

## BREVE TUTORIAL SULL'USO DI PASS PER ALLINEAMENTI

### Informazioni generali sul programma:

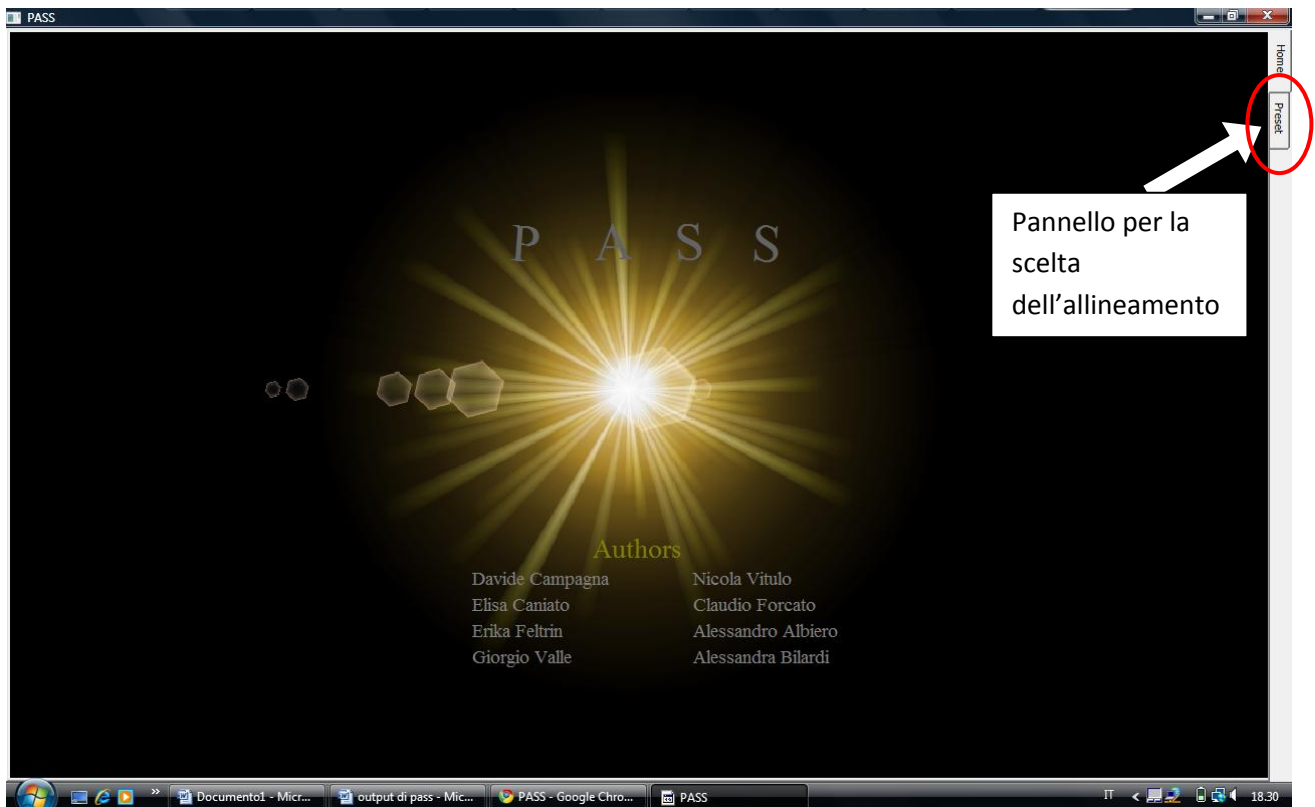
PASS è un programma che permette di allineare corte sequenze (reads) contro una sequenza reference iniziale in modo veloce e anche con introduzione di gap. Può essere usato con tutti i maggiori sistemi operativi e liberamente scaricato dalla rete. È in grado di maneggiare un elevatissimo numero di reads allineate poi alla reference e consente all'utente di modificare la sensibilità dell'allineamento stesso semplicemente settando i diversi parametri. Inoltre permette di utilizzare diversi formati per i dati in input e diverse modalità di visualizzazione per quelli in output.

I principali vantaggi di PASS sono:

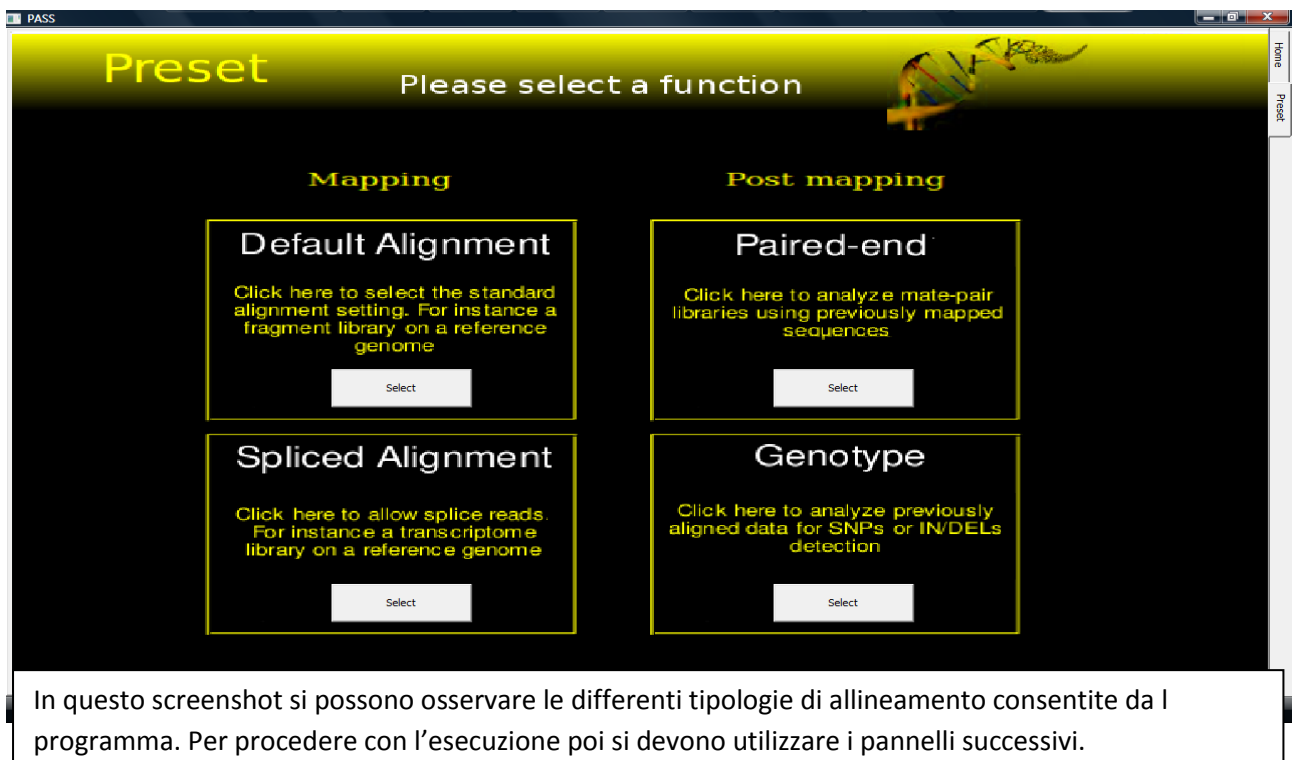
- la velocità nell'esecuzione dell'allineamento;
- i diversi parametri che permettono di modulare la sensibilità;
- la totale compatibilità con le sequenze ottenute con Solexa, Solid e 454;
- altre funzioni per allineamenti locali e supporti per allineamenti con long e spliced reads.

### Descrizione di PASS:

L'interfaccia grafica di PASS è di semplice utilizzo e consente di scegliere la tipologia di allineamento preferita con il pannello **Preset**.



Una volta stabilito quale allineamento utilizzare vengono inseriti dei parametri di default, mentre altri devono essere correttamente compilati in modo da definire i dati in input, le condizioni dell'allineamento e il formato dell'output. Questo viene fatto procedendo nei tre pannelli successivi: **Page1**, **Page2** e **Analysis**.

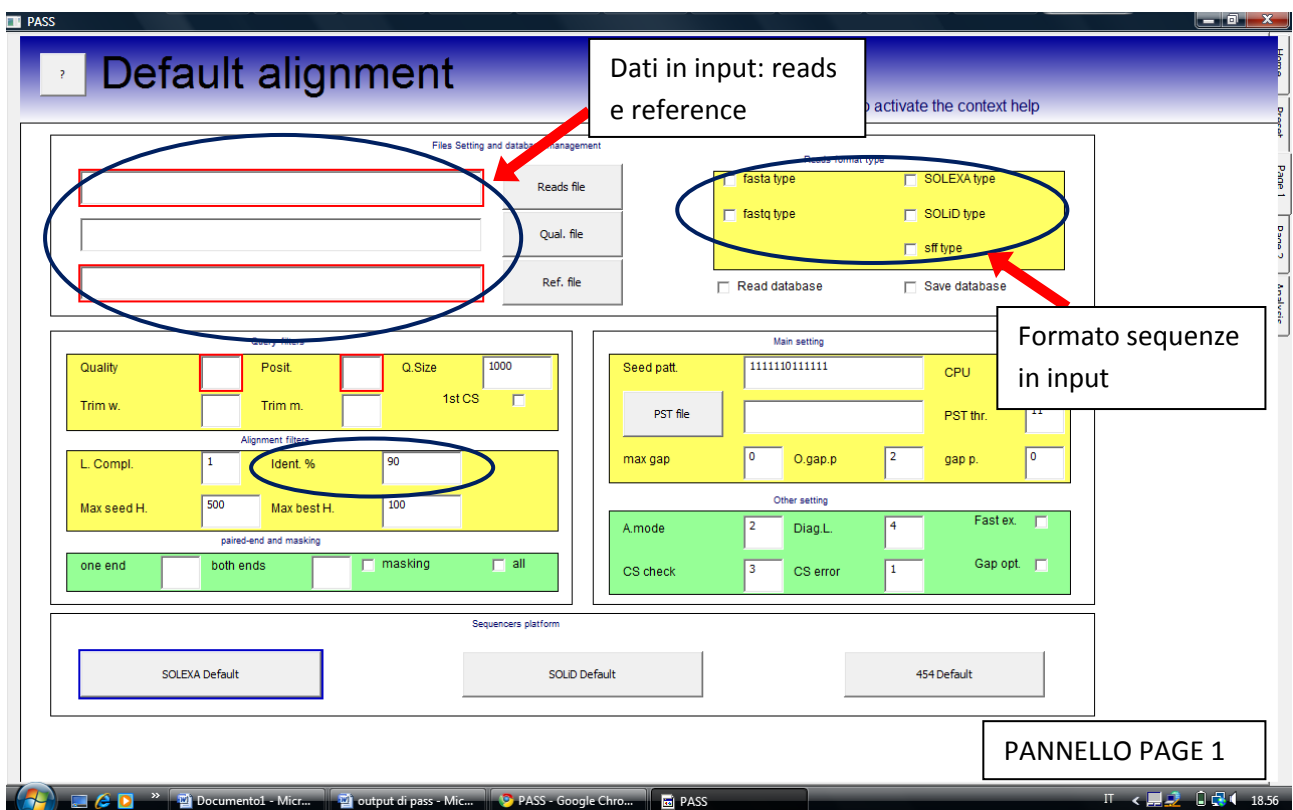


Per poter dunque meglio descrivere tutti i passaggi chiave di PASS si può considerare un esempio di allineamento:

allineamento del genoma di Drosophila contro un dataset di cromatogrammi con identità del 90%.

Si tratta di utilizzare un file in formato multifasta costruito con tutte le brevi sequenze dei cromatogrammi usati ad esercitazione, che costituiranno in questo caso le reads. Queste dunque dovranno essere allineate contro l'intero genoma, ma data la sua estensione nello specifico si è scelto di utilizzare la sola sequenza del cromosoma 2L, poiché era noto che le reads erano tutte state ottenute con sequenziamento del gene tim di Drosophila melanogaster situato in esso.

Perciò dal pannello **Preset** si seleziona l'allineamento "Default Alignment" e compare una finestra dove si notano dei parametri già settati di default e dei campi in rosso da compilare.



Ci si trova quindi nel pannello **Page1** dove per prima cosa si deve definire il formato dei dati in input. In questo caso si tratta del formato fasta perciò si seleziona la casella corrispondente "fasta type" sulla destra. Ora si devono compilare i campi "reads file" e "ref. file" inserendo rispettivamente il multifasta e la sequenza genomica detta reference. Sotto ad essi poi si va ad impostare la percentuale di identità voluta, in questo caso del 90%, inserendo il valore nel campo "Ident. %".

Compilati correttamente tutti i campi, si può passare al secondo pannello **Page2** dove si indica il nome e la posizione nel computer del file da salvare contenente i risultati calcolati. Questo va inserito nel campo "Save file".

Default alignment

Press on the question mark to activate the context help

Output file

Local alignment

Local align.  Align. len. 30 Ext. limit 3 Int. limit 3

Output options

Query block size 100000

GFF output  GFF info

Best hits  Unique

Repeat

sequence to g

aligned

Spliced alignment

Enable spliced alignment

splicing default

splic.fle max intr. 50000

splic.hits 10 i score 2 a prob 0.9

overlap 10 logo report

Sequencers platform

SOLEXA Default SOLID Default 454 Default

Indica dove sarà salvato il file dei risultati

Formato output

PANNELLO PAGE 2

In questo pannello viene anche indicato il formato con cui si visualizzano i risultati, in questo caso si mantiene la tipologia GFF (General Feature Format); in alternativa è possibile visualizzare i risultati con formato "full alignment" come BLAST.

Terminata la compilazione si può passare all'ultimo pannello **Analysis** dove si dà l'avvio all'allineamento impostato tramite il tasto "start" per poi poter visualizzare i risultati.

Default alignment

Press on the question mark to activate the context help

Summary and Statistics

<<< >>>

x 10

1/4

INDEXING

0

Start analysis

Start Stop

Virtual Memory Usage

Progressi dei risultati e grafico andamento

PANNELLO ANALYSIS

Il calcolo degli allineamenti fatto da PASS prevede che vengano allineate tutte le reads, ciascuna contro la reference. In questo modo si arriva ad ottenere i risultati in un documento leggibile tramite il programma Crimson Editor SVN286, dove compare la schermata contenente circa 15 colonne di numeri coma previsto dal formato GFF.

The screenshot shows the Crimson Editor interface with a GFF file open. The file content is as follows:

```

gi|113194944|gb|AE014134.5|pass|match|3503853|3504448|596|+|.ID=2058767:0:0;Name=A01_4.1;P="1-596";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503849|3504448|600|+|.ID=2058767:0:1;Name=A02_4.9;P="1-600";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503849|3504448|600|+|.ID=2058767:0:2;Name=A08_4.66;P="1-600";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503849|3504448|600|+|.ID=2058767:0:3;Name=B02_4.52;P="1-600";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503849|3504448|600|+|.ID=2058767:0:4;Name=B03_4.18;P="1-600";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503844|3504448|605|+|.ID=2058767:0:5;Name=A03_4.46;P="1-605";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503858|3504450|593|+|.ID=2058767:0:6;Name=A05_4.17;P="1-593";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503858|3504447|590|+|.ID=2058767:0:7;Name=A04_4.25;P="1-590";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503858|3504455|598|+|.ID=2058767:0:8;Name=A05_4.33;P="1-598";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503857|3504447|591|+|.ID=2058767:0:9;Name=A06_4.41;P="1-591";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503857|3504457|601|+|.ID=2058767:0:10;Name=A07_4.58;P="1-601";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503850|3504454|605|+|.ID=2058767:0:11;Name=A09_4.77;P="1-605";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503858|3504446|589|+|.ID=2058767:0:12;Name=B01_4.2;P="1-589";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3504282|3504811|530|+|.ID=2058767:0:13;Name=B02_3.10;P="530-1";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503849|3504451|603|+|.ID=2058767:0:14;Name=B02_4.10;P="1-603";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3504282|3504801|520|+|.ID=2058767:0:15;Name=B03_3.18;P="520-1";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503849|3504446|598|+|.ID=2058767:0:16;Name=B04_4.26;P="1-598";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503882|3504458|577|+|.ID=2058767:0:17;Name=B05_4.53;P="1-577";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3504282|3504808|527|+|.ID=2058767:0:18;Name=B06_3.46;P="527-1";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503853|3504446|593|+|.ID=2058767:0:19;Name=B06_4.42;P="1-594";Note="M:1 -> 211/211 T/G,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503853|3504447|595|+|.ID=2058767:0:20;Name=B07_4.59;P="1-595";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503858|3504450|602|+|.ID=2058767:0:21;Name=B08_4.67;P="1-602";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503858|3504451|594|+|.ID=2058767:0:22;Name=C01_4.3;P="1-594";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503853|3504446|594|+|.ID=2058767:0:23;Name=C02_4.11;P="1-594";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503853|3504446|594|+|.ID=2058767:0:24;Name=D04_4.28;P="1-594";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503853|3504446|594|+|.ID=2058767:0:25;Name=D06_4.44;P="1-594";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503853|3504447|595|+|.ID=2058767:1:5;Name=C03_3.20;P="520-1";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503853|3504447|595|+|.ID=2058767:1:6;Name=C07_3.66;P="520-1";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503853|3504447|595|+|.ID=2058767:1:7;Name=D01_3.4;P="520-1";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503853|3504447|595|+|.ID=2058767:1:8;Name=C04_4.27;P="1-599";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503853|3504447|595|+|.ID=2058767:1:9;Name=C05_3.37;P="535-1";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503853|3504447|595|+|.ID=2058767:1:10;Name=C05_4.35;P="1-590";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503853|3504447|595|+|.ID=2058767:1:11;Name=D01_4.4;P="1-590";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503853|3504447|595|+|.ID=2058767:1:12;Name=C07_4.60;P="1-598";Note="M:2 -> 206/206 T/G 595/595 N/T,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503856|3504447|592|+|.ID=2058767:1:13;Name=D02_4.12;P="1-592";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503858|3504446|589|+|.ID=2058767:1:14;Name=D03_4.20;P="1-589";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3504282|3504811|530|+|.ID=2058767:1:15;Name=D04_3.29;P="530-1";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3504282|3504811|530|+|.ID=2058767:1:16;Name=D05_3.38;P="530-1";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503849|3504445|597|+|.ID=2058767:1:17;Name=D04_4.34;P="1-597";Note="M:0,G:0";Hits=1;
gi|113194944|gb|AE014134.5|pass|match|3503849|3504462|614|+|.ID=2058767:1:18;Name=D04_4.70;P="1-614";Note="M:0,G:0";Hits=1;
  
```

Annotations in the image:

- reference**: Points to the second column (gb|AE014134.5).
- multifasta**: Points to the third column (pass|match).
- Posizione di inizio e fine regione allineata in reference**: Points to the fifth and sixth columns (3503853|3504448).
- Score allineamento**: Points to the seventh column (596).
- Name**: A blue circle highlights the eleventh column (Name=A01\_4.1).

I risultati visualizzati devono però essere interpretati. Dall'allineamento tra reads e genoma (cromosoma 2L) di Drosophila si osserva che: la prima colonna contiene il nome con cui è indicata la reference (sequenza genomica) e la seconda il nome con cui appaiono le sequenze brevi o reads (multifasta); la quinta e la sesta colonna indicano le posizioni di inizio e di fine dell'allineamento sulla sequenza reference perciò la regione sul cromosoma dove si trova la read corrispondente; la settima colonna riporta il punteggio dell'allineamento (abbastanza elevato per tutti i casi); la colonna otto definisce su quale filamento della reference si sono allineate le singole reads (5'→3' (+) o 3'→5'(-)); la colonna 11 contenente il campo "Name" indica il nome di ciascuna reads a cui di fatto vanno a corrispondere tutte le voci lungo la riga; la colonna successiva invece riporta le coordinate dei nucleotidi di inizio e fine dell'allineamento delle reads e si vede che tutte sono totalmente allineate; i due campi successivi indicano rispettivamente se sono presenti mismatch "M" specificando poi le posizioni dove si trovano e i gap "G" non trovati con queste reads; l'ultimo campo infine definisce gli "Hits" stabilendo dunque quante volte la reads considerata è presente e nel caso specifico si nota che tutti i valori danno 1 cioè tutte le sequenze non vengono ripetute.

Concludendo dunque si può affermare che tramite PASS è possibile fare degli allineamenti con brevi sequenze contro ad esempio interi genomi ottenendo risultati in tempi ristretti e allo stesso tempo accurati. Ciò significa che si è in grado di stabilire quanti e quali mismatch sono presenti, se sono stati inseriti gap e gli hits che si sono verificati.