

# Time-Optimal Construction of String Synchronizing Sets

Jonas Ellert and Tomasz Kociumaka

## Some gems of stringology...

- exact pattern matching...

$$T = \boxed{\overset{1}{\text{s}} \text{ t a c k s a t s t a c s a t t r a c k a} \overset{n}{\text{a}}}$$

$$P = \boxed{\overset{1}{\text{a}} \text{ c } \overset{m}{\text{k}}}$$



## Some gems of stringology...

- exact pattern matching...

$$T = \overset{1}{\text{s}} \overset{3}{\text{t a c k s}} \overset{18}{\text{a t s t a c s a t t r a c k}} \overset{n}{\text{a}}$$

a c k

a c k

$$P = \overset{1}{\text{a}} \overset{m}{\text{c k}}$$

## Some gems of stringology...

- exact pattern matching...

$$T = \overset{1}{\text{s}} \overset{3}{\text{t a c k}} \text{s a t s t a c s a t t r a c k a} \overset{18}{\text{a c k}} \overset{n}{\text{a}}$$

$$P = \overset{1}{\text{a c k}} \overset{m}{\text{k}} \quad \text{occ} = \{3, 18\}$$



## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$T =$ 

1	3																18	$n$		
s	t	a	c	k	s	a	t	s	t	a	c	s	a	t	t	r	a	c	k	a
					a	c	k											a	c	k

$P =$ 

1	$m$	
a	c	k

 $\text{occ} = \{3, 18\}$

- Lempel-Ziv compression...

$T =$ 

1																				$n$
s	t	a	c	k	s	a	t	s	t	a	c	s	a	t	t	r	a	c	k	a

## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$$T = \overset{1}{\text{s}} \overset{3}{\text{t a c k s}} \overset{18}{\text{a t s t a c s a t t r a c k}} \overset{n}{\text{a}}$$

$$P = \overset{1}{\text{a}} \overset{m}{\text{c k}} \quad \text{occ} = \{3, 18\}$$

a c k
a c k

- Lempel-Ziv compression...

$$T = \overset{1}{\text{s}} \overset{n}{\text{t a c k s a t s t a c s a t t r a c k a}}$$

s

## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$T =$ 

1	3															18	$n$			
s	t	a	c	k	s	a	t	s	t	a	c	s	a	t	t	r	a	c	k	a
					a	c	k											a	c	k

$P =$ 

1	$m$	
a	c	k

 $\text{occ} = \{3, 18\}$

- Lempel-Ziv compression...

$T =$ 

1																			$n$	
s	t	a	c	k	s	a	t	s	t	a	c	s	a	t	t	r	a	c	k	a
s	t																			

  
 $s, t$

## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$$T = \overset{1}{\text{s}} \overset{3}{\text{t a c k s}} \overset{18}{\text{a t s t a c s a t t r a c k}} \overset{n}{\text{a}}$$

$$P = \overset{1}{\text{a}} \overset{m}{\text{c k}} \quad \text{occ} = \{3, 18\}$$

- Lempel-Ziv compression...

$$T = \overset{1}{\text{s}} \overset{2}{\text{t}} \overset{3}{\text{a}} \text{c k s a t s t a c s a t t r a c k a} \overset{n}{\text{a}}$$

$$\text{s, t, a}$$



## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$$T = \overset{1}{\text{s}} \overset{3}{\text{t a c k}} \text{s a t s t a c s a t t r a c k a} \overset{18}{\text{a c k}} \overset{n}{\text{a}}$$

$$P = \overset{1}{\text{a}} \overset{m}{\text{c k}} \quad \text{occ} = \{3, 18\}$$

- Lempel-Ziv compression...

$$T = \overset{1}{\text{s}} \overset{2}{\text{t}} \overset{3}{\text{a}} \overset{4}{\text{c}} \overset{5}{\text{k}} \text{s a t s t a c s a t t r a c k a} \overset{n}{\text{a}}$$

s, t, a, c, k



## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$$T = \overset{1}{\text{s}} \overset{3}{\text{t a c k}} \text{s a t s t a c s a t t r a c k a} \overset{18}{\text{a c k}} \overset{n}{\text{a}}$$

$$P = \overset{1}{\text{a}} \overset{m}{\text{c k}} \quad \text{occ} = \{3, 18\}$$

- Lempel-Ziv compression...

$$T = \overset{1}{\text{s}} \overset{2}{\text{t}} \overset{3}{\text{a}} \overset{4}{\text{c}} \overset{5}{\text{k}} \overset{6}{\text{s}} \overset{7}{\text{a}} \text{t s t a c s a t t r a c k a} \overset{n}{\text{a}}$$

s, t, a, c, k, s, a

## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$$T = \overset{1}{\text{s}} \overset{3}{\text{t a c k}} \text{s a t s t a c s a t t r a c k a} \overset{18}{\text{a c k}} \overset{n}{\text{a}}$$

$$P = \overset{1}{\text{a}} \overset{m}{\text{c k}} \quad \text{occ} = \{3, 18\}$$

- Lempel-Ziv compression...

$$T = \overset{1}{\text{s}} \overset{2}{\text{t}} \overset{3}{\text{a}} \overset{4}{\text{c}} \overset{5}{\text{k}} \overset{6}{\text{s}} \overset{7}{\text{a}} \overset{8}{\text{t}} \text{s t a c s a t t r a c k a} \overset{n}{\text{a}}$$

s, t, a, c, k, s, a, t

## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$T =$ 

1	3															18	$n$						
s	t	a	c	k	s	a	t	s	t	a	c	s	a	t	t	r	a	c	k	a			
				a	c	k															a	c	k

$P =$ 

1	$m$	
a	c	k

 $\text{occ} = \{3, 18\}$

- Lempel-Ziv compression...

$T =$ 

1																				$n$															
s	t	a	c	k	s	a	t	s	t	a	c	s	a	t	t	r	a	c	k	a															
s				t				a				c				k				s				a				t				(1, 4)			

## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$T =$ 

1	3															18	$n$						
s	t	a	c	k	s	a	t	s	t	a	c	s	a	t	t	r	a	c	k	a			
				a	c	k															a	c	k

$P =$ 

1	$m$	
a	c	k

 $\text{occ} = \{3, 18\}$

- Lempel-Ziv compression...

$T =$ 

1																				$n$
s	t	a	c	k	s	a	t	s	t	a	c	s	a	t	t	r	a	c	k	a
				s, t, a, c, k, s, a, t,				(1, 4),				(6, 3)								

## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$T =$ 

1	3															18	$n$						
s	t	a	c	k	s	a	t	s	t	a	c	s	a	t	t	r	a	c	k	a			
				a	c	k															a	c	k

$P =$ 

1	$m$	
a	c	k

 $\text{occ} = \{3, 18\}$

- Lempel-Ziv compression...

$T =$ 

1																				$n$
s	t	a	c	k	s	a	t	s	t	a	c	s	a	t	t	r	a	c	k	a

  
 $s, t, a, c, k, s, a, t, (1, 4), (6, 3), t$

## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$$T = \overset{1}{\text{s}} \overset{3}{\text{t a c k}} \text{s a t s t a c s a t t r a c k a} \overset{18}{\text{a c k}} \overset{n}{\text{a}}$$

$$P = \overset{1}{\text{a}} \overset{m}{\text{c k}} \quad \text{occ} = \{3, 18\}$$

- Lempel-Ziv compression...

$$T = \overset{1}{\text{s}} \text{t a c k s a t s t a c s a t t r a c k a} \overset{n}{\text{a}}$$

s, t, a, c, k, s, a, t, (1, 4), (6, 3), t, r

## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$T =$ 

1	3															18	$n$						
s	t	a	c	k	s	a	t	s	t	a	c	s	a	t	t	r	a	c	k	a			
				a c k																	a c k		

$P =$ 

1	$m$	
a	c	k

 $\text{occ} = \{3, 18\}$

- Lempel-Ziv compression...

$T =$ 

1																				$n$
s	t	a	c	k	s	a	t	s	t	a	c	s	a	t	t	r	a	c	k	a
s, t, a, c, k, s, a, t, (1, 4), (6, 3), t, r, (3, 3)																				

## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$$T = \overset{1}{\text{s}} \overset{3}{\text{t a c k}} \text{s a t s t a c s a t t r a c k a} \overset{18}{\text{a c k}} \overset{n}{\text{a}}$$

$$P = \overset{1}{\text{a}} \overset{m}{\text{c k}} \quad \text{occ} = \{3, 18\}$$

- Lempel-Ziv compression...

$$T = \overset{1}{\text{s}} \overset{4}{\text{t a c k}} \overset{5}{\text{s a t}} \overset{11}{\text{s t a c}} \overset{17}{\text{s a t}} \overset{20}{\text{t r a c k}} \overset{24}{\text{a}}$$

s, t, a, c, k, s, a, t, (1, 4), (6, 3), t, r, (3, 3), a

## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$$T = \overset{1}{\text{s}} \overset{3}{\text{t a c k}} \text{s a t s t a c s a t t r a c k a} \overset{18}{\text{a c k}} \overset{n}{\text{a}} \quad P = \overset{1}{\text{a}} \overset{m}{\text{c k}} \quad \text{occ} = \{3, 18\}$$

- Lempel-Ziv compression...  $\mathcal{O}(n)$  time [Ziv+Lempel '77; Storer+Szymanski '82]

$$T = \overset{1}{\text{s}} \overset{4}{\text{t a c k}} \overset{8}{\text{s a t}} \overset{12}{\text{s t a c}} \overset{16}{\text{s a t}} \overset{20}{\text{t r a c}} \overset{24}{\text{k a}}$$

s, t, a, c, k, s, a, t, (1, 4), (6, 3), t, r, (3, 3), a

## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$$T = \overset{1}{\text{s}} \overset{3}{\text{t a c k s}} \text{ a t s t a c s a t t r a c k a} \overset{18}{\text{a c k}} \overset{n}{\text{a}}$$

$$P = \overset{1}{\text{a}} \overset{m}{\text{c k}} \quad \text{occ} = \{3, 18\}$$

- Lempel-Ziv compression...  $\mathcal{O}(n)$  time [Ziv+Lempel '77; Storer+Szymanski '82]

$$T = \overset{1}{\text{s}} \overset{n}{\text{t a c k s a t s t a c s a t t r a c k a}}$$

s, t, a, c, k, s, a, t, (1, 4), (6, 3), t, r, (3, 3), a

- Longest Common Extension (LCE) data structure...

- query in  $\mathcal{O}(1)$  time: given  $i, j$ , return maximal  $\ell$  s.t.  $T[i..i+\ell] = T[j..j+\ell]$

$$T = \overset{1}{\text{s}} \overset{n}{\text{t a c k s a t s t a c s a t t r a c k a}}$$



## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$$T = \overset{1}{\text{s}} \overset{3}{\text{t a c k s}} \overset{18}{\text{a t s t a c s a t t r a c k}} \overset{n}{\text{a}}$$

$$P = \overset{1}{\text{a}} \overset{m}{\text{c k}} \quad \text{occ} = \{3, 18\}$$

a c k
a c k

- Lempel-Ziv compression...  $\mathcal{O}(n)$  time [Ziv+Lempel '77; Storer+Szymanski '82]

$$T = \overset{1}{\text{s}} \overset{n}{\text{t a c k s a t s t a c s a t t r a c k a}}$$

s, t, a, c, k, s, a, t, (1, 4), (6, 3), t, r, (3, 3), a

- Longest Common Extension (LCE) data structure...

- query in  $\mathcal{O}(1)$  time: given  $i, j$ , return maximal  $\ell$  s.t.  $T[i..i+\ell] = T[j..j+\ell]$

$$T = \overset{1}{\text{s}} \overset{6}{\text{t a c k s a t}} \overset{13}{\text{s t a c s a t t r a c k}} \overset{n}{\text{a}}$$

s a t
s a t

## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$$T = \overset{1}{\text{s}} \overset{3}{\text{t a c k s}} \overset{18}{\text{a t s t a c s a t t r a c k}} \overset{n}{\text{a}}$$

a c k
a c k

$$P = \overset{1}{\text{a}} \overset{m}{\text{c k}} \quad \text{occ} = \{3, 18\}$$

- Lempel-Ziv compression...  $\mathcal{O}(n)$  time [Ziv+Lempel '77; Storer+Szymanski '82]

$$T = \overset{1}{\text{s}} \overset{n}{\text{t a c k s a t s t a c s a t t r a c k a}}$$

s, t, a, c, k, s, a, t, (1, 4), (6, 3), t, r, (3, 3), a

- Longest Common Extension (LCE) data structure...

- query in  $\mathcal{O}(1)$  time: given  $i, j$ , return maximal  $\ell$  s.t.  $T[i..i+\ell] = T[j..j+\ell]$

$$T = \overset{1}{\text{s}} \overset{6}{\text{t a c k s a t s}} \overset{13}{\text{t a c s a t t r a c k}} \overset{n}{\text{a}}$$

s a t s
s a t t

## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$$T = \overset{1}{\text{s}} \overset{3}{\text{t a c k s}} \text{ a t s t a c s a t t r a c k a} \overset{18}{\text{a c k}} \overset{n}{\text{a}}$$

$$P = \overset{1}{\text{a}} \overset{m}{\text{c k}} \quad \text{occ} = \{3, 18\}$$

- Lempel-Ziv compression...  $\mathcal{O}(n)$  time [Ziv+Lempel '77; Storer+Szymanski '82]

$$T = \overset{1}{\text{s}} \overset{4}{\text{t a c k s}} \overset{8}{\text{a t s t a c s}} \overset{12}{\text{a t t r a c k}} \overset{n}{\text{a}}$$

s, t, a, c, k, s, a, t, (1, 4), (6, 3), t, r, (3, 3), a

- Longest Common Extension (LCE) data structure...

- query in  $\mathcal{O}(1)$  time: given  $i, j$ , return maximal  $\ell$  s.t.  $T[i..i+\ell] = T[j..j+\ell]$

$$T = \overset{1}{\text{s}} \overset{6}{\text{t a c k s a t s}} \overset{13}{\text{t a c s a t t r a c k}} \overset{n}{\text{a}}$$

$$\text{LCE}(6, 13) = 3$$

## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n)$  time [Knuth et al. '77]

$$T = \overset{1}{\text{s}} \overset{3}{\text{t a c k s}} \text{ a t s t a c s a t t r a c k a} \overset{18}{\text{a c k}} \overset{n}{\text{a}}$$

$$P = \overset{1}{\text{a}} \overset{m}{\text{c k}} \quad \text{occ} = \{3, 18\}$$

- Lempel-Ziv compression...  $\mathcal{O}(n)$  time [Ziv+Lempel '77; Storer+Szymanski '82]

$$T = \overset{1}{\text{s}} \overset{4}{\text{t a c k s}} \overset{8}{\text{a t s t a c s}} \overset{12}{\text{a t t r a c k}} \overset{n}{\text{a}}$$

s, t, a, c, k, s, a, t, (1, 4), (6, 3), t, r, (3, 3), a

- Longest Common Extension (LCE) data structure...  $\mathcal{O}(n)$  construction time [Farach '97]

- query in  $\mathcal{O}(1)$  time: given  $i, j$ , return maximal  $\ell$  s.t.  $T[i..i+\ell] = T[j..j+\ell]$

$$T = \overset{1}{\text{s}} \overset{6}{\text{t a c k s a t s}} \overset{13}{\text{t a c s a t t r a c k}} \overset{n}{\text{a}}$$

s a t s
s a t t

$$\text{LCE}(6, 13) = 3$$

## Some gems of stringology...

- often  $\mathcal{O}(n)$  time solutions on input of length  $n$
- ... surely that must be optimal, right?

## Some gems of stringology...

- often  $\mathcal{O}(n)$  time solutions on input of length  $n$
- ... surely that must be optimal, right?
  
- **most algorithms in word RAM model** [Hagerup '98]

## Some gems of stringology...

- often  $\mathcal{O}(n)$  time solutions on input of length  $n$
- ... surely that must be optimal, right?
  
- **most algorithms in word RAM model** [Hagerup '98]
- memory words and CPU registers consist of  $w$  bits
- representing integers in  $[0 . . 2^w)$

## Some gems of stringology...

- often  $\mathcal{O}(n)$  time solutions on input of length  $n$
- ... surely that must be optimal, right?
  
- **most algorithms in word RAM model** [Hagerup '98]
- memory words and CPU registers consist of  $w$  bits
- representing integers in  $[0..2^w)$

### Operations:

arithmetic (+ - × /)

bit-wise logic (NOT,AND,OR,XOR)

bit-shifting ( $\ll$ ,  $\gg$ )

## Some gems of stringology...

- often  $\mathcal{O}(n)$  time solutions on input of length  $n$
- ... surely that must be optimal, right?
  
- **most algorithms in word RAM model** [Hagerup '98]
- memory words and CPU registers consist of  $w$  bits
- representing integers in  $[0..2^w)$
- assumption:  $w = \Omega(\log n)$ , usually  $w = \Theta(\log n)$

### Operations:

arithmetic (+ - × /)

bit-wise logic (NOT,AND,OR,XOR)

bit-shifting ( $\ll$ ,  $\gg$ )

## Some gems of stringology...

- often  $\mathcal{O}(n)$  time solutions on input of length  $n$
- ... surely that must be optimal, right?
- **most algorithms in word RAM model** [Hagerup '98]
- memory words and CPU registers consist of  $w$  bits
- representing integers in  $[0..2^w)$
- assumption:  $w = \Omega(\log n)$ , usually  $w = \Theta(\log n)$

### Operations:

arithmetic (+ - × /)

bit-wise logic (NOT,AND,OR,XOR)

bit-shifting ( $\ll$ ,  $\gg$ )

year	data pp (MB)	$w$	$w$ (SIMD)
1986	539	32	–
1993	2,866	32	–
2000	8,988	32–64	128
2007	44,716	64	128

## Some gems of stringology...

- often  $\mathcal{O}(n)$  time solutions on input of length  $n$
- ... surely that must be optimal, right?
- **most algorithms in word RAM model** [Hagerup '98]
- memory words and CPU registers consist of  $w$  bits
- representing integers in  $[0..2^w)$
- assumption:  $w = \Omega(\log n)$ , usually  $w = \Theta(\log n)$

### Operations:

arithmetic (+ - × /)

bit-wise logic (NOT,AND,OR,XOR)

bit-shifting ( $\ll$ ,  $\gg$ )

year	data pp (MB)	$w$	$w$ (SIMD)
1986	539	32	–
1993	2,866	32	–
2000	8,988	32–64	128
2007	44,716	64	128

Implications for, e.g., a binary string:

$$T = \overset{1}{\boxed{1}} \overset{w}{\boxed{0111000}} \boxed{00010101} \boxed{01000110} \boxed{10111000} \boxed{00011010} \boxed{11100111} \boxed{10111010} \overset{n}{\boxed{0}}$$

## Some gems of stringology...

- often  $\mathcal{O}(n)$  time solutions on input of length  $n$
- ... surely that must be optimal, right?
- **most algorithms in word RAM model** [Hagerup '98]
- memory words and CPU registers consist of  $w$  bits
- representing integers in  $[0..2^w)$
- assumption:  $w = \Omega(\log n)$ , usually  $w = \Theta(\log n)$

### Operations:

arithmetic (+ - × /)

bit-wise logic (NOT,AND,OR,XOR)

bit-shifting ( $\ll$ ,  $\gg$ )

year	data pp (MB)	$w$	$w$ (SIMD)
1986	539	32	–
1993	2,866	32	–
2000	8,988	32–64	128
2007	44,716	64	128

Implications for, e.g., a binary string:

$$T = \overset{1}{\boxed{10111000}} \overset{w}{\boxed{00010101}} \boxed{01000110} \boxed{10111000} \boxed{00011010} \boxed{11100111} \boxed{10111010} \overset{n}{}$$

... only  $\approx n/w = \Theta(n/\log n)$  words. Why not process in  $\mathcal{O}(n/\log n)$  time?

## Some gems of stringology...

- often  $\mathcal{O}(n)$  time solutions on input of length  $n$
- ... surely that must be optimal, right?
- most algorithms in word RAM model [Hagerup '98]**
- memory words and CPU registers consist of  $w$  bits
- representing integers in  $[0..2^w)$
- assumption:  $w = \Omega(\log n)$ , usually  $w = \Theta(\log n)$

### Operations:

arithmetic (+ - × /)

bit-wise logic (NOT,AND,OR,XOR)

bit-shifting ( $\ll$ ,  $\gg$ )

year	data pp (MB)	$w$	$w$ (SIMD)
1986	539	32	–
1993	2,866	32	–
2000	8,988	32–64	128
2007	44,716	64	128

More generally, consider  $T \in [0.. \sigma)^n$

$$T = \overset{1}{\boxed{10111000}} \overset{w}{\boxed{00010101}} \overset{w}{\boxed{01000110}} \overset{w}{\boxed{10111000}} \overset{w}{\boxed{00011010}} \overset{n \cdot \lceil \log_2 \sigma \rceil}{\boxed{11100111}}$$

## Some gems of stringology...

- often  $\mathcal{O}(n)$  time solutions on input of length  $n$
- ... surely that must be optimal, right?
- most algorithms in word RAM model** [Hagerup '98]
- memory words and CPU registers consist of  $w$  bits
- representing integers in  $[0..2^w)$
- assumption:  $w = \Omega(\log n)$ , usually  $w = \Theta(\log n)$

### Operations:

arithmetic (+ - × /)

bit-wise logic (NOT,AND,OR,XOR)

bit-shifting ( $\ll$ ,  $\gg$ )

year	data pp (MB)	$w$	$w$ (SIMD)
1986	539	32	–
1993	2,866	32	–
2000	8,988	32–64	128
2007	44,716	64	128

More generally, consider  $T \in [0.. \sigma]^n$

$T =$ 

1	$w$	$n \cdot \lceil \log_2 \sigma \rceil$
1 0 1 1 1 0 0 0	0 0 0 1 0 1 0 1	0 1 0 0 0 1 1 0
1 0 1 1 1 0 0 0	0 0 0 1 1 0 1 0	1 1 1 0 0 1 1 1
s	t	a
c	k	s
a	t	s
t	a	c
s	n	o
w		

$s$	$\text{bin}(s)$
a	000
c	001
k	010
n	011
o	100
s	101
t	110
w	111

## Some gems of stringology...

- often  $\mathcal{O}(n)$  time solutions on input of length  $n$
- ... surely that must be optimal, right?
- most algorithms in word RAM model** [Hagerup '98]
- memory words and CPU registers consist of  $w$  bits
- representing integers in  $[0..2^w)$
- assumption:  $w = \Omega(\log n)$ , usually  $w = \Theta(\log n)$

### Operations:

arithmetic (+ - × /)

bit-wise logic (NOT,AND,OR,XOR)

bit-shifting ( $\ll$ ,  $\gg$ )

year	data pp (MB)	$w$	$w$ (SIMD)
1986	539	32	–
1993	2,866	32	–
2000	8,988	32–64	128
2007	44,716	64	128

More generally, consider  $T \in [0.. \sigma]^n$

$T =$ 

1	$w$	$n \cdot \lceil \log_2 \sigma \rceil$
1 0 1 1 1 0 0 0	0 0 0 1 0 1 0 1	0 1 0 0 0 1 1 0
1 0 1 1 1 0 0 0	0 0 0 1 1 0 1 0	1 1 1 0 0 1 1 1
s	t	a
c	k	s
a	t	s
t	a	c
s	n	o
w		

...  $\approx n \cdot \lceil \log_2 \sigma \rceil / w$  words. **Truly linear time:**  $\mathcal{O}(n \log \sigma / \log n) = \mathcal{O}(n / \log_\sigma n)$

$s$	$\text{bin}(s)$
a	000
c	001
k	010
n	011
o	100
s	101
t	110
w	111







## Some gems of stringology...

- exact pattern matching...  $\mathcal{O}(n/\log_\sigma n)$  time [Bille '11]

$$T = \overset{1}{\text{s}} \overset{3}{\text{t a c k}} \text{s a t s t a c s a t t r a c k a} \overset{18}{\text{a c k}} \overset{n}{\text{a}} \quad P = \overset{1}{\text{a}} \overset{m}{\text{c k}} \quad \text{occ} = \{3, 18\}$$

- Lempel-Ziv compression...  $\mathcal{O}(n/\log_\sigma n)$  or  $\mathcal{O}(n \log \sigma / \sqrt{\log n})$  time [Ellert '23; Kempa+Kociumaka '24]

$$T = \overset{1}{\text{s}} \overset{n}{\text{t a c k s a t s t a c s a t t r a c k a}}$$

s, t, a, c, k, s, a, t, (1, 4), (6, 3), t, r, (3, 3), a

- Longest Common Extension (LCE) data structure...  $\mathcal{O}(n/\log_\sigma n)$  time [Kempa+Kociumaka '19]

- query in  $\mathcal{O}(1)$  time: given  $i, j$ , return maximal  $\ell$  s.t.  $T[i..i+\ell] = T[j..j+\ell]$

$$T = \overset{1}{\text{s}} \overset{6}{\text{t a c k s a t s}} \overset{13}{\text{t a c s a t t r a c k a}} \overset{n}{\text{a}}$$

s a t s      s a t t

LCE(6, 13) = 3

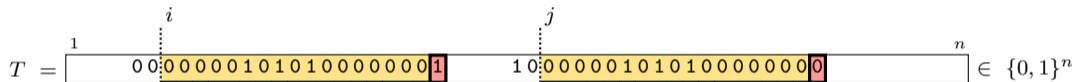
## String synchronizing sets

- would like a powerful tool for designing truly linear time algorithms
- should consider words not symbols or bits



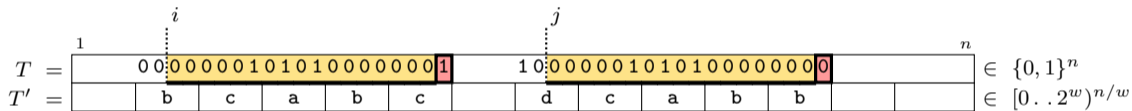
## String synchronizing sets

- would like a powerful tool for designing truly linear time algorithms
- should consider words not symbols or bits
- example for LCE data structure (for binary alphabet)



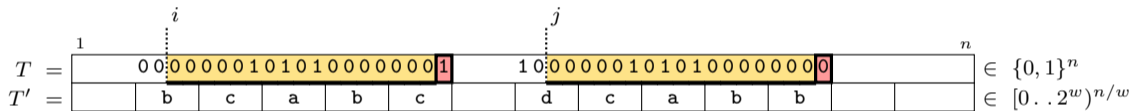
## String synchronizing sets

- would like a powerful tool for designing truly linear time algorithms
- should consider words not symbols or bits
- example for LCE data structure (for binary alphabet)



## String synchronizing sets

- would like a powerful tool for designing truly linear time algorithms
- should consider words not symbols or bits
- example for LCE data structure (for binary alphabet)



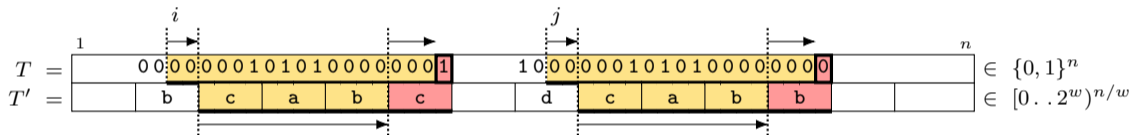
- compute LCE data structure for  $T'$  in  $\mathcal{O}(n/w) = \mathcal{O}(n/\log n)$  time





## String synchronizing sets

- would like a powerful tool for designing truly linear time algorithms
- should consider words not symbols or bits
- example for LCE data structure (for binary alphabet)



- compute LCE data structure for  $T'$  in  $\mathcal{O}(n/w) = \mathcal{O}(n/\log n)$  time



## String synchronizing sets

- fix parameter  $\tau$ , typically  $\tau = \Theta(w)$
- synchronizing function  $\text{sync} : \{0, 1\}^{2\tau} \rightarrow \{\text{true}, \text{false}\}$

## String synchronizing sets

- fix parameter  $\tau$ , typically  $\tau = \Theta(w)$
- synchronizing function  $\text{sync} : \{0, 1\}^{2\tau} \rightarrow \{\text{true}, \text{false}\}$
- position  $i$  is synchronizing if  $\text{sync}(T[i..i + 2\tau])$



## String synchronizing sets

- fix parameter  $\tau$ , typically  $\tau = \Theta(w)$
- synchronizing function  $\text{sync} : \{0, 1\}^{2\tau} \rightarrow \{\text{true}, \text{false}\}$
- position  $i$  is synchronizing if  $\text{sync}(T[i..i+2\tau])$

$$T = \boxed{\begin{array}{c} 1 \quad \text{-----} \quad n \\ x \\ \hline T' = \boxed{a} \end{array}} \in \{0, 1\}^n$$





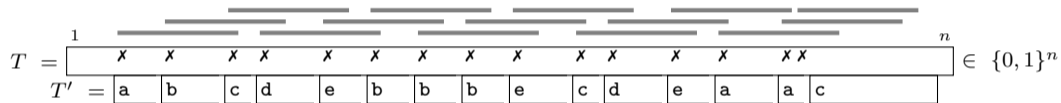
## String synchronizing sets

- fix parameter  $\tau$ , typically  $\tau = \Theta(w)$
- synchronizing function  $\text{sync} : \{0, 1\}^{2\tau} \rightarrow \{\text{true}, \text{false}\}$
- position  $i$  is synchronizing if  $\text{sync}(T[i..i+2\tau])$



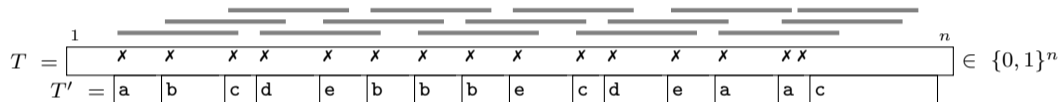
## String synchronizing sets

- fix parameter  $\tau$ , typically  $\tau = \Theta(w)$
- synchronizing function  $\text{sync} : \{0, 1\}^{2\tau} \rightarrow \{\text{true}, \text{false}\}$
- position  $i$  is synchronizing if  $\text{sync}(T[i..i+2\tau])$



## String synchronizing sets

- fix parameter  $\tau$ , typically  $\tau = \Theta(w)$
- synchronizing function  $\text{sync} : \{0, 1\}^{2\tau} \rightarrow \{\text{true}, \text{false}\}$
- position  $i$  is synchronizing if  $\text{sync}(T[i..i+2\tau])$

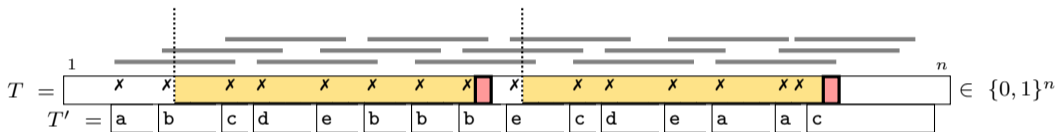


### $\tau$ -synchronizing set of $T$ :

- only  $\mathcal{O}(n/\tau)$  synchronizing positions, so  $|T'| = \mathcal{O}(n/\tau) = \mathcal{O}(n/\log n)$
- all of  $T$  overlapped by synchronizing fragments

## String synchronizing sets

- fix parameter  $\tau$ , typically  $\tau = \Theta(w)$
- synchronizing function  $\text{sync} : \{0, 1\}^{2\tau} \rightarrow \{\text{true}, \text{false}\}$
- position  $i$  is synchronizing if  $\text{sync}(T[i..i+2\tau])$

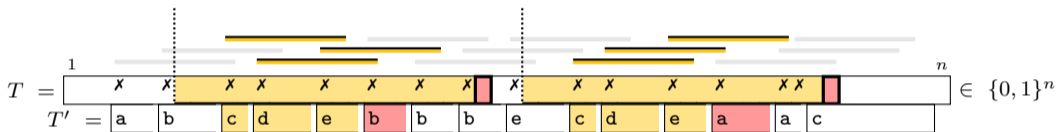


### $\tau$ -synchronizing set of $T$ :

- only  $\mathcal{O}(n/\tau)$  synchronizing positions, so  $|T'| = \mathcal{O}(n/\tau) = \mathcal{O}(n/\log n)$
- all of  $T$  overlapped by synchronizing fragments

## String synchronizing sets

- fix parameter  $\tau$ , typically  $\tau = \Theta(w)$
- synchronizing function  $\text{sync} : \{0, 1\}^{2\tau} \rightarrow \{\text{true}, \text{false}\}$
- position  $i$  is synchronizing if  $\text{sync}(T[i..i+2\tau])$

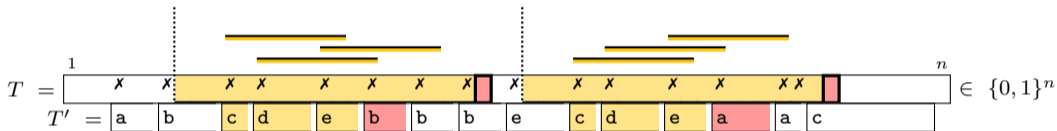


### $\tau$ -synchronizing set of $T$ :

- only  $\mathcal{O}(n/\tau)$  synchronizing positions, so  $|T'| = \mathcal{O}(n/\tau) = \mathcal{O}(n/\log n)$
- all of  $T$  overlapped by synchronizing fragments

## String synchronizing sets

- fix parameter  $\tau$ , typically  $\tau = \Theta(w)$
- synchronizing function  $\text{sync} : \{0, 1\}^{2\tau} \rightarrow \{\text{true}, \text{false}\}$
- position  $i$  is synchronizing if  $\text{sync}(T[i..i+2\tau])$

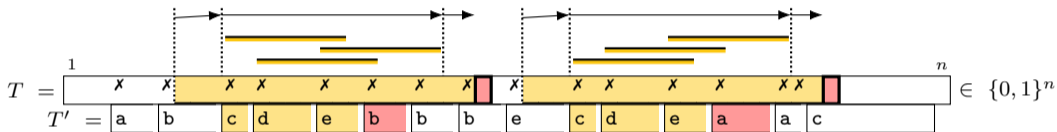


### $\tau$ -synchronizing set of $T$ :

- only  $\mathcal{O}(n/\tau)$  synchronizing positions, so  $|T'| = \mathcal{O}(n/\tau) = \mathcal{O}(n/\log n)$
- all of  $T$  overlapped by synchronizing fragments

## String synchronizing sets

- fix parameter  $\tau$ , typically  $\tau = \Theta(w)$
- synchronizing function  $\text{sync} : \{0, 1\}^{2\tau} \rightarrow \{\text{true}, \text{false}\}$
- position  $i$  is synchronizing if  $\text{sync}(T[i..i+2\tau])$

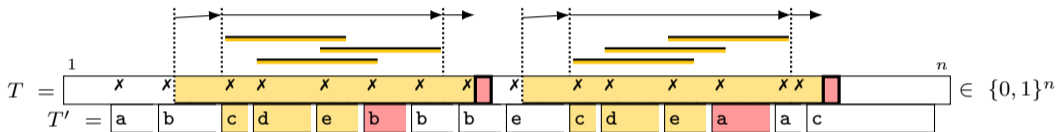


### $\tau$ -synchronizing set of $T$ :

- only  $\mathcal{O}(n/\tau)$  synchronizing positions, so  $|T'| = \mathcal{O}(n/\tau) = \mathcal{O}(n/\log n)$
- all of  $T$  overlapped by synchronizing fragments

## String synchronizing sets

- fix parameter  $\tau$ , typically  $\tau = \Theta(w)$
- synchronizing function  $\text{sync} : \{0, 1\}^{2\tau} \rightarrow \{\text{true}, \text{false}\}$
- position  $i$  is synchronizing if  $\text{sync}(T[i..i+2\tau])$

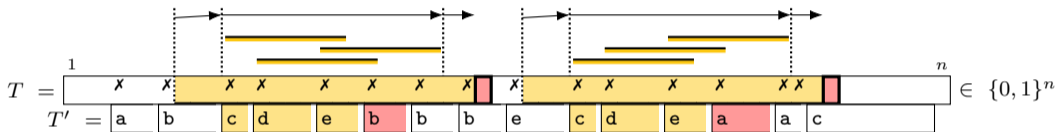


### $\tau$ -synchronizing set of $T$ :

- only  $\mathcal{O}(n/\tau)$  synchronizing positions, so  $|T'| = \mathcal{O}(n/\tau) = \mathcal{O}(n/\log n)$
- all of  $T$  overlapped by synchronizing fragments

## String synchronizing sets

- fix parameter  $\tau$ , typically  $\tau = \Theta(w)$
- synchronizing function  $\text{sync} : \{0, 1\}^{2\tau} \rightarrow \{\text{true}, \text{false}\}$
- position  $i$  is synchronizing if  $\text{sync}(T[i..i+2\tau])$



### $\tau$ -synchronizing set of $T$ :

- only  $\mathcal{O}(n/\tau)$  synchronizing positions, so  $|T'| = \mathcal{O}(n/\tau) = \mathcal{O}(n/\log n)$
- all of  $T$  overlapped by synchronizing fragments
- some extra conditions for periodic fragments apply

## String synchronizing sets

	preprocessing (independent of $\tau$ )	construction (depends on $\tau$ )	
[Kempa+Kociumaka '19]	none	$\mathcal{O}(n/\log_\sigma n)$	only for $\tau < \frac{1}{5} \log_\sigma n$

## String synchronizing sets

	preprocessing (independent of $\tau$ )	construction (depends on $\tau$ )	
[Kempa+Kociumaka '19] [Kociumaka et al. '24]	none $\mathcal{O}(n)$	$\mathcal{O}(n/\log_\sigma n)$ $\mathcal{O}(n/\tau)$	only for $\tau < \frac{1}{5} \log_\sigma n$

## String synchronizing sets

	preprocessing (independent of $\tau$ )	construction (depends on $\tau$ )	
[Kempa+Kociumaka '19]	none	$\mathcal{O}(n/\log_\sigma n)$	only for $\tau < \frac{1}{5} \log_\sigma n$
[Kociumaka et al. '24]	$\mathcal{O}(n)$	$\mathcal{O}(n/\tau)$	
<b>this paper</b>	$\mathcal{O}(n/\log_\sigma n)$	$\mathcal{O}(n/\tau)$	simple representation

## String synchronizing sets

	preprocessing (independent of $\tau$ )	construction (depends on $\tau$ )	
[Kempa+Kociumaka '19]	none	$\mathcal{O}(n/\log_\sigma n)$	only for $\tau < \frac{1}{5} \log_\sigma n$
[Kociumaka et al. '24]	$\mathcal{O}(n)$	$\mathcal{O}(n/\tau)$	
<b>this paper</b>	$\mathcal{O}(n/\log_\sigma n)$ $\mathcal{O}(n/\log_\sigma n)$	$\mathcal{O}(n/\tau)$ $\mathcal{O}(n \log \tau / (\tau \log n))$	simple representation compressed representation

## String synchronizing sets

	preprocessing (independent of $\tau$ )	construction (depends on $\tau$ )	
[Kempa+Kociumaka '19]	none	$\mathcal{O}(n/\log_\sigma n)$	only for $\tau < \frac{1}{5} \log_\sigma n$
[Kociumaka et al. '24]	$\mathcal{O}(n)$	$\mathcal{O}(n/\tau)$	
<b>this paper</b>	$\mathcal{O}(n/\log_\sigma n)$ $\mathcal{O}(n/\log_\sigma n)$	$\mathcal{O}(n/\tau)$ $\mathcal{O}(n \log \tau / (\tau \log n))$	simple representation compressed representation

- most applications use  $\tau = \Theta(\log_\sigma n)$
- some applications can benefit from arbitrary  $\tau$

## String synchronizing sets

	preprocessing (independent of $\tau$ )	construction (depends on $\tau$ )	
[Kempa+Kociumaka '19]	none	$\mathcal{O}(n/\log_\sigma n)$	only for $\tau < \frac{1}{5} \log_\sigma n$
[Kociumaka et al. '24]	$\mathcal{O}(n)$	$\mathcal{O}(n/\tau)$	
<b>this paper</b>	$\mathcal{O}(n/\log_\sigma n)$ $\mathcal{O}(n/\log_\sigma n)$	$\mathcal{O}(n/\tau)$ $\mathcal{O}(n \log \tau / (\tau \log n))$	simple representation compressed representation

- most applications use  $\tau = \Theta(\log_\sigma n)$
- some applications can benefit from arbitrary  $\tau$

### Example: Space-time trade-off for LCEs:

- LCE reduces to scanning  $\mathcal{O}(\tau)$  symbols and computing one LCE in  $T'$

## String synchronizing sets

	preprocessing (independent of $\tau$ )	construction (depends on $\tau$ )	
[Kempa+Kociumaka '19] [Kociumaka et al. '24]	none $\mathcal{O}(n)$	$\mathcal{O}(n/\log_\sigma n)$ $\mathcal{O}(n/\tau)$	only for $\tau < \frac{1}{5} \log_\sigma n$
<b>this paper</b>	$\mathcal{O}(n/\log_\sigma n)$ $\mathcal{O}(n/\log_\sigma n)$	$\mathcal{O}(n/\tau)$ $\mathcal{O}(n \log \tau / (\tau \log n))$	simple representation compressed representation

- most applications use  $\tau = \Theta(\log_\sigma n)$
- some applications can benefit from arbitrary  $\tau$

### Example: Space-time trade-off for LCEs:

- LCE reduces to scanning  $\mathcal{O}(\tau)$  symbols and computing one LCE in  $T'$
- scan takes  $\mathcal{O}(1 + \tau/\log_\sigma n)$  time in a word-wise manner

## String synchronizing sets

	preprocessing (independent of $\tau$ )	construction (depends on $\tau$ )	
[Kempa+Kociumaka '19] [Kociumaka et al. '24]	none $\mathcal{O}(n)$	$\mathcal{O}(n/\log_\sigma n)$ $\mathcal{O}(n/\tau)$	only for $\tau < \frac{1}{5} \log_\sigma n$
<b>this paper</b>	$\mathcal{O}(n/\log_\sigma n)$ $\mathcal{O}(n/\log_\sigma n)$	$\mathcal{O}(n/\tau)$ $\mathcal{O}(n \log \tau / (\tau \log n))$	simple representation compressed representation

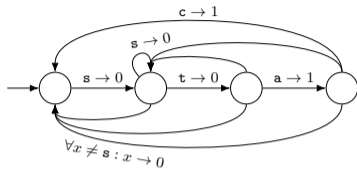
- most applications use  $\tau = \Theta(\log_\sigma n)$
- some applications can benefit from arbitrary  $\tau$

### Example: Space-time trade-off for LCEs:

- LCE reduces to scanning  $\mathcal{O}(\tau)$  symbols and computing one LCE in  $T'$
- scan takes  $\mathcal{O}(1 + \tau/\log_\sigma n)$  time in a word-wise manner
- LCE data structure of  $T'$  is of size  $\mathcal{O}(n/\tau)$  words

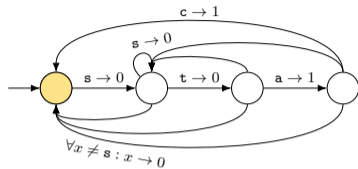
## Some new tools we developed

$T =$  s t a r t s t a c k s a t s t a c s



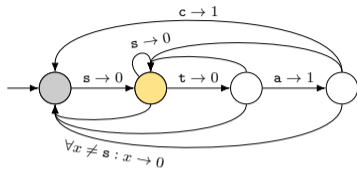
## Some new tools we developed

$T =$  s t a r t s t a c k s a t s t a c s



## Some new tools we developed

$T =$  **s** t a r t s t a c k s a t s t a c s  
 0

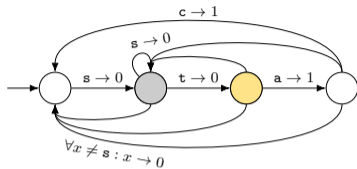


## Some new tools we developed

$T =$ 

s	t	a	r	t	s	t	a	c	k	s	a	t	s	t	a	c	s
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

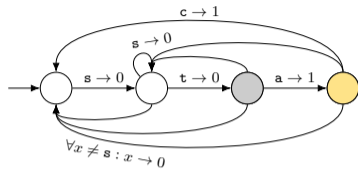
  
 0 0



## Some new tools we developed

$T =$ 

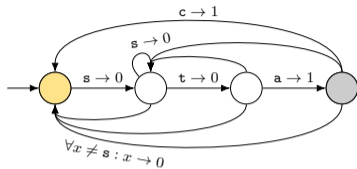
s	t	a	r	t	s	t	a	c	k	s	a	t	s	t	a	c	s
0	0	0															



## Some new tools we developed

$T =$ 

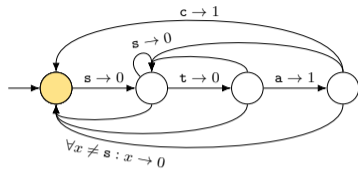
s	t	a	r	t	s	t	a	c	k	s	a	t	s	t	a	c	s
0	0	0	0														



## Some new tools we developed

$T =$ 

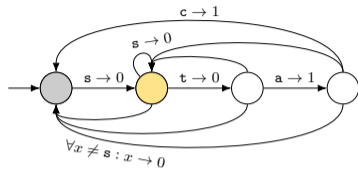
s	t	a	r	t	s	t	a	c	k	s	a	t	s	t	a	c	s
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0



## Some new tools we developed

$T =$ 

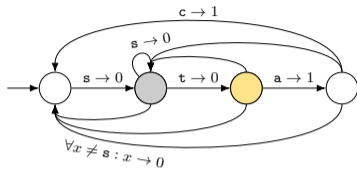
s	t	a	r	t	s	t	a	c	k	s	a	t	s	t	a	c	s
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0



## Some new tools we developed

$T =$ 

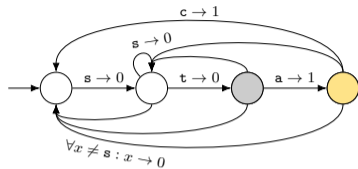
s	t	a	r	t	s	t	a	c	k	s	a	t	s	t	a	c	s
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0



## Some new tools we developed

$T =$ 

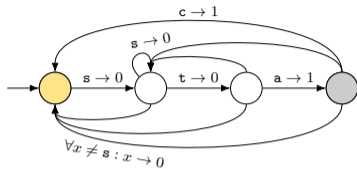
s	t	a	r	t	s	t	a	c	k	s	a	t	s	t	a	c	s
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0



## Some new tools we developed

$T =$ 

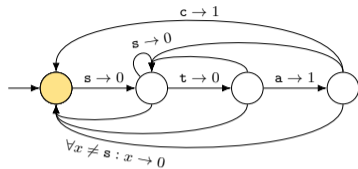
s	t	a	r	t	s	t	a	c	k	s	a	t	s	t	a	c	s
0	0	0	0	0	0	0	0	0	1								



## Some new tools we developed

$T =$ 

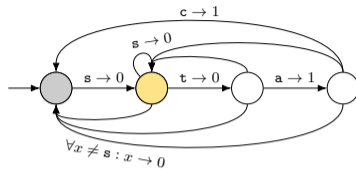
s	t	a	r	t	s	t	a	c	k	s	a	t	s	t	a	c	s
0	0	0	0	0	0	0	0	0	1	0							



## Some new tools we developed

$T =$ 

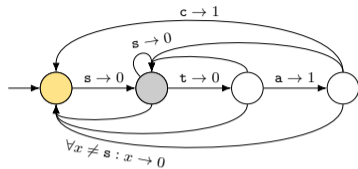
s	t	a	r	t	s	t	a	c	k	s	a	t	s	t	a	c	s
0	0	0	0	0	0	0	0	0	0	1	0	0					



## Some new tools we developed

$T =$ 

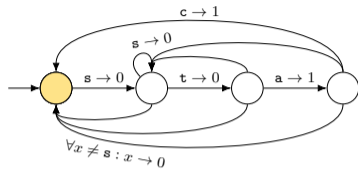
s	t	a	r	t	s	t	a	c	k	s	a	t	s	t	a	c	s
0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0



## Some new tools we developed

$T =$ 

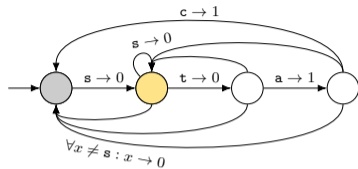
s	t	a	r	t	s	t	a	c	k	s	a	t	s	t	a	c	s
0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0



## Some new tools we developed

$T =$ 

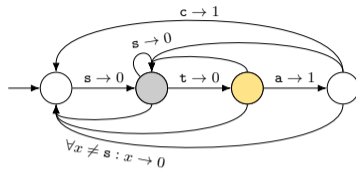
s	t	a	r	t	s	t	a	c	k	s	a	t	s	t	a	c	s				
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0



## Some new tools we developed

$T =$ 

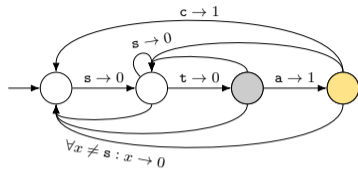
s	t	a	r	t	s	t	a	c	k	s	a	t	s	<b>t</b>	a	c	s				
0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0



## Some new tools we developed

$T =$ 

s	t	a	r	t	s	t	a	c	k	s	a	t	s	t	a	c	s								
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0

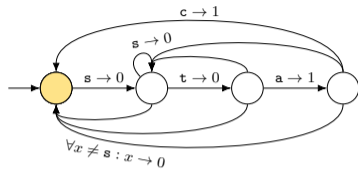




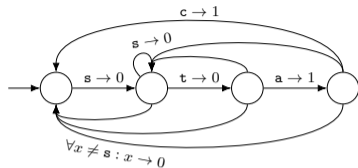
## Some new tools we developed

$T =$ 

s	t	a	r	t	s	t	a	c	k	s	a	t	s	t	a	c	s	
0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	1	0

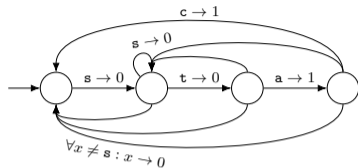


## Some new tools we developed

$$T = \begin{array}{|cccccccccccccccc} \hline \text{s} & \text{t} & \text{a} & \text{r} & \text{t} & \text{s} & \text{t} & \text{a} & \text{k} & \text{s} & \text{a} & \text{t} & \text{s} & \text{t} & \text{a} & \text{c} & \text{s} \\ \hline 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ \hline \end{array}$$


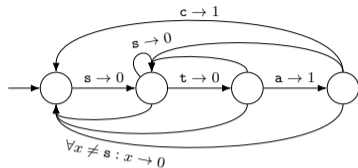
- preprocess  $q$ -state transducer in  $\mathcal{O}(qn^\varepsilon)$  time

## Some new tools we developed

$$T = \begin{array}{|cccccccccccccccc|} \hline \text{s t a r t s t a c k s a t s t a c s} \\ \hline 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 1 0 \\ \hline \end{array}$$


- preprocess  $q$ -state transducer in  $\mathcal{O}(qn^\epsilon)$  time
- run transducer in  $\mathcal{O}(n/\log_\sigma n)$  time

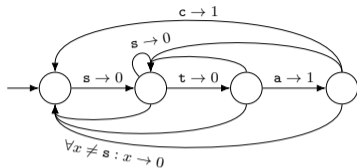
## Some new tools we developed

$$T = \begin{array}{|cccccccccccccccc|} \hline \text{s t a r t s t a c k s a t s t a c s} \\ \hline 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 1 0 \\ \hline \end{array}$$


- preprocess  $q$ -state transducer in  $\mathcal{O}(qn^\epsilon)$  time
- run transducer in  $\mathcal{O}(n/\log_\sigma n)$  time
- also for multiple input/output streams

## Some new tools we developed

$$T = \begin{array}{|cccccccccccccccc|} \hline \text{s t a r t s t a c k s a t s t a c s} \\ \hline 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 1 0 \\ \hline \end{array}$$

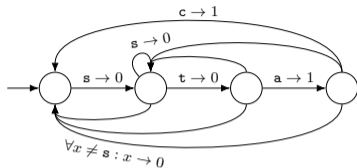


- preprocess  $q$ -state transducer in  $\mathcal{O}(qn^\epsilon)$  time
- run transducer in  $\mathcal{O}(n/\log_\sigma n)$  time
- also for multiple input/output streams

Custom sparse encoding:

## Some new tools we developed

$$T = \begin{array}{|cccccccccccccccc|} \hline \text{s t a r t s t a c k s a t s t a c s} \\ \hline 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 1 0 \\ \hline \end{array}$$



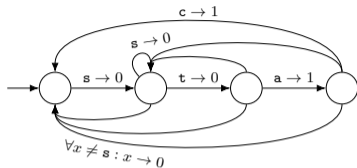
- preprocess  $q$ -state transducer in  $\mathcal{O}(qn^\varepsilon)$  time
- run transducer in  $\mathcal{O}(n/\log_\sigma n)$  time
- also for multiple input/output streams

### Custom sparse encoding:

$$T = \begin{array}{|cccccccccccccccc|} \hline 0 0 0 0 2 0 0 6 0 0 0 1 0 0 0 0 0 9 0 0 4 4 \\ \hline \end{array}$$

## Some new tools we developed

$$T = \begin{array}{|cccccccccccccccc|} \hline \text{s t a r t s t a c k s a t s t a c s} \\ \hline 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 1 0 \\ \hline \end{array}$$



- preprocess  $q$ -state transducer in  $\mathcal{O}(qn^\varepsilon)$  time
- run transducer in  $\mathcal{O}(n/\log_\sigma n)$  time
- also for multiple input/output streams

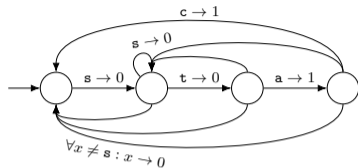
### Custom sparse encoding:

$$T = \begin{array}{|cccccccccccccccc|} \hline 0 0 0 0 2 0 0 6 0 0 0 1 0 0 0 0 0 9 0 0 4 4 \\ \hline \underbrace{\hspace{4em}}_4 \quad \underbrace{\hspace{2em}}_2 \quad \underbrace{\hspace{3em}}_3 \quad \underbrace{\hspace{5em}}_5 \quad \underbrace{\hspace{2em}}_2 \\ \hline \end{array}$$

## Some new tools we developed

$$T = \boxed{\text{s t a r t s t a c k s a t s t a c s}}$$

$$0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0$$



- preprocess  $q$ -state transducer in  $\mathcal{O}(qn^\varepsilon)$  time
- run transducer in  $\mathcal{O}(n/\log_\sigma n)$  time
- also for multiple input/output streams

### Custom sparse encoding:

$$T = \boxed{0 \ 0 \ 0 \ 0 \ 2 \ 0 \ 0 \ 6 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 9 \ 0 \ 0 \ 4 \ 4}$$

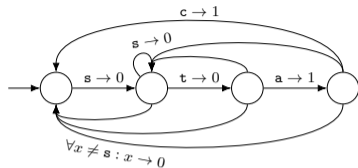
4
2
3
5
2

$$\text{enc}(T) = 0 \cdot \delta(4) \cdot 1 \cdot \delta(2) \cdot 0 \cdot \delta(2) \cdot 1 \cdot \delta(6) \cdot 0 \cdot \delta(3) \cdot 1 \cdot \delta(1) \cdot 0 \cdot \delta(5) \cdot 1 \cdot \delta(9) \cdot 0 \cdot \delta(2) \cdot 1 \cdot \delta(4) \cdot 1 \cdot \delta(4)$$

## Some new tools we developed

$$T = \boxed{\text{s t a r t s t a c k s a t s t a c s}}$$

$$0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0$$



- preprocess  $q$ -state transducer in  $\mathcal{O}(qn^\varepsilon)$  time
- run transducer in  $\mathcal{O}(n/\log_\sigma n)$  time or  $\mathcal{O}(|\text{enc}(T)|/\log n)$
- also for multiple input/output streams

### Custom sparse encoding:

$$T = \boxed{0 \ 0 \ 0 \ 0 \ 2 \ 0 \ 0 \ 6 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 9 \ 0 \ 0 \ 4 \ 4}$$

4
2
3
5
2

$$\text{enc}(T) = 0 \cdot \delta(4) \cdot 1 \cdot \delta(2) \cdot 0 \cdot \delta(2) \cdot 1 \cdot \delta(6) \cdot 0 \cdot \delta(3) \cdot 1 \cdot \delta(1) \cdot 0 \cdot \delta(5) \cdot 1 \cdot \delta(9) \cdot 0 \cdot \delta(2) \cdot 1 \cdot \delta(4) \cdot 1 \cdot \delta(4)$$

## Conclusion

- string synchronizing sets in  $\mathcal{O}(n/\log_{\sigma} n)$  time for all  $\tau$
- truly linear time transducers
- sparse encoding for even faster transducers







## Conclusion

- string synchronizing sets in  $\mathcal{O}(n/\log_\sigma n)$  time for all  $\tau$
- truly linear time transducers
- sparse encoding for even faster transducers






### Open questions:

- complexity gap  $\mathcal{O}(n/\log_\sigma n)$  vs.  $\mathcal{O}(n \log \sigma / \sqrt{\log n})$ , see [Kempa+Kociumaka '25]
- faster algorithms for  $w \gg \log n$ , can we get  $\mathcal{O}(n \log \sigma / w)$  time?

## References

-  Bille, Philip (2011). “Fast searching in packed strings”. In: *J. Discrete Algorithms* 9.1, pp. 49–56. DOI: 10.1016/J.JDA.2010.09.003.
-  Ellert, Jonas (2023). “Sublinear Time Lempel-Ziv (LZ77) Factorization”. In: *String Processing and Information Retrieval - 30th International Symposium, SPIRE 2023, Pisa, Italy, September 26-28, 2023, Proceedings*. Ed. by Franco Maria Nardini, Nadia Pisanti, and Rossano Venturini. Vol. 14240. Lecture Notes in Computer Science. Springer, pp. 171–187. DOI: 10.1007/978-3-031-43980-3\\_14.
-  Farach, Martin (1997). “Optimal Suffix Tree Construction with Large Alphabets”. In: *38th Annual Symposium on Foundations of Computer Science, FOCS 1997*. IEEE Computer Society, pp. 137–143. DOI: 10.1109/SFCS.1997.646102.
-  Hagerup, Torben (1998). “Sorting and Searching on the Word RAM”. In: *Proceedings of the 15th Annual Symposium on Theoretical Aspects of Computer Science (STACS 1998)*. Vol. 1373. Lecture Notes in Computer Science. Springer, pp. 366–398. DOI: 10.1007/BFb0028578.
-  Kempa, Dominik and Tomasz Kociumaka (2019). “String synchronizing sets: sublinear-time BWT construction and optimal LCE data structure”. In: *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing, STOC 2019, Phoenix, AZ, USA, June 23-26, 2019*. Ed. by Moses Charikar and Edith Cohen. ACM, pp. 756–767. DOI: 10.1145/3313276.3316368.
-  — (2024). “Lempel-Ziv (LZ77) Factorization in Sublinear Time”. In: *65th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2024, Chicago, IL, USA, October 27-30, 2024*. IEEE, pp. 2045–2055. DOI: 10.1109/FOCS61266.2024.00122.

## References

-  Kempa, Dominik and Tomasz Kociumaka (2025). “On the Hardness Hierarchy for the  $O(n\sqrt{\log n})$  Complexity in the Word RAM”. In: *Proceedings of the 57th Annual ACM Symposium on Theory of Computing, STOC 2025, Prague, Czechia, June 23-27, 2025*. Ed. by Michal Koucký and Nikhil Bansal. ACM, pp. 290–300. DOI: 10.1145/3717823.3718291.
-  Knuth, Donald E., James H. Morris Jr., and Vaughan R. Pratt (1977). “Fast Pattern Matching in Strings”. In: *SIAM Journal on Computing* 6.2, pp. 323–350. DOI: 10.1137/0206024.
-  Kociumaka, Tomasz et al. (2024). “Internal Pattern Matching Queries in a Text and Applications”. In: *SIAM J. Comput.* 53.5, pp. 1524–1577. DOI: 10.1137/23M1567618.
-  Storer, James A. and Thomas G. Szymanski (1982). “Data compression via textual substitution”. In: *Journal of the ACM* 29.4, pp. 928–951. DOI: 10.1145/322344.322346.
-  Ziv, Jacob and Abraham Lempel (1977). “A universal algorithm for sequential data compression”. In: *IEEE Transactions on Information Theory* 23.3, pp. 337–343. DOI: 10.1109/TIT.1977.1055714.