


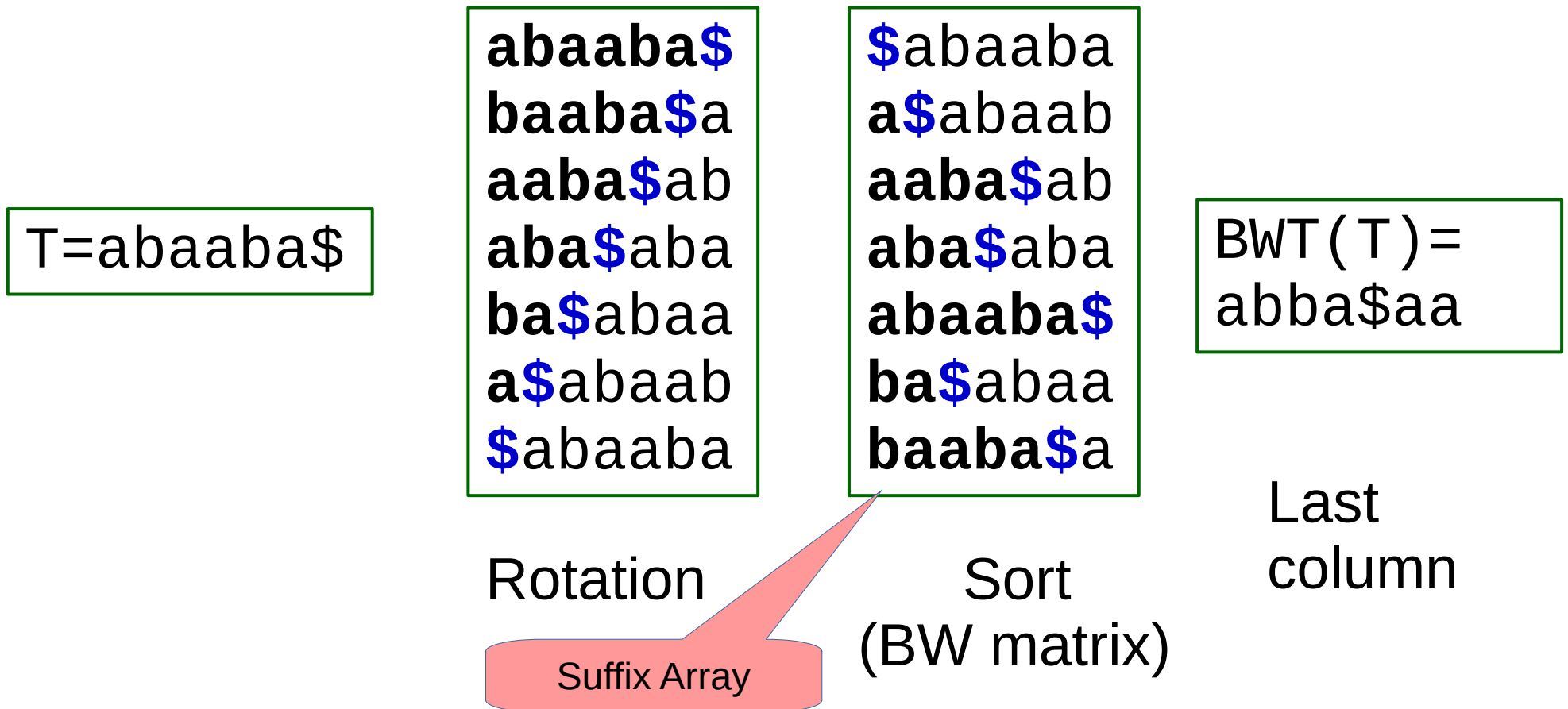
Exercises

- 1) For each mapping approach:
 - i) Why not use the naive, brute force sliding window approach? **time**
 - ii) Why not use hashes? **memory, variable length strings**
 - iii) Why not use suffix trees? **memory**
 - iv) Why not use suffix arrays? **memory**
- 2) Why are data structures important for the string search problem? **affects time and memory**
- 3) Which data structure would you use to search for substrings of varying lengths? **suffix trees or suffix arrays, depending on memory constraints**

Today's Objectives

- Evaluating strategies for short read alignment
 - String search
 - Linear (time constraints)
 - Hash
 - Suffix tree
 - Suffix array
 - Burrows Wheeler Transform (BWT)
 - FM index search of BWT
- 
- memory constraints

Burrows-Wheeler Transform (BWT)



BWT allows **Lossless** compression - the original data can be perfectly reconstructed from the compressed data.

BWT is **reversible**

'Compressed Suffix Array'

BWT enables compression

How is it compressed?

```
BWT("Tomorrow_and_tomorrow_and_tomorrow$")  
='w$wwdd__nnooaaattTmmmrrrrrrrooo__ooo'
```

BWT (also called block-sorting compression) rearranges a character string into runs of similar characters. This is useful for compression, since it tends to be easy to compress a string that has **runs of repeated characters** by techniques such as move-to-front transform and run-length encoding.

The transformation is **reversible**. **Lossless** compression is a class of data compression algorithms that allows the original data to be perfectly reconstructed from the compressed data. Lossy compression only allows approximate reconstruction of the data

Lossless compression algorithms: ZIP, GNU tool gzip, PNG and GIF

Compression

Run length encoding

runs of data are stored as a single data value and count

AAAAAATTTTTTTCTTTTTTGGGG

6A7T1C6T4G

Arithmetic coding

string of characters such as the word "java" is represented using a fixed number of bits per character

j = 01, a = 11, v = 10

java = 01111011

1 character = 8 bits (1 byte)

By using smaller bit codes for common words/characters, we can get compression.

Huffman coding

| Char | Code | Freq | Bits |
|------|-------|------|------|
| E | 10 | 15 | 30 |
| P | 01 | 13 | 26 |
| A | 110 | 10 | 30 |
| S | 11111 | 3 | 15 |

BWT

How is it an index, or how do we reverse it?

\$ a b a a b a
a \$ a b a a b
a a b a \$ a b
a b a \$ a b a
a b a a b a \$
b a \$ a b a a
b a a b a \$ a

BWM(T)

| | |
|---|----------------|
| 6 | \$ |
| 5 | a \$ |
| 2 | a a b a \$ |
| 3 | a b a \$ |
| 0 | a b a a b a \$ |
| 4 | b a \$ |
| 1 | b a a b a \$ |

SA(T)

T=abaaba\$

BWT(T)=
abba\$a

\$abaaba
a\$abaab
aaba\$ab
aba\$aba
abaaba\$
ba\$abaa
baaba\$a

Reversing BWT(T)

a
b
b
a
\$
a
a

BWT(T) Sort

\$
a
a
a
a
b
b

Prepend
BWT(T)

a\$
ba
ba
aa
\$a
ab
ab

\$a
a\$
aa
ab
ab
ba
ba

Sort

a\$a
ba\$
baa
aab
\$ab
aba
aba

Prepend
BWT(T)

\$aba
a\$ab
aaba
aba\$
abaa
ba\$a
baab

Sort

a\$aba
ba\$ab
baaba
aaba\$
\$abaa
aba\$a
abaab

Prepend

\$abaa
a\$aba
aaba\$
aba\$a
abaab
ba\$ab
baaba

Sort

a\$abaa
ba\$aba
baaba\$
aaba\$a
\$abaab
aba\$ab
abaaba

Prepend

\$abaab
a\$abaa
aaba\$a
aba\$ab
abaaba
ba\$aba
baaba\$

Sort

a\$abaab
ba\$abaa
baaba\$a
aaba\$ab
\$abaaba
aba\$aba
abaaba\$

Prepend

\$abaaba
a\$abaab
aaba\$ab
aba\$aba
abaaba\$
ba\$abaa
baaba\$a

Sort

T=abaaba\$

BWT(T)=
abba\$a

Seems like a lot of
computing, is there an
easier way?

\$abaaba
a\$abaab
aaba\$ab
aba\$aba
abaaba\$
ba\$abaa
baaba\$a

Reversing BWT with FM index and LF mapping

FM index (Full test Minute space)

T-ranking

Give each character in T a rank, equal to # times the character occurred previously in T. (see structure)

a₀ b₀ a₁ a₂ b₁ a₃ \$

B-ranking

Give each character in BWT(T) a rank, equal to # times the character occurred. (Use for LF mapping)

T-ranking

T-ranking

Give each character in T a rank, equal to # times the character occurred previously in T.

a₀ b₀ a₁ a₂ b₁ a₃ \$

| | | | | | | | |
|----------------|----|----|----|----|----|----|--------------|
| \$ | a | b | a | a | b | a | ₃ |
| a ₃ | \$ | a | b | a | a | b | ₁ |
| a ₁ | a | b | a | \$ | a | b | ₀ |
| a ₂ | b | a | \$ | a | b | a | ₁ |
| a ₀ | b | a | a | b | a | \$ | |
| b ₁ | a | \$ | a | b | a | a | ₂ |
| b ₀ | a | a | b | a | \$ | a | ₀ |

BWM with T-ranking

\$ a b a a b a
 a **\$** a b a a b
 a a b a **\$** a b
 a b a **\$** a b a
 a b a a b a **\$**
 b a **\$** a b a a
 b a a b a **\$** a

BWM(T)

First

Last

F

L

| | | | | | | |
|----------------------|----------------|----------------|----------------|----------------|----------------|----------------------|
| \$ | a ₀ | b ₀ | a ₁ | a ₂ | b ₁ | a₃ |
| a₃ | \$ | a ₀ | b ₀ | a ₁ | a ₂ | b ₁ |
| a₁ | a ₂ | b ₁ | a ₃ | \$ | a ₀ | b ₀ |
| a₂ | b ₁ | a ₃ | \$ | a ₀ | b ₀ | a₁ |
| a₀ | b ₀ | a ₁ | a ₂ | b ₁ | a ₃ | \$ |
| b ₁ | a ₃ | \$ | a ₀ | b ₀ | a ₁ | a₂ |
| b ₀ | a ₁ | a ₂ | b ₁ | a ₃ | \$ | a₀ |

BWM(T) with T-ranking
 Order does not change!

BWM with B-ranking

B-ranking

Give each character in BWT(T) a rank, equal to # times the character occurred.

| <i>F</i> | | | | | | | <i>L</i> |
|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------|
| \$ | a ₀ | b ₀ | a ₁ | a ₂ | b ₁ | a ₃ | |
| a ₃ | \$ | a ₀ | b ₀ | a ₁ | a ₂ | b ₁ | |
| a ₁ | a ₂ | b ₁ | a ₃ | \$ | a ₀ | b ₀ | |
| a ₂ | b ₁ | a ₃ | \$ | a ₀ | b ₀ | a ₁ | |
| a ₀ | b ₀ | a ₁ | a ₂ | b ₁ | a ₃ | \$ | |
| b ₁ | a ₃ | \$ | a ₀ | b ₀ | a ₁ | a ₂ | |
| b ₀ | a ₁ | a ₂ | b ₁ | a ₃ | \$ | a ₀ | |

BWM(T) with T-ranking
Order does not change!

| <i>F</i> | | | | | | | <i>L</i> |
|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------|
| \$ | a ₃ | b ₁ | a ₁ | a ₂ | b ₀ | a ₀ | |
| a ₀ | \$ | a ₃ | b ₁ | a ₁ | a ₂ | b ₀ | |
| a ₁ | a ₂ | b ₀ | a ₃ | \$ | a ₃ | b ₁ | |
| a ₂ | b ₀ | a ₀ | \$ | a ₃ | b ₁ | a ₁ | |
| a ₃ | b ₁ | a ₁ | a ₂ | b ₀ | a ₀ | \$ | |
| b ₀ | a ₀ | \$ | a ₃ | b ₁ | a ₁ | a ₂ | |
| b ₁ | a ₁ | a ₂ | b ₀ | a ₀ | \$ | a ₃ | |

Ascending rank

BWM(T) with B-ranking
Order of a_i and b_i does
not change!

Last First (LF) mapping B-ranking

| <i>F</i> | | | | | | | <i>L</i> |
|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------|
| \$ | a ₃ | b ₁ | a ₁ | a ₂ | b ₀ | a ₀ | |
| a ₀ | \$ | a ₃ | b ₁ | a ₁ | a ₂ | b ₀ | |
| a ₁ | a ₂ | b ₀ | a ₃ | \$ | a ₃ | b ₁ | |
| a ₂ | b ₀ | a ₀ | \$ | a ₃ | b ₁ | a ₁ | |
| a ₃ | b ₁ | a ₁ | a ₂ | b ₀ | a ₀ | \$ | |
| b ₀ | a ₀ | \$ | a ₃ | b ₁ | a ₁ | a ₂ | |
| b ₁ | a ₁ | a ₂ | b ₀ | a ₀ | \$ | a ₃ | |

BWM(T) with B-ranking

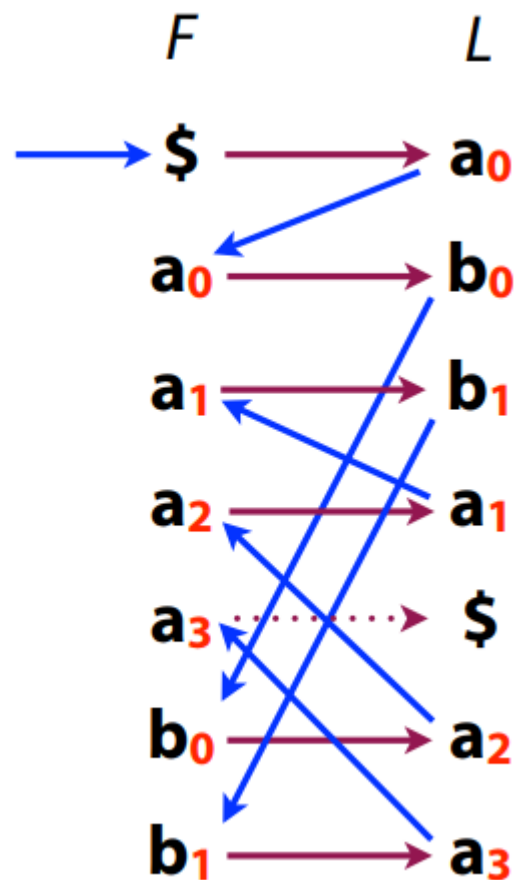
F can be represented
by n integers, n =
alphabet size
4 (a), 2 (b)

- To start at the end (\$) we go to the first position in (F).
- What is before \$? Go to same position in (L) a₀
- What is before a₀? Find the row that begins with a₀. Skip \$ when counting.
- b₀ is before a₀
- What is before b₀? Skip \$ (1), a (4)
- a₂ is before b₀

| <i>F</i> | <i>L</i> | |
|----------------|----------------|---|
| \$ | a ₀ | |
| a ₀ | b ₀ | |
| a ₁ | b ₁ | ← Which BWM row <i>begins</i> with b ₁ ? |
| a ₂ | a ₁ | Skip row starting with \$ (1 row) |
| a ₃ | \$ | Skip rows starting with a (4 rows) |
| b ₀ | a ₂ | Skip row starting with b ₀ (1 row) |
| b ₁ | a ₃ | Answer: row 6 |

row 6 →

Reverse BWT



Given *L*, and positions of *F* rows, we can reverse BW(*T*) to *T*.

Reverse of chars we visited = **a**₃ **b**₁ **a**₁ **a**₂ **b**₀ **a**₀ \$ = *T*

Storage

F
L
\$ a b a a b a
a \$ a b a a b
a a b a \$ a b
a b a \$ a b a
a b a a b a \$
b a \$ a b a a
b a a b a \$ a

| | |
|---|----------------|
| 6 | \$ |
| 5 | a \$ |
| 2 | a a b a \$ |
| 3 | a b a \$ |
| 0 | a b a a b a \$ |
| 4 | b a \$ |
| 1 | b a a b a \$ |

FM-index is a compressed full-text substring index based on the Burrows-Wheeler transform

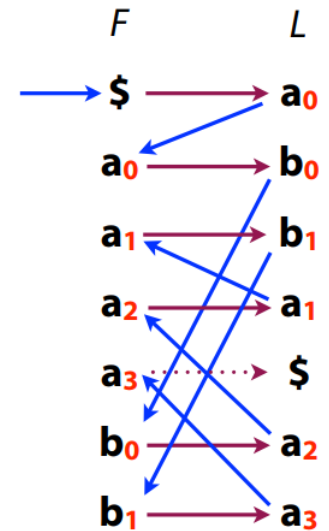
- F can be represented by n integers, n = alphabet size
- L can be compressed

Suffix Array

- Array of integers
- Array of characters
- no compression

Exercises

- 1) Construct BWT(T), with $T = \text{'CATGCAT'}$
- 2) Draw out LF mapping to show how to reverse BWT(T) using arrows and LF columns as shown on the right:
- 3) Given $\text{BWT}(T) = \text{ATGCA\$A}$, what is F, draw LF mapping and give T.



FM index search of BWT

FM index Query

Search the reverse of the string

$P = \mathbf{ab}\mathbf{a}$

Easy to find all the
rows beginning with
a, thanks to F 's
simple structure

| F | | | | | | | L |
|-----------------------|----|----|----|----|----|----|-----------------------|
| \$ | a | b | a | a | b | a | a ₃ |
| a ₀ | \$ | a | b | a | a | b | b ₁ |
| a ₁ | a | b | a | \$ | a | b | b ₀ |
| a ₂ | b | a | \$ | a | b | a | a ₁ |
| a ₃ | b | a | a | b | a | \$ | |
| b ₀ | a | \$ | a | b | a | a | a ₂ |
| b ₁ | a | a | b | a | \$ | a | a ₀ |

FM index Query

$P = \mathbf{aba}$

| | F | | L |
|-------|-----|-----------|-------|
| | \$ | a b a a b | a_0 |
| a_0 | \$ | a b a a | b_0 |
| a_1 | a | b a \$ a | b_1 |
| a_2 | b | a \$ a b | a_1 |
| a_3 | b | a a b a | \$ |
| b_0 | a | \$ a b a | a_2 |
| b_1 | a | a b a \$ | a_3 |

Look at those rows in L .

b₀, b₁ are **b**s occuring just to left.

Use LF Mapping. Let new range delimit those **bs**

$P = \mathbf{a} \mathbf{b} \mathbf{a}$

| | F | | L |
|----------------------|-----|-----------|-------------------------|
| | \$ | a b a a b | a₀ |
| a₀ | \$ | a b a a | b₀ |
| a₁ | a | b a \$ a | b₁ |
| a₂ | b | a \$ a b | a₁ |
| a₃ | b | a a b a | \$ |
| b₀ | a | \$ a b a | a₂ |
| b₁ | a | a a b a | \$ a₃ |

Now we have the rows with prefix **ba**

FM index Query

$P = \mathbf{aba}$

| F | | | | | | | L |
|-------|----|----|----|----|----|-------|-----|
| \$ | a | b | a | a | b | a_0 | |
| a_0 | \$ | a | b | a | a | b_0 | |
| a_1 | a | b | a | \$ | a | b_1 | |
| a_2 | b | a | \$ | a | b | a_1 | |
| a_3 | b | a | a | b | a | \$ | |
| b_0 | a | \$ | a | b | a | a_2 | |
| b_1 | a | a | b | a | \$ | a_3 | |

← a_2, a_3 occur just to left.

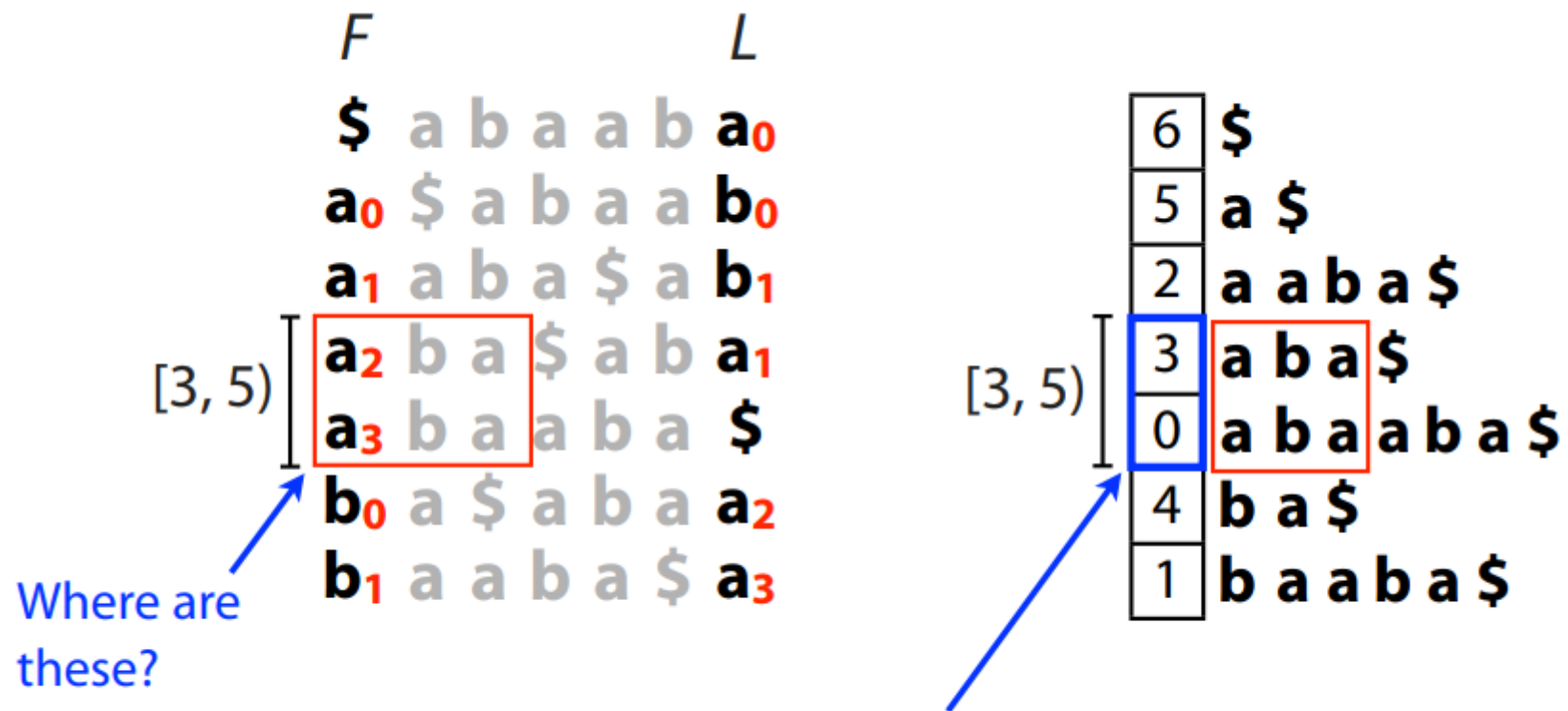
$P = \mathbf{aba}$

Use LF Mapping →

| F | | | | | | | L |
|-------|----|----|----|----|----|-------|-----|
| \$ | a | b | a | a | b | a_0 | |
| a_0 | \$ | a | b | a | a | b_0 | |
| a_1 | a | b | a | \$ | a | b_1 | |
| a_2 | b | a | \$ | a | b | a_1 | |
| a_3 | b | a | a | b | a | \$ | |
| b_0 | a | \$ | a | b | a | a_2 | |
| b_1 | a | a | b | a | \$ | a_3 | |

Now we have the rows with prefix **aba**

Location of matches



Issues to resolve

Scanning for preceding character is slow

| | | | | | | |
|----------------|----|----|----|----|----|----------------|
| \$ | a | b | a | a | b | a ₀ |
| a ₀ | \$ | a | b | a | a | b ₀ |
| a ₁ | a | b | a | \$ | a | b ₁ |
| a ₂ | b | a | \$ | a | b | a ₁ |
| a ₃ | b | a | a | b | a | \$ |
| b ₀ | a | \$ | a | b | a | a ₂ |
| b ₁ | a | a | b | a | \$ | a ₃ |

$O(m)$ scan

Need way to find where matches occur in T :

Where?

| | | | | | | |
|----------------|----|----|----|----|----|----------------|
| \$ | a | b | a | a | b | a ₀ |
| a ₀ | \$ | a | b | a | a | b ₀ |
| a ₁ | a | b | a | \$ | a | b ₁ |
| a ₂ | b | a | \$ | a | b | a ₁ |
| a ₃ | b | a | a | b | a | \$ |
| b ₀ | a | \$ | a | b | a | a ₂ |
| b ₁ | a | a | b | a | \$ | a ₃ |

Storing ranks takes too much space

Pre-tally ranks

Idea: pre-calculate # **a**s,
bs in L up to every row:

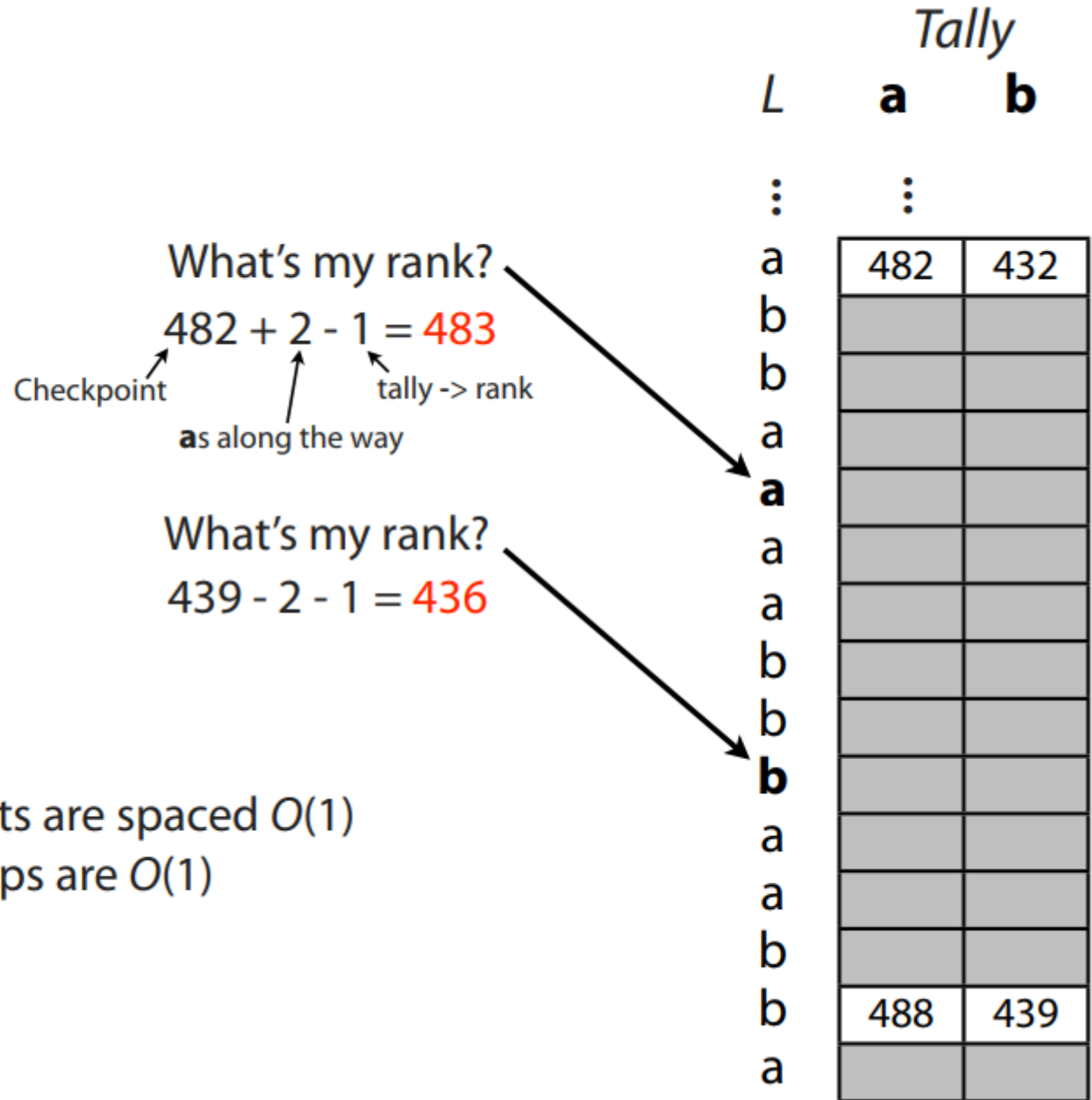
| | | <i>Tally</i> | |
|----------|----------|--------------|----------|
| F | L | a | b |
| \$ | a | 1 | 0 |
| a | b | 1 | 1 |
| a | b | 1 | 2 |
| a | a | 2 | 2 |
| a | \$ | 2 | 2 |
| b | a | 3 | 2 |
| b | a | 4 | 2 |

We infer **b**₀ and **b**₁
appear in L in this range

$O(1)$ time, but requires
 $m \times |\Sigma|$ integers

Exchange memory for speed using checkpoints

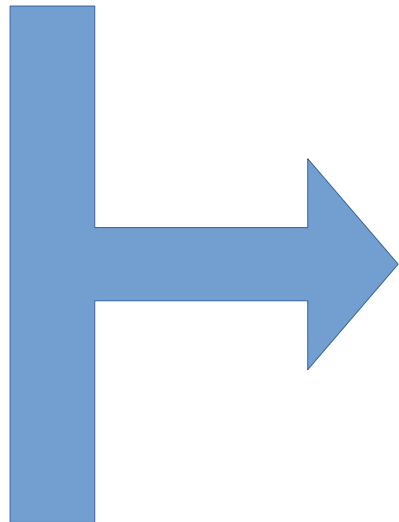
| | | <i>Tally</i> | | |
|----------|----------|--------------|----------|---------------------------------|
| <i>F</i> | <i>L</i> | a | b | |
| \$ | a | 1 | 0 | ← Lookup here succeeds as usual |
| a | b | | | |
| a | b | | | |
| a | a | | | |
| a | \$ | | | ← Oops: not a checkpoint |
| b | a | 3 | 2 | ← But there's one nearby |
| b | a | | | |



Assuming checkpoints are spaced $O(1)$
distance apart, lookups are $O(1)$

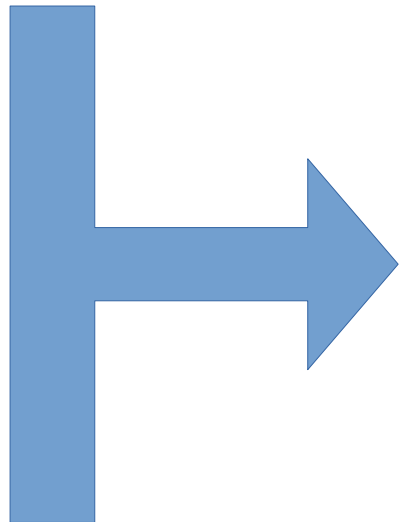
| | |
|---------|------|
| letter: | aaaa |
| Tally: | 1234 |
| Rank: | 0123 |

| <i>L</i> | <i>Tally</i> | |
|----------|--------------|----------|
| | a | b |
| : | : | |
| a | 482 | 432 |
| b | | |
| b | | |
| a | | |
| a | | |
| a | | |
| a | | |
| b | | |
| b | | |
| b | | |
| a | | |
| a | | |
| b | | |
| b | 488 | 439 |
| a | | |



Given this range, find the
tally of 'a': 483-486
What about 'b': 434-437

| <i>Tally</i> | |
|--------------|-----------------|
| <i>L</i> | a b |
| ⋮ | ⋮ |
| a | 482 432 |
| b | |
| b | |
| a | |
| a | |
| a | |
| a | |
| b | |
| b | |
| b | |
| a | |
| a | |
| b | |
| b | 488 439 |
| a | |

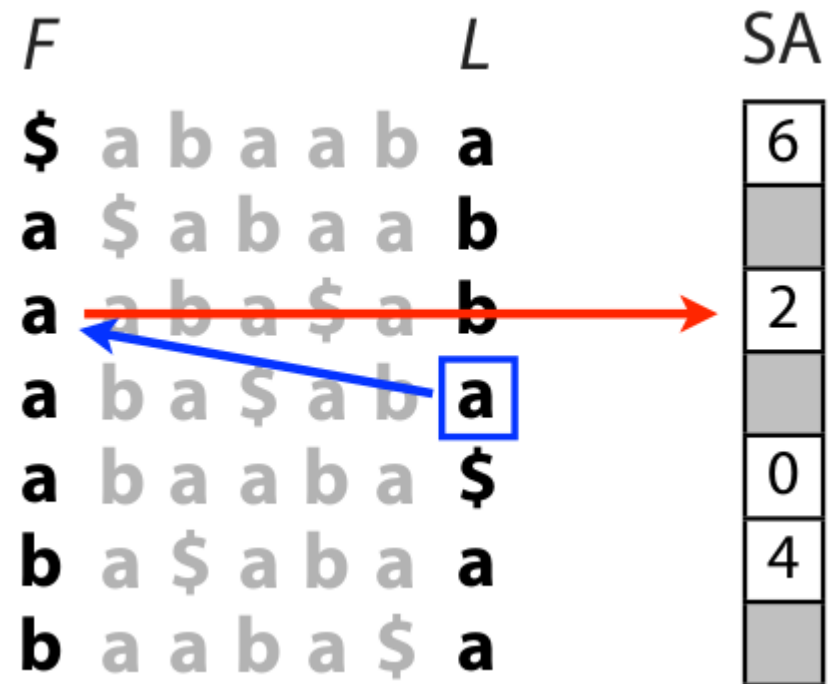
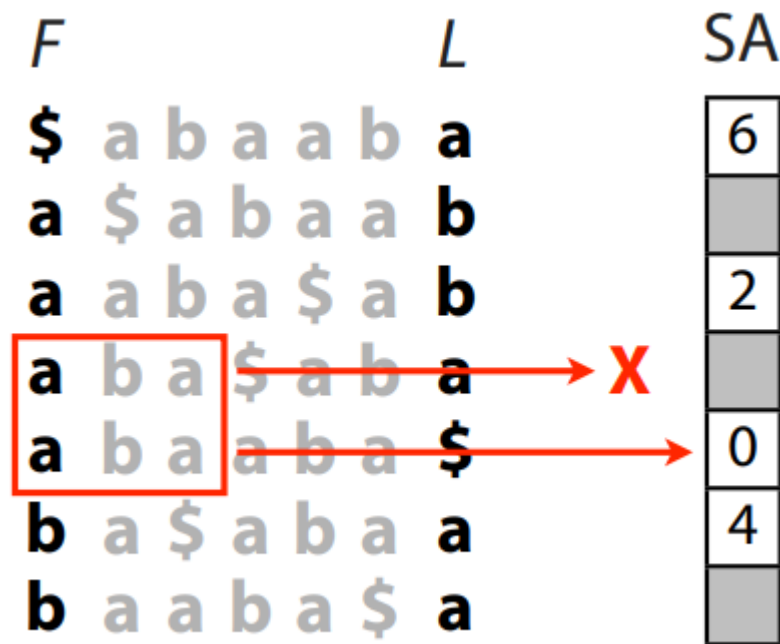


Given this range, find the
tally of 'a': 483-486

What about 'b': 434-437

Only need to count to the
nearest checkpoint

Positions with subset of Suffix Array



If $a[1]$ is at position 2 in SA
 Then $a[3]$ is at position 3 in SA
 aba is at positions [0] and [3]

Components of the FM Index:

First column (F): $\sim |\Sigma|$ integers

Last column (L): m characters

SA sample: $m \cdot a$ integers, where a is fraction of rows kept

Checkpoints: $m \times |\Sigma| \cdot b$ integers, where b is fraction of rows checkpointed

Example: DNA alphabet (2 bits per nucleotide), T = human genome,
 $a = 1/32$, $b = 1/128$

First column (F): 16 bytes

Last column (L): 2 bits * 3 billion chars = 750 MB

SA sample: 3 billion chars * 4 bytes/char / 32 = \sim 400 MB

Checkpoints: 3 billion * 4 bytes/char / 128 = \sim 100 MB

Total < 1.5 GB

Problems and solutions

We now have fast, memory efficient string search

What if our reads are not a perfect match (sequencing errors or SNPs)

- Seed and extend
- Seed and align (extend with gaps)

Solution: Seeding is needed to avoid mismatches (errors & SNPs)

How many seeds? If the read has 30 characters, and seed length is 10, and the seed interval is 6, the seeds extracted will be:

| | |
|-------------------|---------------------------------------|
| Read: | TAGCTACGCTCTACGCTATCATGCATAAAC |
| Seed 1 fw: | TAGCTACGCT |
| Seed 1 rc: | AGCGTAGCTA |
| Seed 2 fw: | CGCTCTACGC |
| Seed 2 rc: | GCGTAGAGCG |
| Seed 3 fw: | ACGCTATCAT |
| Seed 3 rc: | ATGATAGCGT |
| Seed 4 fw: | TCATGCATAA |
| Seed 4 rc: | TTATGCATGA |

Whats the largest number of mismatches that can occur and still **guarantee** at least one seed? **1**

Solution: Seeding is needed to avoid mismatches (errors & SNPs)

How many seeds? If the read has 30 characters, and seed length is 10, and the seed interval is 6, the seeds extracted will be:

| | | |
|------------|--------------------------------|--|
| Read: | TAGCTACGCTCTACGCTATCATGCATAAAC | |
| Seed 1 fw: | TAGCTACGCT | |
| Seed 1 rc: | AGCGTAGCTA | |
| Seed 2 fw: | CGCTCTACGC | |
| Seed 2 rc: | GCGTAGAGCG | |
| Seed 3 fw: | ACGCTATCAT | |
| Seed 3 rc: | ATGATAGCGT | |
| Seed 4 fw: | TCATGCATAA | |
| Seed 4 rc: | TTATGCATGA | |

Whats the largest number of mismatches that can occur and still **guarantee** at least one seed? **1**

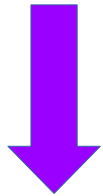
Solution: seed and extend

Read: TAGCTACGCTCTACGCT**G**TCATGCATAAAC

Seed: TAGCTACGCT

Genome:

Seed



ATCGATTAGCTACGCTCTACGCTATCATGCATAAAC CAGCATCGCA



Extend

Mismatches are ok

What about gaps

Read: TAGCTACGCTCTACGCT**G**TCATGCATAAAC

Seed: TAGCTACGCT

Genome:

Seed

ATCGATTAGCTACGCTCTACGCTATCATGCATACCAGCATCGCA

TAAAC

TAAAC

TAAAC

TACCA

T--ACCA

T-ACCA

3S

0S

1S

0G

2G

1G

Extend

Gaps are not ok

Exercises

- 1) Find positions of F for 'baa', 'aab' and 'ab' using FM index query.

| F | L | Position of F |
|----|----|---------------|
| \$ | b | x[0] |
| a | b | x[1] |
| a | b | x[3] |
| a | b | x[4] |
| a | \$ | x[5] |
| a | a | x[6] |
| a | a | x[7] |
| a | a | x[8] |
| b | b | x[9] |
| b | a | x[10] |
| b | b | x[11] |
| b | a | x[12] |
| b | a | x[13] |
| b | a | x[14] |

For example: 'ab' must start with 'a' and so could be in positions 1-8 of x (which stores positions). The question is which of these 'a' positions are followed by 'b'. Remember to search for the reverse of the string.