



# BIOINFORMATICS

## LAB2: 1st Exercise

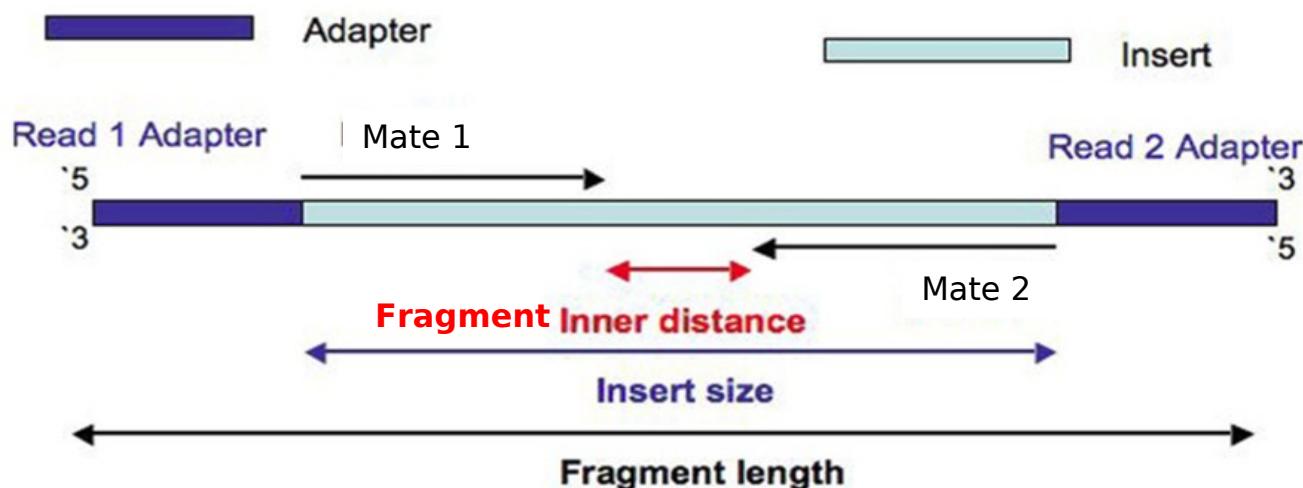
PACIELLO Giulia  
May 5-6th 2016

# REMINDER (1)

- READ: Elementary sequence output of Next Generation Sequencing Machines.

ATCGTCGATCCAGTCCC GG TACGT CGT CCC GGG CT

- SINGLE-END SEQUENCING: Allows users to sequence only one end of a DNA/RNA fragment.
- PAIRED-END SEQUENCING: Allows users to sequence both ends of a DNA/RNA fragment. Paired-end sequencing facilitates detection of genomic rearrangements, as well as gene fusions and novel transcripts.



# REMINDER (2)

- **FASTQ FILE:** Text-based format file storing biological sequences and their corresponding quality scores.

## Fastq File storing mates 1

- **FASTA FILE**: Text-based format file storing biological sequences without quality information.

>HWI-1KL152:175:C3D02ACXX:1:1101:1206:2168/1  
GGGCACCAAAACTTATTGCCAGTGGTAGTTCTGGCAAGGCCATCCAAATAGCAGGTGAAGGCACCTGGCTGACCATCAATGCTTCCACATTGTA  
>HWI-1KL152:175:C3D02ACXX:1:1101:1302:2077/1  
TGCCTAGGGGGAGGGGGCTAGGGACTAGGATGATGGGGGGCAGGATAGTTCAGACGGTTCTATTCCTGAGCGTCTAGATGTTAGCTGGAGAG  
>HWI-1KL152:175:C3D02ACXX:1:1101:1375:2092/1  
GGGGCCCCCTCAGAATGATATTGGCTCACGGGGAGACATAGCCTATGAGGCTGTTCTATAGTTCAGCAGCAGGAGATAATGCCATGTTTCA  
>HWI-1KL152:175:C3D02ACXX:1:1101:1422:2104/1  
GGTGGTTCTGGCTCTGCTGACTGATAACCTTGGCTCAGTTCATCTAACATGATCTTCCCTCTAAATCCAGATCTTGATGCTGGGGCTGTC  
>HWI-1KL152:175:C3D02ACXX:1:1101:1358:2111/1  
GGGGATCTGGTTAACTGGGGGAAGTTCAGCTTCTGCTGATCTAACATGGAGAGGAATAAAGTGGGTAGTAGGGAGGAGCATTCTGAAAGCTGAG

## Fasta File storing mates 1

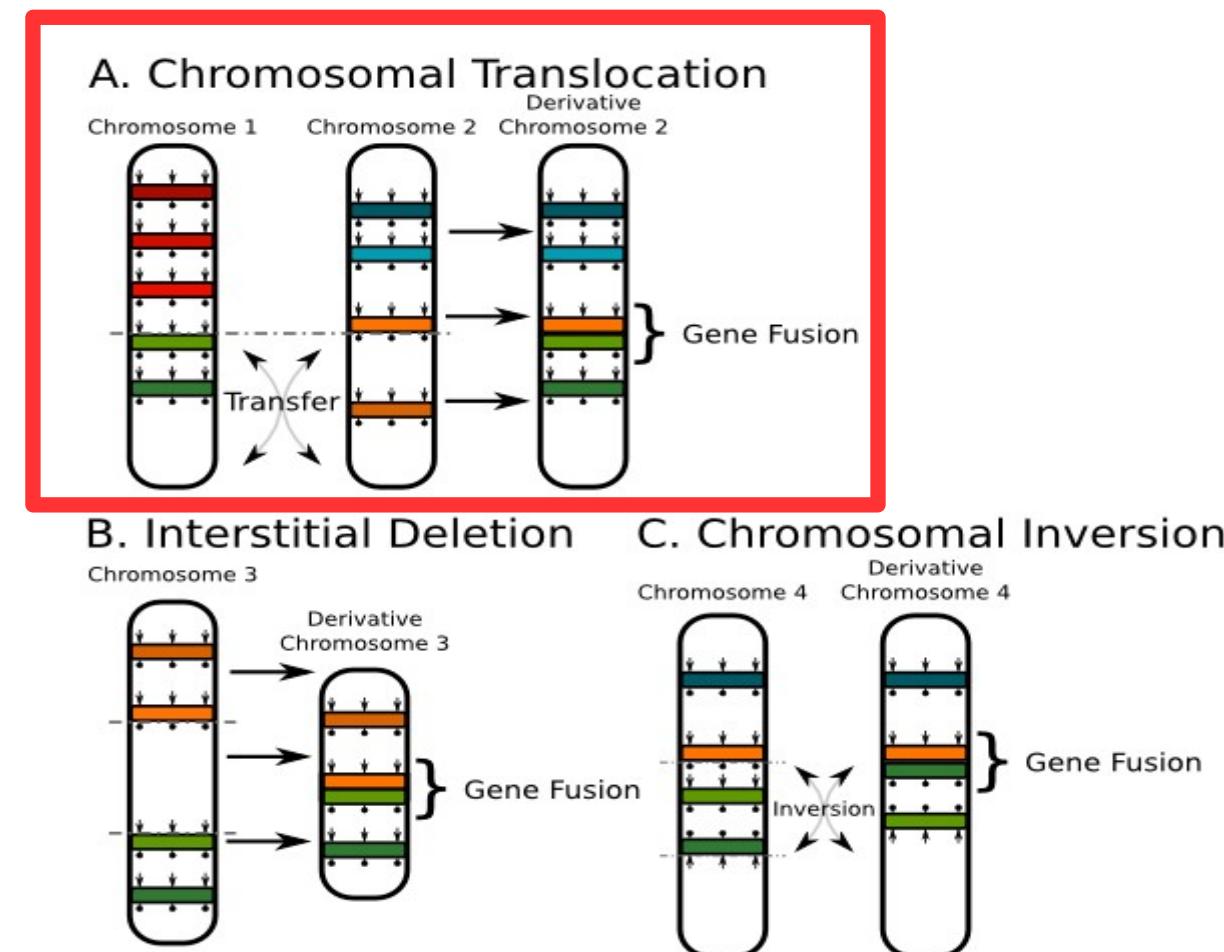
## Fastq File storing mates 2

```
>HWI-1KL152:175:C3D02ACXX:1:1101:1206:2168/2
ATTTGTGACATTGCTTATCCCGTAGAGGGGGATATGATAGTCGCCGAGCGTATGCACACGAACGCCAAATATGGTGTGAAGGTTGGCCTGACAAA
@HWI-1KL152:175:C3D02ACXX:1:1101:1302:2077/2
TAACATCTCACGCTCAGGGAACTACGCCCTGCTAACATCTCCCTGGCCCATCATCCTAGTCCTATGCCCTCCATCCCTACGCAAGATCGAAGCG
>HWI-1KL152:175:C3D02ACXX:1:1101:1375:2092/2
GCAAATAACCCCCCTAATAAAATTAACTCATTCACGCCCTCCACCCCATCCACATCTCGCATGTAACACTCGCTACTCTTGGC
>HWI-1KL152:175:C3D02ACXX:1:1101:1422:2104/2
GCTCTGGAGGCAGGATGGCCAGGCTATGGGTATGGGACTCAACGAAAGGAAACACCTTACACGCTAGATGGTGGGACATCATCACGCCCTGTGCTT
>HWI-1KL152:175:C3D02ACXX:1:1101:1358:2111/2
CCAGTGGGTTATCAGGGGACCAAGGGGGCTCAGGCTTACAGTGAATGGCTCCCTACTCACCCACTTATCCCTCCATGTAATTCAAGGACAA
```

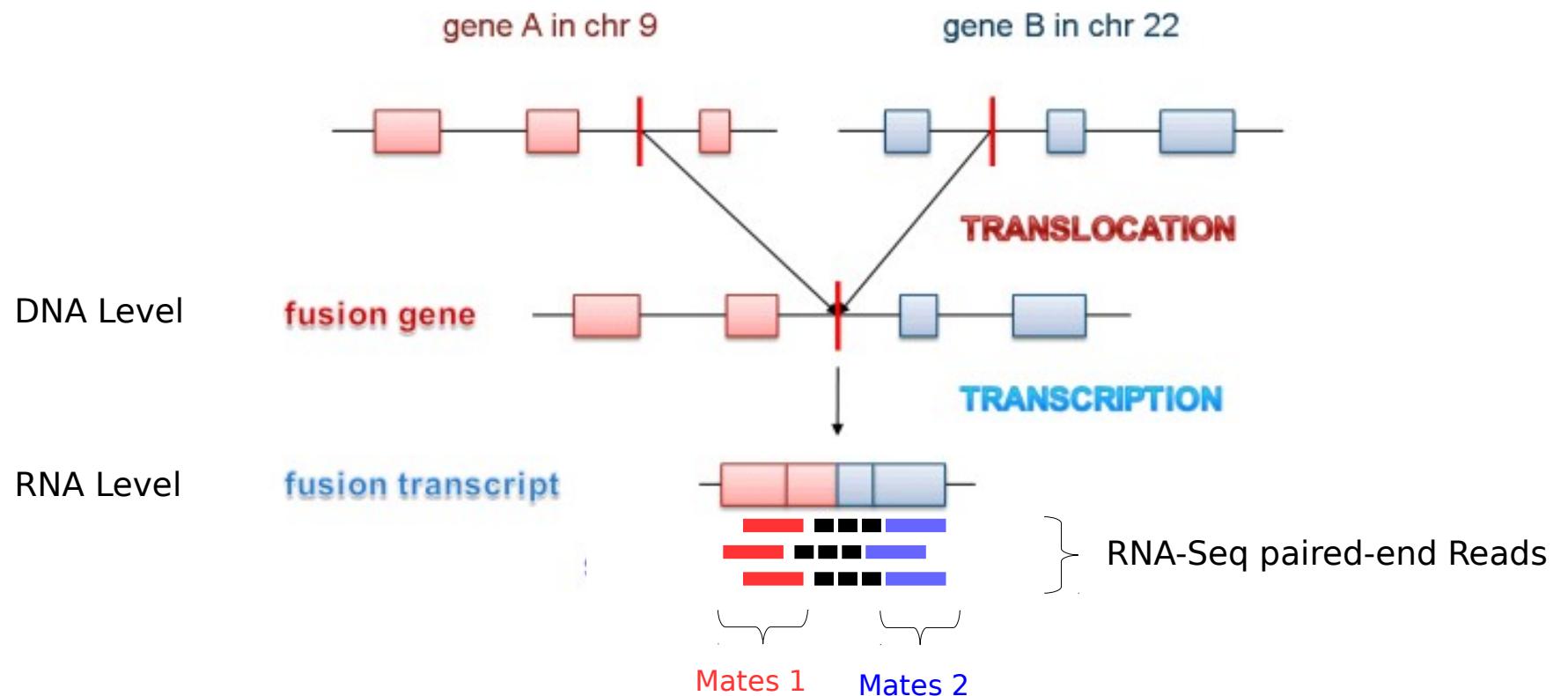
## Fasta File storing mates 2

# REMINDER (3)

- **FUSION GENE:** A fusion gene is a hybrid gene formed by two previously separated genes. At the DNA level these genes can originate from three different mechanisms.



# REMINDER (4)



# CHIMERASCAN WORKFLOW (1)

- CREATE ALIGNMENT INDEX

Chimerascan indexer uses the bowtie-build indexing scheme from Bowtie to build the index of the reference genome and transcriptome on which reads will be mapped.

- PREPARE READS FOR ALIGNMENT:

1) All reads quality scores are converted from Phred to Sanger format

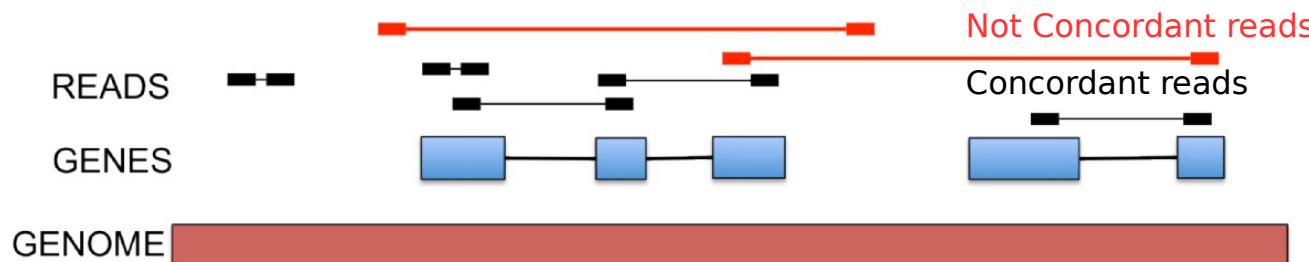
$$\text{Phred} = -10 \log_{10} P \quad \text{where } P \text{ is the base-calling error probability.}$$

$$\text{Sanger} = \text{Phred} + 22$$

2) All reads identifier are converted from an arbitrarily long string to a number.

- ALIGN PAIRED-END READS:

Paired-end reads are aligned to the genome and the transcriptome reference taking advantage of Bowtie aligner. Both the mates belonging to a pair-end read must align on the reference within a distance fixed by the user (Concordant reads).



Concordant reads: Align as a pair on the genome (within a fixed distance) or on the transcriptome within a single gene.

Not Concordant reads: Span multiple genes or large genomic distance.

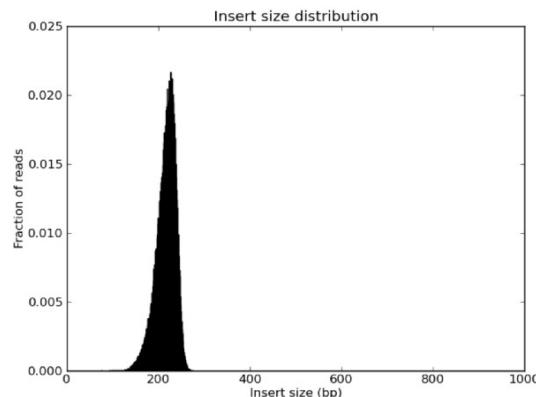
# CHIMERASCAN WORKFLOW (2)

- CREATE A SORTED/INDEXED BAM FILE

Alignments of concordant reads, originally stored in a SAM file, are sorted and compressed to produce an indexed BAM file.

- ESTIMATE FRAGMENT SIZE DISTRIBUTION

Aligned reads are parsed in order to compute the empirical distribution of inner distances. Only uniquely mapped reads are used to sample this distribution.



- REALIGN INITIALLY UNMAPPED READS

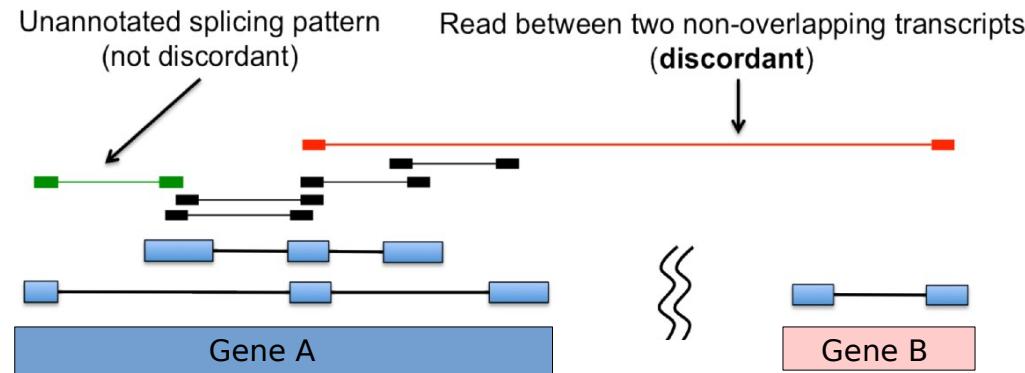
Reads that have not been aligned concordantly are trimmed into smaller segments and realigned on the reference genome and transcriptome.

# CHIMERASCAN WORKFLOW (3)

- DISCOVER DISCORDANT READS

The trimmed alignments are scanned for evidence of discordant reads. A discordant read is defined if:

- 1)The two mates do not align to the genome within a user-specified inner distance;
- 2)The pair does not align to a single transcript;
- 3)The two mates do not map on different transcripts of the same gene.



- NOMINATE CHIMERAS

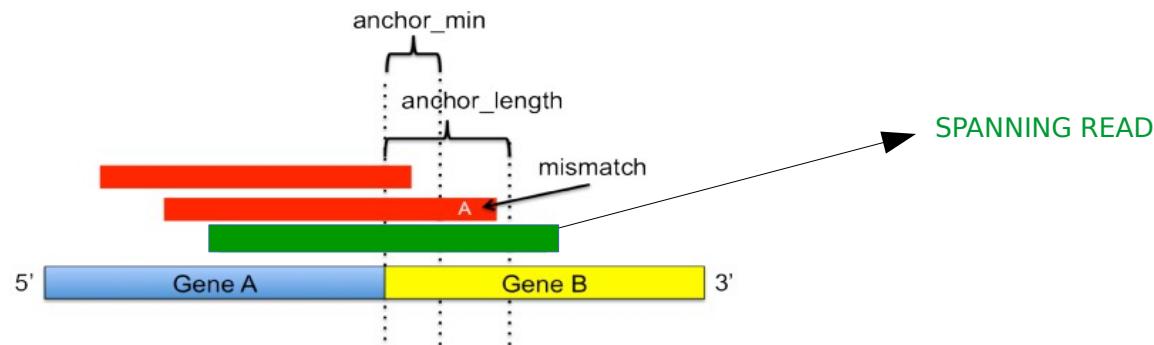
ChimeraScan clusters discordant reads producing a list of putative chimeric candidates (A). For each putative chimeric transcript ChimeraScan build a virtual reference by extracting from the genome its sequence (B).



# CHIMERASCAN WORKFLOW (4)

- CANDIDATE SPANNING READS REALIGNMENT AND SPANNING READS DEFINITION

Discordant reads with trimmed alignments bordering a junction or unmapped reads whose mates align to a predicted chimera are realigned on the obtained putative chimeric transcripts sequences. Spanning reads are defined accordingly to a minimum anchor length and a maximum number of mismatches in the anchor region.



- FILTERING CHIMERAS

In order to ensure high reliability of the detection:

- 1) Chimeras with very low coverage are filtered;
- 2) Chimeras with fragment sizes far outside the range of the distribution are filtered;
- 3) Known false positives are filtered.

- REPORTING CHIMERAS

ChimeraScan produces a tabular text file describing each chimera, and optionally generates a user-friendly HTML page with links to detailed descriptions of the chimeric genes.

# EXERCISE OBJECTIVES

- 1) Identify gene fusions in the sample dataset\_mate\*.fq using Chimerascan Tool;
- 2) Check on COSMIC database if the identified gene fusions have been previously reported:

<http://cancer.sanger.ac.uk/cancergenome/projects/cosmic/>

- 3) Reconstruct, starting from the so called spanning reads provided by Chimerascan tool for the detected gene fusions, the longest consensus region sequence for the fusion. This activity has to be performed by writing an ad-hoc python script.

# DATA (downloaded from the Course Material page)

## 1) The INDEX Folder containing:

- *align\_index.fa*: The reference file. It is in fasta format and stores some portions of sequences belonging to human chr 3, 11, 18 and 21.  
For computational reasons we previously produced all the index files needed for Chimerascan run

(*align\_index.1.ebwt,align\_index.2.ebwt,align\_index.3.ebwt,align\_index.4.ebwt,align\_index.fa.fai,align\_index.rev.1.ebwt,align\_index.rev.2.ebwt*)  
thanks to Chimerascan\_index.py.

## 2) The input files:

- *dataset\_mate1.fq* and *dataset\_mate2.fq*: These files contain respectively mate1 and mate2 sequences of the sample under investigation.

# CHIMERASCAN RUN

1)Create a folder directory to store Chimerascan output (chimera\_out dir)

2)Run chimerascan\_run.py program:

```
python /opt/chimerascan-0.4.5/chimerascan/chimerascan_run.py  
path_to_indexfolder/INDEX/ /path_to_datasetfolder/dataset_mate1.fq  
/path_to_datasetfolder/dataset_mate2.fq  
/path_to_outputdirectory/chimera_out
```

3) Try to change Chimerascan running parameters:

<https://code.google.com/p/chimerascan/wiki/Running>

# COSMIC DB SEARCH

Insert here Gene Name

The screenshot shows the COSMIC v72 homepage. At the top, there is a search bar with the placeholder text "eg: Braf, COLO-829, Carcinoma, V600E, BRCA-UK, Campbell" and a red "SEARCH" button. Above the search bar is a navigation bar with links for Home, About, Licensing, Data Download, News, Help, and Sign in. A large arrow points from the "Insert here Gene Name" text towards the search bar. To the right of the search bar is a circular "Genomic Landscape of Cancer" visualization, which is a multi-layered plot showing genomic data across chromosomes 1 through 22 and X. Below the search bar, there are two sections: "Resources" and "Tools". The "Resources" section lists the Cell Lines Project, COSMIC Whole Genomes, Cancer Gene Census, Drug Sensitivity, Mutational Signatures (marked as New), and GRCh37 Cancer Archive (also marked as New). The "Tools" section lists the Cancer Browser, Genome Browser, CONAN, and COSMIC Mart.

cancer.sanger.ac.uk/cosmic

Google

**COSMIC**  
Catalogue of somatic mutations in cancer

Home | About | Licensing | Data Download | News | Help | Sign in

COSMIC v72

eg: Braf, COLO-829, Carcinoma, V600E, BRCA-UK, Campbell

SEARCH

R Resources

Key COSMIC resources

- Cell Lines Project
- COSMIC Whole Genomes
- Cancer Gene Census
- Drug Sensitivity
- Mutational Signatures New
- GRCh37 Cancer Archive New

T Tools

Additional tools to explore COSMIC

- Cancer Browser
- Genome Browser
- CONAN
- COSMIC Mart

Genomic Landscape of Cancer

# CHIMERASCAN OUTPUT ELABORATION

Starting from `chimeras.bedpe` file write a program to elaborate the output in order to obtain from the spanning reads a breakpoint consensus sequence:

## BEFORE

```
GC GGAGGCGGGAGGGCGAGGGGCGGGGGAGCGCCGCCCTGGAGCGCGGCAG|GAAGC  
CGGAGGCGGGAGGGCGAGGGGCGGGGGAGCGCCGCCCTGGAGCGCGGCAG|GAAGCC  
CGGAGGCGGGAGGGCGAGGGGCGGGGGAGAGCCGCCCTGGAGCGCGGCAG|GAAGCC  
CGGCGGGCGGGAGGGCGAGGGGCGGGGGAGCGCCGCCCTGGAGCGCGGCCG|GAAGCC  
GGAGGCGGGAGGGCGAGGGGCGGGGGAGAGCCGCCCTGGAGCGCGGCAG|GAAGCCT  
GGAGGCGGGAGGGCGAGGGGCGGGGGAGAGCCGCCCTGGAGCGCGGCAG|GAAGCCT  
GGTGGCGGGAGGGCGAGGGGCGGGGGAGCGCCGCCCTGGAGCGCGGCAG|GAAGCCT  
GGAGGCGGGAGGGCGAGGGGCGGGGGAGCGCCGCCCTGGAGCGCGGCAT|GAAGCCT  
GGAGGCGGGAGGGCGAGGGGCGGGCGAGCGCCGCCCTGGAGCGCGGCAG|GGAGCCT  
GAGGCGGGAGGGCGAGGGGCGGGGGAGCGCCGCCCTGGAGCGCGGCAG|GAAGCCTT  
AGGCAGGGAGGGCGAGGGGCGGGGGAGCGCCGCCCTGGAGCGCGGCAG|GAAGCCTTA  
CGGGAGGGCGAGGGGCGGGGGAGCGCCGCCCTGGAGCGCGGCAG|GAAGCCTTATC  
CGGAGGGCGAGGGGCGGGGGAGCGCCGCCCTGGAGCGCGGCAG|GAAGCCTTATCA  
GGAGGGCGAGGGGCGGGGGAGCGCCGCCCTGGAGCGCGGCAG|GAAGCCTTATCAG  
AGGGCGAGGGGCGGGGGAGCGCCGCCCTGGAGCGCGGCAG|GAAGCCTTATCAGTT  
GGCGAGGGGCGGGGGAGCGCCGCCCTGGAGCGCGGCAG|GAAGCCTTATCAGGTC  
GCGAGGGGCGGGGGAGCGCCGCCCTGGAGCGCGGCAG|GAAGCCTTATCAGTTGTG  
GGCGGGGGAGCGCCGCCCTGGAGCGCGGCAG|GAAGCCTTATCAGTTGTGAGTGAG  
GGCGGGGGAGCGCCGCCCTGGAGCGCGGCAG|GAAGCCTTATCAGTTGTGAGTGAG
```

## AFTER

```
GC GGAGGCGGGAGGGCGAGGGGCGGGGGAGCGCCGCCCTGGAGCGCGGCAG|  
GAAGCCTTATCAGTTGTGAGTGAG
```