

# IV107 Bioinformatika I

## Přednáška 10

Katedra informačních technologií  
Masarykova Univerzita Brno

Jaro 2010

### Řetězce a algoritmy na řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu hledaného motivu

Algoritmus využívající analýzu prohledávaného řetězce

### Hledání opakování

Tandemové opakování

Palindromy

### Srovnávání dvou sekvencí

Vylepšení pro maximálně  $k$  chyb

### Příště

video HHMI

## Řetězce a algoritmy na řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu hledaného motivu

Algoritmus využívající analýzu prohledávaného řetězce

## Hledání opakování

Tandemové opakování

Palindromy

## Srovnávání dvou sekvencí

DP - Needleman-Wunsch

Vylepšení pro maximálně  $k$  chyb

video HHMI

### Řetězce a algoritmy na řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu hledaného motivu

Algoritmus využívající analýzu prohledávaného řetězce

### Hledání opakování

Tandemové opakování

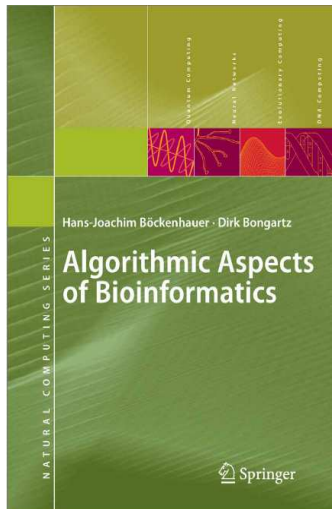
Palindromy

### Srovnávání dvou sekvencí

Vylepšení pro maximálně  $k$  chyb

### Příště

video HHMI



## Řetězce a algoritmy na řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu hledaného motivu

Algoritmus využívající analýzu prohledávaného řetězce

## Hledání opakování

Tandemové opakování

Palindromy

## Srovnávání dvou sekvencí

Vylepšení pro maximálně  $k$  chyb

## Příště

video HHMI

abeceda	$\{\epsilon, a, c, g, t\}$
podřetězec	aag <b>gtacg</b> cgt
prefix	<b>gtacg</b> cgtggt
suffix	cgtat <b>gtacg</b>

## Řetězce a algoritmy na řetězcích

### Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu hledaného motivu

Algoritmus využívající analýzu prohledávaného řetězce

### Hledání opakování

Tandemové opakování

Palindromy

### Srovnávání dvou sekvencí

Vylepšení pro maximálně  $k$  chyb

### Příště

video HHMI

konkatenace  
průnik  
sjednocení

$x=cgcat$   $y=att$   $x \cdot y=cgcatatt$

$x=cgcat$   $y=att$   $Over(x, y)=at$

$x=cgcat$   $y=att$   $\langle x, y \rangle=cgcat$

Řetězce a algoritmy na  
řetězcích

## Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu  
hledaného motivu

Algoritmus využívající analýzu  
prohledávaného řetězce

## Hledání opakování

Tandemové opakování

Palindromy

## Srovnávání dvou sekvencí

Vylepšení pro maximálně  $k$  chyb

## Příště

video HHMI

Cílem je zjistit všechny pozice delšího řetězce, na kterých se vyskytuje kratší řetězec

- ▶ přesný výskyt
- ▶ přibližný výskyt

řetězec  $t$  dlouhý ( $n$ ), např genomová sekvence  
motiv  $p$  krátký ( $m$ ), např `cgcgggctgggtggctcg`

Řetězce a algoritmy na řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu hledaného motivu

Algoritmus využívající analýzu prohledávaného řetězce

Hledání opakování

Tandemové opakování

Palindromy

Srovnávání dvou sekvencí

Vylepšení pro maximálně  $k$  chyb

Příště

video HHMI

## Řetězce a algoritmy na řetězcích

Základní pojmy

### Základní algoritmy

Algoritmus využívající analýzu hledaného motivu

Algoritmus využívající analýzu prohledávaného řetězce

## Hledání opakování

Tandemové opakování

Palindromy

## Srovnávání dvou sekvencí

Vylepšení pro maximálně  $k$  chyb

## Příště

video HHMI

```
a c t g t g t a t g a a a t c g c
1..n → t g t c a
      1..m →
```

Složitost:  $O(mn)$

Řetězce a algoritmy na  
řetězcích

Základní pojmy

**Základní algoritmy**

Algoritmus využívající analýzu  
hledaného motivu

Algoritmus využívající analýzu  
prohledávaného řetězce

Hledání opakování

Tandemové opakování

Palindromy

Srovnávání dvou sekvencí

Vylepšení pro maximálně  $k$  chyb

Příště

video HHMI

a c t g t g t a t g a a a t c g c

→ g a t c a t

x ↑ ↑ ←

*máme v motivu další t?*

a c t g t g t a t g a a a t c g c

+1 → g a **t** c a t

*kde máme v motivu další výskyt suffixu at?*

a c t g t g t a t g a a a t c g c

+3 → g **a t** c a t

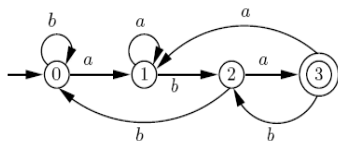
Realizujeme krok, který je větší

Složitost

konstrukce:  $O(\|abeceda\|.m)$

hledání:  $O(mn)$  (v praxi ale blíže k  $O(n)$ )





Automat vytvořen z motivu  $p$  postupně čte symboly z řetězce  $m$ . Koncový stav automatu dosáhneme po načtení celého hledaného motivu.

## Řetězce a algoritmy na řetězcích

Základní pojmy

Základní algoritmy

**Algoritmus využívající analýzu hledaného motivu**

Algoritmus využívající analýzu prohledávaného řetězce

## Hledání opakování

Tandemové opakování

Palindromy

## Srovnávání dvou sekvencí

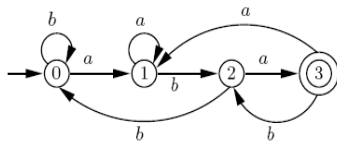
Vylepšení pro maximálně  $k$  chyb

## Příště

video HHMI

t=bababaa p=aba

$\epsilon$	0
b	0
ba	1
bab	2
<b>baba</b>	<b>3</b>
babab	2
bab <b>aba</b>	<b>3</b>
bababaa	1



## Složitost

konstrukce: naivní  $O(m^3)$ ; optimální  $O(\|abceda\|.m)$

hledání:  $O(n)$

Řetězce a algoritmy na  
řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu  
hledaného motivu

Algoritmus využívající analýzu  
prohledávaného řetězce

Hledání opakování

Tandemové opakování

Palindromy

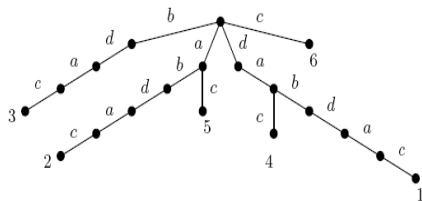
Srovnávání dvou sekvencí

Vylepšení pro maximálně  $k$  chyb

Příště

video HHMI

# Suffixový strom pro řetězec dabdac



## Řetězce a algoritmy na řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu hledaného motivu

**Algoritmus využívající analýzu prohledávaného řetězce**

## Hledání opakování

Tandemové opakování

Palindromy

## Srovnávání dvou sekvencí

Vylepšení pro maximálně  $k$  chyb

## Příště

video HHMI

# Kompaktní suffixový strom pro řetězec

aaabbbbc

Řetězce a algoritmy na řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu hledaného motivu

Algoritmus využívající analýzu prohledávaného řetězce

Hledání opakování

Tandemové opakování

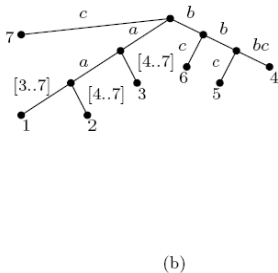
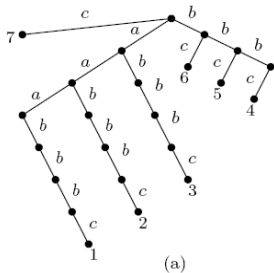
Palindromy

Srovnávání dvou sekvencí

Vylepšení pro maximálně k chyb

Příště

video HHMI



Konstrukce:  $O(n \cdot \log n)$   
Hledání:  $O(m \cdot \|abeceda\| + k)$

# Sufixové pole - ukazovatele na polohy suffixů seřazené lexikograficky

Dlouho bylo považováno za méně kvalitní datovou strukturu, protože neobsahuje přímo informace o společných prefixech. Ty lze však spočítat do lcp pole (least common prefix) tak, že konstrukce pole i stromu má stejnou složitost.

$t = \text{dabdac}$

$sa(t) = 6, 1, 4, 2, 5, 0, 3$

$rank(t) = 5, 1, 3, 6, 2, 4, 0$

$lcp(t) = 0, 0, 1, 0, 0, 0, 2$

6 0

1 0 abdac

4 1 ac

2 0 bdac

5 0 c

0 0 dabdac

3 2 dac

Řetězce a algoritmy na  
řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu  
hledaného motivu

Algoritmus využívající analýzu  
prohledávaného řetězce

Hledání opakování

Tandemové opakování

Palindromy

Srovnávání dvou sekvencí

Vylepšení pro maximálně  $k$  chyb

Příště

video HHMI

## Řetězce a algoritmy na řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu hledaného motivu

Algoritmus využívající analýzu prohledávaného řetězce

## Hledání opakování

Tandemové opakování

Palindromy

## Srovnávání dvou sekvencí

DP - Needleman-Wunsch

Vylepšení pro maximálně  $k$  chyb

video HHMI

Řetězce a algoritmy na  
řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu  
hledaného motivu

Algoritmus využívající analýzu  
prohledávaného řetězce

**Hledání opakování**

Tandemové opakování

Palindromy

**Srovnávání dvou sekvencí**

Vylepšení pro maximálně  $k$  chyb

**Příště**

video HHMI

# Tandemová a palindromická opakování nesou biologický i praktický význam

**palindrom** možná sekundární struktura DNA nebo RNA  
**tandem** regulace genů, telomery, identifikace jedinců  
z DNA

Řetězce a algoritmy na  
řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu  
hledaného motivu

Algoritmus využívající analýzu  
prohledávaného řetězce

Hledání opakování

**Tandemové opakování**

Palindromy

Srovnávání dvou sekvencí

Vylepšení pro maximálně k chyb

Příště

video HHMI

# Nejdelší společný prefix dvou pozic

t g c a g a a g c a g a t c c t g a c g  
↑ ↑

Složitost naivního algoritmu  $O(n^3)$

Řetězce a algoritmy na  
řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu  
hledaného motivu

Algoritmus využívající analýzu  
prohledávaného řetězce

Hledání opakování

**Tandemové opakování**

Palindromy

Srovnávání dvou sekvencí

Vylepšení pro maximálně  $k$  chyb

Příště

video HHMI





- ▶ konstrukce stromu:  $O(n \cdot \log n)$
- ▶ hledání lcp pro dvě konkrétní pozice  $O(n \cdot \log n)$
- ▶ Prohledávání sekvence

Složitost:  $O(n \cdot (\log n)^2 + p)$

Řetězce a algoritmy na  
řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu  
hledaného motivu

Algoritmus využívající analýzu  
prohledávaného řetězce

Hledání opakování

**Tandemové opakování**

Palindromy

Srovnávání dvou sekvencí

Vylepšení pro maximálně  $k$  chyb

Příště

video HHMI

# Palindromy - nejdelší společný prefix mezi originální a komplementární sekvencí umožňuje urychlení hledání podobně jako pro tandemové opakování

↓ 8

t	g	c	a	g	a	a	g	c	t	t	c	t	g	t	c	t	g	a	c	g	
a	c	g	t	c	t	t	c	g	a	a	g	a	c	a	g	a	c	t	g	c	
								↑ 9*													

Složitost naivního algoritmu  $O(n^3)$

Složitost naivního algoritmu  $O(nl)$  (pro omezenou vzdálenost a délku)

Složitost s použitím suffixových struktur  $O(n)$

Řetězce a algoritmy na řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu hledaného motivu

Algoritmus využívající analýzu prohledávaného řetězce

Hledání opakování

Tandemové opakování

Palindromy

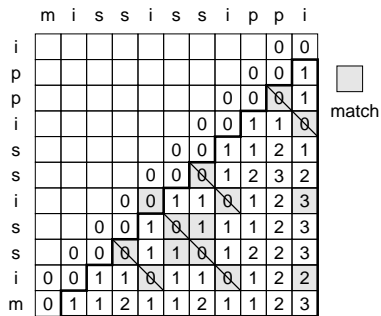
Srovnávání dvou sekvencí

Vylepšení pro maximálně  $k$  chyb

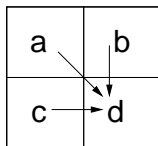
Příště

video HHMI

# Využití DP pro identifikaci palindromů



a)



$$d = \min \begin{cases} a & \text{match} \\ a+1 & \text{mismatch} \\ b+1 & \text{insertion} \\ c+1 & \text{deletion} \end{cases}$$

b)

Řetězce a algoritmy na řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu hledaného motivu

Algoritmus využívající analýzu prohledávaného řetězce

Hledání opakování

Tandemové opakování

Palindromy

Srovnávání dvou sekvencí

Vylepšení pro maximálně k chyb

Příště

video HHMI

# Využití SA a LCP k rychlému postupu po diagonále

Řetězce a algoritmy na řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu hledaného motivu

Algoritmus využívající analýzu prohledávaného řetězce

Hledání opakování

Tandemové opakování

Palindromy

Srovnávání dvou sekvencí

Vylepšení pro maximálně k chyb

Příště

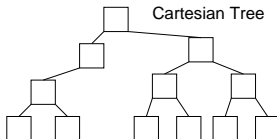
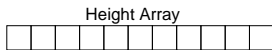
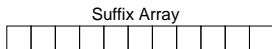
video HHMI

	m	i	s	s	i	s	s	i	p	p	i
i										0	0
p									0	0	1
p								0	0	0	1
i							0	0	1	1	0
s						0	0	1	1	2	1
s					0	0	0	1	1	2	3
i				0	0	1	0	1	1	2	3
s			0	0	1	0	1	1	1	2	3
s		0	0	0	1	0	0	1	1	2	3
i	0	0	1	0	1	0	1	0	1	2	2
m	0	1	1	2	1	1	2	1	1	2	3

→ Longest Common Prefix

- - → Neighboring cells with worse scores

a)



b)

## Řetězce a algoritmy na řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu hledaného motivu

Algoritmus využívající analýzu prohledávaného řetězce

## Hledání opakování

Tandemové opakování

Palindromy

## Srovnávání dvou sekvencí

DP - Needleman-Wunsch

Vylepšení pro maximálně  $k$  chyb

video HHMI

Řetězce a algoritmy na  
řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu  
hledaného motivu

Algoritmus využívající analýzu  
prohledávaného řetězce

Hledání opakování

Tandemové opakování

Palindromy

Srovnávání dvou sekvencí

Vylepšení pro maximálně  $k$  chyb

Příště

video HHMI

# Výpočet omezeného počtu buněk v tabulce DP

## Řetězce a algoritmy na řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu hledaného motivu

Algoritmus využívající analýzu prohledávaného řetězce

## Hledání opakování

Tandemové opakování

Palindromy

## Srovnávání dvou sekvencí

Vylepšení pro maximálně  $k$  chyb

## Příště

video HHMI

Stačí počítat  $2k+1$  diagonál bez ohledu na délku sekvencí

Složitost:  $O(kn)$  (naproti  $O(mn)$ )

# Tabulka pro algoritmus dynamického programování

## Řetězce a algoritmy na řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu hledaného motivu

Algoritmus využívající analýzu prohledávaného řetězce

## Hledání opakování

Tandemové opakování

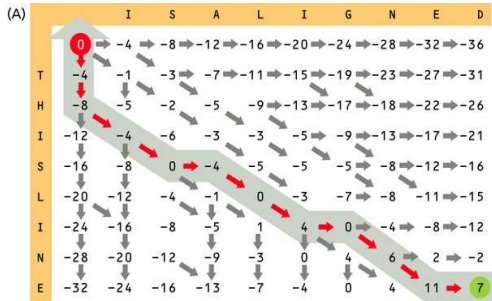
Palindromy

## Srovnávání dvou sekvencí

Vylepšení pro maximálně k chyb

## Příště

video HHMI



(B) THIS-LI-NE-  
--ISALIGNED



# Využití SA a LCP k rychlému postupu po diagonále

## Řetězce a algoritmy na řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu hledaného motivu

Algoritmus využívající analýzu prohledávaného řetězce

## Hledání opakování

Tandemové opakování

Palindromy

## Srovnávání dvou sekvencí

**Vylepšení pro maximálně  $k$  chyb**

## Příště

video HHMI

Složitost:  $O(k^2)$  (viz pou/vzit/i pro palindromy)

## Video HHMI

### Řetězce a algoritmy na řetězcích

Základní pojmy

Základní algoritmy

Algoritmus využívající analýzu hledaného motivu

Algoritmus využívající analýzu prohledávaného řetězce

### Hledání opakování

Tandemové opakování

Palindromy

### Srovnávání dvou sekvencí

Vylepšení pro maximálně  $k$  chyb

### Příště

video HHMI

Dodatek

Dodatek

For Further Reading

# For Further Reading

Dodatek

For Further Reading

X