



Bit-Parallel Algorithms for Exact Search of Extended Patterns

Algorithms for Sequence Analysis

Sven Rahmann

Summer 2024

Review and Today's Lecture

Reminder

- Horspool's Algorithm:
 $O(mn)$ time worst-case, $O(n/m)$ best-case with long shifts
- Shift-And Algorithm:
 $O(mn/w)$ time, **bit-parallel**, w is register width (64 bits)

Review and Today's Lecture

Reminder

- Horspool's Algorithm:
 $O(mn)$ time worst-case, $O(n/m)$ best-case with long shifts
- Shift-And Algorithm:
 $O(mn/w)$ time, **bit-parallel**, w is register width (64 bits)

This lecture: More on bit-parallel algorithms

- 1 How to get **longer shifts** than Horspool's algorithm?
→ **BNDM algorithm** (backward non-deterministic DAWG matching)
- 2 Bit-parallel algorithms for more **general patterns**

Backward Non-Deterministic DAWG Matching

Reminder: Horspool's Algorithm

Horspool shift function

Text: ??????A????? ??????B????? ??????C??????
 └──────────┘ └────────┘ └────────┘ └────────┘ └────────┘
Pattern: BAAAAAB BAAAAAB BAAAAAB

Approach

- Compare characters **from right to left** in current window
- Shift window based on **last character only**

Weak points

Small alphabets lead to short shifts (especially bad for long patterns).

Substring-based Shifts

Ideas

- Read from right to left (like Horspool)
- **Read on** after mismatch to achieve longer shifts
- When current **substring** of window is not a **substring** of pattern, window can be shifted beyond that **substring**.
- Tracking **window suffixes** that are **pattern prefixes** further increases shifts.

Sought: Data Structure

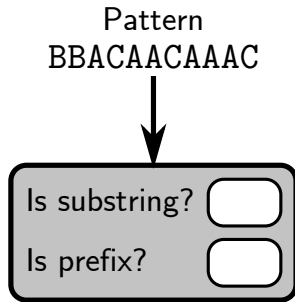
Requirements for supported queries

- Add characters from right to left
- Is window read so far a **substring** of the pattern?
- Is window read so far a **prefix** of the pattern?

Sought: Data Structure

Requirements for supported queries

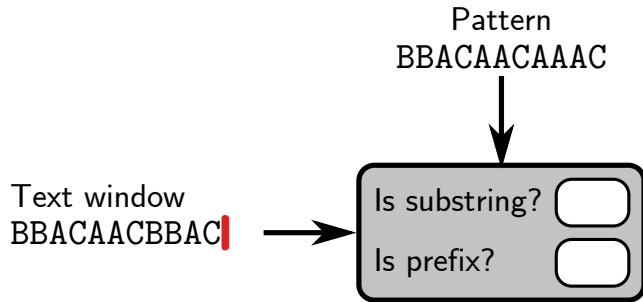
- Add characters from right to left
- Is window read so far a **substring** of the pattern?
- Is window read so far a **prefix** of the pattern?



Sought: Data Structure

Requirements for supported queries

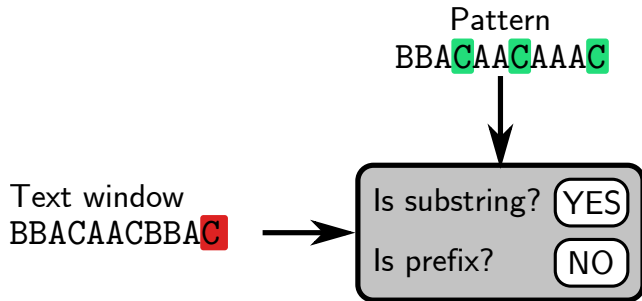
- Add characters from right to left
- Is window read so far a **substring** of the pattern?
- Is window read so far a **prefix** of the pattern?



Sought: Data Structure

Requirements for supported queries

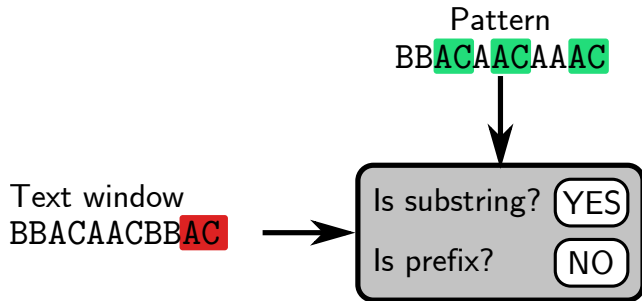
- Add characters from right to left
- Is window read so far a **substring** of the pattern?
- Is window read so far a **prefix** of the pattern?



Sought: Data Structure

Requirements for supported queries

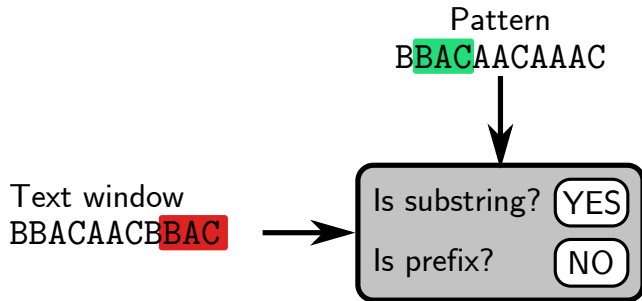
- Add characters from right to left
- Is window read so far a **substring** of the pattern?
- Is window read so far a **prefix** of the pattern?



Sought: Data Structure

Requirements for supported queries

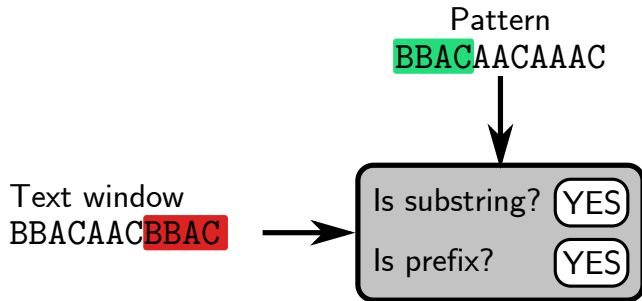
- Add characters from right to left
- Is window read so far a **substring** of the pattern?
- Is window read so far a **prefix** of the pattern?



Sought: Data Structure

Requirements for supported queries

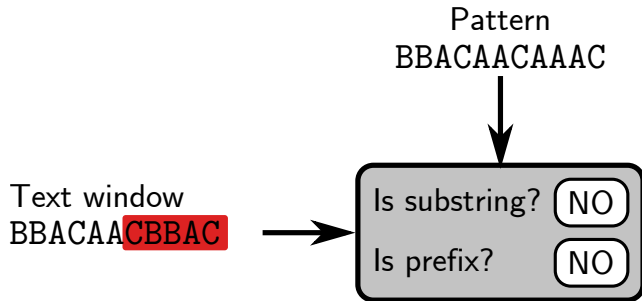
- Add characters from right to left
- Is window read so far a **substring** of the pattern?
- Is window read so far a **prefix** of the pattern?



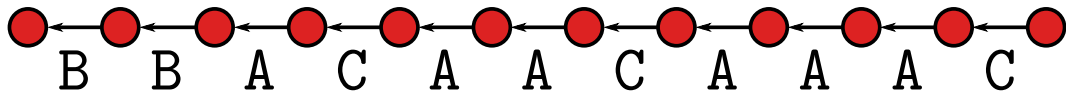
Sought: Data Structure

Requirements for supported queries

- Add characters from right to left
- Is window read so far a **substring** of the pattern?
- Is window read so far a **prefix** of the pattern?



Solution: Non-deterministic suffix automaton



Pattern
BBACAACAAAC

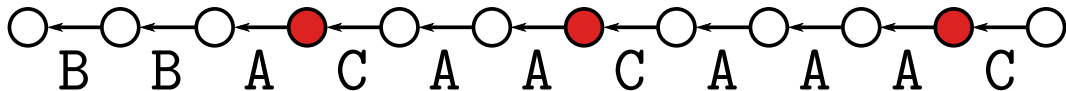
Text window
BBACAACBBAC



Is substring? ☐

Is prefix? ☐

Solution: Non-deterministic suffix automaton



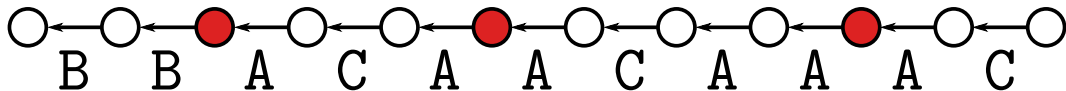
Pattern
BBA**CAAC**AAAC

Text window
BBACAACBBAC**C**



Is substring? YES
Is prefix? NO

Solution: Non-deterministic suffix automaton



Pattern
BBACAACAAAC

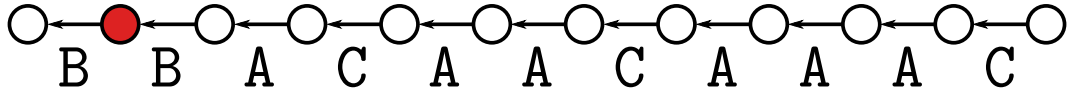
Text window
BBACAACBBAC



Is substring? YES

Is prefix? NO

Solution: Non-deterministic suffix automaton



Pattern
B **BAC** AACAAC

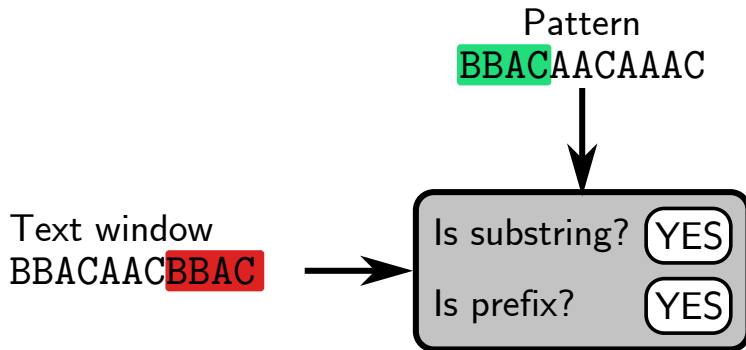
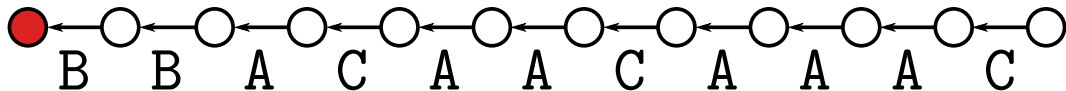
Text window
BBACAACB **BAC**



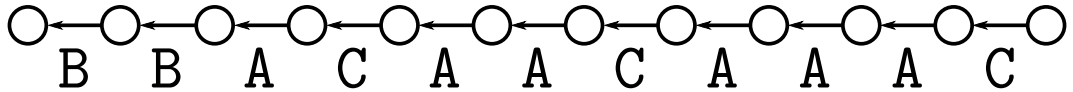
Is substring? YES

Is prefix? NO

Solution: Non-deterministic suffix automaton



Solution: Non-deterministic suffix automaton



Pattern
BBACAACAAAC

Text window

BBACAA**CBBAC**



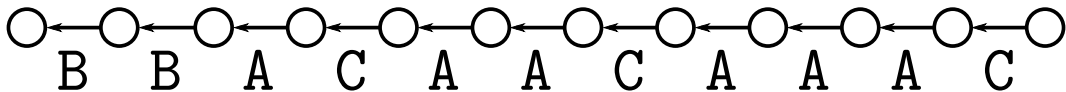
Is substring?

NO

Is prefix?

NO

Non-Deterministic Suffix Automaton



Construction

- Construct Pattern Search NFA of **reverse pattern**.
- Do not use a self-loop inside the first (rightmost) state.
- **All** states are start states.

Usage

- Use Shift-And approach to maintain set of active states
- **Any** state is active \Leftrightarrow substring occurs in pattern
- **Accept** state is active \Leftrightarrow found prefix
- **Accept** state is active **and** $|P|$ characters processed \Leftrightarrow found match

BNDM Algorithm

BNDM Algorithm Outline

For **each window**:

- 1 **Initialize** suffix automaton (all states active)
- 2 Read window from **right to left** until no states active or full window read.
- 3 Keep track of **longest window suffix** that is a **pattern prefix**.
- 4 **Shift** window to **align** this suffix with pattern prefix

BNDM Algorithm: Code

```
def BNDM(P, T):
    masks, accept_state = compute_masks(P[::-1])
    n, m, pos = len(T), len(P), len(P)
    while pos <= n:
        j, lastsuffix, A = 1, 0, (1 << m) - 1
        while A != 0:
            A &= masks[T[pos-j]]          # update A
            if A & accept_state != 0:      # accept state?
                if j == m:                  # full pattern found?
                    yield (pos - m, pos)
                    break
                else:                       # no, found proper prefix
                    lastsuffix = j          # store suffix
            j += 1; A = A << 1             # go to next window character
        pos += m - lastsuffix              # shift window
```

Deterministic Counterpart: BDM

Names

- **BDM**: Backward DAWG Matching
- **BNDM**: Backward Non-deterministic DAWG Matching
- **DAWG**: Directed Acyclic Word Graph

BDM (Backward DAWG Matching) Algorithm

- As before, we could transform the NFA into a DFA
→ **deterministic suffix automaton** (a DAWG)
- Use the subset construction (can be inefficient),
use a suffix tree (later),
or really just used BNDM instead of BDM (if $|P| \leq 64$)

Bit-Parallel Algorithms for Extended Patterns

Overview

So far, patterns were simple strings, $P \in \Sigma^*$.

For several applications (e.g., transcription factor binding sites on DNA), it is necessary to consider patterns that allow

- different characters (some subset of Σ) at some positions,
- variable-length runs of arbitrary characters,
- optional characters at some positions.

All of these patterns are subsets of **regular expressions**, which can be searched for by DFAs.

Overview

So far, patterns were simple strings, $P \in \Sigma^*$.

For several applications (e.g., transcription factor binding sites on DNA), it is necessary to consider patterns that allow

- different characters (some subset of Σ) at some positions,
- variable-length runs of arbitrary characters,
- optional characters at some positions.

All of these patterns are subsets of **regular expressions**, which can be searched for by DFAs.

Why talk about it?

Specialized **bit-parallel** implementations for each pattern class are more efficient. All of the above patterns can be recognized by variations of the **Shift-And** algorithm.

Generalized Strings

- A **generalized string** over Σ is a string over $2^\Sigma \setminus \{\emptyset\}$, i.e., a string whose characters are non-empty subsets of Σ .

Generalized Strings

- A **generalized string** over Σ is a string over $2^\Sigma \setminus \{\emptyset\}$, i.e., a string whose characters are non-empty subsets of Σ .
- **Example:** Consider the set $\{\text{Maier, Meier, Mayer, Meyer}\}$. It can be written as a single generalized string: $\{M\}\{a,e\}\{i,y\}\{e\}\{r\}$. Shorthand: $M[ae][iy]er$

Generalized Strings

- A **generalized string** over Σ is a string over $2^\Sigma \setminus \{\emptyset\}$, i.e., a string whose characters are non-empty subsets of Σ .
- **Example:** Consider the set $\{\text{Maier, Meier, Mayer, Meyer}\}$.
It can be written as a single generalized string: $\{M\}\{a,e\}\{i,y\}\{e\}\{r\}$.
Shorthand: $M[ae][iy]er$
- **Notation:** Singleton sets are represented by their unique element.
Larger sets are represented by square brackets: $[ae]$ for $\{a,e\}$.
We write $\#$ for $\Sigma \in 2^\Sigma$ (“any character”).

Generalized Strings

- A **generalized string** over Σ is a string over $2^\Sigma \setminus \{\emptyset\}$, i.e., a string whose characters are non-empty subsets of Σ .
- **Example:** Consider the set { Maier, Meier, Mayer, Meyer }.
It can be written as a single generalized string: { M } { a,e } { i,y } { e } { r }.
Shorthand: M[ae][iy]er
- **Notation:** Singleton sets are represented by their unique element.
Larger sets are represented by square brackets: [ae] for { a,e }.
We write # for $\Sigma \in 2^\Sigma$ (“any character”).
- In DNA sequences, the IUPAC code specifies a one-letter code for each subset:
size 1: ACGT; size 2: SWRYKM; size 3: BDHV; size 4: N.

The Shift-And Algorithm for Generalized Strings

- Recall the Shift-And algorithm with active state bits D :
 $D \leftarrow ((D \ll 1) | 1) \& \text{mask}(c)$
- The Shift-And algorithm can process generalized strings without modifications.
- The bit masks tell which characters are allowed at which position.
It is no problem that more than one bit is set at some positions.

The Shift-And Algorithm for Generalized Strings

- Recall the Shift-And algorithm with active state bits D :
 $D \leftarrow ((D \ll 1) | 1) \& \text{mask}(c)$
- The Shift-And algorithm can process generalized strings without modifications.
- The bit masks tell which characters are allowed at which position.
It is no problem that more than one bit is set at some positions.
- **Example:** $P = \text{abba\#b}$ over $\Sigma = \{a, b\}$.

	b#abba (reversed because of bit numbers)
<i>mask</i> [a]	011001
<i>mask</i> [b]	110110

Bounded-length Runs of Arbitrary Characters

- A **run of arbitrary characters** is a sequence of Σ s (written as #s) in a generalized string.
- We allow **variable run lengths**, but with fixed **lower and upper bounds**.
- **Notation:** $\#(L, U)$ with lower bound L and upper bound U
- **Example:** $P = \text{bba}\#(1,3)\text{a}$:
After bba, we have one to three arbitrary characters, followed by a.

Bounded-length Runs of Arbitrary Characters

- A **run of arbitrary characters** is a sequence of Σ s (written as #s) in a generalized string.
- We allow **variable run lengths**, but with fixed **lower and upper bounds**.
- **Notation:** $\#(L, U)$ with lower bound L and upper bound U
- **Example:** $P = \text{bba}\#(1,3)\text{a}$:
After bba, we have one to three arbitrary characters, followed by a.
- Three restrictions:
 - 1 An element $\#(L, U)$ does not appear first or last in the pattern.
(We could remove them without substantially changing the pattern.)
 - 2 No two such elements appear next to each other.
(No problem, just add them: $\#(L, U)\#(L', U') \hat{=} \#(L + L', U + U')$.)
 - 3 We require $1 \leq L \leq U$.
(Allowing $L = 0$ is technically more challenging.)

An NFA for Bounded-length Runs of Arbitrary Characters

- Before considering a bit-parallel implementation, we design an NFA.
- We need ϵ -transitions, an extension of the standard NFA definition: ϵ -transitions happen instantaneously, without consuming a character.

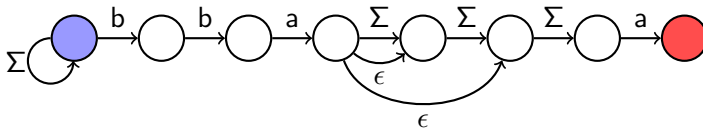
An NFA for Bounded-length Runs of Arbitrary Characters

- Before considering a bit-parallel implementation, we design an NFA.
- We need ϵ -transitions, an extension of the standard NFA definition: ϵ -transitions happen instantaneously, without consuming a character.
- The ϵ -transitions allow us to skip the optional characters.
For technical reasons, they **exit the initial state** of the run;
the **first** #s in each run are optional.
(One could do it differently, but that would be harder to implement!)

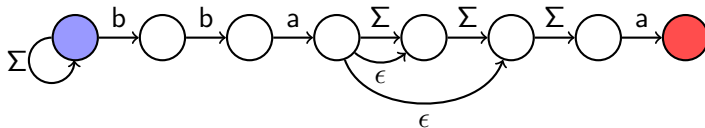
An NFA for Bounded-length Runs of Arbitrary Characters

- Before considering a bit-parallel implementation, we design an NFA.
- We need ϵ -transitions, an extension of the standard NFA definition: ϵ -transitions happen instantaneously, without consuming a character.
- The ϵ -transitions allow us to skip the optional characters.
For technical reasons, they **exit the initial state** of the run;
the **first** #s in each run are optional.
(One could do it differently, but that would be harder to implement!)

Example: $P = \text{bba}\#(1,3)\text{a}$:

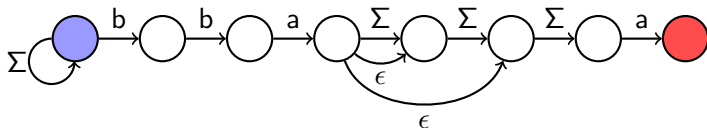


Bit-parallel Implementation



- We use the Shift-And algorithm on the maximal-length pattern as a basis. Then we additionally implement the ϵ -transitions.
- Masks are constructed as before (for #: 1-bits for each character).

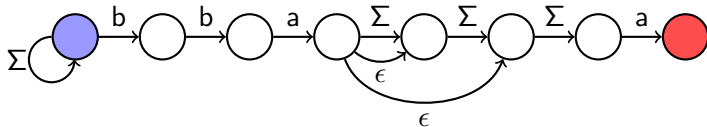
Bit-parallel Implementation



- We use the Shift-And algorithm on the maximal-length pattern as a basis. Then we additionally implement the ϵ -transitions.
- Masks are constructed as before (for #: 1-bits for each character).
- **Example:** $P = \text{bba}\#(1,3)\text{a}$ with $\Sigma = \{a, b, c\}$:

	a###abb
$\text{mask}[a]$	1111100
$\text{mask}[b]$	0111011
$\text{mask}[c]$	0111000

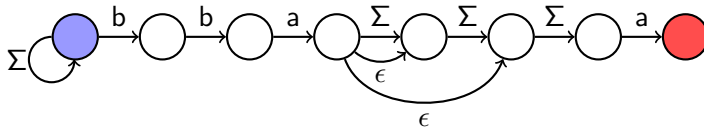
Implementation of ϵ -Transitions



- ϵ -transitions are instantaneous:

Whenever a state with outgoing ϵ -transitions becomes active (1-bit), this is immediately propagated to the targets of the outgoing ϵ -edges; these are by construction adjacent to the source state.

Implementation of ϵ -Transitions



- ϵ -transitions are instantaneous:
Whenever a state with outgoing ϵ -transitions becomes active (1-bit), this is immediately propagated to the targets of the outgoing ϵ -edges; these are by construction adjacent to the source state.
- The actual propagation of 1-bits will be achieved by subtraction (next slide).
- We use two additional bit masks:
 - Bit mask I marks states with outgoing ϵ -transitions.
 - Bit mask F marks the state after the target of the last ϵ -transition of each run.

a###abb

F 0100000

I 0000100

Propagation of Ones

- Let A be the bit mask of active states. Then $A \& I$ selects active I -states.
- Subtraction $F - (A \& I)$ propagates 1-bits, zeroes F -bits.

$$\begin{array}{r} F \quad 0100000 \\ A \& I \quad 0000100 \\ \hline - \quad 0011100 \end{array}$$

Propagation of Ones

- Let A be the bit mask of active states. Then $A \& I$ selects active I -states.
- Subtraction $F - (A \& I)$ propagates 1-bits, zeroes F -bits.

$$\begin{array}{r} F \quad 0100000 \\ A \& I \quad 0000100 \\ \hline - \quad 0011100 \end{array}$$

- **Problem:** Inactive I -states keep corresponding F -bit set:

$$\begin{array}{r} F \quad 010000100000 \\ A \& I \quad 000000000100 \\ \hline - \quad 010000011100 \end{array}$$

Propagation of Ones (Continued)

- **Problem:** Inactive I -states keep corresponding F -bit set.
- **Solution:** Zero out F -bits by a bitwise and with the negation of F :

F	010000100000
$A \& I$	000000000100
<hr/>	
$F - (A \& I)$	010000011100

Propagation of Ones (Continued)

- **Problem:** Inactive I -states keep corresponding F -bit set.
- **Solution:** Zero out F -bits by a bitwise and with the negation of F :

F	010000100000
$A \& I$	000000000100
<hr/>	
$F - (A \& I)$	010000011100
$\sim F$	101111011111
<hr/>	
$(F - (A \& I)) \& \sim F$	000000011100

Propagation of Ones (Continued)

- **Problem:** Inactive I -states keep corresponding F -bit set.
- **Solution:** Zero out F -bits by a bitwise and with the negation of F :

F	010000100000
$A \& I$	000000000100
<hr/>	
$F - (A \& I)$	010000011100
$\sim F$	101111011111
<hr/>	
$(F - (A \& I)) \& \sim F$	000000011100

- The resulting modified Shift-And update is thus:

- 1 Apply standard Shift-And update:

$$A = ((A \ll 1) \mid 1) \& \text{mask}[c]$$

- 2 Propagate active I -states along ϵ -transitions:

$$A = A \mid ((F - (A \& I)) \& \sim F)$$

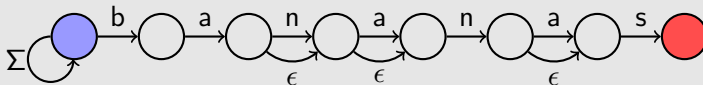
Patterns with Optional Characters

- Another modification of the Shift-And algorithm allows optional characters.
- **Notation:** Write ? after the optional character.
- **Example:** The set {color, colour} becomes $P = \text{colou?r}$.
- Consecutive optional characters (“blocks”) are allowed.

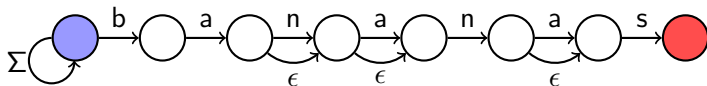
Patterns with Optional Characters

- Another modification of the Shift-And algorithm allows optional characters.
- **Notation:** Write ? after the optional character.
- **Example:** The set {color, colour} becomes $P = \text{colou?r}$.
- Consecutive optional characters (“blocks”) are allowed.

Larger example: $P = \text{ban?a?na?s}$ and $T = \text{banabanns}$



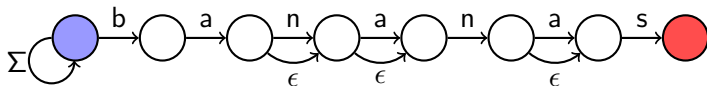
Bit-Parallel Implementation of Optional Characters



<i>I</i> :	0	1	0	0	1	0	0
<i>F</i> :	0	0	0	1	0	1	0
<i>O</i> :	0	0	1	1	0	1	0

- Three bit masks:
I: block start; *O*: targets of ϵ -transitions; *F*: block end
- Note: actual bit patterns are reversed (bit numbering vs. state numbering).

Bit-Parallel Implementation of Optional Characters



<i>I</i> :	0	1	0	0	1	0	0
<i>F</i> :	0	0	0	1	0	1	0
<i>O</i> :	0	0	1	1	0	1	0

- Three bit masks:
I: block start; *O*: targets of ϵ -transitions; *F*: block end
- Note: actual bit patterns are reversed (bit numbering vs. state numbering).
- Activity of any state within a block must be propagated to the block's end.

Bit-Parallel Implementation of Optional Characters (Continued)

- Activity of any state within a block must be propagated to the block's end:
Propagate the lowest active bit within a block up to the F -bit.

Bit-Parallel Implementation of Optional Characters (Continued)

- Activity of any state within a block must be propagated to the block's end: Propagate the lowest active bit within a block up to the F -bit.
- Consider how 1-bit propagation via subtraction works:

$$\begin{array}{r} 1101010000 \\ - \quad \quad \quad 1 \\ \hline 1101001111 \end{array}$$

$$\begin{array}{r} 1101011000 \\ - \quad \quad \quad 100 \\ \hline 1101010100 \end{array}$$

- Bits to the left (green) and to the right (black) are unchanged; only bits between the rightmost ones in the current block change (red).

Bit-Parallel Implementation of Optional Characters (Continued)

- Activity of any state within a block must be propagated to the block's end: Propagate the lowest active bit within a block up to the F -bit.
- Consider how 1-bit propagation via subtraction works:

$$\begin{array}{r} 1101010000 \\ - \quad \quad \quad 1 \\ \hline 1101001111 \end{array}$$

$$\begin{array}{r} 1101011000 \\ - \quad \quad \quad 100 \\ \hline 1101010100 \end{array}$$

- Bits to the left (green) and to the right (black) are unchanged; only bits between the rightmost ones in the current block change (red).
- We develop the machinery by example:

A .0010100.

I .0000001.

O .1111110.

F .1000000.

|

>> .1111100.

A .0010100.

A|F .1010100.

(A|F)-I .1010011.

((A|F)-I)=(A|F) .1111000.

O&((A|F)-I)=(A|F) .1111000.

A|(O&((A|F)-I)=(A|F)) .1111100.

Bit-Parallel Implementation of Optional Characters (Conclusion)

- **Note:** Bitwise equality $X = Y$ can be implemented as $\sim(X \oplus Y)$.
- Full implementation:
 - 1 Create masks for all characters;
treat optional characters as regular characters.
 - 2 Standard Shift-And update of active states A :
$$A = ((A \ll 1) \mid 1) \& \text{mask}[c]$$
 - 3 Propagate active states over optional characters:
$$AF = A \mid F$$
$$A = A \mid (0 \& (\sim(AF - I) \wedge AF))$$

(Here, \wedge denotes the xor-operation.)

Summary

Topic

Bit-parallel methods for exact pattern matching of single patterns without text indexing

Properties of bit-parallel algorithms

- Typically only applicable if an “almost linear” NFA recognizes the pattern, and if this NFA has at most 64 (register width) states
- Shift-And approach is simple and very flexible, extends to general patterns; running time is always $O(n)$ for constant $|P| < 64$.
- BNDM approach is also simple and flexible; may pathologically use $O(mn)$ time even for constant $m = |P| < 64$, but has best-case running time of $O(m + n/m)$.

Possible exam questions

- Explain the idea of bit-parallel simulation of NFAs.
- Explain the suffix automaton and the BNDM algorithm.
- What are the advantages of BNDM over Horspool's algorithm?
- What are the advantages of BNDM over the Shift-And algorithm?
- What is a generalized string?
- How does the Shift-And algorithm change when you allow generalized strings?
- Why would you want to use the Shift-And algorithm for runs with bounded length, when the algorithms for optional characters is more general ($\#(3, 5) = \#?#?###$)?
- How do you implement bit-parallel propagation of an active state?

Conclusion

Topic

Exact pattern matching of single patterns without text indexing

Strengths of different algorithms

- **Shift-And:** simple, applicable if $|P| \leq 64$
 - **B(N)DM:** for $|P| < 64$; best case of $O(m + n/m)$;
long shifts even for small alphabet + long pattern
 - **Horspool:** best case of $O(m + n/m)$ for large alphabet + long pattern
 - **Knuth-Morris-Pratt (KMP):** best worst-case time of $O(m + n)$
-
- Automata theory was very useful.
 - Next topic: index structures (i.e. preprocessing the text)