

Using the command line tools

Tim Millar

2017-07-31

Contents

TEFingerprint	1
Extracting informative reads	2
Initial alignment	2
Extract informative reads	3
Fingerprinting one or more samples	4
Filtering GFF Output	6

TEFingerprint

TEFingerprint is a collection of related tools including:

- `tefingerprint`
- `tef-extract-informative`
- `tef-filter-gff`

The help text can be displayed for any of these tools using the `-h` or `--help` flags e.g.:

```
tefingerprint -h
```

Running the full TEFingerprint pipeline requires the following input data:

- paired-end reads in fastq format (2 files per sample).
- a reference genome in fasta format.
- a reference set of repeat elements in fasta format.
- optionally, an annotation of repeat elements within the reference genome in gff3 format.

The TEFingerprint pipeline requires some consistency with names of repeat elements. It is recommended that you use a naming convention similar to `<super-family>_<family>_<additional-info>` e.g.

Gypsy_26_chr15_18793972 or Gypsy_Gypsy26_chr15_18793972. The chosen naming convention should be used in both the reference set of repeat elements (a fasta file) and the reference annotation of repeat elements (a gff3 file in column 3, the “type” column).

When identifying clusters of selected categories of repeat elements, the start of repeat names from input files is matched to the specified categories of interest. For example if two repeats (transposons) are named **Gypsy_2_more-info** and **Gypsy_7_more-info** they will both be identified as being in the **Gypsy** category. However, they will be placed in separate categories if **Gypsy_2** and **Gypsy_7** are specified as categories. If a repeat name matches two or more categories it will be placed in the one with the longest name. For example a repeat named **Gypsy_27_more-info** will be categorised as **Gypsy_27** in preference of **Gypsy_2**.

Extracting informative reads

Identifying informative reads requires a paired-end reads in fastq format, a library of transposable elements in fasta format and a reference genome in fasta format. Fasta files should be pre indexed using BWA. This pipeline also requires that BWA and Samtools are installed and on the users \$PATH. The following describes the process to extract informative reads from a single sample and will need to be repeated per sample.

Paired end reads are initially aligned to the library of transposable elements manually with bwa mem. The **tef-extract-informative** tool is then used to identify and extract unmapped reads with mapped mates (informative reads) which are mapped against the reference genome and tagged with the ID of their mates transposable element.

The following file names are used in this example:

- **reads1.fastq** and **reads2.fastq**: paired-end reads.
- **reference.fasta**: reference genome.
- **repeats.fasta**: reference set of repeat elements.

Initial alignment

Create index fasta files with bwa:

```
bwa index -a bwtsw reference.fasta
bwa index -a is repeats.fasta
```

Map paired end reads to repeat elements using **bwa mem** and convert to a sorted bam file with **samtools**:

```
bwa mem -t 8 \
    repeats.fasta \
```

```
reads1.fastq \
reads2.fastq \
| samtools view -Su - \
| samtools sort - \
-o reads_mapped_to_repeats.bam
```

Index the resulting bam file with `samtools`:

```
samtools index reads_mapped_to_repeats.bam
```

Extract informative reads

Runing the `tef-extract-informative` program on the indexed bam file:

```
tef-extract-informative \
  reads_mapped_to_repeats.bam \
  --reference reference.fasta \
  --output informative.bam \
  --threads 8
```

Arguments:

- Input bam file
- `-r/--reference` reference genome to align informative reads to.
- `-o/--output` the name of the output (indexed) bam file.
- `-t/--threads` specifies how many threads to use for alignment in bwa (python components of the pipeline are single threaded).

Additional arguments:

- `--exclude-tails` This option will exclude the soft-clipped tails of reads mapped to transposons. This option may increase the number of full length reads used in the analysis which are only included if no soft-clips are used from their mate read.
- `--exclude-tips` This option will exclude the soft-clipped tips of reads mapped to transposons. This option may increase the number of full length reads used in the analysis which are only included if no soft-clips are used from their mate read.
- `--soft-clip-minimum-length` the minimum length allowed for soft-clipped tips or tails to be included (defaults to 38). Exclusion of a soft-clip may result in the inclusion of a (less precise) full-length informative read. A smaller length limit may increase the precision of results at the expense of accuracy and a larger value may improve accuracy at the loss of some precision and signal depth. See the method documentation for further details.
- `--exclude-full-length-reads` exclude full-length informative reads. This option can improve cluster precision by only using soft-clipped tips and tails. However this will result in much lower cluster depth and should

only be used with high read depths and high quality reference transposon sequences.

- `--mate-element-tag` by default, the sam-tag used to store repeat element names associated with each read is **ME** (Mate Element). An alternative tag can be specified with this option.
- `--keep-temp-files` Temporary intermediate files will not be deleted.
- `--use-os-temp` by default, the intermediate files are written to the same location as the output and then removed when the pipeline is completed. This flag will instead write temporary files to a location provided by the operating system which may be faster on some systems. This flag will ignore the `--keep-temp-files` flag.

The output file **informative.bam** contains informative reads mapped to the reference genome. These reads use the same name as they had in the original fastq file with an additional suffix indicating their alignment when mapping to the reference set of transposons. The suffix begins with **:R1** or **:R2** indicating which read of the pair it is. Reads that originate from soft-clips have a further suffix indicating their alignment orientation against a transposon ('+' or '-') and the end of the read which they originate from ('3' or '5' prime ends). For example a read originating from a soft-clip might have a name with the suffix **:R1:+3**. Full length informative reads did not align to any of the transposons (their mate did) and therefore don't include the second part of the suffix. Each of these reads is also tagged with the repeat element that their pair was mapped to using the 'ME' tag in the bam file.

Fingerprinting one or more samples

The **tefingerprint** tool can be used to create a fingerprint of a single sample, or to create a comparative fingerprint of multiple samples (i.e. comparing the fingerprints of more than one sample)

Example usage for comparing two or more bam files:

```
tefingerprint informative1.bam informative2.bam ... \  
  -f family1 family2 ... \  
  -m 10 \  
  -e 350 \  
  --gff fingerprint.gff.gz \  
  --csv fingerprint.csv.gz
```

Or when specifying most common parameters:

```
tefingerprint informative1.bam informative2.bam ... \  
  -a annotation.gff \  
  -r chr chr2 chr3 ... \  
  -f family1 family2 ... \  

```

```

-m 10 \
-e 350 \
-b 25 \
-j 50 \
-n 3 \
-q 30 \
-t 4 \
--gff fingerprint.gff.gz \
--tsv fingerprint.tsv.gz

```

Where `informative1.bam ...` are the bam file(s) being fingerprinted, and `fingerprint.gff.gz` and `fingerprint.tsv.gz` are respectively the output in (block-gzip compressed) gff3 and tab delimited text formats. Output files will be uncompressed if the `.gz` suffix is not specified.

If block-gzipped, the gff output can be tabix indexed with `tabix -p gff fingerprint.gff.gz` and the tab delimited text output can be tabix indexed with `tabix -s 1 -b 4 -e 5 -S 1 fingerprint.tsv.gz`.

Arguments:

- A single bam file to be fingerprinted or multiple bam files for a comparative fingerprint.
- **-a/--annotation-of-known-elements** An optional annotation of known repeat elements in gff (3) format for matching to identified insertions. Known elements are also used for joining pairs of clusters either side of an insertion. Known elements are also used for joining pairs of clusters either side of an insertion. This gff file may be compressed with gzip or bz2.
- **-r/--references** The reference sequence(s) (e.g. chromosomes) to be fingerprinted. If left blank (None) all references sequences in the input file will be used. *Default = None.*
- **-f/--families** Repeat element/transposon categories to be used. These must be exact string match's to start of read name and are used to split reads into categories for analysis. Not specifying at least one valid category will result in empty output files. *Default = None.*
- **-m/--minimum-reads** Minimum number of read tips required to form a cluster. *Default = 10.*
- **-e/--epsilon** The maximum allowable distance among a set of read tips required to form a cluster. This should be approximately equal to the insert size of paired reads. *Default = 250.*
- **-s/--splitting-method** Method used for splitting proximate clusters. One of "IDBCAN", "SIDBCAN" or "IDBCAN-aggressive" (deprecated). See the full documentation for details. *Default = "SIDBCAN".*
- **-b/--buffer-fingerprints** Additional buffer to be added to margins of fingerprints. This is used avoid identifying small clusters as unique, when there is only slight miss-match in read positions across samples (i.e. false positives). It also improves the robustness of comparisons by allowing more reads to be included in each bin. The buffer is trimmed back to the

extent of the furthestmost read tips it contains. *Default = 25*.

- **-j/--join-distance** Used to try and match clusters of informative reads to a known repeat-element (if provided) as well as joining pairs of clusters at either end of a repeat insertion. This represents the maximum distance to search for a known repeat from the end of each cluster. If no known repeat is present (or none are provided) then clusters will be paired if they are within twice this distance of one another. *Default = 25*.
- **-n/--number-of-common-elements** The number of most common repeat elements contributing to each cluster that are counted. *Default = 3*.
- **-q/--mapping-quality** Minimum allowed mapping quality for informative reads mapped to the reference genome. *Default = 30*.
- **-t/--threads** Maximum number of cpu threads to be used. The maximum number of threads that can be utilised is the number of reference molecules to be fingerprinted. *Default = 1*.
- **--gff** File name for GFF output. Compression will be applied by extension e.g. ".gz" (block-gzip) or ".bz2". Output may be written to standard output using "-". *Default = None*.
- **--csv** File name for CSV output. Compression will be applied by extension e.g. ".gz" (block-gzip) or ".bz2". Output may be written to standard output using "-". *Default = None*.
- **--tsv** File name for tab delimited text output. Compression will be applied by extension e.g. ".gz" (block-gzip) or ".bz2". Output may be written to standard output using "-". *Default = None*.

Additional arguments:

- **--minimum-epsilon** Minimum epsilon values used when calculating support for clusters. This is only used in hierarchical clustering and should usually be left as 0. *Default = 0*. This option may be deprecated in a future release.
- **--mate-element-tag** Sam-tag used in bam file to indicate the repeat element matched to each the mate read. *Default = "ME"*.
- **--no-max-count-proportion** Switch to disable calculation and colour coding of gff output based on max count proportion. This metric is useful for identifying potentially unique insertions but disabling it may improve performance.

Filtering GFF Output

The **tef-filter-gff** script can be used to filter down the gff formatted results of **tefingerprint**. Filters can be applied to attributes in the attribute column or to the first 8 standard gff3 columns. The first 8 standard gff3 columns are respectively named "seqid", "source", "type", "start", "end", "score", "strand" and "phase".

Filters take the form '<column/attribute><operator><value>' where:

- <column/attribute> is the name of the column or attribute that the filter is applied to.
- <operator> is one of the following operators =, ==, !=, < >, >=, <= that describes the comparison being performed.
- <value> is the value the each feature is compared to.

Filters should be contained within quotes ' ' so that the operator is not interpreted as a shell command.

The following operators are only used for numerical comparisons: < >, >=, <=.

The operators =, == and != will try to compare values as numerical (floating points) but will also check for equivalence or non-equivalence of string values. Note that =, == are identical.

Multiple filters may be combined within an “all” or “any” context. I.e. in an “all” context each feature must match all of the filters to be kept. In an “any” context each feature must only one of the filters to be kept. If both an “all” and an “any” context are used then they are evaluated separately before being combined in an additional “all” context.

Unix style wildcards may be used and will expand to match all possible column and attribute fields that they can. The resulting set of filters will then be evaluated within the context of the original filter.

Example usage with one column filter and two attribute filters:

```
tcf-filter-gff fingerprint.gff.gz \
  --all 'seqid=chr1' 'start>=1000' 'stop<9000' \
  --any 'sample_?_count>100' \
  -o fingerprint_filtered.gff.gz
```

Where `fingerprint.gff.gz` is a gff file compressed with gzip and `fingerprint_filtered.gff.gz` is a filtered version of that file.

The above example is evaluated as follows: the “all” context will select only feature from chromosome 1 that are in the interval 1000-8999. The “any” context contains a filter with the wildcard “?” which will expand the filter to match multiple samples and evaluate each of the resulting filters e.g.: with three samples it would expand to the equivalent of `--any 'sample_0_count>100' 'sample_1_count>100' 'sample_2_count>100'`. Therefore the full command would select features where any one of the samples contains more than 100 reads, from within the interval chr:1000-8999.

Arguments:

- A gff (3) file to be filtered. If - is specified then gff lines will be read (in plain text) from standard in.
- `--all` filters to apply to apply in the “all” context. These should take the form '<column><operator><value>'

- `--any` filters to apply to apply in the “any” context. These should take the form '`<column><operator><value>`'
- `-o/--output` a file to write the data to. By default - the data are written to standard out. If the data are written to a file, this file will be compressed with gzip or bz2 based on the files extension e.g. `.gff.gz` or `.gff.bz2`.