



# Announcements

---

1. **NO QUIZ THIS WEEK!**
2. Programming Assignment Posted
  - **TESTING CODE UPDATED**
    - small bug in tritonic array generation
    - download new version
  - Due Wednesday, 13 November
3. Attendance Code:



289105

# Meeting Goals

---

Discuss String Matching procedures:

- Brute Force
- DFA procedure
- Knuth-Morris-Pratt

# String Matching

# The String Matching Problem

---

**Input:**


$\Sigma = \text{alphabet}$

- A **text**  $T \in \Sigma^*$  of length  $n$
- A **pattern**  $P \in \Sigma^*$  of length  $m$  (typically  $m \ll n$ )

**Output:**

- The index of the **first occurrence** of  $P$  in  $T$ , or  $-1$  if  $T$  does not contain  $P$  as a substring:
  - $\min \{i \mid T[i, i+m) = P\}$

**Example.**

- $T = 101100$ 1101 $1101$    
0 1 2 3 4 5 6
- $P_1 = 1101$ 
  - Output:  $i \leftarrow 6$
- $P_2 = 000$ 
  - Output:  $i \leftarrow -1$

# Brute Force Matching

# Brute Force Matching

---

**Guess** an index  $i$  where a match might occur

- Possible guesses  $i = 0, 1, \dots, n - m - 1$

**Check** if match at  $i$ :

- is  $T[\underline{i}, i + m) = \underline{P}$ ?
- verify each character individually

**Cost** = number of comparisons made

# Brute Force Matching

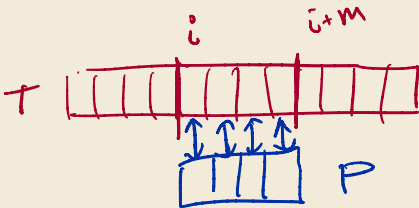
**Guess** an index  $i$  where a match might occur

- Possible guesses  $i = 0, 1, \dots, n - m - 1$

**Check** if match at  $i$ :

- is  $T[i, i + m) = P$ ?
- verify each character individually

```
1: procedure VERIFYMATCH( $T, P, i$ )  
2:    $j \leftarrow 0$   
3:   while  $j < m$  do  
4:     if  $T[i + j) \neq P[j)$  then  
5:       return FALSE  
6:     end if  
7:      $j \leftarrow j + 1$   
8:   end while  
9:   return TRUE  
10: end procedure
```



**Cost** = number of comparisons made



# Brute Force Matching

**Guess** an index  $i$  where a match might occur

- Possible guesses  $i = 0, 1, \dots, n - m - 1$

**Check** if match at  $i$ :

- is  $T[i, i + m) = P$ ?
- verify each character individually

```
1: procedure VERIFYMATCH( $T, P, i$ )
2:    $j \leftarrow 0$ 
3:   while  $j < m$  do
4:     if  $T[i + j) \neq P[j)$  then
5:       return FALSE
6:     end if
7:      $j \leftarrow j + 1$ 
8:   end while
9:   return TRUE
10: end procedure
```

**Cost** = number of comparisons made

Parameters:  
 $n$  = size of  $T$   
 $m$  = size of  $P$

## PollEverywhere Question

What are the worst case and best case running times of VERIFYMATCH?



[pollev.com/comp526](https://pollev.com/comp526)

# Brute Force Matching

**Guess** an index  $i$  where a match might occur

- Possible guesses  $i = 0, 1, \dots, n - m - 1$

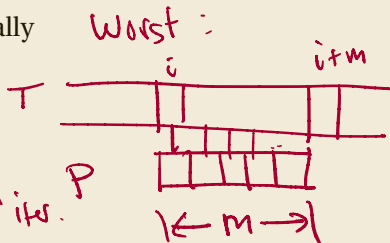
**Check** if match at  $i$ :

- is  $T[i, i + m) = P$ ?
- verify each character individually

```
1: procedure VERIFYMATCH( $T, P, i$ )
2:    $j \leftarrow 0$ 
3:   while  $j < m$  do
4:     if  $T[i + j] \neq P[j]$  then
5:       return FALSE
6:     end if
7:      $j \leftarrow j + 1$ 
8:   end while
9:   return TRUE
10: end procedure
```

**Cost** = number of comparisons made

**Best and Worst Cases:**



$\Theta(m) = \text{worst}$

Best  $O(1)$

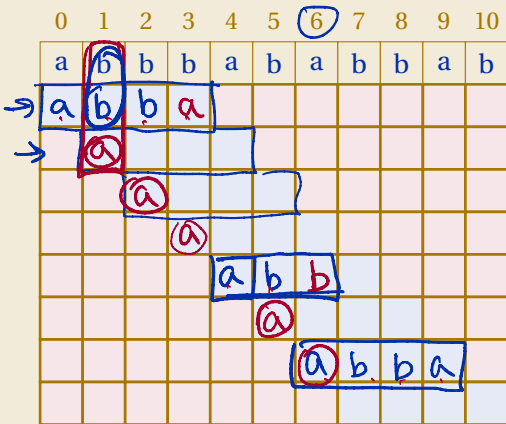
$\rightarrow T[i] \neq P[0]$



# Brute Force Example

## Example

- $T = abbbababbab$
- $P = \underline{abba}$



**procedure**

BRUTEFORCEMATCH( $T, P$ )

**for**  $i = 0, 1, \dots, n - m - 1$  **do**

**if** VERIFYMATCH( $T, P, i$ ) **then**

**return**  $i$

**end if**

**end for**

**return**  $-1$

**end procedure**

← return 6

# Brute Force Efficiency

---

The **worst case** complexity of brute force search is  $\Theta(nm)$ ...  
...but when is this **actually** achieved?

$T = a a a a a \dots a$

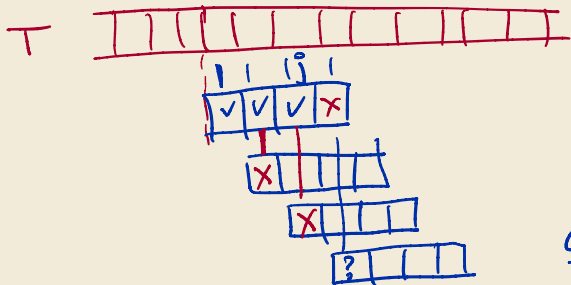
$P = a a a \dots a b$

# Brute Force Efficiency

The **worst case** complexity of brute force search is  $\Theta(nm)$ ...  
...but when is this **actually** achieved?

**Example.** Consider the case where  $P$  contains *no repeated characters*.

$$P = P_1 P_2 \dots P_m \neq$$



In this case  
never do  $> 2n$   
comparisons

$\leq 2n$  comparisons

# Brute Force Efficiency

---

The **worst case** complexity of brute force search is  $\Theta(nm)$ ...  
...but when is this **actually** achieved?

**Example.** Consider the case where  $P$  contains *no repeated characters*.

- Claim: brute force search running time is now  $O(n)$ 
  - In fact, at most  $2n$  comparisons made!
  - *Why?*

# Brute Force Efficiency

---

The **worst case** complexity of brute force search is  $\Theta(nm)$ ...  
...but when is this **actually** achieved?

**Example.** Consider the case where  $P$  contains *no repeated characters*.

- Claim: brute force search running time is now  $O(n)$ 
  - In fact, at most  $2n$  comparisons made!
  - *Why?*
- Which of these comparisons were unnecessary?
  - How can you search with fewer comparisons?



# Brute Force Efficiency

---

The **worst case** complexity of brute force search is  $\Theta(nm)$ ...  
...but when is this **actually** achieved?

**Example.** Consider the case where  $P$  contains *no repeated characters*.

- Claim: brute force search running time is now  $O(n)$ 
  - In fact, at most  $2n$  comparisons made!
  - *Why?*
- Which of these comparisons were unnecessary?
  - How can you search with fewer comparisons?

**More generally:** How can we use results of *previous comparisons* to avoid making unnecessary comparisons in the future?

- Goal: never re-read a character from  $T$ !

# Matching with a DFA

# Sliding Comparisons

---

## Example

- $T = \text{aabababbabacaa}$
- $P = \underline{\text{ababaca}}$

*a a b a b a b b a b a b a c a a*

## Idea:

- Scan through  $T$  keeping track of current matches
- Each new character  $T$  read, compare it to next character of  $P$
- If mismatch slide  $P$  so that **longest prefix** of  $P$  matches

# Sliding Comparisons

---

## Example

- $T = \text{aabababbabacaa}$
- $P = \text{ababaca}$

$T$      $a$   $a$   $b$   $a$   $b$   $a$   $b$   $b$   $a$   $b$   $a$   $b$   $a$   $c$   $a$   $a$   
 $P$      $a$   $b$   ~~$a$   $b$   $a$   $c$   $a$~~

## Idea:

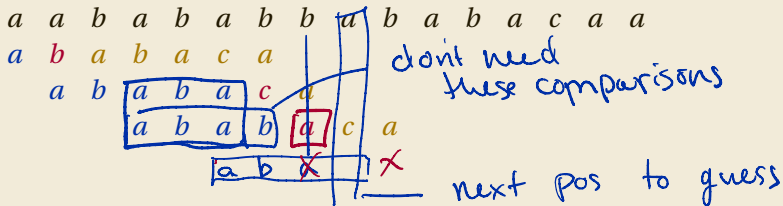
- Scan through  $T$  keeping track of current matches
- Each new character  $T$  read, compare it to next character of  $P$
- If mismatch slide  $P$  so that **longest prefix** of  $P$  matches



# Sliding Comparisons

## Example

- $T = \text{aabababbabacaa}$
- $P = \text{ababaca}$



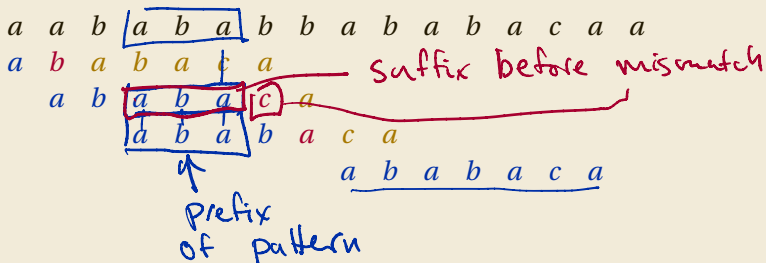
## Idea:

- Scan through  $T$  keeping track of current matches
- Each new character  $T$  read, compare it to next character of  $P$
- If mismatch slide  $P$  so that **longest prefix** of  $P$  matches

# Sliding Comparisons

## Example

- $T = \text{aabababbabacaa}$
- $P = \text{ababaca}$



## Idea:

- Scan through  $T$  keeping track of current matches
- Each new character  $T$  read, compare it to next character of  $P$
- If mismatch slide  $P$  so that **longest prefix** of  $P$  matches

# Representing States and Matches

---

**Question.** What information do we need to compute and store to determine next comparison?



# Representing States and Matches

---

**Question.** What information do we need to compute and store to determine next comparison?

- How many matches in  $P$  have we made so far?
- What what is the longest matching prefix for each possible next character in  $T$ 
  - if we read character  $x$ , how far do we need to “slide”  $P$  to match a prefix?

# Representing States and Matches

---

**Question.** What information do we need to compute and store to determine next comparison?

- How many matches in  $P$  have we made so far?
- What what is the longest matching prefix for each possible next character in  $T$ 
  - if we read character  $x$ , how far do we need to “slide”  $P$  to match a prefix?

**Information** to store

- **states** that represent number of matches with current prefix of  $P$
- **transitions** from current state to next states, depending on next character read from  $T$

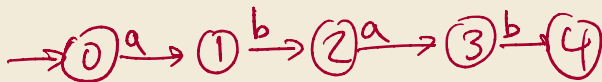
**Note.** This information depends *only* on the pattern  $P$ , not the text  $T$ .

# DFAs

A **Deterministic Finite Automaton (DFA)** consists of:

- A finite set  $Q$  of **states**
- A finite **alphabet**  $\Sigma$
- A **transition function**  $\delta : Q \times \Sigma \rightarrow Q$  —
- An **initial state**  $q_0 \in Q$
- A set  $F \subseteq Q$  of **accepting states**

rule: if  
in state  $q$   
and read  
character  $c$   
(in  $T$ ) which  
state to



move  
to

# DFAs

---

A **Deterministic Finite Automaton (DFA)** consists of:

- A finite set  $Q$  of **states**
- A finite **alphabet**  $\Sigma$
- A **transition function**  $\delta : Q \times \Sigma \rightarrow Q$
- An **initial state**  $q_0 \in Q$
- A set  $F \subseteq Q$  of **accepting states**

**Interpretation.** A DFA is used to determine if a string (text)  $T$  has some property (e.g., containing a pattern  $P$ ):

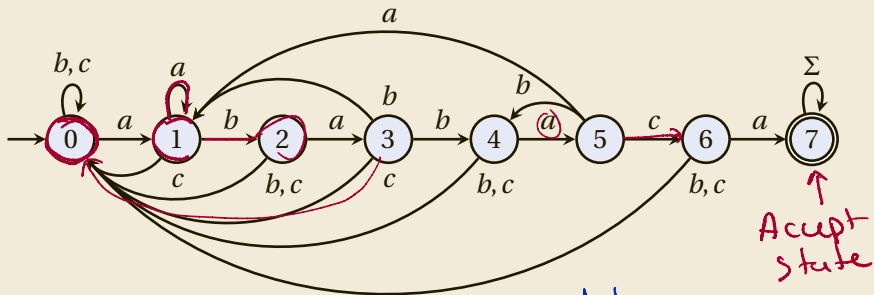
- Start at the state  $q_0$
- Read characters from  $T$  sequentially
  - if in state  $q$  and read character  $c$ , move to state  $\delta(q, c)$
- Return TRUE if end in “accepting” state

← transition function

# DFA Example

## Example

- T = aabacaababacaa ←
- P = ababaca



match

text	a	a	b	a	c	a	a	b	a	b	a	c	a	a
state	1	1	2	3	0	1	2	3	4	5	6	7	7	7

# DFA Efficiency

## PollEverywhere Question

Given a DFA for matching  $P[0, m)$  in  $T[0, n)$ , what is the running time of applying the DFA? Assume following links is  $O(1)$  time.

1.  $\Theta(nm)$
2.  $\Theta(n \log m)$
3.  $\Theta(n + m)$
4.  $\Theta(n)$



[pollev.com/comp526](https://pollev.com/comp526)

# DFA Efficiency

**Observe:** If we are given a DFA, executing it

- reads each character of  $T$  once
- updates state once per character

⇒ running time  $O(n)$

**So** the overall running time for pattern matching with a DFA is

$O(n)$  + time to build DFA

- assuming computation of  $\delta$  is  $O(1)$ .

## PollEverywhere Question

Given a DFA for matching  $P[0, m)$  in  $T[0, n)$ , what is the running time of applying the DFA? Assume following links is  $O(1)$  time.

1.  $\Theta(nm)$
2.  $\Theta(n \log m)$
3.  $\Theta(n + m)$
4.  $\Theta(n)$



[pollev.com/comp526](https://pollev.com/comp526)

# DFA Efficiency

**Observe:** If we are *given* a DFA, executing it

- reads each character of  $T$  once
- updates state once per character

⇒ running time  $O(n)$

**So** the overall running time for pattern matching with a DFA is  $O(n)$  + time to build DFA

- assuming computation of  $\delta$  is  $O(1)$ .

**But** how do we build the DFA?

???

## PollEverywhere Question

Given a DFA for matching  $P[0, m]$  in  $T[0, n]$ , what is the running time of applying the DFA? Assume following links is  $O(1)$  time.

1.  $\Theta(nm)$
2.  $\Theta(n \log m)$
3.  $\Theta(n + m)$
4.  $\Theta(n)$



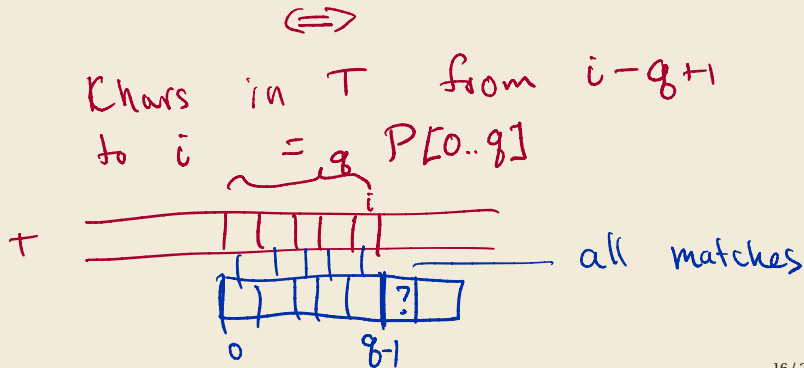
[pollev.com/comp526](https://pollev.com/comp526)



# DFA Interpretation & Construction

**Semantic Question.** What does it *mean* to be in state  $q$ ?

In state  $q$  after reading  $T[i]$



# DFA Interpretation & Construction

---

**Semantic Question.** What does it *mean* to be in state  $q$ ?

- Current position in  $T$  matches  $P$  up to the first  $q$  characters
- Symbolically  $T[j - q + 1, j] = P[0, q)$

**Question.** What happens when we read  $T[j + 1]$ ?

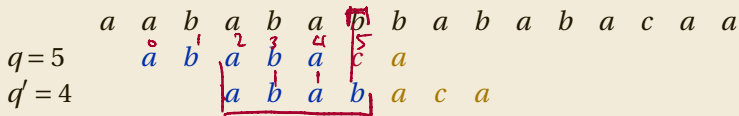
# DFA Interpretation & Construction

**Semantic Question.** What does it *mean* to be in state  $q$ ?

- Current position in  $T$  matches  $P$  up to the first  $q$  characters
- Symbolically  $T[j - q + 1, j] = P[0, q]$

**Question.** What happens when we read  $T[j + 1]$ ?

- If  $T[j + 1] = P[q]$ , transition to state  $q + 1$
- If  $T[j + 1] \neq P[q]$ , find the length  $q' \leq q$  of the longest prefix of  $P$  that matches  $T[j - q' + 1, j]$  that matches  $P[0, q']$



# DFA Interpretation & Construction

---

**Semantic Question.** What does it *mean* to be in state  $q$ ?

- Current position in  $T$  matches  $P$  up to the first  $q$  characters
- Symbolically  $T[j - q + 1, j] = P[0, q]$

**Question.** What happens when we read  $T[j + 1]$ ?

- If  $T[j + 1] = P[q]$ , transition to state  $q + 1$
- If  $T[j + 1] \neq P[q]$ , find the length  $q' \leq q$  of the longest prefix of  $P$  that matches  $T[j - q' + 1, j + 1]$  that matches  $P[0, q']$

	<i>a</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>	<i>a</i>
$q = 5$		<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>								
$q' = 4$				<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>						

- **Insight:** if  $T[j + 1] = c$  this is the same as matching  $P[0..q]$  against  $P[1..q]$ 
  - we can use the DFA constructed so far to find this!

# DFA Interpretation & Construction

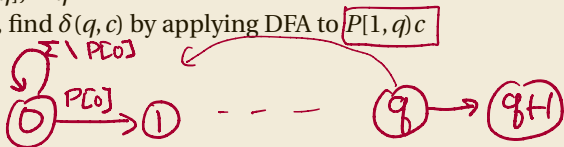
- **Insight:** if  $T[j+1] = c$  this is the same as matching  $P[0..q]$  against  $P[1..q]c$ 
  - we can use the DFA constructed so far to find this!

## Inductive Construction.

- Start with states 0 and 1 with

$$\delta(0, c) = \begin{cases} 1 & \text{if } P[0] = c \\ 0 & \text{otherwise.} \end{cases}$$

- Once we've constructed DFA up to state  $q$ :
  - take  $\delta(q, P[q]) = q+1$
  - for  $c \neq P[q]$ , find  $\delta(q, c)$  by applying DFA to  $P[1, q]c$



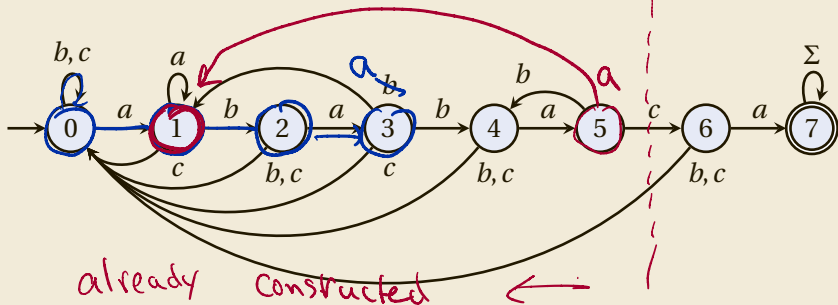
# DFA Interpretation & Construction

Todo: Fix img!

- **Insight:** if  $T[j+1] = c$  this is the same as matching  $P[0..q]$  against  $P[1..q]c$ 
  - we can use the DFA constructed so far to find this!
- Once we've constructed DFA up to state  $q$ :
  - take  $\delta(q, P[q]) = q+1$
  - for  $c \neq P[q]$ , find  $\delta(q, c)$  by applying DFA to  $P[1, q]c$

~~bababa~~

**Example.** Compute  $\delta(5, a)$  for  $P = ababaca$ .



# DFA Interpretation & Construction

---

- **Insight:** if  $T[j+1] = c$  this is the same as matching  $P[0..q]$  against  $P[1..q]c$ 
  - we can use the DFA constructed so far to find this!

## Inductive Construction.

- Start with states 0 and 1 with

$$\delta(0, c) = \begin{cases} 1 & \text{if } P[0] = c \\ 0 & \text{otherwise.} \end{cases}$$

- Once we've constructed DFA up to state  $q$ :
  - take  $\delta(q, P[q]) = q+1$
  - for  $c \neq P[q]$ , find  $\delta(q, c)$  by applying DFA to  $P[1, q]c$

## Analysis (idea).

- Argue by induction on  $q$  that the DFA enters state  $q$  on reading  $T[j]$  if and only if  $q$  is the largest number such that  $T[j-q+1, j] = P[0, q]$ .

# DFA Lookup Table Construction

---

DFA *diagrams* are great for humans, but not so great for computers...



# DFA Lookup Table Construction

---

DFA *diagrams* are great for humans, but not so great for computers...

## Problems.

1. How do we represent the DFA in a computer friendly format?
2. How do construct the DFA in that format efficiently?

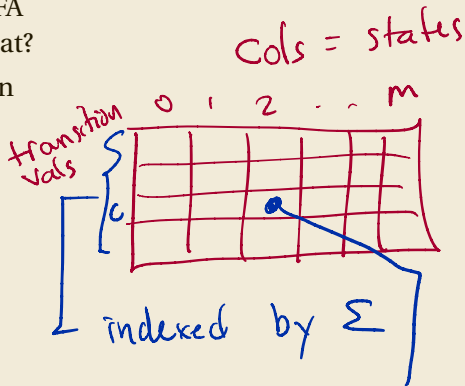
# DFA Lookup Table Construction

## Problems.

1. How do we represent the DFA in a computer friendly format?
2. How do construct the DFA in that format efficiently?

## Solutions.

1. Store a **lookup table**  $\delta[][]$ 
  - columns = states, rows = characters
  - $\delta[q][c] \leftarrow \delta(q, c)$



$\delta(2, c)$   
in state 2,  
read val 'c'  
 $\Rightarrow$  next state.

# DFA Lookup Table Construction

---

## Problems.

1. How do we represent the DFA in a computer friendly format?
2. How do construct the DFA in that format efficiently?

## Solutions.

1. Store a **lookup table**  $\delta[][]$ 
  - columns = states, rows = characters
  - $\delta[q][c] \leftarrow \delta(q, c)$
2. Compute column by column
  - trick: keep track of state for  $P[1, q]$  because we'll reuse this for each  $P[1, q]c$

# DFA Lookup Table Construction

## Solutions.

1. Store a **lookup table**  $\delta[][]$

- columns = states, rows = characters
- $\delta[q][c] \leftarrow \delta(q, c)$

2. Compute column by column

- trick: keep track of state for  $P[1, q]$  because we'll reuse this for each  $P[1, q]c$
- $x$  is largest value with

$$\boxed{P[0, x]} = \boxed{P[q-x, q]}$$

```
1: procedure CONSTRUCTDFA( $P[0..m]$ )
2:   for  $c \in \Sigma$  do
3:      $\delta[0][c] \leftarrow 0$ 
4:   end for
5:    $\delta[0][P[0]] \leftarrow 1$ 
6:    $x \leftarrow 0$ 
7:   for  $q = 1, 2, \dots, m-1$  do
8:     for  $c \in \Sigma$  do
9:        $\delta[q][c] \leftarrow \delta[x][c]$ 
10:    end for
11:     $\delta[q][P[q]] \leftarrow q+1$ 
12:     $x \leftarrow \delta[x][P[q]]$ 
13:  end for
14: end procedure
```

*init first col*

*iterate over columns*

# DFA Lookup Table Construction

**Example.**  $P[0..6] = ababaca$

$\delta(c, q)$	0	1	2	3	4	5	6
$P[q]$	a	b	a	b	a	c	a
a	1	1	3				
b	0	2	0	...			
c	0	0	0				
x	0	0	1				

```

1: procedure CONSTRUCTDFA( $P[0..m]$ )
2:   for  $c \in \Sigma$  do
3:      $\delta[0][c] \leftarrow 0$ 
4:   end for
5:    $\delta[0][P[0]] \leftarrow 1$ 
6:    $x \leftarrow 0$ 
7:   for  $q = 1, 2, \dots, m-1$  do
8:     for  $c \in \Sigma$  do
9:        $\delta[q][c] \leftarrow \delta[x][c]$ 
10:    end for
11:     $\delta[q][P[q]] \leftarrow q+1$ 
12:     $x \leftarrow \delta[x][P[q]]$ 
13:  end for
14: end procedure
  
```

Copy  
col x  
to col q

# DFA Lookup Table Construction

## PollEverywhere Question

What is the running time of CONSTRUCTDFA when  $P$  has length  $m$  and  $|\Sigma| = s$ ?



[pollev.com/comp526](http://pollev.com/comp526)

```
1: procedure CONSTRUCTDFA( $P[0..m]$ )
2:   for  $c \in \Sigma$  do
3:      $\delta[0][c] \leftarrow 0$ 
4:   end for
5:    $\delta[0][P[0]] \leftarrow 1$ 
6:    $x \leftarrow 0$ 
7:   for  $q = 1, 2, \dots, m-1$  do
8:     for  $c \in \Sigma$  do
9:        $\delta[q][c] \leftarrow \delta[x][c]$ 
10:    end for
11:     $\delta[q][P[q]] \leftarrow q+1$ 
12:     $x \leftarrow \delta[x][P[q]]$ 
13:  end for
14: end procedure
```

$\mathcal{O}(s)$

$\mathcal{O}(1)$

$s$

$\mathcal{O}(1)$

$m-2$  iter.  $\rightarrow \mathcal{O}(ms)$

# DFA Lookup Table Application

## Pitting it Together

- Construct the DFA
- Apply the DFA

start at  
state 0

```
1: procedure APPLYDFA( $T[0..n], \delta, m$ )
2:    $q \leftarrow 0$ 
3:   for  $i = 0, 1, \dots, n-1$  do
4:      $q \leftarrow \delta[q][T[i]]$ 
5:     if  $q = m$  then
6:       return  $i$ 
7:     end if
8:   end for
9:   return  $-1$ 
10: end procedure
11: procedure DFAMATCH( $P[0..m], T[0..n]$ )
12:    $\delta \leftarrow$  CONSTRUCTDFA( $P, T$ )
13:   return APPLYDFA( $T, \delta, m$ )
14: end procedure
```

$O(n)$   
time

# DFA Lookup Table Application

## Pitting it Together

- Construct the DFA
- Apply the DFA
- Running time is  $\Theta(\underline{n} + \underline{m}|\Sigma|)$ 
  - $\Theta(m|\Sigma|)$  for making DFA
  - $\Theta(n)$  for applying DFA
- Additional space overhead:  $\Theta(m|\Sigma|)$ 
  - store the DFA

could  
be a lot of space

```
1: procedure APPLYDFA( $T[0..n], \delta, m$ )
2:    $q \leftarrow 0$ 
3:   for  $i = 0, 1, \dots, n - 1$  do
4:      $q \leftarrow \delta[q][T[i]]$ 
5:     if  $q = m$  then
6:       return  $i$ 
7:     end if
8:   end for
9:   return  $-1$ 
10: end procedure
11: procedure DFAMATCH( $P[0..m], T[0..n]$ )
12:    $\delta \leftarrow$  CONSTRUCTDFA( $P, T$ )
13:   return APPLYDFA( $T, \delta, m$ )
14: end procedure
```



# Knuth-Morris- Pratt

# Failure Link Automaton

---

## DFA efficiency.

- Space/time to build DFA:  $\Theta(m|\Sigma|)$
- Time to execute DFA:  $\Theta(n)$

⇒ Overall time is  $\Theta(n + m|\Sigma|)$

- additional space overhead is  $\Theta(m|\Sigma|)$

**Question.** Can we perform string matching in time  $O(n)$  with *less space overhead*?

# Failure Link Automaton

---

## DFA efficiency.

- Space/time to build DFA:  $\Theta(m|\Sigma|)$
- Time to execute DFA:  $\Theta(n)$

⇒ Overall time is  $\Theta(n + m|\Sigma|)$

- additional space overhead is  $\Theta(m|\Sigma|)$

**Question.** Can we perform string matching in time  $O(n)$  with *less space overhead*?

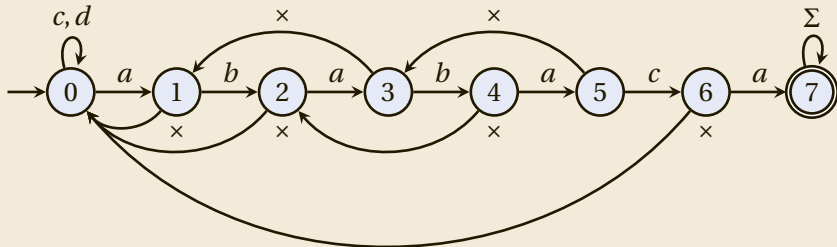
**Idea.** When comparison fails, don't have a separate transition for each failing character

- Just record failure and “shift” pattern as far forward as possible

# Failure Link Automaton

## Example

- T = aabacaababacaa
- P = ababaca



<b>text</b>	a	a	b	a	c	a	a	b	a	b	a	c	a	a
<b>states</b>														

# FLA Execution

---

A **Failure Link Automaton (FLA)**

consists of:

- A finite set  $Q$  of **states**
- A finite **alphabet**  $\Sigma$
- A **transition function**  
 $\varphi : Q \times (\Sigma \cup \{x\}) \rightarrow Q$
- An **initial state**  $q_0 \in Q$
- A set  $F \subseteq Q$  of **accepting states**

# FLA Execution

---

A **Failure Link Automaton (FLA)**

consists of:

- A finite set  $Q$  of **states**
- A finite **alphabet**  $\Sigma$
- A **transition function**  
 $\varphi : Q \times (\Sigma \cup \{\times\}) \rightarrow Q$
- An **initial state**  $q_0 \in Q$
- A set  $F \subseteq Q$  of **accepting states**

**Execution.** To apply and FLA to  $T$

- Start at the state  $q_0$
- Read characters from  $T$  sequentially
  - if in state  $q$  and read character  $c$ :
    - if  $\varphi(q, c)$  is defined, move to state  $\varphi(q, c)$
    - otherwise move to state  $\varphi(q, \times)$  and **re-read**  $c$
- Return TRUE if end in “accepting” state

# FLA Execution

## PollEverywhere Question

Given an FLA for a pattern  $P$  of length  $m$ , how many times could we follow failure links for a single character  $c$  read from  $T$  in the worst case?



[pollev.com/comp526](https://pollev.com/comp526)

**Execution.** To apply and FLA to  $T$

- Start at the state  $q_0$
- Read characters from  $T$  sequentially
  - if in state  $q$  and read character  $c$ :
    - if  $\varphi(q, c)$  is defined, move to state  $\varphi(q, c)$
    - otherwise move to state  $\varphi(q, \times)$  and **re-read**  $c$
- Return TRUE if end in “accepting” state

# FLA Execution

---

**Execution.** To apply and FLA to  $T$

- Start at the state  $q_0$
- Read characters from  $T$  sequentially
  - if in state  $q$  and read character  $c$ :
    - if  $\varphi(q, c)$  is defined, move to state  $\varphi(q, c)$
    - otherwise move to state  $\varphi(q, \times)$  and **re-read**  $c$
- Return TRUE if end in “accepting” state



# FLA Running Time

---

## More careful analysis

- If we match up to  $P[j]$ , then we can only follow up to  $j$  back links
- In order to witness  $j$  failures, must have witnessed  $j$  successes!

# FLA Running Time

---

## More careful analysis

- If we match up to  $P[j]$ , then we can only follow up to  $j$  back links
- In order to witness  $j$  failures, must have witnessed  $j$  successes!

**Amortized cost** of each character read from  $T$

- If read character  $c$  is a **match**:
  - pay 1 for comparison
  - put 1 unit cost in the **bank**
- If read character  $c$  is a **mismatch**
  - *withdraw* 1 from the bank
- By analysis above account balance is always non-negative

⇒ amortized cost of each comparison is 2

⇒ hence overall running time of execution is  $O(n)$

# FLA Construction

---

**Observation.** Each state  $q$  has

- 1 forward link to state  $q + 1$
- 1 fail link

Given  $P$ , we don't need to store forward link label:

- forward link label from  
 $P[q] = P[q + 1]$

Only need to store fail link state!

- this can be stored as a single array of size  $m$

⇒ only  $O(m)$  space overhead

# FLA Construction

---

**Definition.** The **failure link array** *fail* of  $P$  the array of  $m$  numbers that stores the (index of) the next state for each failure

- How do we construct it?

# FLA Construction

**Definition.** The **failure link array** *fail* of *P* the array of *m* numbers that stores the (index of) the next state for each failure

- How do we construct it?
- Again *x* is length of largest prefix that matches a suffix of  $P[1, q]$

**Example.**  $P[0..6] = \text{ababaca}$

<i>q</i>	0	1	2	3	4	5	6
$P[q]$	a	b	a	b	a	c	a
<i>fail</i> [ <i>q</i> ]							

```
1: procedure FAILURELINK( $P[0, m]$ )
2:   fail[0]  $\leftarrow$  0
3:   x  $\leftarrow$  0
4:   for  $j = 1, 2, \dots, m - 1$  do
5:     fail[j]  $\leftarrow$  x
6:     while  $P[x] \neq P[j]$  do
7:       if  $x = 0$  then
8:         x  $\leftarrow$  -1
9:         break
10:      else
11:        x  $\leftarrow$  fail[x]
12:      end if
13:    end while
14:    x  $\leftarrow$  x + 1
15:  end for
16: end procedure
```

# FLA Construction

---

**Question.** What is the running time of FAILURELINK on input of size  $m$ ?

```
1: procedure FAILURELINK( $P[0, m]$ )
2:    $fail[0] \leftarrow 0$ 
3:    $x \leftarrow 0$ 
4:   for  $j = 1, 2, \dots, m - 1$  do
5:      $fail[j] \leftarrow x$ 
6:     while  $P[x] \neq P[j]$  do
7:       if  $x = 0$  then
8:          $x \leftarrow -1$ 
9:         break
10:      else
11:         $x \leftarrow fail[x]$ 
12:      end if
13:    end while
14:     $x \leftarrow x + 1$ 
15:  end for
16: end procedure
```

# FLA Construction

---

**Question.** What is the running time of FAILURELINK on input of size  $m$ ?

## Observations.

- $x$  incremented once per  $j$
- $fail[x] < x$
- Each “while” iteration decrements  $x$

So at most  $2m$  updates to  $x$

- cf. amortized analysis
- $x =$  bank balance

```
1: procedure FAILURELINK( $P[0, m]$ )
2:    $fail[0] \leftarrow 0$ 
3:    $x \leftarrow 0$ 
4:   for  $j = 1, 2, \dots, m - 1$  do
5:      $fail[j] \leftarrow x$ 
6:     while  $P[x] \neq P[j]$  do
7:       if  $x = 0$  then
8:          $x \leftarrow -1$ 
9:         break
10:      else
11:         $x \leftarrow fail[x]$ 
12:      end if
13:    end while
14:     $x \leftarrow x + 1$ 
15:  end for
16: end procedure
```

# KMP Algorithm

---

**Question.** How do we apply the failure link array to find a match?



# KMP Algorithm

---

**Question.** How do we apply the failure link array to find a match?

- Scan along  $T[0, n)$ 
  - index  $i$
- Maintain position in  $P[0, m)$ 
  - index  $j$
  - current prefix match
- When  $T[i] = P[j]$ , increment  $i$  and  $j$
- Otherwise,  $j \leftarrow fail[j]$ 
  - unless  $j = 0$ , then  $i \leftarrow i + 1$

# KMP Algorithm

---

**Question.** How do we apply the failure link array to find a match?

- Scan along  $T[0, n)$ 
  - index  $i$
- Maintain position in  $P[0, m)$ 
  - index  $j$
  - current prefix match
- When  $T[i] = P[j]$ , increment  $i$  and  $j$
- Otherwise,  $j \leftarrow fail[j]$ 
  - unless  $j = 0$ , then  $i \leftarrow i + 1$

```
1: procedure KMP( $T[0..n), P[0..m)$ )
2:    $fail \leftarrow FAILURELINK(P)$ 
3:    $i \leftarrow 0$ 
4:    $j \leftarrow 0$ 
5:   while  $i < n$  do
6:     if  $T[i] = P[j]$  then
7:        $i \leftarrow i + 1, j \leftarrow j + 1$ 
8:       if  $j = m$  then return  $i - j$ 
9:     else
10:      if  $j \geq 1$  then
11:         $j \leftarrow fail[j]$ 
12:      else
13:         $i \leftarrow i + 1$ 
14:      end if
15:    end if
16:  end while
17: end procedure
```

# KMP Algorithm

---

## Analysis:

- Running time  $O(n + m)$ 
  - $O(m)$  to build *fail*
  - $O(n)$  to apply KMP
  - analysis uses **amortized analysis**
- Additional space  $O(m)$ 
  - just need to store *fail* and indices

```
1: procedure KMP( $T[0..n], P[0..m]$ )
2:    $fail \leftarrow \text{FAILURELINK}(P)$ 
3:    $i \leftarrow 0$ 
4:    $j \leftarrow 0$ 
5:   while  $i < n$  do
6:     if  $T[i] = P[j]$  then
7:        $i \leftarrow i + 1, j \leftarrow j + 1$ 
8:       if  $j = m$  then return  $i - j$ 
9:     else
10:      if  $j \geq 1$  then
11:         $j \leftarrow fail[j]$ 
12:      else
13:         $i \leftarrow i + 1$ 
14:      end if
15:    end if
16:  end while
17: end procedure
```

# KMP Algorithm

---

## Analysis:

- Running time  $O(n + m)$ 
  - $O(m)$  to build *fail*
  - $O(n)$  to apply KMP
  - analysis uses **amortized analysis**
- Additional space  $O(m)$ 
  - just need to store *fail* and indices

## Clean Takeaway:

*fail*[*j*] is the length of the longest prefix of  $P[0..j]$  that is a suffix of  $P[1..j]$

```
1: procedure KMP( $T[0..n], P[0..m]$ )
2:   fail ← FAILURELINK( $P$ )
3:    $i$  ← 0
4:    $j$  ← 0
5:   while  $i < n$  do
6:     if  $T[i] = P[j]$  then
7:        $i$  ←  $i + 1, j$  ←  $j + 1$ 
8:       if  $j = m$  then return  $i - j$ 
9:     else
10:      if  $j \geq 1$  then
11:         $j$  ← fail[ $j$ ]
12:      else
13:         $i$  ←  $i + 1$ 
14:      end if
15:    end if
16:  end while
17: end procedure
```

# DFA vs FLA

---

**Question.** Which is better? DFA matching or KMP algorithm?

- KMP has overall running time  $O(n + m)$ 
  - amortized 2 comparisons per  $T$  access
- DFA has overall running time  $O(n + m|\Sigma|)$ 
  - 1 comparison per  $T$  access
  - $|\Sigma|$  dependence

# Next Time

---

More String Matching!

# Scratch Notes

---