

Suffix arrays, BWT and FM-index

Alan Medlar

Wednesday 16th March 2016

Outline

- **Lecture:** Technical background for read mapping tools used in this course
 - Suffix array
 - Burrows-Wheeler transform (BWT)
 - FM-index
- **Lab session:** Using BWA to map paired-end data against the human genome, SAM/BAM files, etc

Read mapping

- Sequencers can generate up to **100 million** reads per sample
- Human genome is **~3 billion** basepairs
- Need to map reads to the genome to discover variants (SNVs, indels), counts (gene expression)

Preliminaries

- String
 - sequence of characters,
 - e.g. "banana", "ATGC", "MDLISTFS"
- Alphabet { A, C, G, T, \$ }, { A-Z, a-z, \$ }
- Lexicographical order
 - $\$ < A < C < G < T$

Preliminaries

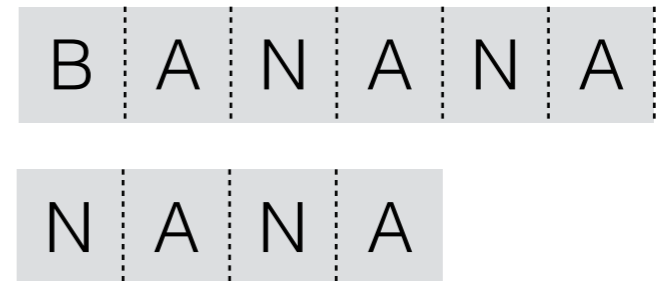
- Prefix
 - non-empty substring that is the beginning of another string (left-to-right)
 - e.g. "**banana**", "**ATGC**", "**MDLISTFS**"
- Suffix
 - non-empty substring that is the ending of another string (right-to-left)
 - e.g. "ba**nana**", "AT**GC**", "**MDLISTFS**"

Naïve exact search

- Text = "banana"
- Query = "nana"
- Linear search

Naïve exact search

- Text = "banana"
- Query = "nana"
- Linear search



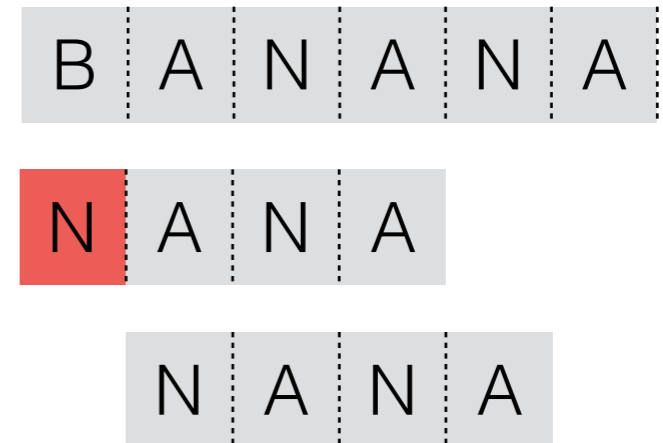
Naïve exact search

- Text = "banana"
- Query = "nana"
- Linear search



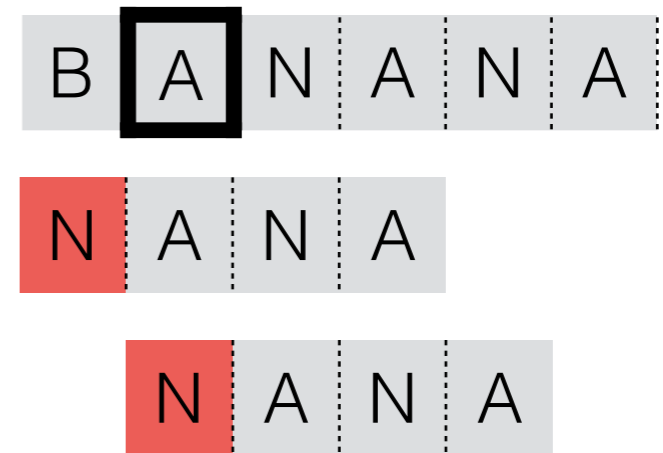
Naïve exact search

- Text = "banana"
- Query = "nana"
- Linear search



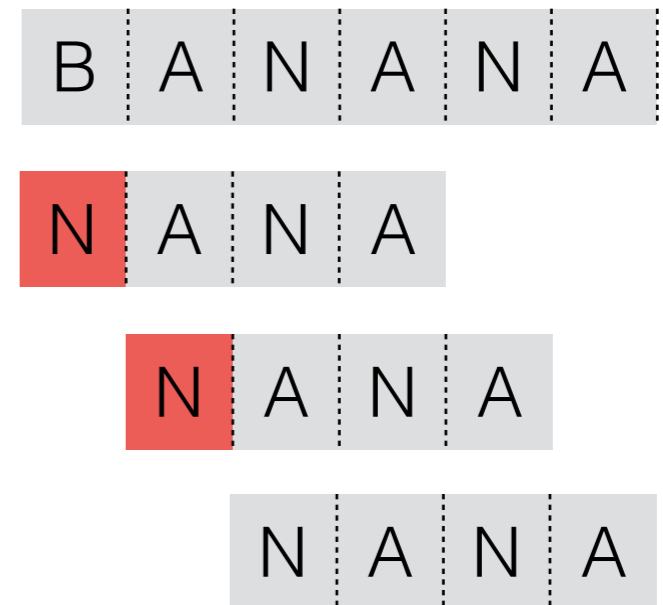
Naïve exact search

- Text = "banana"
- Query = "nana"
- Linear search



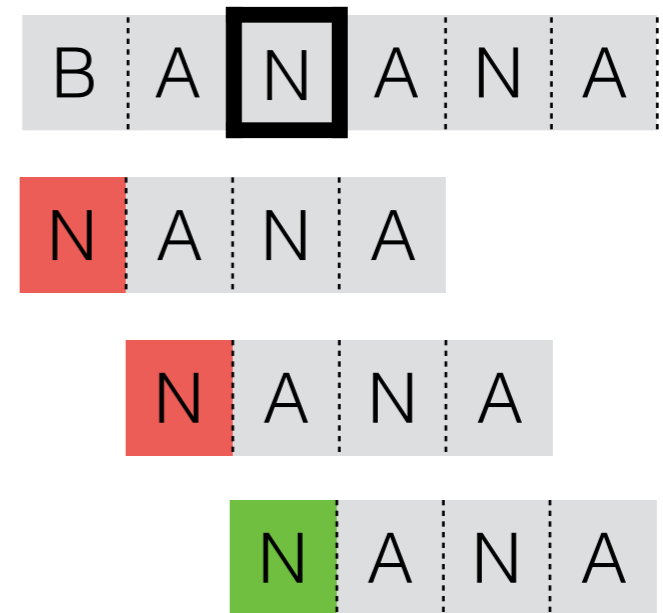
Naïve exact search

- Text = "banana"
- Query = "nana"
- Linear search



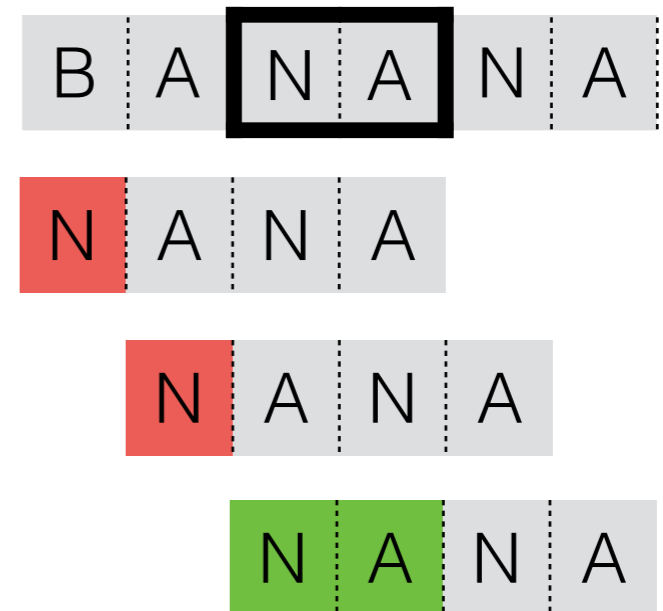
Naïve exact search

- Text = "banana"
- Query = "nana"
- Linear search



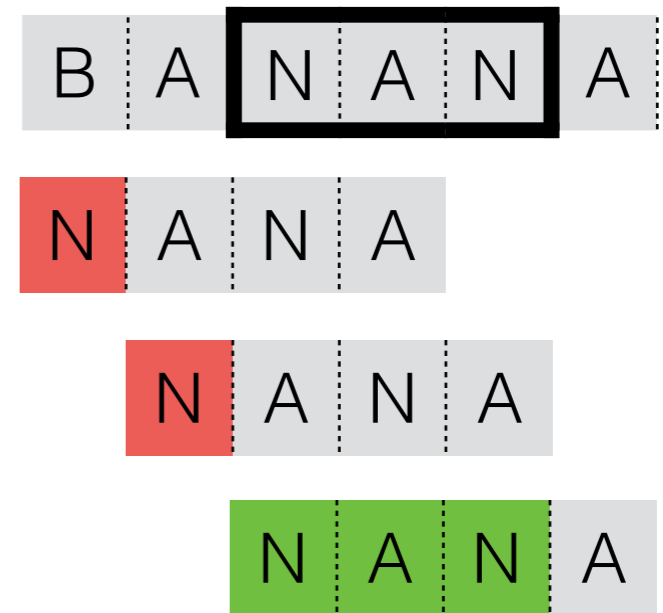
Naïve exact search

- Text = "banana"
- Query = "nana"
- Linear search



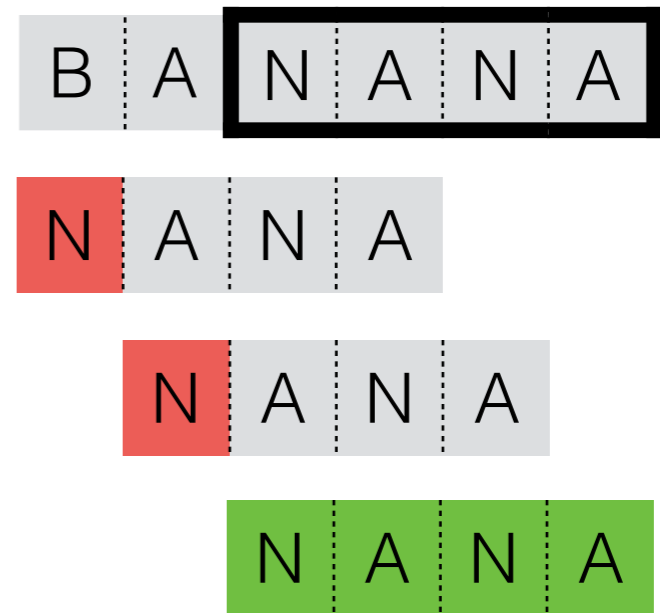
Naïve exact search

- Text = "banana"
- Query = "nana"
- Linear search



Naïve exact search

- Text = "banana"
- Query = "nana"
- Linear search



Naïve search is too slow

- Human genome ~3 billion basepairs
- Read 100 basepairs
- Complexity of search scales **linearly** with the length of the text!

Suffix array

- Introduced by Manber and Myers (1990) as a space efficient alternative to suffix tree (independently by Gonnet (1987))
- Sorted array of all suffixes of a given text
- Allows fast search of very large texts (e.g. genomes)

SA: building

B	A	N	A	N	A	\$

\$ is lexicographically **lower** than all other characters in the alphabet and cannot appear in the text otherwise

SA: building

B	A	N	A	N	A	\$	
A	N	A	N	A	\$		
N	A	N	A	\$			
A	N	A	\$				
N	A	\$					
A	\$						
\$							

SA: building

B	A	N	A	N	A	\$	0
A	N	A	N	A	\$		1
N	A	N	A	\$			2
A	N	A	\$				3
N	A	\$					4
A	\$						5
\$							6

SA: building

B	A	N	A	N	A	\$	0
A	N	A	N	A	\$		1
N	A	N	A	\$			2
A	N	A	\$				3
N	A	\$					4
A	\$						5
\$							6

\$							6
A	\$						5
A	N	A	\$				3
A	N	A	N	A	\$		1
B	A	N	A	N	A	\$	0
N	A	\$					4
N	A	N	A	\$			2

SA: building

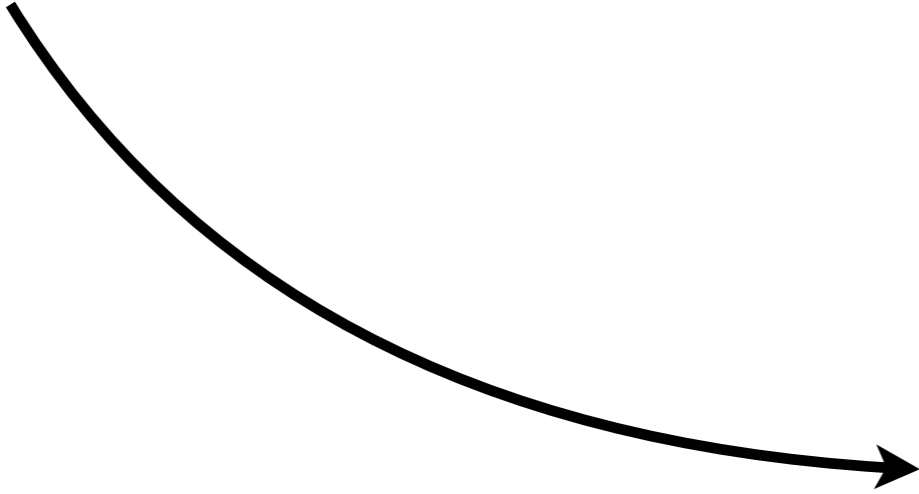
B	A	N	A	N	A	\$	0
A	N	A	N	A	\$		1
N	A	N	A	\$			2
A	N	A	\$				3
N	A	\$					4
A	\$						5
\$							6

\$							6
A	\$						5
A	N	A	\$				3
A	N	A	N	A	\$		1
B	A	N	A	N	A	\$	0
N	A	\$					4
N	A	N	A	\$			2

SA: querying

- Search for **prefixes** in the suffix array that match our query string
- SA is sorted, so we can use binary search!

N	A	N	A
---	---	---	---



\$							6
A	\$						5
A	N	A	\$				3
A	N	A	N	A	\$		1
B	A	N	A	N	A	\$	0
N	A	\$					4
N	A	N	A	\$			2

N	A	N	A
---	---	---	---

\$							6
A	\$						5
A	N	A	\$				3
A	N	A	N	A	\$		1
B	A	N	A	N	A	\$	0
N	A	\$					4
N	A	N	A	\$			2

N A N A

\$							6
A	\$						5
A	N	A	\$				3
A	N	A	N	A	\$		1
B	A	N	A	N	A	\$	0
N	A	\$					4
N	A	N	A	\$			2

N	A	N	A
---	---	---	---

\$							6
A	\$						5
A	N	A	\$				3
A	N	A	N	A	\$		1
B	A	N	A	N	A	\$	0
N	A	\$					4
N	A	N	A	\$			2

N | A | N | A

\$							6
A	\$						5
A	N	A	\$				3
A	N	A	N	A	\$		1
B	A	N	A	N	A	\$	0
N	A	\$					4
N	A	N	A	\$			2

N A N A

\$							6
A	\$						5
A	N	A	\$				3
A	N	A	N	A	\$		1
B	A	N	A	N	A	\$	0
N	A	\$					4
N	A	N	A	\$			2

N A N A

\$							6
A	\$						5
A	N	A	\$				3
A	N	A	N	A	\$		1
B	A	N	A	N	A	\$	0
N	A	\$					4
N	A	N	A	\$			2

N A N A

\$							6
A	\$						5
A	N	A	\$				3
A	N	A	N	A	\$		1
B	A	N	A	N	A	\$	0
N	A	\$					4
N	A	N	A	\$			2

N	A	N	A
---	---	---	---

\$							6
A	\$						5
A	N	A	\$				3
A	N	A	N	A	\$		1
B	A	N	A	N	A	\$	0
N	A	\$					4
N	A	N	A	\$			2

N | A | N | A

\$							6
A	\$						5
A	N	A	\$				3
A	N	A	N	A	\$		1
B	A	N	A	N	A	\$	0
N	A	\$					4
N	A	N	A	\$			2

N A N A

\$							6
A	\$						5
A	N	A	\$				3
A	N	A	N	A	\$		1
B	A	N	A	N	A	\$	0
N	A	\$					4
N	A	N	A	\$			2

N A N A

\$							6
A	\$						5
A	N	A	\$				3
A	N	A	N	A	\$		1
B	A	N	A	N	A	\$	0
N	A	\$					4
N	A	N	A	\$			2

N A N A

\$							6
A	\$						5
A	N	A	\$				3
A	N	A	N	A	\$		1
B	A	N	A	N	A	\$	0
N	A	\$					4
N	A	N	A	\$			2

N A N A

\$							6
A	\$						5
A	N	A	\$				3
A	N	A	N	A	\$		1
B	A	N	A	N	A	\$	0
N	A	\$					4
N	A	N	A	\$			2

SA vs. naïve search

- Searching the human genome (~3 billion basepairs, n) for a single-end read (100 basepairs, m)
- Naïve search $O(mn)$
- Suffix array search $O(m \log(n))$

SA vs. naïve search

- Searching the human genome (~3 billion basepairs, n) for a single-end read (100 basepairs, m)
- Naïve search $O(mn)$
- Suffix array search $O(m \log(n))$

n	$O(n)$	$O(\log(n))$
8	8	3
16	16	4
32	32	5
64	64	6
128	128	7
256	256	8
512	512	9

Good enough for read mapping?

- Human genome is ~3 billion basepairs
- Assume 5 bytes per basepair (1 byte characters, 4 byte integers) = **~14 GB**
- NGS data really hit in 2009 (16 GB RAM at the time was a luxury!)

Burrows-Wheeler transform

- Invented by Burrows and Wheeler (1994) while working at DEC
- Used in compression (.bz2 files)
- Interested in three things:
 - how to perform BWT
 - why BWT is useful for compression
 - how to reverse BWT

BWT

B	A	N	A	N	A	\$
A	N	A	N	A	\$	

BWT

B	A	N	A	N	A	\$
A	N	A	N	A	\$	B

BWT

B	A	N	A	N	A	\$
A	N	A	N	A	\$	B
N	A	N	A	\$		

BWT

B	A	N	A	N	A	\$
A	N	A	N	A	\$	B
N	A	N	A	\$	B	A

BWT

B	A	N	A	N	A	\$
A	N	A	N	A	\$	B
N	A	N	A	\$	B	A
A	N	A	\$	B	A	N
N	A	\$	B	A	N	A
A	\$	B	A	N	A	N
\$	B	A	N	A	N	A

BWT

B	A	N	A	N	A	\$
A	N	A	N	A	\$	B
N	A	N	A	\$	B	A
A	N	A	\$	B	A	N
N	A	\$	B	A	N	A
A	\$	B	A	N	A	N
\$	B	A	N	A	N	A

\$	B	A	N	A	N	A
A	\$	B	A	N	A	N
A	N	A	\$	B	A	N
A	N	A	N	A	\$	B
B	A	N	A	N	A	\$
N	A	\$	B	A	N	A
N	A	N	A	\$	B	A

BWT

\$	B	A	N	A	N	A
A	\$	B	A	N	A	N
A	N	A	\$	B	A	N
A	N	A	N	A	\$	B
B	A	N	A	N	A	\$
N	A	\$	B	A	N	A
N	A	N	A	\$	B	A

BWT compression

- $T = \text{"banana\$"}$
- $\text{BWT}(T) = \text{"annb\$aa"}$

BWT compression

- T =
"peter_piper_picked_a_peck_of_pickled_peppers_a_peck_of
_pickled_peppers_peter_piper_picked_if_peter_piper_picked
_a_peck_of_pickled_peppers_wheres_the_peck_of_pickled_
peppers_peter_piper_picked"
- BWT(T) =
"ddsddkkkkaeaadddsfsrrrrffffrrrrss___eeeeiiiiiiiieeeeeeeehpp
ppkkkkllllppppppppptttthpppprpppppioooootwpppppppppp_pppp
ccccccccccckkkk_____iiiipppp_____eee
eeeeeeeeeeeeerrrereeee_"

Relation to suffix array

- BWT matrix truncated at "\$" in each row **is** the suffix array of the same text
- BWT can be computed directly from the suffix array

\$	B	A	N	A	N	A
A	\$	B	A	N	A	N
A	N	A	\$	B	A	N
A	N	A	N	A	\$	B
B	A	N	A	N	A	\$
N	A	\$	B	A	N	A
N	A	N	A	\$	B	A

Reverse BWT

- It not very useful to compress something if we cannot get the original text back!
- $BWT'(BWT(T)) = T$

LF mapping (T-rank)

B | A | N | A | N | A | \$

LF mapping (T-rank)

B A N A N A \$

T-rank

B_0 A_0 N_0 A_1 N_1 A_2 \$

LF mapping (T-rank)

T-rank

B	A	N	A	N	A	\$
B ₀	A ₀	N ₀	A ₁	N ₁	A ₂	\$

F						L
\$	B ₀	A ₀	N ₀	A ₁	N ₁	A ₂
A ₂	\$	B ₀	A ₀	N ₀	A ₁	N ₁
A ₁	N ₁	A ₂	\$	B ₀	A ₀	N ₀
A ₀	N ₀	A ₁	N ₁	A ₂	\$	B ₀
B ₀	A ₀	N ₀	A ₁	N ₁	A ₂	\$
N ₁	A ₂	\$	B ₀	A ₀	N ₀	A ₁
N ₀	A ₁	N ₁	A ₂	\$	B ₀	A ₀

LF mapping (T-rank)

T-rank

B	A	N	A	N	A	\$
B ₀	A ₀	N ₀	A ₁	N ₁	A ₂	\$

F						L
\$	B ₀	A ₀	N ₀	A ₁	N ₁	A₂
A₂	\$	B ₀	A ₀	N ₀	A ₁	N₁
A₁	N ₁	A ₂	\$	B ₀	A ₀	N₀
A₀	N ₀	A ₁	N ₁	A ₂	\$	B₀
B₀	A ₀	N ₀	A ₁	N ₁	A ₂	\$
N₁	A ₂	\$	B ₀	A ₀	N ₀	A₁
N₀	A ₁	N ₁	A ₂	\$	B ₀	A₀

LF mapping (T-rank)

T-rank

B	A	N	A	N	A	\$
B ₀	A ₀	N ₀	A ₁	N ₁	A ₂	\$

F						L
\$	B ₀	A ₀	N ₀	A ₁	N ₁	A₂
A₂	\$	B ₀	A ₀	N ₀	A ₁	N₁
A₁	N ₁	A ₂	\$	B ₀	A ₀	N₀
A₀	N ₀	A ₁	N ₁	A ₂	\$	B₀
B₀	A ₀	N ₀	A ₁	N ₁	A ₂	\$
N₁	A ₂	\$	B ₀	A ₀	N ₀	A₁
N₀	A ₁	N ₁	A ₂	\$	B ₀	A₀

LF mapping (T-rank)

T-rank

B	A	N	A	N	A	\$
B ₀	A ₀	N ₀	A ₁	N ₁	A ₂	\$

F						L
\$	B ₀	A ₀	N ₀	A ₁	N ₁	A₂
A₂	\$	B ₀	A ₀	N ₀	A ₁	N₁
A₁	N ₁	A ₂	\$	B ₀	A ₀	N₀
A₀	N ₀	A ₁	N ₁	A ₂	\$	B₀
B₀	A ₀	N ₀	A ₁	N ₁	A ₂	\$
N₁	A ₂	\$	B ₀	A ₀	N ₀	A₁
N₀	A ₁	N ₁	A ₂	\$	B ₀	A₀

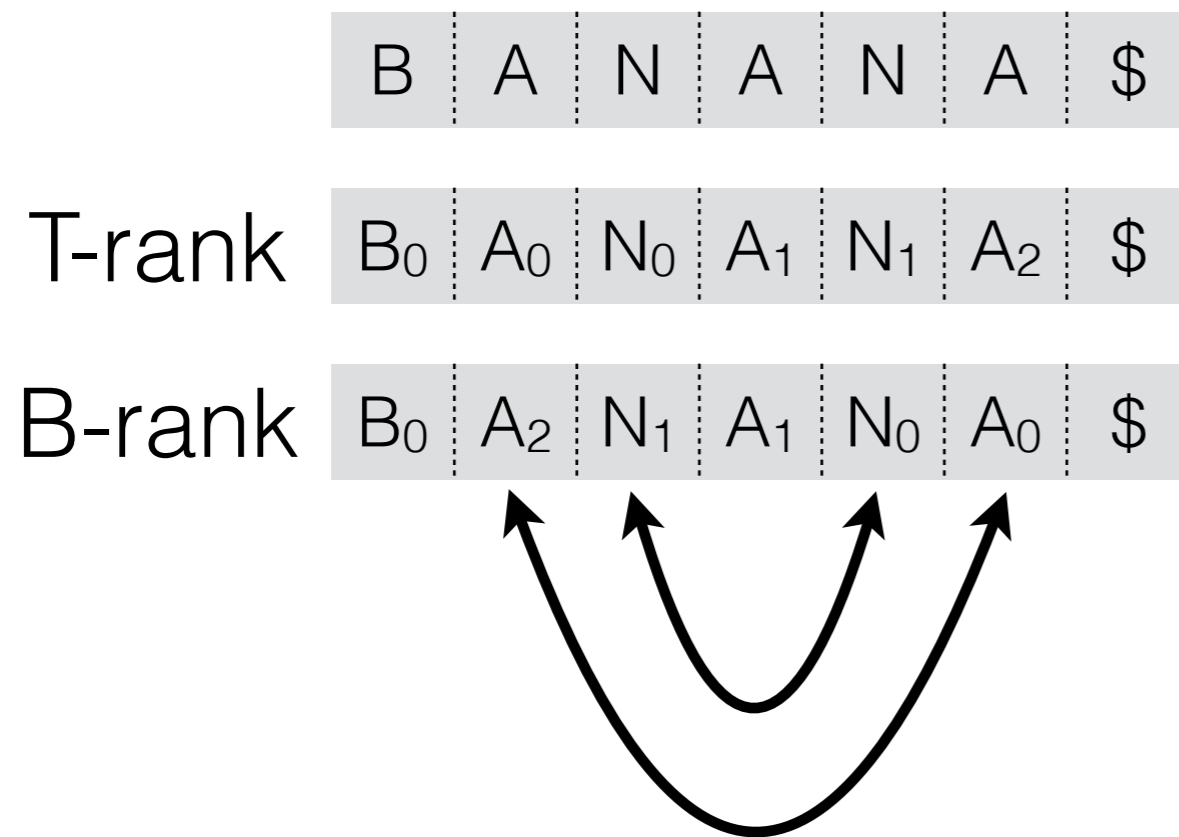
LF mapping (T-rank)

	B	A	N	A	N	A	\$
T-rank	B ₀	A ₀	N ₀	A ₁	N ₁	A ₂	\$

F							L
\$	B ₀	A ₀	N ₀	A ₁	N ₁	A ₂	A₂
A₂	\$	B ₀	A ₀	N ₀	A ₁	A ₂	N₁
A₁	N ₁	A ₂	\$	B ₀	A ₀	A ₁	N₀
A₀	N ₀	A ₁	N ₁	A ₂	\$	A ₂	B₀
B₀	A ₀	N ₀	A ₁	N ₁	A ₂	A ₂	\$
N₁	A ₂	\$	B ₀	A ₀	N ₀	A ₁	A₁
N₀	A ₁	N ₁	A ₂	\$	B ₀	A ₂	A₀

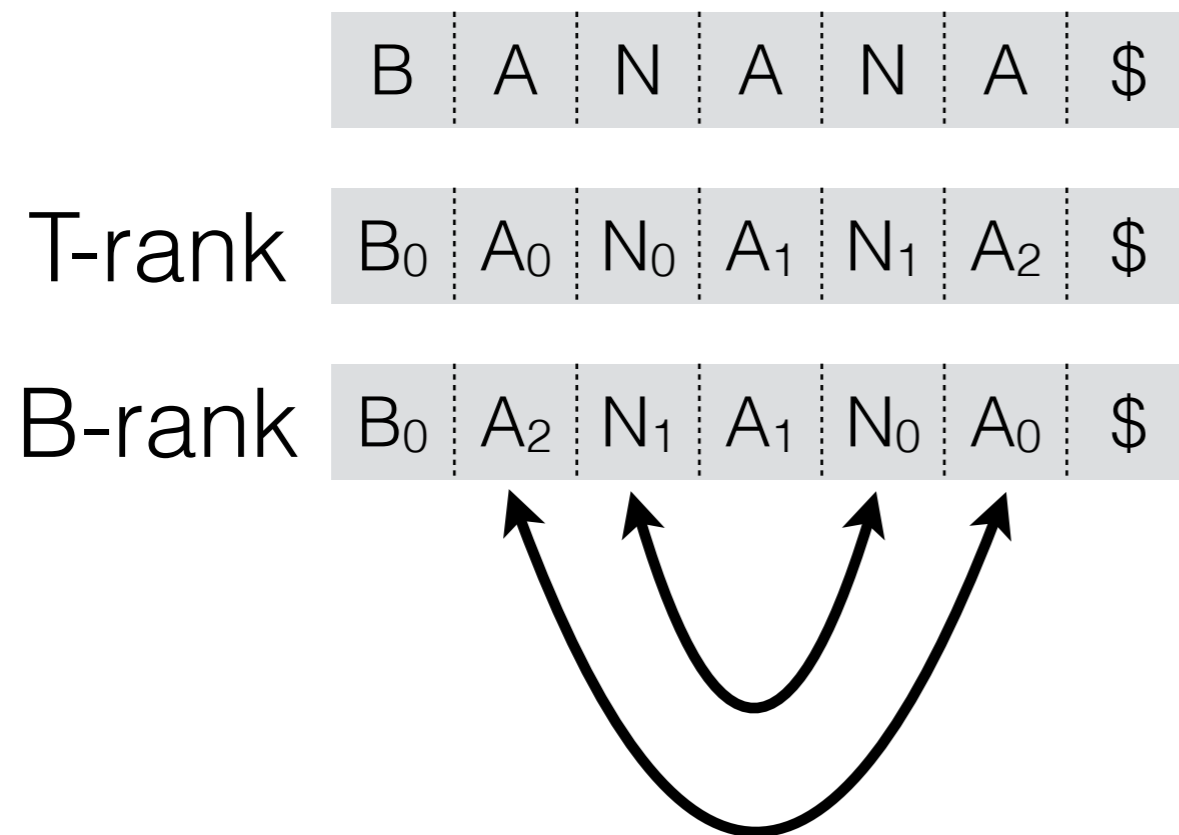
Ns in the L column are sorted by their "**right context**", same as Ns in F column!

LF mapping (B-rank)



F						L
\$	B ₀	A ₂	N ₁	A ₁	N ₀	A ₀
A ₀	\$	B ₀	A ₂	N ₁	A ₁	N ₀
A ₁	N ₀	A ₀	\$	B ₀	A ₂	N ₁
A ₂	N ₁	A ₁	N ₀	A ₀	\$	B ₀
B ₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$
N ₀	A ₀	\$	B ₀	A ₂	N ₁	A ₁
N ₁	A ₁	N ₀	A ₀	\$	B ₀	A ₂

LF mapping (B-rank)



F						L
\$	B ₀	A ₂	N ₁	A ₁	N ₀	A₀
A₀	\$	B ₀	A ₂	N ₁	A ₁	N₀
A₁	N ₀	A ₀	\$	B ₀	A ₂	N₁
A₂	N ₁	A ₁	N ₀	A ₀	\$	B₀
B₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$
N₀	A ₀	\$	B ₀	A ₂	N ₁	A₁
N₁	A ₁	N ₀	A ₀	\$	B ₀	A₂

LF mapping (B-rank)

- F column contains very little information, just counts of each character

F						L
\$	B ₀	A ₂	N ₁	A ₁	N ₀	A₀
A₀	\$	B ₀	A ₂	N ₁	A ₁	N₀
A₁	N ₀	A ₀	\$	B ₀	A ₂	N₁
A₂	N ₁	A ₁	N ₀	A ₀	\$	B₀
B₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$
N₀	A ₀	\$	B ₀	A ₂	N ₁	A₁
N₁	A ₁	N ₀	A ₀	\$	B ₀	A₂

LF mapping (B-rank)

L
A₀
N₀
N₁
B₀
\$
A₁
A₂

{ \$:1, A:3, B:1, N:2 }



Which row contains N₁ in the F column?

LF mapping (B-rank)

L
A₀
N₀
N₁
B₀
\$
A₁
A₂

{ \$:1, A:3, B:1, N:2 }



Which row contains N₁ in the F column?

- Skip \$ (+1)
- Skip As (+3)
- Skip Bs (+1)
- Skip first N (+1) = 6

LF mapping (B-rank)

	F	L
0	\$	A ₀
1	A ₀	N ₀
2	A ₁	N ₁
3	A ₂	B ₀
4	B ₀	\$
5	N ₀	A ₁
6	N ₁	A ₂

{ \$:1, A:3, B:1, N:2 }

← Which row contains N₁ in the F column?

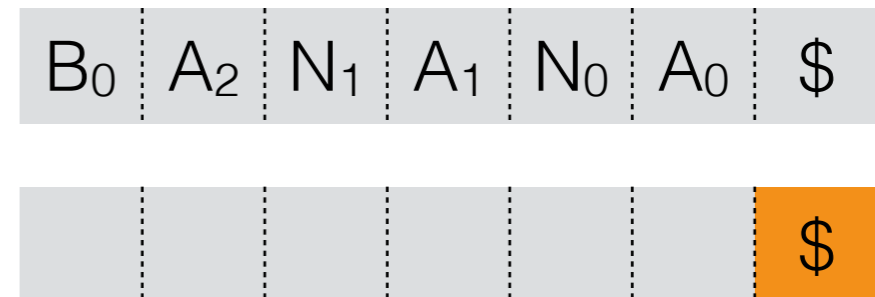
- Skip \$ (+1)
- Skip As (+3)
- Skip Bs (+1)
- Skip first N (+1) = 6

Reverse BWT

- Use B-ranking to reverse BWT, recreating the text T from right-to-left

{ \$:1, A:3, B:1, N:2 }

	F	L
0	\$	A ₀
1	A ₀	N ₀
2	A ₁	N ₁
3	A ₂	B ₀
4	B ₀	\$
5	N ₀	A ₁
6	N ₁	A ₂

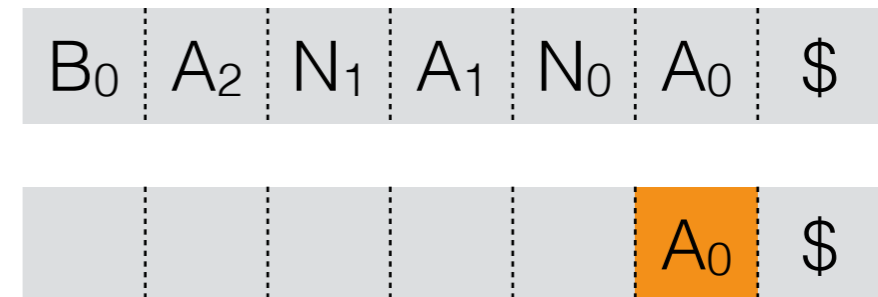


Reverse BWT

- Use B-ranking to reverse BWT, recreating the text T from right-to-left

{ \$:1, A:3, B:1, N:2 }

	F	L
0	\$	A ₀
1	A ₀	N ₀
2	A ₁	N ₁
3	A ₂	B ₀
4	B ₀	\$
5	N ₀	A ₁
6	N ₁	A ₂

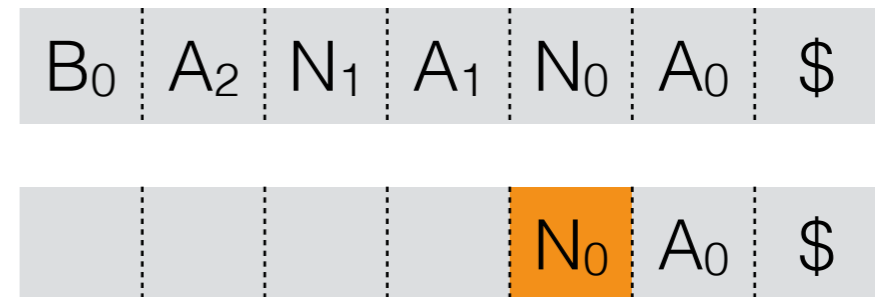


Reverse BWT

- Use B-ranking to reverse BWT, recreating the text T from right-to-left

{ \$:1, A:3, B:1, N:2 }

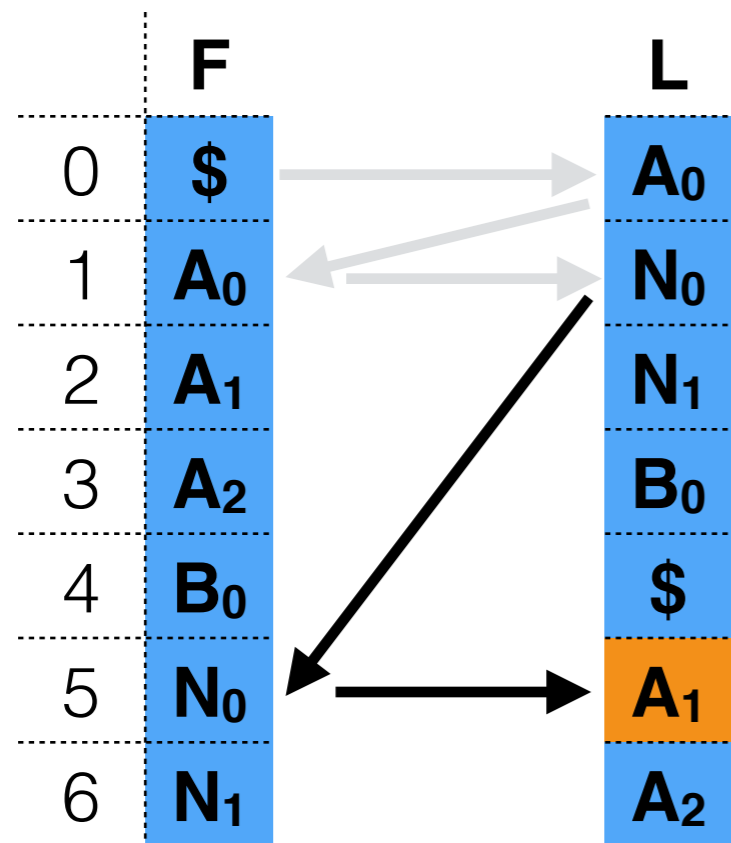
	F	L
0	\$	A ₀
1	A ₀	N ₀
2	A ₁	N ₁
3	A ₂	B ₀
4	B ₀	\$
5	N ₀	A ₁
6	N ₁	A ₂



Reverse BWT

- Use B-ranking to reverse BWT, recreating the text T from right-to-left

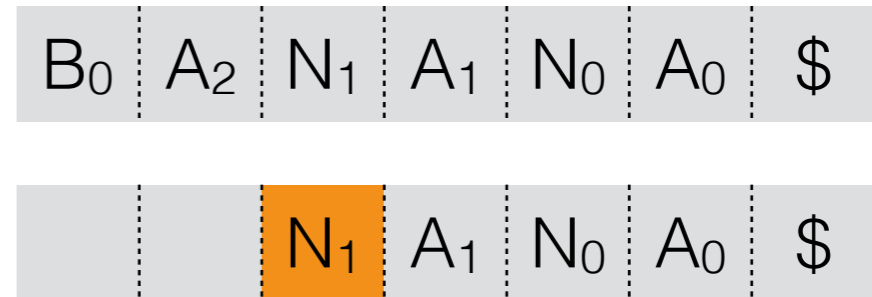
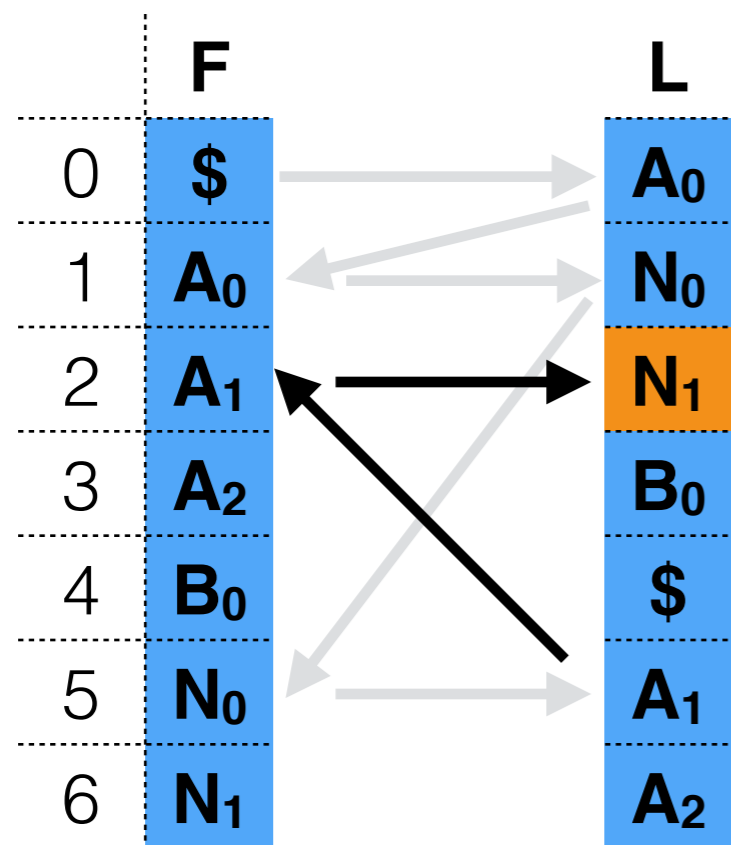
{ \$:1, A:3, B:1, N:2 }



Reverse BWT

- Use B-ranking to reverse BWT, recreating the text T from right-to-left

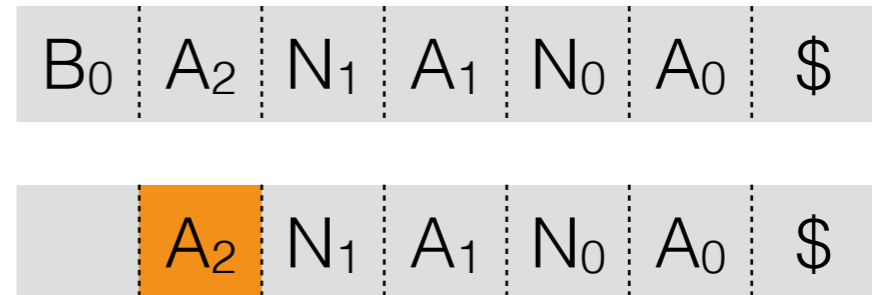
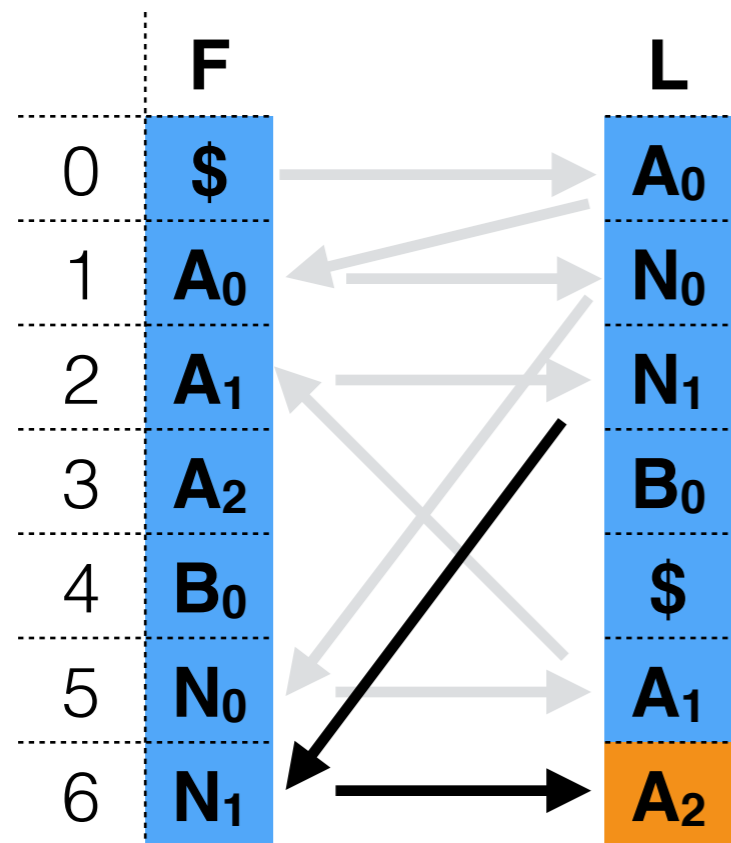
{ \$:1, A:3, B:1, N:2 }



Reverse BWT

- Use B-ranking to reverse BWT, recreating the text T from right-to-left

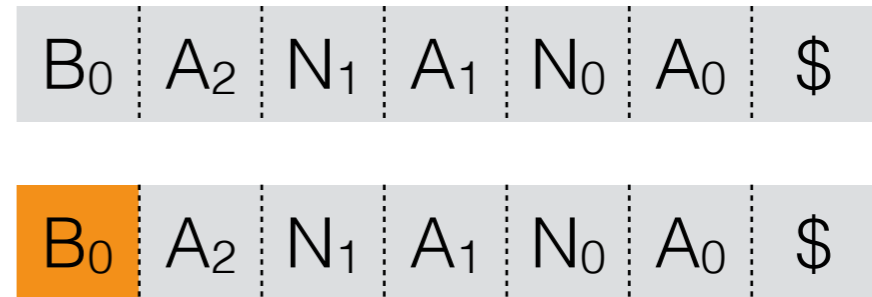
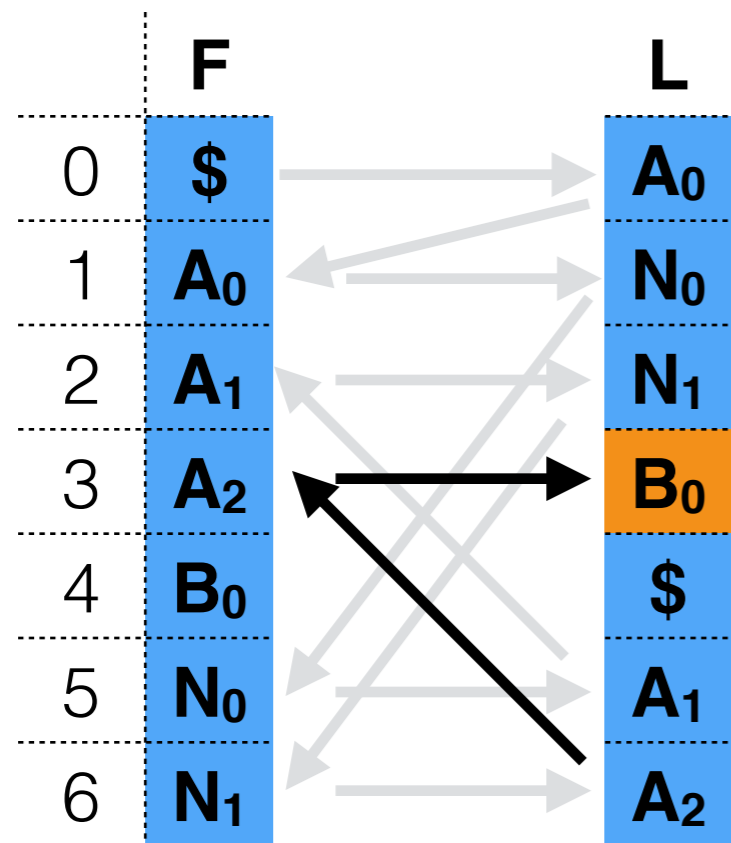
{ \$:1, A:3, B:1, N:2 }



Reverse BWT

- Use B-ranking to reverse BWT, recreating the text T from right-to-left

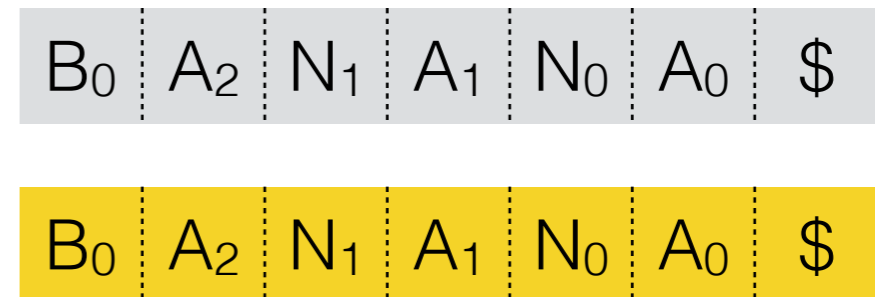
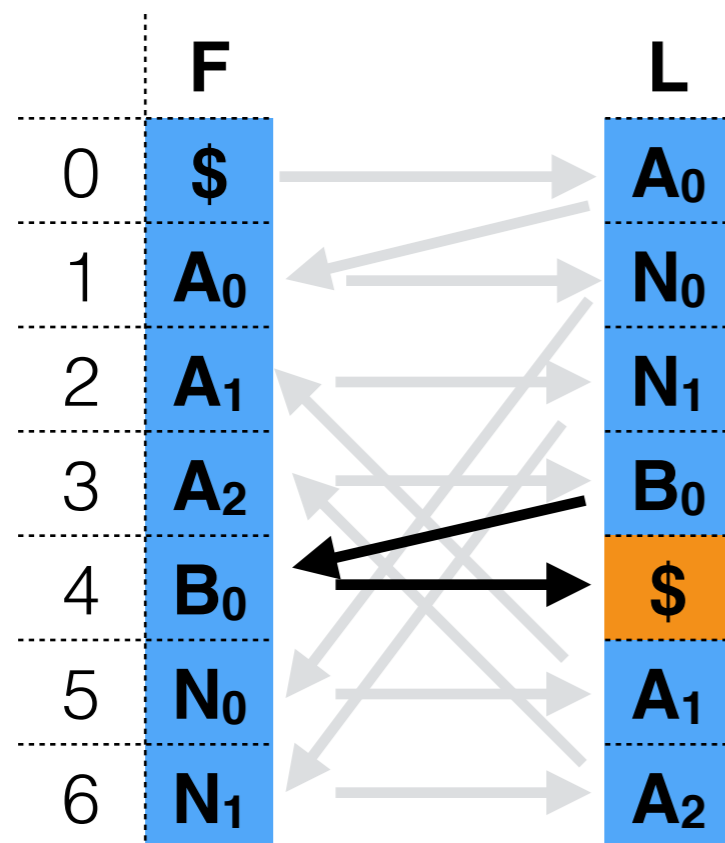
{ \$:1, A:3, B:1, N:2 }



Reverse BWT

- Use B-ranking to reverse BWT, recreating the text T from right-to-left

{ \$:1, A:3, B:1, N:2 }



FM-index

- All BWT allows us to do is compress text
- Ferragina and Manzini (2000)
"Full-text index in Minute space"
- Combine BWT with other auxiliary data structures to get an index
- Space savings: e.g. Human genome (3 billion bp)
 - SA = ~14 GB (5 bytes/bp)
 - FM = ~1.5 GB (2 bits/bp)

Cannot search BWT like SA

- Rotation matrix contains the suffix array

\$	B	A	N	A	N	A
A	\$	B	A	N	A	N
A	N	A	\$	B	A	N
A	N	A	N	A	\$	B
B	A	N	A	N	A	\$
N	A	\$	B	A	N	A
N	A	N	A	\$	B	A

Cannot search BWT like SA

- Rotation matrix contains the suffix array
- But we only store F and L columns, so binary search of prefixes not possible

\$	B	A	N	A	N	A
A	\$	B	A	N	A	N
A	N	A	\$	B	A	N
A	N	A	N	A	\$	B
B	A	N	A	N	A	\$
N	A	\$	B	A	N	A
N	A	N	A	\$	B	A

BWT search

- In SA, we matched successively longer **prefixes** (left-to-right) of query string (binary search)
- In BWT, we will match successively longer **suffixes** (right-to-left) of query string (reverse BWT transform)

BWT search

N A N A

F						L
\$	B ₀	A ₂	N ₁	A ₁	N ₀	A₀
A₀	\$	B ₀	A ₂	N ₁	A ₁	N₀
A₁	N ₀	A ₀	\$	B ₀	A ₂	N₁
A₂	N ₁	A ₁	N ₀	A ₀	\$	B₀
B₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$
N₀	A ₀	\$	B ₀	A ₂	N ₁	A₁
N₁	A ₁	N ₀	A ₀	\$	B ₀	A₂

BWT search

N A N A

F						L
\$	B ₀	A ₂	N ₁	A ₁	N ₀	A₀
A₀	\$	B ₀	A ₂	N ₁	A ₁	N₀
A₁	N ₀	A ₀	\$	B ₀	A ₂	N₁
A₂	N ₁	A ₁	N ₀	A ₀	\$	B₀
B₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$
N₀	A ₀	\$	B ₀	A ₂	N ₁	A₁
N₁	A ₁	N ₀	A ₀	\$	B ₀	A₂

BWT search

N A N A

F						L
\$	B ₀	A ₂	N ₁	A ₁	N ₀	A₀
A₀	\$	B ₀	A ₂	N ₁	A ₁	N₀
A₁	N ₀	A ₀	\$	B ₀	A ₂	N₁
A₂	N ₁	A ₁	N ₀	A ₀	\$	B₀
B₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$
N₀	A ₀	\$	B ₀	A ₂	N ₁	A₁
N₁	A ₁	N ₀	A ₀	\$	B ₀	A₂

BWT search

N A N A

F						L
\$	B ₀	A ₂	N ₁	A ₁	N ₀	A₀
A₀	\$	B ₀	A ₂	N ₁	A ₁	N₀
A₁	N ₀	A ₀	\$	B ₀	A ₂	N₁
A₂	N ₁	A ₁	N ₀	A ₀	\$	B₀
B₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$
N₀	A ₀	\$	B ₀	A ₂	N ₁	A₁
N₁	A ₁	N ₀	A ₀	\$	B ₀	A₂

BWT search

N A N A

F						L
\$	B ₀	A ₂	N ₁	A ₁	N ₀	A₀
A₀	\$	B ₀	A ₂	N ₁	A ₁	N₀
A₁	N ₀	A ₀	\$	B ₀	A ₂	N₁
A₂	N ₁	A ₁	N ₀	A ₀	\$	B₀
B₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$
N₀	A ₀	\$	B ₀	A ₂	N ₁	A₁
N₁	A ₁	N ₀	A ₀	\$	B ₀	A₂

BWT search

N A N A

F						L
\$	B ₀	A ₂	N ₁	A ₁	N ₀	A₀
A₀	\$	B ₀	A ₂	N ₁	A ₁	N₀
A₁	N ₀	A ₀	\$	B ₀	A ₂	N₁
A₂	N ₁	A ₁	N ₀	A ₀	\$	B₀
B₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$
N₀	A ₀	\$	B ₀	A ₂	N ₁	A₁
N₁	A ₁	N ₀	A ₀	\$	B ₀	A₂

BWT search

N A N A

F						L
\$	B ₀	A ₂	N ₁	A ₁	N ₀	A₀
A₀	\$	B ₀	A ₂	N ₁	A ₁	N₀
A₁	N ₀	A ₀	\$	B ₀	A ₂	N₁
A₂	N ₁	A ₁	N ₀	A ₀	\$	B₀
B₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$
N₀	A ₀	\$	B ₀	A ₂	N ₁	A₁
N₁	A ₁	N ₀	A ₀	\$	B ₀	A₂

BWT search

N A N A

We know BWT **contains** the query, but unlike SA, we do not know the **location** of the match in T!

F						L
\$	B ₀	A ₂	N ₁	A ₁	N ₀	A₀
A₀	\$	B ₀	A ₂	N ₁	A ₁	N₀
A₁	N ₀	A ₀	\$	B ₀	A ₂	N₁
A₂	N ₁	A ₁	N ₀	A ₀	\$	B₀
B₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$
N₀	A ₀	\$	B ₀	A ₂	N ₁	A₁
N₁	A ₁	N ₀	A ₀	\$	B ₀	A₂

BWT search

N A N A

F						L								
\$	B ₀	A ₂	N ₁	A ₁	N ₀	A₀	\$						6	
A₀	\$	B ₀	A ₂	N ₁	A ₁	N₀	A	\$					5	
A₁	N ₀	A ₀	\$	B ₀	A ₂	N₁	A	N	A	\$			3	
A₂	N ₁	A ₁	N ₀	A ₀	\$	B₀	A	N	A	N	A	\$	1	
B₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$	B	A	N	A	N	A	\$	0
N₀	A ₀	\$	B ₀	A ₂	N ₁	A₁	N	A	\$				4	
N₁	A ₁	N ₀	A ₀	\$	B ₀	A₂	N	A	N	A	\$		2	

Idea: just store SA as well?

BWT search

N A N A

F						L								
\$	B ₀	A ₂	N ₁	A ₁	N ₀	A₀	\$							6
A₀	\$	B ₀	A ₂	N ₁	A ₁	N₀	A	\$						
A₁	N ₀	A ₀	\$	B ₀	A ₂	N₁	A	N	A	\$				3
A₂	N ₁	A ₁	N ₀	A ₀	\$	B₀	A	N	A	N	A	\$		
B₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$	B	A	N	A	N	A	\$	0
N₀	A ₀	\$	B ₀	A ₂	N ₁	A₁	N	A	\$					
N₁	A ₁	N ₀	A ₀	\$	B ₀	A₂	N	A	N	A	\$			

Idea 2: store part of SA?

BWT search

N A N A

	F						L										
	\$	B ₀	A ₂	N ₁	A ₁	N ₀	A₀		\$								6
	A₀	\$	B ₀	A ₂	N ₁	A ₁	N₀		A	\$							
	A₁	N ₀	A ₀	\$	B ₀	A ₂	N₁		A	N	A	\$					3
+1	A₂	N ₁	A ₁	N ₀	A ₀	\$	B₀		A	N	A	N	A	\$			
	B₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$		B	A	N	A	N	A	\$		0
	N₀	A ₀	\$	B ₀	A ₂	N ₁	A₁		N	A	\$						
	N₁	A ₁	N ₀	A ₀	\$	B ₀	A₂		N	A	N	A	\$				

... and walk backwards through the BWT!

BWT search

N A N A

	F					L								
	\$	B ₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$					6	
	A ₀	\$	B ₀	A ₂	N ₁	A ₁	N ₀	A	\$					
	A ₁	N ₀	A ₀	\$	B ₀	A ₂	N ₁	A	N	A	\$		3	
+1	A ₂	N ₁	A ₁	N ₀	A ₀	\$	B ₀	A	N	A	N	A	\$	
+1	B ₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$	B	A	N	A	N	A	\$
	N ₀	A ₀	\$	B ₀	A ₂	N ₁	A ₁	N	A	\$				
	N ₁	A ₁	N ₀	A ₀	\$	B ₀	A ₂	N	A	N	A	\$		

... and walk backwards through the BWT!

BWT search

N A N A

	F					L								
	\$	B ₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$						6
	A ₀	\$	B ₀	A ₂	N ₁	A ₁	N ₀	A	\$					
	A ₁	N ₀	A ₀	\$	B ₀	A ₂	N ₁	A	N	A	\$			3
+1	A ₂	N ₁	A ₁	N ₀	A ₀	\$	B ₀	A	N	A	N	A	\$	
+1	B ₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$	B	A	N	A	N	A	\$
	N ₀	A ₀	\$	B ₀	A ₂	N ₁	A ₁	N	A	\$				
	N ₁	A ₁	N ₀	A ₀	\$	B ₀	A ₂	N	A	N	A	\$		

... and walk backwards through the BWT!

BWT search

N A N A

	F					L									
	\$	B ₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$						6	6
	A ₀	\$	B ₀	A ₂	N ₁	A ₁	N ₀	A	\$						5
	A ₁	N ₀	A ₀	\$	B ₀	A ₂	N ₁	A	N	A	\$			3	3
+1	A ₂	N ₁	A ₁	N ₀	A ₀	\$	B ₀	A	N	A	N	A	\$		1
+1	B ₀	A ₂	N ₁	A ₁	N ₀	A ₀	\$	B	A	N	A	N	A	\$	0
	N ₀	A ₀	\$	B ₀	A ₂	N ₁	A ₁	N	A	\$					4
	N ₁	A ₁	N ₀	A ₀	\$	B ₀	A ₂	N	A	N	A	\$			2

... and walk backwards through the BWT!

BWT search

- Finding location takes constant time if the offsets into T are **evenly spaced in T** , not in the SA!
- Make tradeoff between space (RAM) and time (how long lookups take)

Things we left out

- Rank calculations in the BWT need to be fast!
Needs another auxiliary data structure
- Only covered exact matching, read alignment requires mismatches (e.g. SNP in read, not in genome)
- Other details: store forwards and backwards indices of genome due to sequencing error profile