



Lecture 12: String Matching III

COMP526: Efficient Algorithms

Updated: November 12, 2024

Will Rosenbaum
University of Liverpool

Announcements

1. Programming Assignment 1 **DUE WEDNESDAY**
 - Use updated testing code (from last Wednesday)
 - Submission through Canvas
 - *Only* submit `pr_tester.py`
 - Late Policy: 5% off per day down to 50%
2. Quiz due Friday
 - Covers string matching
 - *including today's lecture*
 - 2 questions (multiple choice)
3. Attendance Code:

Meeting Goals

Discuss String Matching procedures:

- Knuth-Morris-Pratt
- Boyer-Moore

The String Matching Problem

Input:

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

Output:

- The index of the **first occurrence** of P in T

The String Matching Problem

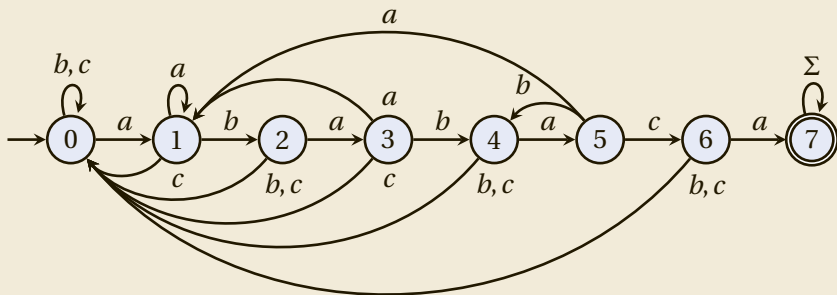
Input:

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

Output:

- The index of the **first occurrence** of P in T

Last Time. Search with DFA



Example: $T = abababac$

The String Matching Problem

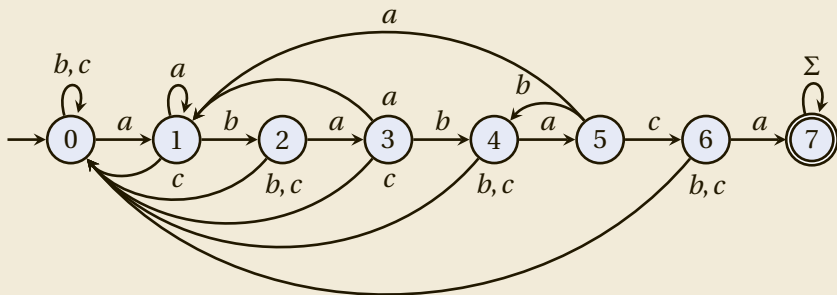
Input:

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

Output:

- The index of the **first occurrence** of P in T

Last Time. Search with DFA



Result: Search in time $\Theta(n + |\Sigma|n)$ with space overhead $|\Sigma|n$.

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*?

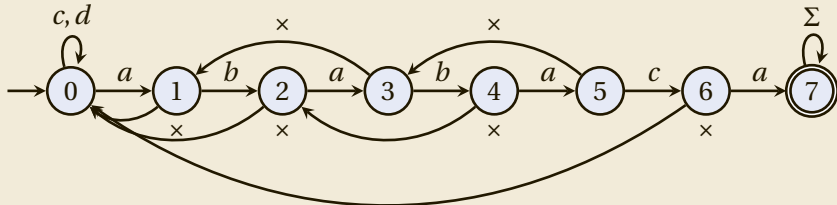
Idea. When comparison fails, don't have a separate transition for each failing character

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

Failure Link Automaton

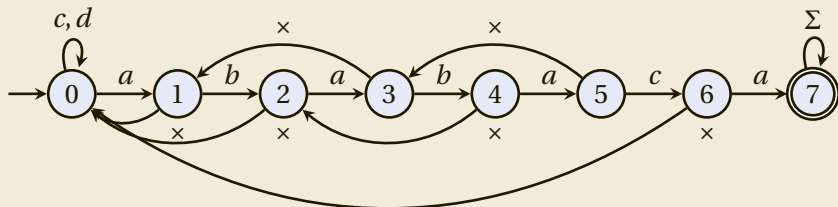
Example

- T = aababaababacaa
- P = ababaca



text	a	a	b	a	b	a	a	b	a	b	a	c	a	a
states														

States and Shifts



<i>a</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</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>
<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</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>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</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>	<i>b</i>	<i>a</i>	<i>b</i>	<i>a</i>	<i>c</i>	<i>a</i>	

Correspondence: matches increment T index i , mismatches shift P

- shift amount aligns largest possible number of matches

FLA Execution

A **Failure Link Automaton (FLA)**

consists of:

- A finite set Q of **states**
- A finite **alphabet** Σ
- 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** Σ
- 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

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 q to $q+1$ is $P[q]$

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

q	0	1	2	3	4	5	6
$fail[q]$							

```
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.

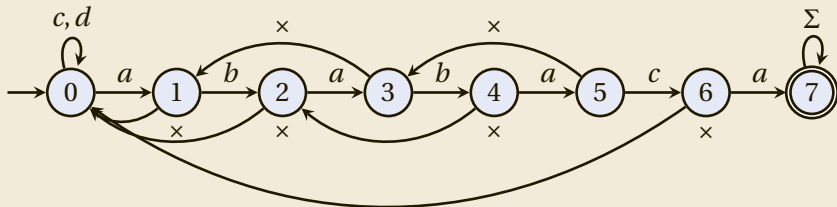
- x incremented once per j
- $fail[x] < x$
- Each “while” iteration decrements x

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
```

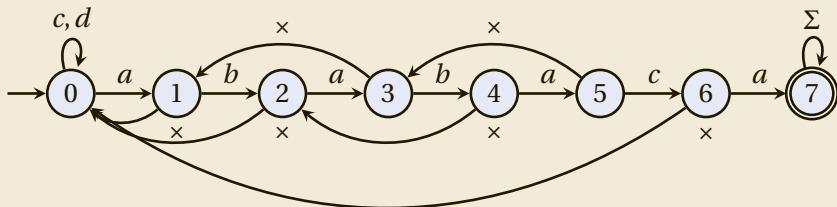
Failure Links: 3 Views



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

q	0	1	2	3	4	5	6
$fail[q]$	0	0	0	1	2	3	0

Failure Links: 3 Views



a b a b a c a
 a b a b a c a
 a b a b a c a

q	0	1	2	3	4	5	6
$fail[q]$	0	0	0	1	2	3	0

$fail[q]$ is

- the max of alignments formed by shifting P if first mismatch at $P[q]$
- longest prefix of $P[0, q]$ that is a suffix of $P[1, q]$

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:    $fail \leftarrow FAILURELINK(P)$ 
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 Example

Example. Find the failure link array for $P[0,8) = \text{BCBABCBA}$.

i	0	1	2	3	4	5	6	7
$fail[i]$								

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

KMP Example

Example. Find the failure link array for $P[0, 8) = \text{BCBABCBA}$.

i	0	1	2	3	4	5	6	7
$fail[i]$	0	0	0	1	0	1	2	3

KMP Example

Example. Find the failure link array for $P[0,8) = \text{BCBABCBA}$.

i	0	1	2	3	4	5	6	7
$fail[i]$	0	0	0	1	0	1	2	3

Interpretation. If $T[i..i+j)$ matches $P[0..j)$, but $T[i+j) \neq P[j)$, then $fail[j]$ is the maximum number matches between $T[i+1, i+j)$ and P .

1	2	3	4	5	6	7						
B	C	B	A	B	C	B	A					
		B	C	B	A	B	C	B	A			
				B	C	B	A	B	C	B	A	

KMP Example

Example. Find the failure link array for $P[0,8) = \text{BCBABCBA}$.

i	0	1	2	3	4	5	6	7
$fail[i]$	0	0	0	1	0	1	2	3

Interpretation. If $T[i..i+j)$ matches $P[0..j)$, but $T[i+j) \neq P[j)$, then $fail[j]$ is the maximum number matches between $T[i+1, i+j)$ and P .

1	2	3	4	5	6	7						
B	C	B	A	B	C	B	A					
		B	C	B	A	B	C	B	A			
				B	C	B	A	B	C	B	A	

Visualization. See website.

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

Boyer-Moore

Beyond Worst-Case Pattern Matching?

A Puzzle. Suppose we have

- $P[0, 4) = \text{AAAA}$
- $T[0, 14) = \text{BBBBBBBBBBBBBB}$

If we know P , what is the fewest number of accesses we can make to T to **certify** that T does not contain P ?

Beyond Worst-Case Pattern Matching?

A Puzzle. Suppose we have

- $P[0, 4) = \text{AAAA}$
- $T[0, 14) = \text{BBBBBBBBBBBBBB}$

If we know P , what is the fewest number of accesses we can make to T to **certify** that T does not contain P ?

?	?	?	B	?	?	?	B	?	?	?	B	?	?
A	A	A	A										
				A	A	A	A						
								A	A	A	A		

Beyond Worst-Case Pattern Matching?

A Puzzle. Suppose we have

- $P[0, 4) = \text{AAAA}$
- $T[0, 14) = \text{BBBBBBBBBBBBBB}$

If we know P , what is the fewest number of accesses we can make to T to **certify** that T does not contain P ?

?	?	?	B	?	?	?	B	?	?	?	B	?	?
A	A	A	A										
				A	A	A	A						
								A	A	A	A		

Observation.

- By starting comparisons from the *end* of P , we could eliminate more possible alignments.

Two Heuristics

Strategy. To test match of $P[0..m)$ with $T[j..j+m)$, perform comparisons from *right to left*

Two Heuristics

Strategy. To test match of $P[0..m)$ with $T[j..j+m)$, perform comparisons from *right to left*

Heuristic 1. If we encounter $T[i]$ that does not occur in P , shift P entirely past index i .

T:	...	A	B	D	C	A	A	C	A	B	C	A	...
P:		C	A	B	C	A							
				→	C	A	B	C	A				

Two Heuristics

Strategy. To test match of $P[0..m)$ with $T[j..j+m)$, perform comparisons from *right to left*

Heuristic 1. If we encounter $T[i]$ that does not occur in P , shift P entirely past index i .

T:	...	A	B	D	C	A	A	C	A	B	C	A	...
P:		C	A	B	C	A							
				→	C	A	B	C	A				

Heuristic 2. If we match on a suffix of P but mismatch at index i , shift P to next alignment of suffix.

T:	...	A	B	D	C	A	A	C	A	B	C	A	...
P:					C	A	B	C	A				
							→	C	A	B	C	A	

Boyer-Moore Algorithm

Combining these heuristics gives the **Boyer-Moore algorithm**

- Compare alignments from right to left
- If we encounter $T[i]$ that does not occur in P , shift P entirely past index i .
- If we match on a suffix of P but mismatch at index i , shift P to next alignment of suffix

T:	...	A	B	D	C	A	A	C	A	B	C	A	...
P:		C	A	B	C	A							
				→	C	A	B	C	A				

T:	...	A	B	D	C	A	A	C	A	B	C	A	...
P:					C	A	B	C	A				
							→	C	A	B	C	A	

Boyer-Moore Algorithm

Combining these heuristics gives the **Boyer-Moore algorithm**

- Compare alignments from right to left
- If we encounter $T[i]$ that does not occur in P , shift P entirely past index i .
- If we match on a suffix of P but mismatch at index i , shift P to next alignment of suffix

Features of this approach:

- Worst-case running time on $P[0..m)$ and $T[0..n)$ is $\Theta(nm)$
 - achieved if all instances of P must be reported
 - can be improved to $\Theta(n + m + |\Sigma|)$ with some care if T does not contain P

Boyer-Moore Algorithm

Combining these heuristics gives the **Boyer-Moore algorithm**

- Compare alignments from right to left
- If we encounter $T[i]$ that does not occur in P , shift P entirely past index i .
- If we match on a suffix of P but mismatch at index i , shift P to next alignment of suffix

Features of this approach:

- Worst-case running time on $P[0..m)$ and $T[0..n)$ is $\Theta(nm)$
 - achieved if all instances of P must be reported
 - can be improved to $\Theta(n + m + |\Sigma|)$ with some care if T does not contain P
- Typical running time can be much better!
 - For some *random* string models, expected running time is $O(n/m)$
 - For English text, typically uses $\sim 0.25n$ comparisons if no match

Boyer-Moore Algorithm

Combining these heuristics gives the **Boyer-Moore algorithm**

- Compare alignments from right to left
- If we encounter $T[i]$ that does not occur in P , shift P entirely past index i .
- If we match on a suffix of P but mismatch at index i , shift P to next alignment of suffix

Features of this approach:

- Worst-case running time on $P[0..m)$ and $T[0..n)$ is $\Theta(nm)$
 - achieved if all instances of P must be reported
 - can be improved to $\Theta(n + m + |\Sigma|)$ with some care if T does not contain P
- Typical running time can be much better!
 - For some *random* string models, expected running time is $O(n/m)$
 - For English text, typically uses $\sim 0.25n$ comparisons if no match
- Space overhead is $\Theta(m + |\Sigma|)$

Summary of String Matching

- **Brute Force:**
 - simplest description
 - $\Theta(nm)$ running time
 - $O(1)$ space overhead

Summary of String Matching

- **Brute Force:**
 - simplest description
 - $\Theta(nm)$ running time
 - $O(1)$ space overhead
- **DFA**
 - few comparisons (worst case)
 - $\Theta(n + m|\Sigma|)$ running time
 - $\Theta(m|\Sigma|)$ space overhead (DFA table)

Summary of String Matching

- **Brute Force:**
 - simplest description
 - $\Theta(nm)$ running time
 - $O(1)$ space overhead
- **DFA**
 - few comparisons (worst case)
 - $\Theta(n + m|\Sigma|)$ running time
 - $\Theta(m|\Sigma|)$ space overhead (DFA table)
- **Knuth-Morris-Pratt**
 - simple description
 - $\Theta(n + m)$ running time (inc. all occurrences)
 - $\Theta(m)$ space overhead (fail array)

Summary of String Matching

- **Brute Force:**

- simplest description
- $\Theta(nm)$ running time
- $O(1)$ space overhead

- **DFA**

- few comparisons (worst case)
- $\Theta(n + m|\Sigma|)$ running time
- $\Theta(m|\Sigma|)$ space overhead (DFA table)

- **Knuth-Morris-Pratt**

- simple description
- $\Theta(n + m)$ running time (inc. all occurrences)
- $\Theta(m)$ space overhead (fail array)

- **Boyer-Moore**

- efficient in practice (English text)
- $\Theta(nm)$ worst case to find all occurrences, can be as small as $O(n/m)$
- $\Theta(m)$ overhead

Summary of String Matching

- **Brute Force:**
 - simplest description
 - $\Theta(nm)$ running time
 - $O(1)$ space overhead
- **DFA**
 - few comparisons (worst case)
 - $\Theta(n + m|\Sigma|)$ running time
 - $\Theta(m|\Sigma|)$ space overhead (DFA table)
- **Knuth-Morris-Pratt**
 - simple description
 - $\Theta(n + m)$ running time (inc. all occurrences)
 - $\Theta(m)$ space overhead (fail array)
- **Boyer-Moore**
 - efficient in practice (English text)
 - $\Theta(nm)$ worst case to find all occurrences, can be as small as $O(n/m)$
 - $\Theta(m)$ overhead
- **Rabin-Karp**
 - based on **hashing**
 - generalizes beyond one-dimensional strings
 - expected running time $O(n + m)$
 - $O(1)$ space overhead

Next Time

Data Compression!

- How much **space** do we need to store our data?

Scratch Notes
