-1, c) = \emptyset$$

Correctness

Lemma: NFA state set invariant

Let $A \subset Q$ be a set of active states of our NFA. Then, $q \in A \setminus \{-1\}$ iff the last $q + 1$ read characters equal the patterns prefix $P[\dots q]$. In particular, state $|P| - 1$ is active iff the last $|P|$ characters equal the full pattern.

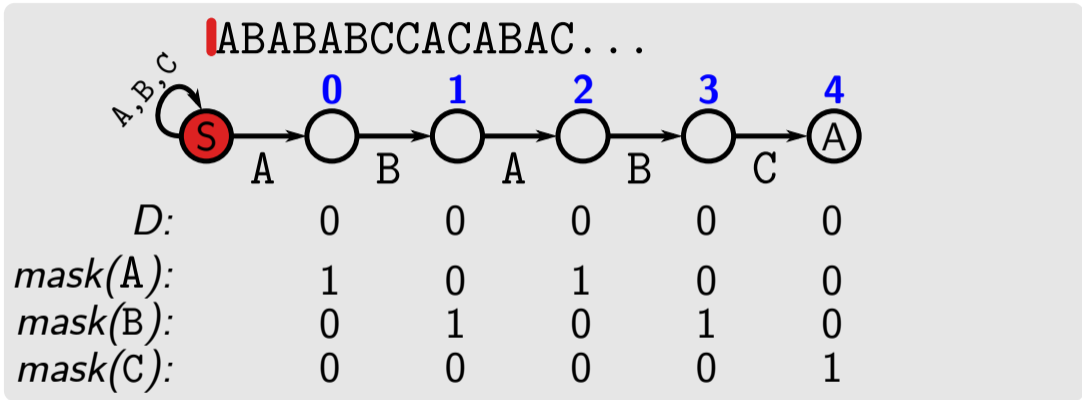
Proof: Follows directly from the NFA definition.

Theorem: Correctness of Pattern Matching NFA

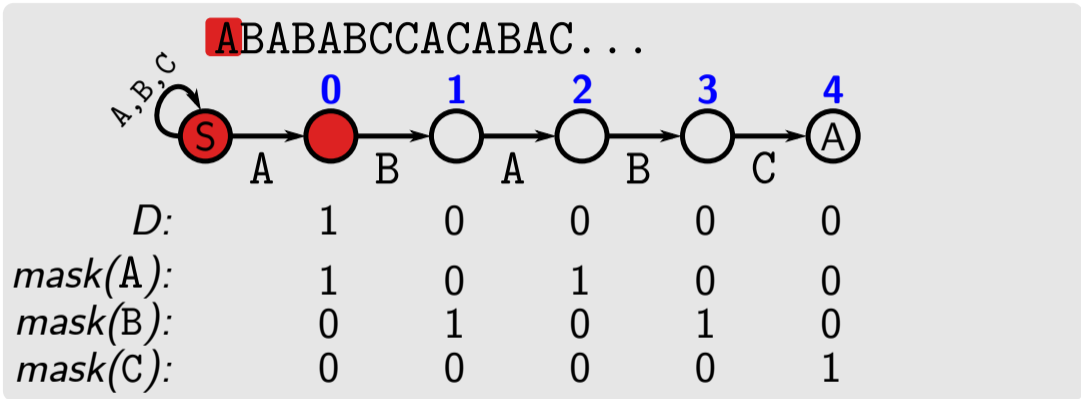
The pattern matching NFA for pattern P accepts exactly the language Σ^*P .

Proof: Follows immediately from the above lemma.

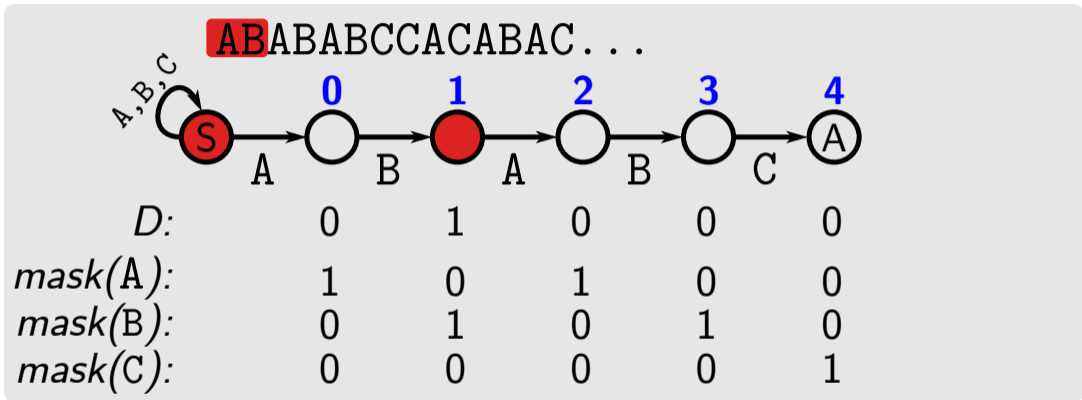
Efficient Implementation: The Shift-And Algorithm



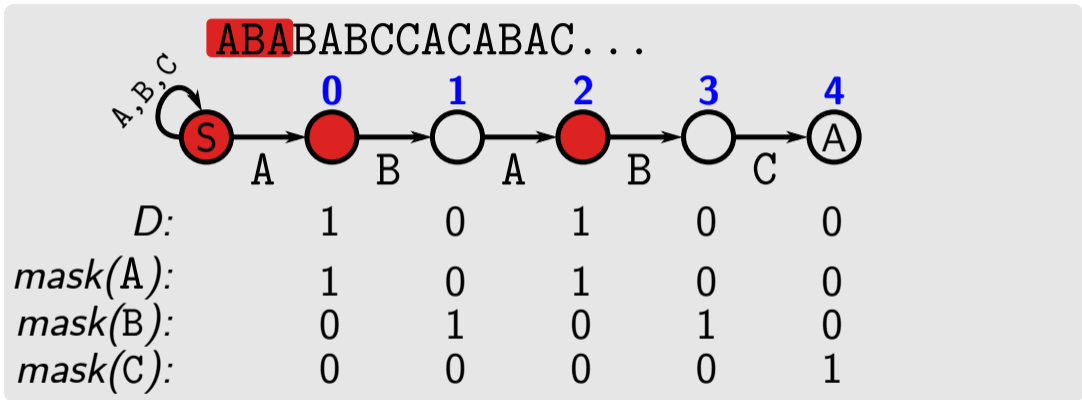
Efficient Implementation: The Shift-And Algorithm



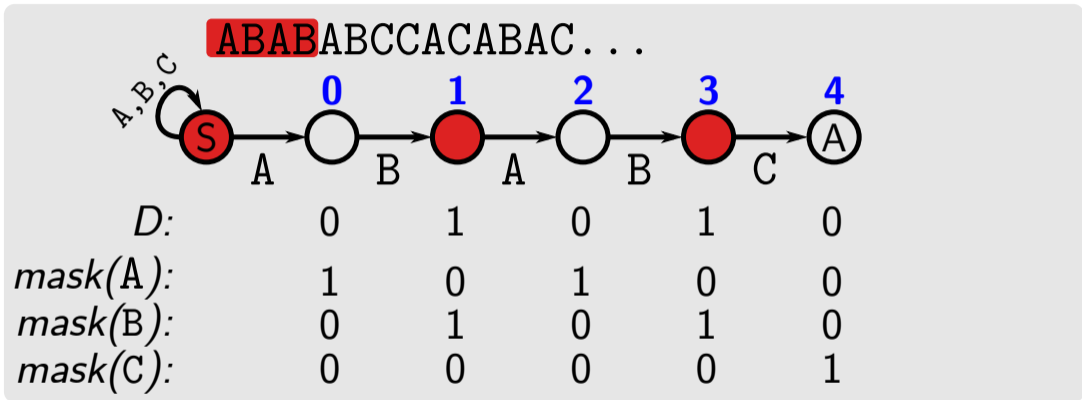
Efficient Implementation: The Shift-And Algorithm



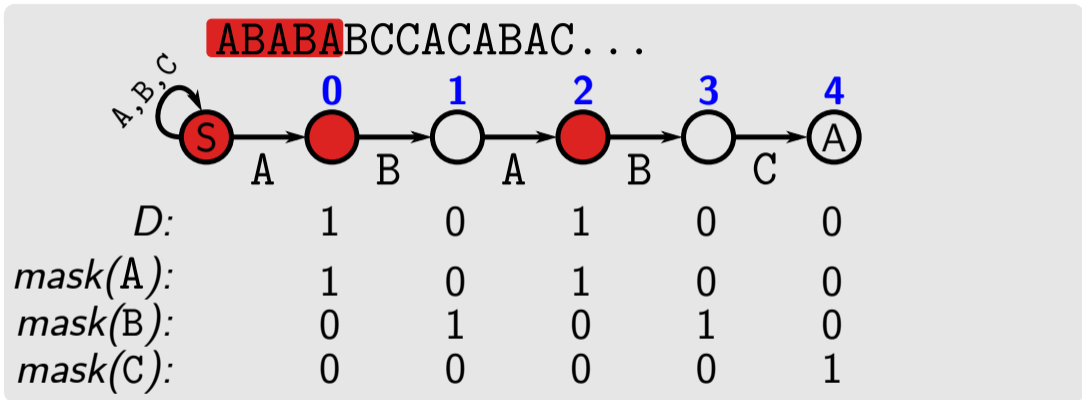
Efficient Implementation: The Shift-And Algorithm



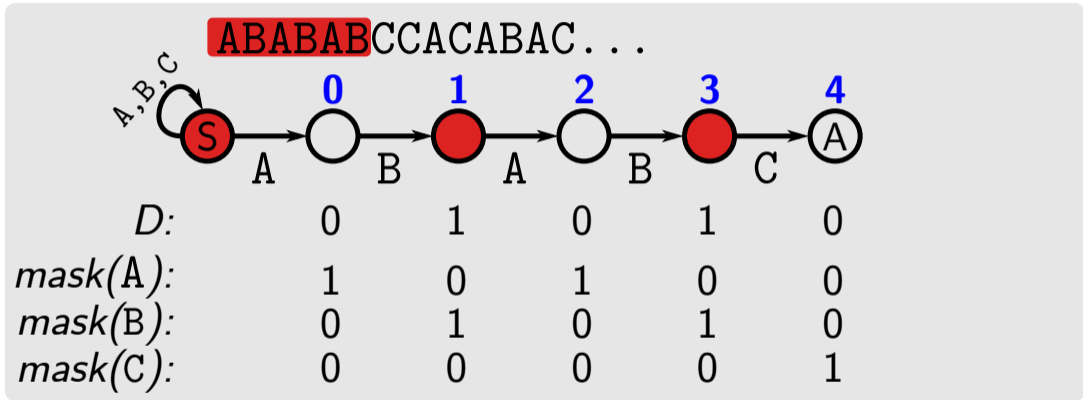
Efficient Implementation: The Shift-And Algorithm



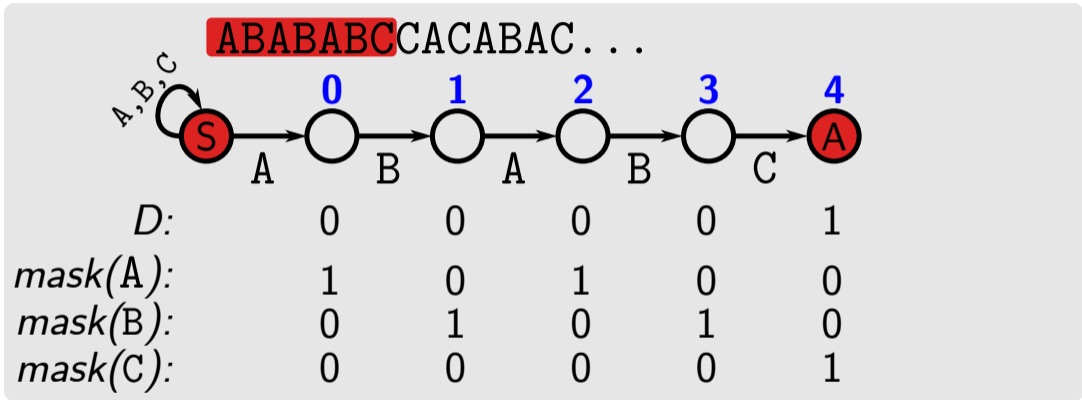
Efficient Implementation: The Shift-And Algorithm



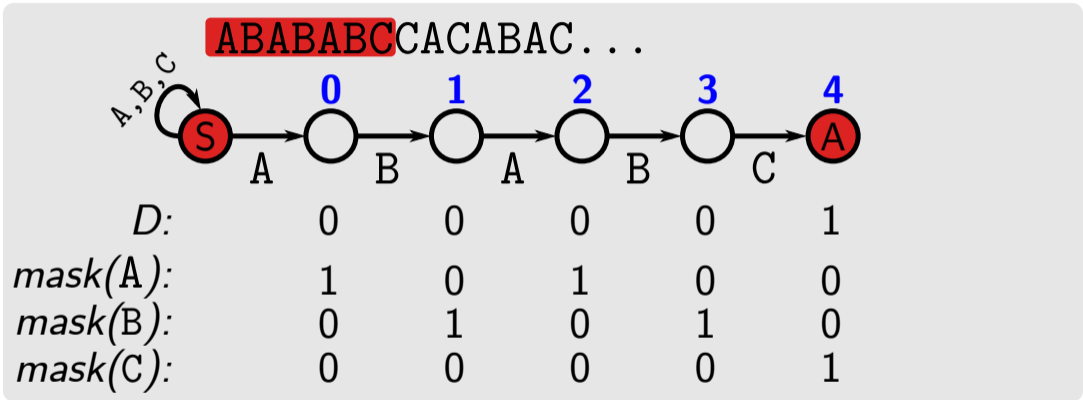
Efficient Implementation: The Shift-And Algorithm



Efficient Implementation: The Shift-And Algorithm

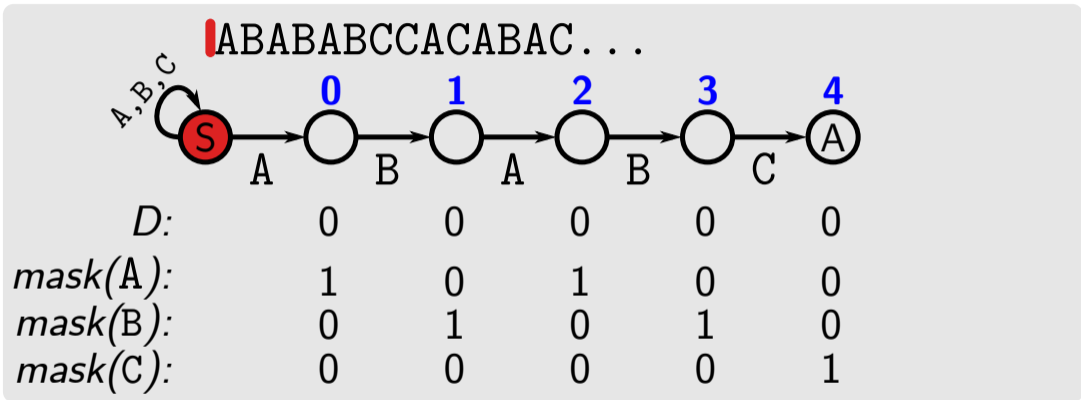


Efficient Implementation: The Shift-And Algorithm



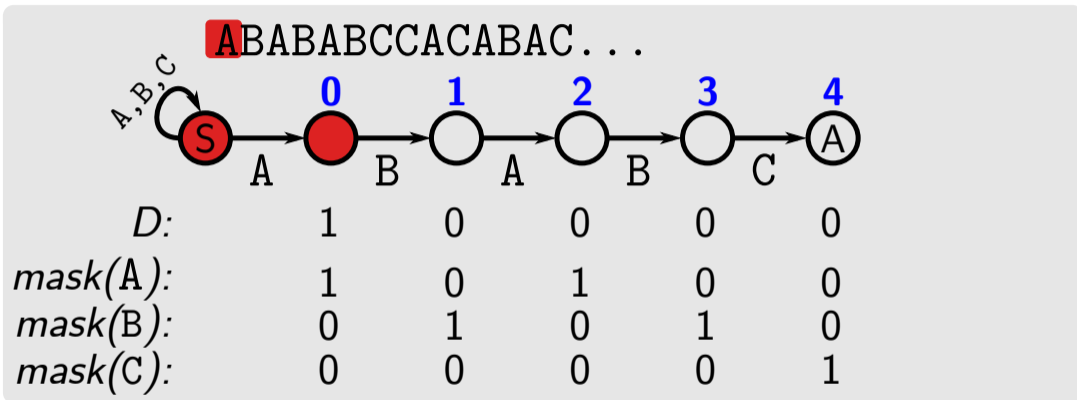
$$D \leftarrow ((D \ll 1) | 1) \& \text{mask}(c)$$

Efficient Implementation: The Shift-And Algorithm



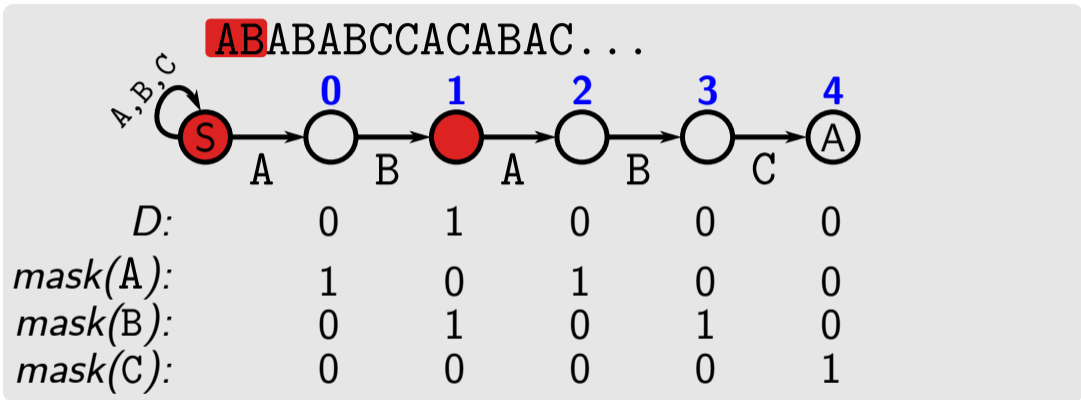
$$D \leftarrow ((D \ll 1) | 1) \& \text{mask}(c)$$

Efficient Implementation: The Shift-And Algorithm



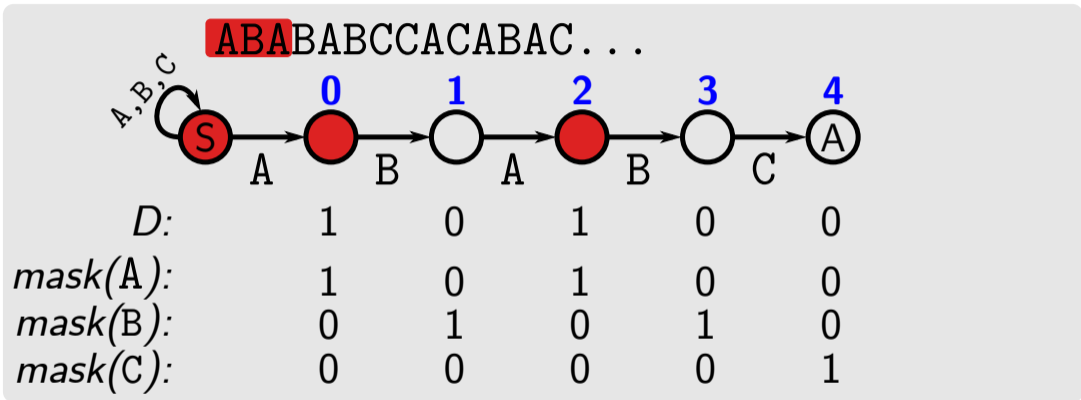
$$D \leftarrow ((D \ll 1) | 1) \& \text{mask}(c)$$

Efficient Implementation: The Shift-And Algorithm



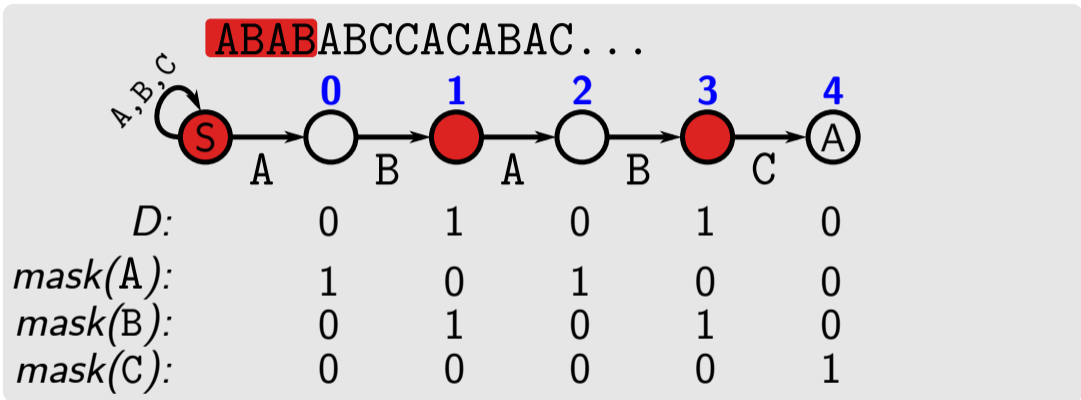
$$D \leftarrow ((D \ll 1) | 1) \& \text{mask}(c)$$

Efficient Implementation: The Shift-And Algorithm



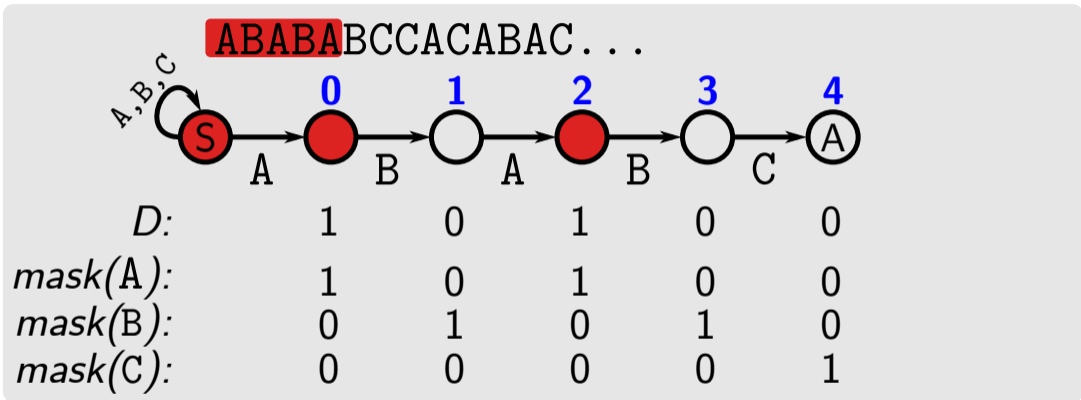
$$D \leftarrow ((D \ll 1) | 1) \& \text{mask}(c)$$

Efficient Implementation: The Shift-And Algorithm



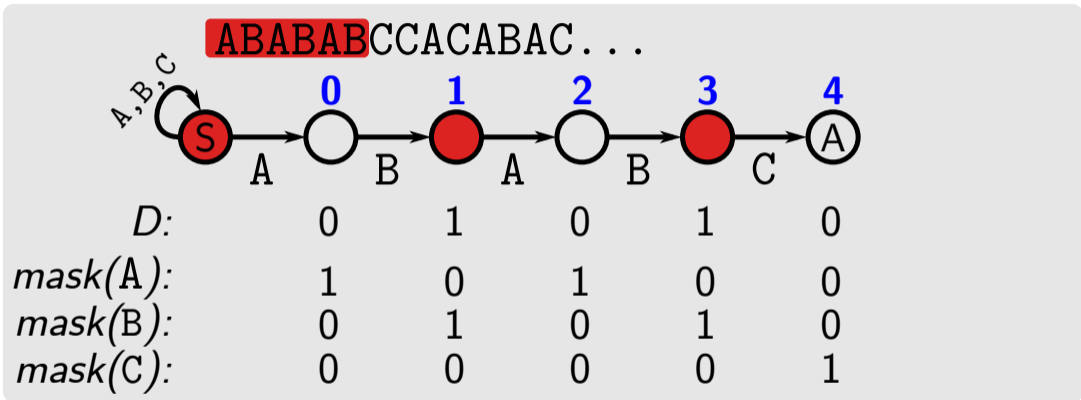
$$D \leftarrow ((D \ll 1) | 1) \& \text{mask}(c)$$

Efficient Implementation: The Shift-And Algorithm



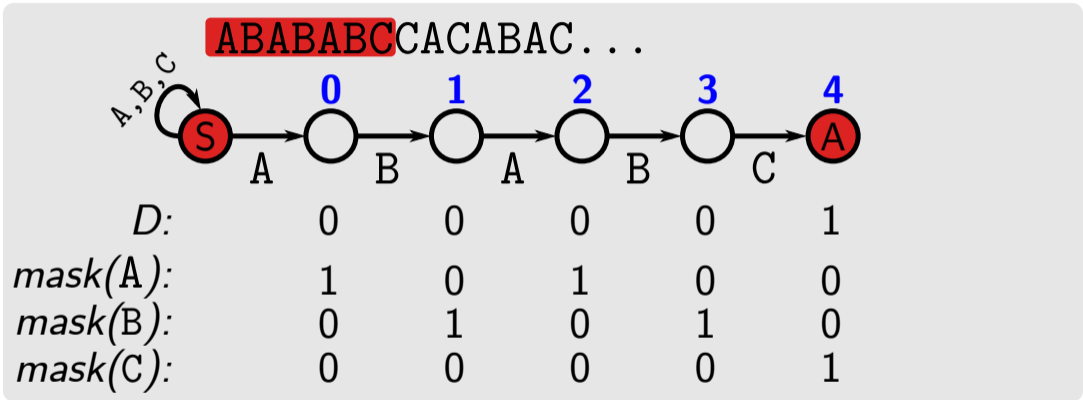
$$D \leftarrow ((D \ll 1) | 1) \& \text{mask}(c)$$

Efficient Implementation: The Shift-And Algorithm



$$D \leftarrow ((D \ll 1) | 1) \& \text{mask}(c)$$

Efficient Implementation: The Shift-And Algorithm



$$D \leftarrow ((D \ll 1) | 1) \& \text{mask}(c)$$

Code for Shift-And Algorithm

```
1 def ShiftAnd(P, T):
2     m = len(P)
3     masks = dict() # empty dictionary
4     bit = 1
5     for c in P:
6         if c not in masks: masks[c] = 0
7         masks[c] = masks[c] | bit
8         bit = bit * 2
9     accept_state = bit // 2
10    D = 0 # bit-mask of active states
11    i = 0
12    for c in T:
13        D = ((D << 1) | 1) & masks[c]
14        if (D & accept_state) != 0:
15            yield i
16        i += 1
```

Running time of Shift-And Algorithm

- m : Pattern length
- n : Text length
- w : Machine register width

Running time

If $m < w$, then the shift-and algorithm runs in $O(m + n)$ time.

Generally (i.e. when m/w is not constant), it takes $O(m + nm/w)$ time.

Conclusions

- Fast when pattern fits into one machine word
- Running time independent of how similar text and pattern are.
- Running time independent of alphabet size.

DFA-based Pattern Matching

Reminder: How to Turn an NFA into a DFA

Definition: equivalence of automata

Two automata are **equivalent** if they accept the same sets of words (languages).

NFA vs. DFA

- **NFA:** Many states can be active at the same time
- **DFA:** Only one state active at any given time

Subset construction (NFA \rightarrow DFA)

- **Idea:** create one DFA state for every **subset** of NFA states
- We can omit states that are **unreachable**

Example: Subset Construction (NFA \rightarrow DFA)

Accept the strings over $\{a, b, c\}$, where 3 or 4 divides the sum of the number of as and bs:

Subset Construction: How Many DFA States?

In general

Subset construction leads to **exponential blow-up** in the number of states:
 $|Q|$ NFA states turn into $2^{|Q|}$ DFA states.

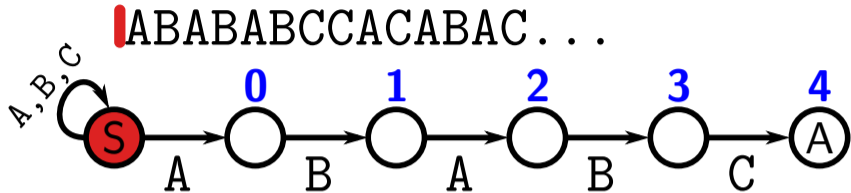
In practice

Blow-up often does not happen (many states are **unreachable**).

For our pattern matching automata

DFA always has the **same** number of states as **NFA**.

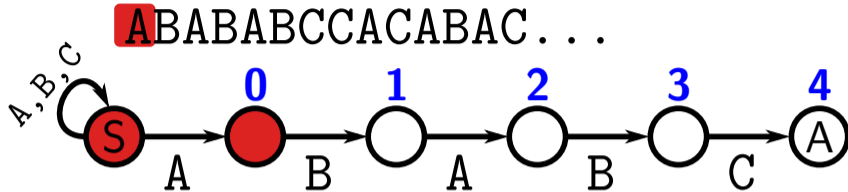
Reachable State Sets in NFA



Lemma (from previous lecture): NFA state set invariant

Let $A \subset Q$ be a set of active states of our NFA. Then, $q \in A \setminus \{-1\}$ iff the last $q + 1$ read characters equal the pattern's prefix $P[. . . q]$. In particular, state $|P| - 1$ is active iff the last $|P|$ characters equal the full pattern.

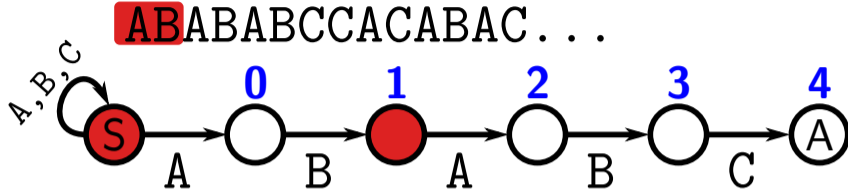
Reachable State Sets in NFA



Lemma (from previous lecture): NFA state set invariant

Let $A \subset Q$ be a set of active states of our NFA. Then, $q \in A \setminus \{-1\}$ iff the last $q + 1$ read characters equal the pattern's prefix $P[. . . q]$. In particular, state $|P| - 1$ is active iff the last $|P|$ characters equal the full pattern.

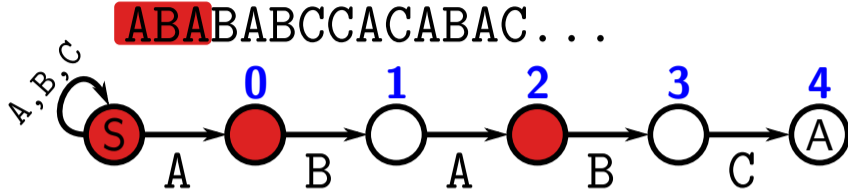
Reachable State Sets in NFA



Lemma (from previous lecture): NFA state set invariant

Let $A \subset Q$ be a set of active states of our NFA. Then, $q \in A \setminus \{-1\}$ iff the last $q + 1$ read characters equal the pattern's prefix $P[. \dots q]$. In particular, state $|P| - 1$ is active iff the last $|P|$ characters equal the full pattern.

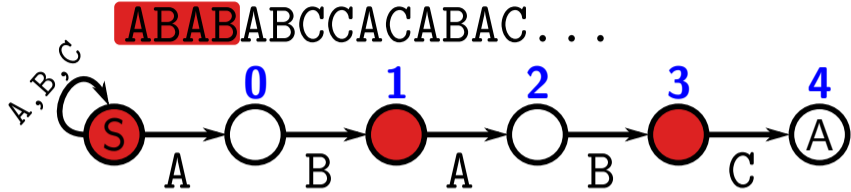
Reachable State Sets in NFA



Lemma (from previous lecture): NFA state set invariant

Let $A \subset Q$ be a set of active states of our NFA. Then, $q \in A \setminus \{-1\}$ iff the last $q + 1$ read characters equal the pattern's prefix $P[. \dots q]$. In particular, state $|P| - 1$ is active iff the last $|P|$ characters equal the full pattern.

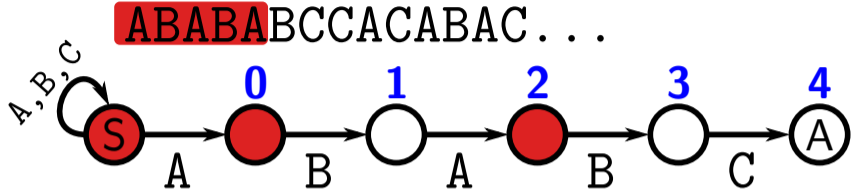
Reachable State Sets in NFA



Lemma (from previous lecture): NFA state set invariant

Let $A \subset Q$ be a set of active states of our NFA. Then, $q \in A \setminus \{-1\}$ iff the last $q + 1$ read characters equal the pattern's prefix $P[. . . q]$. In particular, state $|P| - 1$ is active iff the last $|P|$ characters equal the full pattern.

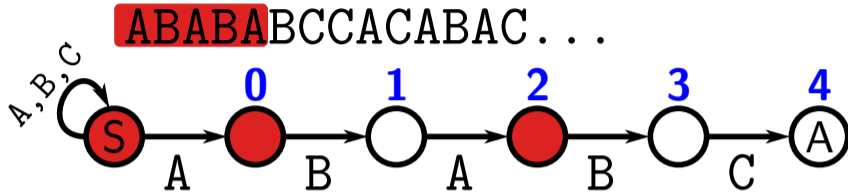
Reachable State Sets in NFA



Lemma (from previous lecture): NFA state set invariant

Let $A \subset Q$ be a set of active states of our NFA. Then, $q \in A \setminus \{-1\}$ iff the last $q + 1$ read characters equal the pattern's prefix $P[. . . q]$. In particular, state $|P| - 1$ is active iff the last $|P|$ characters equal the full pattern.

Reachable State Sets in NFA



Lemma (from previous lecture): NFA state set invariant

Let $A \subset Q$ be a set of active states of our NFA. Then, $q \in A \setminus \{-1\}$ iff the last $q + 1$ read characters equal the pattern's prefix $P[. \dots q]$. In particular, state $|P| - 1$ is active iff the last $|P|$ characters equal the full pattern.

Lemma

Let A be the set of active states of a pattern matching NFA. Now, let $a^* := \max A$. Then, A is completely determined by a^* .

Consequences of State Set Lemma

Lemma (from previous slide)

Let A be the set of active states of a pattern matching NFA. Now, let $a^* := \max A$. Then, A is completely determined by a^* .

Consequences

For each possible a^* from -1 to $m - 1$, there is exactly one reachable set of active states.

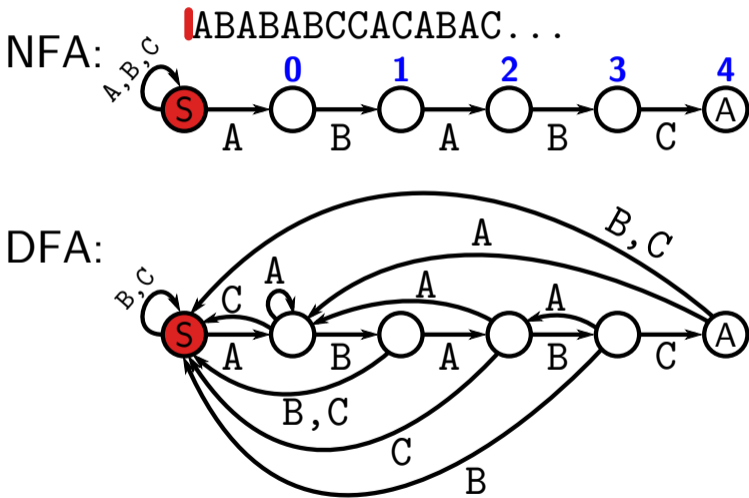
⇒ there are exactly $m + 1$ possible sets of active states.

⇒ there is an equivalent DFA with $m + 1$ states.

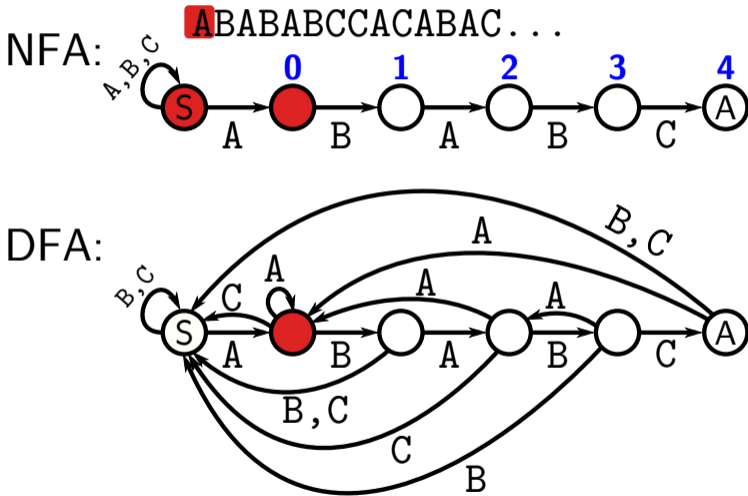
DFA state set

- We use same state set $Q = \{-1, \dots, m - 1\}$ for DFAs.
- Being in state $q \in Q$ in DFA \iff NFA has active states set A with $\max A = q$.

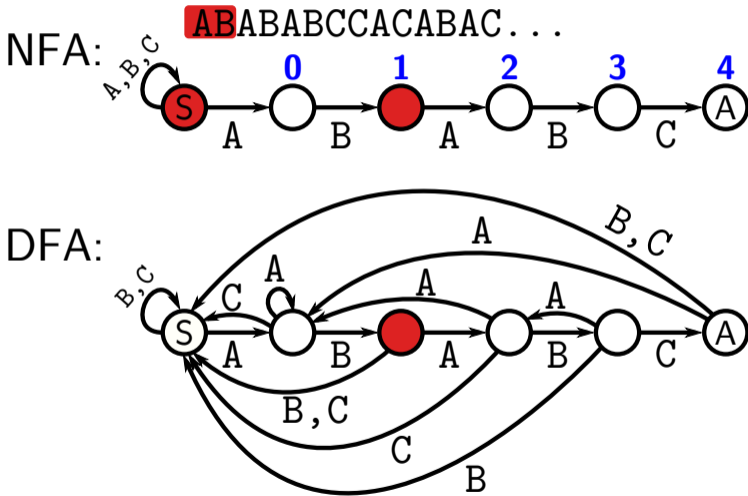
DFA-based vs. NFA-based Pattern Matching



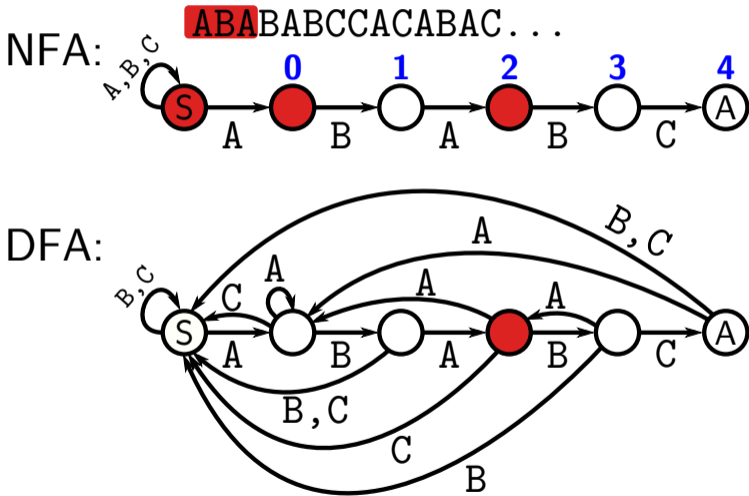
DFA-based vs. NFA-based Pattern Matching



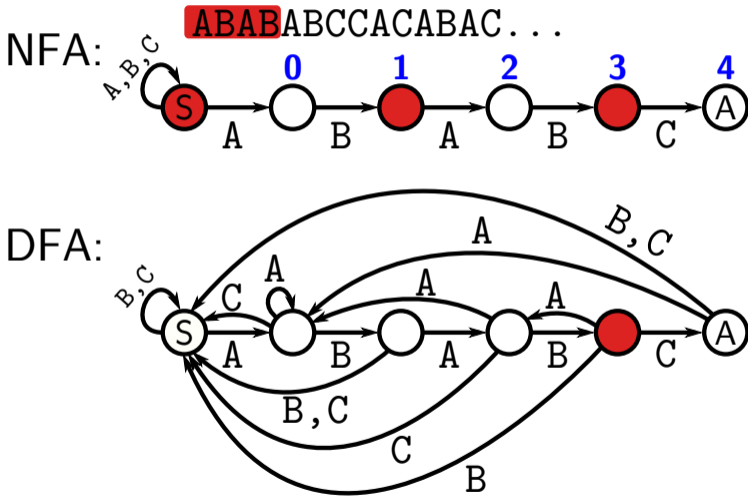
DFA-based vs. NFA-based Pattern Matching



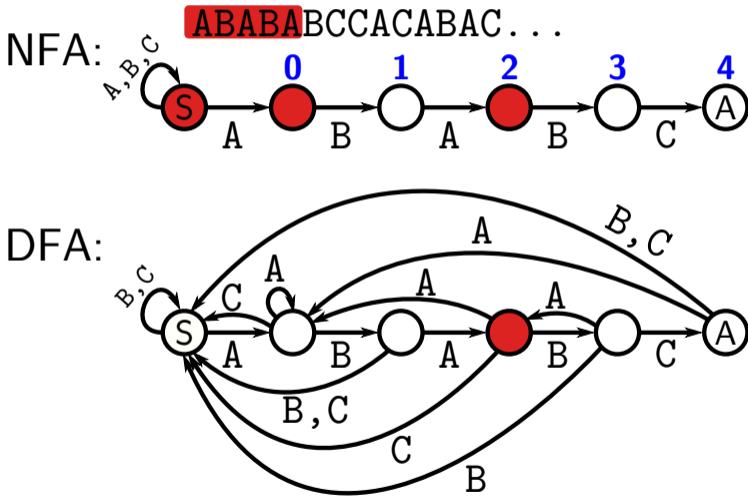
DFA-based vs. NFA-based Pattern Matching



DFA-based vs. NFA-based Pattern Matching



DFA-based vs. NFA-based Pattern Matching



DFA-based Pattern Matching: Overview

Algorithm overview

- 1 Construct pattern matching NFA in $O(m)$ time.
- 2 Construct DFA by computing the set of active states for every possible a^* in $O(m^2)$ time.
- 3 Use DFA for searching text in $O(n)$ time.

DFA-based Pattern Matching: Overview

Algorithm overview

- 1 Construct pattern matching NFA in $O(m)$ time.
- 2 Construct DFA by computing the set of active states for every possible a^* in $O(m^2)$ time.
- 3 Use DFA for searching text in $O(n)$ time.

Resulting **total time**: $O(m^2 + n)$

DFA-based Pattern Matching: Code

```
1 def DFA_with_delta(m, delta, T):
2     q = -1
3     for i in range(len(T)):
4         q = delta(q, T[i])
5         if q == m - 1:
6             yield (i-m+1, i+1)
7
8 def DFA(P, T):
9     delta = DFA_delta_table(P)
10    return DFA_with_delta(len(P), delta, T)
```


DFA-based Pattern Matching: Code

```
1 def DFA_with_delta(m, delta, T):
2     q = -1
3     for i in range(len(T)):
4         q = delta(q, T[i])
5         if q == m - 1:
6             yield (i-m+1, i+1)
7
8 def DFA(P, T):
9     delta = DFA_delta_table(P)
10    return DFA_with_delta(len(P), delta, T)
```

Summary

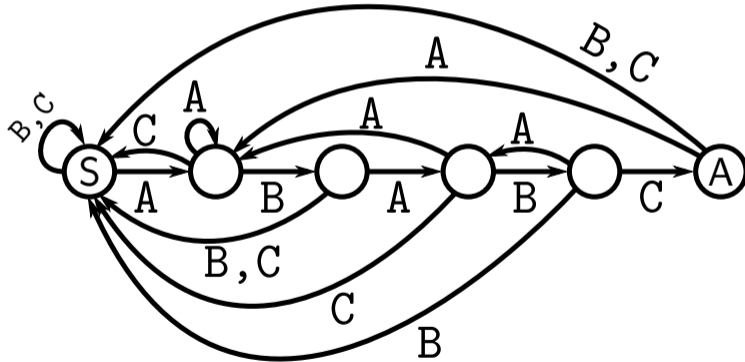
- Total time: $O(m^2 + n)$
- How do we get to $O(m + n)$?

Knuth-Morris-Pratt Algorithm

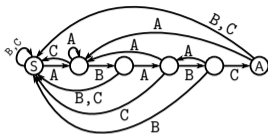
How to Compute the DFA Transition Function?

DFA transition function δ

- Matching character: just move right (easy)
- Non-matching character: move left (but how far?)



Example: Construction of DFA Transition Table (ABABC)



Using the *lps*-Table to Generate δ

Given: state q and character c

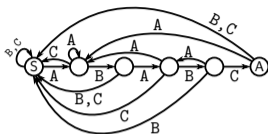
Approach to compute $\delta(q, c)$

- If $c = P[q + 1]$, then $q \mapsto q + 1$
- If not, try again for q' , where q' is “the next NFA state that would be active” (if this was an NFA and not a DFA)

The *lps*-function: $lps : \{0, \dots, m - 1\} \rightarrow \mathbb{N}$

- Recall: state q corresponds to prefix $P[\dots q]$ of length $q + 1$
- $lps(q)$ is the length of the **longest prefix** of P that is a true **suffix** of $P[\dots q]$.
- Then, we can get q' through $q' = lps(q) - 1$.

Example: *lps*-Table



Using lps to Compute δ

```
1 def DFA_delta_lps(q, c, P, lps):  
2     """state q, character c, pattern P, lps function/table"""  
3     m = len(P)  
4     while q == m-1 or (P[q+1] != c and q > -1):  
5         q = lps[q] - 1  
6     if P[q+1] == c: q += 1  
7     return q
```

Running time analysis

- Still takes $O(m)$ time in the worst case (while loop in line 4), leading to $O(mn)$ time for pattern matching over the whole text.

Using lps to Compute δ

```
1 def DFA_delta_lps(q, c, P, lps):
2     """state q, character c, pattern P, lps function/table"""
3     m = len(P)
4     while q == m-1 or (P[q+1] != c and q > -1):
5         q = lps[q] - 1
6     if P[q+1] == c: q += 1
7     return q
```

Running time analysis

- Still takes $O(m)$ time in the worst case (while loop in line 4), leading to $O(mn)$ time for pattern matching over the whole text.
- But m iterations cannot happen for all n text characters!
Amortized analysis: How many times can lines 4, 5 be executed in total?

Using lps to Compute δ

```
1 def DFA_delta_lps(q, c, P, lps):
2     """state q, character c, pattern P, lps function/table"""
3     m = len(P)
4     while q == m-1 or (P[q+1] != c and q > -1):
5         q = lps[q] - 1
6     if P[q+1] == c: q += 1
7     return q
```

Running time analysis

- Still takes $O(m)$ time in the worst case (while loop in line 4), leading to $O(mn)$ time for pattern matching over the whole text.
- But m iterations cannot happen for all n text characters!
Amortized analysis: How many times can lines 4, 5 be executed in total?
- Line 5 decreases q , but q cannot drop below -1 .
Only line 6 can ever increase q , at most once per iteration.

Knuth-Morris-Pratt Algorithm

```
1 def DFA_delta_lps(q, c, P, lps):
2     """state q, character c, pattern P, lps function/table"""
3     m = len(P)
4     while q == m-1 or (P[q+1] != c and q > -1):
5         q = lps[q] - 1
6     if P[q+1] == c: q += 1
7     return q
```

```
1 def KMP(P, T):
2     lps = compute_lps(P)
3     m, q = len(P), -1
4     for i in range(len(T)):
5         q = DFA_delta_lps(q, T[i], P, lps)
6         if q == m - 1: yield i
```

Running time: $O(n + m)$ since `DFA_delta_lps` takes **amortized** constant time.

Computing the lps-Table

```
1 def compute_lps(P):
2     m = len(P)
3     q = -1
4     lps = [0] * m # lps[0] = 0 is correct
5     for i in range(1, m):
6         while q > -1 and P[q+1] != P[i]:
7             q = lps[q] - 1
8         if P[q+1] == P[i]: q += 1
9         # Invariant (Q) holds here
10        lps[i] = q+1
11    return lps
```

Invariant (Q): $q = \max \{k < i : P[i - k \dots i] = P[0 \dots k]\}$

Summary: Knuth-Morris-Pratt Algorithm

Knuth-Morris-Pratt Algorithm

- lps-function gives a succinct representation of δ .
- Using lps to evaluate δ takes **amortized** constant time.
- Constructing lps-table works similar and takes $O(m)$ time.
- KMP algorithm has optimal running time of $O(m + n)$.

Historical Note

In the original paper (1977), the algorithm is not presented in terms of DFAs. However, the authors point out that *“it was still legitimate to conclude that automata theory had actually been helpful in this practical problem.”*

Summary

- Today's topic: **Exact Pattern Matching** (for single patterns without index)
- Reminder: NFAs and DFAs
- Avoid reading text characters more than once
 - **NFA-based pattern matching**
- Efficient bit-parallel implementation of pattern matching NFA
 - **Shift-And algorithm**
- Best asymptotic worst-case time:
 - **DFA-based algorithms**
- Compact representation of DFA transitions:
 - **Knuth-Morris-Pratt algorithm**

Possible exam questions

- How can finite automata be used to solve the pattern matching problem?
- What is the difference between NFAs and DFAs?
- What running times can be achieved in NFA/DFA based pattern matching?
- How is the set of active NFA states related to the read text so far?
- Give the formal definition of a pattern matching NFA and explain it.
- Explain the Shift-And algorithm.
- Explain the subset construction (from NFA to DFA).
- Why do the special NFAs studied here have the same number of states as the corresponding DFAs?
- Explain the Knuth-Morris-Pratt (KMP) algorithm and its relation to DFAs.
- How can one construct the *lps* function and in what time?
- What is the running time of KMP (worst case / best case)?