

Corso di Bioinformatica

http://compngen.bio.unipd.it/~stefania/Didattica/AA2011-2012/Bioinformatica_BTS/Bioinfo_BTS.html

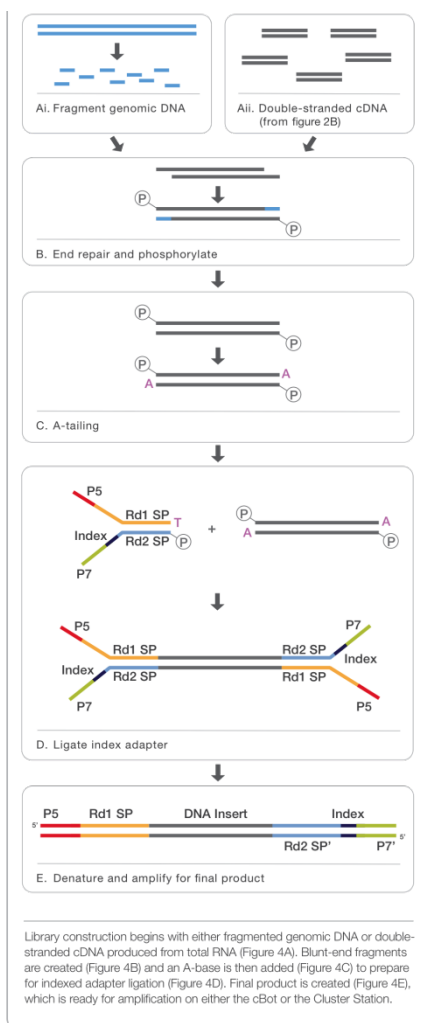
Per ogni punto dell'esercitazione copiare e incollare le informazioni richieste in un file sotto forma di breve testo con figure, in modo da assemblare un traccia di relazione con commenti sintetici.
Alla fine inviare via mail la relazione prodotta nel corso dell'esercitazione (stefania.bortoluzzi@unipd.it).
Salvare il file della relazione ogni 5'!

IV ESERCITAZIONE

Analisi del trascrittoma in organismi modello a partire da dati di deep sequencing: mappaggio delle reads al genoma ed analisi dei risultati per l'identificazione di geni espressi.

In questa esercitazione svolgeremo una parte delle analisi necessarie per identificare i geni espressi in un campione biologico di un organismo modello (genoma disponibile) a partire da dati di sequenziamento massivo del trascrittoma (RNA-seq) con tecnologia Illumina. In particolare, usando un web server molto potente (Galaxy), effettueremo il mappaggio delle reads di sequenza al DNA genomico con due programmi diversi ed integreremo i risultati con alcune track di annotazione del genoma mediante il sistema IGV (Integrative Genomics Viewer) per identificare esoni espressi in geni noti, e scoprire possibili nuovi geni.

Figura 1. Schema della costruzione di una libreria Paired End



1. Scaricare i dati necessari per l'esercitazione

Aprire il browser Firefox, all'indirizzo compngen.bio.unipd.it/~enrico scaricare il file .zip disponibile online, salvandolo in una cartella col vostro nome in "Documenti".

Ci vorranno circa 5 minuti.

Decomprimendo il file troverete una cartella contenente:

Raw_data:

SRR167673_chr10_1.fastq.zip

SRR167673_chr10_2.fastq.zip

Si tratta dei dati grezzi di sequenziamento RNA-seq Illumina con metodica "paired end a partire da RNA di tessuto adiposo di maiale. Nei due files (in quanto il sequenziamento è stato fatto "paired end") sono incluse reads in formato FASTQ.

Si tratta di circa 500.000 reads (solo una parte del dataset originale) di 90 basi ciascuna per file. Le reads appaiate distano 200 sulla sequenza originale.

I file sono zippati.

Mapping_results:

bowtie_sscg9.65_chr10_mapped_reads.bam

bowtie_sscg9.65_chr10_mapped_reads.bam.bai

tophat_sscg9.65_chr10_mapped_reads.bam

tophat_sscg9.65_chr10_mapped_reads.bam.bai

Risultati del mappaggio, da usare più tardi, o "di scorta".

2. Usare Galaxy project

<http://g2.bx.psu.edu/>

Data intensive biology *for everyone*.

Galaxy is an open, web-based platform for data intensive biomedical research. Whether on the free public server or your own instance, you can perform, reproduce, and share complete analyses.



Data Upload in Galaxy: scegliere “Get data” → “Upload data from your computer”

Scegliere il file da caricare (entrambi i fastq zippati, uno alla volta) e cliccare su “Execute”.

Attenzione!!! **Please do not use your browser's "stop" or "reload" buttons until the upload is complete, or it may be interrupted.**

Il sistema prima esegue l'upload, poi dezippa automaticamente i files. Quando appare il colore verde l'operazione e' completata. Cliccando “sull'occhio” potete vedere la parte iniziale del file.

Copiare ed incollare le prime due reads e commentare il formato.

Cliccando “sulla matitina” è possibile editare le proprietà del file.

Cambiare il formato scegliendo “fastqsanger” per entrambi i files.

Mappaggio delle reads al genoma di maiale usando Bowtie.

Sulla sinistra, scrivere “bowtie”: il sistema fornisce diverse opzioni, scegliere “Map with Bowtie for Illumina”.

Compilare il form per lanciare il programma di mappaggio.

Scegliere come genoma di riferimento “Pig (Sus scrofa): Sscrofa9.58 (SGSC)” (dato che e' disponibile non e' necessario fornire un indice della sequenza genomica).

Specificare che la libreria e' “mate paired” e usare il file #1 come forward e il #2 come reverse.

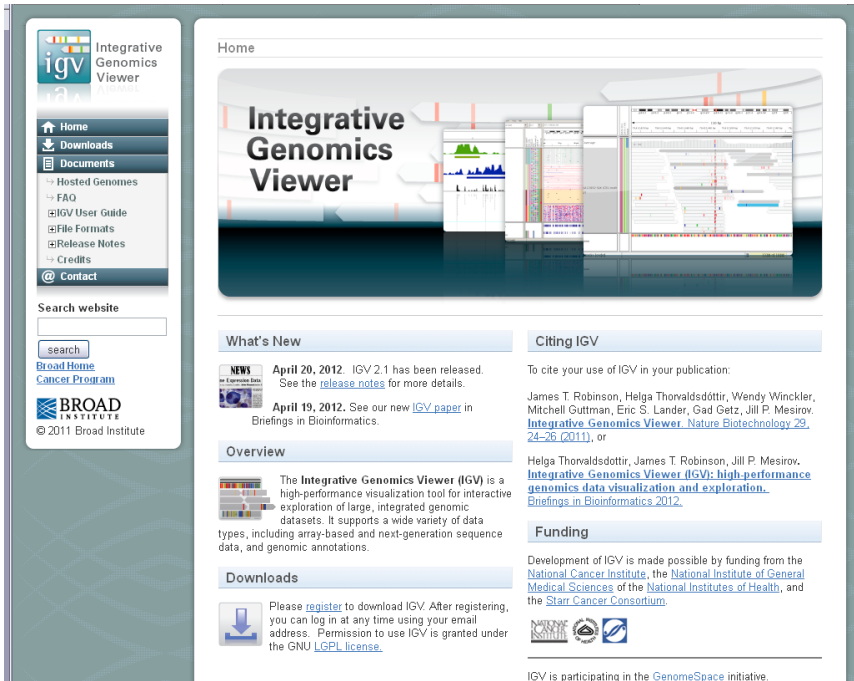
Cliccare su “Execute”.

Il mappaggio e' un processo computazionalmente intensivo, data la numerosità delle reads e la dimensione del genoma. Il programma Bowtie e' uno dei programmi di mappaggio più veloci, ma non permette la divisione delle reads (unspliced mapping)

The screenshot shows the Galaxy web interface. The main panel displays the 'Map with Bowtie for Illumina (version 1.1.2)' tool configuration. The 'Tools' sidebar on the left shows 'bowtie' selected. The configuration form includes options for 'Reference genome' (Pig (Sus scrofa): Sscrofa9.58 (SGSC)), 'Library type' (Paired-end), and 'FASTQ files' (2: SRR167673_chr10_2.fastq). The 'Execute' button is visible at the bottom of the form. The 'History' panel on the right shows the execution results, including the file 'SRR167673_chr10_2.fastq' and its contents.

3. Visualizzazione e “analisi” del mappaggio con IGV (Integrative genomics Viewer).

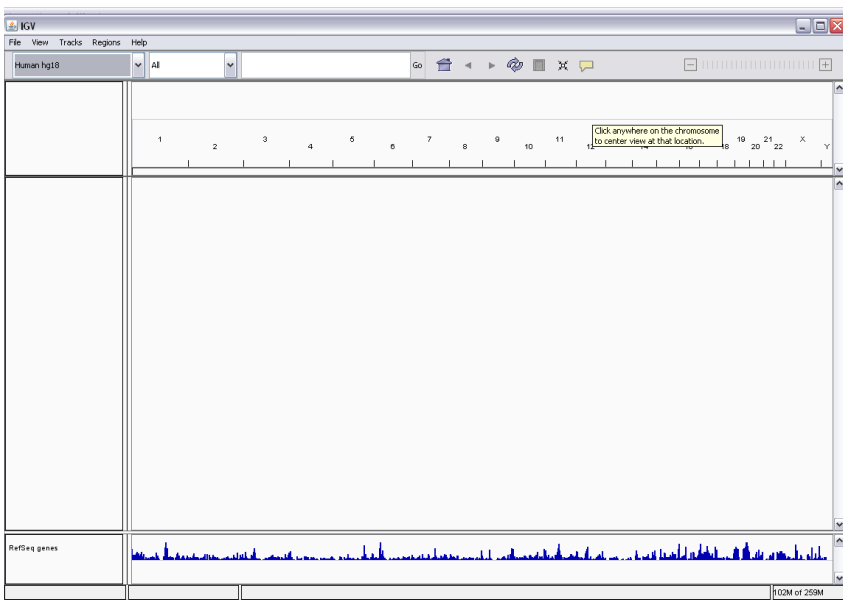
The Integrative Genomics Viewer (IGV) is a high-performance visualization tool for interactive exploration of large, integrated genomic datasets. It supports a wide variety of data types, including array-based and next-generation sequence data, and genomic annotations.



Sulla sinistra scegliere “Downloads” e usare l’indirizzo mail compngen.unipd@gmail.com già preparato.

Prima di tutto lanciare IGV usando JAVA Web Start, scegliere “Launch with 750 MB” (minima “potenza” permessa, ma la massima utilizzabile con i PC disponibili), aprire il file e dare “esegui”.

Quando si entra in IGV, si visualizza di default il genoma umano (24 cromosomi), con già caricata la track delle annotazioni corrispondente, quale la track “RefSeq” genes in basso. Dal menu a tendina è possibile scegliere un cromosoma oppure il genoma mitocondriale.



Nel menu in lato a sinistra scegliere il genoma di maiale “Sus scrofa 9.56”.

Quanti cromosomi vedete?

IGV ci servirà per visualizzare gli allineamenti ottenuti tra le reads RNA-seq, che in ultima analisi rappresentano frammenti di RNA trascritti in tessuto adiposo di maiale, sulla sequenza genomica, e in integrazione con alcune annotazioni della stessa, quali le posizioni e la struttura dei geni noti.

Espandere la track “Gene in

basso” usando il tasto destro.

Caricare il file degli allineamenti ottenuti con Bowtie, in formato BAM sorted: “Load from File”. Caricare anche il file `tophat_sscg9.65_chr10_mapped_reads.bam` presente nella cartella fornita.

Scegliere cromosoma 10, verrà visualizzata una schermata con l’intero cromosoma, con le reads non visibili. **Quante Mb è lungo il cromosoma 10?**

Aumentando lo zoom visualizziamo gli allineamenti e contemporaneamente vediamo i geni noti in maggiore dettaglio. Oltre un certo “ingrandimento” visualizziamo la sequenza genomica come una barra con i nucleotidi colorati.

Usando lo zoom o le coordinate visualizzare una finestra che comprenda circa 20 Kb.

Scorrere il cromosoma a partire dall’inizio (ptel; coordinate: 10:1-20,001).

Nella parte superiore vedrete due tracks per ciascun mappaggio:

- Coverage: diagramma a barre che indica il numero di sequenze allineate per nucleotide
- Allineamenti: dettaglio degli allineamenti reads-genoma con gli SNPs colorati.

Confrontare i risultati di mappaggio ottenuti con Bowtie e con Tophat.

Secondo il mappaggio con TopHat:

- trovare un gene espresso (coverage superiore a 10) e uno “non espresso”.
- Esistono regioni “introniche” che risultano espresse? Cercarne una lunga almeno 50 nt.
- Cercare un esone con almeno uno SNP e visualizzarlo insieme con la traduzione del messaggero nella regione intorno.