# Short Read Alignment
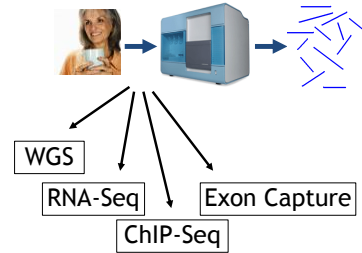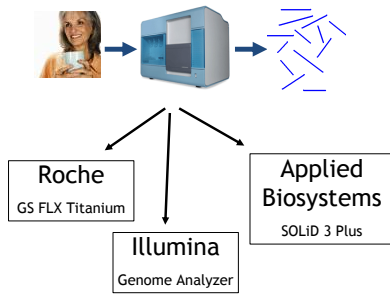
Tobias Rausch
7th June 2010
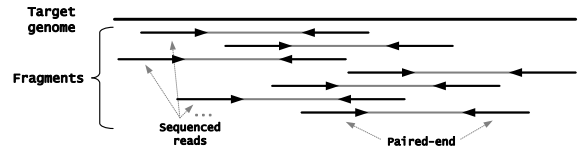
# Sequencing
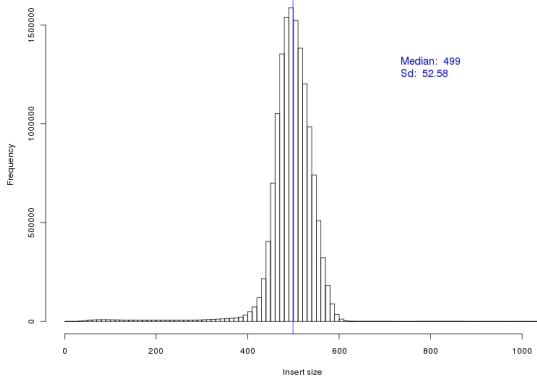


# Sequencing
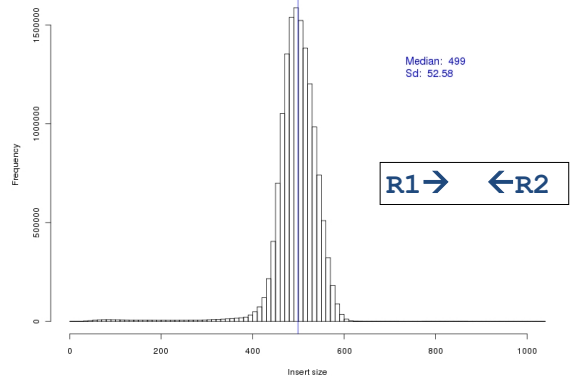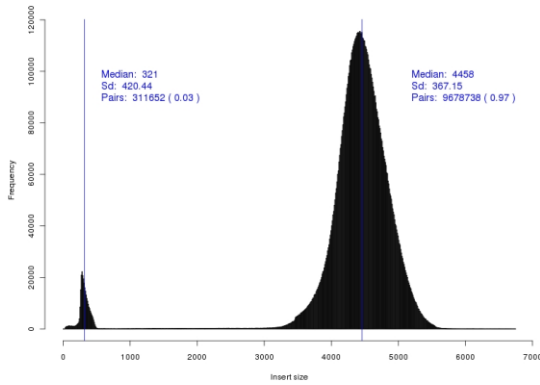


# Paired-End Sequencing



# Paired-End Libraries



Median: 499
Sd: 52.58

# Paired-End Libraries



Median: 499
Sd: 52.58

R1→   ←R2

## Mate-Pair Libraries



Median: 321
Sd: 420.44
Pairs: 311652 ( 0.03 )

Median: 4458
Sd: 367.15
Pairs: 9678738 ( 0.97 )

## Mate-Pair Libraries



Median: 321
Sd: 420.44
Pairs: 311652 ( 0.03 )

Median: 4458
Sd: 367.15
Pairs: 9678738 ( 0.97 )

←R2    R1→

R1→    ←R2

## Data Analysis



## Data Analysis



**Read Mapping**

Reference

**Assembly**

Overlap

Layout

Consensus

## Assembly

- String Graph Assembler
  - Overlap - Layout - Consensus assemblers
  - Examples
    - *Celera Assembler, Arachne, Atlas*

- De-Bruijn Graph Assembler
  - Short-read assemblers
  - Examples:
    - *Velvet, Abyss, SOAPdenovo*
  - Transcriptome assembly: *Oases*

## Read Mapping



Reference

## Read Mapping

Reference



Reference Genome
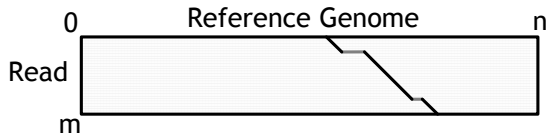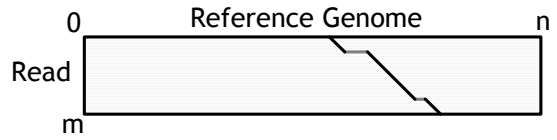
0    n

Read

m

## Read Mapping

Reference



Reference Genome

0    n

Read

m

- Quadratic algorithm
  - Requires $O(m*n)$ time and space
- Infeasible for millions of short reads

## Filtering

*Genome*

Preprocess

Index

## Filtering

*Genome*

Preprocess

*Read*

Filter Algorithm ← Index

*Filtration Phase*

*Potential Matches*

## Filtering

*Genome*

Preprocess

*Read*

Filter Algorithm ← Index

*Filtration Phase*

*Potential Matches*

Exact Algorithm

*Verification Phase*

*True Matches*
*False Matches*

## Simple k-mer Index, k=3

S = ACGAAAACTCGATTACTCGACC

|     | Hitlist |     | Hitlist |     | Hitlist |
|-----|---------|-----|---------|-----|---------|
| AAA |         | ACC |         | CGA |         |
| AAC |         | ACG |         | ... |         |
| AAG |         | ACT |         | GAA |         |
| AAT |         | AGA |         | ... |         |
| ACA |         | ... |         | TTT |         |

- Size of that table: $4^3 = 64$ entries = $|\Sigma|^k$

3

## Simple k-mer Index, k=3

S = <span style="color:red">ACG</span>AAAACTCGATTACTCGACC

|      | Hitlist |      | Hitlist |      | Hitlist |
|------|---------|------|---------|------|---------|
| AAA  |         | ACC  |         | CGA  |         |
| AAC  |         | ACG  | 0       | …    |         |
| AAG  |         | ACT  |         | GAA  |         |
| AAT  |         | AGA  |         | …    |         |
| ACA  |         | …    |         | TTT  |         |

- Size of that table: $4^3 = 64$ entries $= |\Sigma|^k$

## Simple k-mer Index, k=3

S = <span style="color:red">ACGA</span>AAACTCGATTACTCGACC

|      | Hitlist |      | Hitlist |      | Hitlist |
|------|---------|------|---------|------|---------|
| AAA  |         | ACC  |         | CGA  | 1       |
| AAC  |         | ACG  | 0       | …    |         |
| AAG  |         | ACT  |         | GAA  |         |
| AAT  |         | AGA  |         | …    |         |
| ACA  |         | …    |         | TTT  |         |

- Size of that table: $4^3 = 64$ entries $= |\Sigma|^k$

## Simple k-mer Index, k=3

S = <span style="color:red">ACGAA</span>AACTCGATTACTCGACC

|      | Hitlist |      | Hitlist |      | Hitlist |
|------|---------|------|---------|------|---------|
| AAA  |         | ACC  |         | CGA  | 1       |
| AAC  |         | ACG  | 0       | …    |         |
| AAG  |         | ACT  |         | GAA  | 2       |
| AAT  |         | AGA  |         | …    |         |
| ACA  |         | …    |         | TTT  |         |

- Size of that table: $4^3 = 64$ entries $= |\Sigma|^k$

## Simple k-mer Index, k=3

S = ACGAAAACTCGATTACTCGACC

|      | Hitlist |      | Hitlist |      | Hitlist |
|------|---------|------|---------|------|---------|
| AAA  | 3,4     | ACC  | 19      | CGA  | 1       |
| AAC  | 5       | ACG  | 0       | …    | …       |
| AAG  | Empty   | ACT  | 6,14    | GAA  | 2       |
| AAT  | Empty   | AGA  | …       | …    | …       |
| ACA  | Empty   | …    | …       | TTT  | Empty   |

## Searching a Read

|      | Hitlist |      | Hitlist |      | Hitlist |
|------|---------|------|---------|------|---------|
| AAA  | 3,4     | ACC  | 19      | CGA  | 1       |
| AAC  | 5       | ACG  | 0       | …    | …       |
| AAG  | Empty   | ACT  | 6,14    | GAA  | 2       |
| AAT  | Empty   | AGA  | …       | …    | …       |
| ACA  | Empty   | …    | …       | TTT  | Empty   |

- Read Sequence: <span style="color:red">ACT</span>G
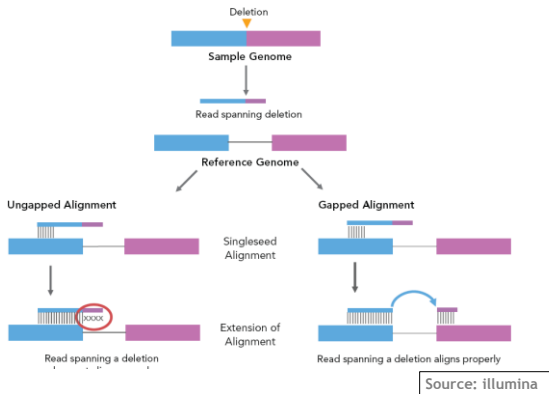  - Potential match at position 6 and 14

## Verification Algorithm
## Banded Dynamic Programming
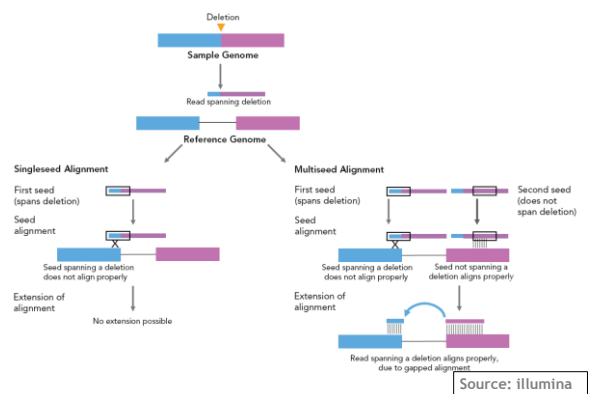
## Techniques

- Index
  - Hash tables, k-mer Index
  - Suffix arrays
  - Burrows-Wheeler-Transformation (BWT) of a suffix array
- Filtering Algorithms
  - Single or multiple seeds
  - Pigeonhole principle
  - Q-gram filtering
- Verification
  - Simple seed-and-extend
  - Banded dynamic programming
  - Quality-based dynamic programming

## Read Mappers



Source: illumina

## ELANDv2 – Gapped Banded Alignment (20bp)



Source: illumina

## ELANDv2 – Multiseed Alignment (Seed max. 2 errors)



Source: illumina

## Parallelization

- Data Decomposition
  - Split the reads
  - Examples: Bowtie, Eland

- Functional Decomposition
  - Separate filtering and verification processes
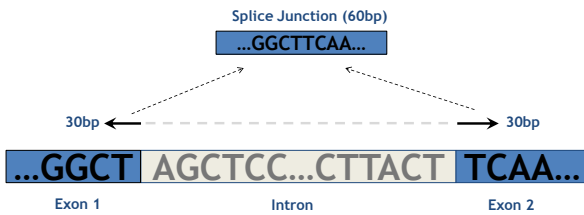
## RNA-Seq

## RNA-Seq



## RNA-Seq

- Read-Mapping Protocol
  - Alignment against contaminants (rRNA)
  - Alignment against splice-junctions
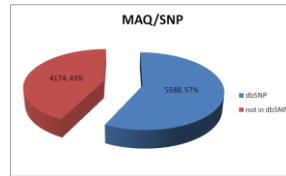  - Alignment against genome

## RNA-Seq

- Read-Mapping Protocol
  - Alignment against contaminants (rRNA)
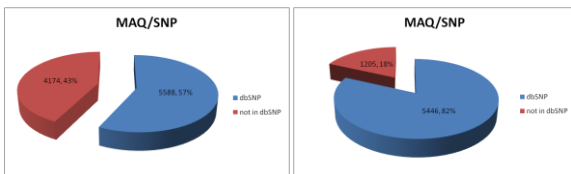  - Alignment against splice-junctions
  - Alignment against genome



## Calling SNPs

- Direct Alignment against hg18



## Calling SNPs

- Direct Alignment against hg18



- Alignment against rRNA (1%) + Alignment against splice junctions (11%)

## SAM/BAM

- Generic format for storing large nucleotide sequence alignments
- SAM Tools
  - Sorting alignments
  - Merging alignments
  - Indexing alignments
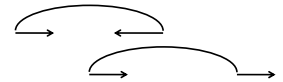  - Viewing alignments

## SAM record

❏Tab-delimited format
  ❏ Field 1: Query name
  ❏ Field 2: Flag
  ❏ Field 3: Reference sequence name
  ❏ Field 4: 1-based leftmost coordinate of the clipped sequence
  ❏ Field 5: Mapping quality
  ❏ Field 6: CIGAR strings
  ❏ Field 7: Mate reference sequence name
  ❏ Field 8: 1-based leftmost coordinate of the clipped sequence
  ❏ Field 9: Insert size (5' to 5')
  ❏ Field 10: Query sequence
  ❏ Field 11: Sequence qualities

## SAM record

❏Tab-delimited format
  ❏ Field 1: Query name
  ❏ Field 2: Flag
  ❏ Field 3: Reference sequence name
  ❏ Field 4: 1-based leftmost coordinate of the clipped sequence
  ❏ Field 5: Mapping quality
  ❏ Field 6: CIGAR strings
  ❏ Field 7: Mate reference sequence name
  ❏ Field 8: 1-based leftmost coordinate of the clipped sequence
  ❏ Field 9: Insert size (5' to 5')
  ❏ Field 10: Query sequence
  ❏ Field 11: Sequence qualities

## Sam / Bam Format



## Sam / Bam Format



- Sequence characters agreeing with the reference are set to " . " or " , " for reads aligned to the forward or reverse strand.

## Sam / Bam Format



CIGAR Strings
- 39M
- 19M1D5M
- 9M1I23M
- 9M2I23M
- 23M2D10M
- 26M
- 12M1D15M
- 16M

- M: Alignment match or mismatch
- I: Insertion to the reference
- D: Deletion from the reference

## Sam / Bam Format



CIGAR Strings
- 39M
- 19M1D5M
- 9M1I1P23M
- 9M2I23M
- 23M2D10M
- 26M
- 12M1D15M
- 16M

- P: Padding (silent deletion)
- This is not even implemented by BWA
  – Because it would require a *de novo local assembler*!

# Sam / Bam Format

- N: Skipped region from the reference
  - For spliced reads:
    - ACATGATA..............................................GAGCTTTA   (Cigar: 8M56N8M)
- Two more CIGAR characters
  - S: Soft clip on the read
  - H: Hard clip on the read

# Flags

Bitwise FLAG:     $f_{15}f_{14}f_{13}f_{12}f_{11}f_{10}f_9f_8f_7f_6f_5f_4f_3f_2f_1f_0$   with  $f_i = \{0,1\}$

$f_0$:  0 = Read is not paired in sequencing, 1 = Read is paired in seq.

$f_1$:  1 = The read is mapped in a proper pair

$f_2$:  1 = The query sequence itself is unmapped

$f_3$:  1 = The mate is unmapped

$f_4$:  0 = forward strand, 1 = reverse strand

...