

Package ‘DBSverify’

August 6, 2021

Title Analyze Next Generation Sequencing Data to Verify Somatic Double-BaseSubstituion (DBS) Calls

Version 0.10

Description DBSverify stands for “Doublet Base Substitution verify”.

A “Doublet Base Substitution” is the simultaneous mutation of adjacent bases, for example mutations from TA to GA and from TC to AA, which are often caused by the carcinogen aristolochic acid. The main function is `Read_DBS_VCF_and_BAMs_to_verify_DBSs`. This function reads a VCF (variant call format) file containing somatic DBS mutations and examines the supporting reads in the corresponding tumor and normal BAM files to assess whether each DBS is likely real, i.e. arose at one time on a single chromosome. This is as opposed to being two adjacent single base mutations on homologous chromosomes, or adjacent mutations that occurred on the same chromosome but at detectably different times. Many somatic mutation callers do not call DBSs at all, and those that do sometimes make obvious errors, such as calling DBSs that consist of a germline SNP next to a somatic single base mutation. Because of the way the functionality of this package is used, the main function for most users will be `Read_DBS_VCF_and_BAMs_to_verify_DBSs`.

License GPL-3

Encoding UTF-8

LazyData true

Roxygen list(markdown = TRUE)

RoxygenNote 7.1.1

Imports stringr,
ICAMS (>= 3.0.0),
data.table

Suggests dplyr,
GenomicRanges,
spelling,
testthat,
tibble

Language en-US

R topics documented:

CategorizeReads	2
DBS_conclusion_1_row	2
GetAllBAMSlicesSamtools	3
PCAWG_read_table_and_evaluate_DBS	4
ReadSamfile	5
Read_DBS_VCF_and_BAMs_to_verify_DBSs	6
Read_SBS_VCF_and_BAMs_to_verify_DBSs	8
Slice2ReadSupport	9
Index	10

CategorizeReads	<i>Categorizes the reads in one sam file.</i>
-----------------	---

Description

Categorizes the reads in one sam file.

Usage

CategorizeReads(sam, POS, REF, ALT)

Arguments

- sam An in-memory version of a sam file as data.frame.
- POS The position in the sam at which the DBS starts.
- REF The reference allele.
- ALT The putative alternative allele.

Value

a named character vector, the names are the read names in the sam file. Each element is one of "Read supports only 1st position" "Read supports only 2nd position", "WT read" "Mut read"
@keywords internal

DBS_conclusion_1_row	<i>Examine 1 row of a DBS-only VCF that was already processed by Slice2ReadSupport to decide if the DBS is likely real.</i>
----------------------	---

Description

Examine 1 row of a DBS-only VCF that was already processed by [Slice2ReadSupport](#) to decide if the DBS is likely real.

Usage

DBS_conclusion_1_row(row, germlineCutOff = 0.1, max.half.support.T.reads = 1)

Arguments

<code>row</code>	One row of a DBS VCF (as a <code>data.frame</code>) already processed by Slice2ReadSupport (so that the fields <code>vcf\$NreadSupport</code> and <code>vcf\$TreadSupport</code> are populated).
<code>germlineCutOff</code>	If this proportion of normal reads show one or the other variant (or both variants), consider this a germline variant or partial germline variant.
<code>max.half.support.T.reads</code>	Do not tolerate more than this number of reads in the tumor that support one but not both mutated positions.

Value

A character string indicating the conclusion about the putative DBS.

`GetAllBAMSlicesSamtools`

Get and save all slices of BAM as specified by positions in a VCF table.

Description

Get and save all slices of BAM as specified by positions in a VCF table.

Usage

```
GetAllBAMSlicesSamtools(
  vcf,
  bam.name,
  padding = 10,
  where.to.put.slices = tempfile(),
  verbose = 0
)
```

Arguments

<code>vcf</code>	A VCF (Variant Call Format) "file" as a <code>data.table</code> or similar.
<code>bam.name</code>	The file name (path) to the BAM file to slice.
<code>padding</code>	How many base pairs on each side of the first base of the DBS to to keep in the BAM slices.
<code>where.to.put.slices</code>	If <code>NULL</code> , create a temporary directory to store the slices. Otherwise, a character string that specifies a directory in which to store the BAM slices. This is directory is created if necessary.
<code>verbose</code>	If <code>> 0</code> print a message when starting the number of slices generated every verbose slices.

Value

A character string that specifies the directory the containing the BAM slices. The slices are stored as SAM files.

PCAWG_read_table_and_evaluate_DBS

Read a table that specifies how to process VCFs and mini-BAMS to evaluate DBS calls, for PCAWG Collaboratory data.

Description

Read a table that specifies how to process VCFs and mini-BAMS to evaluate DBS calls, for PCAWG Collaboratory data.

Usage

```
PCAWG_read_table_and_evaluate_DBS (
    in.table,
    in.vcf.dir,
    minibam.dir,
    out.vcf.dir = in.vcf.dir,
    bam.suffix = "_dbs_srt",
    verbose = 1
)
```

Arguments

<code>in.table</code>	the file path of the table to process; in production, <code>.../DBSverify/data-raw/collaboratory_bams_full_20</code> for testing <code>.../DBSverify/data-raw/short_collaboratory_bams.csv</code> .
<code>in.vcf.dir</code>	The path to the directory containing the DBS VCF files.
<code>minibam.dir</code>	The path to the directory containing the mini BAMs.
<code>out.vcf.dir</code>	The path to the directory in which to put the "evaluated" DBS VCF files.
<code>bam.suffix</code>	String to add to end of BAM file name; depends on the conventions used by the script (run on the Collaboratory) that generated the miniBAMs.
<code>verbose</code>	If <code>> 0</code> generate some progress messages.

Details

This is a specialized function for processing PCAWG data from the ICGC (International Cancer Genome) "Collaboratory" cloud computing system, once the miniBAMs have been created in the Collaboratory and downloaded. The `in.table` and associated BED files were used to specify the contents of the miniBAMs. The result consists of the "evaluated" DBS VCF files. The naming of the input and output VCF files and the mini BAMs is governed by the contents of `in.table`, with the VCF file names incorporating the "aliquot_id" and the miniBAM names based on the `icgc_donor_id` and the T_Specimen ID and N_Specimen ID.

ReadSamfile	<i>Read a "SAM file", discarding some reads that cannot be interpreted for our purposes.</i>
-------------	--

Description

Read a "SAM file", discarding some reads that cannot be interpreted for our purposes.

Usage

```
ReadSamfile(filename)
```

Arguments

filename	The name of the SAM file to read.
----------	-----------------------------------

Details

SAM stands for "Sequence Alignment Map", a text file that represents aligned next-generation sequencing reads ([https://en.wikipedia.org/wiki/SAM_\(file_format\)](https://en.wikipedia.org/wiki/SAM_(file_format))). Only keep reads that satisfy certain conditions:

- Mate pair maps to same chromosome
- Mapping quality ≥ 30
- "FLAG" < 256 (see info on the return value element `reads.with.bad.FLAG`, below).
- The CIGAR string is only `\d+M` (one or more digits followed M, and nothing else before or after). This means there are no insertions or deletions in the read versus the reference and there is no soft clipping. Insertions or deletions or clipping shift the read within the SAM file, after which this function cannot keep track of the DBS location.

Value

A list with the elements:

- `good.reads`, `data.frame` with with column names for the first 11 columns as specified in [https://en.wikipedia.org/wiki/SAM_\(file_format\)](https://en.wikipedia.org/wiki/SAM_(file_format)) with one row per read. The result contains only the reads that meet certain conditions.
- `total.depth`, the initial depth including "bad" reads
- `reads.with.bad.FLAG`. These are reads with FLAG ≥ 256 , which marks reads that (i) are "not primary alignment" (ii) failed vendor QC (iii) are PCR or optical duplicates (iv) are supplementary alignments (e.g. split, split /inverted read). See <https://broadinstitute.github.io/picard/explain-flags.html>
- `reads.with.bad.CIGAR`, reads with CIGAR string that indicates an indel in the read or clipping.
- `reads.with.bad.MAPQ`, reads with MAPQ < 30 .
- `reads.with.bad.Mate_CHROM`, reads with a mate on a different chromosome.

Read_DBS_VCF_and_BAMs_to_verify_DBSs

Determine whether sequencing reads in fact support (candidate) DBSs present in a VCF file.

Description

Determine whether sequencing reads in fact support (candidate) DBSs present in a VCF file.

Usage

```
Read_DBS_VCF_and_BAMs_to_verify_DBSs (
  input.vcf,
  Nbam.name,
  Tbam.name,
  N.slice.dir = tempfile(),
  T.slice.dir = tempfile(),
  unlink.slice.dir = TRUE,
  exclude.SBSs = TRUE,
  verbose = 0,
  outfile = NULL,
  filter.status = "PASS"
)
```

Arguments

<code>input.vcf</code>	If a character string, then the path to a VCF file; otherwise A a single VCF "file" as a data.frame or similar object.
<code>Nbam.name</code>	The name of the BAM file for the normal sample corresponding to <code>vcf.name</code> .
<code>Tbam.name</code>	The name of the BAM file for the tumor sample corresponding to <code>vcf.name</code> .
<code>N.slice.dir</code>	Directory for the slices of the normal BAM. Created if necessary.
<code>T.slice.dir</code>	Directory for the slices of the tumor BAM. Created if necessary. Must be different than <code>N.slice.dir</code> .
<code>unlink.slice.dir</code>	If TRUE unlink <code>N.slice.dir</code> and <code>T.slice.dir</code> before return.
<code>exclude.SBSs</code>	If TRUE silently filter out (exclude) SBSs in the input VCF. This makes sense if the the VCF is from a caller (like Mutect or the Hartwig Medical Foundation caller) that calls both SBSs and DBS.
<code>verbose</code>	If > 0 print a message when starting the number of slices generated every verbose slices.
<code>outfile</code>	If not NULL then write the "evaluated" VCF to <code>outfile</code> ; otherwise write it to <code>paste0(input.vcf(vcf.name, "_evaluated.vcf"))</code> . Must be non-NULL if <code>input.vcf</code> is not a file path.
<code>filter.status</code>	If not NULL only keep rows where the FILTER column in the VCF is equal to <code>filter.status</code> .

Details

Creates a new VCF file. This VCF file has no data rows if there were no DBSs to analyze. Otherwise, this VCF contains some additional columns. Any SBSs or indels in the input are silently ignored, and no attempt is made to merge adjacent SBSs.

1. `NreadSupport` With regard to the two positions of the DBS in the normal BAM, a string with 4 numbers separated by ":", with the numbers indicating respectively:
 - the number of reads that are reference sequence at both positions of the DBS,
 - the number of reads that have the alternative allele only at the 1st position of the DBS,
 - the number of reads that have the alternative allele only at the second position of the DBS, and
 - the number of reads that have the alternative alleles at both positions of the DBS.
2. `TreadSupport` Information analogous to that in `NreadSupport`, for the tumor BAM.
3. `num_bad_mapped_reads` The total number of tumor reads with `MAPQ < 30` or with a mate on a different chromosome. If there are many badly mapped reads in the slice the slice may represent a segmental duplication.
4. `num_bad_mapped_DBS_reads` The number of tumor reads with the putative DBS but with `MAPQ < 30` or a mate on a different chromosome. If many badly mapped reads support DBSs the DBS might results from mismapped reads in a segmental duplication.
5. `DBSconclusion` A string that describes whether the DBSs is believable ("True DBS"), or if the DBS is not believable, a string that describes why not.

The decision in `DBSconclusion` is based on multiple criteria. I also suggest relying on any available upstream filtering of SBSs that get merged into DBSs, as well as upstream filtering of DBSs. It is difficult to capture all the possible characteristics of likely miscalled DBSs, especially if they stem from mismapped reads that nevertheless have high `MAPQ` (mapping quality). The code that implements these criteria is in `DBS_conclusion_1_row`, which depends on classification of individual reads in `ReadSamfile`. The filtering of reads in the input SAM files is described in the documentation for `ReadSamfile`.

Once the reads to analyze are selected, additional criteria include

- There must be ≥ 5 normal reads at the site of the putative tumor DBS.
- At each separate position of the DBS, the normal reads must have $< 10\%$ of the variant in the DBS at that position.
- < 2 normal reads support the DBS.
- At least 2 tumor reads support the DBS.
- There are more well-mapped tumor reads than badly mapped tumor reads at the site of the DBS.
- There are more well-mapped tumor reads than badly mapped tumor reads that contain the DBS. (`ReadSamfile` keeps track of the well-mapped and badly mapped reads).
- If 1 normal read supports the DBS, then there must be a statistically greater proportion of tumor reads supporting the DBS (by Fisher's test).
- There must be ≥ 5 normal reads at the site of the putative tumor DBS.
- At each separate position of the DBS, the normal reads must have $< 10\%$ of the variant in the DBS at that position.
- < 2 normal reads support the DBS.
- At least 2 tumor reads support the DBS.

- There are more well-mapped tumor reads than badly mapped tumor reads at the site of the DBS.
- There are more well-mapped tumor reads than badly mapped tumor reads that contain the DBS. (`ReadSamfile` keeps track of the well-mapped and badly mapped reads).
- If 1 normal read supports the DBS, then there must be a statistically greater proportion of tumor reads supporting the DBS (by Fisher's test).

Value

Invisibly, a list with the elements

1. The name of the DBS-only VCF file created.
2. The in-memory representation of the DBS VCF as a `data.table`.
3. The name of the directory with the normal SAM slices, if `unlink.slice.dir` is `FALSE`.
4. The name of the directory with the tumor SAM slices, if `unlink.slice.dir` is `FALSE`.

`Read_SBS_VCF_and_BAMs_to_verify_DBSs`

Aggressively merge SBSs into DBSs then determine whether sequencing reads in fact support DBSs.

Description

Aggressively merge SBSs into DBSs then determine whether sequencing reads in fact support DBSs.

Usage

```
Read_SBS_VCF_and_BAMs_to_verify_DBSs (
  input.vcf,
  Nbam.name,
  Tbam.name,
  variant.caller,
  N.slice.dir = tempfile(),
  T.slice.dir = tempfile(),
  unlink.slice.dir = TRUE,
  outfile = NULL,
  verbose = 0
)
```

Arguments

<code>input.vcf</code>	If a character string, then the path to a VCF file; otherwise A a single VCF "file" as a <code>data.frame</code> or similar object.
<code>Nbam.name</code>	The name of the BAM file for the normal sample corresponding to <code>vcf.name</code> .
<code>Tbam.name</code>	The name of the BAM file for the tumor sample corresponding to <code>vcf.name</code> .
<code>variant.caller</code>	One of "strelka", "PCAWG", or "unknown". Merging adjacent SBS is done by <code>ReadAndSplitVCFs</code> in the ICAMS package.

<code>N.slice.dir</code>	Directory for the slices of the normal BAM. Created if necessary.
<code>T.slice.dir</code>	Directory for the slices of the tumor BAM. Created if necessary. Must be different than <code>N.slice.dir</code> .
<code>unlink.slice.dir</code>	If TRUE unlink <code>N.slice.dir</code> and <code>T.slice.dir</code> before return.
<code>outfile</code>	If not NULL then write the "evaluated" VCF to <code>outfile</code> ; otherwise write it to <code>paste0(input.vcf(vcf.name, "_evaluated.vcf"))</code> . Must be non-NULL if <code>input.vcf</code> is not a file path.
<code>verbose</code>	If > 0 print a message when starting the number of slices generated every verbose slices.

Details

Note: argument `input.vcf` must be a file path. This function creates a new VCF file. See [Read_DBS_VCF_and_BAMs_to_verify_DBSs](#) for details.

Value

Same as [Read_DBS_VCF_and_BAMs_to_verify_DBSs](#).

<code>Slice2ReadSupport</code>	<i>Calculate the support for a putative DBS from the SAM slice containing the overlapping reads.</i>
--------------------------------	--

Description

Calculate the support for a putative DBS from the SAM slice containing the overlapping reads.

Usage

```
Slice2ReadSupport(slice.dir, CHROM, POS, REF, ALT)
```

Arguments

<code>slice.dir</code>	The directory containing the SAM slices.
<code>CHROM</code>	The chromosome identifier.
<code>POS</code>	The first position of the DBS.
<code>REF</code>	The reference variant.
<code>ALT</code>	The alternate variant.

Index

CategorizeReads, [2](#)

DBS_conclusion_1_row, [2](#), [7](#)

GetAllBAMSlicesSamtools, [3](#)

PCAWG_read_table_and_evaluate_DBS,
[4](#)

Read_DBS_VCF_and_BAMs_to_verify_DBSs,
[6](#), [9](#)

Read_SBS_VCF_and_BAMs_to_verify_DBSs,
[8](#)

ReadSamfile, [5](#), [7](#), [8](#)

Slice2ReadSupport, [2](#), [3](#), [9](#)