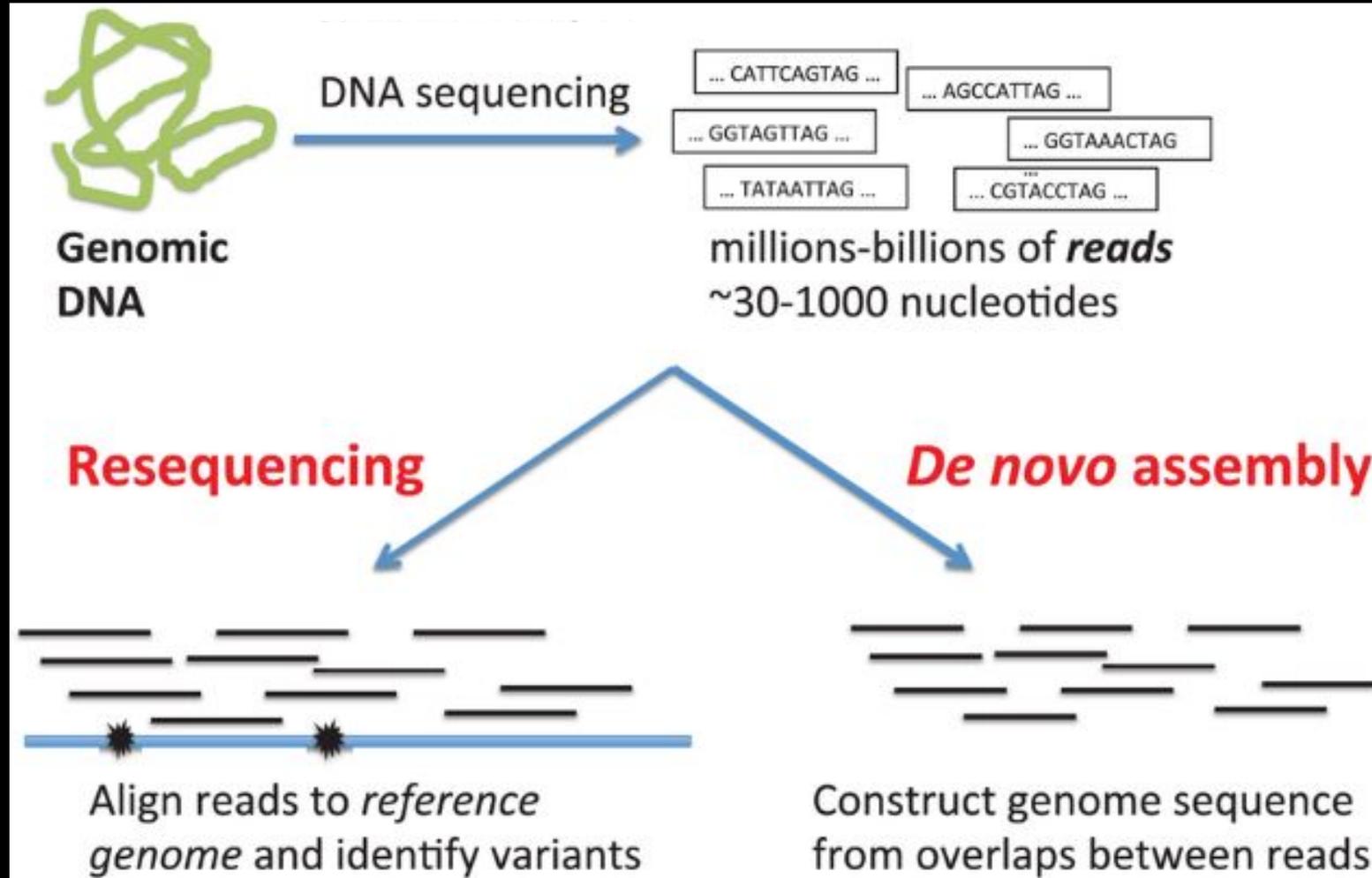


# Comparing Highly Similar Sequences: Suffix Trees, Suffix Arrays and the Burrows-Wheeler Transform

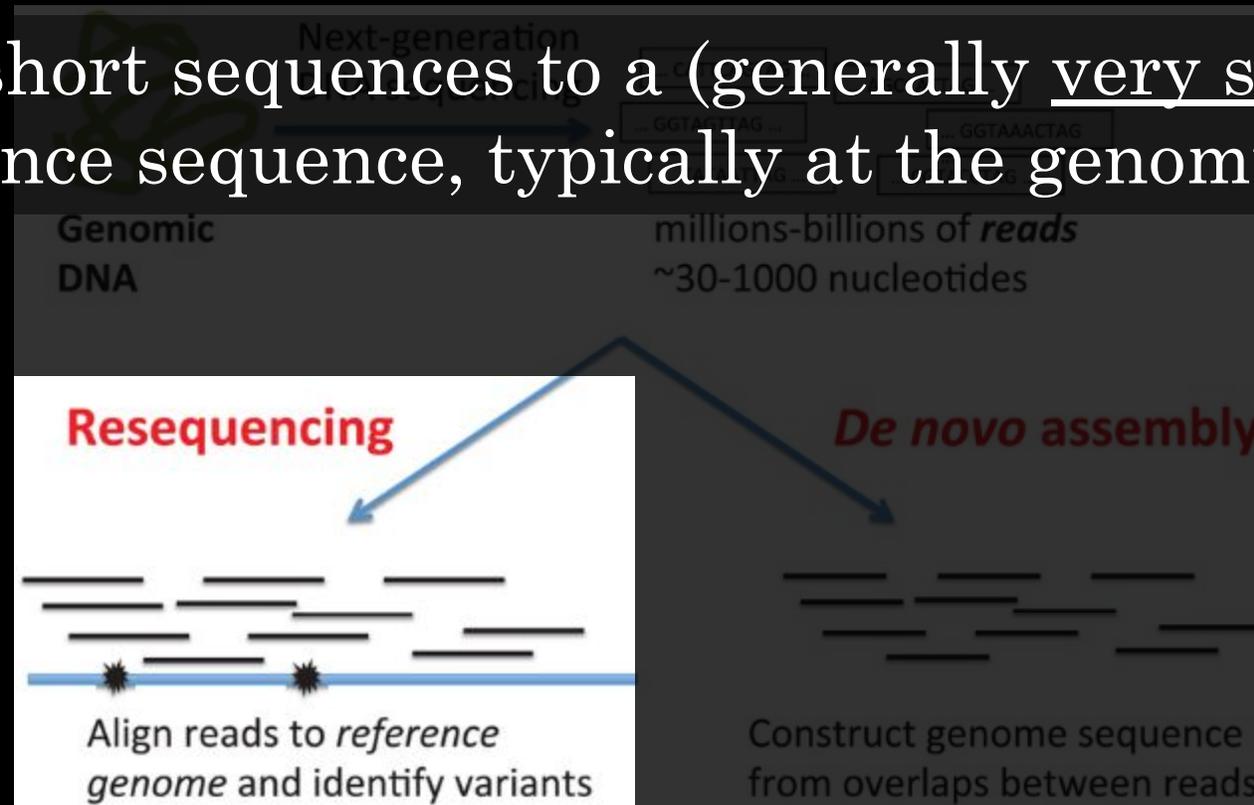


# Resequencing: A Different(-ish) Use Case



# Resequencing: A Different(-ish) Use Case

Map short sequences to a (generally very similar) reference sequence, typically at the genomic level



# Contrasting Cases

- MEGABLAST: **Somewhat similar** sequences
- BLAST: **Remote** homology – lots of sequence differences, possibly a lot of gaps too
- PSI-BLAST / hidden Markov models: **Super remote** homology
- Mapping: **Very similar**, few to no differences / gaps

# Resequencing



UK

100K

RARE GENETIC VARIANTS IN HEALTH AND DISEASE

2001

## articles

# Initial sequencing and analysis of the human genome

International Human Genome Sequencing Consortium\*

\* A partial list of authors appears on the opposite page. Affiliations are listed at the end of the paper.

The human genome holds an extraordinary trove of information about human development, physiology, medicine and evolution. Here we report the results of an international collaboration to produce and make freely available a draft sequence of the human genome. We also present an initial analysis of the data, describing some of the insights that can be gleaned from the sequence.

## THE HUMAN GENOME

## The Sequence of the Human Genome

J. Craig Venter,<sup>1\*</sup> Mark D. Adams,<sup>1</sup> Eugene W. Myers,<sup>1</sup> Peter W. Li,<sup>1</sup> Richard J. Mural,<sup>1</sup> Granger G. Sutton,<sup>1</sup> Hamilton O. Smith,<sup>1</sup> Mark Yandell,<sup>1</sup> Cheryl A. Evans,<sup>1</sup> Robert A. Holt,<sup>1</sup> Jeannine D. Gocayne,<sup>1</sup> Peter Amanatides,<sup>1</sup> Richard M. Ballew,<sup>1</sup> Daniel H. Huson,<sup>1</sup> Jennifer Russo Wortman,<sup>1</sup> Qing Zhang,<sup>1</sup> Chinnappa D. Kodira,<sup>1</sup> Xiangqun H. Zheng,<sup>1</sup> Lin Chen,<sup>1</sup> Marian Skupski,<sup>1</sup> Gangadharan Subramanian,<sup>1</sup> Paul D. Thomas,<sup>1</sup> Jinghui Zhang,<sup>1</sup> George L. Gabor Miklos,<sup>2</sup> Catherine Nelson,<sup>3</sup> Samuel Broder,<sup>1</sup> Andrew G. Clark,<sup>4</sup> Joe Nadeau,<sup>5</sup> Victor A. McKusick,<sup>6</sup> Norton Zinder,<sup>7</sup> Arnold J. Levine,<sup>7</sup> Richard J. Roberts,<sup>8</sup> Mel Simon,<sup>9</sup> Carolyn Slayman,<sup>10</sup> Michael Hunkapiller,<sup>11</sup> Randall Bolanos,<sup>1</sup> Arthur Delcher,<sup>1</sup> Ian Dew,<sup>1</sup> Daniel Fasulo,<sup>1</sup> Michael Flanigan,<sup>1</sup> Liliana Florea,<sup>1</sup> Aaron Halpern,<sup>1</sup> Sridhar Hannenhalli,<sup>1</sup> Saul Kravitz,<sup>1</sup> Samuel Levy,<sup>1</sup> Clark Mobarry,<sup>1</sup> Knut Reinert,<sup>1</sup> Karin Remington,<sup>1</sup> Jane Abu-Threideh,<sup>1</sup> Ellen Beasley,<sup>1</sup> Kendra Biddick,<sup>1</sup> Vivien Bonazzi,<sup>1</sup> Rhonda Brandon,<sup>1</sup> Michele Cargill,<sup>1</sup> Ishwar Chandramouliswaran,<sup>1</sup> Rosane Charlab,<sup>1</sup> Kabir Chaturvedi,<sup>1</sup> Zuoming Deng,<sup>1</sup> Valentina Di Francesco,<sup>1</sup> Patrick Dunn,<sup>1</sup> Karen Eilbeck,<sup>1</sup> Carlos Evangelista,<sup>1</sup> Andrei E. Gabrielian,<sup>1</sup> Weiniu Gan,<sup>1</sup> Wangmao Ge,<sup>1</sup> Fangcheng Gong,<sup>1</sup> Zhiping Gu,<sup>1</sup> Ping Guan,<sup>1</sup> Thomas J. Heiman,<sup>1</sup> Maureen E. Higgins,<sup>1</sup> Rui-Ru Ji,<sup>1</sup> Zhaoxi Ke,<sup>1</sup> Karen A. Ketchum,<sup>1</sup> Zhongwu Lai,<sup>1</sup> Yiding Lei,<sup>1</sup> Zhenya Li,<sup>1</sup> Jiayin Li,<sup>1</sup> Yong Liang,<sup>1</sup> Xiaoying Lin,<sup>1</sup> Fu Lu,<sup>1</sup> Gerald V. Mackay,<sup>1</sup> Natalia Milshina,<sup>1</sup> Helen M. Moore,<sup>1</sup> Achyuth Kumar K. Nair,<sup>1</sup>

2022

## HUMAN GENOMICS

# The complete sequence of a human genome

Sergey Nurk<sup>1</sup>†, Sergey Koren<sup>1</sup>†, Arang Rhie<sup>1</sup>†, Mikko Rautiainen<sup>1</sup>†, Andrey V. Bzikadze<sup>2</sup>, Alla Mikheenko<sup>3</sup>, Mitchell R. Vollger<sup>4</sup>, Nicolas Altemose<sup>5</sup>, Lev Uralsky<sup>6,7</sup>, Ariel Gershman<sup>8</sup>, Sergey Aganezov<sup>9</sup>†, Savannah J. Hoyt<sup>10</sup>, Mark Diekhans<sup>11</sup>, Glennis A. Logsdon<sup>4</sup>, Michael Alonge<sup>9</sup>, Stylianos E. Antonarakis<sup>12</sup>, Matthew Borchers<sup>13</sup>, Gerard G. Bouffard<sup>14</sup>, Shelise Y. Brooks<sup>14</sup>, Gina V. Caldas<sup>15</sup>, Nae-Chyun Chen<sup>9</sup>, Haoyu Cheng<sup>16,17</sup>, Chen-Shan Chin<sup>18</sup>, William Chow<sup>19</sup>, Leonardo G. de Lima<sup>13</sup>, Philip C. Dishuck<sup>4</sup>, Richard Durbin<sup>19,20</sup>, Tatiana Dvorkina<sup>3</sup>, Ian T. Fiddes<sup>21</sup>, Giulio Formenti<sup>22,23</sup>, Robert S. Fulton<sup>24</sup>, Arkarachai Fungtammasan<sup>18</sup>, Erik Garrison<sup>11,25</sup>, Patrick G. S. Grady<sup>10</sup>, Tina A. Graves-Lindsay<sup>26</sup>, Ira M. Hall<sup>27</sup>, Nancy F. Hansen<sup>28</sup>, Gabrielle A. Hartley<sup>10</sup>, Marina Haukness<sup>11</sup>, Kerstin Howe<sup>19</sup>, Michael W. Hunkapiller<sup>29</sup>, Chirag Jain<sup>1,30</sup>, Miten Jain<sup>11</sup>, Erich D. Jarvis<sup>22,23</sup>, Peter Kerpedjiev<sup>31</sup>, Melanie Kirsche<sup>9</sup>, Mikhail Kolmogorov<sup>32</sup>, Jonas Korfach<sup>29</sup>, Milinn Kremitzki<sup>26</sup>, Heng Li<sup>16,17</sup>, Valerie V. Maduro<sup>33</sup>, Tobias Marschall<sup>34</sup>, Ann M. McCartney<sup>1</sup>, Jennifer McDaniel<sup>35</sup>, Danny E. Miller<sup>4,36</sup>, James C. Mullikin<sup>14,28</sup>, Eugene W. Myers<sup>37</sup>, Nathan D. Olson<sup>35</sup>, Benedict Paten<sup>11</sup>, Paul Peluso<sup>29</sup>, Pavel A. Pevzner<sup>32</sup>, David Porubsky<sup>4</sup>, Tamara Potapova<sup>13</sup>, Evgeny I. Rogaev<sup>6,7,38,39</sup>, Jeffrey A. Rosenfeld<sup>40</sup>, Steven L. Salzberg<sup>9,41</sup>, Valerie A. Schneider<sup>42</sup>, Fritz J. Sedlazeck<sup>43</sup>, Kishwar Shafin<sup>11</sup>, Colin J. Shew<sup>44</sup>, Alaina Shumate<sup>41</sup>, Ying Sims<sup>19</sup>, Arian F. A. Smit<sup>45</sup>, Daniela C. Soto<sup>44</sup>, Ivan Sović<sup>29,46</sup>, Jessica M. Storer<sup>45</sup>, Aaron Streets<sup>5,47</sup>, Beth A. Sullivan<sup>48</sup>, Françoise Thibaud-Nissen<sup>42</sup>, James Torrance<sup>19</sup>, Justin Wagner<sup>35</sup>, Brian P. Walenz<sup>1</sup>, Aaron Wenger<sup>29</sup>, Jonathan M. D. Wood<sup>19</sup>, Chunlin Xiao<sup>42</sup>, Stephanie M. Yan<sup>49</sup>, Alice C. Young<sup>14</sup>, Samantha Zarate<sup>9</sup>, Urvashi Surti<sup>50</sup>, Rajiv C. McCoy<sup>49</sup>, Megan Y. Dennis<sup>44</sup>, Ivan A. Alexandrov<sup>3,7,51</sup>, Jennifer L. Gerton<sup>13,52</sup>, Rachel J. O'Neill<sup>10</sup>, Winston Timp<sup>8,41</sup>, Justin M. Zook<sup>35</sup>, Michael C. Schatz<sup>9,49</sup>, Evan E. Eichler<sup>4,53</sup>\*, Karen H. Miga<sup>11,54</sup>\*, Adam M. Phillippy<sup>1</sup>\*

Since its initial release in 2000, the human reference genome has covered only the euchromatic fraction of the genome, leaving important heterochromatic regions unfinished. Addressing the remaining 8% of the genome, the Telomere-to-Telomere (T2T) Consortium presents a complete 3.055 billion–base pair sequence of a human genome, T2T-CHM13, that includes gapless assemblies for all chromosomes except Y, corrects errors in the prior references, and introduces nearly 200 million base pairs of sequence containing 1956 gene predictions, 99 of which are predicted to be protein coding. The completed regions include all centromeric satellite arrays, recent segmental duplications, and the short arms of all five acrocentric chromosomes, unlocking these complex regions of the genome to variational and functional studies.

Lander et al. (2001) *Nature*Venter et al. (2001) *Science*Nurk et al. (2022) *Science*

“Landmark human genome papers”

<https://doe-humangenomeproject.ornl.gov/landmark-hgp-papers/>

# The Critical Question

- How can we exploit the high expected degree of similarity to accelerate the search?
- Key considerations:
  - DNA is a four-letter alphabet, there will be lots of **redundant sequences** in the genome
  - We can use approaches that treat mismatches and gaps less efficiently than dynamic programming

# Preliminaries: strings

## Offsets

S ← AGGCCTA

0123456

S[3]=C

S[5:]=TA

## Substring

AGCGCTA

AGCGCTA

GCGC

CTA

A

GC

GC

## Terminator

AGCGCTA\$

\$<A<C<G<T

AA\$<AAA\$

BEE<BEER

# Preliminaries: strings

## Prefixes

AGCGCTA\$

AGCGCTA

AGCGCT

AGCGC

AGCG

AGC

AG

A

## Suffixes

AGCGCTA\$

AGCGCTA\$

GCGCTA\$

CGCTA\$

GCTA\$

CTA\$

TA\$

A\$

\$

# Naive Searching

Find places where *pattern*  $P$  occurs as a substring of *text*  $T$ .  
Each such place is an *occurrence* or *match*.

Let  $n = |P|$ , and let  $m = |T|$     Assume  $n \leq m$

*Alignment*: a way of putting  $P$ 's characters opposite  $T$ 's.  
May or may not correspond to an match.

$P$ : word

$T$ : There would have been a time for such a word

Alignment 1: word

Alignment 2: word

# Naive Searching

P: word

T: There would have been a time for such a word

word word word word word word word word word **word**

word word word word word word word word word

word word word word word word word word word

word word word word word word word word word

word word word word word word word word word

← One  
occurrence

Try all possible alignments. For each, check if it matches.

This is the *naïve algorithm*.

Not good: back to  $m - n + 1$  possible matches or  $\sim 4^n$  with gaps!

# Skip fruitless alignments: Boyer-Moore

*P*: word

*T*: There would have been a time for such a word

.....word----->  
----->

*u* doesn't occur in *P*, so skip next two alignments

*P*: word

*T*: There would have been a time for such a word

.....word----->  
word skip!  
word skip!  
word

extending to  
approximate  
matching is  
awkward.

# Scaling in genome size will not work!

$P$ : word

$T$ : There would have been a time for such a word

-----word----->

As  $m$  &  $n$  grow, # characters comparisons grows with...

$|P| = n$     $|T| = m$

	Naïve matching	Boyer-Moore
Worst case	$m \cdot n$	$m$
Best case	$m$	$m / n$

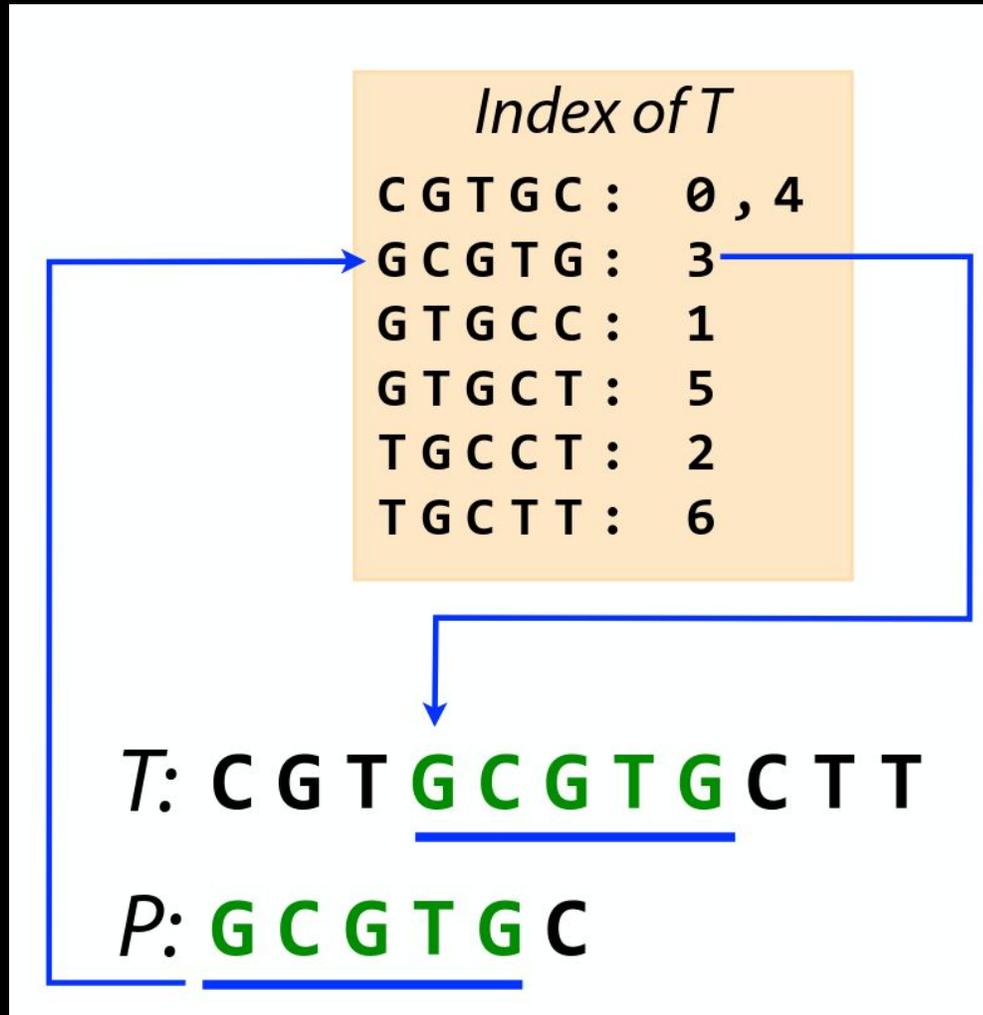
Read mapping:

- Many short reads (small  $n$ ) vs.
- Large genome (big  $m$ ).
- Multiple matches possible

BLAST has same issue

Can we use some indexing tricks to scale in  $n$  rather than  $m$ ?

# K-mer Indexing



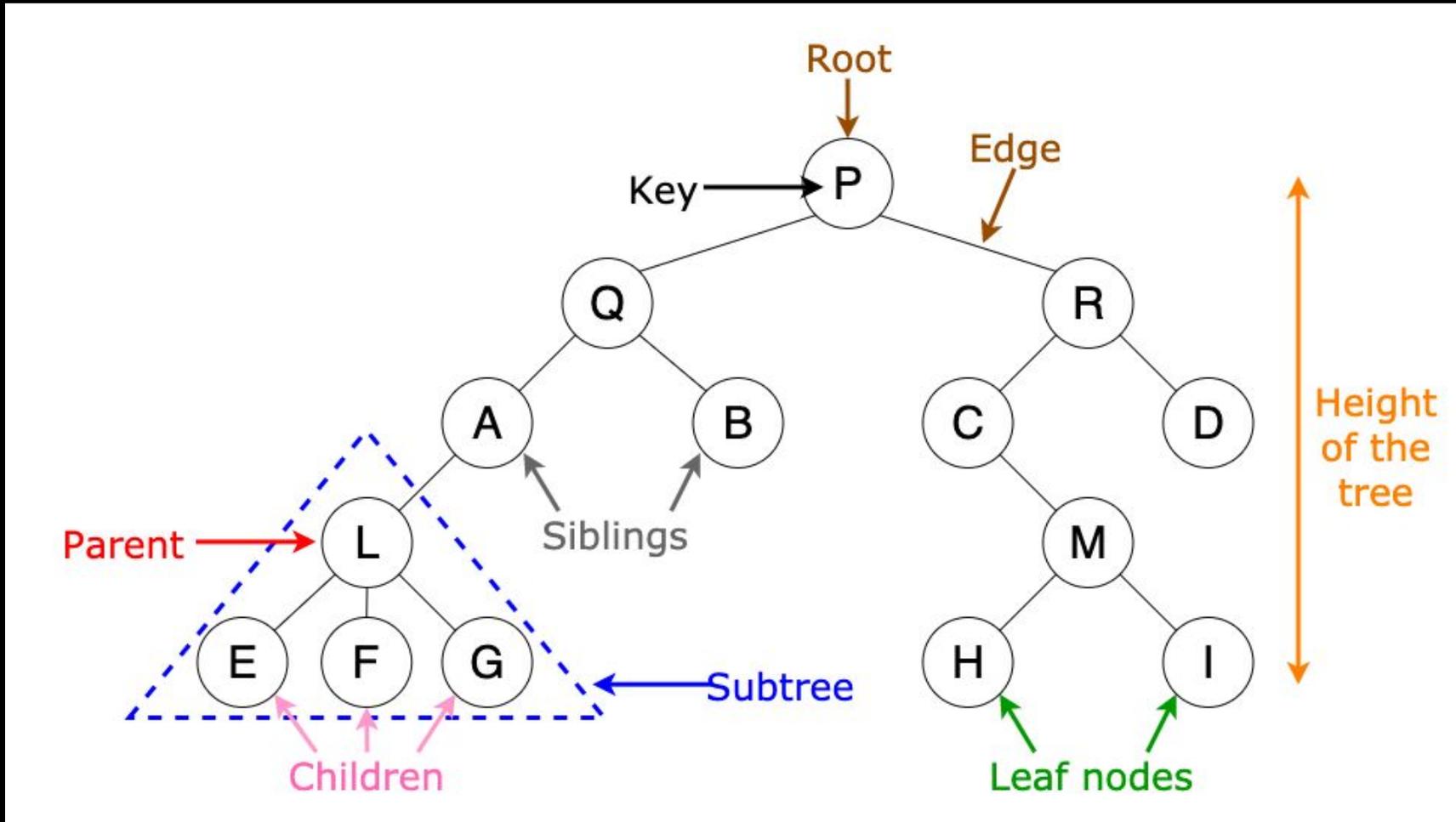
Reference indexed once and  $O(1)$  lookup

**BUT:**

- K-mers shorter than reads
- Extensions requires many lookups (neighbourhood enumeration) or expensive DP
- Repetitive K-mers mean many candidates
- Index doesn't easily capture context (adjacent K-mers)

Let's try a different way  
of indexing a genome

# Trees in Computer Science



# A trie is a special type of tree

Keys: instant, internal, internet

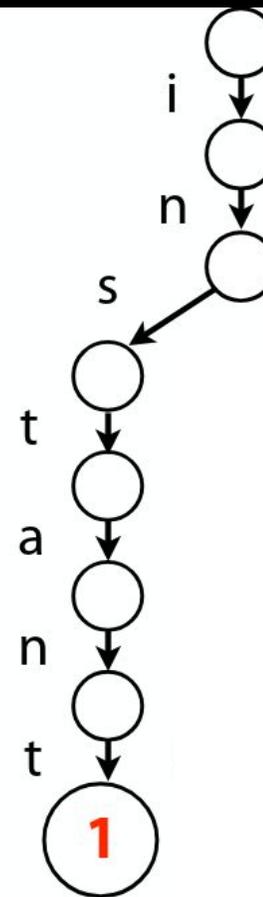
Key	Value
instant	1
internal	2
internet	3

Smallest tree such that:

Each edge is labeled with a character  $c \in \Sigma$

For given node, at most one child edge has label  $c$ , for any  $c \in \Sigma$

Each key is "spelled out" along some path starting at root



# A trie is a special type of tree

Keys: instant, internal, internet

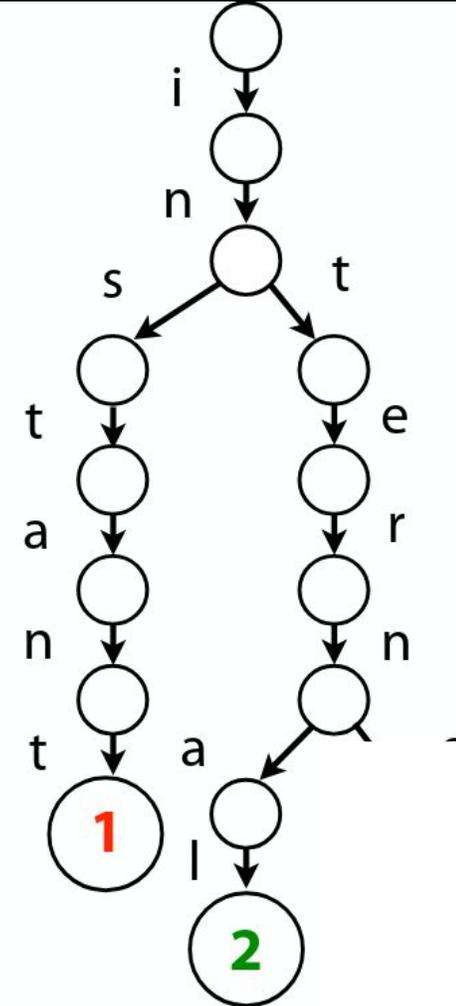
Key	Value
instant	1
internal	2
internet	3

Smallest tree such that:

Each edge is labeled with a character  $c \in \Sigma$

For given node, at most one child edge has label  $c$ , for any  $c \in \Sigma$

Each key is "spelled out" along some path starting at root



# A trie is a special type of tree

Keys: instant, internal, internet

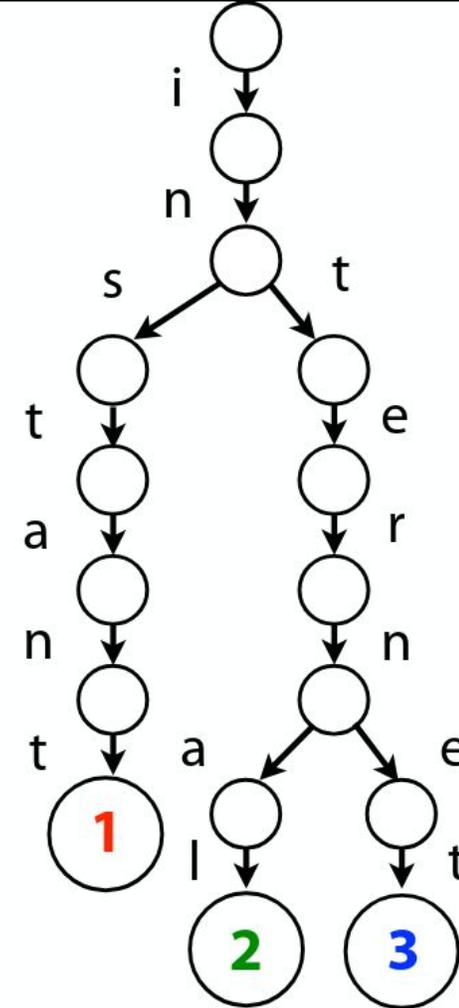
Key	Value
instant	1
internal	2
internet	3

Smallest tree such that:

Each edge is labeled with a character  $c \in \Sigma$

For given node, at most one child edge has label  $c$ , for any  $c \in \Sigma$

Each key is "spelled out" along some path starting at root



# Suffix trie: trie containing all suffixes

T = aba\$

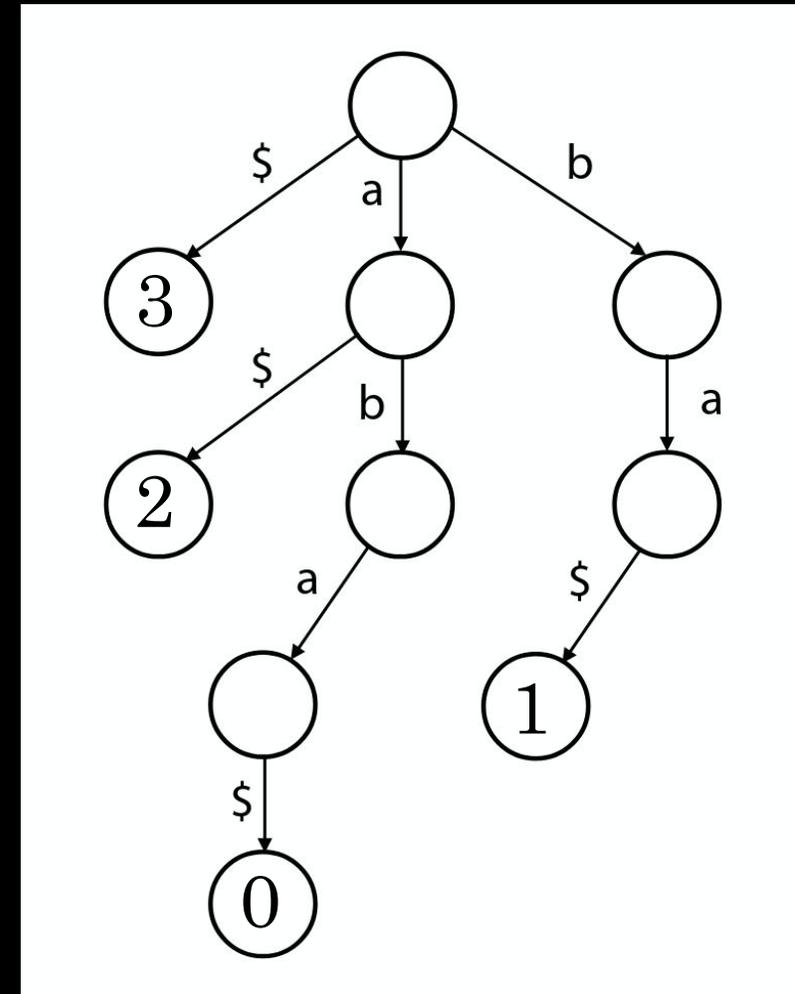
aba\$

ba\$

a\$

\$

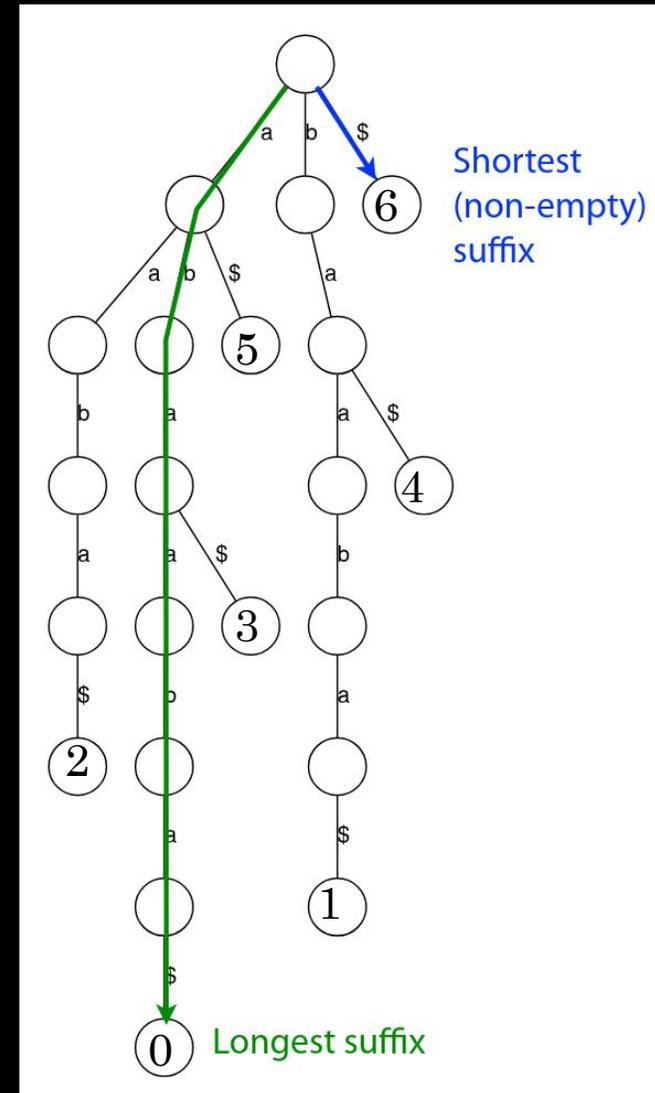
Each path from root to leaf represents a suffix; each suffix is represented by some path from root to leaf



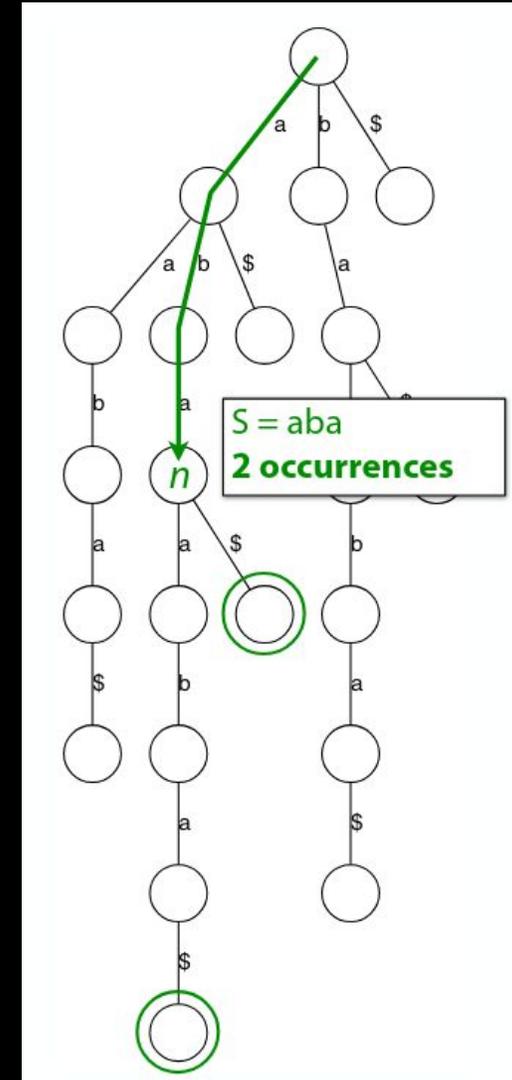
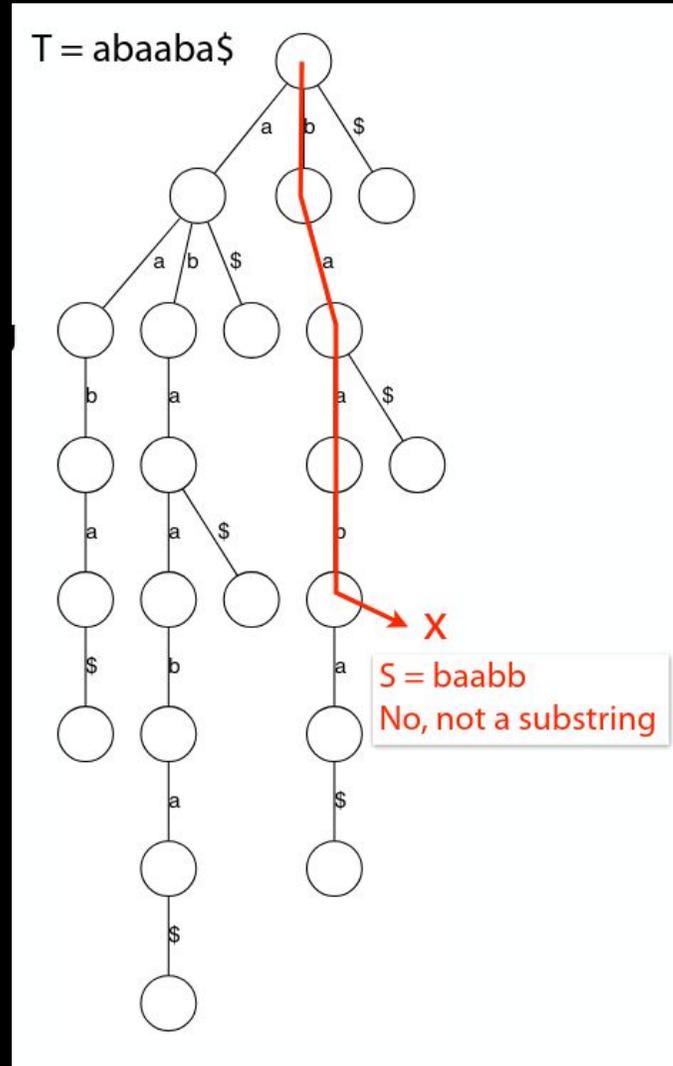
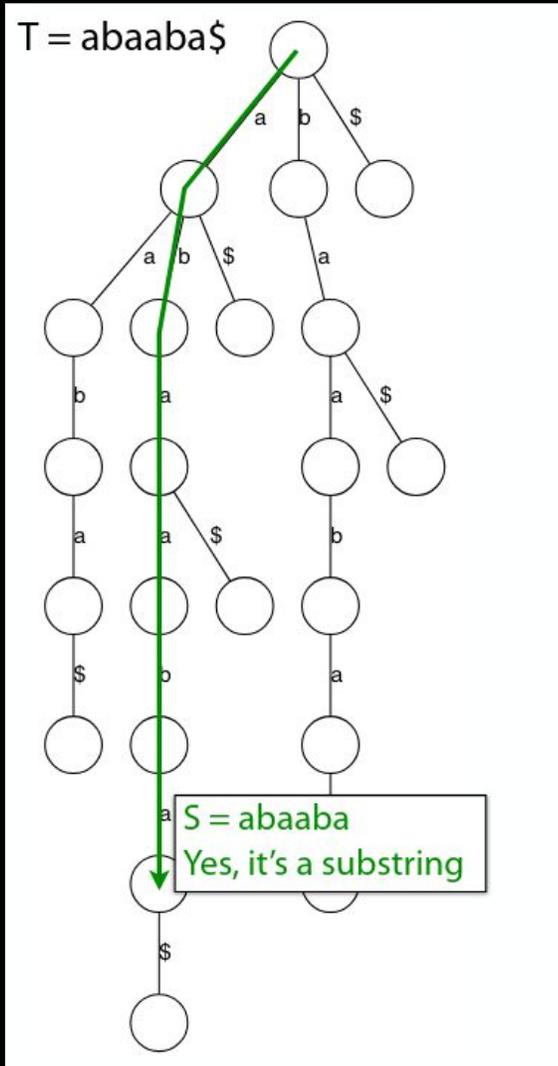
# Suffix trie: trie containing all suffixes

T = abaaba\$  
abaaba\$  
baaba\$  
aaba\$  
aba\$  
ba\$  
a\$  
\$

Each path from root to leaf represents a suffix; each suffix is represented by some path from root to leaf

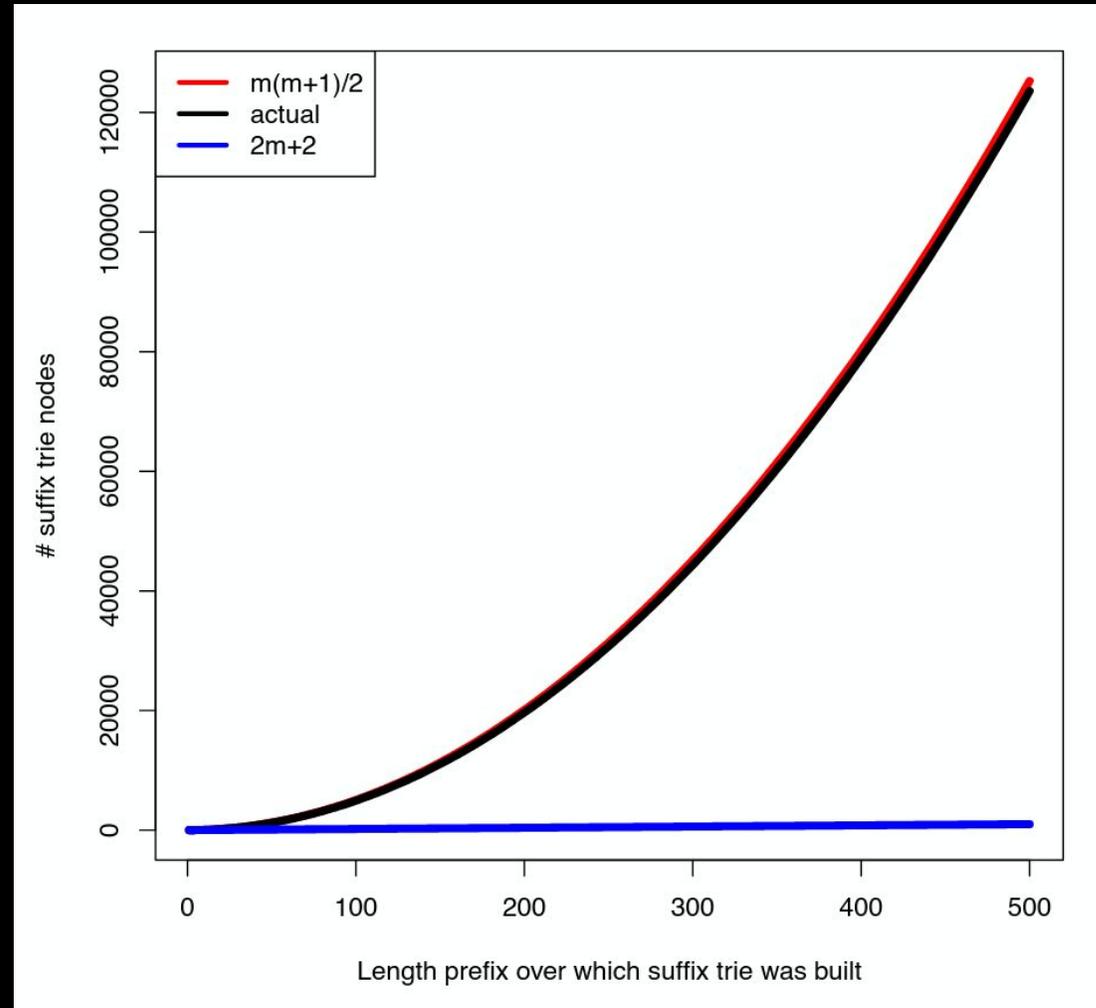


# Why? Search and count in $O(|S|)$

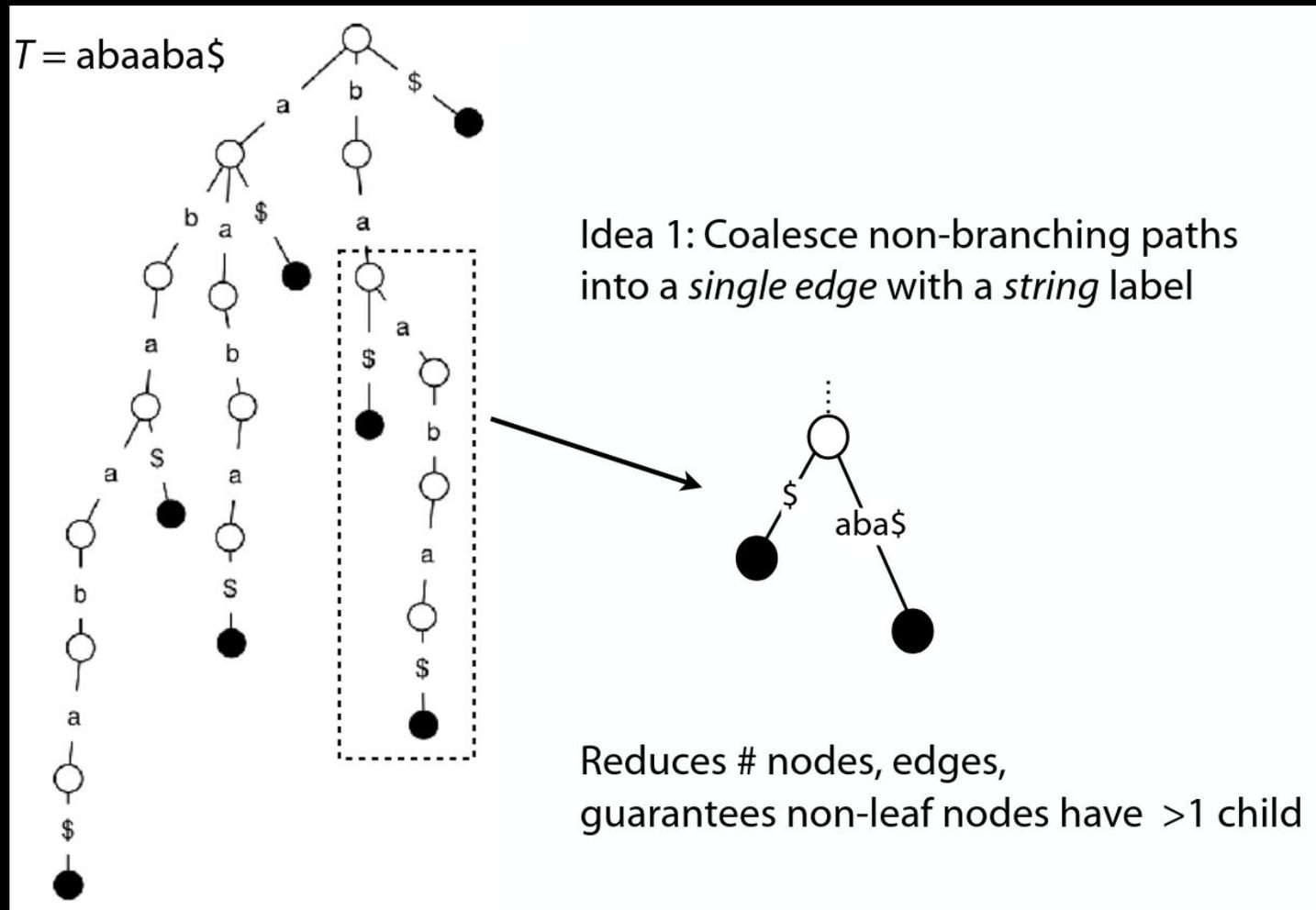


# But Suffix Trie $\sim O(m^2 / 2)$ space

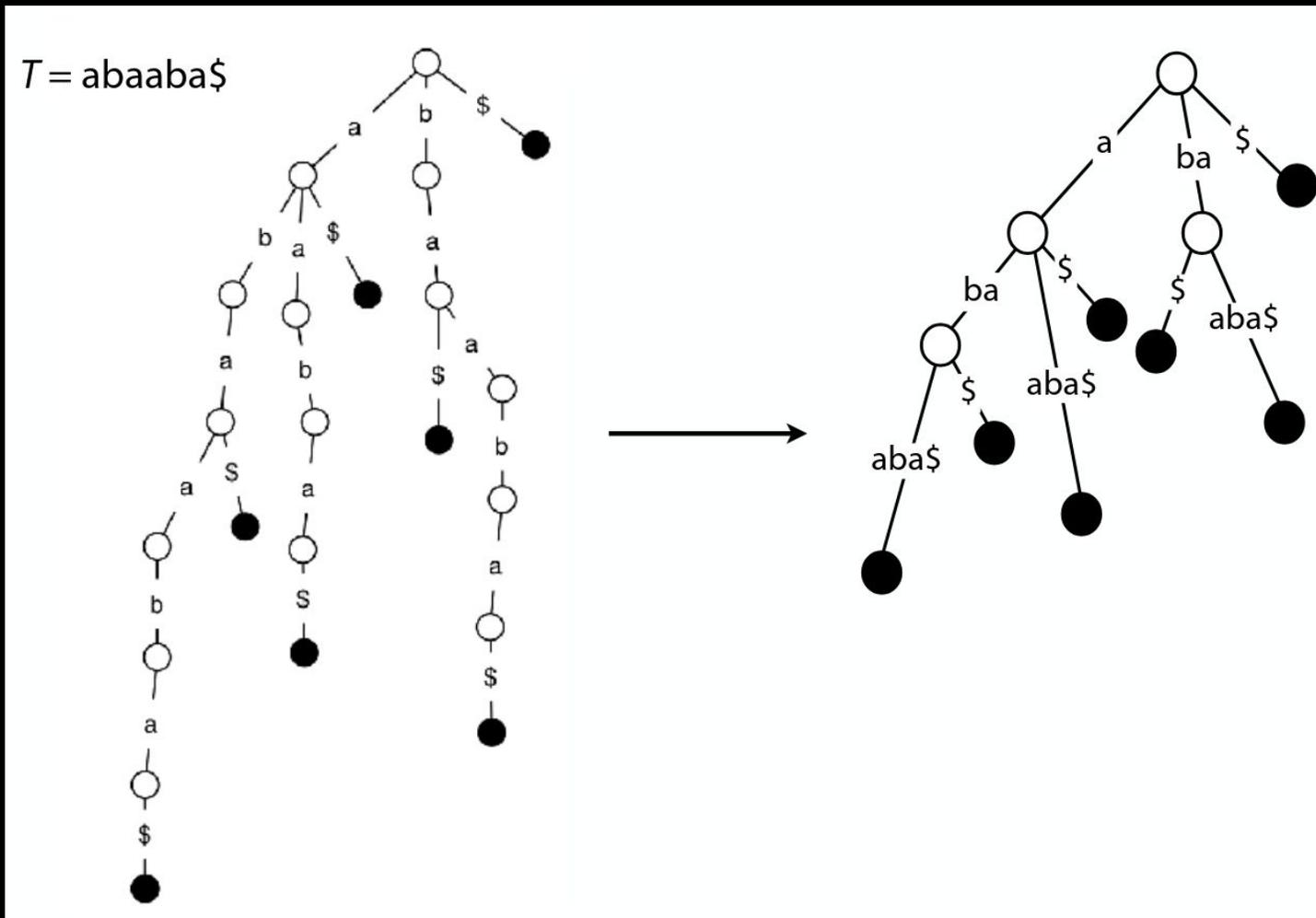
T = abaaba\$  
abaaba\$  
baaba\$  
aaba\$  
aba\$  
ba\$  
a\$  
\$



# Compressing a Suffix Trie into a Tree



# Compressing a Suffix Trie into a Tree

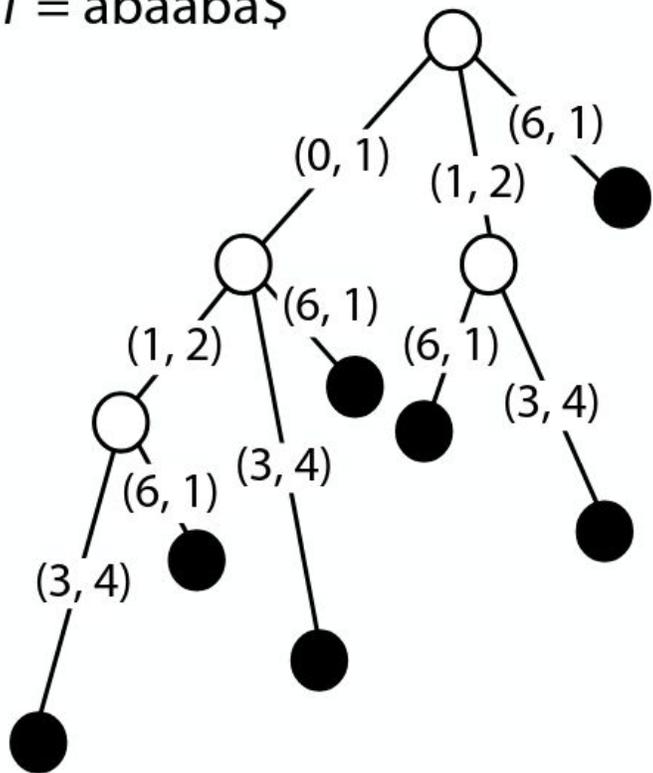


Storing suffix strings still memory intensive

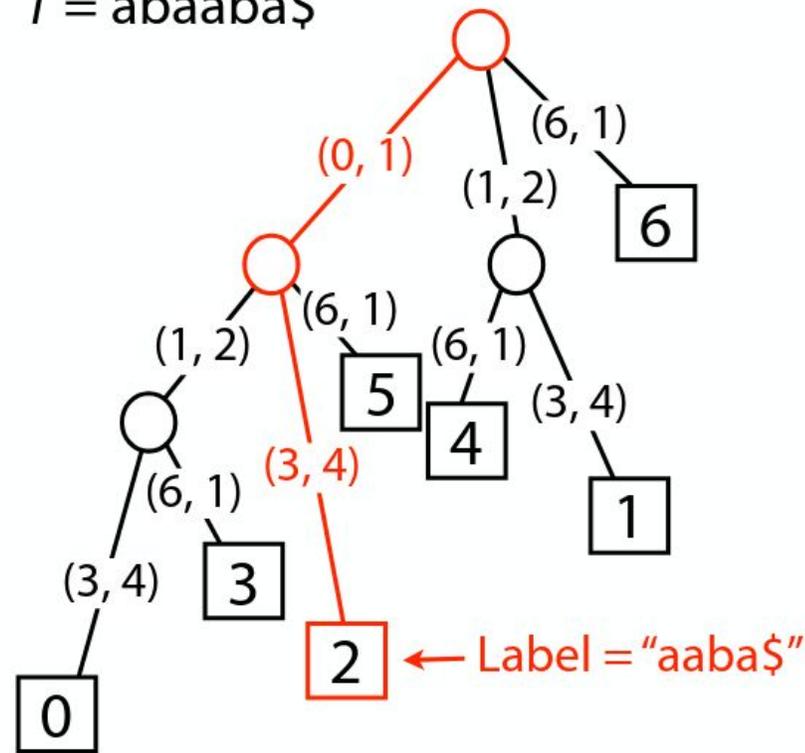


# Label leaves with offsets

$T = \text{abaaba}\$$



$T = \text{abaaba}\$$



# Retain the same operations

Suffix Trie  $\sim O(m^2/2)$  space

Suffix Tree  $\sim O(m)$  space!

BUT... don't forget constant factors!

MUMmer  $\approx 15.76$  bytes per nt

Suffix tree of human genome will still be  $>45\text{GB}$  (or larger depending on exact data structures underlying suffix tree nodes/edges)

	Suffix tree
Time: Does P occur?	$O(n)$
Time: Count $k$ occurrences of P	$O(n + k)$
Time: Report $k$ locations of P	$O(n + k)$
Space	$O(m)$

$m = |T|, n = |P|, k = \# \text{ occurrences of } P \text{ in } T$

# From suffix trees to suffix arrays

Idx	Suffixes	SA-Idx	Idx	Sorted Suffix
0	BANANA\$	0	6	\$
1	ANANA\$	1	5	A\$
2	NANA\$	2	3	ANA\$
3	ANA\$	3	1	ANANA\$
4	NA\$	4	0	BANANA\$
5	A\$	5	4	NA\$
6	\$	6	2	NANA\$

Every suffix, sorted lexicographically  
(\$ is smallest)

Index value remains the same!

# What is the point of *this*?

Idx	Suffixes	SA-Idx	Idx	Sorted Suffix
0	BANANA\$	0	6	\$
1	ANANA\$	1	5	A\$
2	NANA\$	2	3	ANA\$
3	ANA\$	3	1	ANANA\$
4	NA\$	4	0	BANANA\$
5	A\$	5	4	NA\$
6	\$	6	2	NANA\$

Suffix Array [6, 5, 3, 1, 0, 4, 2]

We now have a compressed representation where lexicographically similar strings **are adjacent to one another**

All suffixes that start with “A” are at positions 1 to 3 (because the first position is numbered 0) in the suffix array

So if we’re searching for “A” we only need to look at this **consecutive block of rows!**

Still  $O(m)$  space BUT much reduced constant factor  $\sim 4$  bytes/base \* 3bn bp  $\approx 12$  GB.  
Suffix tree is  $>45$  GB.

# Does my sequence match?

Search “ANA” against the suffix array

Idx	Suffixes	SA-Idx	Idx	Sorted Suffix
0	BANANA\$	0	6	\$
1	ANANA\$	1	5	A\$
2	NANA\$	2	3	ANA\$
3	ANA\$	3	1	ANANA\$
4	NA\$	4	0	BANANA\$
5	A\$	5	4	NA\$
6	\$	6	2	NANA\$

Suffix Array [6, 5, 3, 1, 0, 4, 2]

Binary search to find first “ANA”

Progress to last “ANA”

There are two matches of “ANA” in “BANANA”

They start at positions 3 and 1

Naive:  $O(n \log m)$ , Skipping:  $O(n + \log m)$  - smaller but slower than suffix array!

BUT... still a large index  
and inexact matching (not  
covered) is a bit painful!

Want something smaller than suffix array but with query  
time of suffix tree (and can handle inexact matches)

# Burrows-Wheeler Transform

**All Rotations**

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

**Sort Rotations**

\$ 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

**Take last column**

\$ 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

$\$ < a < b < n$

$\text{BWT}(\text{BANANA}\$) = \text{ANNB}\$AA$

# Why Do Such a Ridiculous Thing?

F		L
\$	BANANA	A
A	\$BANAN	
A	ANA\$BAN	
A	NANA\$B	
B	BANANA\$	
N	A\$BANA	
N	ANA\$BA	

If we're doing string searching, we need to explore different rows of the matrix

The end of one row should point to the row that starts with the **same exact character** in the string. This is LF mapping

# Why Do Such a Ridiculous Thing?

F            L  
\$BANANA A  
A\$BANAN  
ANA\$BAN  
ANANA\$B  
BANANA\$  
NA\$BANA  
NANA\$BA

The first “A” in L is also the first “A” in F

# Why Do Such a Ridiculous Thing?

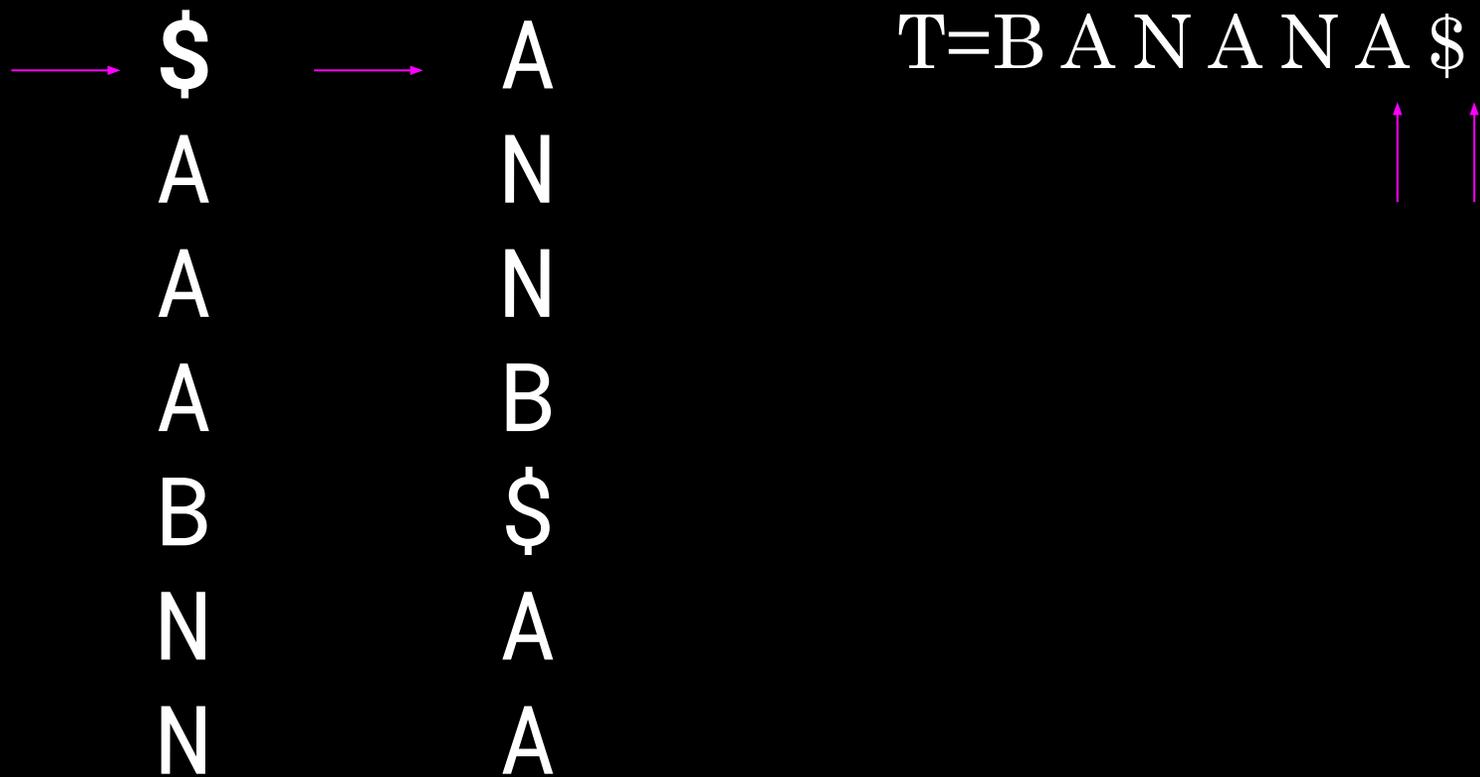
F	L
\$BANANA	A
A\$BANAN	N
ANA\$BAN	
ANANA\$B	
BANANA\$	
NA\$BANA	
NANA\$BA	

The first “N” in L is also the first “N” in F

# LF-Mapping allows reversal of BWT

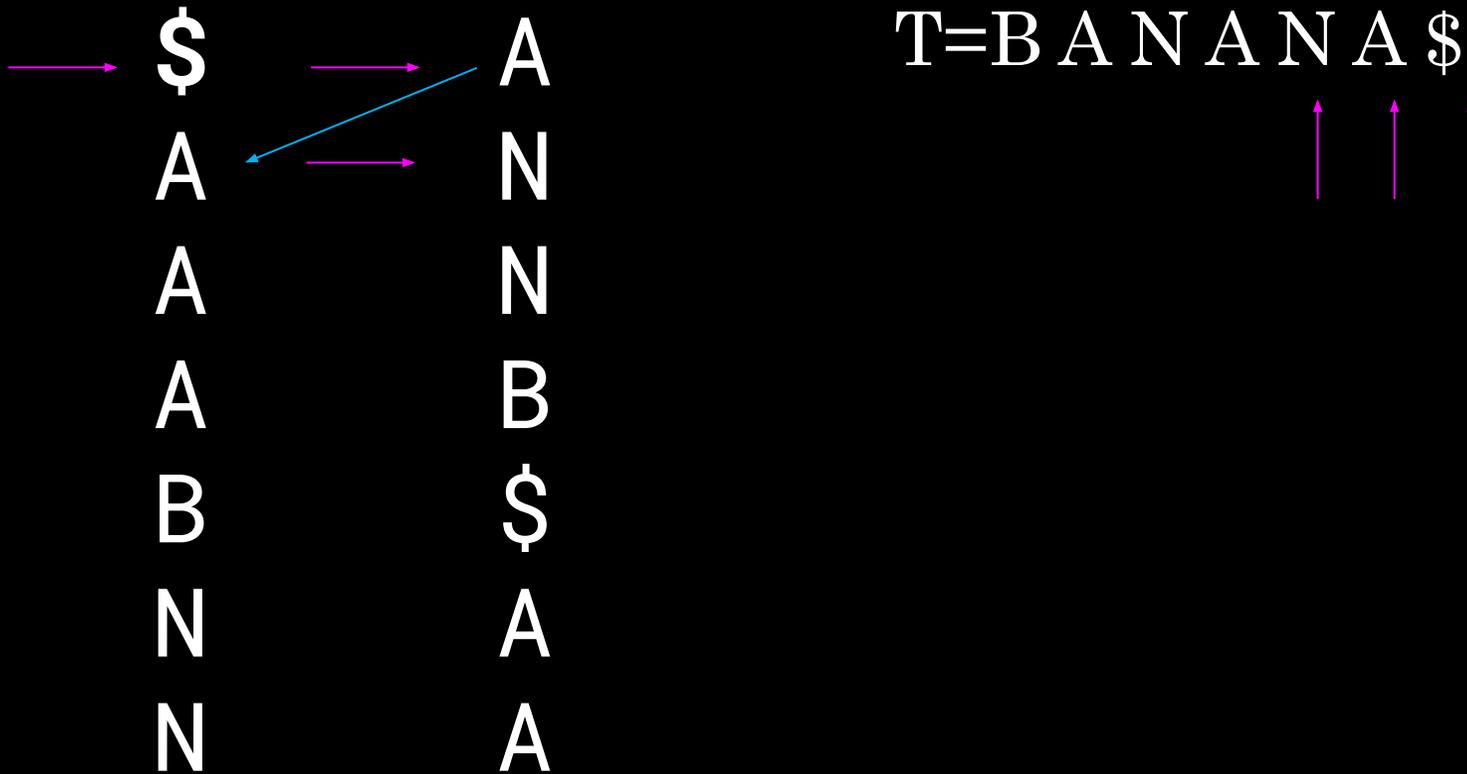
\$BANANA  
A\$BANAN  
ANA\$BAN  
ANANA\$B  
BANANA\$  
NA\$BANA  
NANA\$BA

# LF-Mapping allows reversal of BWT



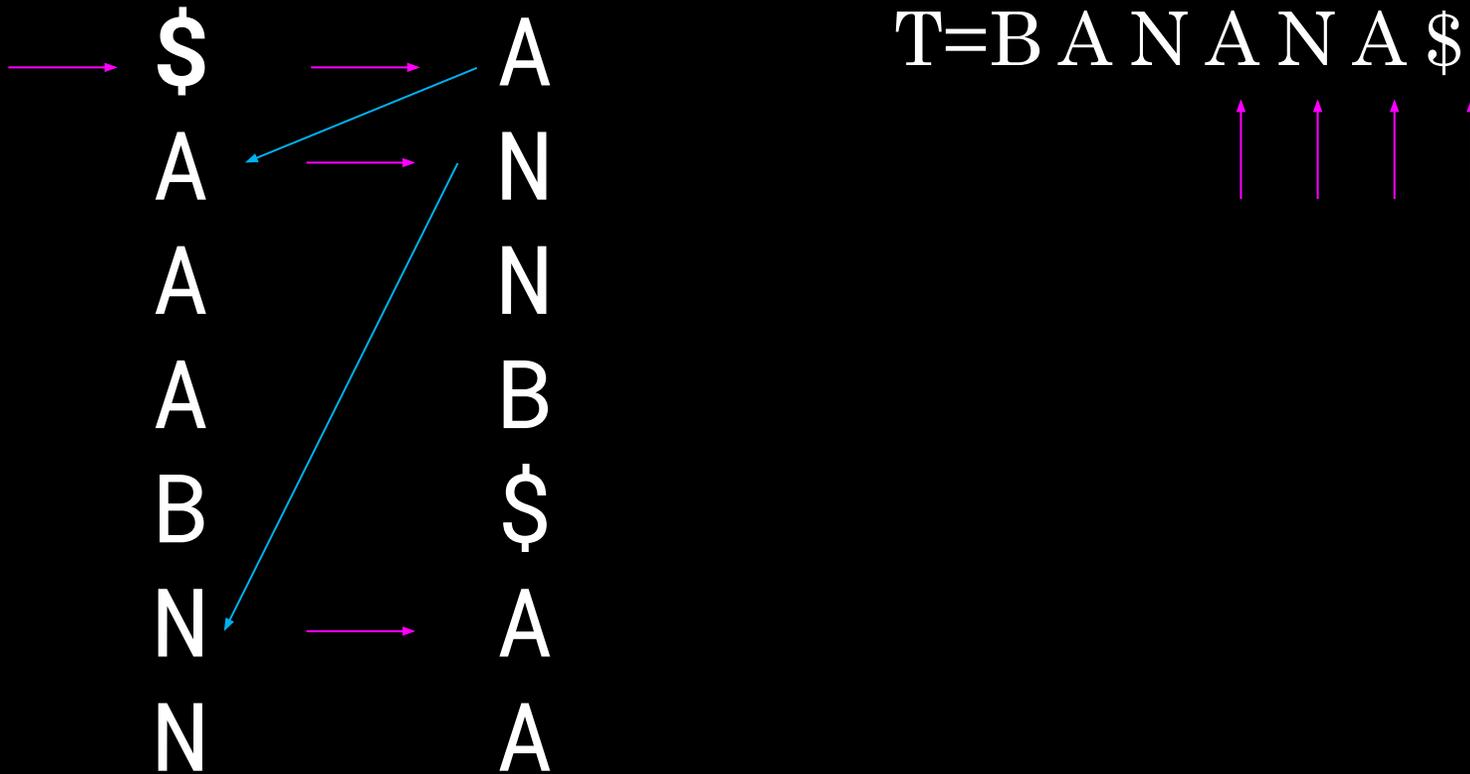
BWT(T)=ANNB\$AA

# LF-Mapping allows reversal of BWT



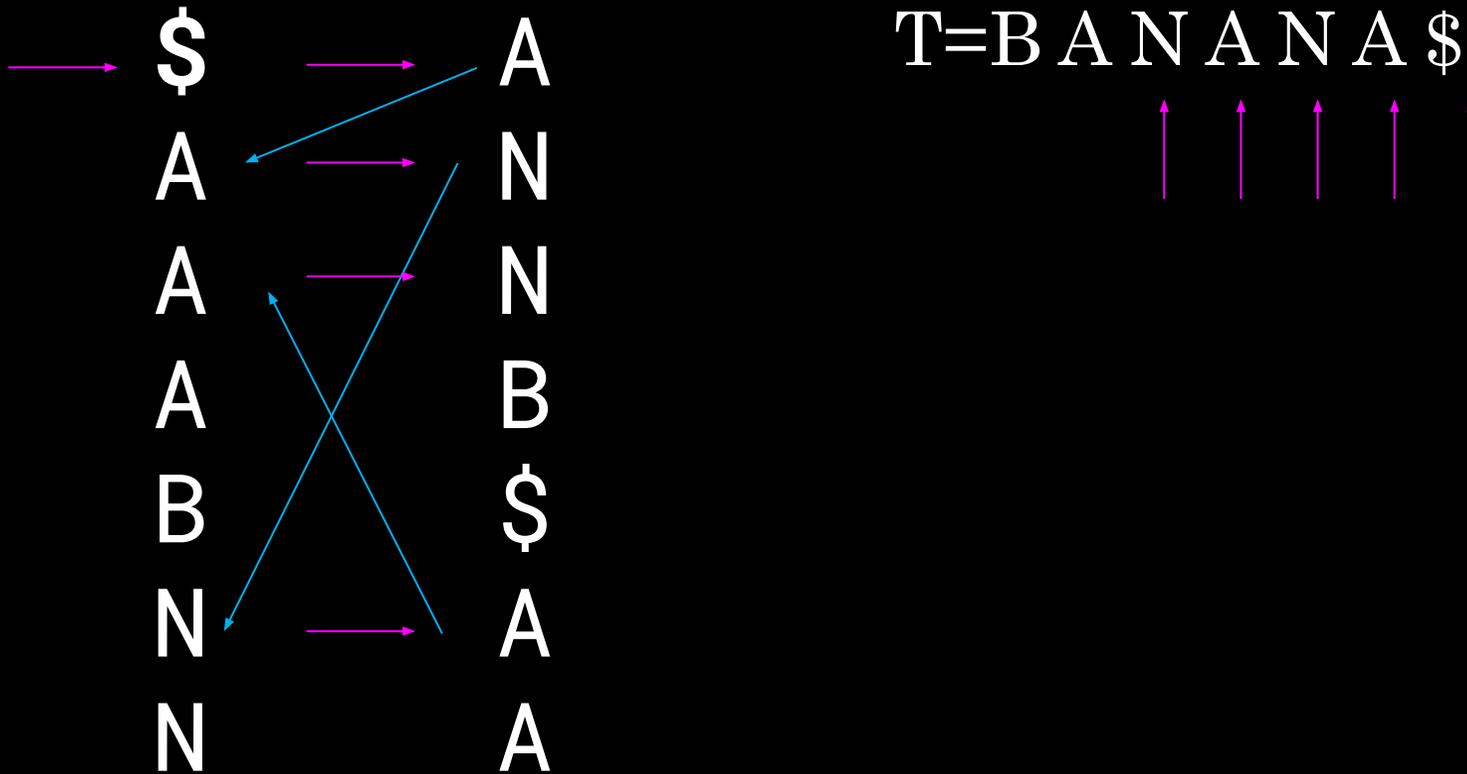
BWT(T)=ANNB\$AA

# LF-Mapping allows reversal of BWT



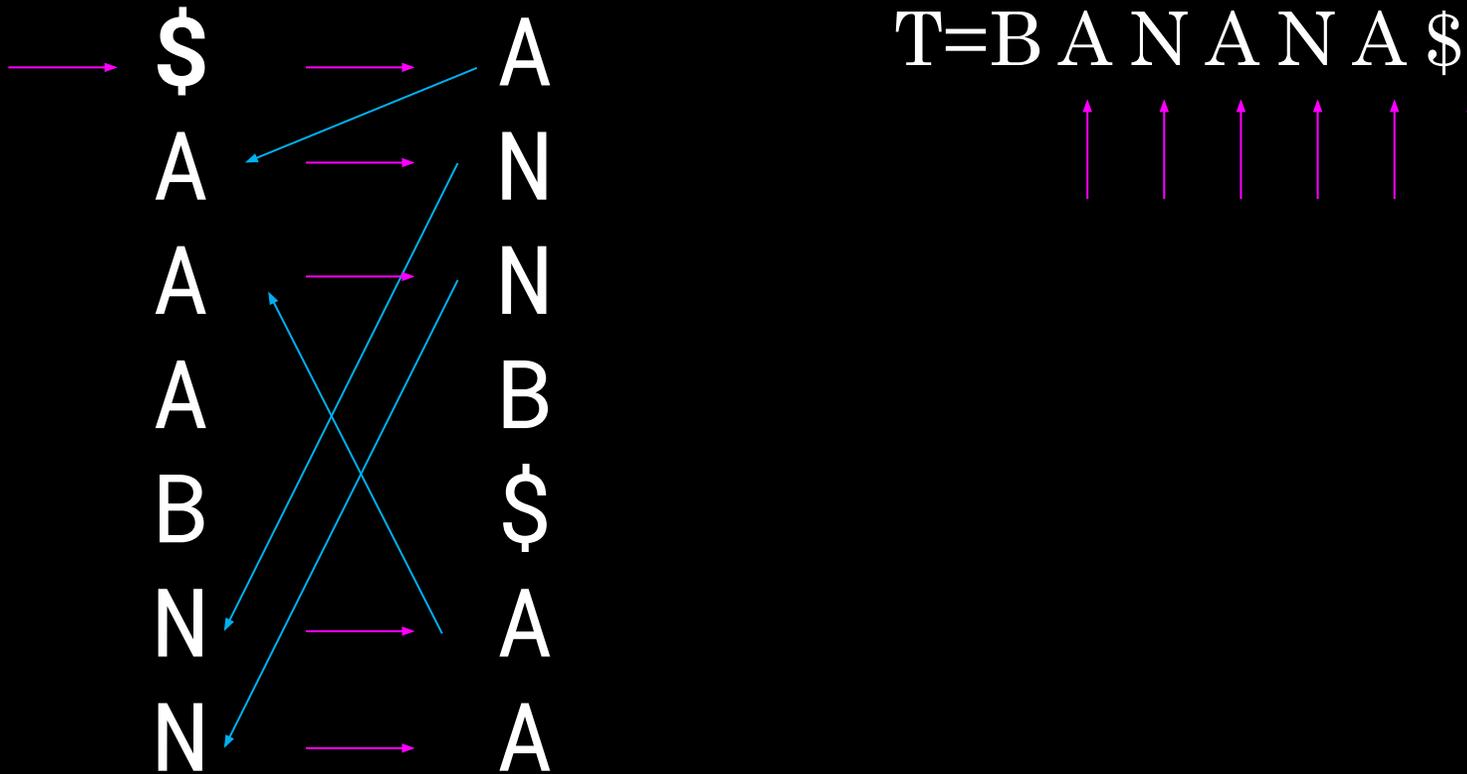
BWT(T) = ANNBS\$AA

# LF-Mapping allows reversal of BWT



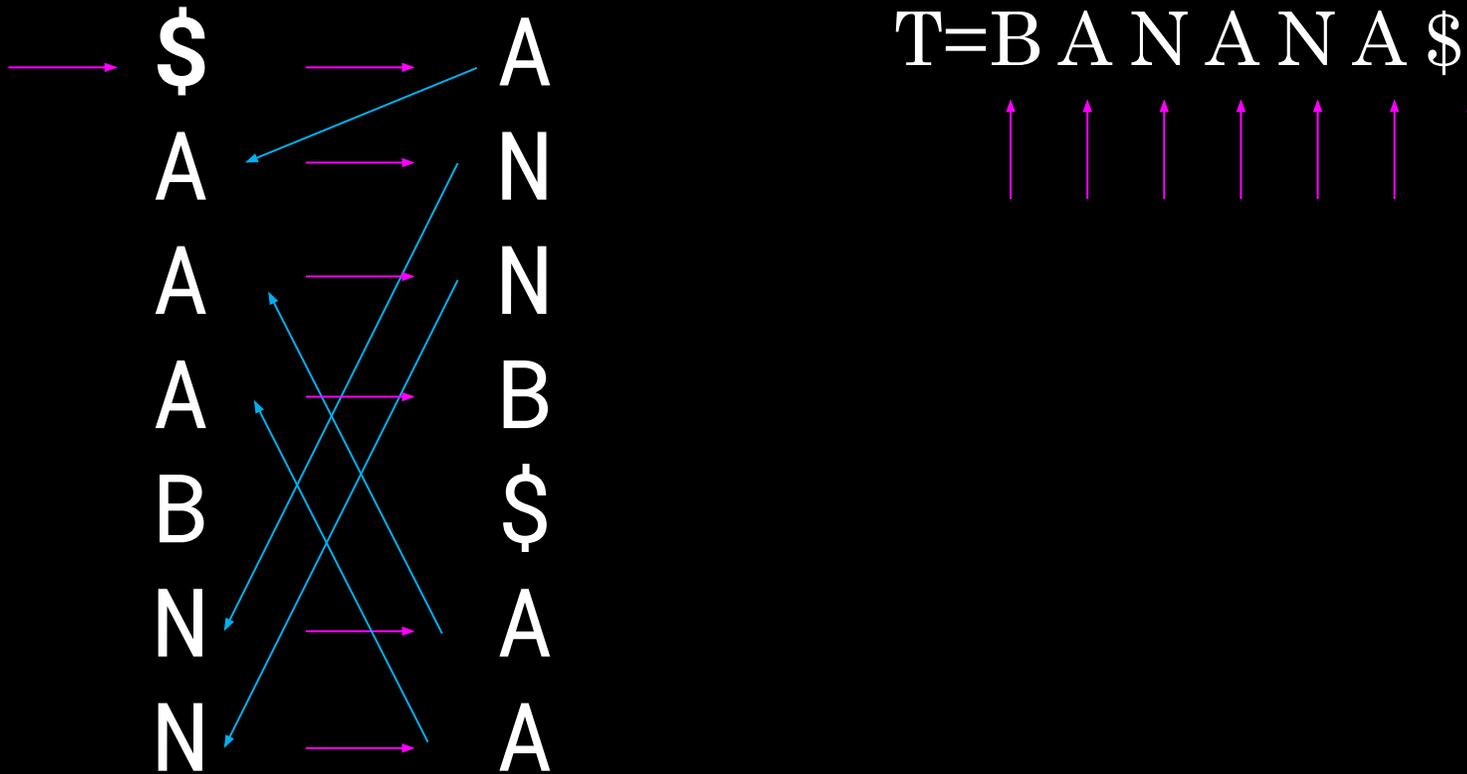
$BWT(T) = ANNB\$AA$

# LF-Mapping allows reversal of BWT



$BWT(T) = ANNB\$AA$

# LF-Mapping allows reversal of BWT



$BWT(T) = ANNBS\$AA$

# Searching BWM for a string

Look for range of rows of BWM(T) with  $P$  as prefix

Start with shortest suffix, then match successively longer suffixes

$P = \mathbf{aba}$

Easy to find all the  
rows beginning with  $\mathbf{a}$

	$F$					$L$	
	\$	a	b	a	a	b	$\mathbf{a_0}$
	$\mathbf{a_0}$	\$	a	b	a	a	$\mathbf{b_0}$
	$\mathbf{a_1}$	a	b	a	\$	a	$\mathbf{b_1}$
	$\mathbf{a_2}$	b	a	\$	a	b	$\mathbf{a_1}$
	$\mathbf{a_3}$	b	a	a	b	a	\$
	$\mathbf{b_0}$	a	\$	a	b	a	$\mathbf{a_2}$
	$\mathbf{b_1}$	a	a	b	a	\$	$\mathbf{a_3}$

# Searching BWM for a string

We have rows beginning with **a**, now we want rows beginning with **ba**

$P = \mathbf{aba}$

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

← Look at those rows in *L*.  
**b<sub>0</sub>**, **b<sub>1</sub>** are **bs** occurring just to left.

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

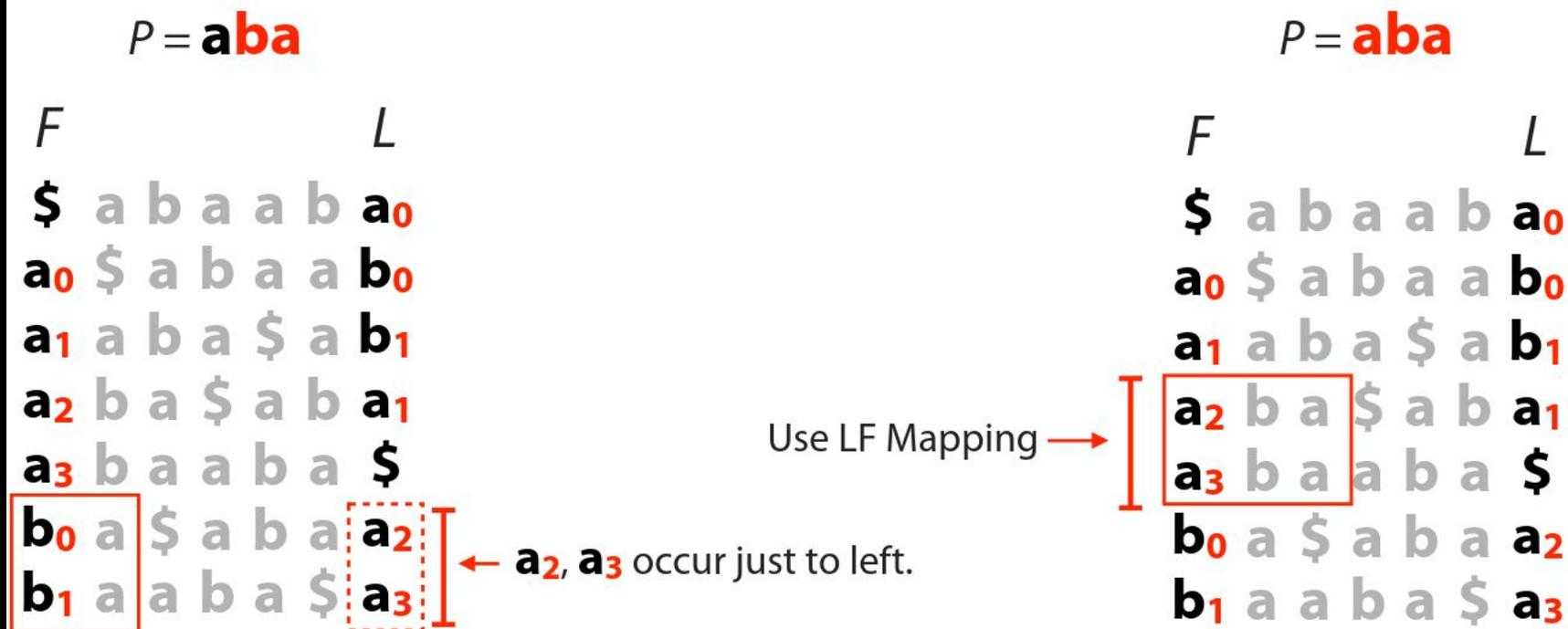
$P = \mathbf{aba}$

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

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

# Searching BWM for a string

We have rows beginning with **ba**, now we seek rows beginning with **aba**

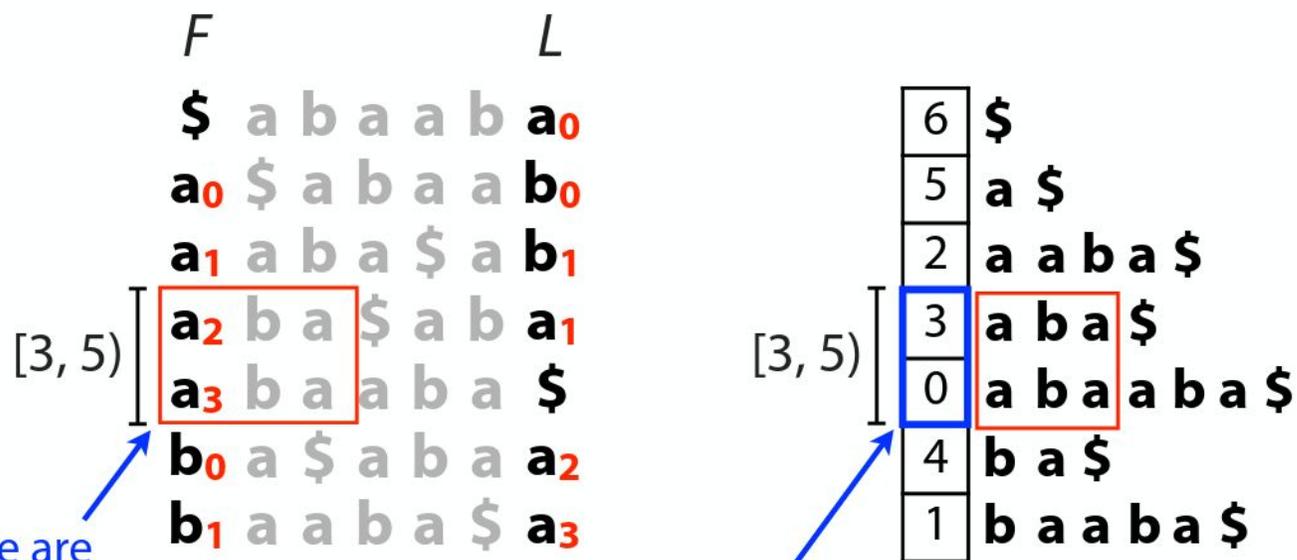


Now we have the rows with prefix **aba**

# Searching BWM for a string

$P = \mathbf{aba}$

Got the same range, [3, 5), we would have got from suffix array

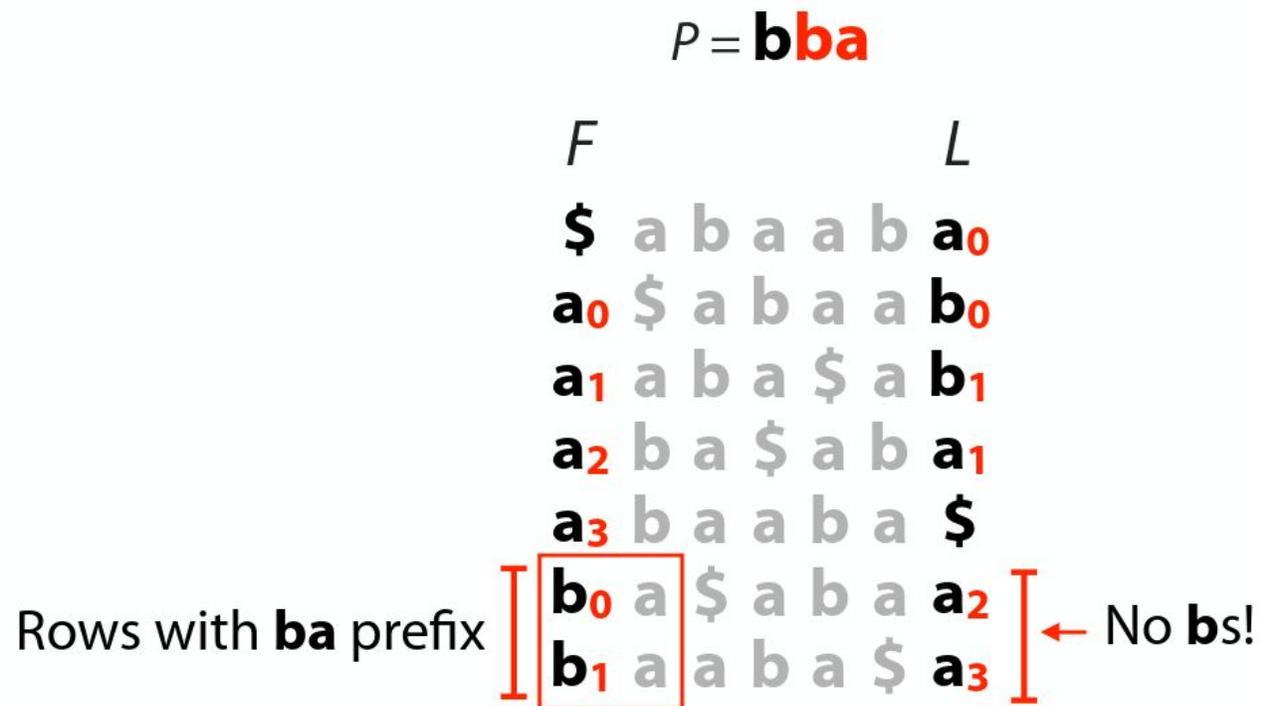


Where are these?

Unlike suffix array, we don't immediately know *where* the matches are in T...

# Searching BWM for a string

When  $P$  does not occur in  $T$ , we eventually fail to find next character in  $L$ :



# BWA: The Burrows-Wheeler Aligner

Idx	Sorted Suffix		
6	\$	\$BANAN	A
5	A\$	A\$BANA	N
3	ANA\$	ANA\$BA	N
1	ANANA\$	ANANA\$	B
0	BANANA\$	BANANA	\$
4	NA\$	NA\$BAN	A
2	NANA\$	NANA\$B	A

Rapid string search using the **sorted index** (equivalent to F) and the **BWT string** (equal to L)

# BWA: The Burrows-Wheeler Aligner

Idx	Sorted Suffix
6	\$
5	A\$
3	ANA\$
1	ANANA\$
0	BANANA\$
4	NA\$
2	NANA\$

\$BANAN	A
A\$BANA	N
ANA\$BA	N
ANANA\$	B
BANANA	\$
NA\$BAN	A
NANA\$B	A

Key parts  
of an  
FM-Index

Rapid string search using the **sorted index**  
(equivalent to F) and the **BWT string** (equal to L)  
And nothing else

**Why this is awesome:** Sequence reads are effectively searched against different parts of the reference genome at the same time

Using the L-F mapping allows us to **remove redundancy** relative to the suffix array

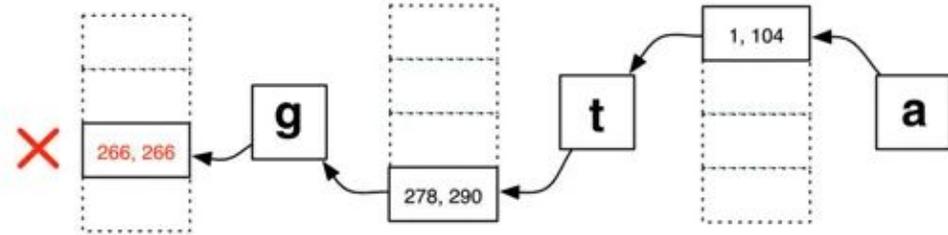
- **Why this is slightly less awesome:** Preprocessing requires many GB of memory
- What about mismatches?

# Backtracking

Search string: GGTA

“Genome” contains no GGTA, but contains GG TG

Exact



# BWA refinements

Allowing mismatches: maximum deviation from the search string

- Store searches in a heap to prioritize the lowest mismatches in the search so far

Custom penalties for mismatches, insertion and deletions

Double indexing: BWTs from both ends meet in the middle (avoids massive amounts of futile backtracing)

Memory refinements: store only parts of the BWT and supplementary data structures, calculate the rest on the fly

# Suffix tree vs Suffix array vs FM-index

	Suffix tree	Suffix array	FM Index
Time: Does P occur?	$O(n)$	$O(n \log m)$	$O(n)$
Time: Count $k$ occurrences of P	$O(n + k)$	$O(n \log m)$	$O(n)$
Time: Report $k$ locations of P	$O(n + k)$	$O(n \log m + k)$	$O(n + k)$
Space	$O(m)$	$O(m)$	$O(m)$
Needs T?	yes	yes	no
Bytes per input character	>15	~4	~0.5

$m = |T|, n = |P|, k = \# \text{ occurrences of } P \text{ in } T$

Program	Single-end			Paired-end		
	Time (s)	Conf (%)	Err (%)	Time (s)	Conf (%)	Err (%)
bowtie-125	1966	88.0	0.07	1701	91.0	0.37
BWA-125	3021	93.0	0.05	3059	97.6	0.04
MAQ-125	17506	92.7	0.08	19388	96.3	0.02
SOAP2-125	555	91.5	0.17	1187	90.8	0.14

One million pairs of 32, 70 and 125 bp reads, respectively, were simulated from the human genome with 0.09% SNP mutation rate, 0.01% indel mutation rate and 2% uniform sequencing base error rate. The insert size of 32 bp reads is drawn from a normal distribution  $N(170,25)$ , and of 70 and 125 bp reads from  $N(500,50)$ . CPU time in seconds on a single core of a 2.5 GHz Xeon E5420 processor (Time), percent confidently mapped reads (Conf) and percent erroneous alignments out of confident mappings (Err) are shown in the table.

SOAP2: VERY approximately 1000x faster than BLAST

**Table 2.** Evaluation on real data

Program	Time (h)	Conf (%)	Paired (%)
Bowtie	5.2	84.4	96.3
BWA	4.0	88.9	98.8
MAQ	94.9	86.1	98.7
SOAP2	3.4	88.3	97.5

The 12.2 million read pairs were mapped to the human genome. CPU time in hours on a single core of a 2.5 GHz Xeon E5420 processor (Time), percent confidently mapped reads (Conf) and percent confident mappings with the mates mapped in the correct orientation and within 300 bp (Paired), are shown in the table.

The

