

Accessing the Suffix Array via ϕ^{-1} -Forest



Christina Boucher¹

¹Herbert Wertheim College of Engineering, University of Florida,
USA



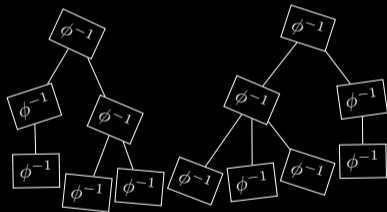
Dominik Köppl²

²M&D Data Science Center, Tokyo Medical and Dental University,
Japan



Herman Perera¹

Massimiliano Rossi¹



extended talk on
https://youtu.be/SjAX1Ru2_gE



what's it all about?

goal

given the r -index, can we practically speed up the time for retrieving a suffix array (SA) entry?

why interesting?

- ▀ r -index is a version of the FM-index with refined SA samples
- ▀ while using less space than the FM-index on repetitive texts, the sparse SA samples make random accesses to $SA[i]$ slow

example

index sequences

- ▮ GATTACAT

- ▮ GATACAT

- ▮ GATTAGATA

for that:

- ▮ concatenate with \$, and

- ▮ use # as terminal symbol

input becomes $T = \text{GATTACAT\$GATACAT\$GATTAGATA\#}$

i	SA	rotation matrix	BWT
1	27	#GATTACAT\$GATACAT\$GATTAGAT	A
2	9	\$GATACAT\$GATTAGATA#GATTACA	T
3	17	\$GATTAGATA#GATTACAT\$GATACA	T
4	26	A#GATTACAT\$GATACAT\$GATTAGA	T
5	5	ACAT\$GATACAT\$GATTAGATA#GAT	T
6	13	ACAT\$GATTAGATA#GATTACAT\$GA	T
7	22	AGATA#GATTACAT\$GATACAT\$GAT	T
8	7	AT\$GATACAT\$GATTAGATA#GATTA	C
9	15	AT\$GATTAGATA#GATTACAT\$GATA	C
10	24	ATA#GATTACAT\$GATACAT\$GATTA	G
11	11	ATACAT\$GATTAGATA#GATTACAT\$	G
12	2	ATTACAT\$GATACAT\$GATTAGATA#	G
13	19	ATTAGATA#GATTACAT\$GATACAT\$	G
14	6	CAT\$GATACAT\$GATTAGATA#GATT	A
15	14	CAT\$GATTAGATA#GATTACAT\$GAT	A
16	23	GATA#GATTACAT\$GATACAT\$GATT	A
17	10	GATACAT\$GATTAGATA#GATTACAT	\$
18	1	GATTACAT\$GATACAT\$GATTAGATA	#
19	18	GATTAGATA#GATTACAT\$GATACAT	\$
20	8	T\$GATACAT\$GATTAGATA#GATTAC	A
21	16	T\$GATTAGATA#GATTACAT\$GATAC	A
22	25	TA#GATTACAT\$GATACAT\$GATTAG	A
23	4	TACAT\$GATACAT\$GATTAGATA#GA	T
24	12	TACAT\$GATTAGATA#GATTACAT\$G	A
25	21	TAGATA#GATTACAT\$GATACAT\$GA	T
26	3	TTACAT\$GATACAT\$GATTAGATA#G	A
27	20	TTAGATA#GATTACAT\$GATACAT\$G	A

FM-index for text $T :=$

GATTACAT\$GATACAT\$GATTAGATA#

- uses BWT and a wavelet tree for pattern matching
- counting pattern occurrences works out of the box
- for locating the pattern occurrences, we need SA
- FM index samples SA by text position distance

i	SA	rotation matrix	BWT
1	27	#GATTACAT\$GATACAT\$GATTAGAT	A
2	9	\$GATACAT\$GATTAGATA#GATTACA	T
3	17	\$GATTAGATA#GATTACAT\$GATACA	T
4	26	A#GATTACAT\$GATACAT\$GATTAGA	T
5	5	ACAT\$GATACAT\$GATTAGATA#GAT	T
6	13	ACAT\$GATTAGATA#GATTACAT\$GA	T
7	22	AGATA#GATTACAT\$GATACAT\$GAT	T
8	7	AT\$GATACAT\$GATTAGATA#GATTA	C
9	15	AT\$GATTAGATA#GATTACAT\$GATA	C
10	24	ATA#GATTACAT\$GATACAT\$GATTA	G
11	11	ATACAT\$GATTAGATA#GATTACAT\$	G
12	2	ATTACAT\$GATACAT\$GATTAGATA#	G
13	19	ATTAGATA#GATTACAT\$GATACAT\$	G
14	6	CAT\$GATACAT\$GATTAGATA#GATT	A
15	14	CAT\$GATTAGATA#GATTACAT\$GAT	A
16	23	GATA#GATTACAT\$GATACAT\$GATT	A
17	10	GATACAT\$GATTAGATA#GATTACAT	\$
18	1	GATTACAT\$GATACAT\$GATTAGATA	#
19	18	GATTAGATA#GATTACAT\$GATACAT	\$
20	8	T\$GATACAT\$GATTAGATA#GATTAC	A
21	16	T\$GATTAGATA#GATTACAT\$GATAC	A
22	25	TA#GATTACAT\$GATACAT\$GATTAG	A
23	4	TACAT\$GATACAT\$GATTAGATA#GA	T
24	12	TACAT\$GATTAGATA#GATTACAT\$G	A
25	21	TAGATA#GATTACAT\$GATACAT\$GA	T
26	3	TTACAT\$GATACAT\$GATTAGATA#G	A
27	20	TTAGATA#GATTACAT\$GATACAT\$G	A

FM-index for text $T :=$

GATTACAT\$GATACAT\$GATTAGATA#

- uses BWT and a wavelet tree for pattern matching
- counting pattern occurrences works out of the box
- for locating the pattern occurrences, we need SA
- FM index samples SA by text position distance

r -index

- only stores SA samples at run boundaries

i	SA	rotation matrix	BWT
1	27	#GATTACAT\$GATACAT\$GATTAGAT	A
2	9	\$GATACAT\$GATTAGATA#GATTACA	T
3	17	\$GATTAGATA#GATTACAT\$GATACA	T
4	26	A#GATTACAT\$GATACAT\$GATTAGA	T
5	5	ACAT\$GATACAT\$GATTAGATA#GAT	T
6	13	ACAT\$GATTAGATA#GATTACAT\$GA	T
7	22	AGATA#GATTACAT\$GATACAT\$GAT	T
8	7	AT\$GATACAT\$GATTAGATA#GATTA	C
9	15	AT\$GATTAGATA#GATTACAT\$GATA	C
10	24	ATA#GATTACAT\$GATACAT\$GATTA	G
11	11	ATACAT\$GATTAGATA#GATTACAT\$	G
12	2	ATTACAT\$GATACAT\$GATTAGATA#	G
13	19	ATTAGATA#GATTACAT\$GATACAT\$	G
14	6	CAT\$GATACAT\$GATTAGATA#GATT	A
15	14	CAT\$GATTAGATA#GATTACAT\$GAT	A
16	23	GATA#GATTACAT\$GATACAT\$GATT	A
17	10	GATACAT\$GATTAGATA#GATTACAT	\$
18	1	GATTACAT\$GATACAT\$GATTAGATA	#
19	18	GATTAGATA#GATTACAT\$GATACAT	\$
20	8	T\$GATACAT\$GATTAGATA#GATTAC	A
21	16	T\$GATTAGATA#GATTACAT\$GATAC	A
22	25	TA#GATTACAT\$GATACAT\$GATTAG	A
23	4	TACAT\$GATACAT\$GATTAGATA#GA	T
24	12	TACAT\$GATTAGATA#GATTACAT\$G	A
25	21	TAGATA#GATTACAT\$GATACAT\$GA	T
26	3	TTACAT\$GATACAT\$GATTAGATA#G	A
27	20	TTAGATA#GATTACAT\$GATACAT\$G	A

SA access

how do we get $SA[i]$ with the r -index?

Toehold Lemma, Gagie+'18

$$BWT[i] = BWT[i + 1] \Rightarrow$$

$$SA[i + 1] - SA[i] = SA[j + 1] - SA[j] \text{ for}$$

$$SA[j] := SA[i] - 1$$

Example

- ▀ $SA[2] = 9, SA[20] = 8$
- ▀ $BWT[2] = BWT[3] = T$
- ▀ $SA[3] - SA[2] = SA[21] - SA[20]$.

i	SA	rotation matrix	BWT
1	27	#GATTACAT\$GATACAT\$GATTAGAT	A
2	9	\$GATACAT\$GATTAGATA#GATTACA	T
3	17	\$GATTAGATA#GATTACAT\$GATACA	T
4	26	A#GATTACAT\$GATACAT\$GATTAGA	T
5	5	ACAT\$GATACAT\$GATTAGATA#GAT	T
6	13	ACAT\$GATTAGATA#GATTACAT\$GA	T
7	22	AGATA#GATTACAT\$GATACAT\$GAT	T
8	7	AT\$GATACAT\$GATTAGATA#GATTA	C
9	15	AT\$GATTAGATA#GATTACAT\$GATA	C
10	24	ATA#GATTACAT\$GATACAT\$GATTA	G
11	11	ATACAT\$GATTAGATA#GATTACAT\$	G
12	2	ATTACAT\$GATACAT\$GATTAGATA#	G
13	19	ATTAGATA#GATTACAT\$GATACAT\$	G
14	6	CAT\$GATACAT\$GATTAGATA#GATT	A
15	14	CAT\$GATTAGATA#GATTACAT\$GAT	A
16	23	GATA#GATTACAT\$GATACAT\$GATT	A
17	10	GATACAT\$GATTAGATA#GATTACAT	\$
18	1	GATTACAT\$GATACAT\$GATTAGATA	#
19	18	GATTAGATA#GATTACAT\$GATACAT	\$
20	8	T\$GATACAT\$GATTAGATA#GATTAC	A
21	16	T\$GATTAGATA#GATTACAT\$GATAC	A
22	25	TA#GATTACAT\$GATACAT\$GATTAG	A
23	4	TACAT\$GATACAT\$GATTAGATA#GA	T
24	12	TACAT\$GATTAGATA#GATTACAT\$G	A
25	21	TAGATA#GATTACAT\$GATACAT\$GA	T
26	3	TTACAT\$GATACAT\$GATTAGATA#G	A
27	20	TTAGATA#GATTACAT\$GATACAT\$G	A

Corollary

- let $k \geq 0$ be largest value such that $\text{BWT}[i'] = \text{BWT}[i' + 1]$ for all i' with $\text{SA}[i'] \in [\text{SA}[i] - k + 1.. \text{SA}[i]]$
- let $\text{SA}[j] := \text{SA}[i] - k$

Then:

- $\text{BWT}[j] \neq \text{BWT}[j + 1]$ but still
- $\text{SA}[i + 1] - \text{SA}[i] = \text{SA}[j + 1] - \text{SA}[j]$.

i	SA	rotation matrix	BWT
1	27	#GATTACAT\$GATACAT\$GATTAGAT	A
2	9	\$GATACAT\$GATTAGATA#GATTACA	T
3	17	\$GATTAGATA#GATTACAT\$GATACA	T
4	26	A#GATTACAT\$GATACAT\$GATTAGA	T
5	5	ACAT\$GATACAT\$GATTAGATA#GAT	T
6	13	ACAT\$GATTAGATA#GATTACAT\$GA	T
7	22	AGATA#GATTACAT\$GATACAT\$GAT	T
8	7	AT\$GATACAT\$GATTAGATA#GATTA	C
9	15	AT\$GATTAGATA#GATTACAT\$GATA	C
10	24	ATA#GATTACAT\$GATACAT\$GATTA	G
11	11	ATACAT\$GATTAGATA#GATTACAT\$	G
12	2	ATTACAT\$GATACAT\$GATTAGATA#	G
13	19	ATTAGATA#GATTACAT\$GATACAT\$	G
14	6	CAT\$GATACAT\$GATTAGATA#GATT	A
15	14	CAT\$GATTAGATA#GATTACAT\$GAT	A
16	23	GATA#GATTACAT\$GATACAT\$GATT	A
17	10	GATACAT\$GATTAGATA#GATTACAT	\$
18	1	GATTACAT\$GATACAT\$GATTAGATA	#
19	18	GATTAGATA#GATTACAT\$GATACAT	\$
20	8	T\$GATACAT\$GATTAGATA#GATTAC	A
21	16	T\$GATTAGATA#GATTACAT\$GATAC	A
22	25	TA#GATTACAT\$GATACAT\$GATTAG	A
23	4	TACAT\$GATACAT\$GATTAGATA#GA	T
24	12	TACAT\$GATTAGATA#GATTACAT\$G	A
25	21	TAGATA#GATTACAT\$GATACAT\$GA	T
26	3	TTACAT\$GATACAT\$GATTAGATA#G	A
27	20	TTAGATA#GATTACAT\$GATACAT\$G	A

Corollary

- let $k \geq 0$ be largest value such that $\text{BWT}[i'] = \text{BWT}[i' + 1]$ for all i' with $\text{SA}[i'] \in [\text{SA}[i] - k + 1.. \text{SA}[i]]$
- let $\text{SA}[j] := \text{SA}[i] - k$

Then:

- $\text{BWT}[j] \neq \text{BWT}[j + 1]$ but still
- $\text{SA}[i + 1] - \text{SA}[i] = \text{SA}[j + 1] - \text{SA}[j]$.

i	SA	rotation matrix	BWT
1	27	#GATTACAT\$GATACAT\$GATTAGAT	A
2	9	\$GATACAT\$GATTAGATA#GATTACA	T
3	17	\$GATTAGATA#GATTACAT\$GATACA	T
4	26	A#GATTACAT\$GATACAT\$GATTAGA	T
5	5	ACAT\$GATACAT\$GATTAGATA#GAT	T
6	13	ACAT\$GATTAGATA#GATTACAT\$GA	T
7	22	AGATA#GATTACAT\$GATACAT\$GAT	T
8	7	AT\$GATACAT\$GATTAGATA#GATTA	C
9	15	AT\$GATTAGATA#GATTACAT\$GATA	C
10	24	ATA#GATTACAT\$GATACAT\$GATTA	G
11	11	ATACAT\$GATTAGATA#GATTACAT\$	G
12	2	ATTACAT\$GATACAT\$GATTAGATA#	G
13	19	ATTAGATA#GATTACAT\$GATACAT\$	G
14	6	CAT\$GATACAT\$GATTAGATA#GATT	A
15	14	CAT\$GATTAGATA#GATTACAT\$GAT	A
16	23	GATA#GATTACAT\$GATACAT\$GATT	A
17	10	GATACAT\$GATTAGATA#GATTACAT	\$
18	1	GATTACAT\$GATACAT\$GATTAGATA	#
19	18	GATTAGATA#GATTACAT\$GATACAT	\$
20	8	T\$GATACAT\$GATTAGATA#GATTAC	A
21	16	T\$GATTAGATA#GATTACAT\$GATAC	A
22	25	TA#GATTACAT\$GATACAT\$GATTAG	A
23	4	TACAT\$GATACAT\$GATTAGATA#GA	T
24	12	TACAT\$GATTAGATA#GATTACAT\$G	A
25	21	TAGATA#GATTACAT\$GATACAT\$GA	T
26	3	TTACAT\$GATACAT\$GATTAGATA#G	A
27	20	TTAGATA#GATTACAT\$GATACAT\$G	A

Corollary

- let $k \geq 0$ be largest value such that $\text{BWT}[i'] = \text{BWT}[i' + 1]$ for all i' with $\text{SA}[i'] \in [\text{SA}[i] - k + 1.. \text{SA}[i]]$
- let $\text{SA}[j] := \text{SA}[i] - k$

Then:

- $\text{BWT}[j] \neq \text{BWT}[j + 1]$ but still
- $\text{SA}[i + 1] - \text{SA}[i] = \text{SA}[j + 1] - \text{SA}[j]$.

i	SA	rotation matrix	BWT
1	27	#GATTACAT\$GATACAT\$GATTAGAT	A
2	9	\$GATACAT\$GATTAGATA#GATTACA	T
3	17	\$GATTAGATA#GATTACAT\$GATACA	T
4	26	A#GATTACAT\$GATACAT\$GATTAGA	T
5	5	ACAT\$GATACAT\$GATTAGATA#GAT	T
6	13	ACAT\$GATTAGATA#GATTACAT\$GA	T
7	22	AGATA#GATTACAT\$GATACAT\$GAT	T
8	7	AT\$GATACAT\$GATTAGATA#GATTA	C
9	15	AT\$GATTAGATA#GATTACAT\$GATA	C
10	24	ATA#GATTACAT\$GATACAT\$GATTA	G
11	11	ATACAT\$GATTAGATA#GATTACAT\$	G
12	2	ATTACAT\$GATACAT\$GATTAGATA#	G
13	19	ATTAGATA#GATTACAT\$GATACAT\$	G
14	6	CAT\$GATACAT\$GATTAGATA#GATT	A
15	14	CAT\$GATTAGATA#GATTACAT\$GAT	A
16	23	GATA#GATTACAT\$GATACAT\$GATT	A
17	10	GATACAT\$GATTAGATA#GATTACAT	\$
18	1	GATTACAT\$GATACAT\$GATTAGATA	#
19	18	GATTAGATA#GATTACAT\$GATACAT	\$
20	8	T\$GATACAT\$GATTAGATA#GATTAC	A
21	16	T\$GATTAGATA#GATTACAT\$GATAC	A
22	25	TA#GATTACAT\$GATACAT\$GATTAG	A
23	4	TACAT\$GATACAT\$GATTAGATA#GA	T
24	12	TACAT\$GATTAGATA#GATTACAT\$G	A
25	21	TAGATA#GATTACAT\$GATACAT\$GA	T
26	3	TTACAT\$GATACAT\$GATTAGATA#G	A
27	20	TTAGATA#GATTACAT\$GATACAT\$G	A

Corollary

- let $k \geq 0$ be largest value such that $BWT[i'] = BWT[i' + 1]$ for all i' with $SA[i'] \in [SA[i] - k + 1..SA[i]]$
- let $SA[j] := SA[i] - k$

Then:

- $BWT[j] \neq BWT[j + 1]$ but still
- $SA[i + 1] - SA[i] = SA[j + 1] - SA[j]$.

i	SA	rotation matrix	BWT
1	27	#GATTACAT\$GATACAT\$GATTAGAT	A
2	9	\$GATACAT\$GATTAGATA#GATTACA	T
3	17	\$GATTAGATA#GATTACAT\$GATTACA	T
4	26	A#GATTACAT\$GATACAT\$GATTAGA	T
5	5	ACAT\$GATACAT\$GATTAGATA#GAT	T
6	13	ACAT\$GATTAGATA#GATTACAT\$GA	T
7	22	AGATA#GATTACAT\$GATACAT\$GAT	T
8	7	AT\$GATACAT\$GATTAGATA#GATTA	C
9	15	AT\$GATTAGATA#GATTACAT\$GATA	C
10	24	ATA#GATTACAT\$GATACAT\$GATTA	G
11	11	ATACAT\$GATTAGATA#GATTACAT\$	G
12	2	ATTACAT\$GATACAT\$GATTAGATA#	G
13	19	ATTAGATA#GATTACAT\$GATACAT\$	G
14	6	CAT\$GATACAT\$GATTAGATA#GATT	A
15	14	CAT\$GATTAGATA#GATTACAT\$GAT	A
16	23	GATA#GATTACAT\$GATACAT\$GATT	A
17	10	GATACAT\$GATTAGATA#GATTACAT	\$
18	1	GATTACAT\$GATACAT\$GATTAGATA	#
19	18	GATTAGATA#GATTACAT\$GATACAT	\$
20	8	T\$GATACAT\$GATTAGATA#GATTAC	A
21	16	T\$GATTAGATA#GATTACAT\$GATAC	A
22	25	TA#GATTACAT\$GATACAT\$GATTAG	A
23	4	TACAT\$GATACAT\$GATTAGATA#GA	T
24	12	TACAT\$GATTAGATA#GATTACAT\$G	A
25	21	TAGATA#GATTACAT\$GATACAT\$GA	T
26	3	TTACAT\$GATACAT\$GATTAGATA#G	A
27	20	TTAGATA#GATTACAT\$GATACAT\$G	A

Corollary

- let $k \geq 0$ be largest value such that $BWT[i'] = BWT[i' + 1]$ for all i' with $SA[i'] \in [SA[i] - k + 1..SA[i]]$
- let $SA[j] := SA[i] - k$

Then:

- $BWT[j] \neq BWT[j + 1]$ but still
- $SA[i + 1] - SA[i] = SA[j + 1] - SA[j]$.

i	SA	rotation matrix	BWT
1	27	#GATTACAT\$GATACAT\$GATTAGAT	A
2	9	\$GATACAT\$GATTAGATA#GATTACA	T
3	17	\$GATTAGATA#GATTACAT\$GATTACA	T
4	26	A#GATTACAT\$GATACAT\$GATTAGAT	T
5	5	ACAT\$GATACAT\$GATTAGATA#GATT	T
6	13	ACAT\$GATTAGATA#GATTACAT\$GATT	T
7	22	AGATA#GATTACAT\$GATACAT\$GATT	T
8	7	AT\$GATACAT\$GATTAGATA#GATTACA	C
9	15	AT\$GATTAGATA#GATTACAT\$GATTACA	C
10	24	ATA#GATTACAT\$GATACAT\$GATTAGAT	G
11	11	ATACAT\$GATTAGATA#GATTACAT\$GATT	G
12	2	ATTACAT\$GATACAT\$GATTAGATA#GATT	G
13	19	ATTAGATA#GATTACAT\$GATACAT\$GATT	G
14	6	CAT\$GATACAT\$GATTAGATA#GATTAGAT	A
15	14	CAT\$GATTAGATA#GATTACAT\$GATTAGAT	A
16	23	GATA#GATTACAT\$GATACAT\$GATTAGAT	A
17	10	GATACAT\$GATTAGATA#GATTACAT\$GATT	A
18	1	GATTACAT\$GATACAT\$GATTAGATA#GATT	A
19	18	GATTAGATA#GATTACAT\$GATACAT\$GATT	A
20	8	T\$GATACAT\$GATTAGATA#GATTACAT\$GATT	A
21	16	T\$GATTAGATA#GATTACAT\$GATACAT\$GATT	A
22	25	TA#GATTACAT\$GATACAT\$GATTAGAT	A
23	4	TACAT\$GATACAT\$GATTAGATA#GATTACA	T
24	12	TACAT\$GATTAGATA#GATTACAT\$GATTACA	T
25	21	TAGATA#GATTACAT\$GATACAT\$GATTAGAT	T
26	3	TTACAT\$GATACAT\$GATTAGATA#GATTAGAT	A
27	20	TTAGATA#GATTACAT\$GATACAT\$GATTAGAT	A

Corollary

- let $k \geq 0$ be largest value such that $BWT[i'] = BWT[i' + 1]$ for all i' with $SA[i'] \in [SA[i] - k + 1..SA[i]]$
- let $SA[j] := SA[i] - k$

Then:

- $BWT[j] \neq BWT[j + 1]$ but still
- $SA[i + 1] - SA[i] = SA[j + 1] - SA[j]$.

why it works:

for each backward step, we move to the preceding character pair in the text

i	SA	rotation matrix	BWT
1	27	#GATTACAT\$GATACAT\$GATTAGAT	A
2	9	\$GATACAT\$GATTAGATA#GATTACA	T
3		\$GATTAGATA#GATTACAT\$GATACA	T
4		A#GATTACAT\$GATACAT\$GATTAGA	T
5		ACAT\$GATACAT\$GATTAGATA#GAT	T
6		ACAT\$GATTAGATA#GATTACAT\$GA	T
7	22	AGATA#GATTACAT\$GATACAT\$GAT	T
8	7	AT\$GATACAT\$GATTAGATA#GATTA	C
9	15	AT\$GATTAGATA#GATTACAT\$GATA	C
10	24	ATA#GATTACAT\$GATACAT\$GATTA	G
11		ATACAT\$GATTAGATA#GATTACAT\$	G
12		ATTACAT\$GATACAT\$GATTAGATA#	G
13	19	ATTAGATA#GATTACAT\$GATACAT\$	G
14	6	CAT\$GATACAT\$GATTAGATA#GATT	A
15		CAT\$GATTAGATA#GATTACAT\$GAT	A
16	23	GATA#GATTACAT\$GATACAT\$GATT	A
17	10	GATACAT\$GATTAGATA#GATTACAT	\$
18	1	GATTACAT\$GATACAT\$GATTAGATA	#
19	18	GATTAGATA#GATTACAT\$GATACAT	\$
20	8	T\$GATACAT\$GATTAGATA#GATTAC	A
21		T\$GATTAGATA#GATTACAT\$GATAC	A
22	25	TA#GATTACAT\$GATACAT\$GATTAG	A
23	4	TACAT\$GATACAT\$GATTAGATA#GA	T
24	12	TACAT\$GATTAGATA#GATTACAT\$G	A
25	21	TAGATA#GATTACAT\$GATACAT\$GA	T
26	3	TTACAT\$GATACAT\$GATTAGATA#G	A
27	20	TTAGATA#GATTACAT\$GATACAT\$G	A

r -index in $\mathcal{O}(r)$ space

store

- ▮ $\mathcal{S}[x]$: sample at start of x -th run

- ▮ $\mathcal{E}[x]$: sample at end of x -th run

where $x \in [1..r]$, and r is the number of character runs in BWT.

interested in following queries on \mathcal{E} :

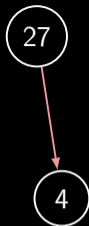
- ▮ $\mathcal{E}.\text{pred}(p) : \max\{q \in \mathcal{E} : q \leq p\}$

- ▮ $\mathcal{E}.\text{succ}(p) : \min\{q \in \mathcal{E} : q > p\}$

for that: build predecessor and successor data structure on \mathcal{E}

$\mathcal{E}.\text{pred}(\mathcal{S}[x + 1])$

x	$\mathcal{S}[x]$	$\mathcal{E}[x]$
1	27	27
2	4	9
3	4	7
4	23	24
5	4	6
6	10	10
7	1	1
8	18	18
9	4	8
10	4	4
11	12	12
12	21	21
13	1	3

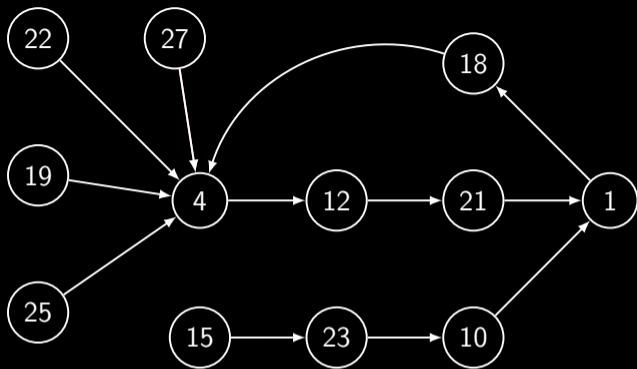


ϕ^{-1} graph

- ▀ each entry of \mathcal{E} is a node
- ▀ create an arc from $\mathcal{E}[x]$ to $\mathcal{E}[y]$ if $\mathcal{E}[y] = \mathcal{E}.\text{pred}(\mathcal{S}[x + 1])$
(sorted $\mathcal{E} = [1, 4, 10, 12, 15, 18, 19, 20, 21, 22, 23, 25, 27]$)

$\mathcal{E}.\text{pred}(\mathcal{S}[x + 1])$

x	$\mathcal{S}[x]$	$\mathcal{E}[x]$
1	27	27
2	4	9
3	4	7
4	23	24
5	4	6
6	10	10
7	1	1
8	18	18
9	4	8
10	4	4
11	12	12
12	21	21
13	1	3

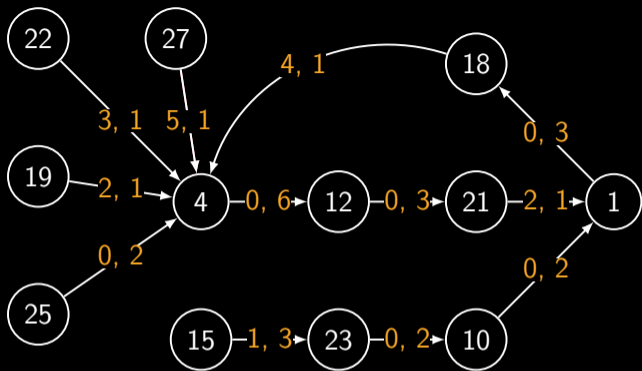


ϕ^{-1} graph

- each entry of \mathcal{E} is a node
- create an arc from $\mathcal{E}[x]$ to $\mathcal{E}[y]$ if $\mathcal{E}[y] = \mathcal{E}.\text{pred}(\mathcal{S}[x + 1])$ (sorted $\mathcal{E} = [1, 4, 10, 12, 15, 18, 19, 20, 21, 22, 23, 25, 27]$)

$\mathcal{E}.pred(\mathcal{S}[x + 1])$

x	$\mathcal{S}[x]$	$\mathcal{E}[x]$	c_x	l_x	
1	27	27	5	1	
2	4	9	3	1	
3	4	7	1	3	
4	23	24	2	1	
5	4	6	0	2	
6	10	10	0	2	
7	1	1	1	0	3
8	18	18	4	1	
9	4	8	25	0	2
10	4	4	4	0	6
11	12	12	12	0	3
12	21	21	21	0	2
13	1	3	20	-	2



▀ label the arc of x -th run with cost c_x and limit l_x

▀ $c_x := \mathcal{S}[x + 1] - \mathcal{E}.pred(\mathcal{S}[x + 1])$

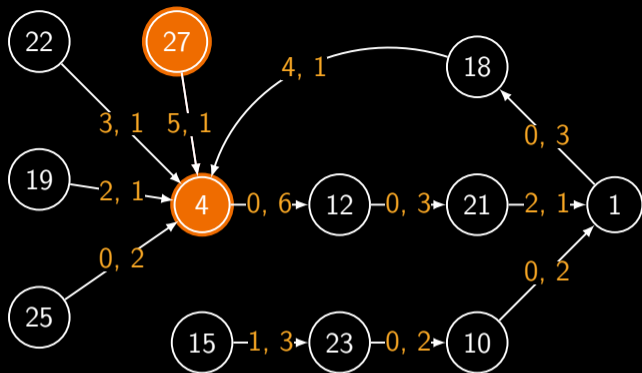
▀ $l_x := \mathcal{E}.succ(\mathcal{E}[x]) - \mathcal{E}[x]$

(sorted $\mathcal{E} =$

[1, 4, 10, 12, 15, 18, 19, 20, 21, 22, 23, 25, 27])

to compute $SA[i + 1]$
from $SA[i]$

1. start node is
 $p := \mathcal{E}.\text{pred}(SA[i])$
2. initial cost c_0 is
 $SA[i] - p$
3. stop if accumulated
cost c is above
limit of arc
4. add cost of arc to
 c and move to next
node v
5. report SA value
 $v + c$
6. goto 3.

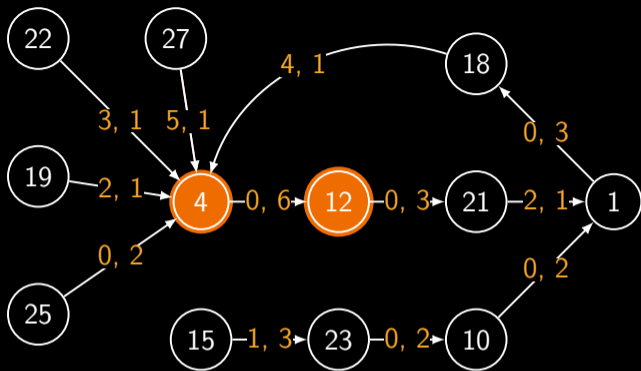


example

- ▀ $SA[1] = 27, p = 27, c \leftarrow c_0 = 0$
- ▀ arc $(27, 4)$ has limit $1 > c \Rightarrow$ move to 4
- ▀ add cost 5 of arc $(27, 4)$ to $c \leftarrow c_0 + 5 = 5$
- ▀ $SA[2] = 9$ is node label 4 plus cost 5

to compute $SA[i + 1]$
from $SA[i]$

1. start node is
 $p := \mathcal{E}.\text{pred}(SA[i])$
2. initial cost c_0 is
 $SA[i] - p$
3. stop if accumulated
cost c is above
limit of arc
4. add cost of arc to
 c and move to next
node v
5. report SA value
 $v + c$
6. goto 3.

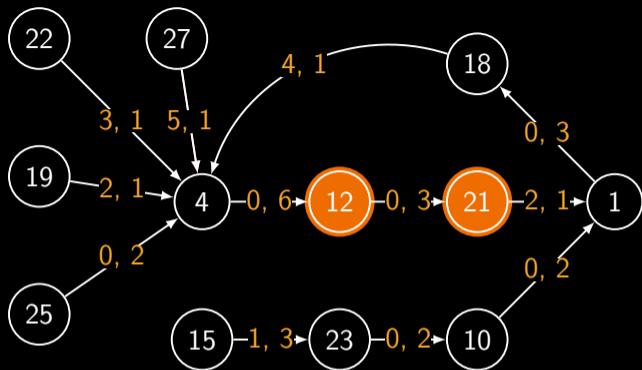


example

- ▀ are at node 4 with accumulated cost $c = 5$
- ▀ arc $(4, 12)$ has limit $6 > c \Rightarrow$ move to 12
- ▀ $SA[3] = 17$ is node label 12 plus cost 5

to compute $SA[i + 1]$
from $SA[i]$

1. start node is
 $p := \mathcal{E}.\text{pred}(SA[i])$
2. initial cost c_0 is
 $SA[i] - p$
3. stop if accumulated
cost c is above
limit of arc
4. add cost of arc to
 c and move to next
node v
5. report SA value
 $v + c$
6. goto 3.

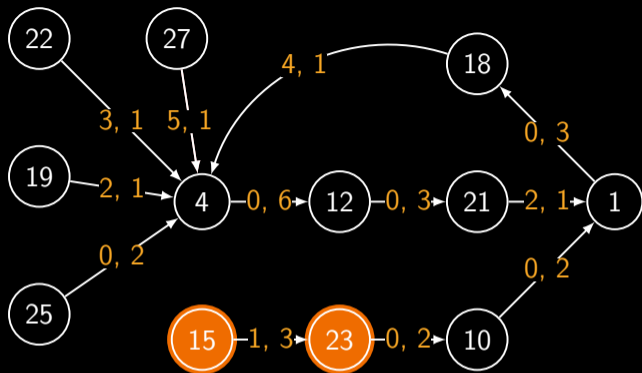


example

- ▀ are at node 12 with accumulated cost $c = 5$
- ▀ arc (12, 21) has limit $3 < c \Rightarrow$ stop
- ▀ continue with $SA[3] = 17$
- ▀ $p = \mathcal{E}.\text{pred}(17) = 15, c_0 = 2.$

to compute $SA[i + 1]$
from $SA[i]$

1. start node is
 $p := \mathcal{E}.\text{pred}(SA[i])$
2. initial cost c_0 is
 $SA[i] - p$
3. stop if accumulated
cost c is above
limit of arc
4. add cost of arc to
 c and move to next
node v
5. report SA value
 $v + c$
6. goto 3.

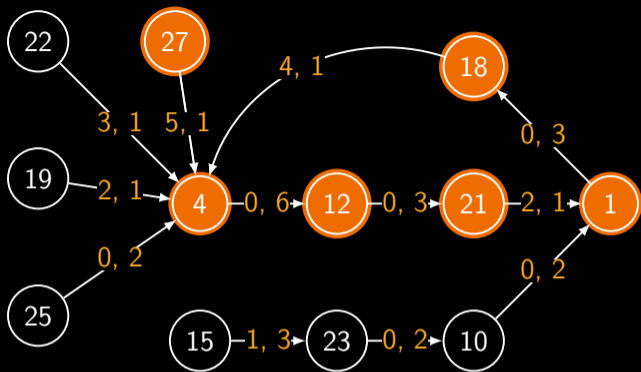


example

- ▀ $p = \mathcal{E}.\text{pred}(17) = 15, c \leftarrow c_0 = 2$
- ▀ arc (15, 23) has limit $3 > c \Rightarrow$ continue
- ▀ add cost 1 of arc (15, 23) to $c \leftarrow c_0 + 1 = 3$
- ▀ $SA[4] = 26$ is node label 23 plus cost 3

to compute $SA[i + 1]$
from $SA[i]$

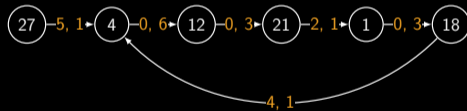
1. start node is
 $p := \mathcal{E}.\text{pred}(SA[i])$
2. initial cost c_0 is
 $SA[i] - p$
3. stop if accumulated
cost c is above
limit of arc
4. add cost of arc to
 c and move to next
node v
5. report SA value
 $v + c$
6. goto 3.



further speedup

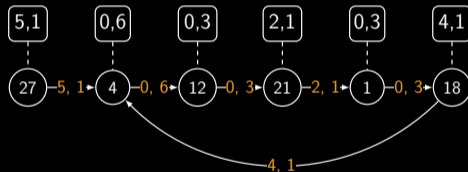
- ▀ for each traversed arc, we can omit a predecessor query on \mathcal{E}
 - ▀ if such traversable paths are long, we add shortcuts
- \Rightarrow build ϕ^{-1} trees on long paths

ϕ^{-1} tree



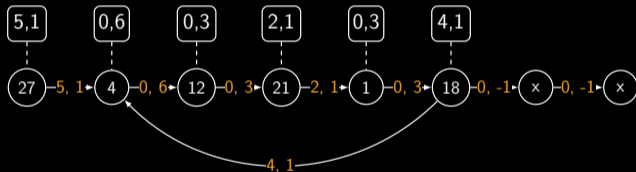
ϕ^{-1} tree

- identify label of out-going arc with node itself



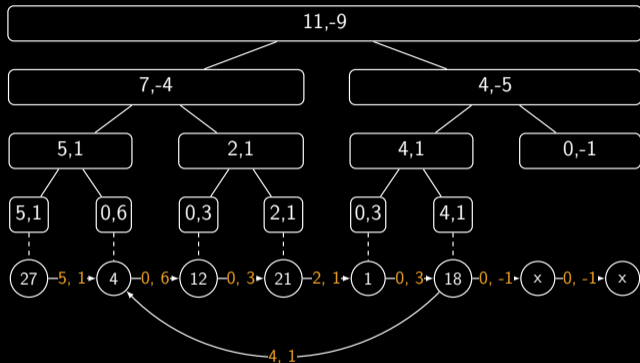
ϕ^{-1} tree

- identify label of out-going arc with node itself
- fill up path with dummy arcs having label 0, -1 (=untraversable)



ϕ^{-1} tree

- identify label of out-going arc with node itself
- fill up path with dummy arcs having label 0, -1 (=untraversable)
- build perfect binary tree on path by partitioning it



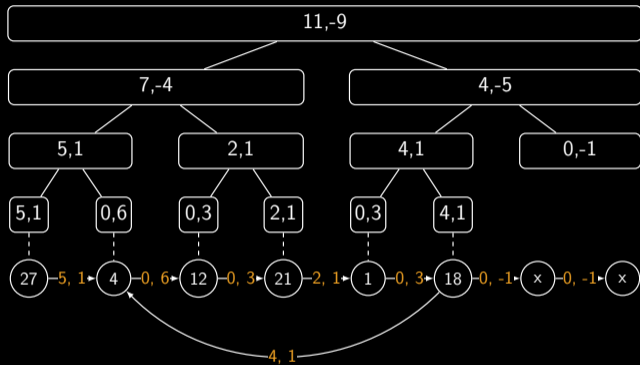
ϕ^{-1} tree

- identify label of out-going arc with node itself
- fill up path with dummy arcs having label $0, -1$ (=untraversable)
- build perfect binary tree on path by partitioning it
- label of internal node is based on the label of its children as below:

$$c_1 + c_2, \min(l_1, l_2 - c_1)$$

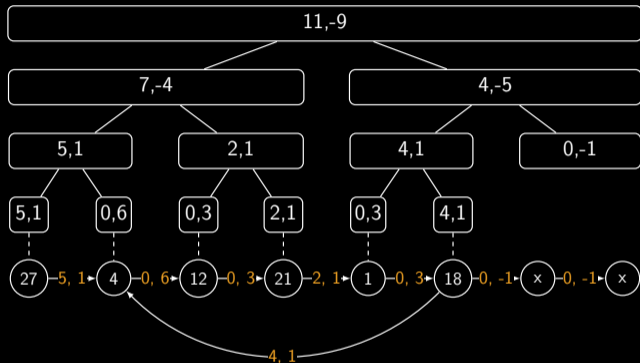
$$c_1, l_1$$

$$c_2, l_2$$



querying ϕ^{-1} tree

1. climb up until exceeding the limit at a node
2. climb down to leaf at which we exceed the limit the first time

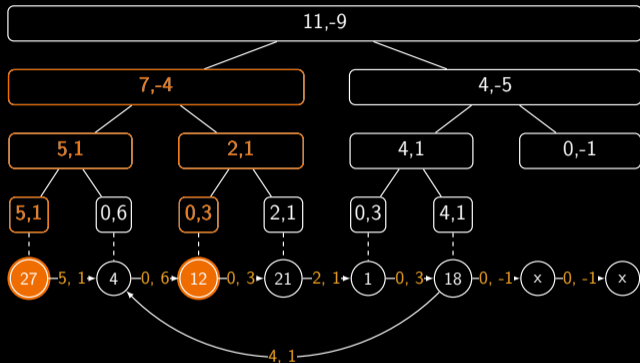


querying ϕ^{-1} tree

1. climb up until exceeding the limit at a node
2. climb down to leaf at which we exceed the limit the first time

example

1. start at 27 with cost 0
2. climb up to 7, -4 with limit $-4 < 0$
3. take cost of left child and descend to 0, 3
4. return 12 with cost 5, having skipped 4



experiments

dataset

- ▮ Chromosome 19 sequences from the 1000 Genomes Project

machine

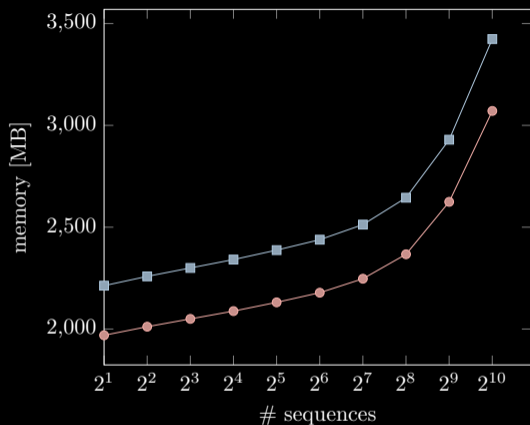
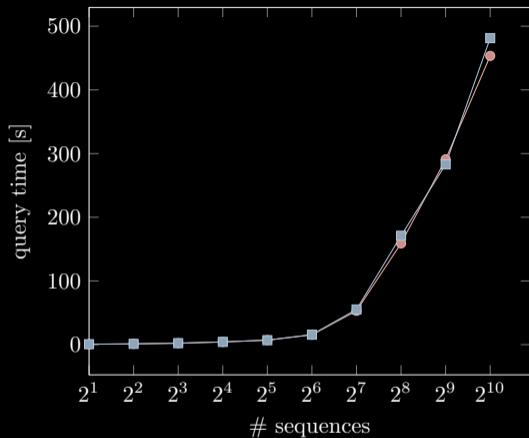
- ▮ AMD EPYC 75F3 32-core processor
- ▮ 512 GB of RAM
- ▮ 64-bit Linux

alternative solutions

- ▮ standard r -index, Gagie+'18
- ▮ sr -index, Cobas+'21
- ▮ RLCSA, Mäkinen+'10

Chromosome 19

—●— r-index —■— ϕ^{-1} forest



ϕ^{-1} forest has no clear advantage

conclusion

introduced ϕ^{-1} forest

- ▀ $\mathcal{O}(r)$ space data structure on top of r -index
- ▀ provides random access to SA
- ▀ query time depends on
 - length of a run
 - values of costs and limits

open problems

- ▀ only trivial bound known
 - $\mathcal{O}(\log r)$ time per ϕ^{-1} tree traversal
 - $\mathcal{O}(\log \log_w(n/r))$ time for a predecessor query
- ▀ need theoretical analysis of number of predecessor calls

Thank you for listening. Any questions are welcome!

i	SA	rotation matrix	BWT
1	27	#GATTACAT\$GATACAT\$GATTAGAT	A
2	9	\$GATACAT\$GATTAGATA#GATTACA	T
3		\$GATTAGATA#GATTACAT\$GATACA	T
4		A#GATTACAT\$GATACAT\$GATTAGA	T
5		ACAT\$GATACAT\$GATTAGATA#GAT	T
6		ACAT\$GATTAGATA#GATTACAT\$GA	T
7	22	AGATA#GATTACAT\$GATACAT\$GAT	T
8	7	AT\$GATACAT\$GATTAGATA#GATTA	C
9	15	AT\$GATTAGATA#GATTACAT\$GATA	C
10	24	ATA#GATTACAT\$GATACAT\$GATTA	G
11		ATACAT\$GATTAGATA#GATTACAT\$	G
12		ATTACAT\$GATACAT\$GATTAGATA#	G
13	19	ATTAGATA#GATTACAT\$GATACAT\$	G
14	6	CAT\$GATACAT\$GATTAGATA#GATT	A
15		CAT\$GATTAGATA#GATTACAT\$GAT	A
16	23	GATA#GATTACAT\$GATACAT\$GATT	A
17	10	GATACAT\$GATTAGATA#GATTACAT	\$
18	1	GATTACAT\$GATACAT\$GATTAGATA	#
19	18	GATTAGATA#GATTACAT\$GATACAT	\$
20	8	T\$GATACAT\$GATTAGATA#GATTAC	A
21		T\$GATTAGATA#GATTACAT\$GATAC	A
22	25	TA#GATTACAT\$GATACAT\$GATTAG	A
23	4	TACAT\$GATACAT\$GATTAGATA#GA	T
24	12	TACAT\$GATTAGATA#GATTACAT\$G	A
25	21	TAGATA#GATTACAT\$GATACAT\$GA	T
26	3	TTACAT\$GATACAT\$GATTAGATA#G	A
27	20	TTAGATA#GATTACAT\$GATACAT\$G	A

how to compute $SA[i + 1]$ from $SA[i]$ with the r -index:

- ▮ if $SA[i] = \mathcal{E}[x]$ for an $x \in [1..r]$
 $\Rightarrow SA[i + 1] = \mathcal{S}[x + 1]$
- ▮ otherwise, take $p = \mathcal{E}.\text{pred}(SA[i])$, and let j be such that $SA[j] = p$
- ▮ apply toehold lemma:
- ▮ $SA[i + 1] - SA[i] = SA[j + 1] - SA[j]$
- ▮ we obtain $SA[j + 1]$ from the above case ($SA[j] \in \mathcal{E}$)

ϕ^{-1}

p	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27
$T[p]$	G	A	T	T	A	C	A	T	\$	G	A	T	A	C	A	T	\$	G	A	T	T	A	G	A	T	A	#
$\phi^{-1}[p]$	18	19	20	12	13	14	15	16	17	1	2	21	22	23	24	25	26	8	6	27	3	7	10	11	4	5	9

computing $SA[i + 1]$ from $SA[i]$ is actually an application of

$$\phi^{-1}[SA[i]] := \begin{cases} 1 & \text{if } i = n, \\ SA[i + 1] & \text{otherwise} \end{cases}$$

then:

- ▀ $\phi^{-1}(\mathcal{E}[x]) = \mathcal{S}[x + 1]$
- ▀ take $p = \mathcal{E}.\text{pred}(SA[i])$, and let j be such that $SA[j] = p$
- ▀ rewrite $SA[i + 1] = SA[j + 1] + SA[i] - SA[j]$ as

$$\phi^{-1}(SA[i]) = SA[i + 1] = \phi^{-1}(p) + SA[i] - p$$

i	SA	rotation matrix	BWT
1	27	#GATTACAT\$GATACAT\$GATTAGAT	A
$i \rightarrow 2$	9	\$GATACAT\$GATTAGATA#GATTACA	T
3		\$GATTAGATA#GATTACAT\$GATACA	T
4		A#GATTACAT\$GATACAT\$GATTAGA	T
5		ACAT\$GATACAT\$GATTAGATA#GAT	T
6		ACAT\$GATTAGATA#GATTACAT\$GA	T
7	22	AGATA#GATTACAT\$GATACAT\$GAT	T
8	7	AT\$GATACAT\$GATTAGATA#GATTA	C
9	15	AT\$GATTAGATA#GATTACAT\$GATA	C
10	24	ATA#GATTACAT\$GATACAT\$GATTA	G
11		ATACAT\$GATTAGATA#GATTACAT\$	G
12		ATTACAT\$GATACAT\$GATTAGATA#	G
13	19	ATTAGATA#GATTACAT\$GATACAT\$	G
14	6	CAT\$GATACAT\$GATTAGATA#GATT	A
15		CAT\$GATTAGATA#GATTACAT\$GAT	A
16	23	GATA#GATTACAT\$GATACAT\$GATT	A
17	10	GATACAT\$GATTAGATA#GATTACAT	\$
18	1	GATTACAT\$GATACAT\$GATTAGATA	#
19	18	GATTAGATA#GATTACAT\$GATACAT	\$
20	8	T\$GATACAT\$GATTAGATA#GATTAC	A
21		T\$GATTAGATA#GATTACAT\$GATAC	A
22	25	TA#GATTACAT\$GATACAT\$GATTAG	A
23	4	TACAT\$GATACAT\$GATTAGATA#GA	T
$\phi^{-1}(4) \rightarrow 24$	12	TACAT\$GATTAGATA#GATTACAT\$G	A
25	21	TAGATA#GATTACAT\$GATACAT\$GA	T
26	3	TTACAT\$GATACAT\$GATTAGATA#G	A
27	20	TTAGATA#GATTACAT\$GATACAT\$G	A

from SA[i] to SA[$i + 1$]

- ▀ $\phi^{-1}(\mathcal{E}[x]) = \mathcal{S}[x + 1]$
- ▀ take $p = \mathcal{E}.\text{pred}(\text{SA}[i])$, and let j be such that $\text{SA}[j] = p$
- ▀ $\phi^{-1}(\text{SA}[i]) = \phi^{-1}(p) + \text{SA}[i] - p$

example

- ▀ $i = 2$, SA[2] = 9 is known
 - ▀ $4 = \mathcal{E}.\text{pred}(9)$
 - ▀ $\phi^{-1}(4) = 12$
- $\Rightarrow \text{SA}[3] = 12 + 9 - 4 = 17$

(sorted

$\mathcal{E} = [1, 4, 10, 12, 15, 18, 19, 20, 21, 22, 23, 25, 27])$

i	SA	rotation matrix	BWT
1	27	#GATTACAT\$GATACAT\$GATTAGAT	A
2	9	\$GATACAT\$GATTAGATA#GATTACA	T
$i \rightarrow 3$	17	\$GATTAGATA#GATTACAT\$GATACA	T
4		A#GATTACAT\$GATACAT\$GATTAGA	T
5		ACAT\$GATACAT\$GATTAGATA#GAT	T
6		ACAT\$GATTAGATA#GATTACAT\$GA	T
7	22	AGATA#GATTACAT\$GATACAT\$GAT	T
8	7	AT\$GATACAT\$GATTAGATA#GATTA	C
9	15	AT\$GATTAGATA#GATTACAT\$GATA	C
$\phi^{-1}(15) \rightarrow 10$	24	ATA#GATTACAT\$GATACAT\$GATTA	G
11		ATACAT\$GATTAGATA#GATTACAT\$	G
12		ATTACAT\$GATACAT\$GATTAGATA#	G
13	19	ATTAGATA#GATTACAT\$GATACAT\$	G
14	6	CAT\$GATACAT\$GATTAGATA#GATT	A
15		CAT\$GATTAGATA#GATTACAT\$GAT	A
16	23	GATA#GATTACAT\$GATACAT\$GATT	A
17	10	GATACAT\$GATTAGATA#GATTACAT	\$
18	1	GATTACAT\$GATACAT\$GATTAGATA	#
19	18	GATTAGATA#GATTACAT\$GATACAT	\$
20	8	T\$GATACAT\$GATTAGATA#GATTAC	A
21		T\$GATTAGATA#GATTACAT\$GATAC	A
22	25	TA#GATTACAT\$GATACAT\$GATTAG	A
23	4	TACAT\$GATACAT\$GATTAGATA#GA	T
24	12	TACAT\$GATTAGATA#GATTACAT\$G	A
25	21	TAGATA#GATTACAT\$GATACAT\$GA	T
26	3	TTACAT\$GATACAT\$GATTAGATA#G	A
27	20	TTAGATA#GATTACAT\$GATACAT\$G	A

from $SA[i]$ to $SA[i + 1]$

- ▀ $\phi^{-1}(\mathcal{E}[x]) = \mathcal{S}[x + 1]$
- ▀ take $p = \mathcal{E}.\text{pred}(SA[i])$, and let j be such that $SA[j] = p$
- ▀ $\phi^{-1}(SA[i]) = \phi^{-1}(p) + SA[i] - p$

example

- ▀ $i = 3$, $SA[3] = 17$ is known
 - ▀ $15 = \mathcal{E}.\text{pred}(17)$
 - ▀ $\phi^{-1}(15) = 24$
- $\Rightarrow SA[4] = 24 + 17 - 15 = 26$

(sorted

$\mathcal{E} = [1, 4, 10, 12, 15, 18, 19, 20, 21, 22, 23, 25, 27])$

i	SA	rotation matrix	BWT
1	27	#GATTACAT\$GATACAT\$GATTAGAT	A
2	9	\$GATACAT\$GATTAGATA#GATTACA	T
3	17	\$GATTAGATA#GATTACAT\$GATACA	T
4		A#GATTACAT\$GATACAT\$GATTAGA	T
5		ACAT\$GATACAT\$GATTAGATA#GAT	T
6		ACAT\$GATTAGATA#GATTACAT\$GA	T
7	22	AGATA#GATTACAT\$GATACAT\$GAT	T
8	7	AT\$GATACAT\$GATTAGATA#GATTA	C
9	15	AT\$GATTAGATA#GATTACAT\$GATA	C
10	24	ATA#GATTACAT\$GATACAT\$GATTA	G
11		ATACAT\$GATTAGATA#GATTACAT\$	G
12		ATTACAT\$GATACAT\$GATTAGATA#	G
13	19	ATTAGATA#GATTACAT\$GATACAT\$	G
14	6	CAT\$GATACAT\$GATTAGATA#GATT	A
15		CAT\$GATTAGATA#GATTACAT\$GAT	A
16	23	GATA#GATTACAT\$GATACAT\$GATT	A
17	10	GATACAT\$GATTAGATA#GATTACAT	\$
18	1	GATTACAT\$GATACAT\$GATTAGATA	#
19	18	GATTAGATA#GATTACAT\$GATACAT	\$
20	8	T\$GATACAT\$GATTAGATA#GATTAC	A
21		T\$GATTAGATA#GATTACAT\$GATAC	A
22	25	TA#GATTACAT\$GATACAT\$GATTAG	A
23	4	TACAT\$GATACAT\$GATTAGATA#GA	T
24	12	TACAT\$GATTAGATA#GATTACAT\$G	A
25	21	TAGATA#GATTACAT\$GATACAT\$GA	T
26	3	TTACAT\$GATACAT\$GATTAGATA#G	A
27	20	TTAGATA#GATTACAT\$GATACAT\$G	A

Observation

- ▀ # predecessor queries is bounded by length of run
- ⇒ practically slow for long runs

Can we compute m iterations of ϕ^{-1} faster?

i	SA	rotation matrix	BWT
1	27	#GATTACAT\$GATACAT\$GATTAGAT	A
2	9	\$GATACAT\$GATTAGATA#GATTACA	T
3	17	\$GATTAGATA#GATTACAT\$GATACA	T
4		A#GATTACAT\$GATACAT\$GATTAGA	T
5		ACAT\$GATACAT\$GATTAGATA#GAT	T
6		ACAT\$GATTAGATA#GATTACAT\$GA	T
7	22	AGATA#GATTACAT\$GATACAT\$GAT	T
8	7	AT\$GATACAT\$GATTAGATA#GATTA	C
9	15	AT\$GATTAGATA#GATTACAT\$GATA	C
10	24	ATA#GATTACAT\$GATACAT\$GATTA	G
11		ATACAT\$GATTAGATA#GATTACAT\$	G
12		ATTACAT\$GATACAT\$GATTAGATA#	G
13	19	ATTAGATA#GATTACAT\$GATACAT\$	G
14	6	CAT\$GATACAT\$GATTAGATA#GATT	A
15		CAT\$GATTAGATA#GATTACAT\$GAT	A
16	23	GATA#GATTACAT\$GATACAT\$GATT	A
17	10	GATACAT\$GATTAGATA#GATTACAT	\$
18	1	GATTACAT\$GATACAT\$GATTAGATA	#
19	18	GATTAGATA#GATTACAT\$GATACAT	\$
20	8	T\$GATACAT\$GATTAGATA#GATTAC	A
21		T\$GATTAGATA#GATTACAT\$GATAC	A
22	25	TA#GATTACAT\$GATACAT\$GATTAG	A
23	4	TACAT\$GATACAT\$GATTAGATA#GA	T
24	12	TACAT\$GATTAGATA#GATTACAT\$G	A
25	21	TAGATA#GATTACAT\$GATACAT\$GA	T
26	3	TTACAT\$GATACAT\$GATTAGATA#G	A
27	20	TTAGATA#GATTACAT\$GATACAT\$G	A

suppose $SA[i] = \mathcal{E}[x]$, and let

$$\blacktriangleright \mathcal{E}[x] := \mathcal{E}.\text{pred}(SA[i])$$

$$\blacktriangleright \mathcal{E}[y] := \mathcal{E}.\text{pred}(SA[i + 1])$$

$$\begin{aligned} \text{then: } \phi^{-1}(SA[i]) &= SA[i + 1] = \mathcal{S}[x + 1] \\ &= \mathcal{E}[y] + \mathcal{S}[x + 1] - \mathcal{E}[y] \\ &= \mathcal{E}[y] + c_x \end{aligned}$$

where $c_x := \mathcal{S}[x + 1] - \mathcal{E}.\text{pred}(\mathcal{S}[x + 1])$ is the cost of x -th run.

example

$$\blacktriangleright i = 1, SA[i] = \mathcal{E}[1] = 27,$$

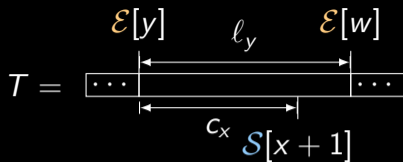
$$\blacktriangleright \mathcal{S}[2] = 9, \mathcal{E}[y] = \mathcal{E}.\text{pred}(\mathcal{S}[2]) = 4$$

$$\blacktriangleright c_1 = \mathcal{S}[2] - \mathcal{E}[y] = 5; \phi^{-1}(27) = 4 + 5 = 9$$

suppose $SA[i + 1] \notin \mathcal{E}$.

Then:

$$\begin{aligned}\phi^{-1}(SA[i + 1]) &= \phi^{-1}(\phi^{-1}(SA[i])) \\ &= \phi^{-1}(\mathcal{S}[x + 1]) \\ &= \phi^{-1}(\mathcal{E}[y] + c_x)\end{aligned}$$



▀ apply toehold lemma for $\mathcal{E}.\text{pred}(\mathcal{S}[x + 1]) = \mathcal{E}[y]$:

$$\Rightarrow \phi^{-1}(\mathcal{E}[y] + c_x) = \phi^{-1}(\mathcal{E}[y]) + c_x$$

▀ finally, $\phi^{-1}(\mathcal{E}[y]) = \mathcal{S}[y + 1] = \mathcal{E}[z] + c_y$,
where $\mathcal{E}[z] := \mathcal{E}.\text{pred}(\mathcal{S}[y + 1])$

▀ total: $\phi^{-1}(SA[i + 1]) = \mathcal{E}[z] + c_x + c_y$

given $\mathcal{E}[w] := \mathcal{E}.\text{succ}(\mathcal{E}[y])$, we assumed that

▀ $c_x = \mathcal{S}[x + 1] - \mathcal{E}[y] < \mathcal{E}[w] - \mathcal{E}[y] = l_y$,

where $l_y := \mathcal{E}.\text{succ}(\mathcal{E}[y]) - \mathcal{E}[y]$ is the **limit** of the y -th run

recursive application

- ▶ given $SA[i]$ with $\mathcal{E}.\text{pred}(SA[i]) = \mathcal{E}[x_1]$,
- ▶ let $c_0 := SA[i] - \mathcal{E}[x_1]$ and
- ▶ x_1, x_2, \dots, x_m be the indices of the runs we visit, such that $\mathcal{E}.\text{pred}(\mathcal{E}[x_j] + \sum_{k=1}^{j-1} c_{x_k} + c_0) = \mathcal{E}[x_{j+1}]$ for all $j \in [2..m-1]$

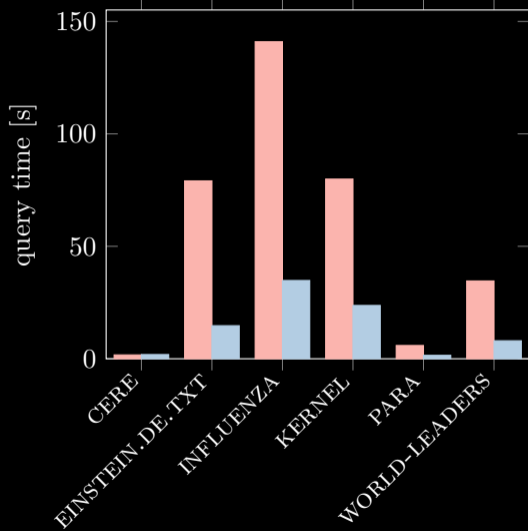
then:

$$\phi^{-m}(SA[i]) = \mathcal{S}[x_m + 1] = \mathcal{E}[x_m] + \sum_{k=1}^{m-1} c_{x_k} + c_0$$

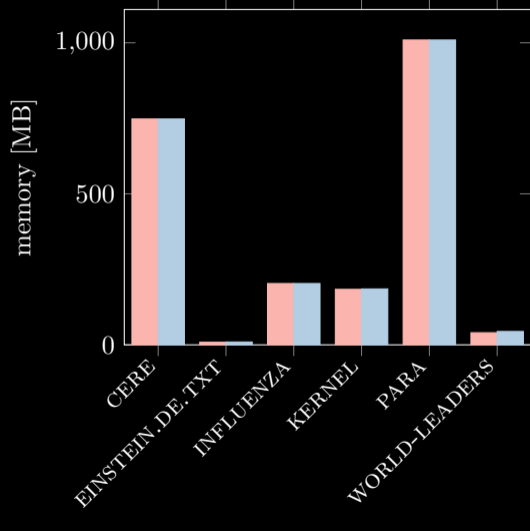
conclusion:

- ▶ just need to sum up $\sum_k c_{x_k}$
- ▶ but check that sum does not exceed the limit
- ▶ can translate this to a path problem on a directed labeled graph

Pizza&Chili experiments



■ r-index ■ ϕ^{-1} forest



ϕ^{-1} forest faster with negligible memory overhead