

Computing Matching Statistics on Wheeler DFAs

Alessio Conte¹ Nicola Cotumaccio^{2,3} Travis Gagie³
Giovanni Manzini¹ Nicola Prezza⁴ Marinella Sciortino⁵

¹University of Pisa, Italy

²GSSI, L'Aquila, Italy

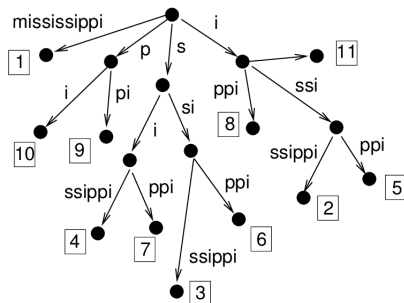
³Dalhousie University, Halifax, Canada

⁴Ca' Foscari University, Venice, Italy

⁵University of Palermo, Italy

Suffix trees

The suffix tree¹ of a string is a general-purpose data structure which is able to efficiently handle a variety of problems (pattern matching, approximate pattern matching, shortest/longest substrings with some desired properties, palindromes, and so on).



¹P. Weiner, *Linear Pattern Matching Algorithms*, FOCS (SWAT) 1973.

- The suffix tree of a string can be compressed by exploiting the repetitiveness of the string².
- It is also possible to build (variants of) suffix trees which are able to solve pattern matching queries where some variables are fixed up to a permutation³.

²T. Gagie, G. Navarro, N. Prezza, *Fully Functional Suffix Trees and Optimal Text Searching in BWT-Runs Bounded Space*, JACM 2020.

³B. S. Baker, *A Theory of Parametrized Pattern Matching: Algorithms and Applications*, STOC 1993.

Suffix trees

- The main limitation of suffix trees is their space consumption.
- As a consequence, suffix trees have been replaced with *suffix arrays*.
- Suffix arrays allow to solve pattern matching queries, but they do not have the whole functionality of suffix trees.

i	Sorted suffixes	SA[i]
1	\$	12
2	i\$	11
3	ippi\$	8
4	issippi\$	5
5	ississippi\$	2
6	mississippi\$	1
7	pi\$	10
8	ppi\$	9
9	sippi\$	7
10	sissippi\$	4
11	ssippi\$	6
12	ssissippi\$	3

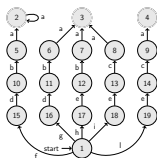
Suffix trees

- Nonetheless, if a suffix array is augmented with the tree topology and the *longest common prefix (LCP) array*, then one retrieves the whole functionality of a suffix tree.
- All these data structures can be compressed, thus leading to *compressed suffix trees*.

i	Sorted suffixes	$LCP[i]$
1	\$	
2	i\$	0
3	ippi\$	1
4	issippi\$	1
5	ississippi\$	4
6	mississippi\$	0
7	pi\$	0
8	ppi\$	1
9	sippi\$	0
10	sissippi\$	2
11	ssippi\$	1
12	ssissippi\$	3

From strings to graphs

- In the meanwhile, the notion of suffix array has been generalized from strings to trees⁴, Wheeler graphs⁵ and arbitrary labeled graphs⁶.
- Since the functionality of a suffix tree can be simulated starting from a suffix array, the natural question is whether it is possible to design suffix trees of graphs.



⁴P. Ferragina, F. Luccio, G. Manzini, S. Muthukrishnan, *Structuring labeled trees for optimal succinctness, and beyond*, FOCS 2005.

⁵T. Gagie, G. Manzini, J. Sirén, *Wheeler graphs: A framework for BWT-based data structures*, TCS 2017.

⁶N. Cotumaccio, N. Prezza, *On indexing and compressing finite automata*, SODA 2021.

Suffix trees of graphs

- In the absence of suffix trees for graphs, some authors have proposed alternative approaches (such as the direct product of graphs), but none of them provides the same flexibility and generality of suffix trees.
- Suffix trees of strings were designed well before the modern development of the data compression field, and the invention of suffix trees of graphs, while elicited by the research on suffix arrays, would have an impact which would go far beyond the applications to compression.
- At the same time, the journey leading to suffix trees of graphs is expected to provide a deeper insight into data compression concepts and techniques (entropy, Lempel-Ziv factorization,...) and their interpretation and application in graph theory.

Our contribution

- The main contribution of our paper is to provide a first clear step towards extending suffix tree functionality to labeled graphs.
- Since in order to simulate the suffix tree of a string we need not only a suffix array, but also an LCP array, we show how to define the LCP array *of a graph*.
- We actually focus on (deterministic) Wheeler graphs (or equivalently, automata), because Wheeler graphs best capture the intuition behind suffix arrays in a graph setting, thus they are expected to be the crucial class of graphs one should work with.

Our contribution

- A classical and useful problem that can be solved using the suffix tree of a string (and in particular the LCP array of a string) is the problem of determining the *matching statistics of a given pattern w.r.t the string*.
- We test our definition of the LCP array for a Wheeler graph by showing that the LCP array can be effectively and efficiently used to determine the *matching statistics of a given pattern w.r.t the Wheeler graph*, which is the natural generalization of the matching statistics problem from a string setting to a graph setting.
- More precisely, we generalize an algorithm⁷ for determining matching statistics from a string setting to a graph setting.

⁷E. Ohlebusch, S. Gog, A. Kügell, *Computing matching statistics and maximal exact matches on compressed full-text indexes*, SPIRE 2010.

Matching statistics of a string π w.r.t a string T

$T = \text{mississippi\$}$
 $\pi = \text{stpissi}$

- ① $l_1 = 1, [l_1, r_1] = [9, 12]$.
- ② $l_2 = 0, [l_2, r_2] = [1, 12]$.
- ③ $l_3 = 2, [l_3, r_3] = [7, 7]$.
- ④ $l_4 = 4, [l_4, r_4] = [4, 5]$.
- ⑤ $l_5 = 3, [l_5, r_5] = [11, 12]$.
- ⑥ $l_6 = 2, [l_6, r_6] = [9, 10]$.
- ⑦ $l_7 = 1, [l_7, r_7] = [2, 5]$.

i	Sorted suffixes
1	\$
2	i\$
3	ippi\$
4	issippi\$
5	ississippi\$
6	mississippi\$
7	pi\$
8	ppi\$
9	sippi\$
10	sissippi\$
11	ssippi\$
12	ssissippi\$

$T = \text{mississippi}\$, \pi = \text{stpissi}$

We start from the end of π , and we repeatedly apply a backward search step (using the FM-index, which relies on the Burrows-Wheeler transform of the string).

- 1 $\ell_1 = 1, [l_1, r_1] = [9, 12]$.
- 2 $\ell_2 = 0, [l_2, r_2] = [1, 12]$.
- 3 $\ell_3 = 2, [l_3, r_3] = [7, 7]$.
- 4 $\ell_4 = 4, [l_4, r_4] = [4, 5]$.
- 5 $\ell_5 = 3, [l_5, r_5] = [11, 12]$.
- 6 $\ell_6 = 2, [l_6, r_6] = [9, 10]$.
- 7 $\ell_7 = 1, [l_7, r_7] = [2, 5]$.

i	Sorted suffixes
1	\$
2	i\$
3	ippi\$
4	issippi\$
5	ississippi\$
6	mississippi\$
7	pi\$
8	ppi\$
9	sippi\$
10	issippi\$
11	ssippi\$
12	ssissippi\$

$T = \text{mississippi}\$, \pi = \text{stpissi}$

$\pi = \text{stpissi}$

- 1 $l_1 = 1, [l_1, r_1] = [9, 12]$.
- 2 $l_2 = 0, [l_2, r_2] = [1, 12]$.
- 3 $l_3 = 2, [l_3, r_3] = [7, 7]$.
- 4 $l_4 = 4, [l_4, r_4] = [4, 5]$.
- 5 $l_5 = 3, [l_5, r_5] = [11, 12]$.
- 6 $l_6 = 2, [l_6, r_6] = [9, 10]$.
- 7 $l_7 = 1, [l_7, r_7] = [2, 5]$.

i	Sorted suffixes	$LCP[i]$	$SA[i]$	$BWT[i]$
1	\$		12	i
2	i\$	0	11	p
3	ippi\$	1	8	s
4	issippi\$	1	5	s
5	ississippi\$	4	2	m
6	mississippi\$	0	1	\$
7	pi\$	0	10	p
8	ppi\$	1	9	i
9	sippi\$	0	7	s
10	sissippi\$	2	4	s
11	ssippi\$	1	6	i
12	ssissippi\$	3	3	i

$T = \text{mississippi}\$, \pi = \text{stpissi}$

$\pi = \text{stpissi}$

① $l_1 = 1, [l_1, r_1] = [9, 12]$.

② $l_2 = 0, [l_2, r_2] = [1, 12]$.

③ $l_3 = 2, [l_3, r_3] = [7, 7]$.

④ $l_4 = 4, [l_4, r_4] = [4, 5]$.

⑤ $l_5 = 3, [l_5, r_5] = [11, 12]$.

⑥ $l_6 = 2, [l_6, r_6] = [9, 10]$.

⑦ $l_7 = 1, [l_7, r_7] = [2, 5]$.

i	Sorted suffixes	$LCP[i]$	$SA[i]$	$BWT[i]$
1	\$		12	i
2	i\$	0	11	p
3	ippi\$	1	8	s
4	issippi\$	1	5	s
5	ississippi\$	4	2	m
6	mississippi\$	0	1	\$
7	pi\$	0	10	p
8	ppi\$	1	9	i
9	sippi\$	0	7	s
10	ssissippi\$	2	4	s
11	ssippi\$	1	6	i
12	ssissippi\$	3	3	i

$T = \text{mississippi}\$, \pi = \text{stpissi}$

$\pi = \text{stpissi}$

① $l_1 = 1, [l_1, r_1] = [9, 12]$.

② $l_2 = 0, [l_2, r_2] = [1, 12]$.

③ $l_3 = 2, [l_3, r_3] = [7, 7]$.

④ $l_4 = 4, [l_4, r_4] = [4, 5]$.

⑤ $l_5 = 3, [l_5, r_5] = [11, 12]$.

⑥ $l_6 = 2, [l_6, r_6] = [9, 10]$.

⑦ $l_7 = 1, [l_7, r_7] = [2, 5]$.

i	Sorted suffixes	$LCP[i]$	$SA[i]$	$BWT[i]$
1	\$		12	i
2	i\$	0	11	p
3	ippi\$	1	8	s
4	issippi\$	1	5	s
5	ississippi\$	4	2	m
6	mississippi\$	0	1	\$
7	pi\$	0	10	p
8	ppi\$	1	9	i
9	sippi\$	0	7	s
10	sissippi\$	2	4	s
11	ssippi\$	1	6	i
12	ssissippi\$	3	3	i

$T = \text{mississippi}\$, \pi = \text{stpissi}$

$\pi = \text{stpissi}$

① $l_1 = 1, [l_1, r_1] = [9, 12]$.

② $l_2 = 0, [l_2, r_2] = [1, 12]$.

③ $l_3 = 2, [l_3, r_3] = [7, 7]$.

④ $l_4 = 4, [l_4, r_4] = [4, 5]$.

⑤ $l_5 = 3, [l_5, r_5] = [11, 12]$.

⑥ $l_6 = 2, [l_6, r_6] = [9, 10]$.

⑦ $l_7 = 1, [l_7, r_7] = [2, 5]$.

i	Sorted suffixes	$LCP[i]$	$SA[i]$	$BWT[i]$
1	\$		12	i
2	i\$	0	11	p
3	ippi\$	1	8	s
4	issippi\$	1	5	s
5	ississippi\$	4	2	m
6	mississippi\$	0	1	\$
7	pi\$	0	10	p
8	ppi\$	1	9	i
9	sippi\$	0	7	s
10	sissippi\$	2	4	s
11	ssippi\$	1	6	i
12	ssissippi\$	3	3	i

$T = \text{mississippi}\$, \pi = \text{stpissi}$

$\pi = \text{stpissi}$

- 1 $\ell_1 = 1, [l_1, r_1] = [9, 12]$.
- 2 $\ell_2 = 0, [l_2, r_2] = [1, 12]$.
- 3 $\ell_3 = 2, [l_3, r_3] = [7, 7]$.
- 4 $\ell_4 = 4, [l_4, r_4] = [4, 5]$.
- 5 $\ell_5 = 3, [l_5, r_5] = [11, 12]$.
- 6 $\ell_6 = 2, [l_6, r_6] = [9, 10]$.
- 7 $\ell_7 = 1, [l_7, r_7] = [2, 5]$.

i	Sorted suffixes	$LCP[i]$	$SA[i]$	$BWT[i]$
1	\$		12	i
2	i\$	0	11	p
3	ippi\$	1	8	s
4	issippi\$	1	5	s
5	issippi\$	4	2	m
6	mississippi\$	0	1	\$
7	pi\$	0	10	p
8	ppi\$	1	9	i
9	sippi\$	0	7	s
10	sissippi\$	2	4	s
11	ssippi\$	1	6	i
12	ssissippi\$	3	3	i

The backward search fails: "pissi" does not occur in $T = \text{mississippi}\$$.

$T = \text{mississippi}\$, \pi = \text{stpissi}$

We use the LCP array to determine the longest prefix of "issi" with *more* occurrences than "issi" in T .

$\pi = \text{stpissi}$

① $\ell_1 = 1, [l_1, r_1] = [9, 12]$.

② $\ell_2 = 0, [l_2, r_2] = [1, 12]$.

③ $\ell_3 = 2, [l_3, r_3] = [7, 7]$.

④ $\ell_4 = 4, [l_4, r_4] = [4, 5]$.

⑤ $\ell_5 = 3, [l_5, r_5] = [11, 12]$.

⑥ $\ell_6 = 2, [l_6, r_6] = [9, 10]$.

⑦ $\ell_7 = 1, [l_7, r_7] = [2, 5]$.

i	Sorted suffixes	$LCP[i]$	$SA[i]$	$BWT[i]$
1	\$		12	i
2	i\$	0	11	p
3	ippi\$	1	8	s
4	issippi\$	1	5	s
5	ississippi\$	4	2	m
6	mississippi\$	0	1	\$
7	pi\$	0	10	p
8	ppi\$	1	9	i
9	sippi\$	0	7	s
10	sissippi\$	2	4	s
11	ssippi\$	1	6	i
12	ssissippi\$	3	3	i

The desired prefix is the one of length $\max\{1, 0\} = 1$, that is, "i".

$T = \text{mississippi}\$, \pi = \text{stpissi}$

We now extend⁸ the interval $[4, 5]$ (suffixes starting with "issi") to $[2, 5]$ (suffixes starting with "i").

$\pi = \text{stpissi}$

- 1 $l_1 = 1, [l_1, r_1] = [9, 12]$.
- 2 $l_2 = 0, [l_2, r_2] = [1, 12]$.
- 3 $l_3 = 2, [l_3, r_3] = [7, 7]$.
- 4 $l_4 = 4, [l_4, r_4] = [4, 5]$.
- 5 $l_5 = 3, [l_5, r_5] = [11, 12]$.
- 6 $l_6 = 2, [l_6, r_6] = [9, 10]$.
- 7 $l_7 = 1, [l_7, r_7] = [2, 5]$.

i	Sorted suffixes	$LCP[i]$	$SA[i]$	$BWT[i]$
1	\$		12	i
2	i\$	0	11	p
3	ippi\$	1	8	s
4	issippi\$	1	5	s
5	ississippi\$	4	2	m
6	mississippi\$	0	1	\$
7	pi\$	0	10	p
8	ppi\$	1	9	i
9	sippi\$	0	7	s
10	sissippi\$	2	4	s
11	ssippi\$	1	6	i
12	ssissippi\$	3	3	i

⁸There exists a data structure efficiently supporting such an extension.

$T = \text{mississippi}\$, \pi = \text{stpissi}$

We try again to perform a backward search step.

$\pi = \text{st~~pi~~ssi}$

① $l_1 = 1, [l_1, r_1] = [9, 12]$.

② $l_2 = 0, [l_2, r_2] = [1, 12]$.

③ $l_3 = 2, [l_3, r_3] = [7, 7]$.

④ $l_4 = 4, [l_4, r_4] = [4, 5]$.

⑤ $l_5 = 3, [l_5, r_5] = [11, 12]$.

⑥ $l_6 = 2, [l_6, r_6] = [9, 10]$.

⑦ $l_7 = 1, [l_7, r_7] = [2, 5]$.

i	Sorted suffixes	$LCP[i]$	$SA[i]$	$BWT[i]$
1	\$		12	i
2	i\$	0	11	p
3	ippi\$	1	8	s
4	issippi\$	1	5	s
5	ississippi\$	4	2	m
6	mississippi\$	0	1	\$
7	pi\$	0	10	p
8	ppi\$	1	9	i
9	sippi\$	0	7	s
10	sissippi\$	2	4	s
11	ssippi\$	1	6	i
12	ssissippi\$	3	3	i

This time the backward search step is successful, and we go on with the algorithm.

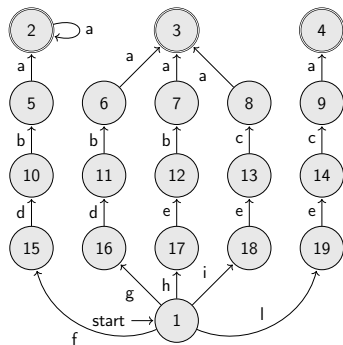
Wheeler graphs

- Wheeler DFAs have a total order on the set of states.
- The set $T(\pi)$ of all states *reached* by a pattern π is always an interval.
- Intuively, the ordering plays the role of the suffix array.

$$T(e) = [12, 14]$$

$$T(ca) = [3, 4]$$

$$T(dba) = [2, 3].$$

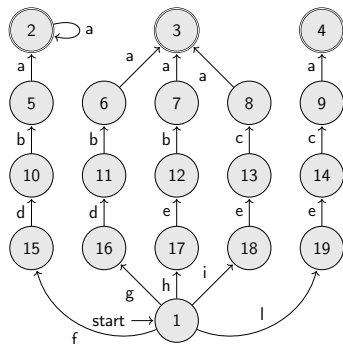


Matching statistics of a string π w.r.t a Wheeler DFA \mathcal{A}

\mathcal{A} = the Wheeler DFA in the figure

$\pi = caa$

- 1 $l_1 = 1, [l_1, r_1] = [8, 9]$.
- 2 $l_2 = 2, [l_1, r_1] = [3, 4]$.
- 3 $l_3 = 2, [l_1, r_1] = [2, 2]$.

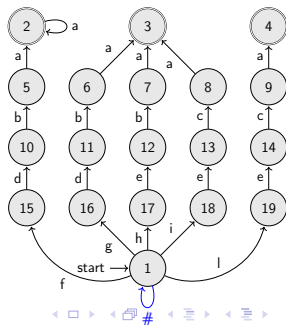


The definition is symmetrical for historical reasons.

LCP array of a Wheeler DFA

- Since in the algorithm on strings we had an LCP array for $T = \text{mississippi}\$,$ now we need to somehow define an LCP array for the Wheeler DFA \mathcal{A} .
- Even though there may be infinitely many strings reaching a state, we show that it suffices to consider the *minimum* and *maximum* such strings.

State	i	Minima and maxima	LCP[i]
1	1	#####...	∞
	2	#####...	
2	3	aaaaa...	0
	4	abdf##...	
3	5	abdg##...	3
	6	acei##...	
4	7	acel##...	3
	8	acel##...	
...



It is not that easy...

- Analogously to the algorithm on strings, given a string π , we will sometimes need to determine the longest suffix π' of π reaching more states than π .
- If $T(\pi) = [r, s]$, we can determine π' by the following formula:

$$|\pi'| = \max \left\{ \min \{ \text{lcp}(\text{max}_{r-1}, \text{min}_r), \text{lcp}(\text{min}_r, \pi^R) \}, \min \{ \text{lcp}(\pi^R, \text{max}_s), \text{lcp}(\text{max}_s, \text{min}_{s+1}) \} \right\}.$$

- **Orange** values are stored in the LCP array (just like in the algorithm on strings).
- **Blue** values were not needed in the algorithm on strings, but they are now needed because some strings reaching $T(\pi)$ may not have π as a suffix. We maintain the blue values, and we show that they can be efficiently updated during the algorithm.

Computing Matching Statistics on Wheeler DFAs

Alessio Conte¹ Nicola Cotumaccio^{2,3} Travis Gagie³
Giovanni Manzini¹ Nicola Prezza⁴ Marinella Sciortino⁵

¹University of Pisa, Italy

²GSSI, L'Aquila, Italy

³Dalhousie University, Halifax, Canada

⁴Ca' Foscari University, Venice, Italy

⁵University of Palermo, Italy