# Molecular Evolution and **Bioinformatics**

Reference based sequence mapping

# Literature

How to map billions of short reads onto
genomes?
Trapnell and  Salzberg Nature Biotechnology 27,
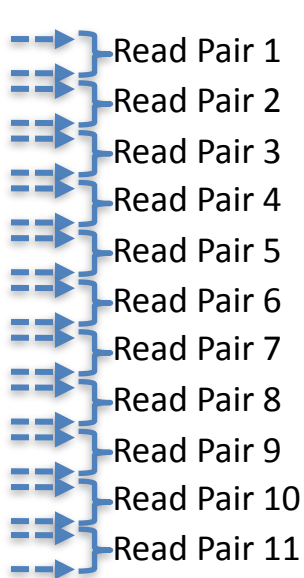455 - 457 (2009)
doi:10.1038/nbt0509-455

Ultrafast and memory-efficient alignment of
short DNA sequences to the human genome.
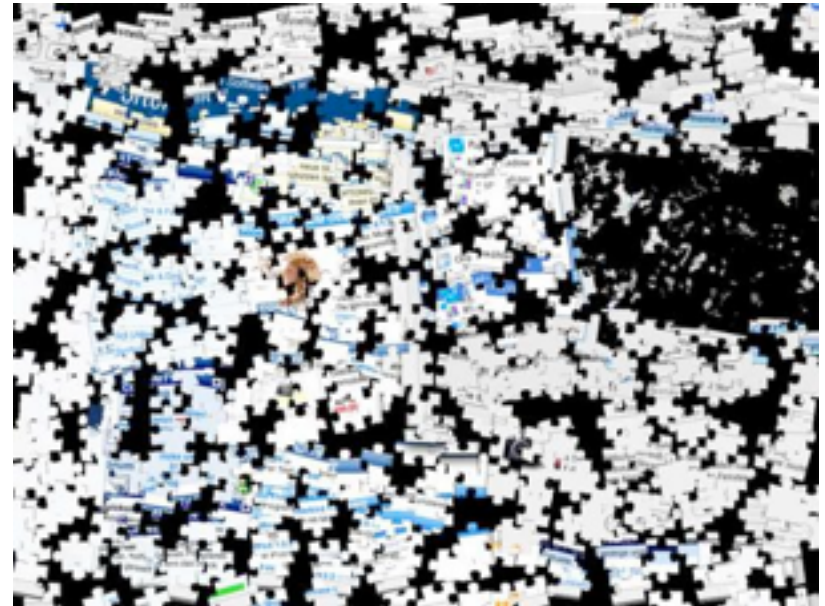Langmead et al. Genome Biology **10**:R25 (2009)
doi:10.1186/gb-2009-10-3-r25

# What do we need short read mapping for?
## Assembly of whole genome shotgun sequencing data

Read Pair 1
Read Pair 2
Read Pair 3
Read Pair 4
Read Pair 5
Read Pair 6
Read Pair 7
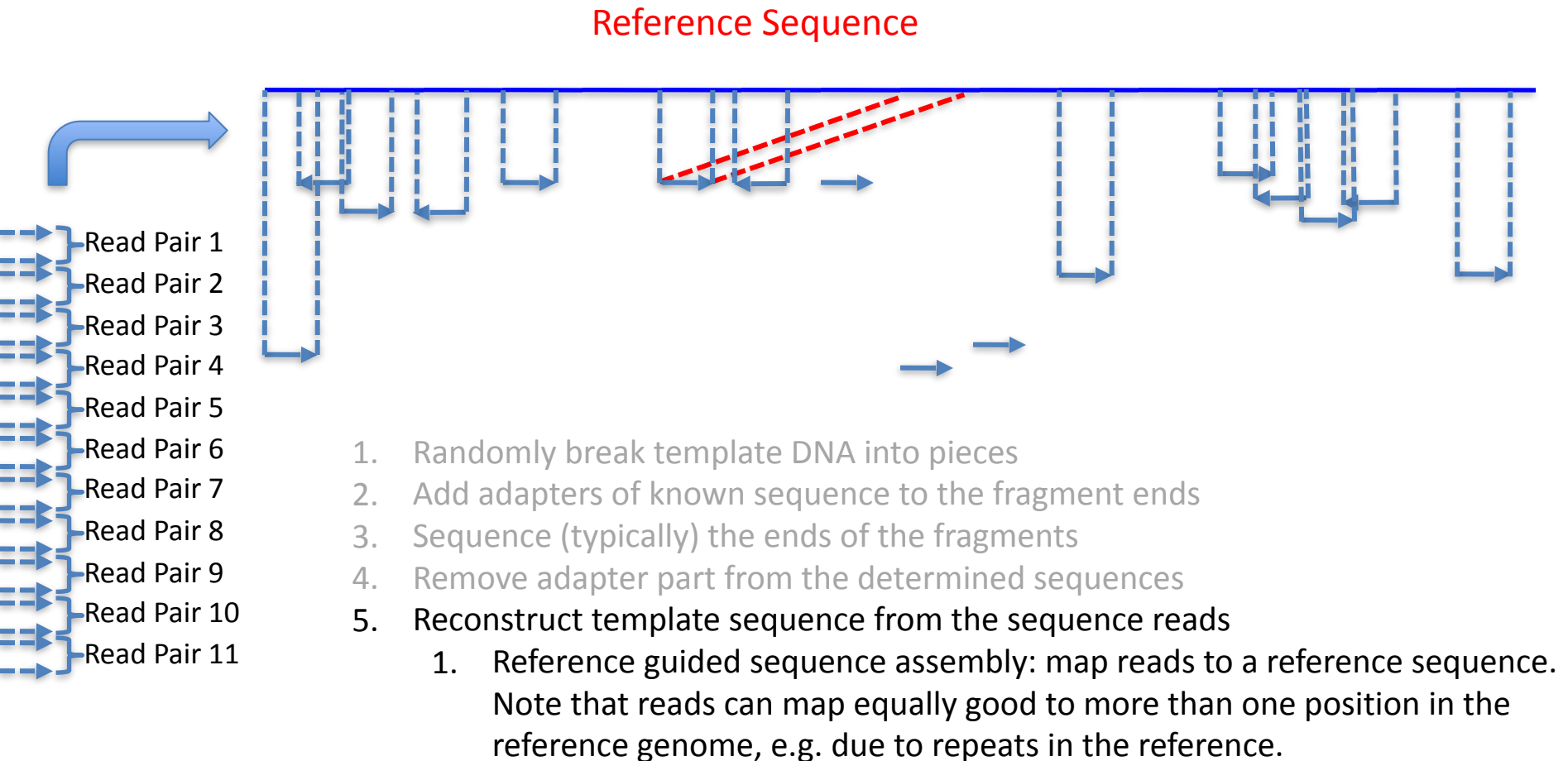Read Pair 8
Read Pair 9
Read Pair 10
Read Pair 11

Reconstruct template

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

# Shotgun sequencing and **reference guided** sequence assembly

Reference Sequence

Read Pair 1
Read Pair 2
Read Pair 3
Read Pair 4
Read Pair 5
Read Pair 6
Read Pair 7
Read Pair 8
Read Pair 9
Read Pair 10
Read Pair 11

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. Remove adapter part from the determined sequences
5. Reconstruct template sequence from the sequence reads
    1. Reference guided sequence assembly: map reads to a reference sequence. Note that reads can map equally good to more than one position in the reference genome, e.g. due to repeats in the reference.
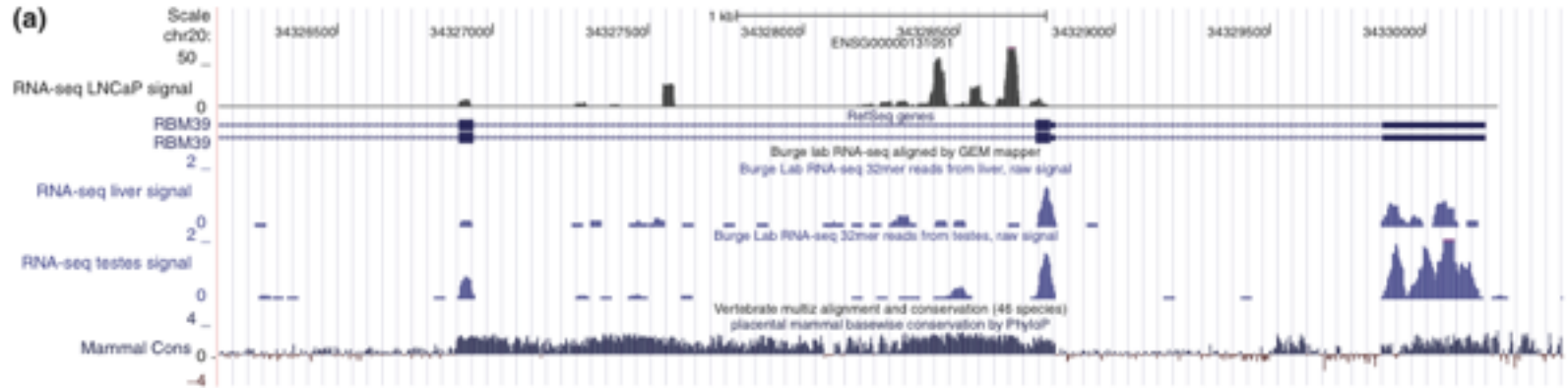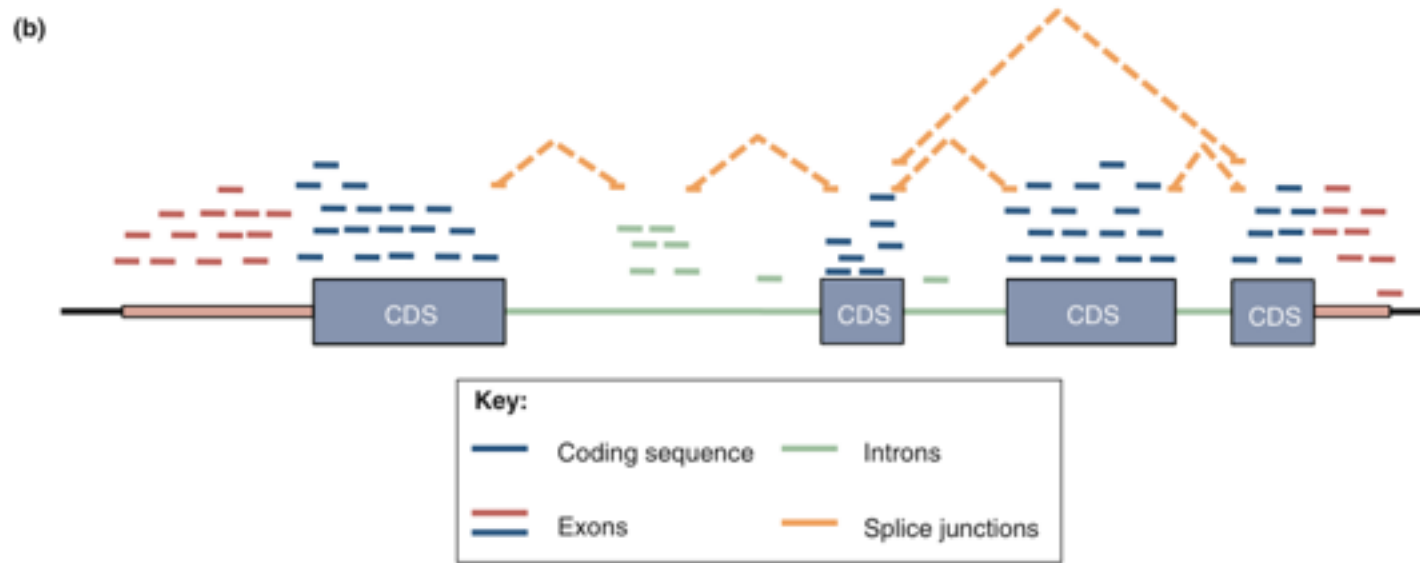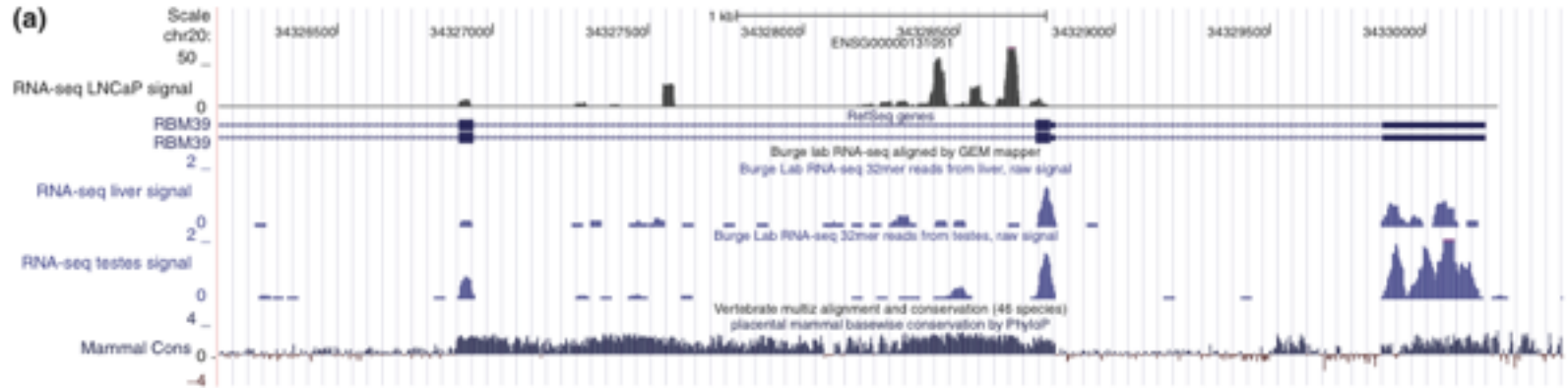
—— Reference sequence
   a) e.g. genome of a different individual from the same species to study species diversity
   b) e.g. genome of a closely related species

# Further short read applications: Mapping of RNA seq

# RNA-seq mapping helps building hypotheses concerning gene structure

# Further short read applications

- ## Genotyping

Goal: identify variations

```
                                               GGTATAC...
...CCATAG    TATGCGCCC      CGGAAATTT CGGTATAC
...CCAT    CTATATGCG      TCGGAAATT  CGGTATAC
...CCAT GGCTATATG       CTATCGGAA    GCGGTATA
...CCA AGGCTATAT      CCTATCGGA    TTGCGGTA    C...
...CCA AGGCTATAT    GCCCTATCG      TTTGCGGT      C...
...CC   AGGCTATAT    GCCCTATCG  AAATTTGC     ATAC...
...CC  TAGGCTATA GCGCCCTA    AAATTTGC  GTATAC...
```
**...CCATAGGCTATATGCGCCCTATCGGCAATTTGCGGTATAC...**

- ## ChIP-seq, CLIP-Seq, Methyl-seq

```
                    GAAATTTGC
                    GGAAATTTG
                   CGGAAATTT
                   CGGAAATTT
                  TCGGAAATT
                 CTATCGGAAA
               CCTATCGGA        TTTGCGGT
             GCCCTATCG      AAATTTGC
             GCCCTATCG      AAATTTGC      ATAC...
...CC
```
**...CCATAGGCTATATGCGCCCTATCGGCAATTTGCGGTATAC...**

Goal: classify, measure significant peaks

**ChIP:** Chromatin Immunoprecipitation, **CLIP:** Cross-linking immunoprecipitation

# 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

- 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

```
...TGATCATA...  better than  ...TGATCATA...
   | | | | |                 | |     | |
   GATCAA                    GAGAAT

...TGATATTA...  better than  ...TGATcaTA...
   | | |   |                 |     | | |
   GATcaT                    GTACAT
```

Finding mapping positions of reads in genomes is, in principle, a pattern matching problem. We search for the occurrence of short words in long sentences.

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

```
Searchstring: AATGAGACATGAA

      Query1: CATG
      Query2: ATGT
```

Finding mapping positions is, in principle, very easy. Just slide the word along the sequence and stop when either end of sequence is reached or mapping position is found

Searchstring: AATGAGACATGAA

Query1: CATG
Query2: ATGT

Naïve approach

```
                      1234567891234
          Searchstring: AATGAGACATGAA
                        CATG
```

Finding mapping positions is, in principle, very easy. Just slide the word along the sequence and stop when either end of sequence is reached or mapping position is found

Searchstring: AATGAGACATGAA

Query1: CATG
Query2: ATGT



Naïve approach

```
           1234567891234
Searchstring: AATGAGACATGAA
           CATG
           ✖➕➕➕
```

Finding mapping positions is, in principle, very easy. Just slide the word along the sequence and stop when either end of sequence is reached or mapping position is found

Searchstring: AATGAGACATGAA

Query1: CATG
Query2: ATGT



Naïve approach

```
             1234567891234
Searchstring: AATGAGACATGAA
             CATG
             ✘✘✘✘
```

Finding mapping positions is, in principle, very easy. Just slide the word along the sequence and stop when either end of sequence is reached or mapping position is found

Searchstring: AATGAGACATGAA

Query1: CATG
Query2: ATGT



Naïve approach

```
             1234567891234
Searchstring: AATGAGACATGAA
                CATG
                ✗✗✗✚
```

Finding mapping positions is, in principle, very easy. Just slide the word along the sequence and stop when either end of sequence is reached or mapping position is found

```
Searchstring: AATGAGACATGAA
```

```
Query1: CATG
Query2: ATGT
```



```
              1234567891234
Searchstring: AATGAGACATGAA
                     CATG
                     ++++
```

Naïve approach

Full match found,
Output result

Query1 maps to position 8-12 in searchstring

Finding mapping positions is, in principle, very easy. Just slide the word along the sequence and stop when either end of sequence is reached or mapping position is found

Searchstring: AATGAGACATGAA

Query1: CATG
Query2: ATGT



Naïve approach

1234567891234
Searchstring: AATGAGACATGAA
TTGT

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.

# Indexing speeds up searches

INDEX

# 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

| Kmer | Position |
|------|----------|
| GAG | 1 |
| AGA | 2 |

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

# Main differences between mapping approaches

1) How does the mapper decide when the query maps multiple times to the sequence?

Mapping position of Q2: 1-4 **AND** 8-11

Relevant e.g. for RNA seq quantifying gene expression and for Chip-Seq as you must not consider a read more than once.

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

Q2

| Kmer | Position |
|------|----------|
| CAT | 1 |
| ATG | 2 |

# Main differences between mapping approaches

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

| Kmer | Position |
|------|----------|
| ATG | 1 |
| TGT | 2 |

One mismatch

The 2nd 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**



**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**

# The "traditional way": Hash tables

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

# 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 $h_P$.



'AGC'

Sum up the values* for each letter in $P$, divide by prime number *mod* and take the rest as hash value

$h_{AGC} = (1+2+3)\%3$

0

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

1  2  1  0  1  1  1
1  2  2  2  2  1  0  2

**ACTTGAACAAGCTTGAGATCAGGAGGGGAGA**

1  1  0  1  2  1  1  0
2  1  1  1  2  2  1  1

1) 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 $h_k$



1) 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!

# MAQ* uses **spaced seeds** 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'

gatgtgacatacctgttctactgaggct     HASH → 2057673064

gatgtgacatacctgttctactgaggct → 3178370917

gatgtgacatacctgttctactgaggct → 773088662

gatgtgacatacctgttctactgaggct → 1856750201

gatgtgacatacctgttctactgaggct → 2510061809

gatgtgacatacctgttctactgaggct → 119777054

hash value for template 1: 832589471

GENOME
ttcagttatccagaaaatgctattttccccaaaatgaaatctaaaatagtaactcaagtg
aaacattgtcaggtgtgtaaggaaggaaaatatgaagggctacccacaaacccacaaac
aacggaaactccaattcctaagtttccgggacacactattcatatagatataatttcta
cagacaaaaacgggtacttacggcaattcacaaattttcaaaacttgcgaaagcaaaaa
taaaaaaaattcaaaatcaatagaagacataagaaaacctttacatgacatcttatttt
attttggagtaccaaaatacgttgtaatggatatagaaaaatcctttaattccgcaaca
acagcctttatgatgaaagaccagctgggcatacaaattttcaaagcatccccttataa
aagttctgtaaacggacaaatagaacggtttcattctatcctcgctgaaattaaaagat
gtttaaaaactaaacaggtacaccgaacatttgaagaacagttcaattcagctgtctag
gaatataactacacaattcaccctgtatcaaaaatacaaaacaaaataacgcccttaaa
aatattttggcagaaatataaccactgatccaggaaaatatgacgaaagcagataagg
caacatcgaaaacctataatcaaaacaggcaacaggcttaaaaaccataacgaaaaaag
caataagatcaaagattatgagccaggacaaacagtttttataaagcaaatcacaaggc
caggttctaagctgtacaagaaagaaacattaaaagaaaacagacaaacatttgttatc
acagaagcaggaagatgtgacatacctgttctactgaaggct

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

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

# Why does MAQ* use **spaced seeds** for mapping?

Reference

AGACTGAGGTACGTAGACCATGATCGATACCCAAAAAGCTAGA

GTACGTAGACGATGATCCATACCCAAAA

Read (28 bp prefix)

MAQ uses seed pairs as it allows per default up to 2 mismatches between seed and reference. As each 28 mere is represented by 4 non-overlapping seeds, we are guaranteed that we **always have at least 2 seeds** that must result in a perfect match to the reference.

# 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.

# The "new way": 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

$   **a c a a c g**   S

a c a a c g $
c a a c g $ a
a a c g $ a c
a c g $ a c a
c g $ a c a a
g $ a c a a c
$ a c a a c g

→ sort lexicographically →

$ a c a a c **g**
a a c g $ a **c**
a c a a c g **$**
a c g $ a c **a**
c a a c g $ **a**
c g $ a c a **a**
g $ a c a a **c**

BWT(S)

# Burrows Wheeler Transform (BWT)



(a)

$$acaacg\$ \rightarrow \begin{array}{l} \$\,a\,c\,a\,a\,c\,g \\ a\,a\,c\,g\,\$\,a\,c \\ a\,c\,a\,a\,c\,g\,\$ \\ a\,c\,g\,\$\,a\,c\,a \\ c\,a\,a\,c\,g\,\$\,a \\ c\,g\,\$\,a\,c\,a\,a \\ g\,\$\,a\,c\,a\,a\,c \end{array} \rightarrow g\,c\,\$\,a\,a\,a\,c$$

BWT(S)

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 ith occurrence of character X in the last column corresponds to the same text character as the ith occurrence of X in the first column

# Burrows Wheeler Transform (BWT)



**(a)**

```
        $ a c a a c g
  1st   a a c g $ a c
        a c a a c g $
a c a a c g $ → a c g $ a c a → g c $ a a a c
        c a a c g $ a
        c g $ a c a a
        g $ a c a a c
```

BWT(S)

**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)



BWT(S)

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.

# How to cope with mismatches?
# Query GGTA



Exact

1, 104

g    t    a

266, 266    278, 290

Empty range:
Abort or backtrack.

Ranges of the matrix
rows beginning with
the suffix observed to
that point.

# Bowtie conducts a 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...

# 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).
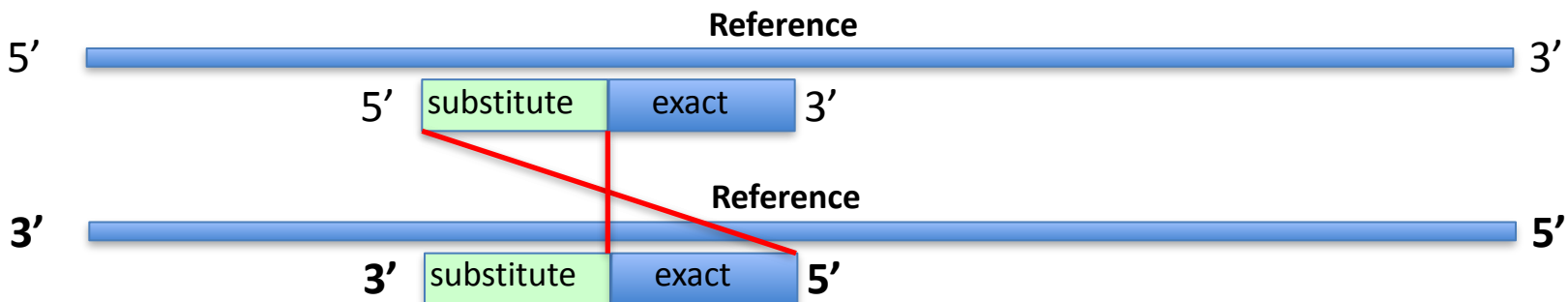
**Solution Part1:** double indexing (similar to MAQ), using two indices for the genome
- Index 1: BWT of the original genome
- Index 2: BWT of the genome with reversed character order (**not reverse complemented!**)

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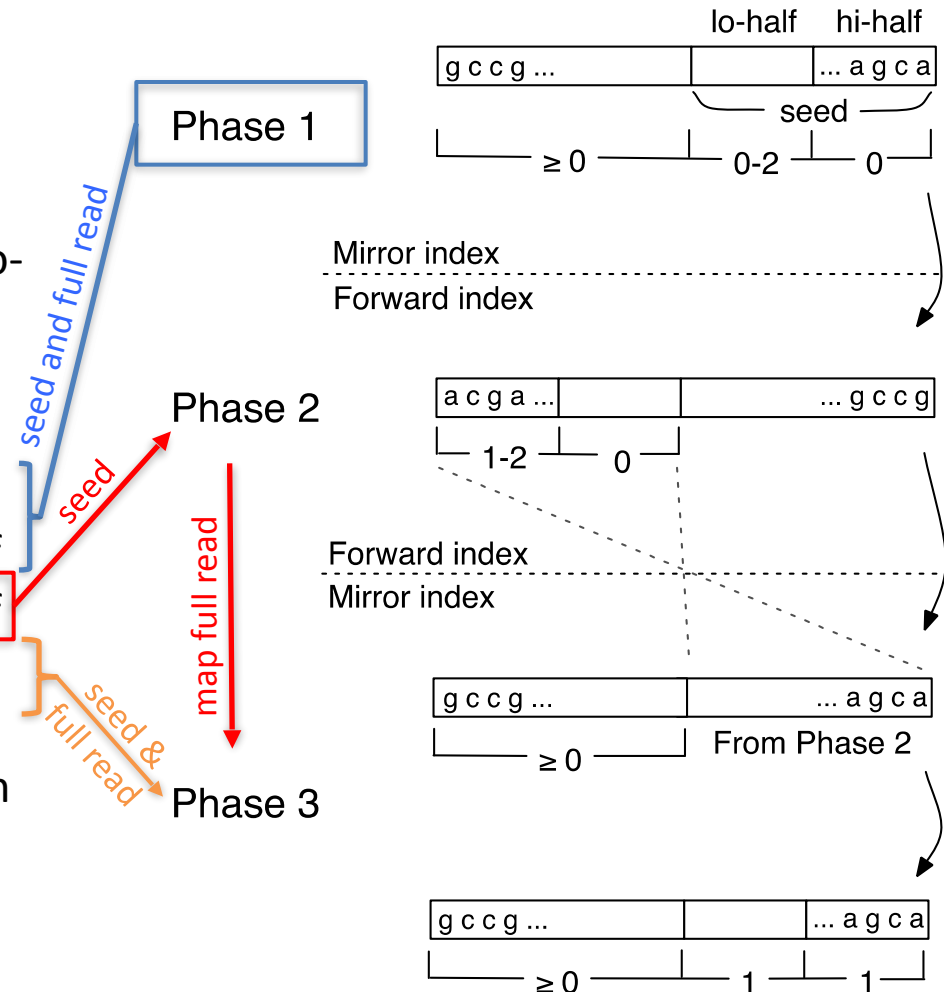
**Solution Part3:** set a hard upper limit of backtracks to be performed.

# 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).

# SAM/BAM format* & SAMtools

Header section

```
@HD    VN:1.0
@SQ    SN:chr20 AS:HG18 LN:62435964
@RG    ID:L1 PU:SC_1_10 LB:SC_1 SM:NA12891
@RG    ID:L2 PU:SC_2_12 LB:SC_2 SM:NA12891
```
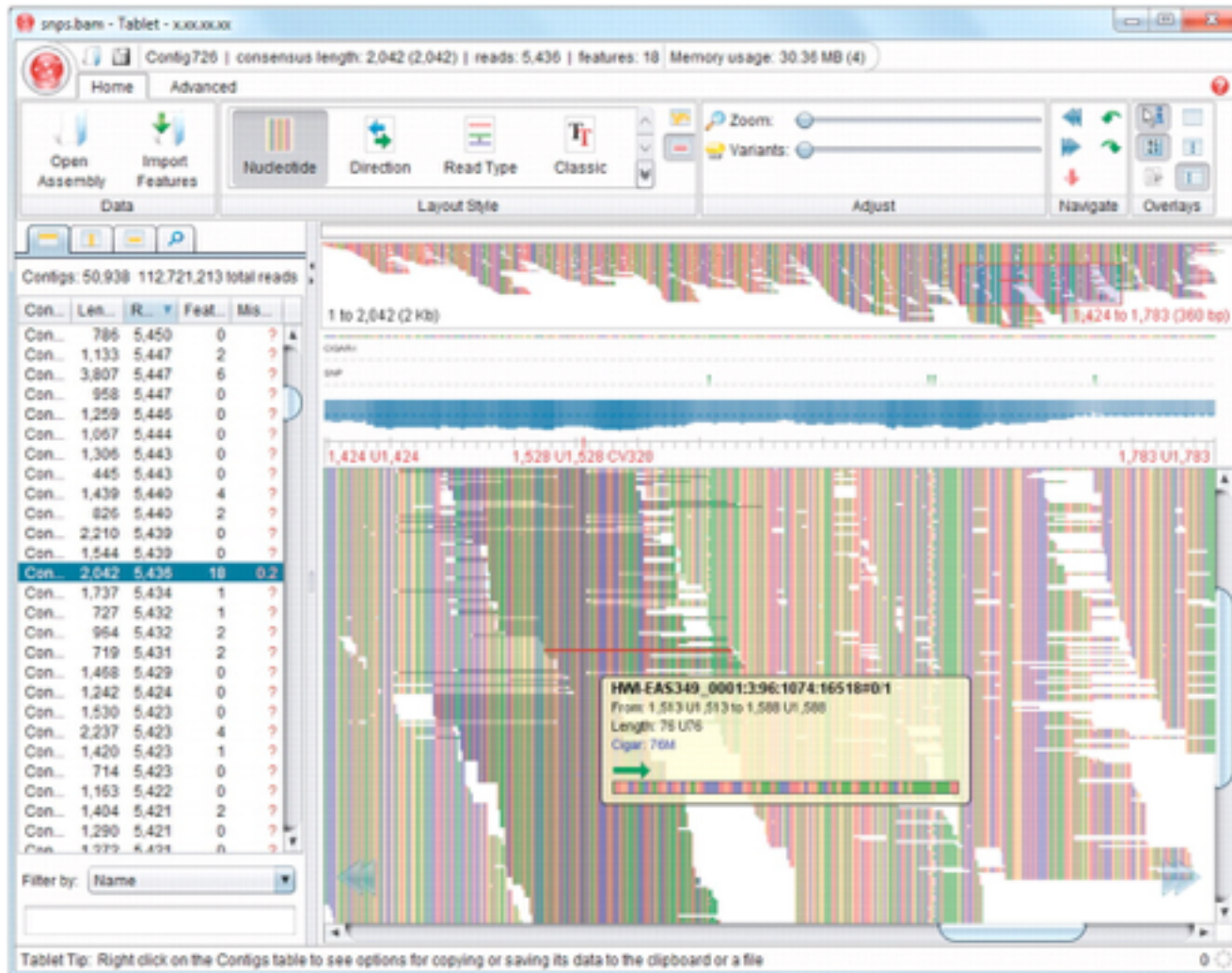
position of alignment

Alignment section

| Query Name | | Ref sequence | | | | | | | query sequence (same strand as ref) | query quality | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| V00-HWI-EAS132:3:38:959:2035#0 | 147 | chr1 | 28 | 255 | 36M | = | 79 | 0 | GATCTGATGGCAGAAAACCCCTCTCAGTCCGTCGTG | aaX`[\`Y^Y^]ZX``\`EV_BBBBBBBBBBBBBBBB | NM:i:1 |
| V00-HWI-EAS132:4:99:122:772#0 | 177 | chr1 | 32 | 255 | 36M | = | 9127 | 0 | AAAGGATCTGATGGCAGAAAACCCCTCTCAGTCCGT | aaaaaa\OWaI_\WL\aa`Xa^]\ZUaa[XWT\^XR | NM:i:1 |
| V00-HWI-EAS132:4:44:473:970#0 | 25 | chr1 | 40 | 255 | 36M | * | 0 | 0 | GTCGTGGTGAAGGATCTGATGGCAGAAAACACCTCT | __YaZ`W[aZNUZ[U[_TL[KVVX^QURUTDRVZBB | NM:i:2 |
| V00-HWI-EAS132:4:29:113:1934#0 | 99 | chr1 | 44 | 255 | 36M | = | 107 | 0 | GGGTTTTCTGCCATCAGATCCTTTACCACGACAGAC | aaaQaa__``]\\_^`¯^a^`a`_^^^_XQ[ZS\XX | NM:i:1 |

- Sequence Alignment/Map & Binary Alignment/Map (Heng Li, et al)

- SAMtools - C or Java
  - convert formats
  - sort/merge alignments
  - generate per-position information
  - call SNPs/indel variants
  - show alignments in a text viewer
  - http://bioinformatics.oxfordjournals.org/cgi/reprint/btp352v1

*for details on the format see http://samtools.sourceforge.net/SAM1.pdf

# Data visualization using Tablet

Milne et al. Brief Bioinform (2013) 14 (2): 193-202.

# Mapping: Watch out for the following…

- Up to 2 mismatches in first 28 bases reported by default (Maq, Bowtie, BWA). Maq additionally reports some 3-mismatches.

- If >=1 exact matches exist for a read, Bowtie will report one of them. If the best match is inexact, bowtie will not always return the highest quality alignment (unless run with --best).

- Bowtie, BWA, and SOAP2 all randomly place repeat reads.

- Although in theory BWA works with arbitrarily long reads,its performance is degraded on long reads (probably with any BWT aligner). Consider using Bowtie2.

- Non-ACTG bases in the genome are converted to random bases in BWA

- How does the mapper deal with paired end reads? Will it report a mapping success only when both reads map in the correct distance (specify insert size) and orientation?