

# Mapping on reference genomes



Stefan Simm

[simm@bio.uni-frankfurt.de](mailto:simm@bio.uni-frankfurt.de)

2019



# Literature

---

Mapping short DNA sequencing reads and calling variants using mapping quality scores. Li et al. Genome Research (2008)  
doi: 10.1101/gr.078212.108 (MAQ)

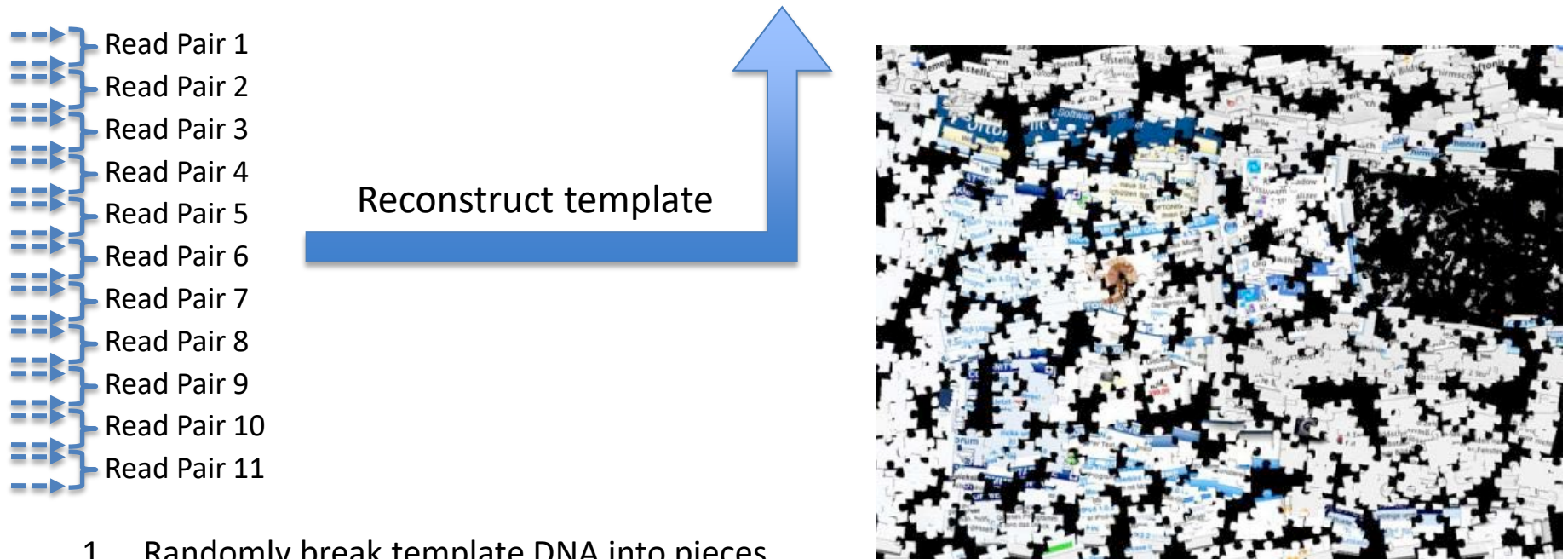
Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. Langmead et al. Genome Biology (2009)  
doi:10.1186/gb-2009-10-3-r25 (BOWTIE)

Mapping Reads on a Genomic Sequence: An Algorithmic Overview and a Practical Comparative Analysis  
Schbath et al. Journal of computational biology (2012)  
doi:10.1089/cmb.2012.0022 (Review)

NextGenMap: fast and accurate read mapping in highly polymorphic genomes. Sedlazeck et al. Bioinformatics (2013)  
doi:10.1093/bioinformatics/btt468 (NextGenMap)

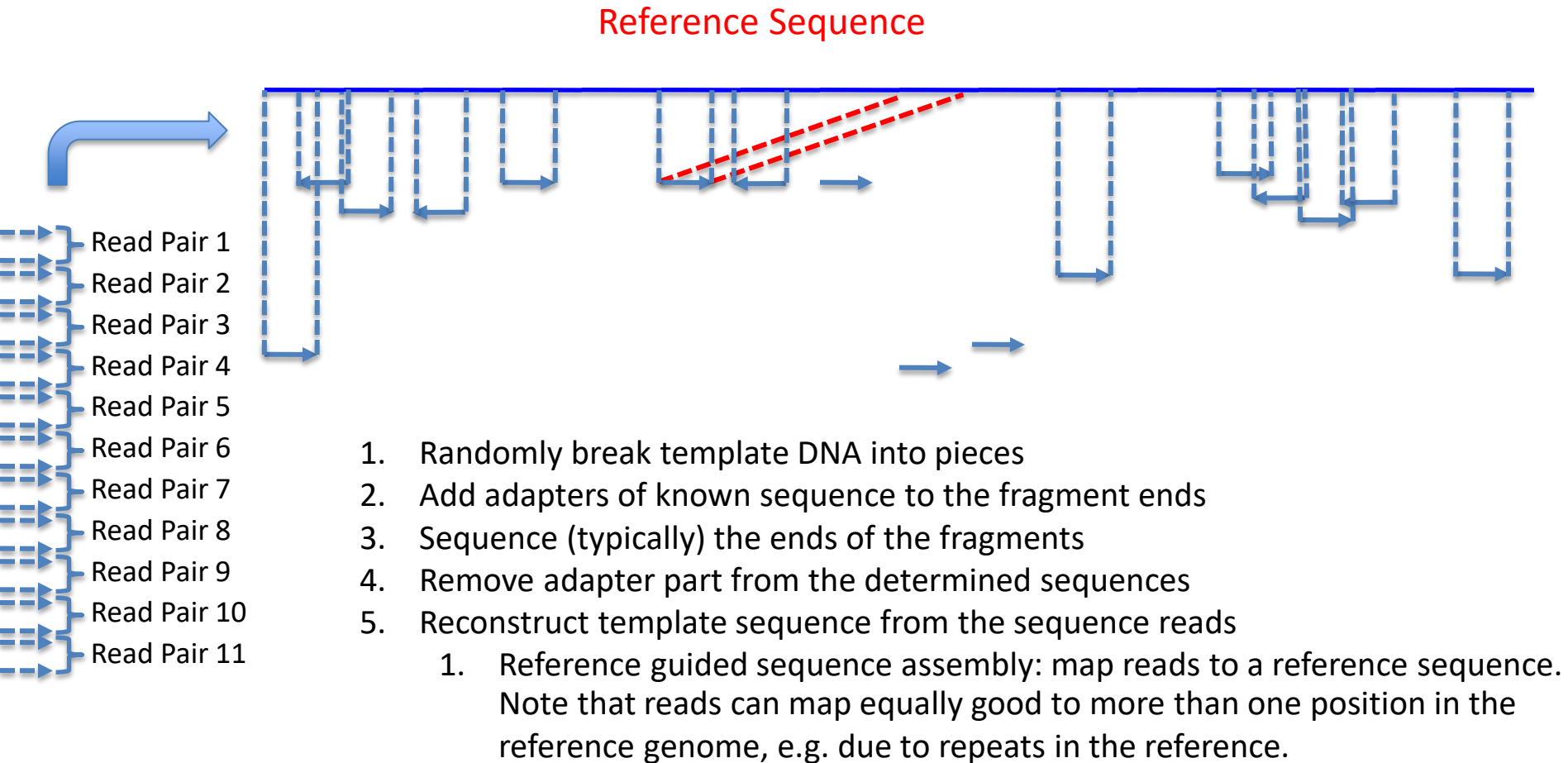
# **REFERENCE BASED MAPPING**

# Strategies to sequence long DNA molecules: Shotgun sequencing



1. Randomly break template DNA into pieces
2. Add adapters of known sequence to the fragment ends
3. Sequence (typically) the ends of the fragments
4. Identify and remove adapter part from the determined sequences
5. Reconstruct template sequence from the sequence reads

# Strategies to sequence long DNA molecules: Shotgun sequencing and **reference guided** sequence assembly

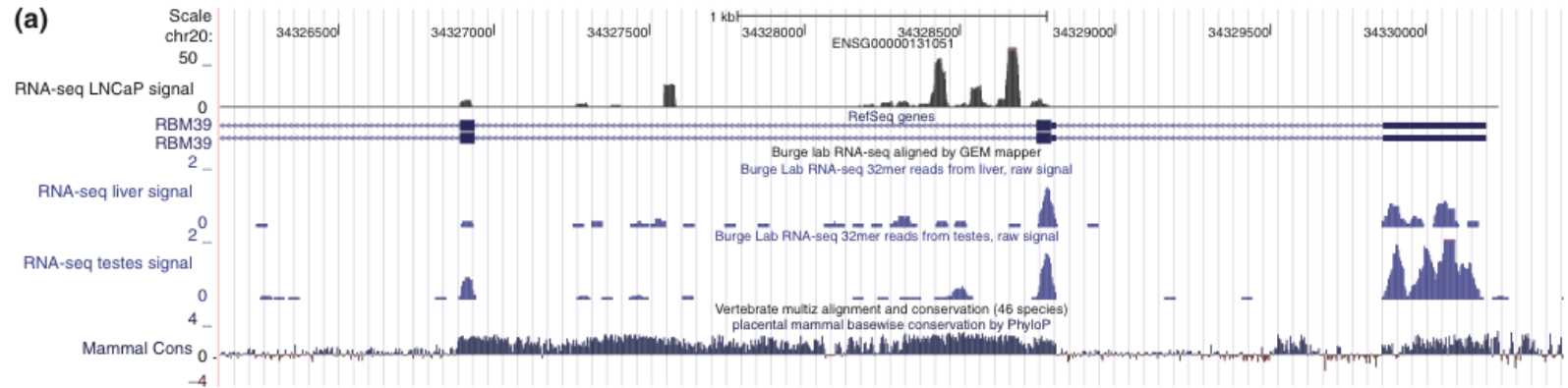


— Reference sequence

- a) e.g. genome of a different individual from the same species to study species diversity
- b) e.g. genome of a closely related species

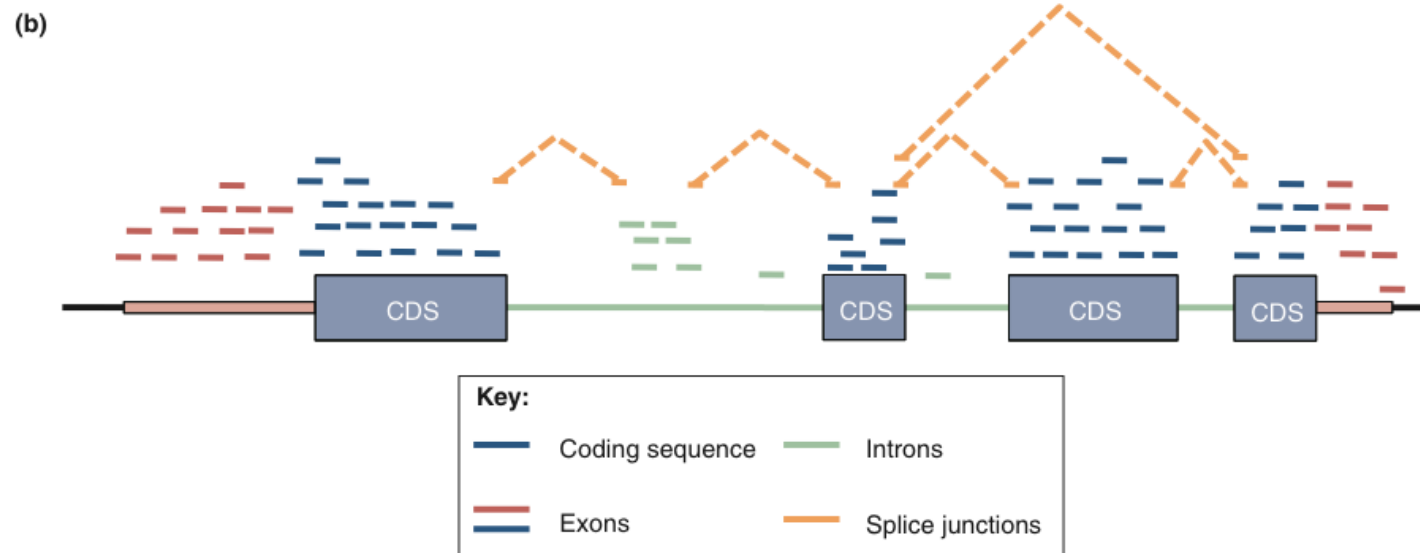
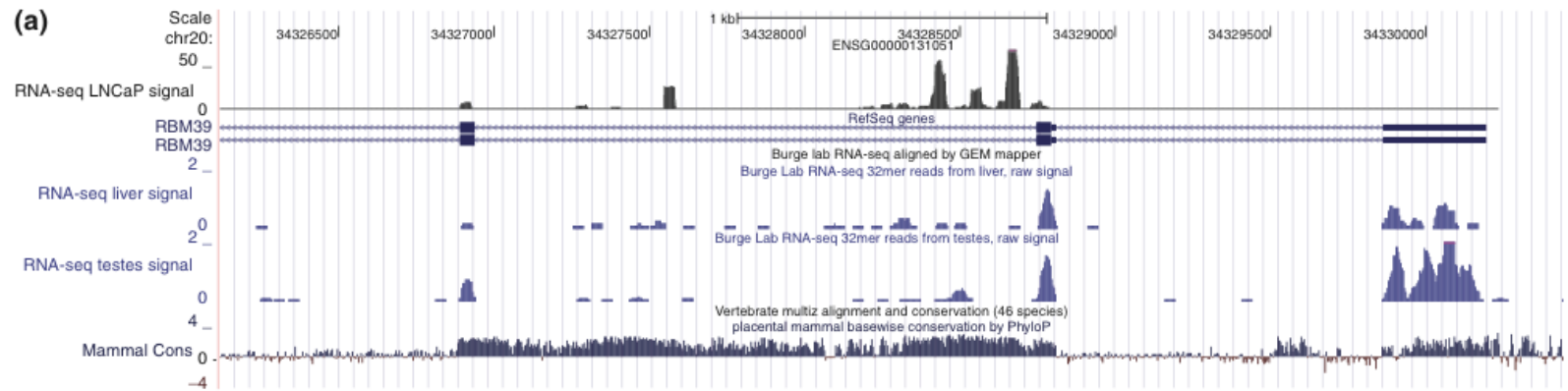


# Short Read Application: RNA seq mapping





# RNA-seq mapping helps building hypotheses concerning gene structure



# Short Read Applications

- Genotyping

Goal: identify variations

GGTATAC...

```
...CCATAG TATGCGCCC CGGA AATTTCGGTATAC
...CCAT CTATATGCG TCGG AATT CGGTATAC
...CCATGGCTATATG CTATCGG AAA GCGGTATA
...CCAAGGCTATAT CCTATCGG A TTGCGGT C...
...CCAAGGCTATAT GCCCTATCG TTTGCGGT C...
...CC AGGCTATAT GCCCTATCG A AATTTC ATAC...
...CCTAGGCTATAGCGCCCTA A AATTTC GTATAC...
```

...CCATAGGCTATATGCGCCCTATCGGCAATTTCGGTA  
TAC...

Goal: classify, measure significant peaks

- RNA-seq, ChIP-seq, Methyl-seq, Ribo-seq

GAAATTTGC  
GGAAATTTG  
CGGAAATTT  
CGGAAATTT  
TCGGAAATT  
CTATCGGAAA  
CCTATCGGA TTTGCGGT  
GCCCTATCG AAATTTGC  
GCCCTATCG AAATTTGC ATAC...

...CC

...CCATAGGCTATATGCGCCCTATCGGCAATTTCGGTATAC...



# Challenges

---

- mapping millions/billions of reads to a large genome is hard:
  - how quickly can we map the reads to the genome?
  - how do we deal with multiple mapping positions?
  - how do we deal with sequencing errors and genetic divergence/diversity
  - how do we deal with reads that span intron-exon boundaries?

# **SHORT READ ALIGNMENT**

# Short Read Alignment

---

- Given a reference and a set of reads, report at least one “good” local alignment for each read if one exists
  - Approximate answer to: from where in genome did the read originate?

## What is “good”?

- Fewer mismatches is better
- Failing to align a low-quality base is better than failing to align a high-quality base

...TGATCA**TA**... better than ...TGAT**TC**ATA...  
    ↓↓↓    ↓  
    GATCA**A**  
  
...TGAT**ATT**A... better than ...TGAT**ca**TA...  
    ↓↓↓    ↓  
    GAT**ca**T  
            ↓  
            G**T**ACAT

Finding mapping positions is, in principle, very easy!

---

Genome to search in:

AATGAGACATGAA

Reads to search:

Query1: CATG

Query2: ATGT

Finding mapping positions is, in principle, very easy.


---

Genome: AATGAGACATGAA

Query1: CATG



Naïve approach

Searchstring: AATGAGACATGAA  
CATG   
✗+++

Just slide the word along the sequence and stop when either end of sequence is reached or mapping position is found.

Finding mapping positions is, in principle, very easy.

---

Genome: AATGAGACATGAA

Query1: CATG



Naïve approach

1 2 3 4 5 6 7 8 9 1 2 3 4  
Searchstring: AATGAGACATGAA  
CATG →  
XXXXX

Just slide the word along the sequence and stop when either end of sequence is reached or mapping position is found.

Finding mapping positions is, in principle, very easy.


---

Genome: AATGAGACATGAA

Query1: CATG



Naïve approach

1 2 3 4 5 6 7 8 9 1 2 3 4  
Searchstring: AATGAGACATGAA  
CATG   
xxxx+

Just slide the word along the sequence and stop when either end of sequence is reached or mapping position is found.



# Finding mapping positions is, in principle, very easy.

---

Genome: AATGAGACATGAA

Query1: CATG



Naïve approach

1 2 3 4 5 6 7 8 9 1 2 3 4

Searchstring: AATGAGACATGAA

CATG →

+++



Query1 maps to position 8-12 in Genome

Full match found,  
Output result

Genome: AATGAGACATGAA

## Query2: TTGT



# Naïve approach

1234567891234

Searchstring: AATGAGACATGAA

TTGT  
XXXXX

At most  $n-k$  comparisons, with  $n$  is the length of the search string, and  $k$  is the query length (read length).

This is not feasible for short read mapping.

## Run Time of the Naïve approach

---

Naive approach:

- $O(L_G L_r N_r)$
- $L_G \rightarrow$  Size of the genome sequence
- $L_r \rightarrow$  Size of the read
- $N_r \rightarrow$  Number of the reads

Allowing gaps:

- Needleman Wunsch as dynamic programming algorithm
- Same complexity  $O(L_G L_r N_r)$

# Indexing speeds up searches

## INDEX

Adams, Jesse	60-61	Collins, Clemmons	26	Foster, Job Wesley	78
Alvarado, Hiram O.	61	Conover, Benjamin Edward	46	Foster, T. W.	79
Arnold, Dan and Benina	61-62	Conover, B. F.	46, 138	Frio County Centennial	5
Arnold, George and		Conover, Freddie Marvin	46, 138	Frio County Historical	
Agatha	62	Conover, Fred N.	46-47	Commission	4
Arnold, Henry and Cora	18-19	Conover, George W.	47	Frio County History	4
Assembly of God Church		Conover, George Washington	47	Frio County Markers	4-5
(Dilley)	17	Conover, Mac D.	47	Frio County Time Capsule	6-7
Assembly of God Church		Conover, Minnie	47	Frio Pioneer Jail Museum	6
(Pearsall)	58	Conover, William O.	47	Frio Public Library	7-8
Avant, Forrest J.	19	County, Roosevelt and Lois	69-70	Fudge, Albert Lester	79-80
Avant, James Ross	19-20	Cowden, George	70-71	Geyer, E. F.	10-11
Avant, Robert F. and		Cowley, W. B.	71-72	Goodman, Nancy Johnson	12
Florrie	20	Cox, Joseph	72	Gross, L. E.	80
Beall, Dr. J. E.	62-63	Crawford, V. T. and Mary	72, 139	Gill, Hugh E.	30-31
Bennett, John	63	Crocker, Minnie	72	Grueser, Winifred	31
Berry, James E.	63-64	Crossette, Joanne	147	Hardcastle, Henry J.	80-81
Berry, Mrs. James E.	64	Cude, Frank	10	Harrell, William D.	11
Betts, Nena Ward	64	Cude, Lowell	10	Harris, W. A. and Nancy	138
Bigfoot Baptist Church	9	Cude, Willis Franklin	72-73	Harris, J. Will	31
Bigfoot Church of		Dalkowitz, Harry and Evelyn	73-74,	Harris and Hinder	81
Christ	9		140	Haynes, Wynn Worsham	81-82
Bigfoot Methodist		Danchak, John Michael	74	Henson, William F.	50
Church	9	Danchak, Joseph Franklin	74-75	Herron, James	82
Bilhartz, August and		Davenport, Mary	75	Higdon, George J.	82-83
Tessie	64	DeVilbiss, Luther	75-76, 150	Higdon, Grady and Ruth	83
Bilhartz, Joseph and		DeWoody, Alonzo	26-27	Higdon, Jasper	83
Hortense B.	64-65	DeWoody, T. V.	27	Hiler, Daniel I.	31-32
Blackaller, James		Dillard, Nathan	27	Hiler, W. S.	83-85
Harrison	65	Dilley Church of Christ	17	Hiler, James H.	85
Boon, William	65-66	Dilley Presbyterian Church	18	Hiler, Lee Reavis	85-86
Boswell, Miss Bird	66	Dromgoole, Glenn A. and		Hinojosa, Miguel	86
Boyd, William A.	9-10	Minnie L.	76	Holland, Albert Green	86-87
Braun, John George	45-46	Dubose, Ruby K.	48	Holmes, Joseph M.	11-12
Breazeale, Morris H.	66-67	Dunn, Claude Bernard	27-28	Howard, Dr. and Mrs. E. M.	87, 149
Brown, Bernard *	20-21, 137	Dunn, Hubert	28	Howard, Earl Winfield	87-88
Brown, Charles A.	21	Dunn, John Anthony	28	Howerton, James	32-33
Brown, Frank S.	21	Dunn, Oscar James	28	Immaculate Heart of Mary	
Brown, James Gideon	21-22	Dunn, Patrick	28	Parish	59
Brown, Milton B.	22	Edwards, Eliza H. Crain	48	Jeffries, William	88
Brown, Paul P.	22	Edwards, Levi J. Wangruder	48	Johnson, Quill	88-89
Busby, Alex and Gladys	22	Elkins, J. W.	28-29	Jones, Charles Calvin	88
Busby, Claude L.	22-23	Ellis, Houston H.	29	Jones, Floyd G.	89
Busby, Robert	147	Ellis, S. H.	29-30	Jones, G. B. and Michealna	89-90
Busby, Roy and Lucille	23	Ellis, W. W.	30, 141	Jones, James Marion	50
Campbell, Alta	67	Fargason, Cloud O.	48-49	Jones, Jeremiah Denman	90
Campo, Luis	67-68	Fargason, John E., Jr. and		Jones, Michael Joe (Mike)	90, 150
Carmichael, Rev. Walter	68	Mildred	49	Jordan, Sam H.	90-91
Carroll, Barney R.	23-24	Fargason, John E., Sr.	49-50	Karnes, Blanche Crawford	91
Carroll, Charles L.	24	Fields, John E.	76	Kemper, John	91-92
Carroll, James Henry	24	First Baptist Church		Kimball, Rev. James Floyd	92-93
Carroll, Josephine and		(Dilley)	17	King, J. J. and Mabel	33-34, 143
Geneva	24	First Baptist Church		Klopek, Albert and Lucille	34, 145
Carroll, Louie Joseph, Jr.	24	(Pearsall)	58	Loc, Charles Edward	93
Carroll, Louie Joseph, Sr.	24-25	First Christian Church		Lester, Owen E.	50-51
Carroll, Paul Eugene	25	(Pearsall)	58	Lindholm, John Nelson	93-94
Carroll, Wesley Raymond	25	First United Methodist		Lindholm, Mady	94
Carroll, Wesley Robert	25	Church (Dilley)	17-18	Lippard, David	94-95
Carroll, William Patrick	25	First United Methodist		Little, Brice	95
Cavender Family	25-26	Church (Pearsall)	59-60	Littleton, Sam and Ora	51
Church of Christ (Pearsall)	58-59	Fitch, R. D.	76-77	Long, R. G.	95-96
Clausewitz, Martin	68	Flores, Ted	77	Lowe, Judge Marcellus	96
Coleman, Robert	68-69	Foreman, Barton and Bessie	77-78	Lyons, Arthur E.	96
Collins, John W. (Bill)	69	Foster, Lorena	78-79, 142	Malone, John J.	97, 148

## Indexing:

- Allow targeted search
- Genome distributed in chapter / keywords

# Indexing speeds up searches

Searchstring: CATGAGACATGAA

Query1: GAGA

Query2: CATG

Query3: ATGT

- 1) Decide on a word length  $k$ , e.g.,  $k=3$
- 2) Build hash table from search string, storing every word occurring in  $S$  together with its start position.
- 3) Process query and search for each word occurring in  $Q1$  whether it is in the hash table.
- 4) Repeat for  $Q2$ .

Kmer	Position
CAT	1,8
ATG	2,9
TGA	3,10
GAG	4
AGA	5
GAC	6
ACA	7
GAA	11

Q1.1

Kmer	Position
GAG	1
AGA	2

Q1.2

Two lookups are sufficient to find  $Q1$  in  $S$

# Indexing handling mismatches?

2) How does the mapper deal with queries that 'almost' match the reference?

Q3 matches the reference with one mismatch

Relevant for sensitivity and specificity of the mapping. Allowing more mismatches increases sensitivity (consider sequencing error and genetic diversity) but decreases specificity (more false positives).

Kmer	Position
CAT	1,8
ATG	2,9
TGA	3,10
GAG	4
AGA	5
GAC	6
ACA	7
GAA	11

Q3.1

Kmer	Position
ATG	1
TGT	2

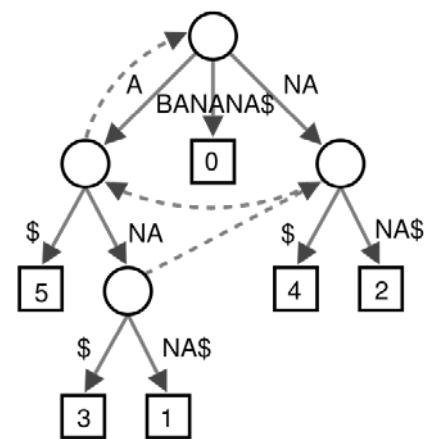
One mismatch

The 2<sup>nd</sup> lookup indicates that Q3 is almost in S

# Main differences between mapping approaches

3) What kind of index does the mapper use?

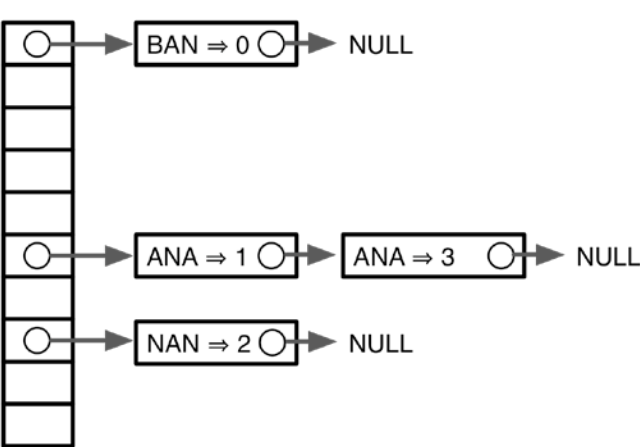
Relevant for speed and memory footprint of the mapper



**Suffix tree**

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

**Suffix array**



**Seed hash tables**  
Many variants, incl. spaced seeds

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

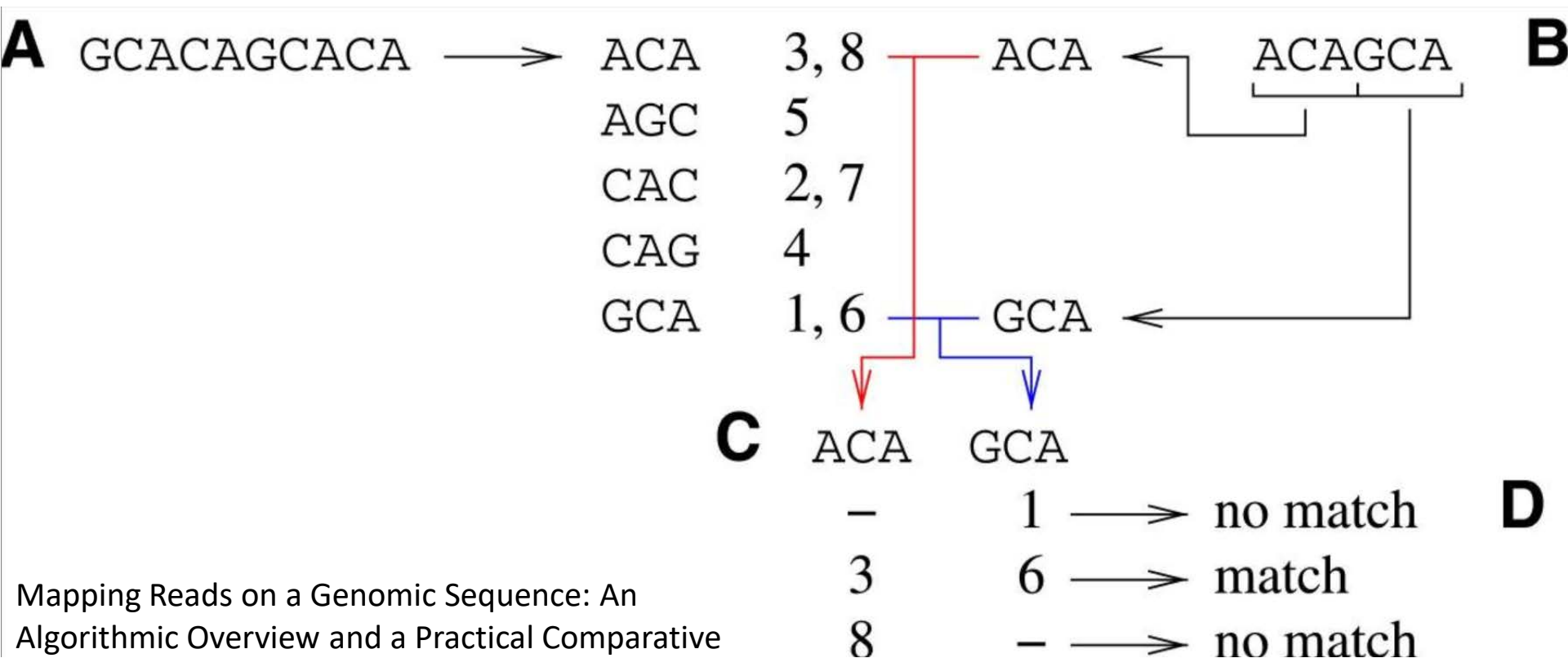
**Burrows Wheeler Transformation**



# **HASH-BASED MAPPING**

# The “traditional way”: Hash tables

- Used by MAQ, Eland, SOAP, SHRiMP, ZOOM, partially by Mosaik, SSAHA2, Stampy

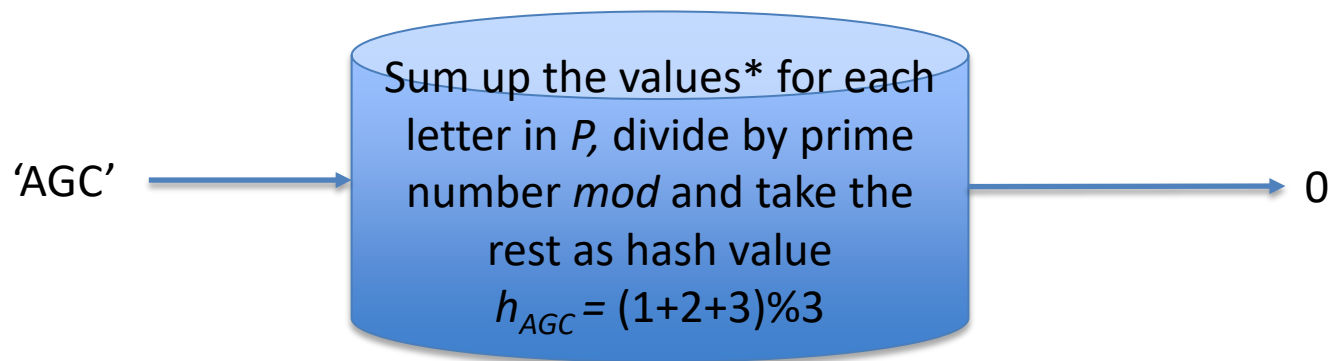


# The use of 'hashing' in exact pattern search (Rabin-Karp; $O(n+m)$ )

---

## Approach:

- 1) Use a 'hash-function' to transform pattern  $P$  into a numerical hash value  $h_P$ .

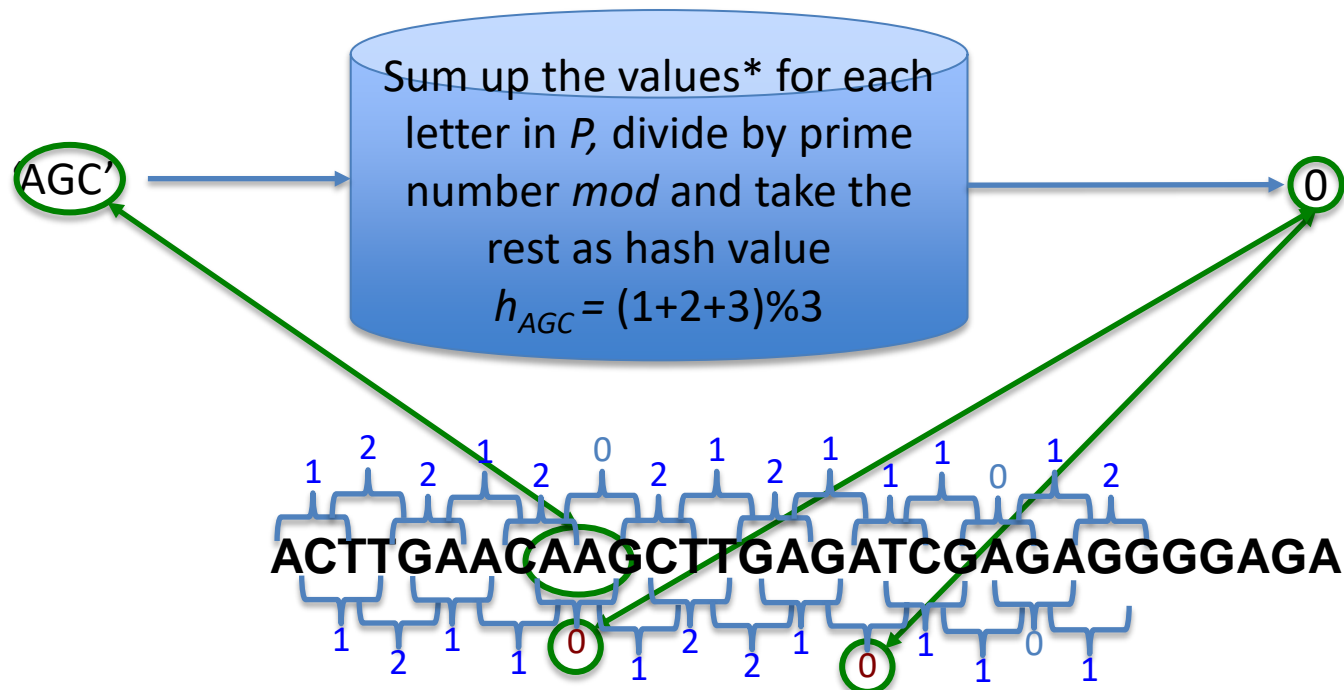


\*JUST AN EXAMPLE:  $v(A)=1$ ,  $v(C)=2$ ,  $v(G)=3$ ,  $v(T)=4$ ;  $mod=3$

# The use of 'hashing' in exact pattern search (Rabin-Karp; $O(n+m)$ )

## Approach:

- 1) Use a 'hash-function' to transform pattern  $P$  into a numerical hash value  $h_P$ .



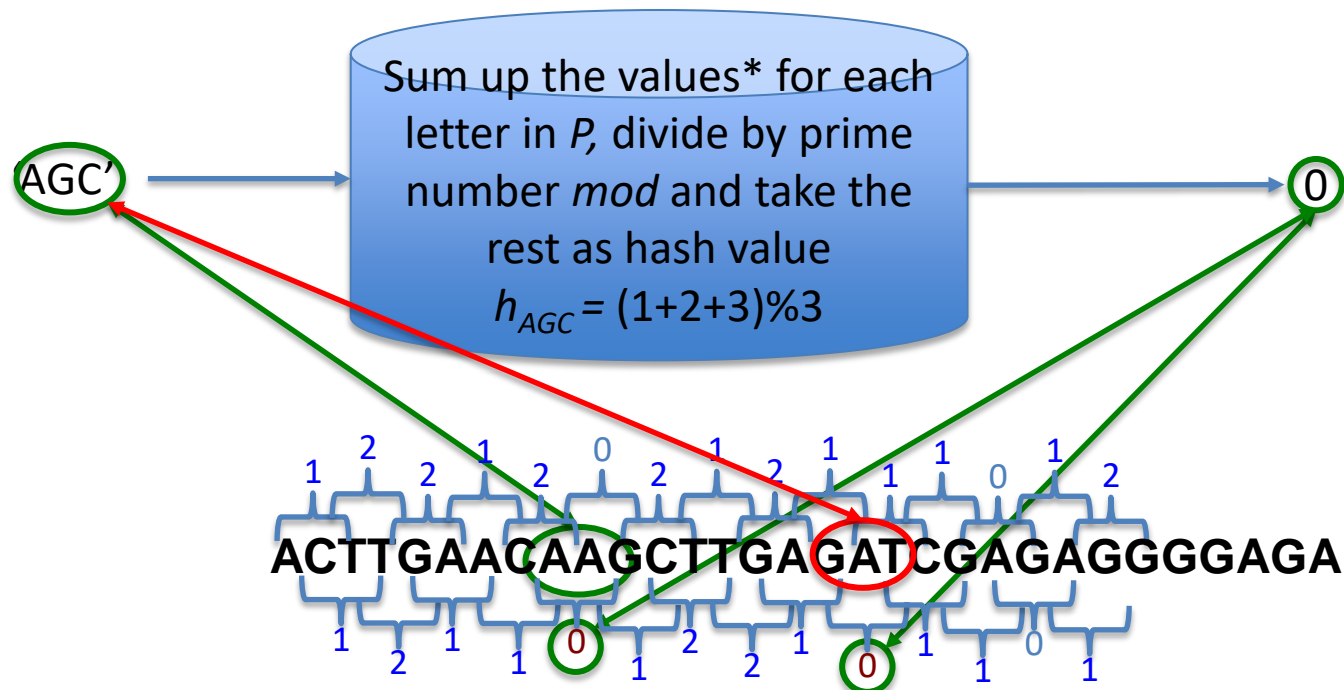
- 2) Search the text  $T$  starting from left for words of length  $|P|$  having the same hash value as  $P$ .

\*JUST AN EXAMPLE:  $v(A)=1$ ,  $v(C)=2$ ,  $v(G)=3$ ,  $v(T)=4$ ;  $mod=3$

# The use of 'hashing' in exact pattern search (Rabin-Karp; $O(n+m)$ )

## Approach:

- 1) Use a 'hash-function' to transform pattern  $P$  into a numerical hash value  $h_P$ .



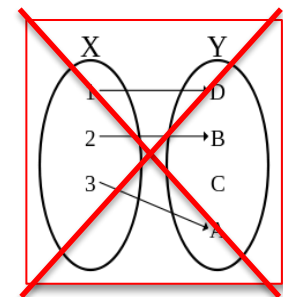
- 2) Search the text  $T$  starting from left for words of length  $|P|$  having the same hash value as  $P$ .

\*JUST AN EXAMPLE:  $v(A)=1$ ,  $v(C)=2$ ,  $v(G)=3$ ,  $v(T)=4$ ;  $mod=3$

# The use of 'hashing' in exact pattern search (Rabin-Karp; $O(n+m)$ )

## Approach:

- 1) Use a 'hash-function' to transform pattern  $P$  into a numerical hash value  $h_P$ .
- 2) Search the text  $T$  starting from left for words of length  $|P|$  having the same hash value as  $P$ .
- 3) Given a word  $K$  with  $h_K = h_P$  was found, perform an exact string comparison to verify that  $K == P$ . (Note, the projection of words with length  $|P|$  in the space of hash values is **not** injective (linkseindeutig!).)

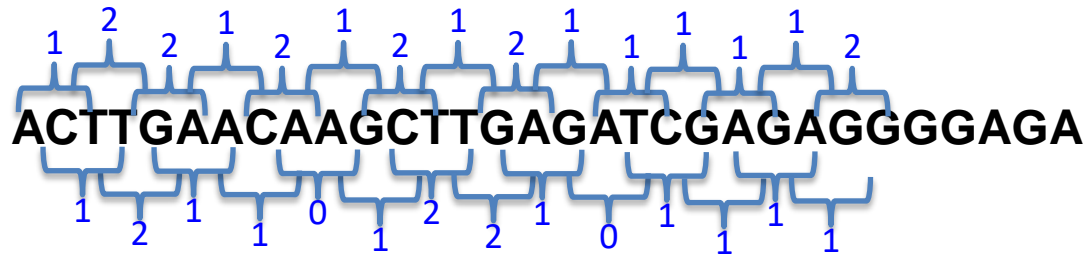


\*JUST AN EXAMPLE:  $v(A)=1$ ,  $v(C)=2$ ,  $v(G)=3$ ,  $v(T)=4$ ;  $mod=3$

# 'hashing' in combination with hash tables help to reduce the average time complexity of the pattern search to $O(1)$ , i.e. constant in time\*

**Idea: Speed up pattern search by creating look-up tables storing the hash values**

- 1) Search the text  $T$  starting from left for words of length  $k$  and compute their hash value  $h_k$



- 2) Store the hash values together with the starting point of the corresponding words in a hash table. Note, a hash table is nothing but a special way of indexing your data, just like a phone book. This will provide direct access to your potential matches in the pattern search once you know the hash value of  $P$ .

Hash value	Position in string
0	10,11,20,25,26
1	1,2,6,7,8, 12,15,18,19,21,22,23,24,27,28
2	3,4,5,9,13,14,16,17,29

\*note that this ignores the time and space you need to populate your hash table!



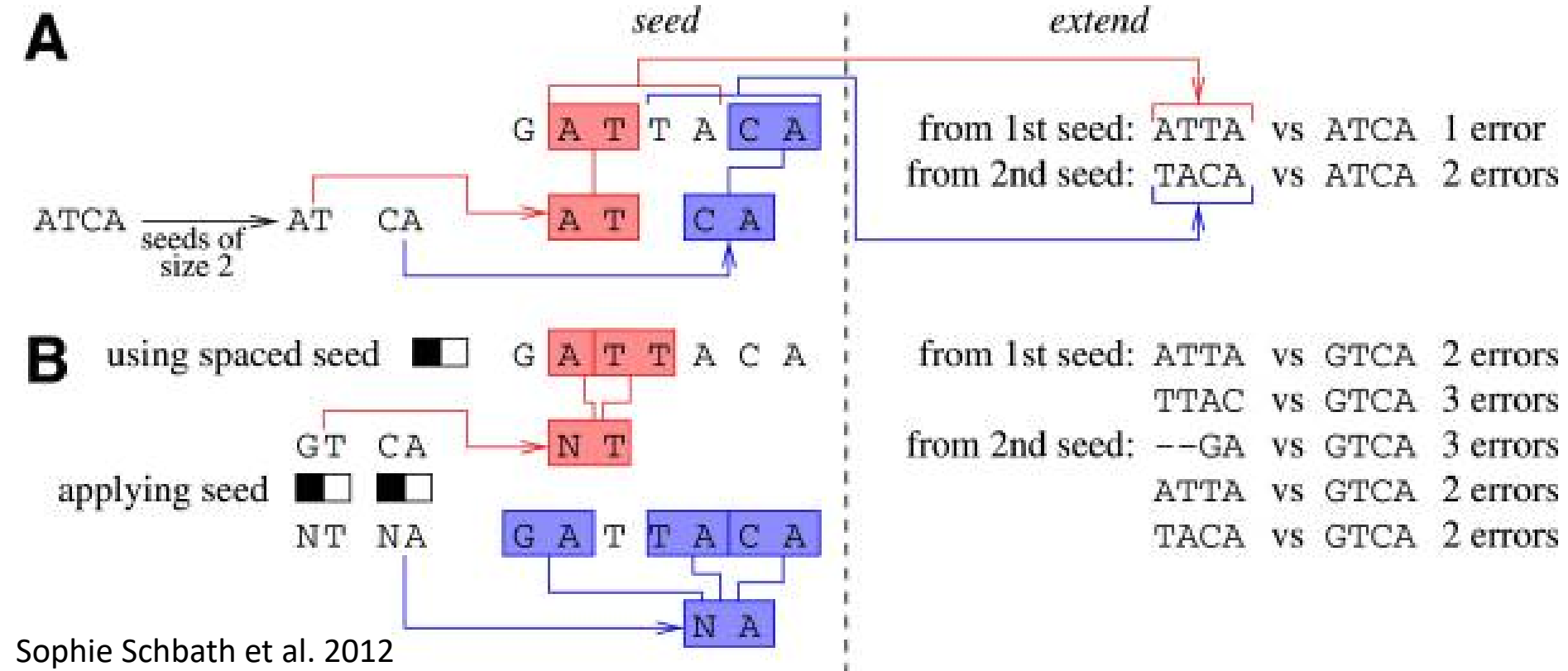
# Seed and extend with **local alignment**

---

- SSAHA & Stampy:
  - Use k-mer (shorter than read) to find it in the genome
  - Seed regions will be extended by Needleman-Wunsch
- Drawback:
  - Many regions have to be analyzed in the extend phase

# Seed and Extend the **pigeon hole principle**

- MAQ, SOAP, RMAP:
  - Chop read in k-mers (allowing errors)
  - All k-mers in the genome + correct order + adjacent to each other → read found



# Seed and Extend the **q**-gram filtering

---

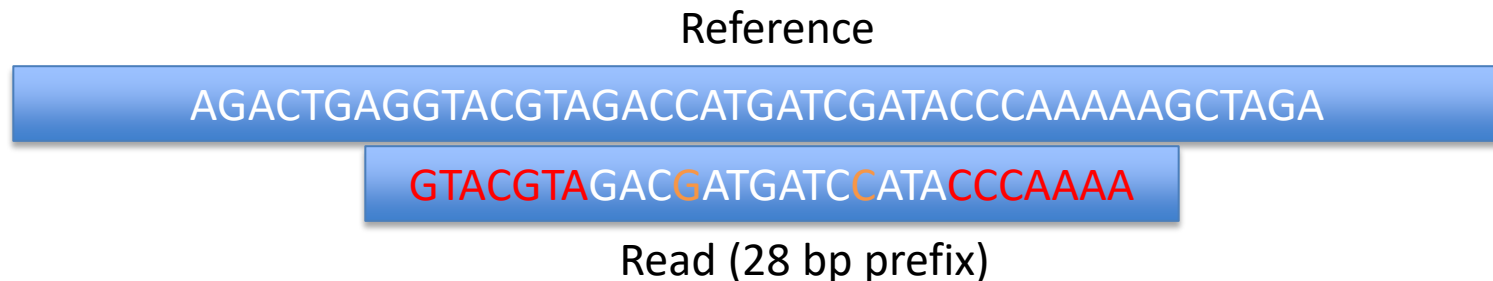
- SHRiMP2 & RazerS:
  - Chop read in k-mers but overlapping
  - If enough k-mers map in a small region a more careful alignment will be done
- Drawback:
  - List too large to be kept in memory
  - Characters stored in 8 bits (2 bits per Nucleotide) but ambiguity code has to be overcome

# Why does MAQ\* use **pigeon hole principle** and **spaced seeds** for mapping?

---

Issues to solve by seed and extend:

- Shorter seeds map more regions on the genome
- Minimum of 10 nucleotides per k-mer
- Allowing “Don’t care” positions to be able to find seeds → **spaced seed approach**



MAQ uses seed pairs as it allows, per default no more than 2 **mismatches** between seed and reference. As each 28 mere is represented by 4 non-overlapping seeds, we **always have at least 2 seeds** that must result in a perfect match to the reference.

# MAQ\*: pigeon hole principle for mapping

Example: read1 - gatgtgacatacctgttctactgaggct



Build **six** hash tables (templates) for the reads (only first 28 bp are considered) **only** from the colored nucleotides

6 'Templates'

gatgtgacatacct	gttctactgaggct	<sup>HASH</sup> →	2057673064
gatgtgacatacct	gttctactgaggct	→	3178370917
gatgtgacatacct	gttctactgaggct	→	773088662
gatgtgacatacct	gttctactgaggct	→	1856750201
gatgtgacatacct	gttctactgaggct	→	2510061809
gatgtgacatacct	gttctactgaggct	→	119777054

hash value for template 1: 832589471

GENOME

ttcagttatccaga<sup>aaatgctattttcc</sup>aaaaatgaaatctaaaatagtaactcaagtgaacattgtcagg  
tgtgtaaggaaggaaaatgaagggctaccacaaacccacaaacaacggaaactccaattcctaagtt  
tccgggacacactattcatatagatataattttacagacaaaaacgggtacttacggcaattcacaattttc  
aaaacttgcgaagcaaaaaataaaaaaattcaaaatcaatagaagacataagaaaacctttacatgac  
atcttatttttttgagtagcaaaaatacgttgtaatggatagaaaaatccttaattccgcaacaacagcctt  
atgatgaaagaccagctgggcatacaaaatttcaaaagcatcccttataaaaagtctgttaaacggacaaata  
gaacgggttcattctatcctcgtgaaatataaagatgttataaaaactaaacagggtacaccgaacattgaaga  
acagttcaattcagctgtctaggaatataactacacaattcaccctgtatcaaaaatacaaaaacaaataac  
gcccttaaaaatatttttggcagaatataaccactgatccaggaaaaatagacgaaagcagataaggcaa  
catcgaaaacctaataatcaaaaacaggcaacaggcttaaaaaccataacgaaaaaagcaataagatcaa  
agattatgaccaggacaaaacagttttataaagcaaatcacaaggccaggttctaagctgtacaagaag  
aaacattaaaaagaaaacagacaaacattgttatcacagaagcaggaagatgtgacatacctgttctactg  
aaggct

compute hash values for spaced seeds in reference (on both strands and for one of the six templates) and perform lookups in hash tables of the reads

# MAQ\*: pigeon hole principle for mapping

Example: read1 - gatgtgacatacctgttctactgaggct



Build **six** hash tables for the reads (only first 28 bp) using **only** the colored nucleotides

gatgtgacatacctgttctactgaggct	<sup>HASH</sup> →	2057673064
gatgtgacatacctgttctactgaggct	→	3178370917
gatgtgacatacctgttctactgaggct	→	773088662
gatgtgacatacctgttctactgaggct	→	1856750201
gatgtgacatacctgttctactgaggct	→	2510061809
gatgtgacatacctgttctactgaggct	→	119777054



continue with next word in the reference from the same template.... until the entire reference sequence has been used.

1258119214

GENOME

ttcaqttatccagaaatgctattttcccaaaatgaatctaaaatagtaactcaagtgaacattgtcagg  
tgtgtaaggaaggaaaatgaagggctaccacaaacccacaaacaacggaactccaattcctaagtt  
tccgggacacactattcatatagatataattttacagacaaaaacgggtacttacggcaattcacaattttc  
aaaactgcgaagcaaaaaataaaaaaattcaaaatcaatagaagacataagaaaacctttacatgac  
atcttattttattttggagtacaaaaatcgttgtaatggatatagaaaaatccttaattccgcaacaacagcctt  
atgatgaaagaccagctgggcatacaaaatttcaaaagcatcccttataaaaagtctgttaaacggacaaata  
gaacggtttcattctatcctcgtgaaatataaagatgttataaaactaaacagggtacaccgaacattgaaga  
acagttcaattcagctgtctaggaatataactacacaattcacccctgtatcaaaaatacaaaaacaaataac  
gcccttaaaaatattttggcagaatataaccactgatccaggaaaaatgatgacgaaagcagataaggcaa  
catcgaaaacctaataatcaaaaacaggcaacaggcttaaaaaccataacgaaaaaagcaataagatcaa  
agattatgacccaggacaaaacagttttataaagcaaatcacaaggccaggttctaagctgtacaagaag  
aaacattaaaaagaaaacagacaaacattgttatcacagaagcaggaagatgtgacatacctgttctactg  
aaggct

\*Li et al *Genome Res.* 2008. 18: 1851-1858



# MAQ \*: pigeon hole principle for mapping

Example: read1 - gatgtgacatacctgttctactgaggct



Build **six** hash tables for the reads (only first 28 bp) using **only** the colored nucleotides

gatgtgacatacctgttctactgaggct	<small>HASH</small> →	2057673064
gatgtgacatacctgttctactgaggct	→	3178370917
gatgtgacatacctgttctactgaggct	→	773088662
gatgtgacatacctgttctactgaggct	→	1856750201
gatgtgacatacctgttctactgaggct	→	2510061809
gatgtgacatacctgttctactgaggct	→	119777054

HIT: Calculate the sum of qualities of mismatched bases  $q$  over the whole length of the read and store together with hit position.

For each read, MAQ stores score and position of only the 2 best hits and the number of 0-, 1-, and 2-mismatch seed positions

## GENOME

```
ttcagttatccagaaaatgctattttccaaaaatgaaatctaaaatagtaactcaagtgaaacattgtcagggtg
gtaaggaaggaaaatgaagggtaccacaaaaccacaaaacaacggaaactccaattcctaagttcc
gggacacactattcatatagatataattttacagacaaaaacgggtactacggcaattcacaattttcaa
acttgcgaaagcaaaaaataaaaaaattcaaatcaatagaagacataagaaaacctttacatgacatctt
attttattttgagtagcaaaaatcgttgtaattgtagataaaaaatccttaattccgcaacaacagccttatg
atgaaagaccagctgggcatacaaatttcaaaagcatcccctataaaaagtctgtaaacggacaaaatagaa
cggtttcattctatcctcgctgaaattaaaagatgtttaaaactaaacaggtacaccgaacatttgagaaca
gttcaattcagctgtctaggaatataactacacaattcacctgtatcaaaaatacaaaaacaaaataacgccc
ttaaaaatattttggcagaatataaccactgatccaggaaaatagacgaaagcagataaggcaacatc
gaaaacctataatcaaaaacaggcaacaggcttaaaaaccataacgaaaaaagcaataagatcaagat
tatgagccaggacaaaacagttttataaagcaaatcacaggccagggttctaagctgtacaagaaagaaac
attaaaagaaaacagacaaaacattgttatcacagaagcaggaagatgtgacatacctgttctactgaa
ggct
```

2057673064

# MAQ\* uses spaced seeds for mapping

Example: read1 - gatgtgacatacctgttctactgaggct

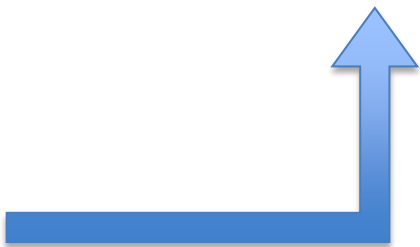


Build **six** hash tables for the reads (only first 28 bp) using **only** the colored nucleotides

gatgtgacatacctgttctactgaggct	<sup>HASH</sup> →	2057673064
gatgtgacatacctgttctactgaggct	→	3178370917
gatgtgacatacctgttctactgaggct	→	773088662
gatgtgacatacctgttctactgaggct	→	1856750201
gatgtgacatacctgttctactgaggct	→	2510061809
gatgtgacatacctgttctactgaggct	→	119777054

hash value for template 2

GENOME  
ttcagttatccaga<sup>aaatgctattttcc</sup>caaaatgaaatctaaaatagtaactcaagtgaacattgtcagg  
tgtgtaaggaaggaaaatgaagggctaccacaaacccacaaacaacggaaactccaattcctaagtt  
tccgggacacactattcatatagatataattttacagacaaaaacgggtacttacggcaattcacaattttc  
aaaactgcgaaagcaaaaaataaaaaaattcaaaatcaatagaagacataagaaaacctttacatgac  
atcttattttattttggagtacaaaaatcgttgtaatggatatagaaaaatccttaattccgcaacaacagcctt  
atgatgaaagaccagctgggcatacaaaatttcaaaagcatccctctataaaaagtctgttaaacggacaaata  
gaacggttcattctatcctcgtctgaaatataaagatgttataaaactaaacagggtacaccgaacatttgaaga  
acagttcaattcagctgtctaggaatataactacacaattcacccctgtatcaaaaatacaaaaacaaataac  
gcccttaaaaatatttttggcagaatataaccactgatccaggaaaaatgatgacgaaagcagataaggcaa  
catcgaaaacctaataatcaaaaacaggcaacaggcttaaaaaccataacgaaaaaagcaataagatcaa  
agattatgagccaggacaaaacagttttataaagcaaatcacaaggccaggttctaagctgtacaagaaag  
aaacattaaaagaaaacagacaaacattgttatcacagaagcagggaagatgtgacatacctgttctactg  
aaggct



compute hash values for spaced seeds in reference (on both strands) and perform lookup in hash tables of the reads



# MAQ\* uses spaced seeds for mapping

Example: read1 - gatgtgacatacctgttctactgaggct



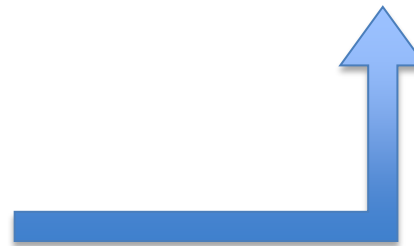
Build **six** hash tables for the reads (only first 28 bp) using **only** the colored nucleotides

gatgtgacatacctgttctactgaggct	<sup>HASH</sup> →	2057673064
gatgtgacatacctgttctactgaggct	→	3178370917
gatgtgacatacctgttctactgaggct	→	773088662
gatgtgacatacctgttctactgaggct	→	1856750201
gatgtgacatacctgttctactgaggct	→	2510061809
gatgtgacatacctgttctactgaggct	→	119777054

hash value 3

GENOME

ttcagttatccagaaaatgctattttccaaaatgaaatctaaaatagtaactcaagtgaacattgtcagg  
tgtgtaagggaaggaaaatgaagggctaccacaaacccacaaacaacggaaactccaattcctaagtt  
tccgggacacactattcatatagatataatttctacagacaaaaacgggtacttacggcaattcacaattttc  
aaaactgcgaagcaaaaaataaaaaaattcaaaatcaatagaagacataagaaaacctttacatgac  
atcttattttatttggagtacaaaaatcgttgtaatggatatagaaaaatccttaattccgcaacaacagcctt  
atgatgaaagaccagctgggcatacaaaatttcaaaagcatccctctataaaaagtctgttaaacggacaaata  
gaacggttcattctatcctcgtctgaaatataaagatgttataaaactaaacagggtacaccgaacatttgaaga  
acagttcaattcagctgtctaggaatataactacacaattcacctgtatcaaaaatacaaaaacaaataac  
gcccttaaaaatatttttggcagaaatataaccactgatccaggaaaaatgatgacgaaagcagataaggcaa  
catcgaaaacctaataatcaaaaacaggcaacaggcttaaaaaccataacgaaaaaagcaataagatcaa  
agattatgagccaggacaaaacagttttataaagcaaatcacaaggccaggttctaagctgtacaagaag  
aaacattaaaagaaaacagacaaacattgttatcacagaagcagggaagatgtgacatacctgttctactg  
aaggct



compute hash values for spaced seeds in reference (on both strands) and perform lookup in hash tables of the reads

# MAQ\* uses spaced seeds for mapping

Example: read1 - gatgtgacatacctgttctactgaggct



Build **six** hash tables for the reads (only first 28 bp) using **only** the colored nucleotides

gatgtgacatacctgttctactgaggct	<sup>HASH</sup> →	2057673064
gatgtgacatacctgttctactgaggct	→	3178370917
gatgtgacatacctgttctactgaggct	→	773088662
gatgtgacatacctgttctactgaggct	→	1856750201
gatgtgacatacctgttctactgaggct	→	2510061809
gatgtgacatacctgttctactgaggct	→	119777054

hash value 4

GENOME  
ttcagttatccaga<sup>aaatgct</sup>attttccaaaatgaaatctaaaatagtaactcaagtgaacattgtcagg  
tgtgtaagggaaggaaaatgaagggctaccacaaacccacaaacaacggaaactccaattcctaagtt  
tccgggacacactattcatatagatataatttctacagacaaaaacgggtacttacggcaattcacaattttc  
aaaactgcgaaagcaaaaaataaaaaaattcaaaatcaatagaagacataagaaaacctttacatgac  
atcttattttatttggagtacaaaaatcgttgtaatggatatagaaaaatccttaattccgcaacaacagcctt  
atgatgaaagaccagctgggcatacaaaatttcaaaagcatccctctataaaaagtctgttaaacggacaaata  
gaacggttcattctatcctcgtctgaaatataaagatgttataaaactaaacagggtacaccgaacatttgaaga  
acagttcaattcagctgtctaggaatataactacacaattcacccctgtatcaaaaatacaaaaacaaataac  
gcccttaaaaatatttttggcagaaatataaccactgatccaggaaaaatgatgacgaaagcagataaggcaa  
catcgaaaacctaataatcaaaaacaggcaacaggcttaaaaaccataacgaaaaagcaataagatcaa  
agattatgagccaggacaaaacagttttataaagcaaatcacaaggccaggttctaagctgtacaagaaag  
aaacattaaaagaaaacagacaaacattgttatcacagaagcagggaagatgtgacatacctgttctactg  
aaggct



compute hash values for spaced seeds in reference (on both strands) and perform lookup in hash tables of the reads

# MAQ\* uses spaced seeds for mapping

Example: read1 - gatgtgacatacctgttctactgaggct



Build **six** hash tables for the reads (only first 28 bp) using **only** the colored nucleotides

gatgtgacatacctgttctactgaggct	<sup>HASH</sup> →	2057673064
gatgtgacatacctgttctactgaggct	→	3178370917
gatgtgacatacctgttctactgaggct	→	773088662
gatgtgacatacctgttctactgaggct	→	1856750201
gatgtgacatacctgttctactgaggct	→	2510061809
gatgtgacatacctgttctactgaggct	→	119777054

hash value 5

GENOME  
ttcagttatccagaaaatgctattttcccaaaatgaaatctaaaatagtaactcaagtgaacattgtcagg  
tgtgtaagggaaggaaaatgaagggctaccacaaacccacaaacaacggaaactccaattcctaagtt  
tccgggacacactattcatatagatataatttctacagacaaaaacgggtacttacggcaattcacaattttc  
aaaactgcgaagcaaaaaataaaaaaattcaaaatcaatagaagacataagaaaacctttacatgac  
atcttattttatttggagtacaaaaatcgttgtaatggatatagaaaaatccttaattccgcaacaacagcctt  
atgatgaaagaccagctgggcatacaaaatttcaaaagcatcccttataaaaagtgtgtaaacggacaaata  
gaacggttcattctatcctcgctgaaatataaagatgttataaaactaaacagggtacaccgaacatttgaaga  
acagttcaattcagctgtctaggaatataactacacaattcacctgtatcaaaaatacaaaaacaaataac  
gcccttaaaaatatttttggcagaaatataaccactgatccaggaaaaatgacgaaagcagataaggcaa  
catcgaaaacctaataatcaaaaacaggcaacaggcttaaaaaccataacgaaaaaagcaataagatcaa  
agattatgagccaggacaaaacagttttataaagcaaatcacaaggccaggttctaagctgtacaagaag  
aaacattaaaaagaaaacagacaaacattgttatcacagaagcagggaagatgtgacatacctgttctactg  
aaggct



compute hash values for spaced seeds in reference (on both strands) and perform lookup in hash tables of the reads

# MAQ\* uses spaced seeds for mapping

Example: read1 - gatgtgacatacctgttctactgaggct



Build **six** hash tables for the reads (only first 28 bp) using **only** the colored nucleotides

gatgtgacatacctgttctactgaggct	<sup>HASH</sup> →	2057673064
gatgtgacatacctgttctactgaggct	→	3178370917
gatgtgacatacctgttctactgaggct	→	773088662
gatgtgacatacctgttctactgaggct	→	1856750201
gatgtgacatacctgttctactgaggct	→	2510061809
gatgtgacatacctgttctactgaggct	→	119777054

hash value 6

GENOME  
ttcagttatccagaaaatgctattttccaaaatgaaatctaaaatagtaactcaagtgaacattgtcagg  
tgtgtaagggaaggaaaatgaagggctaccacaaacccacaaacaacggaactccaattcctaagtt  
tccgggacacactattcatatagatataattttacagacaaaaacgggtacttacggcaattcacaaatttc  
aaaactgcgaaagcaaaaaataaaaaaattcaaaatcaatagaagacataagaaaacctttacatgac  
atcttattttatttggagtacaaaaatcgttgtaatggatatagaaaaatccttaattccgcaacaacagcctt  
atgatgaaagaccagctgggcatacaaaatttcaaaagcatcccttataaaaagtctgtaaacggacaaata  
gaacggttcattctatcctcgtgtaaatataaagatgttataaaactaaacagggtacaccgaacatttgaaga  
acagttcaattcagctgtctaggaatataactacacaattcacctgtatcaaaaatacaaaaacaaataac  
gcccttaaaaatatttttgcagaaatataaccactgatccaggaaaatgatgacgaaagcagataaggcaa  
catcgaaaacctaataatcaaaaacaggcaacaggcttaaaaaccataacgaaaaagcaataagatcaa  
agattatgagccaggacaaaacagttttataaagcaaatcacaaggccaggttctaagctgtacaagaag  
aaacattaaaagaaaacagacaaacattgttatcacagaagcagggaagatgtgacatacctgttctactg  
aaggct



compute hash values for spaced seeds in reference (on both strands) and perform lookup in hash tables of the reads

# MAQ: An overview

---

1. At the alignment stage, MAQ first searches for the **ungapped match** with lowest mismatch score, defined as the sum of qualities at mismatching bases.
2. MAQ only considers positions that have **two or fewer mismatches in the first 28 bp** (default parameters; speed-up).
3. Sequences that fail to reach a mismatch score threshold but whose mate pair is mapped are **searched with a gapped alignment algorithm** in the regions defined by the mate pair.
4. To evaluate the reliability of alignments, MAQ assigns each individual alignment a **phred-scaled quality score** (capped at 99), which measures the probability that the true alignment is not the one found by MAQ.
5. MAQ **always reports a single alignment**, and if a read can be aligned equally well to multiple positions, MAQ will **randomly pick** one position and give it a mapping quality zero. Because their mapping score is set to zero, reads that are mapped equally well to multiple positions will not contribute to variant calling.

# NextGenMap

---

- Hash-based read mapper bridging:
  1. speed
  2. ability to map reads in highly polymorphic regions
- pitfalls:
  1. Mixing single end and paired end reads is not supported → corrupt the mapping results
  2. Not recommended to change parameters –kmer and kmer-skip
  3. Values between 0.5 and 0.8 give the best trade-off between speed and sensitivity

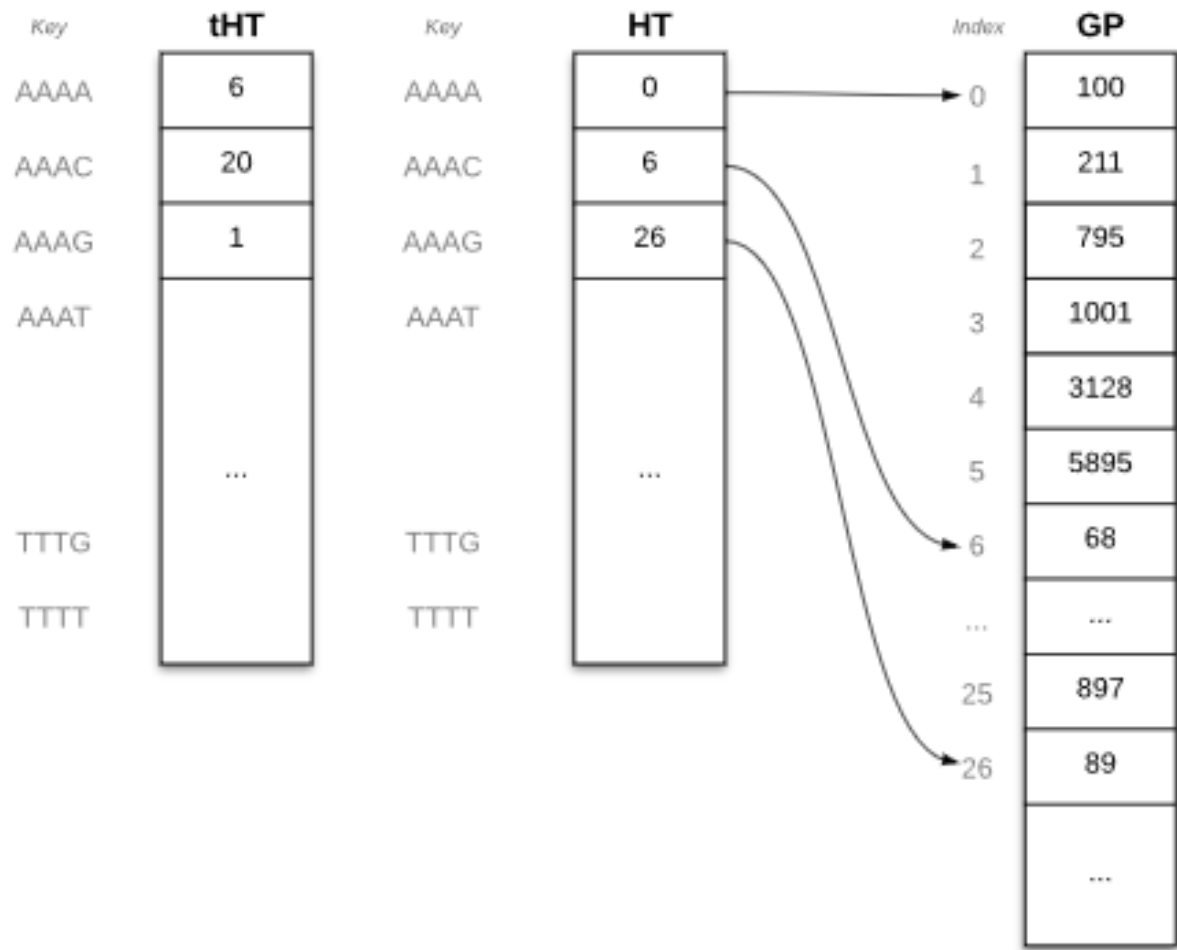
# Algorithm Workflow

---

1. Indexing the reference genome (only once)
2. Identification of possible mapping regions on the reference genome (candidate mapping region (CMR) search)
3. Computation of alignment scores for all CMRs found in step 2
4. Computation of the full alignment for the best scoring CMR from step 3

# Indexing the reference genome

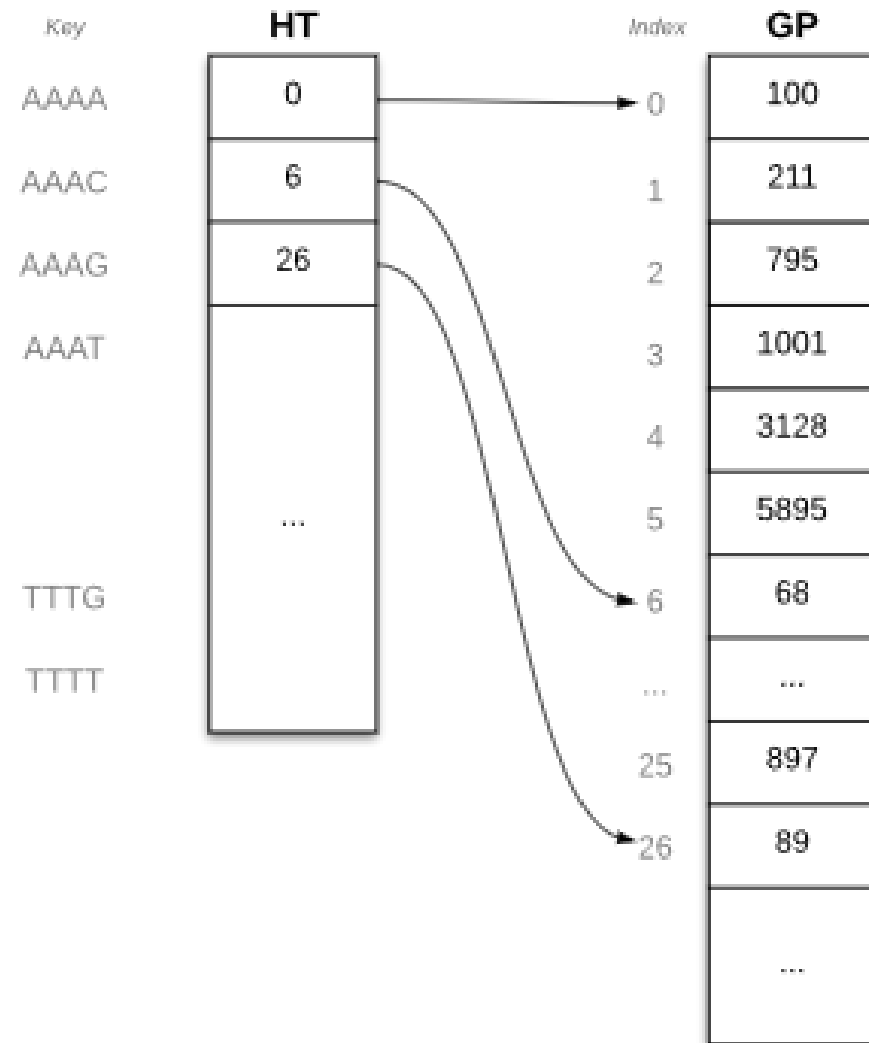
- k-mer = 13 bp
- Every 3. position
- Only A, C, G, T
- Only + strand



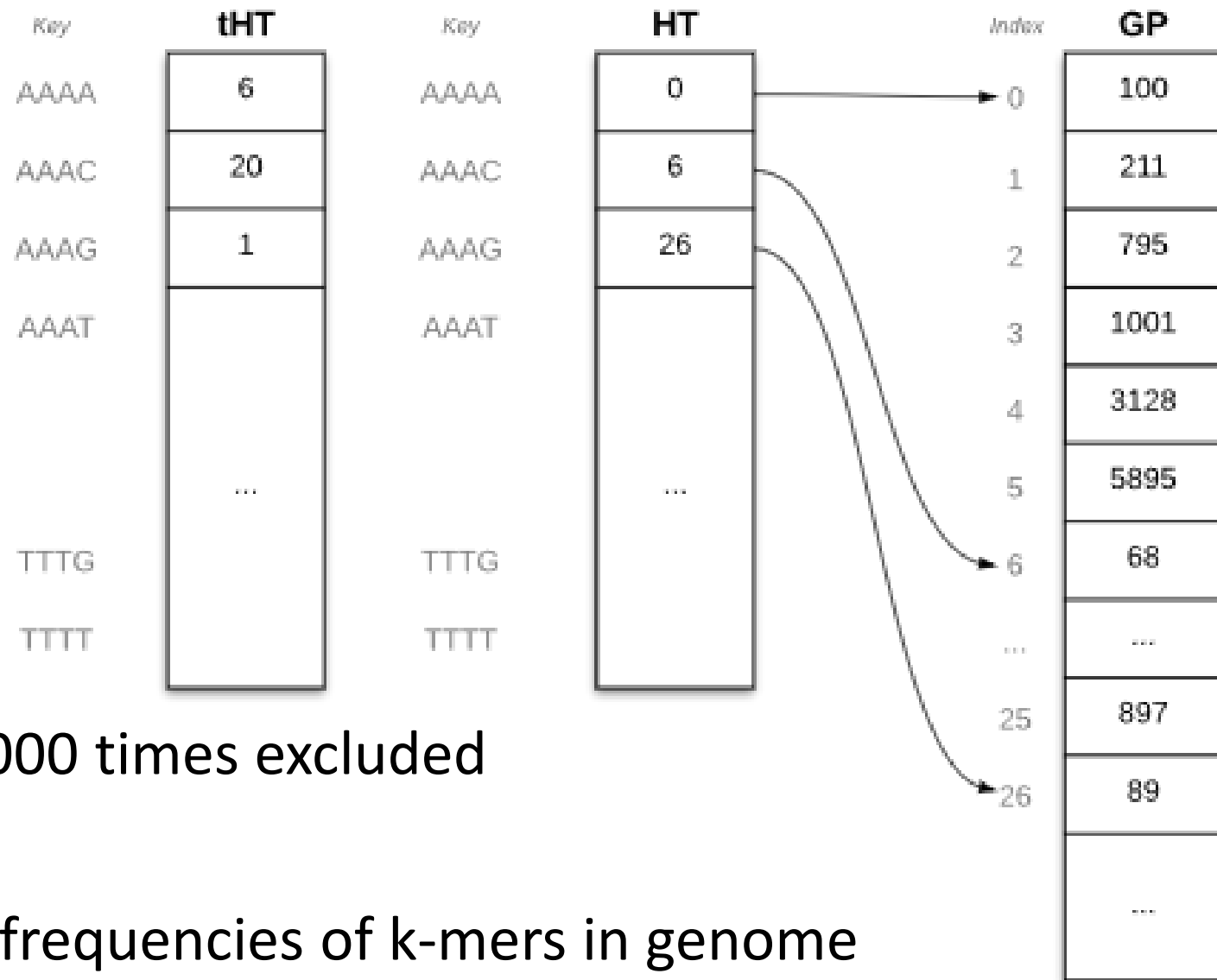


# Indexing the reference genome

- GP = genomic positions
  - Consecutively k-mer blocks
  - Saving start position
- HT = Hash table of k-mers
  - Lexicographic ordered
  - Frequency (numbered serially)

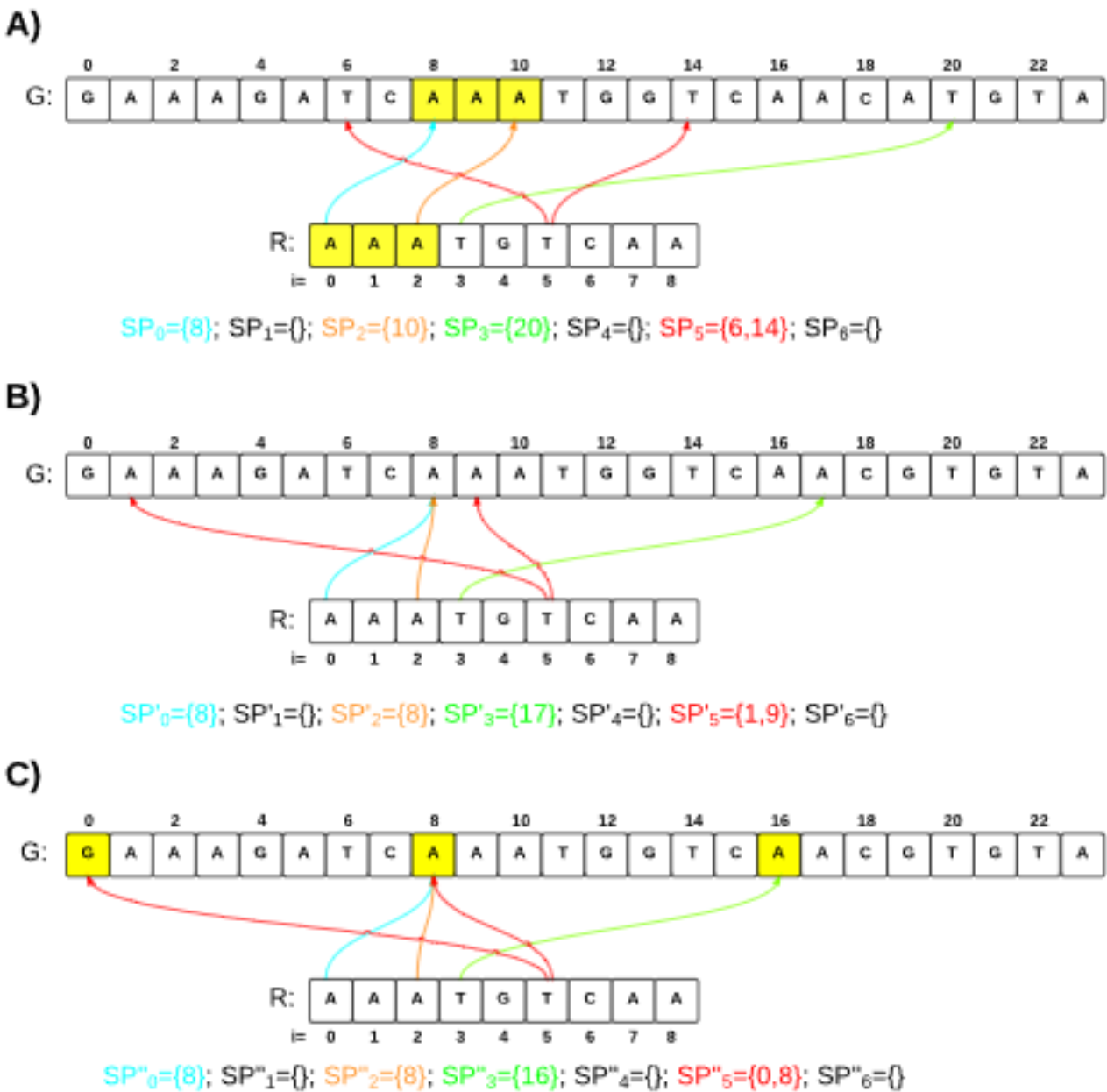


# Indexing the reference genome



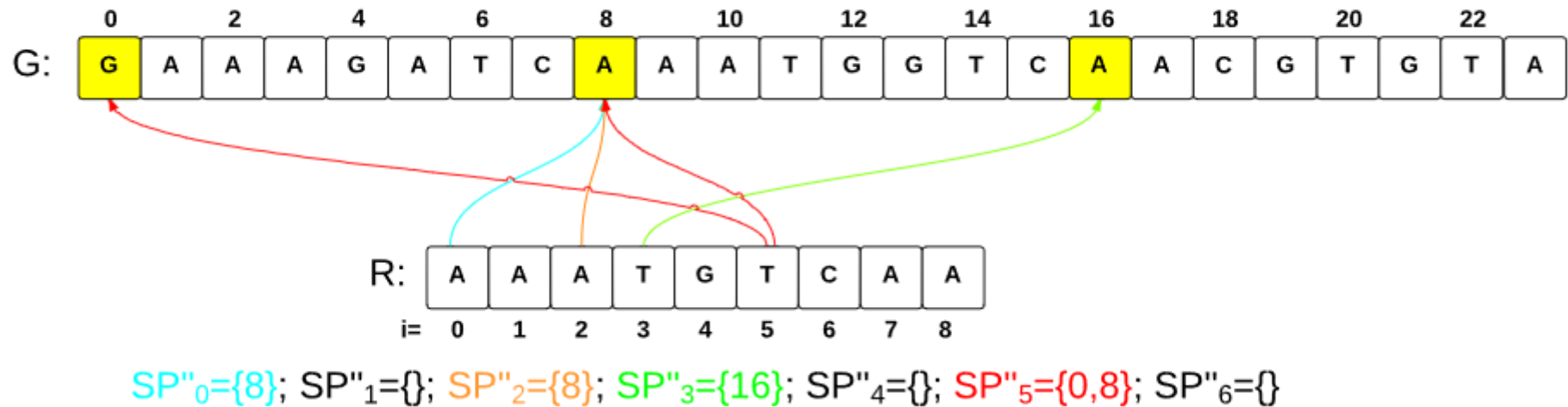
# Identification of CMR

- R = read of length 9bp
- G = reference Genome
- Step-size = 2
- k-mer size = 3bp



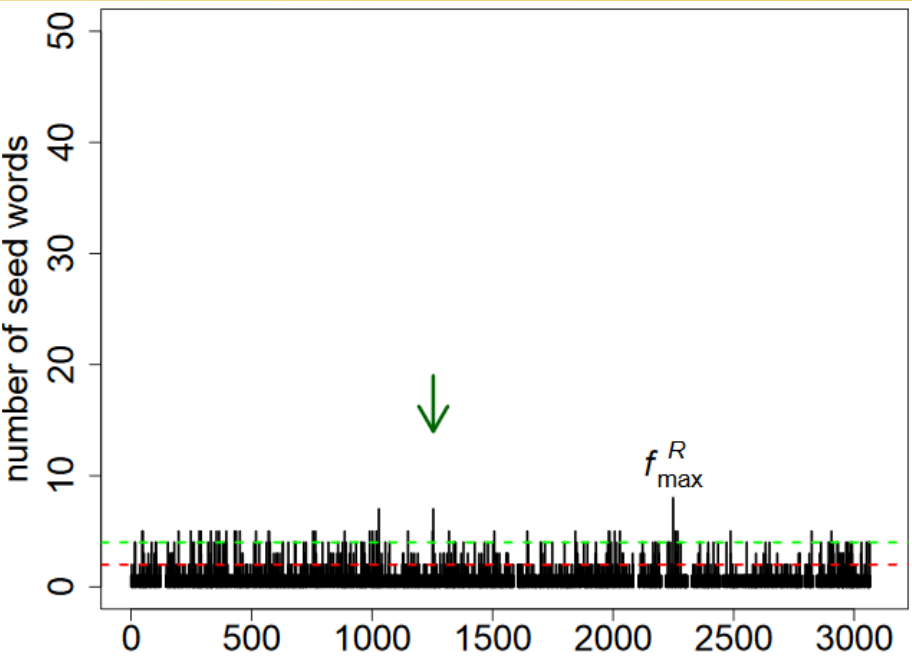
# Identification of CMR

c)



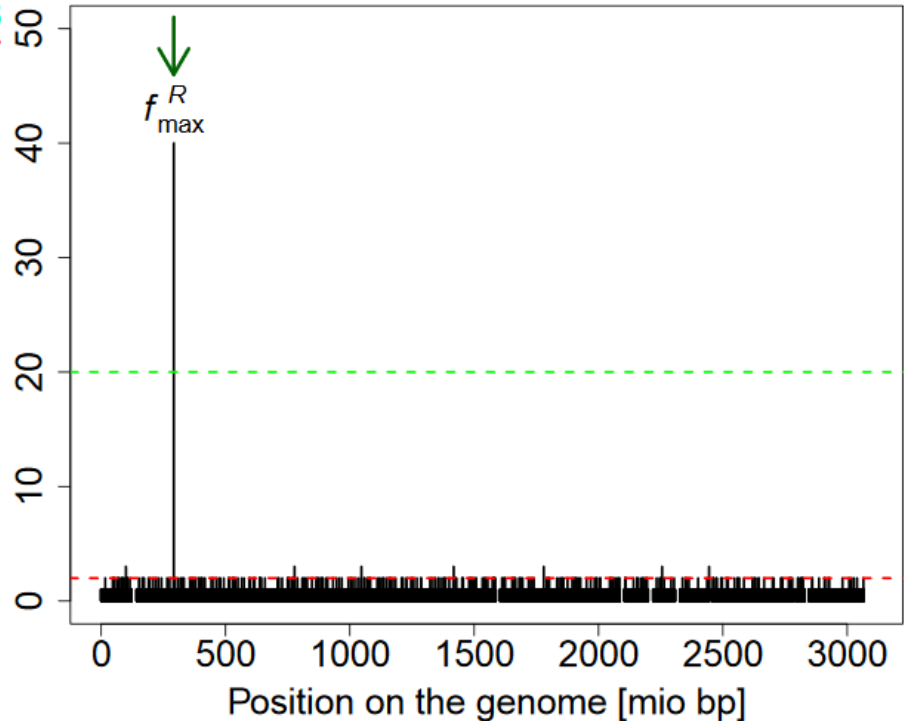
- Position 1 not found due to step-size = 2
- Saving startpoints in reference genome
- Shift of readstart to startposition
- Allow modulo 8 (bit-shift operation) due to polymorphic regions to reduce start positions

# Computation of CMRs



informative

uninformative



Position on the genome [mio bp]



# Computation of CMRs

---

$$\sigma = \frac{\overline{F_{\max}} = \frac{1}{B} \sum_{j=1}^B \max\{F_{R_j}\}}{S_{\max} = \left\lceil \frac{l - k + 1}{\Delta} \right\rceil}$$

- Calculating the ratio of
  - Average read seed count
  - Perfectly matching read
- $\sigma = 1 \rightarrow$  perfect reads
- $\sigma = 0 \rightarrow$  very different reads
- Calculating genomic start points
  - $F_R \rightarrow$  Frequency distribution
  - Genomic position with seed word count above  $\Theta_R$

$$\Theta_R = \sigma \max\{F_R\}$$

# Computation of alignment scores

---

1. Calculating the read alignment score

$$p'' - c/2 \text{ to } p'' + l + c/2$$

2.  $c$  defines consecutive insertions deletions  
dependent on the read length ( $l$ )

$$c = 5 + \bar{l} * 0.15$$

3. MASon (Rescheneder et al. 2012) for  
alignment

# **BURROWS-WHEELER TRANSFORM BASED ALGORITHMS**



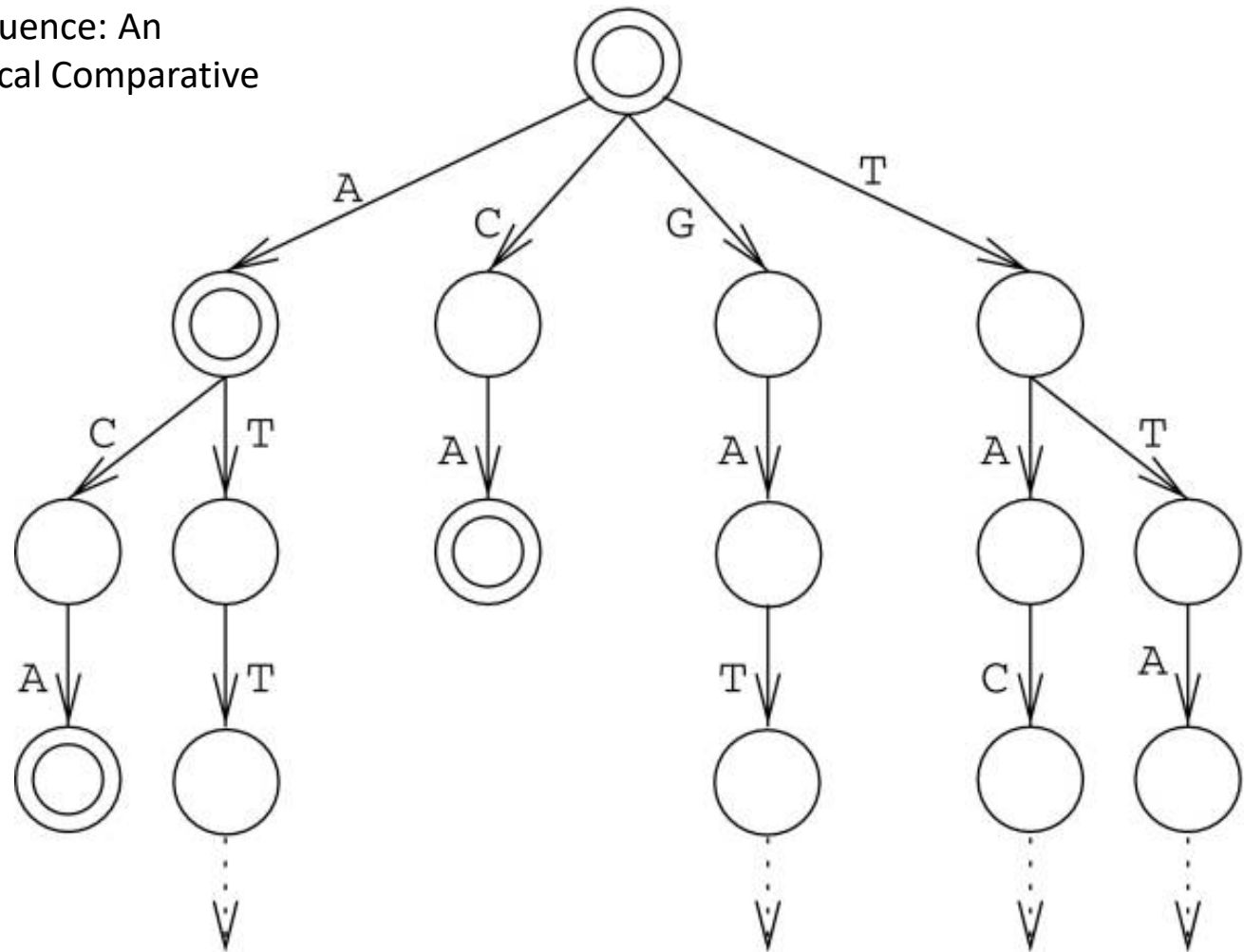
# Suffix Trees / Suffix Tries

---

- **Advantage:**
  - **Hashing performs poor for repeating regions**
- Each suffix of a word is represented as path from leaf to root
  - Size of tree proportional to size of genome
  - Building time proportional to size of genome
  - Search time  $O(L_r)$

# Suffix Trees

Mapping Reads on a Genomic Sequence: An  
Algorithmic Overview and a Practical Comparative  
Analysis (2012) Sophie Schbath



- **Advantage:**
  - Hashing performs poor for repeating regions
  - Repeating regions are squeezed into one path

# Suffix Arrays

- Problem of the trees and tries:
  - Large genomes will not fit in the RAM
- Array is the set of suffixes sorted lexicographically
  - Trick 1: save the startposition of the suffix

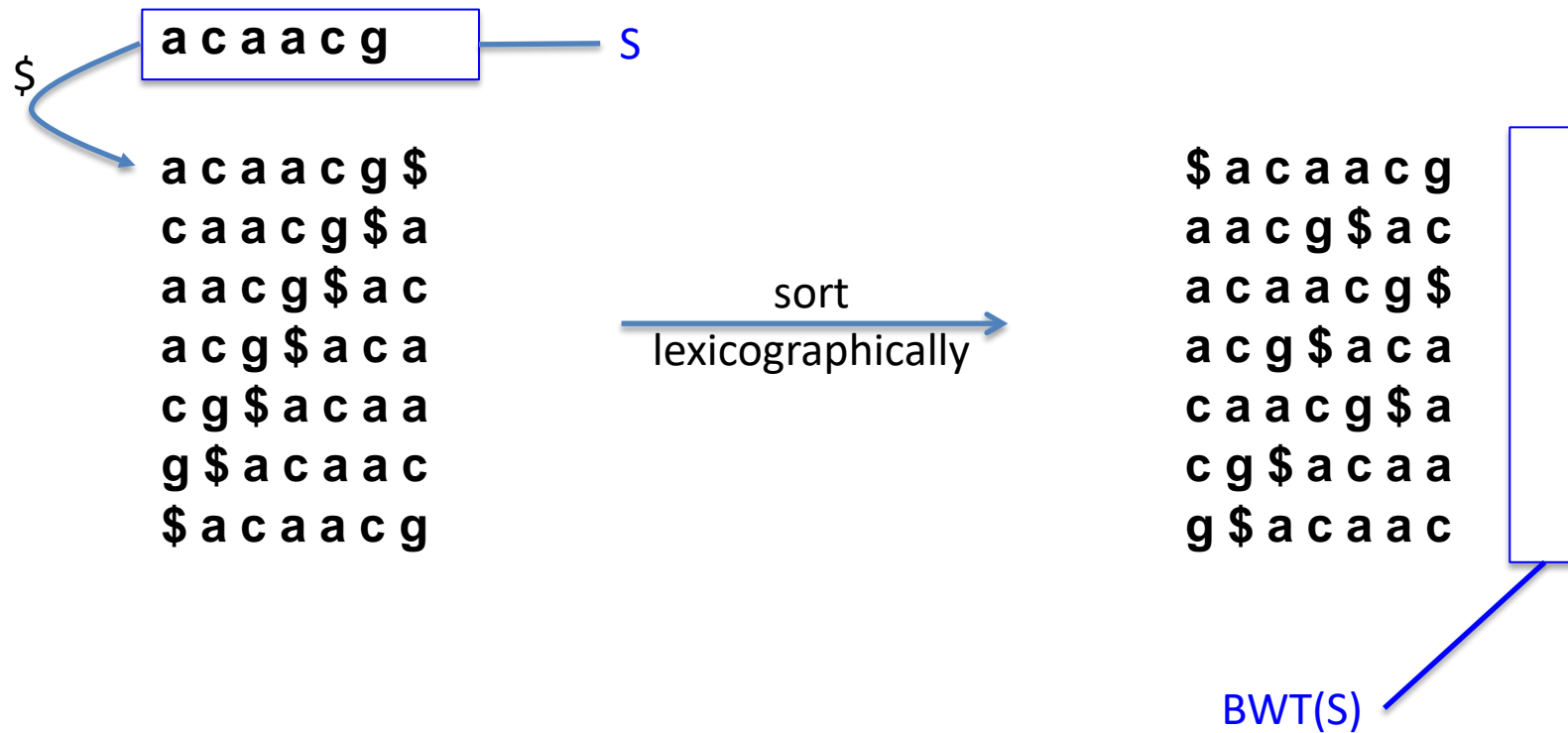
genome	suffixes	sorted suffixes	positions	cylinder suffix array	B–W
GATTACA\$	GATTACA\$	ACA\$	4	ACA\$GATT	T
	ATTACA\$	ATTACA\$	2	ATTACA\$G	G
	TTACA\$	A\$	6	A\$GATTAC	C
	TACA\$	CA\$	5	CA\$GATTA	A
	ACA\$	GATTACA\$	1	GATTACA\$	\$
	CA\$	TACA\$	4	TACA\$GAT	T
	A\$	TTACA\$	3	TTACA\$GA	A
	\$	\$	7	\$GATTACA	A

## Trick 2: The Burrows-Wheeler Transform

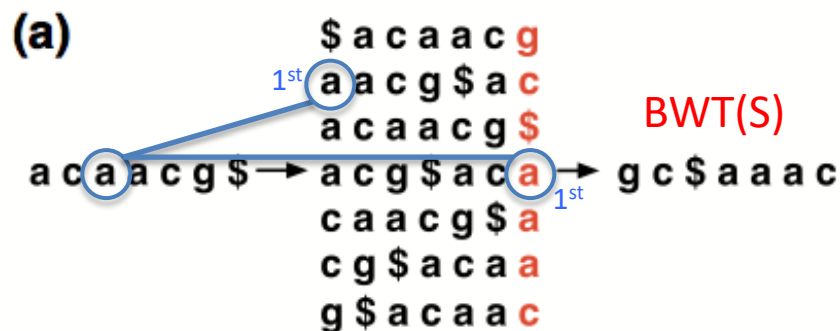
---

- Invented by David Wheeler in 1983 (bell labs), pub. 1994
- Used in data compression (bzip2)
- Used in 2003 on the human genome to define exact word matches (originally for microarray probe design)
- First used for short read alignment by bowtie, now adopted by bwa (maq author) and SOAP2

# The Burrows-Wheeler Transform



# Burrows Wheeler Transform (BWT)

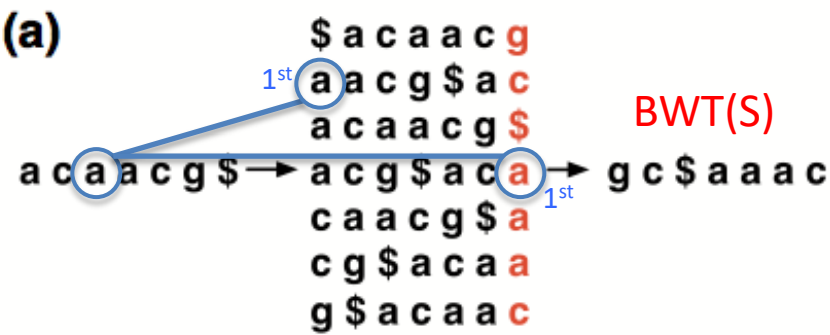


Generate matrix by

1. Appending a \$ to the end of the string  $S$  that should be indexed. \$ should have 2 properties
  1. it must not occur in the string
  2. it should be lexicographically smaller than any character in  $S$
2. generate all cyclic permutations of  $S$
3. sort the resulting matrix lexicographically (the line beginning with the \$ is the first to occur in the matrix).

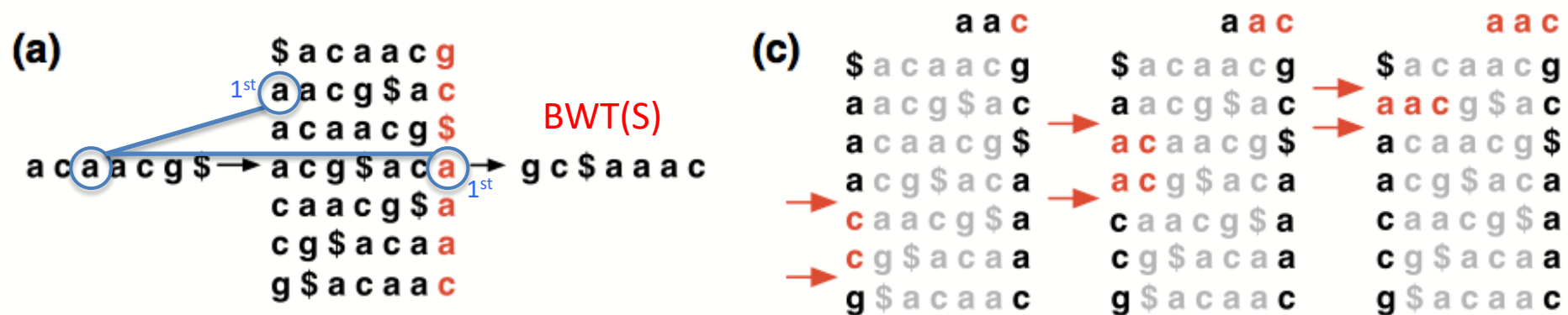
**The matrix has the property of last first (LF) mapping:** The  $i$ th occurrence of character  $X$  in the last column corresponds to the same text character as the  $i$ th occurrence of  $X$  in the first column

# Burrows Wheeler Transform (BWT)



The matrix has the property of **last first (LF) mapping**: This can be used to reconstruct the original text from BWT(S) using the UNPERMUTE algorithm.

# Burrows Wheeler Transform (BWT)



**The matrix has the property of last first (LF) mapping:** This can be used to search for a text within BWT(S) using the EXACTMATCH algorithm.

Key aspects:

- 1) Matrix is sorted lexicographically. Thus, rows beginning with a given sequence appear **consecutively**.
- 2) EXACTMATCH algorithm calculates the range of matrix rows beginning with **successively longer suffixes** of the query.
- 3) At each step, the size of the **range either shrinks or remains the same**.
- 4) When the algorithm completes, rows beginning with  $S_0$  (the entire query) correspond to **exact occurrences of the query** in the text.





Original string: ctgaaactggt

Put a \$ on the end

create cyclic rotations of the string...

---

c	t	g	a	a	a	c	t	g	g	t	\$
t	g	a	a	a	c	t	g	g	t	\$	c
g	a	a	a	c	t	g	g	t	\$	c	t
a	a	a	c	t	g	g	t	\$	c	t	g
a	a	c	t	g	g	t	\$	c	t	g	a
a	c	t	g	g	t	\$	c	t	g	a	a
c	t	g	g	t	\$	c	t	g	a	a	a
t	g	g	t	\$	c	t	g	a	a	a	c
g	g	t	\$	c	t	g	a	a	a	c	t
g	t	\$	c	t	g	a	a	a	c	t	g
t	\$	c	t	g	a	a	a	c	t	g	g
\$	c	t	g	a	a	a	c	t	g	g	t

Alphabetically sort the permuted strings, first column is the “genome dictionary” last column is the Burrows-wheeler transformation

---

\$	c	t	g	a	a	a	c	t	g	g	t	11
a	a	a	c	t	g	g	t	\$	c	t	g	3
a	a	c	t	g	g	t	\$	c	t	g	a	4
a	c	t	g	g	t	\$	c	t	g	a	a	5
c	t	g	a	a	a	c	t	g	g	t	\$	0
c	t	g	g	t	\$	c	t	g	a	a	a	6
g	a	a	a	c	t	g	g	t	\$	c	t	2
g	g	t	\$	c	t	g	a	a	a	c	t	8
g	t	\$	c	t	g	a	a	a	c	t	g	9
t	\$	c	t	g	a	a	a	c	t	g	g	10
t	g	a	a	a	c	t	g	g	t	\$	c	1
t	g	g	t	\$	c	t	g	a	a	a	c	7

Look up ctgg: start at the end with g, lookup in genome dictionary

$\text{top}(g) = 6$ ;  $\text{bot}(g) = 8$

---

0	\$	c	t	g	a	a	a	c	t	g	g	t
1	a	a	a	c	t	g	g	t	\$	c	t	g
2	a	a	c	t	g	g	t	\$	c	t	g	a
3	a	c	t	g	g	t	\$	c	t	g	a	a
4	c	t	g	a	a	a	c	t	g	g	t	\$
5	c	t	g	g	t	\$	c	t	g	a	a	a
6	g	a	a	a	c	t	g	g	t	\$	c	t
7	g	g	t	\$	c	t	g	a	a	a	c	t
8	g	t	\$	c	t	g	a	a	a	c	t	g
9	t	\$	c	t	g	a	a	a	c	t	g	g
10	t	g	a	a	a	c	t	g	g	t	\$	c
11	t	g	g	t	\$	c	t	g	a	a	a	c

Does gg exist, and what are top(gg) and bot(gg)?

Yes, gg exists.

$$\text{top}(gg) = \text{top}(g) + \#g \text{ before } g\text{-block in bwt} = 6 + 1 = 7$$

$$\text{bot}(gg) = \text{top}(gg) + \# \text{ of } gg \text{ in genome} - 1 = 7 + 1 - 1 = 7$$

---

0	\$	c	t	g	a	a	a	c	t	g	g	t
1	a	a	a	c	t	g	g	t	\$	c	t	g
2	a	a	c	t	g	g	t	\$	c	t	g	a
3	a	c	t	g	g	t	\$	c	t	g	a	a
4	c	t	g	a	a	a	c	t	g	g	t	\$
5	c	t	g	g	t	\$	c	t	g	a	a	a
6	g	a	a	a	c	t	g	g	t	\$	c	t
7	g	g	t	\$	c	t	g	a	a	a	c	t
8	g	t	\$	c	t	g	a	a	a	c	t	g
9	t	\$	c	t	g	a	a	a	c	t	g	g
10	t	g	a	a	a	c	t	g	g	t	\$	c
11	t	g	g	t	\$	c	t	g	a	a	a	c

Does tgg exist, and what are top(tgg) and bot(tgg)?

Yes, tgg exists.

$\text{top}(\text{tgg}) = \text{top}(\text{t}) + \#t \text{ before gg-block in bwt} = 9 + 2 = 11$

$\text{bot}(\text{tgg}) = \text{top}(\text{tgg}) + \# \text{ of tgg in genome} - 1 = 11 + 1 - 1 = 11$

---

0	\$	c	t	g	a	a	a	c	t	g	g	t
1	a	a	a	c	t	g	g	t	\$	c	t	g
2	a	a	c	t	g	g	t	\$	c	t	g	a
3	a	c	t	g	g	t	\$	c	t	g	a	a
4	c	t	g	a	a	a	c	t	g	g	t	\$
5	c	t	g	g	t	\$	c	t	g	a	a	a
6	g	a	a	a	c	t	g	g	t	\$	c	t
7	g	g	t	\$	c	t	g	a	a	a	c	t
8	g	t	\$	c	t	g	a	a	a	c	t	g
9	t	\$	c	t	g	a	a	a	c	t	g	g
10	t	g	a	a	a	c	t	g	g	t	\$	c
11	t	g	g	t	\$	c	t	g	a	a	a	c

Does ctgg exist, and what are top(ctgg) and bot(ctgg)?

Yes, ctgg exists.

$\text{top}(\text{ctgg}) = \text{top}(\text{c}) + \# \text{c before tgg-block in bwt} = 4 + 1 = 5$

$\text{bot}(\text{ctgg}) = \text{top}(\text{ctgg}) + \# \text{ of ctgg in genome} - 1 = 5 + 1 - 1 = 5$

---

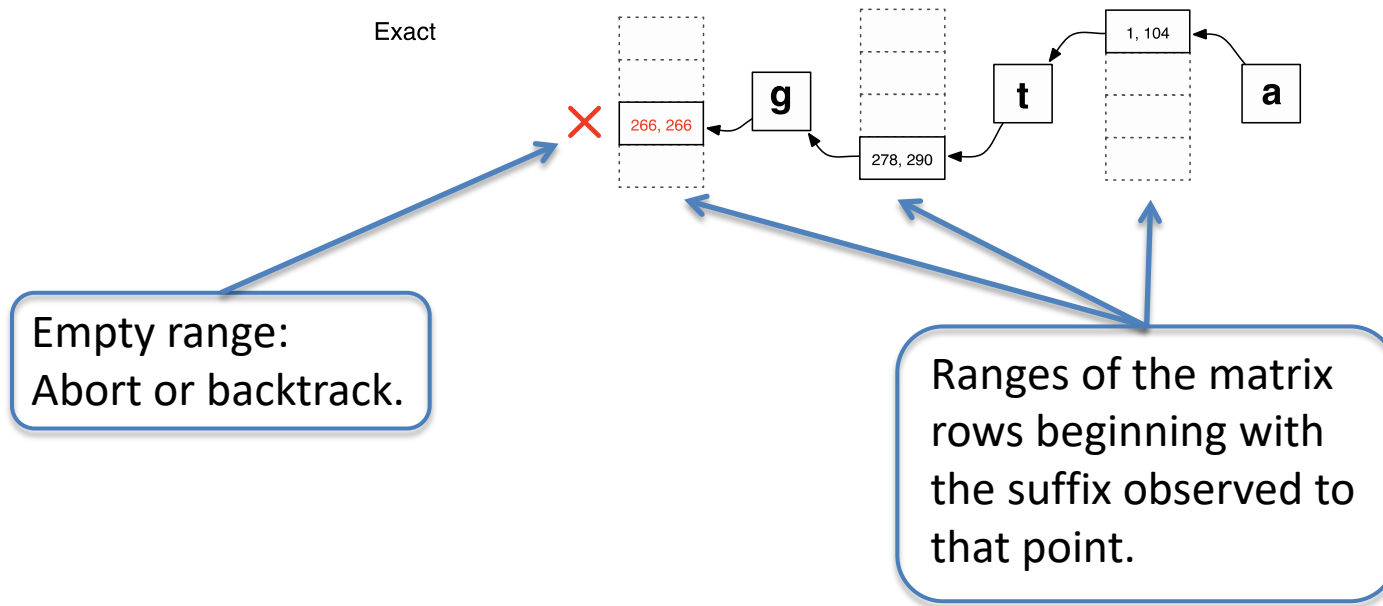
0	\$	c	t	g	a	a	a	c	t	g	g	t
1	a	a	a	c	t	g	g	t	\$	c	t	g
2	a	a	c	t	g	g	t	\$	c	t	g	a
3	a	c	t	g	g	t	\$	c	t	g	a	a
4	c	t	g	a	a	a	c	t	g	g	t	\$
5	c	t	g	g	t	\$	c	t	g	a	a	a
6	g	a	a	a	c	t	g	g	t	\$	c	t
7	g	g	t	\$	c	t	g	a	a	a	c	t
8	g	t	\$	c	t	g	a	a	a	c	t	g
9	t	\$	c	t	g	a	a	a	c	t	g	g
10	t	g	a	a	a	c	t	g	g	t	\$	c
11	t	g	g	t	\$	c	t	g	a	a	a	c

Found ctgg at position 5, which is position 6 in original string ctgaaactggt

0	\$	c	t	g	a	a	a	c	t	g	g	t	11
1	a	a	a	c	t	g	g	t	\$	c	t	g	3
2	a	a	c	t	g	g	t	\$	c	t	g	a	4
3	a	c	t	g	g	t	\$	c	t	g	a	a	5
4	c	t	g	a	a	a	c	t	g	g	t	\$	0
5	c	t	g	g	t	\$	c	t	g	a	a	a	6
6	g	a	a	a	c	t	g	g	t	\$	c	t	2
7	g	g	t	\$	c	t	g	a	a	a	c	t	8
8	g	t	\$	c	t	g	a	a	a	c	t	g	9
9	t	\$	c	t	g	a	a	a	c	t	g	g	10
10	t	g	a	a	a	c	t	g	g	t	\$	c	1
11	t	g	g	t	\$	c	t	g	a	a	a	c	7



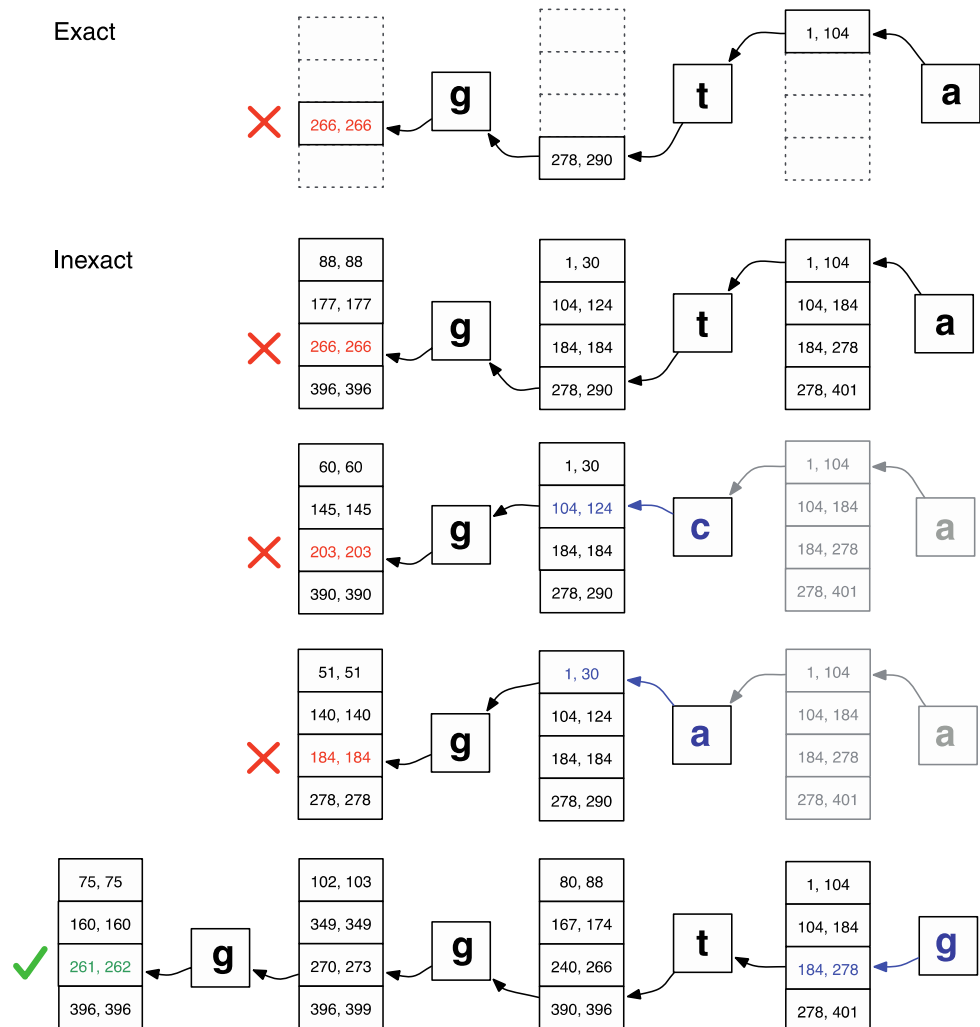
# Coping with mismatches: Query GGTA







# Quality-aware, greedy, randomized, depth-first search through the space of possible alignments.



1. The search proceeds similarly to EXACTMATCH
2. If range becomes empty (suffix does not occur in text), the algorithm backtracks and selects an already-matched query position and substitutes a different base there. The EXACTMATCH algorithm resumes from this modified position.
3. The algorithm allows only substitutions that yield a modified suffix that occurs at least once in the text. If there are multiple candidate substitution positions, then the algorithm greedily selects a position with a **minimal quality** value.
4. Because search is greedy, the first valid alignment is not necessarily the best (Number of mismatches and quality). Bowtie has parameters to cope with this (--best or -all (all alignments)).
5. Excessive Backtracking should be avoided. Note, we start from the **low quality** end...

# Coping with mismatches

---

- Our alignment policy allows a limited number of mismatches and prefers alignments where the sum of the quality values at all mismatched positions is low.
- The search proceeds similarly to EXACTMATCH, calculating matrix ranges for successively longer query suffixes.
- If the range becomes empty (a suffix does not occur in the text), then the algorithm may select an already-matched query position and substitute a different base there. The EXACTMATCH algorithm resumes from this position.
- The algorithm selects only those substitutions that yield a modified suffix that occurs at least once in the text. If there are multiple candidate substitution positions, then the algorithm greedily selects a position with a minimal quality value.

# Bowtie: Avoidance of excessive backtracking (assuming 1 mismatch)

**Problem:** The aligner spends most of its effort fruitlessly backtracking to positions close to the 3' end of the query (error prone).

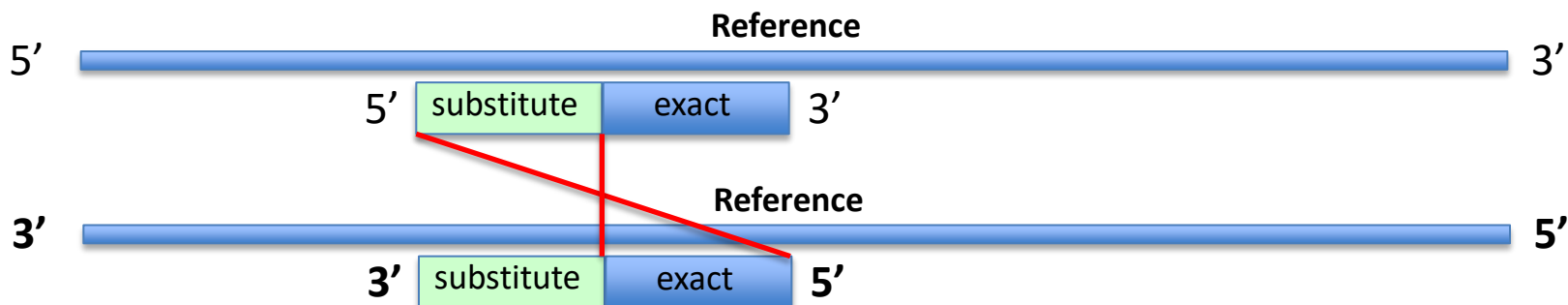
**Solution Part1:** double indexing (similar to MAQ), using two indices for the genome

- Index 1: BWT of the original genome
- Index 2: BWT of the genome with **reversed character order** (**not reverse complemented!**)

**Solution Part2:** The aligner is invoked twice

- First round: Index 1 is used, and the aligner is started with original read with the constraint that it must not substitute a position in the query's right half (3' end).
- Second round: Index 2 is used, and the aligner is started with the reversed read, again with the constraint that it must not substitute a position in the reversed query's right half (originally and still the 5' end).

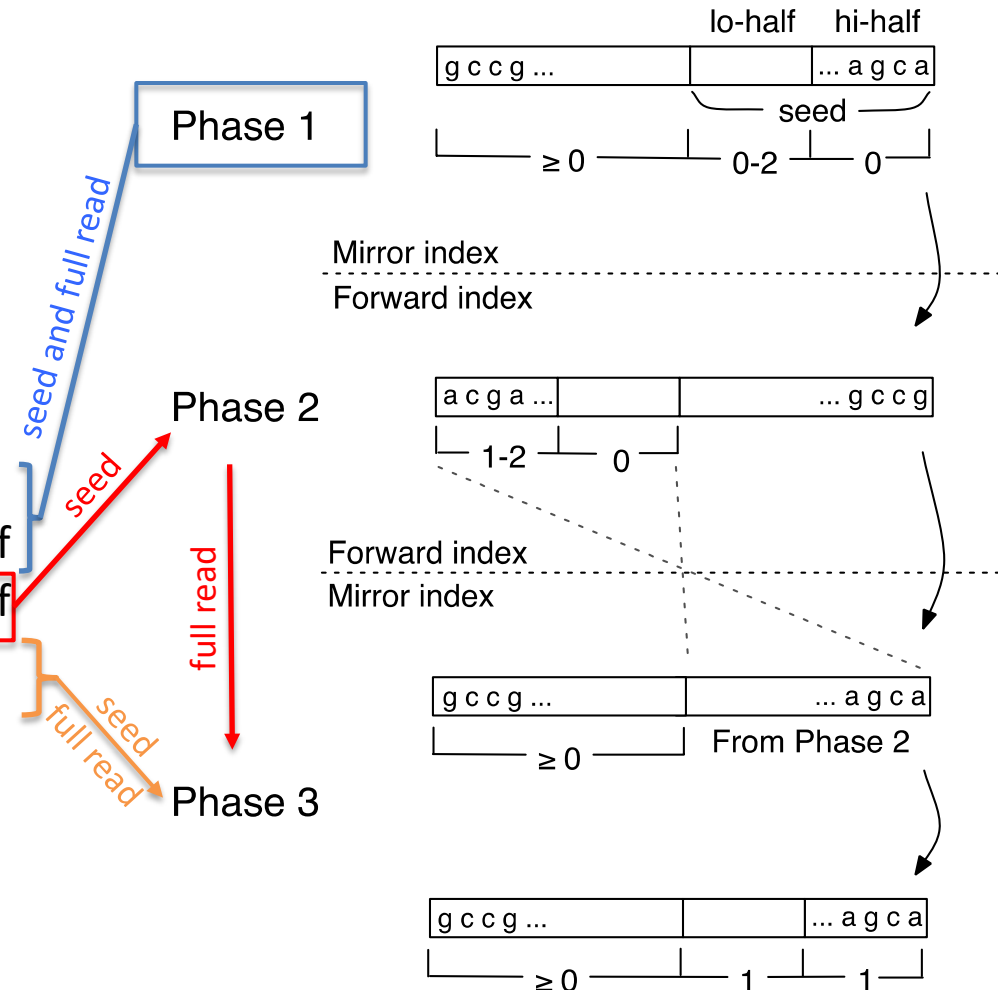
**Solution Part3:** set a hard upper limit of backtracks to be performed.



# The three phases of Bowtie

**In the case of 2 (or more mismatches):**

- Bowtie uses the first 28 bp as seed
- The seed is split into a **high quality 5'** half (hi-half) and a **low quality 3'** half (lo-half)
- For up to 2 mismatches we have four scenarios:
  1. no mismatches in seed
  2. 1-2 mismatches only in the lo-half
  3. 1-2 mismatches only in the hi-half
  4. 1 mismatch each in hi- and low-half
- Any number of mismatches can occur in non-seed part (subject to other thresholds).



# Changes in Bowtie2

---

Supports gapped, local, and paired-end alignment modes:

- For reads longer than about 50 bp Bowtie 2 is generally faster, more sensitive, and uses less memory.
- Bowtie 2 supports [local alignment](#), which doesn't require reads to align end-to-end. Local alignments might be “trimmed” (“soft clipped”) at one or both extremes in a way that optimizes alignment score.
- There is no upper limit on read length in Bowtie 2.
- Bowtie 2 allows alignments to [overlap ambiguous characters](#) (e.g. Ns) in the reference.

# Bowtie2: Seed extraction & alignment

---

## **Seed extraction:**

- Substrings of the read (“seed strings”) are extracted at regular intervals along the read and its reverse complement.
- Seed strings are contiguous (i.e. they are not spaced seeds) and may or may not overlap each other.

## **FM Index assisted seed alignment :**

- Find ungapped alignments for each based on Bowtie1
- Seed strings can be aligned with up to 1 mismatch.

# Bowtie2: Seed alignment prioritization & alignment

---

## **Prioritization:**

- “seed-hit range” → A seed-hit range describes a range of rows in the Burrows-Wheeler matrix that begin with a reference substring that is within 0 or 1 mismatches of the seed substring.
- Bowtie 2 proceeds by repeatedly selecting a row in a random weighted fashion using these weights.

## **Alignment:**

- Bowtie 2 extracts flanking characters from the reference
- Solves a rectangular dynamic programming problem to find high-scoring full alignments in the vicinity of the seed hit.



## Watch out for the following...

---

- Up to 2 mismatches in first 28 bases reported by Maq (default) and 1-2 mismatches by Bowtie1/2 in the seed.
- If  $\geq 1$  matching seed regions exist for one read Bowtie2 has a dynamic programming effort limit of 15. After 15 attempts not performing better than the best so far.
- Bowtie outperforms Bowtie2 in case of short reads ( $\leq 50$  nt) in some cases.
- Bowtie convert N to A, C, G or T randomly / Bowtie2 accepts N's



**THANK YOU**

**GRACIAS**  
**ARIGATO**  
**SHUKURIA**  
**JUSPAXAR**  
**DANKSCHEEN**  
**TASHAKKUR ATU**  
**YAQHANYELAY**  
**SUKSAMA**  
**EKHMET**  
**TINGKI**  
**BI'YAN**  
**SHUKRIA**  
**GOZAIMASHITA**  
**EFCHARISTO**  
**KOMAPSUMINIDA**  
**MAAKE**  
**GRAZIE**  
**MEHRBANI**  
**PALDIES**  
**BOLZIN**  
**MERCI**