

# CS 240 – Data Structures and Data Management

## Module 9: String Matching

O. Veksler

Based on lecture notes by many previous cs240 instructors

**David R. Cheriton School of Computer Science, University of Waterloo**

Spring 2023

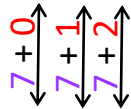
# Outline

- String Matching
  - Introduction
  - Karp-Rabin Algorithm
  - Boyer-Moore Algorithm
  - Suffix Trees

# Pattern Matching Definitions [1]

- Search for a string (pattern) in a large body of text
- $T[0 \dots n - 1]$  **text** (or **haystack**) being searched
- $P[0 \dots m - 1]$  **pattern** (or **needle**) being searched for
- Strings over **alphabet**  $\Sigma$
- Return the first occurrence of  $P$  in  $T$
- Example

$T =$     Little piglets cooked for mother pig



$P =$     pig

$n = 36, m = 3, i = 7$

- return smallest  $i$  such that

$$T[i + j] = P[j] \text{ for } 0 \leq j \leq m - 1$$

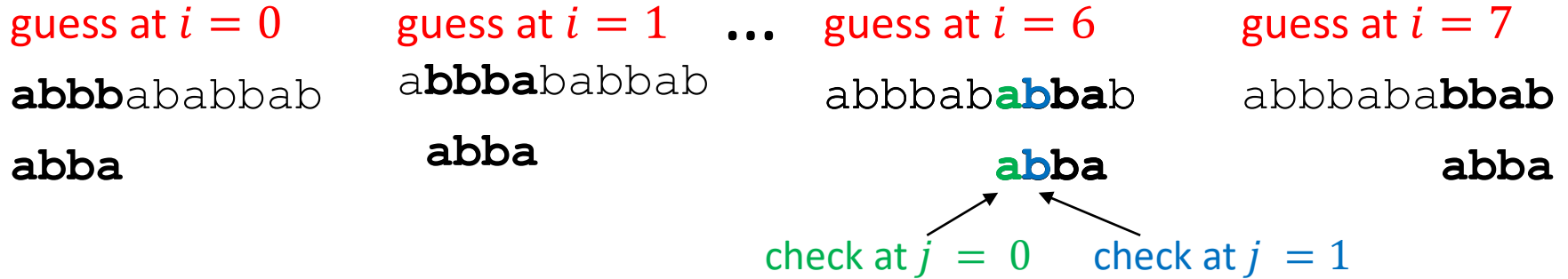
- If  $P$  does not occur in  $T$ , return FAIL
- Applications
  - information retrieval (text editors, search engines), bioinformatics, data mining

# More Definitions [2]

## antid~~estab~~lishmentarianism

- **Substring**  $T[i..j]$   $0 \leq i \leq j < n$  is a string consisting of characters  $T[i], T[i + 1], \dots, T[j]$ 
  - length is  $j - i + 1$
- **Prefix** of  $T$  is a substring  $T[0..i]$  of  $T$  for some  $0 \leq i \leq n - 1$
- **Suffix** of  $T$  is a substring  $T[i..n - 1]$  of  $T$  for some  $0 \leq i \leq n - 1$
- With this definition, prefix and suffix are never empty strings
  - sometimes want to allow empty string prefix and suffix

# General Idea of Algorithms

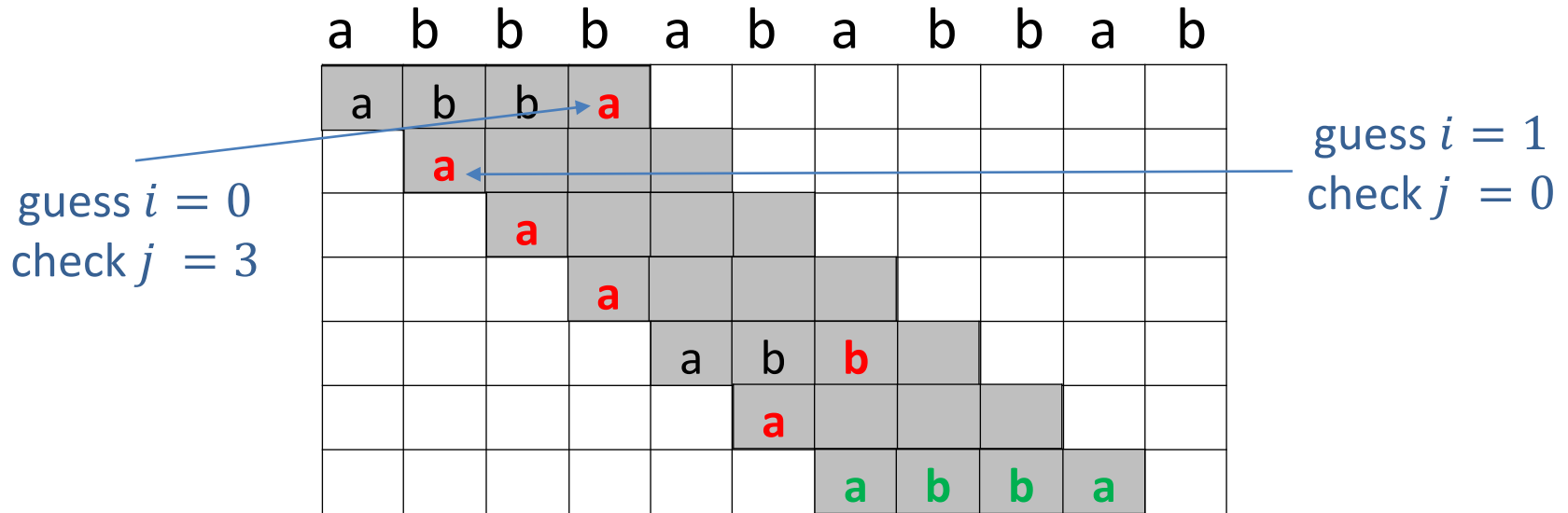


- Pattern matching algorithms consist of **guesses** and **checks**
  - a **guess** or **shift** is a position  $i$  such that  $P$  might start at  $T[i]$
  - valid guesses (initially) are  $0 \leq i \leq n - m$
  - a **check** of a guess is a single position  $j$  with  $0 \leq j < m$  where we compare  $T[i + j]$  to  $P[j]$ 
    - must perform  $m$  checks of a single **correct** guess
  - may make fewer checks of an **incorrect** guess



# Brute-Force Algorithm: Example

Example:  $T = \text{abbababbab}$ ,  $P = \text{abba}$



- Worst possible input
  - $P = \underbrace{a \dots ab}_{m-1 \text{ times}}, T = \underbrace{aaaaaaaaa \dots aaaaaaaaa}_{n \text{ times}}$
- Have to perform  $(n - m + 1)m$  checks, which is  $\Theta((n - m)m)$  runtime
  - this is  $\Theta(nm)$  if  $m \leq n/2$
  - worst running time if  $m = n/2$ 
    - $\Theta(n^2)$

# Brute-force Algorithm

- Checks every possible guess

```
Bruteforce::PatternMatching( $T [0..n - 1]$ ,  $P[0..m - 1]$ )  
 $T$  : String of length  $n$  (text),  $P$ : String of length  $m$  (pattern)  
  for  $i \leftarrow 0$  to  $n - m$  do  
    if strcmp( $T [i \dots i + m - 1]$ ,  $P$ ) = 0  
      return "found at guess  $i$ "  
  return FAIL
```

- Note: *strcmp* takes  $\Theta(m)$  time

```
strcmp( $T [i \dots i + m - 1]$ ,  $P[0..m - 1]$ )  
for  $j \leftarrow 0$  to  $m - 1$  do  
  if  $T [i + j]$  is before  $P[j]$  in  $\Sigma$  then return -1  
  if  $T [i + j]$  is after  $P[j]$  in  $\Sigma$  then return 1  
return 0
```



# How to improve?

- Extra **preprocessing** on pattern  $P$ 
  - **Karp-Rabin**
  - **KMP**
  - **Boyer-Moore**
  - **Eliminate guesses** based on completed matches and mismatches
- Do extra **preprocessing** on the text  $T$ 
  - **Suffix-trees**
  - **Suffix-arrays**
  - **Create a data structure** to find matches easily

# Outline

- String Matching
  - Introduction
  - Karp-Rabin Algorithm
  - Boyer-Moore Algorithm
  - Suffix Trees

# Karp-Rabin Fingerprint Algorithm: Idea

- Hash functions are useful not just for hash tables!
- **Idea:** use hashing to eliminate guesses faster
  - compute hash function for each guess, compare with pattern hash
    - if values are unequal, then current guess cannot match the pattern
    - if values are equal, **verify** that pattern actually matches text
      - equal hash value does not guarantee equal keys
      - although if hash function is good, most likely keys are equal
      - $O(m)$  time to verify, but happens rarely, and most likely only for true match
  - Example:  $P = 5\ 9\ 2\ 6\ 5$ ,  $T = 3\ 1\ 4\ 1\ 5\ 9\ 2\ 6\ 5\ 3\ 5$ 
    - standard hash function: flattening + modular (radix  $R = 10$ ):

$$h(59265) = (5 \cdot 10^4 + 9 \cdot 10^3 + 2 \cdot 10^2 + 6 \cdot 10^1 + 5) \bmod 97 = 59265 \bmod 97 = 95$$

3	1	4	1	5	9	2	6	5	3	5	
hash-value 84											$h(31415) = 84$
	hash-value 94										$h(14159) = 94$
		hash-value 76									$h(41592) = 76$
			hash-value 18								$h(15926) = 18$
				hash-value 95							$h(59265) = 95$

# Karp-Rabin Fingerprint Algorithm – First Attempt

```
Karp-Rabin-Simple::patternMatching( $T, P$ )
```

```
 $h_P \leftarrow h(P[0..m-1])$ 
```

```
for  $i \leftarrow 0$  to  $n - m$ 
```

```
     $h_T \leftarrow h(T[i..i+m-1])$ 
```

```
    if  $h_T = h_P$ 
```

```
        if strcmp( $T[i..i+m-1], P) = 0$ 
```

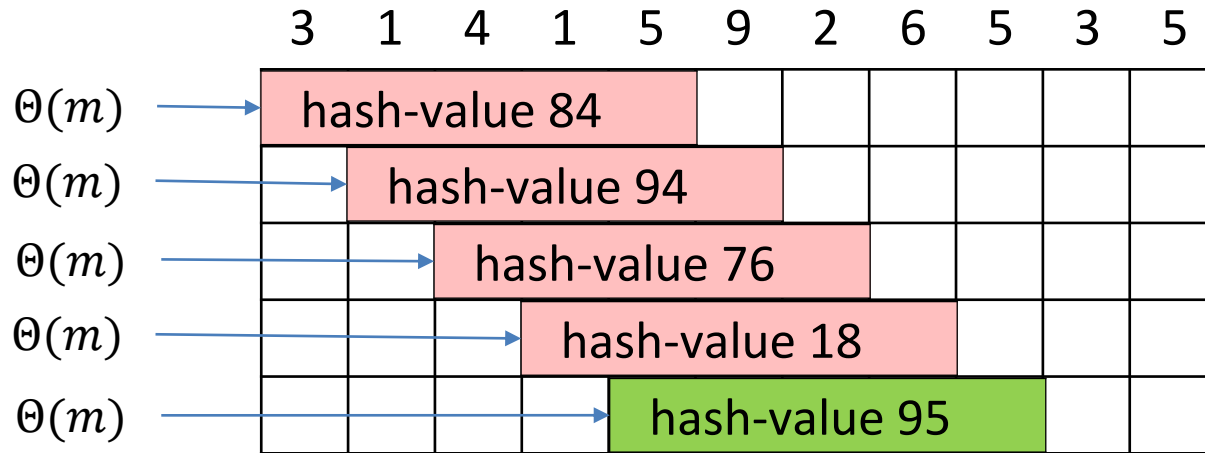
```
            return “found at guess  $i$ ”
```

```
    return FAIL
```

$\Theta(m)$

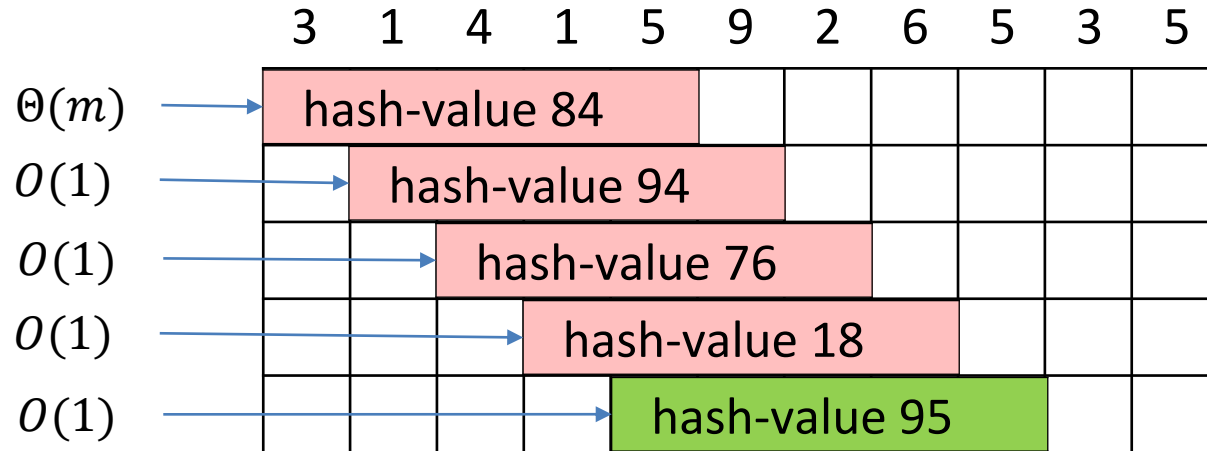
- Algorithm correctness: match is not missed
  - $h(T[i..i+m-1]) \neq h(P) \Rightarrow$  guess  $i$  is not  $P$
- What about running time?

# Karp-Rabin Fingerprint Algorithm: First Attempt



- For each shift,  $\Theta(m)$  time to compute hash value
  - since  $h(T[i \dots i + m - 1])$  depends on all  $m$  characters
  - worse than brute-force!
    - it is possible for brute force matching to use less than  $\Theta(m)$  per shift, as it stops at the first mismatched character
- $n - m + 1$  shifts in text to check
- Total time is  $\Theta(mn)$  if pattern not in text
  - how can we improve this?

# Karp-Rabin Fingerprint Algorithm: Idea



- Idea: compute next hash from previous one in  $O(1)$  time
- $n - m + 1$  shifts in text to check
- $\Theta(m)$  to compute the first hash value
- $O(1)$  to compute all other hash values
- $\Theta(n + m)$  expected time
  - recall that we still need to check if the pattern actually matches text whenever hash value of text is equal to the hash value of pattern
  - if hash function is good, then whenever hash values are equal, pattern most likely matches the text

# Karp-Rabin Fingerprint Algorithm – Fast Rehash

- For historical reasons, hashes are called **fingerprints**
- Insight: can update a fingerprint from previous fingerprint in constant time
  - $O(1)$  time to compute any hash, except first one

- **Example**

$$T = 4\ 1\ 5\ 9\ 2\ 6\ 5\ 3\ 5, \quad P = 5\ 9\ 2\ 6\ 5$$

- Initialization of the algorithm

1. compute first hash:  $h(41592) = 41592 \bmod 97 = 76$  [ $\Theta(m)$  time]
2. also compute  $10000 \bmod 97 = 9$

- Main loop: repeatedly compute next hash from the previous one

- Example : 15926  $\bmod 97$  from 41592  $\bmod 97$

- get rid of the old **first digit** and add new **last digit**

$$41592 \xrightarrow{-4 \cdot 10000} 1592 \xrightarrow{\times 10} 15920 \xrightarrow{+6} 15926$$

- Algebraically,

$$(41592 - (4 \cdot 10000)) \cdot 10 + 6 = 15926$$

# Karp-Rabin Fingerprint Algorithm – Fast Rehash

- Insight: can update a fingerprint from previous fingerprint in constant time

- Example**

$$T = 415926535, \quad P = 59265$$

- Initialization of the algorithm

- compute first hash:  $h(41592) = 41592 \bmod 97 = 76$  [ $\Theta(m)$  time]
- also compute  $10000 \bmod 97 = 9$

- Main loop: repeatedly compute next hash from the previous one

- Example:  $15926 \bmod 97$  from  $41592 \bmod 97$

$$(41592 - (4 \cdot 10000)) \cdot 10 + 6 = 15926$$

$$((41592 - (4 \cdot 10000)) \cdot 10 + 6) \bmod 97 = 15926 \bmod 97$$

$$((41592 \bmod 97 - (4 \cdot (10000 \bmod 97))) \cdot 10 + 6) \bmod 97 = 15926 \bmod 97$$

previous hash

precomputed

$$\left( (76 - (4 \cdot 9)) \cdot 10 + 6 \right) \bmod 97 = 15926 \bmod 97$$

constant number of operations, independent of  $m$

$$18 = 15926 \bmod 97$$



# Karp-Rabin Fingerprint Algorithm – Conclusion

*Karp-Rabin-RollingHash::PatternMatching*( $T, P$ )

$M \leftarrow$  suitable prime number

$h_P \leftarrow h(P[0..m-1])$

$h_T \leftarrow h(T[0..m-1])$

$s \leftarrow 10^{m-1} \bmod M$

**for**  $i \leftarrow 0$  to  $n - m$

**if**  $h_T = h_P$

**if** *strcmp*( $T[i..i+m-1], P) = 0$

**return** “found at guess  $i$ ”

**if**  $i < n - m$  // compute hash-value for next guess

$h_T \leftarrow ((h_T - T[i] \cdot s) \cdot 10 + T[i+m]) \bmod M$

**return** FAIL

- Choose “table size”  $M$  at **random** to be prime in  $\{2, \dots, mn^2\}$
- Expected running time is  $O(m + n)$
- $\Theta(mn)$  worst-case, but this extremely is unlikely
- Improvement: reset  $M$  if no match at  $h_T = h_P$

# Outline

- String Matching
  - Introduction
  - Karp-Rabin Algorithm
  - Boyer-Moore Algorithm
  - Suffix Trees

# Boyer-Moore Algorithm Motivation

- Fastest pattern matching in practice on English Text
- Important components
  - **Reverse-order searching**
    - compare  $P$  with a guess moving *backwards*
  - When a mismatch occurs choose the better option among the two below
    1. **Bad character heuristic**
      - eliminate shifts based on mismatched character of  $T$
    2. **Good suffix heuristic**
      - eliminate shifts based on the matched part (i.e.) suffix of  $P$

# Reverse Searching vs. Forward Searching

$T = \text{whereiswaldo}$ ,  $P = \text{aldo}$

w	h	e	r	e	i	s	w	a	l	d	o
			o								
							o				
								a	l	d	o

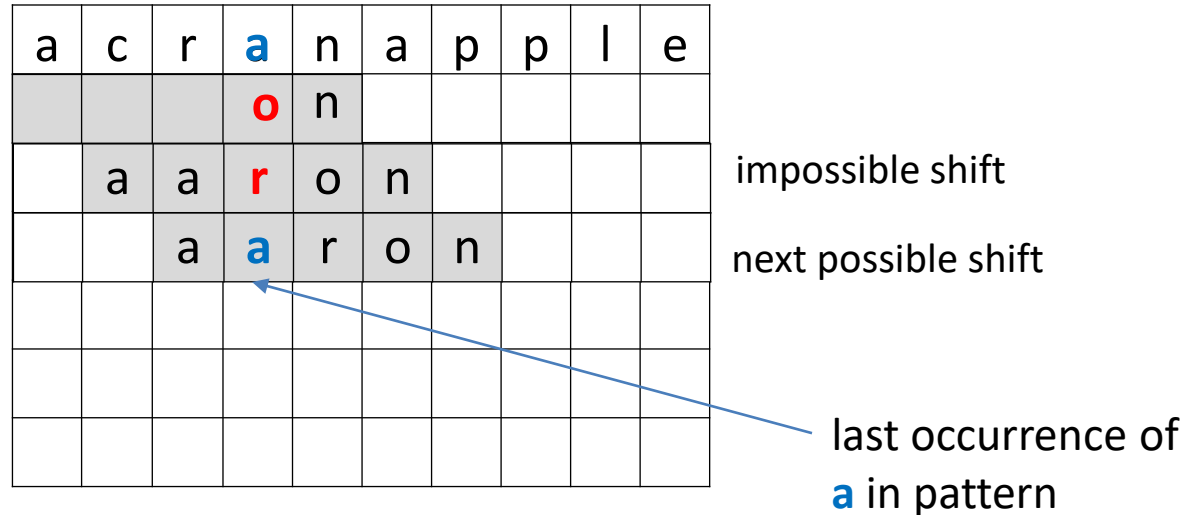
w	h	e	r	e	i	s	w	a	l	d	o
a											

- **r** does not occur in  $P = \text{aldo}$
- shift pattern past **r**
- **w** does not occur in  $P = \text{aldo}$
- shift pattern past **w**
- **bad character heuristic** can rule out many shifts with reverse searching

- **w** does not occur in  $P = \text{aldo}$
- move pattern past **w**
- the first shift moves pattern past **w**
- no shifts are ruled out
- **bad character heuristic** does not rule out any shifts with forward searching when the first character of the pattern is mismatched

# What if Mismatched Text Character Occurs in $P$ ?

$T = \text{acranapple}$ ,  $P = \text{aaron}$



- Mismatched character in the text is **a**
- Find **last** occurrence of **a** in  $P$
- Shift the pattern to the right until **last a** in  $P$  aligns with **a** in text
  - all smaller shifts are impossible since they do not match **a**
- Precompute last occurrence of any letter before matching starts

# Bad Character Heuristic: Side Note

$T = \text{acranapple}$ ,  $P = \text{aaron}$

a	c	r	a	n	a	p	p	l	e
			o	n					
		a	a	r	o	n			
			a	a	r	o	n		

next possible shift

also a valid shift

- If we shifted until the **first** a in P aligns with a in text
  - this would give a possible shift, but misses a previous possible shift, possibly leading to a missed pattern

# Bad Character Heuristic: Full Version

- Extends to the case when mismatched text character does occur in  $P$

$T = \text{acranapple}$ ,  $P = \text{aaron}$

a	c	r	a	n	a	p	p	l	e
			o	n					
			[a]						

- Mismatched character in the text is **a**
- Shift the pattern to the right so that the last **a** in  $P$  aligns with **a** in text
- Continue matching the pattern (in reverse)

# Bad Character Heuristic: Full Version

- Extends to the case when mismatched text character does occur in  $P$

$T = \text{acranapple}$ ,  $P = \text{aaron}$

a	c	r	a	n	a	p	p	l	e
			o	n					
			[a]			n			

- Mismatched character in the text is **a**
- Shift the pattern to the right so that the last **a** in  $P$  aligns with **a** in text
- Continue matching the pattern (in reverse)



# Bad Character Heuristic: Last Occurrence Array

- Compute the **last occurrence array**  $L(c)$  of any character in the alphabet
  - $L(c) = -1$  if character  $c$  does not occur in  $P$ , otherwise
  - $L(c) =$  largest index  $j$  such that  $P[j] = c$
- Example:  $P = \text{aaron}$

- initialization

<i>char</i>	a	n	o	r	all others
$L(c)$	-1	-1	-1	-1	-1

this means:

a	b	c	d	e	f	...	x	y	z
-1	-1	-1	-1	-1	-1		-1	-1	-1

in actual implementation:

0	1	2	3	4	5	...	24	25	26
-1	-1	-1	-1	-1	-1		-1	-1	-1

# Bad Character Heuristic: Last Occurrence Array

- Compute the **last occurrence array**  $L(c)$  of any character in the alphabet
  - $L(c) = -1$  if character  $c$  does not occur in  $P$ , otherwise
  - $L(c) =$  largest index  $j$  such that  $P[j] = c$
- Example:  $P =$  aaron

- computation

**a**aaron

$i = 0$

<i>char</i>	a	n	o	r	all others
$L(c)$	0	-1	-1	-1	-1

$L$  is valid for  $P =$  **a**

# Bad Character Heuristic: Last Occurrence Array

- Compute the **last occurrence array**  $L(c)$  of any character in the alphabet
  - $L(c) = -1$  if character  $c$  does not occur in  $P$ , otherwise
  - $L(c) =$  largest index  $j$  such that  $P[j] = c$
- Example:  $P = \text{aaron}$

- computation

**a**aron

$i = 1$

<i>char</i>	a	n	o	r	all others
$L(c)$	1	-1	-1	-1	-1

$L$  is valid for  $P = \text{aa}$

# Bad Character Heuristic: Last Occurrence Array

- Compute the **last occurrence array**  $L(c)$  of any character in the alphabet
  - $L(c) = -1$  if character  $c$  does not occur in  $P$ , otherwise
  - $L(c) =$  largest index  $j$  such that  $P[j] = c$
- Example:  $P =$  aaron

- computation

aaron

$i = 2$

<i>char</i>	a	n	o	r	all others
$L(c)$	1	-1	-1	2	-1

$L$  is valid for  $P =$  aar

# Bad Character Heuristic: Last Occurrence Array

- Compute the **last occurrence array**  $L(c)$  of any character in the alphabet
  - $L(c) = -1$  if character  $c$  does not occur in  $P$ , otherwise
  - $L(c) =$  largest index  $j$  such that  $P[j] = c$
- Example:  $P =$  aaron

- computation

aaron

$i = 3$

<i>char</i>	a	n	o	r	all others
$L(c)$	1	-1	3	2	-1

$L$  is valid for  $P =$  aar**o**

# Bad Character Heuristic: Last Occurrence Array

- Compute the **last occurrence array**  $L(c)$  of any character in the alphabet
  - $L(c) = -1$  if character  $c$  does not occur in  $P$ , otherwise
  - $L(c) =$  largest index  $j$  such that  $P[j] = c$
- Example:  $P =$  aaron

- computation

aaron

$i = 4$

<i>char</i>	a	n	o	r	all others
$L(c)$	1	4	3	2	-1

$L$  is valid for  $P =$  aaron

- Total time is  $O(m + |\Sigma|)$

# Boyer-More vs. Brute-force Indexing

$$P = cad$$

	$j=0$	$j=1$	$j=2$		
$T$	c	a	b	a	b
$i=0$	c	a	b		

	$j=0$	$j=1$	$j=2$		
	$i=0$	$i=1$	$i=2$		
$T$	c	a	b	a	b
	c	a	b		

## ■ Brute-force

- maintain variables  $i$  and  $j$
- $j$  is the position in the pattern
- $i$  is equal to the current shift
- check is performed by determining if  $T[i + j] = P[j]$

## ■ Boyer-More

- maintain variables  $i$  and  $j$
- $j$  is the position in the pattern
- $i$  is the position in the text where we do the next check
- check is performed by determining if  $T[i] = P[j]$
- current shift is  $i - j$

# Bad Character Heuristic: Shifting Formula

<i>char</i>	a	n	o	r	all others
$L(c)$	1	4	3	2	-1

$T = \text{acranapple}$ ,  $P = \text{aaron}$

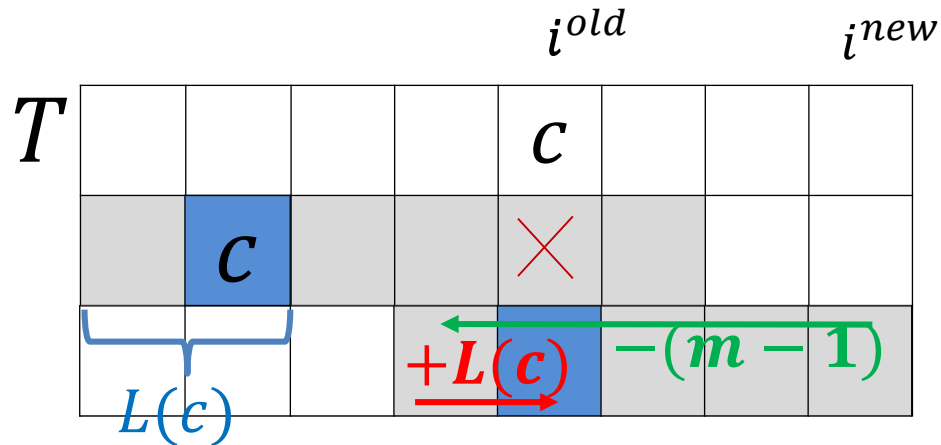
			$j=3$ $i=3$	$j=4$ $i=6$					
a	c	r	a	n	a	p	p	l	e
			o	n					
			[a]			n			

- Let  $L(c)$  be the last occurrence of character  $c$  in  $P$ 
  - $L(\mathbf{a}) = 1$  in our example
- When mismatch occurs at text position  $i$ , pattern position  $j$ , update
  - $j = m - 1$ 
    - start matching at the end of the pattern
  - $i = i + m - 1 - L(c)$
  - for our example
    - $j = 5 - 1 = 4$
    - $i = 3 + 5 - 1 - 1 = 6$



# Bad Character Heuristic: Shifting Formula Explained

- Text character is  $c$  at the mismatch position  $i$  in the text
- $i = i + m - 1 - L(c)$



$$i^{new} - (m - 1) + L(c) = i^{old}$$

$$i^{new} = i^{old} + m - 1 - L(c)$$

$$i = i + m - 1 - L(c)$$

# Bad Character Heuristic: Important Use Condition

- Text character is  $c$  at the mismatch position  $i$  in the text
- $i = i + m - 1 - L(c)$
- Old shift:  $i - j$
- New shift:  $i + (m - 1) - L(c) - (m - 1) = i - L(c)$
- If  $L(c) > j$ , new shift  $<$  old shift, shifts  $P$  in the wrong direction, not useful
  - we already ruled that shift out, no point to come back to it

Example:  $T = \text{acranapple}$ ,  $P = \text{reroa}$

								$j=3$		
								$i=8$		
c	a	c	r	w	a	a	p	a	a	e
				a						
								o	a	
							o	a		

$$L(\mathbf{a}) = 4$$

$$L(\mathbf{a}) > j = 3$$

$$\text{old shift: } i - j = 8 - 3 = 5$$

$$i = 8 + 5 - 1 - 4 = 8$$

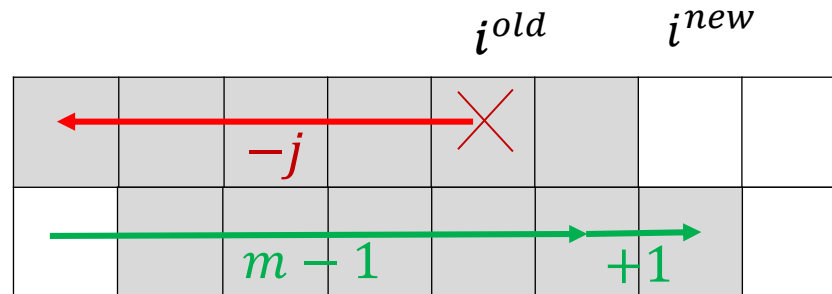
$$j = 5 - 1 = 4$$

$$\text{new shift: } i - j = 8 - 4 = 4$$

- bad character heuristic makes sense to used only if  $L(c) < j$** 
  - $L(c) \neq j$  in case of a mismatch

# Bad Character Heuristic: Brute-Force Step

- If  $L(c) > j$ 
  - pattern would shift in wrong direction if used bad character heuristic
  - therefore, do brute-force step
  - $j = m - 1$
  - $i = i - j + m$



$$i^{old} - j + m - 1 + 1 = i^{new}$$

$$i^{new} = i^{old} - j + m$$

$$i = i - j + m$$

# Bad Character Heuristic: Unified Formula

- If  $L(c) < j$ 
  - $j = m - 1$
  - $i = i + m - 1 - L(c)$
- If  $L(c) > j$ 
  - $j = m - 1$
  - $i = i - j + m$
- Unified formula for  $i$  that works in all cases
$$i = i + m - 1 - \min\{L(c), j - 1\}$$

# Boyer-More Example

<i>char</i>	a	e	p	r	others
$L(c)$	1	3	2	4	-1

$P = \text{paper}$

				$j=4$ $i=4$				$j=4$ $i=7$		$j=4$ $i=9$				$j=3$ $i=13$	$j=4$ $i=14$	$j=4$ $i=15$			
$T$	f	e	e	d	a	l	l	p	o	o	r	p	a	r	r	o	t	s	
				r															$i=7$
				[a]				r											$i=9$
								[p]		r									$i=14$
														e	r				$i=15$
																	r		$i=20$

- Unified formula for  $i$  that works in all cases

$$i = i + m - 1 - \min\{L(c), j - 1\}$$

not found!

# Boyer-Moore Algorithm

*BoyerMoore*( $T, P$ )

$L \leftarrow$  last occurrence array computed from  $P$

$j \leftarrow m - 1$

$i \leftarrow m - 1$

**while**  $i < n$  and  $j \geq 0$  **do** //current guess begins at index  $i - j$

**if**  $T[i] = P[j]$  **then**

$i \leftarrow i - 1$

$j \leftarrow j - 1$

**else**

$i \leftarrow i + m - 1 - \min\{L(c), j - 1\}$

$j \leftarrow m - 1$

**if**  $j = -1$  **return** “found at shift  $i + 1$ ” //  $i$  moved one position to  
// the left of the first char in  $T$

**else return FAIL**

# Good Suffix Heuristic

$P = \text{onobobo}$

$j=3$   
 $i=3$

$T$	o	n	o	o	o	b	o	o	o	i	b	b	o	u	n	d	a	r	y
				b	o	b	o												
			o	n	o	b	o	b	o										

- Text has letters **obo**
- Do the smallest shift so that **obo** fits
- Can precompute this from the pattern itself, before matching starts
  - 'if failure at  $j = 3$ , shift pattern by 2'
- Continue matching from the end of the new shift
- Will not study the precise way to do it

# Boyer-Moore Algorithm with Good Suffix

*BoyerMoore*( $T, P$ )

$L \leftarrow$  last occurrence array computed from  $P$

$S \leftarrow$  good suffix array computed from  $P$

$j \leftarrow m - 1$

$i \leftarrow m - 1$

**while**  $i < n$  and  $j \geq 0$  **do** //current guess begins at index  $i - j$

**if**  $T[i] = P[j]$  **then**

$i \leftarrow i - 1$

$j \leftarrow j - 1$

**else**

$i \leftarrow i + m - 1 - \min\{L(T[i]), S[j]\}$

$j \leftarrow m - 1$

**if**  $j = -1$  **return** “found at shift  $i + 1$ ”

**else return FAIL**



# Boyer-Moore Summary

- Boyer-Moore performs very well, even when using only bad character heuristic
- Worst case run time is  $O(nm)$  with bad character heuristic only, but in practice much faster
- On typical English text, Boyer-Moore looks only at  $\approx 25\%$  of text  $T$
- With good suffix heuristic, can ensure  $O(n + m + |\Sigma|)$  run time
  - no details

# Outline

- String Matching
  - Introduction
  - Karp-Rabin Algorithm
  - Boyer-Moore Algorithm
  - Suffix Trees

# Suffix Tree: Trie of Suffixes

- What if we search for **many patterns**  $P$  within the same **fixed text**  $T$ ?
- **Idea:** preprocess the text  $T$  rather than pattern  $P$
- **Observation:**  $P$  is a substring of  $T$  if and only if  $P$  is a prefix of some suffix of  $T$
- Example:  $P = \text{ish}$

$T = \text{establishment}$

The diagram illustrates the text  $T = \text{establishment}$ . The substring **ish** is highlighted in red. A red bracket above **ish** is labeled "prefix", and a blue bracket below **ishment** is labeled "suffix".

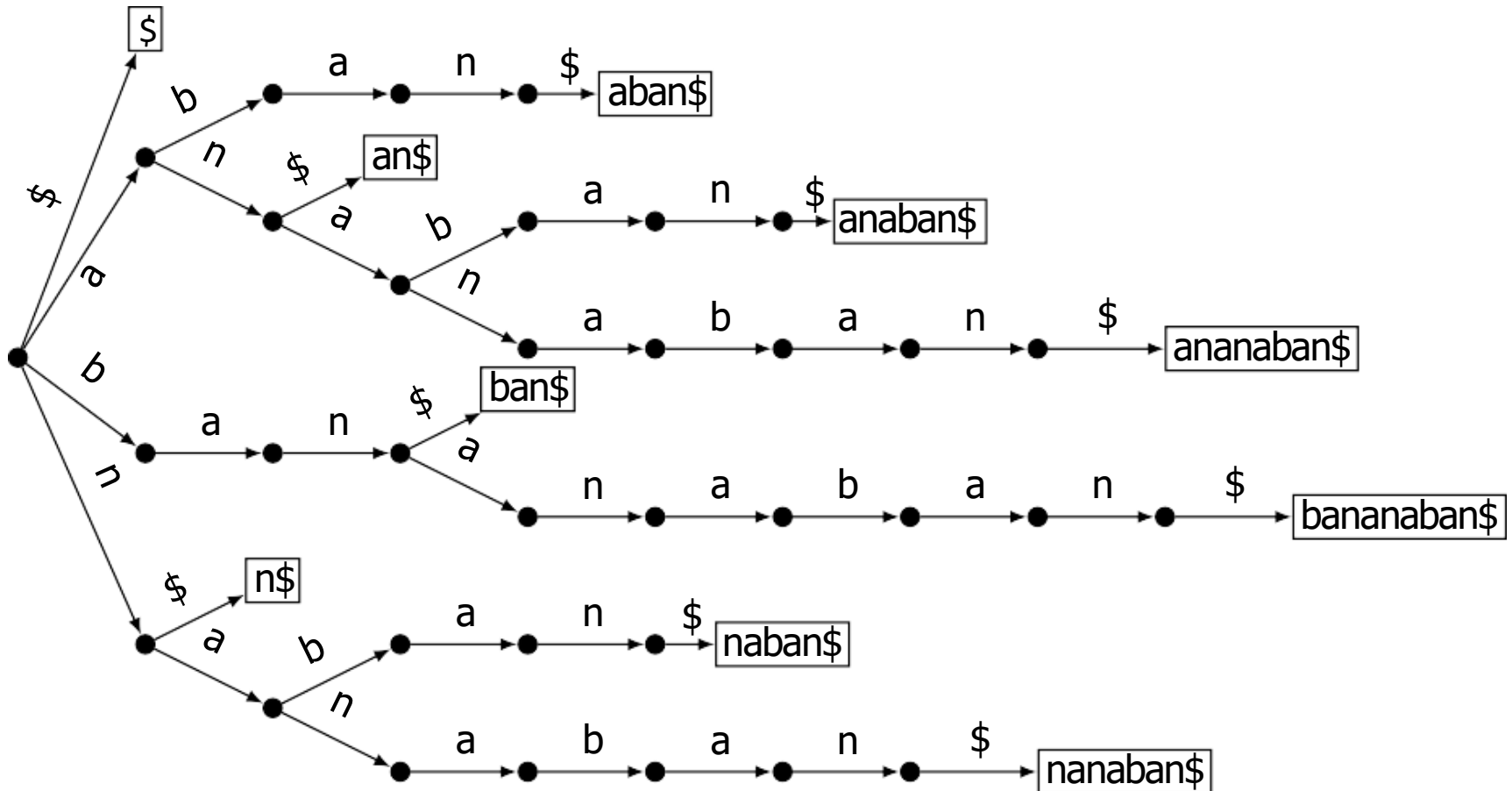
- Store all suffixes of  $T$  in a trie
- To save space
  - use compressed trie
  - store suffixes implicitly via indices into  $T$
- This is called a **suffix tree**

# Trie of suffixes: Example

- $T = \text{bananaban}$

Suffixes = {bananaban, ananaban, nanaban, anaban, naban, aban, ban, an, n,  $\Lambda$ }

$S = \{\text{bananaban}\$, \text{ananaban}\$, \text{nanaban}\$, \text{anaban}\$, \text{naban}\$, \dots, \text{ban}\$, \text{n}\$, \$\}$

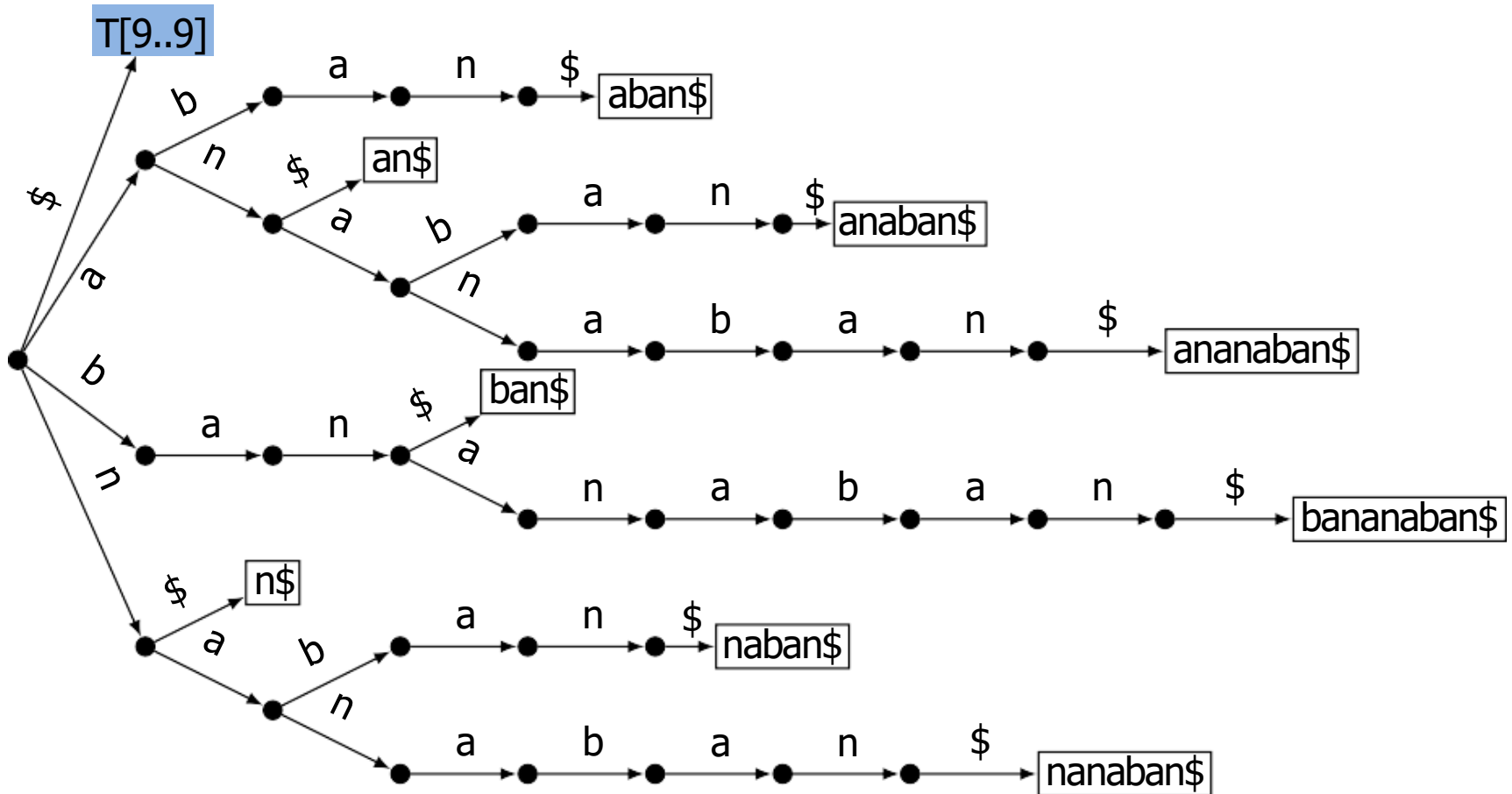


# Trie of suffixes: Example

$T =$ 

<b>0</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>7</b>	<b>8</b>	<b>9</b>
b	a	n	a	n	a	b	a	n	\$

- Store suffixes via indices

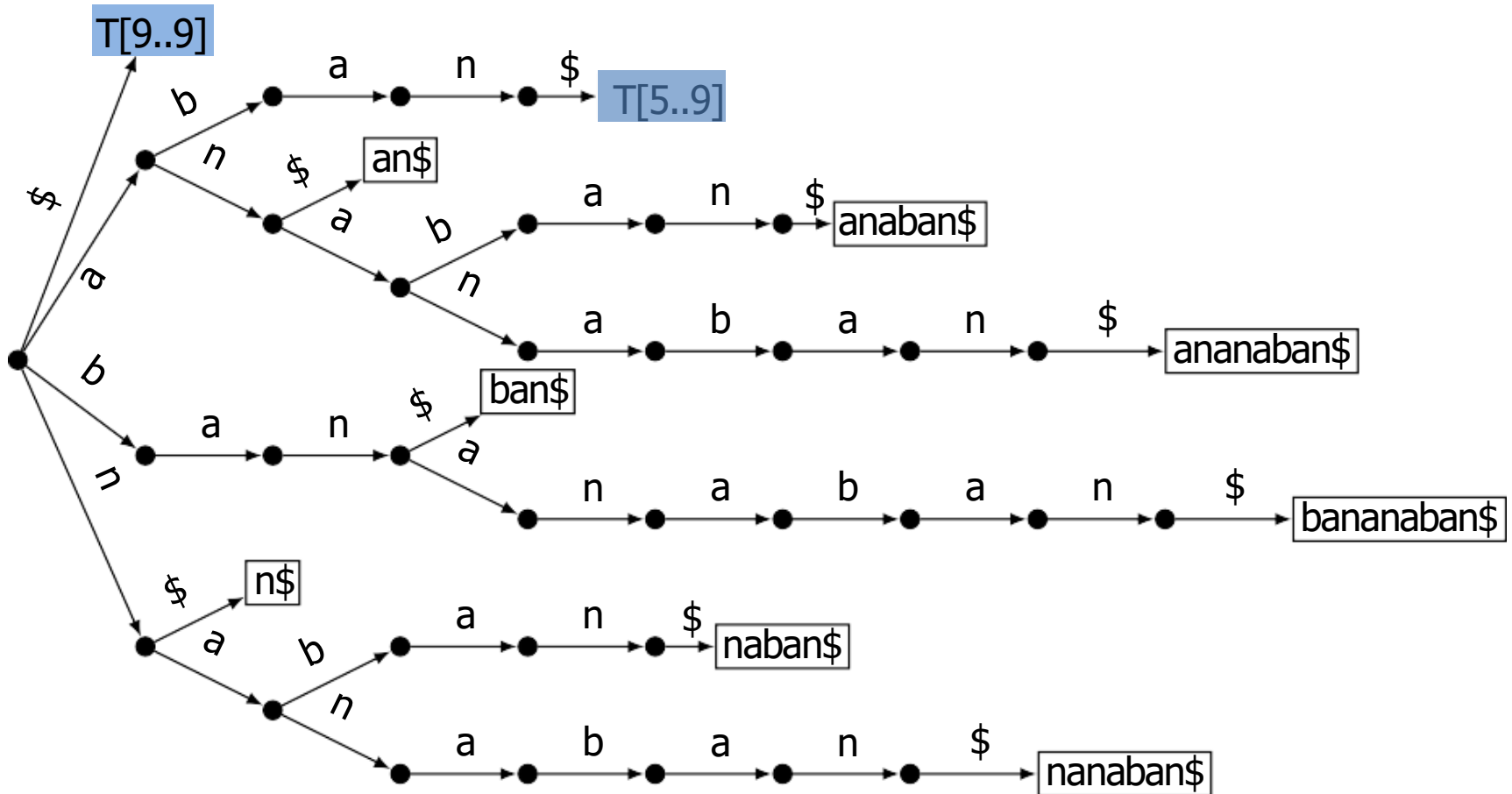


# Trie of suffixes: Example

$T =$ 

0	1	2	3	4	5	6	7	8	9
b	a	n	a	n	a	b	a	n	\$

- Store suffixes via indices

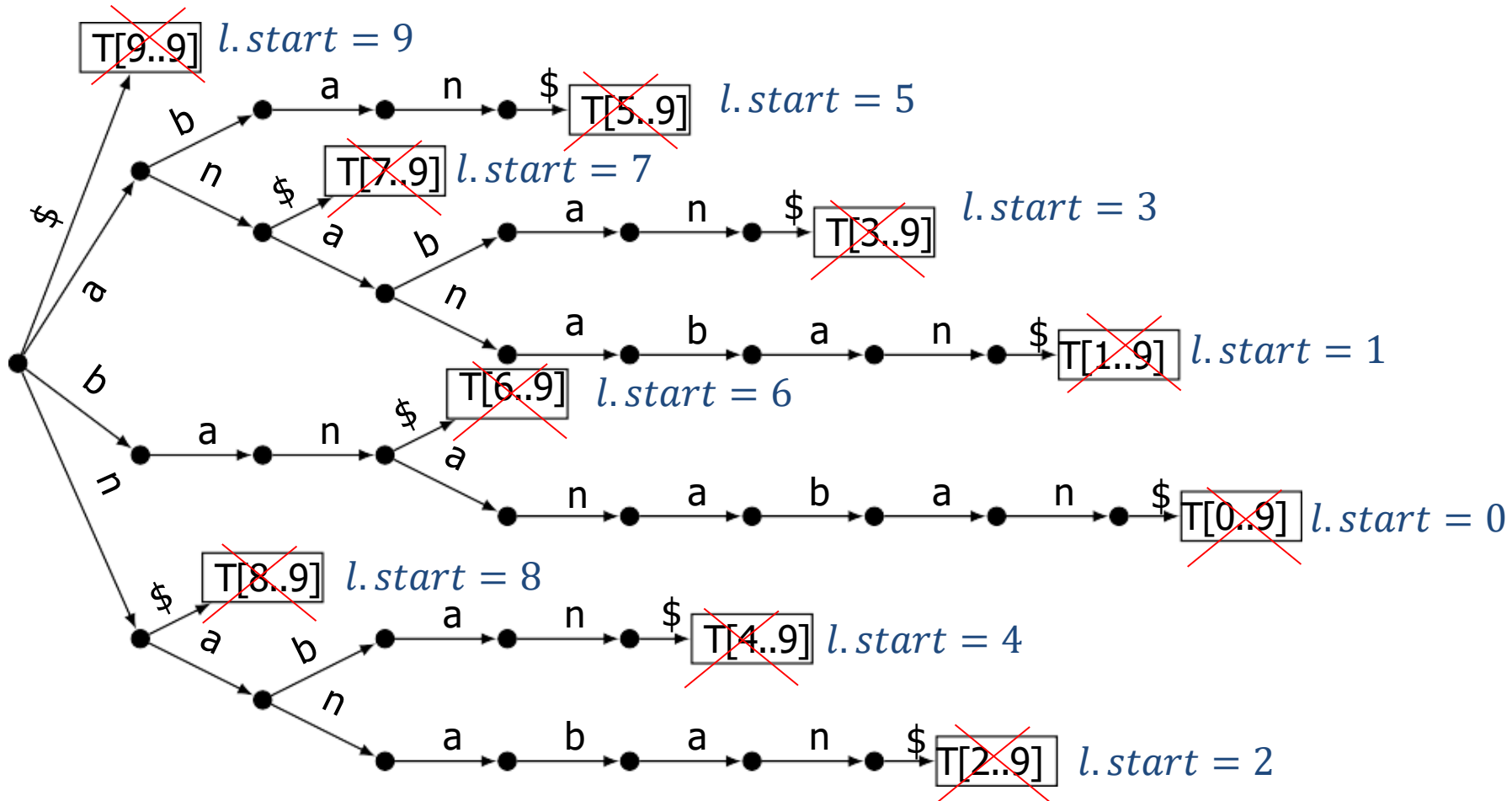


# Tries of suffixes

- In actual implementation, each leaf  $l$  stores the start of its suffix in variable  $l.start$

$T =$

0	1	2	3	4	5	6	7	8	9
b	a	n	a	n	a	b	a	n	\$

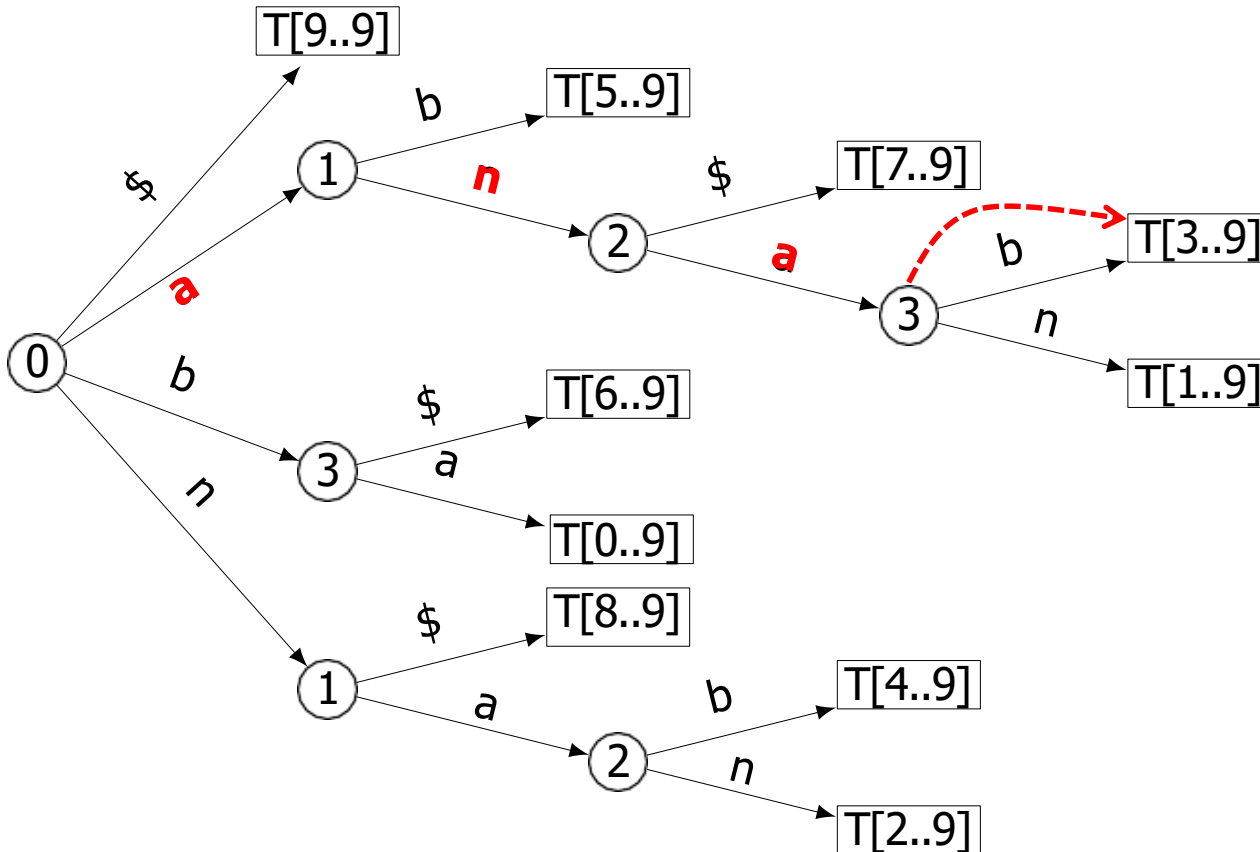


# Suffix tree

$T =$ 

0	1	2	3	4	5	6	7	8	9
b	a	n	a	n	a	b	a	n	\$

- **Suffix tree**: compressed trie of suffixes
- If  $P$  occurs in the text, it is a prefix of one (or more) strings stored in the trie
- Have to modify search in a trie to allow search for a prefix





# Building Suffix Tree

- Building
  - text  $T$  has  $n$  characters and  $n + 1$  suffixes
  - can build suffix tree by inserting each suffix of  $T$  into compressed trie
    - takes  $\Theta(|\Sigma|n^2)$  time
  - there is a way to build a suffix tree of  $T$  in  $\Theta(|\Sigma|n)$  time
    - beyond the course scope
- Pattern Matching
  - essentially search for  $P$  in compressed trie
    - some changes needed, since  $P$  may only be prefix of stored word
  - run-time is
    - $O(|\Sigma|m)$ , assuming each node stores children in a linked list
    - $O(m)$ , assuming each node stores children in an array
- Summary
  - theoretically good, but construction is slow or complicated and lots of space-overhead
  - rarely used in practice