

Estruturas de dados e algoritmos para buscas de padrões

Marco Alves de Alcantara

14 de Dezembro, 2022

Orientadora: Cristina Gomes Fernandes
IME USP



Roteiro

① Introdução

② Vetor de sufixos

③ Vetor LCP

Extensão do vetor LCP

④ Sobre as buscas

⑤ Árvore de sufixos

⑥ Resultados

1 Introdução

2 Vetor de sufixos

3 Vetor LCP

Extensão do vetor LCP

4 Sobre as buscas

5 Árvore de sufixos

6 Resultados

- Queremos buscar as ocorrências de palavras em um texto T

- Queremos buscar as ocorrências de palavras em um texto T
- Estamos pensando em aplicações onde T é um texto bastante longo e fixo (imutável)

- Queremos buscar as ocorrências de palavras em um texto T
- Estamos pensando em aplicações onde T é um texto bastante longo e fixo (imutável)
- Algoritmos como KMP levam tempo $\mathcal{O}(|T|)$ por busca

- Queremos buscar as ocorrências de palavras em um texto T
- Estamos pensando em aplicações onde T é um texto bastante longo e fixo (imutável)
- Algoritmos como KMP levam tempo $\mathcal{O}(|T|)$ por busca
- Queremos algo mais eficiente do que isso

- Queremos buscar as ocorrências de palavras em um texto T
- Estamos pensando em aplicações onde T é um texto bastante longo e fixo (imutável)
- Algoritmos como KMP levam tempo $\mathcal{O}(|T|)$ por busca
- Queremos algo mais eficiente do que isso
 - Iremos pré-processar o texto T antes de realizar as buscas

- Queremos buscar as ocorrências de palavras em um texto T
- Estamos pensando em aplicações onde T é um texto bastante longo e fixo (imutável)
- Algoritmos como KMP levam tempo $\mathcal{O}(|T|)$ por busca
- Queremos algo mais eficiente do que isso
 - ▶ Iremos pré-processar o texto T antes de realizar as buscas
 - » o pré-processamento é relativamente lento (linear em $|T|$)

- Queremos buscar as ocorrências de palavras em um texto T
- Estamos pensando em aplicações onde T é um texto bastante longo e fixo (imutável)
- Algoritmos como KMP levam tempo $\mathcal{O}(|T|)$ por busca
- Queremos algo mais eficiente do que isso
 - ▶ Iremos pré-processar o texto T antes de realizar as buscas
 - » o pré-processamento é relativamente lento (linear em $|T|$)
 - » mas é realizado uma única vez
 - » e permite agilizar **todas** as buscas por qualquer palavra P nesse texto

① Pré-processamento:

- ① Acrescentar um caractere especial \$ ao final de T

① Pré-processamento:

- ① Acrescentar um caractere especial \$ ao final de T
- ② Construir o **vetor de sufixos** (VS) de T

① Pré-processamento:

- ① Acrescentar um caractere especial \$ ao final de T
- ② Construir o **vetor de sufixos** (VS) de T
- ③ Construir o **vetor LCP** de T a partir do VS

① Pré-processamento:

- ① Acrescentar um caractere especial \$ ao final de T
- ② Construir o **vetor de sufixos** (VS) de T
- ③ Construir o **vetor LCP** de T a partir do VS
 - » *e estender o vetor LCP*

① Pré-processamento:

- ① Acrescentar um caractere especial \$ ao final de T
- ② Construir o **vetor de sufixos** (VS) de T
- ③ Construir o **vetor LCP** de T a partir do VS
 - » *e estender o vetor LCP*

② Em seguida, realizamos buscas binárias pelas palavras P no vetor de sufixos de T usando os LCPs estendidos

① Introdução

② **Vetor de sufixos**

③ Vetor LCP

Extensão do vetor LCP

④ Sobre as buscas

⑤ Árvore de sufixos

⑥ Resultados

Vetor de sufixos – definição

Uma estrutura de dados importante para as buscas é o vetor de sufixos:

Vetor de sufixos

- Um **vetor de sufixos** é um vetor de índices $[0..|T|-1]$, onde cada índice representa um sufixo de T
- Se $VS[i] = k$, então $VS[i]$ representa $T[k..]$
- VS é ordenado lexicograficamente com base nos sufixos representados

(https://en.wikipedia.org/wiki/Suffix_array)

Vetor de sufixos – exemplo

Exemplo de VS para o texto $T = \text{“abracadabra$”}$

$T[i..]$	i	$VS[i]$	$T[VS[i]..]$
abracadabra\$	0	11	\$
bracadabra\$	1	10	a\$
racadabra\$	2	7	abra\$
acadabra\$	3	0	abracadabra\$
cadabra\$	4	3	acadabra\$
adabra\$	5	5	adabra\$
dabra\$	6	8	bra\$
abra\$	7	1	bracadabra\$
bra\$	8	4	cadabra\$
ra\$	9	6	dabra\$
a\$	10	9	ra\$
\$	11	2	racadabra\$

Vetor de sufixos – ideia da construção

- Utiliza-se um algoritmo recursivo de divisão e conquista

Vetor de sufixos – ideia da construção

- **Utiliza-se um algoritmo recursivo de divisão e conquista**
 - ▶ O texto T é dividido em 3, alternando qual sufixo fica em qual parte

Vetor de sufixos – ideia da construção

- **Utiliza-se um algoritmo recursivo de divisão e conquista**
 - ▶ O texto T é dividido em 3, alternando qual sufixo fica em qual parte
 - ▶ É construído um texto auxiliar T_{12} , contendo os 3 primeiros caracteres de cada sufixo mod 1 e mod 2

Vetor de sufixos – ideia da construção

- **Utiliza-se um algoritmo recursivo de divisão e conquista**
 - ▶ O texto T é dividido em 3, alternando qual sufixo fica em qual parte
 - ▶ É construído um texto auxiliar T_{12} , contendo os 3 primeiros caracteres de cada sufixo mod 1 e mod 2
 - ▶ Cada tripla é ordenada com o RadixSort, e corresponde a um caractere num alfabeto maior

Vetor de sufixos – ideia da construção

- **Utiliza-se um algoritmo recursivo de divisão e conquista**
 - ▶ O texto T é dividido em 3, alternando qual sufixo fica em qual parte
 - ▶ É construído um texto auxiliar T_{12} , contendo os 3 primeiros caracteres de cada sufixo mod 1 e mod 2
 - ▶ Cada tripla é ordenada com o RadixSort, e corresponde a um caractere num alfabeto maior
 - ▶ Caso haja triplas repetidas, é necessário construir o vetor de sufixos de T_{12} recursivamente, caso contrário, o RadixSort já nos dá esse vetor de sufixos

Vetor de sufixos – ideia da construção

- **Utiliza-se um algoritmo recursivo de divisão e conquista**
 - ▶ O texto T é dividido em 3, alternando qual sufixo fica em qual parte
 - ▶ É construído um texto auxiliar T_{12} , contendo os 3 primeiros caracteres de cada sufixo mod 1 e mod 2
 - ▶ Cada tripla é ordenada com o RadixSort, e corresponde a um caractere num alfabeto maior
 - ▶ Caso haja triplas repetidas, é necessário construir o vetor de sufixos de T_{12} recursivamente, caso contrário, o RadixSort já nos dá esse vetor de sufixos
- **Depois do processo de divisão, fazemos a parte da conquista...**

Vetor de sufixos – ideia da construção – parte 2

- Usamos $VS(T_1)$ para ordenar os sufixos de T_0 (obter $VS(T_0)$)

Vetor de sufixos – ideia da construção – parte 2

- Usamos $VS(T_1)$ para ordenar os sufixos de T_0 (obter $VS(T_0)$)
- Construimos certos pares e triplas formados por 1 (ou 2) caracteres de T_0, T_1 ou T_2 e com 1 rank de $VS(T_{12})$

Vetor de sufixos – ideia da construção – parte 2

- Usamos $VS(T_1)$ para ordenar os sufixos de T_0 (obter $VS(T_0)$)
- Construimos certos pares e triplas formados por 1 (ou 2) caracteres de T_0, T_1 ou T_2 e com 1 rank de $VS(T_{12})$
- Fazemos um procedimento semelhante ao do algoritmo MergeSort para inicializar $VS(T)$, comparando pares ou triplas elemento-por-elemento

Vetor de sufixos – ideia da construção – parte 2

- Usamos $VS(T_1)$ para ordenar os sufixos de T_0 (obter $VS(T_0)$)
- Construimos certos pares e triplas formados por 1 (ou 2) caracteres de T_0 , T_1 ou T_2 e com 1 rank de $VS(T_{12})$
- Fazemos um procedimento semelhante ao do algoritmo MergeSort para inicializar $VS(T)$, comparando pares ou triplas elemento-por-elemento
- Por fim, obtém-se $VS(T)$ em tempo linear em $|T|$.

(<https://gist.github.com/markormesher/59b990fba09972b4737e7ed66912e044>)

(<http://www.cs.cmu.edu/~guyb/realworld/papersS04/KaSa03.pdf>)

① Introdução

② Vetor de sufixos

③ **Vetor LCP**

Extensão do vetor LCP

④ Sobre as buscas

⑤ Árvore de sufixos

⑥ Resultados

Vetor LCP – definição

Além do vetor de sufixos, para que seja possível fazer a busca eficientemente, é necessário também o vetor LCP do texto T :

Vetor LCP

- O **vetor LCP** é um vetor de tamanho $|T| - 1$, indexado a partir de 1, que mede o comprimento do prefixo comum mais longo entre dois sufixos de T
- Se $LCP[i] = k$, então os sufixos de índices $VS[i - 1]$ e $VS[i]$ tem um prefixo comum mais longo de comprimento k

(https://en.wikipedia.org/wiki/LCP_array)

Vetor LCP – exemplo

Vetor LCP para o texto $T = \text{“abracadabra$”}$

i	$VS[i]$	$LCP[i]$	$T[VS[i]..]$
0	11		\$
1	10	0	a\$
2	7	1	<u>a</u> bra\$
3	0	4	<u>abra</u> cadabra\$
4	3	1	<u>a</u> cadabra\$
5	5	1	<u>a</u> dabra\$
6	8	0	bra\$
7	1	3	<u>bra</u> cadabra\$
8	4	0	cadabra\$
9	6	0	dabra\$
10	9	0	ra\$
11	2	2	<u>ra</u> cadabra\$

Vetor LCP – algoritmo de construção

O vetor LCP pode ser construído a partir do VS com o Algoritmo de Kasai

Vetor LCP – algoritmo de construção

O vetor LCP pode ser construído a partir do VS com o Algoritmo de Kasai

```
1  $i \leftarrow 0$ 
2 enquanto  $i < |T|$  faça
3   | rank[VS[i]]  $\leftarrow i$ 
4   |  $i \leftarrow i + 1$ 
5  $h, i \leftarrow 0$ 
6 enquanto  $i < |T|$  faça
7   | se rank[i] > 0 então
8     |  $k \leftarrow \text{VS}[\text{rank}[i] - 1]$ 
9     | enquanto  $T[i + h] = T[k + h]$  faça
10    |   |  $h \leftarrow h + 1$ 
11    |   | LCP[rank[i]]  $\leftarrow h$ 
12    |   | se  $h > 0$  então
13    |   |   |  $h \leftarrow h - 1$ 
14    |   |  $i \leftarrow i + 1$ 
15 retorna LCP
```

- O vetor LCP por si só não ajuda muito nas buscas, já que ele só nos dá informações sobre dois sufixos consecutivos em VS

- O vetor LCP por si só não ajuda muito nas buscas, já que ele só nos dá informações sobre dois sufixos consecutivos em VS
 - ▶ É necessário **estender** o LCP para nos dar informações sobre um **intervalo** de sufixos

- O vetor LCP por si só não ajuda muito nas buscas, já que ele só nos dá informações sobre dois sufixos consecutivos em VS
 - ▶ É necessário **estender** o LCP para nos dar informações sobre um **intervalo** de sufixos
 - ▶ Mais precisamente, sobre os $\mathcal{O}(|T|)$ intervalos possíveis de aparecerem durante uma busca binária em VS

- O vetor LCP por si só não ajuda muito nas buscas, já que ele só nos dá informações sobre dois sufixos consecutivos em VS
 - ▶ É necessário **estender** o LCP para nos dar informações sobre um **intervalo** de sufixos
 - ▶ Mais precisamente, sobre os $\mathcal{O}(|T|)$ intervalos possíveis de aparecerem durante uma busca binária em VS
 - » *por isso, também iremos **pré-processar o vetor LCP**, simulando todos os caminhos possíveis em uma busca binária, e guardar as informações obtidas sobre os LCPs desses intervalos*

- O vetor LCP por si só não ajuda muito nas buscas, já que ele só nos dá informações sobre dois sufixos consecutivos em VS
 - ▶ É necessário **estender** o LCP para nos dar informações sobre um **intervalo** de sufixos
 - ▶ Mais precisamente, sobre os $\mathcal{O}(|T|)$ intervalos possíveis de aparecerem durante uma busca binária em VS
 - » *por isso, também iremos **pré-processar o vetor LCP**, simulando todos os caminhos possíveis em uma busca binária, e guardar as informações obtidas sobre os LCPs desses intervalos*
 - » *durante as buscas, essas informações poderão ser obtidas em tempo $\mathcal{O}(1)$, já que esse pré-processamento só depende de T*

Roteiro

① Introdução

② Vetor de sufixos

③ Vetor LCP

Extensão do vetor LCP

④ Sobre as buscas

⑤ Árvore de sufixos

⑥ Resultados

- Para buscar as ocorrências de uma palavra P , na verdade são necessárias duas buscas binárias

- Para buscar as ocorrências de uma palavra P , na verdade são necessárias duas buscas binárias
 - ▶ Primeiro, busca-se o **predecessor** de P no vetor de sufixos, ou seja, o índice do maior sufixo de T que é lexicograficamente menor que P

- **Para buscar as ocorrências de uma palavra P , na verdade são necessárias duas buscas binárias**
 - ▶ Primeiro, busca-se o **predecessor** de P no vetor de sufixos, ou seja, o índice do maior sufixo de T que é lexicograficamente menor que P
 - ▶ O predecessor serve como um limitante inferior para um intervalo que contenha as ocorrências de P

- **Para buscar as ocorrências de uma palavra P , na verdade são necessárias duas buscas binárias**
 - ▶ Primeiro, busca-se o **predecessor** de P no vetor de sufixos, ou seja, o índice do maior sufixo de T que é lexicograficamente menor que P
 - ▶ O predecessor serve como um limitante inferior para um intervalo que contenha as ocorrências de P
 - ▶ Para o limitante superior, precisamos de um índice a partir do qual P não é mais um prefixo dos sufixos de VS

- **Para buscar as ocorrências de uma palavra P , na verdade são necessárias duas buscas binárias**
 - ▶ Primeiro, busca-se o **predecessor** de P no vetor de sufixos, ou seja, o índice do maior sufixo de T que é lexicograficamente menor que P
 - ▶ O predecessor serve como um limitante inferior para um intervalo que contenha as ocorrências de P
 - ▶ Para o limitante superior, precisamos de um índice a partir do qual P não é mais um prefixo dos sufixos de VS
 - » *podemos criar uma nova palavra Q , cujo predecessor é o limitante superior do intervalo*

- **Para buscar as ocorrências de uma palavra P , na verdade são necessárias duas buscas binárias**
 - ▶ Primeiro, busca-se o **predecessor** de P no vetor de sufixos, ou seja, o índice do maior sufixo de T que é lexicograficamente menor que P
 - ▶ O predecessor serve como um limitante inferior para um intervalo que contenha as ocorrências de P
 - ▶ Para o limitante superior, precisamos de um índice a partir do qual P não é mais um prefixo dos sufixos de VS
 - » *podemos criar uma nova palavra Q , cujo predecessor é o limitante superior do intervalo*
 - » *Q pode ser obtida, por exemplo, incrementando o último caractere de P*

Resumo da busca – exemplo

Suponha, ainda para $T = \text{“abracadabra$”}$, que estejamos procurando as ocorrências de $P = \text{“abra”}$.

Resumo da busca – exemplo

Suponha, ainda para $T = \text{“abracadabra$”}$, que estejamos procurando as ocorrências de $P = \text{“abra”}$.

$T[i..]$	i	$VS[i]$	$T[VS[i]..]$
abracadabra\$	0	11	\$
bracadabra\$	1	10	a\$
racadabra\$	2	7	<u>abra</u> \$
acadabra\$	3	0	<u>abracadabra</u> \$
cadabra\$	4	3	acadabra\$
adabra\$	5	5	adabra\$
dabra\$	6	8	bra\$
abra\$	7	1	bracadabra\$
bra\$	8	4	cadabra\$
ra\$	9	6	dabra\$
a\$	10	9	ra\$
\$	11	2	racadabra\$

Resumo da busca – exemplo

Suponha, ainda para $T = \text{“abracadabra$”}$, que estejamos procurando as ocorrências de $P = \text{“abra”}$.

$T[i..]$	i	$VS[i]$	$T[VS[i]..]$	
abracadabra\$	0	11	\$	
bracadabra\$	1	10	a\$	← $\text{pred}(P = \text{“abra”}) = 1$
racadabra\$	2	7	<u>abra</u> \$	
acadabra\$	3	0	<u>ab</u> racadabra\$	
cadabra\$	4	3	acadabra\$	
adabra\$	5	5	adabra\$	
dabra\$	6	8	bra\$	
abra\$	7	1	bracadabra\$	
bra\$	8	4	cadabra\$	
ra\$	9	6	dabra\$	
a\$	10	9	ra\$	
\$	11	2	racadabra\$	

Resumo da busca – exemplo

Suponha, ainda para $T = \text{“abracadabra$”}$, que estejamos procurando as ocorrências de $P = \text{“abra”}$.

$T[i..]$	i	$VS[i]$	$T[VS[i]..]$	
abracadabra\$	0	11	\$	
bracadabra\$	1	10	a \$	$\leftarrow \text{pred}(P = \text{“abra”}) = 1$
racadabra\$	2	7	<u>abra</u> \$	
acadabra\$	3	0	<u>abra</u> cadabra\$	$\leftarrow \text{pred}(Q = \text{“abrb”}) = 3$
cadabra\$	4	3	acadabra\$	
adabra\$	5	5	adabra\$	
dabra\$	6	8	bra\$	
abra\$	7	1	bracadabra\$	
bra\$	8	4	cadabra\$	
ra\$	9	6	dabra\$	
a\$	10	9	ra\$	
\$	11	2	racadabra\$	

Resumo da busca – exemplo – conclusões

- P ocorre em T se, e somente se, $\text{pred}(P) \neq \text{pred}(Q)$

Resumo da busca – exemplo – conclusões

- P ocorre em T se, e somente se, $\text{pred}(P) \neq \text{pred}(Q)$
- O número de ocorrências é $\text{pred}(Q) - \text{pred}(P)$

Resumo da busca – exemplo – conclusões

- P ocorre em T se, e somente se, $\text{pred}(P) \neq \text{pred}(Q)$
- O número de ocorrências é $\text{pred}(Q) - \text{pred}(P)$
- As posições onde P ocorre em T são os sufixos em $\text{VS}[\text{pred}(P) + 1..\text{pred}(Q)]$

Resumo da busca – exemplo – conclusões

- P ocorre em T se, e somente se, $\text{pred}(P) \neq \text{pred}(Q)$
- O número de ocorrências é $\text{pred}(Q) - \text{pred}(P)$
- As posições onde P ocorre em T são os sufixos em $\text{VS}[\text{pred}(P) + 1..\text{pred}(Q)]$

$T = \text{“}\underline{\text{abracad}}\underline{\text{abra}}\text{”}$
 $\uparrow 0$ $\uparrow 7$

Busca – visão mais detalhada

- Na busca binária para encontrar o predecessor, além dos extremos L e R da busca, também guardamos índices l e r com o número de caracteres casados entre P e os sufixos dos extremos L e R

- Na busca binária para encontrar o predecessor, além dos extremos L e R da busca, também guardamos índices l e r com o número de caracteres casados entre P e os sufixos dos extremos L e R
 - ▶ A finalidade desses índices l e r , bem como do pré-processamento realizado no vetor LCP, é minimizar o número de comparações diretas entre caracteres de P e T , a fim de agilizar a busca

Busca – visão mais detalhada – caso 0

Antes de fazer a busca binária, precisamos verificar se a palavra é maior que o último sufixo do VS. Se for, o predecessor já é o último sufixo.

Busca – visão mais detalhada – caso 0

Antes de fazer a busca binária, precisamos verificar se a palavra é maior que o último sufixo do VS. Se for, o predecessor já é o último sufixo.

i	$VS[i]$	$T[VS[i]..]$	
0	11	\$	
1	10	a\$	
2	7	abra\$	
3	0	abracadabra\$	
4	3	acadabra\$	
5	5	adabra\$	
6	8	bra\$	
7	1	bracadabra\$	
8	4	cadabra\$	
9	6	dabra\$	
10	9	ra\$	
11	2	rac adabra\$	$\leftarrow \text{pred}(P = \text{"rat"}) = 2$

Busca – visão mais detalhada – caso 1

Se $l = r$, então é necessário comparar caracteres de P e T para decidir como prosseguir com a busca (exemplo para $P = \text{“AAGG”}$)

Busca – visão mais detalhada – caso 1

Se $l = r$, então é necessário comparar caracteres de P e T para decidir como prosseguir com a busca (exemplo para $P = \text{“AAGG”}$)

	T[VS[i]..]		T[VS[i]..]
	(...)		(...)
L	<u>AA</u> CCT...		L <u>AA</u> CCT
	<u>AA</u> CTAC...		AACTAC...
	<u>AA</u> GAT...	⇒	M <u>AA</u> GAT...
	<u>AA</u> GGCT...		<u>AA</u> GGCT...
M	<u>AA</u> GGTTAC...		R <u>AA</u> GGTTAC...
	<u>AA</u> GGTTATG...		(...)
	<u>AA</u> GTC...		.
	<u>AA</u> GTT...		.
R	<u>AA</u> T...		
	(...)		

- Se $l > r$, há 3 sub-casos possíveis dependendo da relação entre l e o LCP do intervalo $[L, M]$

- Se $l > r$, há 3 sub-casos possíveis dependendo da relação entre l e o LCP do intervalo $[L, M]$
 - ▶ O caso 3 ($l < r$) não será analisado aqui por ser simétrico ao caso 2 baseado na relação entre r e o LCP do intervalo $[M, R]$

Busca – visão mais detalhada – caso 2

- Se $l > r$, há 3 sub-casos possíveis dependendo da relação entre l e o LCP do intervalo $[L, M]$
 - ▶ O caso 3 ($l < r$) não será analisado aqui por ser simétrico ao caso 2 baseado na relação entre r e o LCP do intervalo $[M, R]$
- Iremos agora ver os 3 sub-casos do caso 2

Busca – visão mais detalhada – sub-caso 2.1

Se $l > r$ e $l < \text{LCP}[L, M]$, é possível avançar a busca sem comparar caracteres (exemplo para $P = \text{“AACT”}$)

Busca – visão mais detalhada – sub-caso 2.1

Se $l > r$ e $l < \text{LCP}[L, M]$, é possível avançar a busca sem comparar caracteres (exemplo para $P = \text{“AACT”}$)

	T[VS[i]..]		T[VS[i]..]
	(...)		
L	<u>AA</u> CAC...		
	<u>AA</u> CCAG...		
	<u>AA</u> CCAT...		
	<u>AA</u> CCCG...		(...)
M	<u>AA</u> CCGT...	L	<u>AA</u> CCGT...
	<u>AA</u> CCTG...		<u>AA</u> CCTG...
	AACTC...	⇒ M	<u>AA</u> CTC...
	AAGTT...		AAGTT...
R	<u>AA</u> T...	R	<u>AA</u> T...
	(...)		(...)

Busca – visão mais detalhada – sub-caso 2.2

Se $l > r$ e $l > \text{LCP}[L, M]$, é possível avançar a busca sem comparar caracteres (exemplo para $P = \text{“AACG”}$)

Busca – visão mais detalhada – sub-caso 2.2

Se $l > r$ e $l > \text{LCP}[L, M]$, é possível avançar a busca sem comparar caracteres (exemplo para $P = \text{“AACG”}$)

	T[VS[i]..]		T[VS[i]..]
	(...)		(...)
L	<u>AAC</u> CAC...		L <u>AAC</u> CAC...
	<u>AACC</u> AG...		<u>AACC</u> AG...
	<u>AACT</u> ...	⇒	M <u>AACT</u> ...
	<u>AAGAC</u> ...		<u>AAGAC</u> ...
M	<u>AAGTG</u> ...		R <u>AA</u> GTG...
	AAGTTG...		(...)
	AAGTTTC...		
	AATC...		
R	<u>AA</u> TG...		
	(...)		

Busca – visão mais detalhada – sub-caso 2.3

Se $l > r$ e $l = \text{LCP}[L, M]$, é necessário comparar caracteres diretamente para avançar (exemplo para $P = \text{“AAT”}$)

Busca – visão mais detalhada – sub-caso 2.3

Se $l > r$ e $l = \text{LCP}[L, M]$, é necessário comparar caracteres diretamente para avançar (exemplo para $P = \text{“AAT”}$)

	T[VS[i]..]		T[VS[i]..]
	(...)		
L	<u>AA</u> CCTACT...		
	<u>AA</u> CCAG...		
	<u>AA</u> CT...		
	<u>AA</u> GACG...		(...)
M	<u>AA</u> GTG...	L	<u>AA</u> GTG...
	AAGTTG...		<u>AA</u> GTTG...
	AAGTTTC... \Rightarrow	M	<u>AA</u> GTTTC...
	AAT...		AAT...
R	<u>A</u> CCTG...	R	<u>A</u> CCTG...
	(...)		(...)

- **Por meio de uma análise amortizada, é possível concluir que a busca pelo predecessor de uma palavra P leva tempo $\mathcal{O}(|P| + \log(|T|))$**
- **Portanto, contar o número de ocorrências de P também leva tempo $\mathcal{O}(|P| + \log(|T|))$**

Busca – complexidade

- Por meio de uma análise amortizada, é possível concluir que a busca pelo predecessor de uma palavra P leva tempo $\mathcal{O}(|P| + \log(|T|))$
- Portanto, contar o número de ocorrências de P também leva tempo $\mathcal{O}(|P| + \log(|T|))$
- Encontrar as ocorrências leva tempo $\mathcal{O}(|P| + \log(|T|) + x)$, onde x é o número de ocorrências de P

Roteiro

① Introdução

② Vetor de sufixos

③ Vetor LCP

Extensão do vetor LCP

④ Sobre as buscas

⑤ **Árvore de sufixos**

⑥ Resultados

Árvore de sufixos – introdução

- Uma alternativa a realizar as buscas no vetor de sufixos é buscar usando a árvore de sufixos de T

Árvore de sufixos – introdução

- **Uma alternativa a realizar as buscas no vetor de sufixos é buscar usando a árvore de sufixos de T**
 - ▶ Uma desvantagem de usar árvore de sufixos é que elas gastam mais espaço

Árvore de sufixos – introdução

- **Uma alternativa a realizar as buscas no vetor de sufixos é buscar usando a árvore de sufixos de T**
 - ▶ Uma desvantagem de usar árvore de sufixos é que elas gastam mais espaço
 - ▶ Uma vantagem é que é feita uma única busca em vez de duas, e, além disso, ela é feita em tempo $\mathcal{O}(|P|)$

Árvore de sufixos – introdução

- **Uma alternativa a realizar as buscas no vetor de sufixos é buscar usando a árvore de sufixos de T**
 - ▶ Uma desvantagem de usar árvore de sufixos é que elas gastam mais espaço
 - ▶ Uma vantagem é que é feita uma única busca em vez de duas, e, além disso, ela é feita em tempo $\mathcal{O}(|P|)$
- **A árvore de sufixos pode ser obtida a partir do vetor de sufixos e LCP em tempo linear em $|T|$ por meio de uma árvore cartesiana**

Árvore de sufixos – introdução

- **Uma alternativa a realizar as buscas no vetor de sufixos é buscar usando a árvore de sufixos de T**
 - ▶ Uma desvantagem de usar árvore de sufixos é que elas gastam mais espaço
 - ▶ Uma vantagem é que é feita uma única busca em vez de duas, e, além disso, ela é feita em tempo $\mathcal{O}(|P|)$
- **A árvore de sufixos pode ser obtida a partir do vetor de sufixos e LCP em tempo linear em $|T|$ por meio de uma árvore cartesiana**
- **Alternativamente, pode ser construída direto de T , sem o vetor de sufixos, por meio do Algoritmo de Ukkonen**

Árvore de sufixos

- Uma **árvore de sufixos** de um texto T é uma árvore tal que:
 - ▶ Cada nó pode ter até $|\Sigma|$ filhos, onde Σ é o alfabeto do texto;

Árvore de sufixos – definição

Árvore de sufixos

- Uma **árvore de sufixos** de um texto T é uma árvore tal que:
 - ▶ Cada nó pode ter até $|\Sigma|$ filhos, onde Σ é o alfabeto do texto;
 - ▶ Cada aresta entre um nó pai e filho contém um **rótulo**, que é um intervalo do texto T ;

Árvore de sufixos – definição

Árvore de sufixos

- Uma **árvore de sufixos** de um texto T é uma árvore tal que:
 - ▶ Cada nó pode ter até $|\Sigma|$ filhos, onde Σ é o alfabeto do texto;
 - ▶ Cada aresta entre um nó pai e filho contém um **rótulo**, que é um intervalo do texto T ;
 - ▶ Cada nó filho é indexado a partir do nó pai pelo primeiro caractere do rótulo da aresta entre eles;

Árvore de sufixos – definição

Árvore de sufixos

- Uma **árvore de sufixos** de um texto T é uma árvore tal que:
 - ▶ Cada nó pode ter até $|\Sigma|$ filhos, onde Σ é o alfabeto do texto;
 - ▶ Cada aresta entre um nó pai e filho contém um **rótulo**, que é um intervalo do texto T ;
 - ▶ Cada nó filho é indexado a partir do nó pai pelo primeiro caractere do rótulo da aresta entre eles;
 - ▶ Os rótulos são o mais longo possível, isso é, se $T \neq \$$, não há nós com exatamente um filho (a árvore é comprimida);

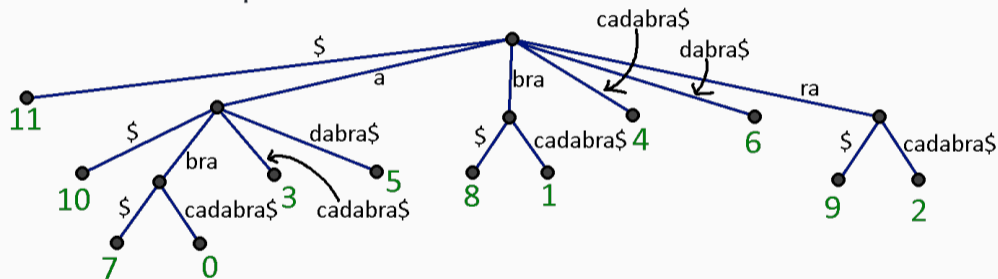
Árvore de sufixos – definição

Árvore de sufixos

- Uma **árvore de sufixos** de um texto T é uma árvore tal que:
 - ▶ Cada nó pode ter até $|\Sigma|$ filhos, onde Σ é o alfabeto do texto;
 - ▶ Cada aresta entre um nó pai e filho contém um **rótulo**, que é um intervalo do texto T ;
 - ▶ Cada nó filho é indexado a partir do nó pai pelo primeiro caractere do rótulo da aresta entre eles;
 - ▶ Os rótulos são o mais longo possível, isso é, se $T \neq \$$, não há nós com exatamente um filho (a árvore é comprimida);
 - ▶ Cada folha contém o índice de um sufixo de T , e todos os sufixos de T são representados por uma folha.

Árvore de sufixos – exemplo

Árvore de sufixos para $T = \text{“abracadabra$”}$



Os índices das folhas aparecem na mesma ordem que no vetor de sufixos

- A busca consiste em simplesmente encontrar o nó menos fundo n cujo caminho até a raiz contém P

- **A busca consiste em simplesmente encontrar o nó menos fundo n cujo caminho até a raiz contém P**
 - ▶ Se n não existir, então P não ocorre em T

- **A busca consiste em simplesmente encontrar o nó menos fundo n cujo caminho até a raiz contém P**
 - ▶ Se n não existir, então P não ocorre em T
 - ▶ Se existir, então o número de ocorrências de P é igual ao número de folhas abaixo de n , e as ocorrências são os índices dessas folhas

Árvore de sufixos – busca

- **A busca consiste em simplesmente encontrar o nó menos fundo n cujo caminho até a raiz contém P**
 - ▶ Se n não existir, então P não ocorre em T
 - ▶ Se existir, então o número de ocorrências de P é igual ao número de folhas abaixo de n , e as ocorrências são os índices dessas folhas
- **Para encontrar n , apenas tente percorrer a árvore a partir do nó atual procurando o nó indexado pelo próximo caractere de P , e verifique se P casa com o rótulo da aresta entre eles**

Roteiro

① Introdução

② Vetor de sufixos

③ Vetor LCP

Extensão do vetor LCP

④ Sobre as buscas

⑤ Árvore de sufixos

⑥ **Resultados**

Resultados – Complexidade assintótica de tempo

Complexidade	KMP	VS e LCP	Árv. suf.
Pré-processamento único	não tem	$\mathcal{O}(T)$	$\mathcal{O}(T)$
Pré-processamento por busca	$\mathcal{O}(P)$	não tem	não tem
Contar as ocorrências	$\mathcal{O}(T)$	$\mathcal{O}(P + \log(T))$	$\mathcal{O}(P)$
Encontrar todas as ocorrências	$\mathcal{O}(T)$	$\mathcal{O}(P + \log(T) + x)$	$\mathcal{O}(P + x)$

Tabela 1: Comparação da complexidade de tempo de diferentes algoritmos de busca, onde x denota o número de ocorrências de P em T .

Buscas de Padrões em Textos

① Introdução

② Vetor de sufixos

③ Vetor LCP

Extensão do vetor LCP

④ Sobre as buscas

⑤ Árvore de sufixos

⑥ Resultados

Perguntas