

Reference Manual: Sasquatch R-Tool



Version:	0.1
Date:	March 14, 2016
Authors:	Ron Schweßinger, Maria Suciú, Simon McGowan, Jelena Telenius & Jim Hughes
Groups:	Genome Biology and Computational Biology Research Group, WIMM, Oxford
Contact:	ron.schwessinger@ndcls.ox.co.uk, genmail@molbiol.ox.ac.uk
Description:	R implementation of Sasquatch, a tool for predicting the impact of sequence variation on DNase I footprinting potential.
Required Packages:	ggplot2 RColorBrewer grid TFBSTools Biostrings
Recommended:	pbapply
License:	...

Contents

CalcSFR	2
CompareSequences	3
DecodeKmer	5
DissectSequence	5
GetCount	6
GetFootprint	6
GetFootprintStrand	7
GetPossibleMutations	8
GetSFR	9
GrepProfile	10
InSilicoMutation	10
MakeRainbowTrackHub	12
PlotOverlap	14
PlotOverlapKmers	16
PlotSingle	17
PlotSingleKmer	18
PlotSingleStrands	19
PruneProfile	20
QueryJaspar	20
QueryJasparBatch	21

QueryLongSequence	22
RainbowPlot	23
RefVarBatch	24
SmoothProfile.	25
SobelBorders	26
Sobeln	27

CalcSFR

Description

Calculate the Shoulder-to-Footprint-Ratio of a given footprint profile and estimated shoulder positions and ranges.

Usage

```
CalcSFR(profile,  
         upstream.shoulder.position,  
         upstream.shoulder.position,  
         upstream.shoulder.range,  
         upstream.shoulder.range)
```

Arguments

profile	input cut profile (should be smoothed)
us.mid	center bp position of the upstream shoulder (as estimated from <i>SobelBorders</i>)
ds.mid	center bp position of the downstream shoulder
range.us	range in bp of the upstream shoulder (as estimated from <i>SobelBorders</i>)
range.ds	range in bp of the downstream shoulder

Value

Returns the SFR as single numeric value.

Examples

```
frag.type <- "DNase"  
  
data.dir <- "./my_sasq_database/human/DNase/"  
  
tissue <- "blood_tissue"  
  
kmer <- "CACGTG"  
  
# get the average profile for a single kmer  
profile<- GetFootprint(  
  kmer=kmer,  
  tissue=tissue,  
  data.dir=data.dir,  
  frag.type=frag.type,  
  smooth=TRUE  
)  
  
# estimate the shoulders from the profile  
sh <- SobelBorders(profile, kl=nchar(kmer))  
  
#calculate the SFR  
CalcSFR(profile, sh$us, sh$ds, sh$range.us, sh$range.ds)
```

CompareSequences

Description

Wrapper function to compare two input sequences on a kmer basis. Both sequence are split up and compared pairwise, calculating the damage. Total damage is calculated either as summed up or as highest pairwise damage, as chosen. Overlay plots for none, all or only the highest scoring pair are created as specified.

Usage

```
CompareSequences(sequence1,  
                  sequence2,  
                  kl,  
                  damage.mode="exhaustive",  
                  tissue,  
                  data.dir,  
                  vocab.flag=FALSE,  
                  vocab.file=paste0(data.dir,"/",tissue,"/vocabulary_",tissue,".txt"),  
                  frag.type,  
                  plots="highest",  
                  smooth=TRUE,  
                  ylim=c(0,0.01),  
                  xlim=c(-125,125),  
                  plot.shoulders=FALSE)
```

Arguments

sequence1	input sequence 1, character string (usually length of kmer or a window over a variant)
sequence2	input sequence 2, character string, required to have same length as sequence 1
kl	length of the kmer windows, to split up the sequences
damage.mode	{ "exhaustive" , "local" } mode to calculate the total damage for the comparison. "exhaustive" sum over all pair-wise comparison. "local" extract the highest pair-wise damage score. [default= "exhaustive"]
tissue	name of tissue directory of interest
data.dir	path to directory containing preprocessed data
vocab.flag	{FALSE, TRUE} flag if a precalculated vocabulary file is present and should be used [default=FALSE] Speeds up the analysis but precalculation is time consuming, therefore has to be done manually after initial processing of new tissue data. Recommended if tissue is frequently used for analysis.
vocab.file	path to precalculated vocabulary file [default=paste0(data.dir,"",tissue,"vocabulary_",tissue,".txt")]
frag.type	{ "DNase" , "ATAC" } fragmentation type.
plots	{FALSE, "highest" , "all" } indicate if and what overlay plots to retrieve [default= "highest"] FALSE no plots, highest only the highest scoring pair, "all" plots for all pairs
smooth	{FALSE, TRUE} smooth the profile [default=TRUE]
ylim	y-limit for plots [default c(0,0.01)]
xlim	x-limit for plots [default c(-125,125)]
plot.shoulders	{FALSE, TRUE} print the estimated shoulders with the overlay plots [default=FALSE]

Value

Returns a list containing two data frames and a plot or list of plots, depending on how the plotting was specified.

...\$summary: a data frame listing the reference and variant sequence with highest scoring kmer pair

...\$df: data frame containing the pair-wise results per row (kmer and SFR from reference and variant and the damage)

...\$plot: Overlay plot(s), either single ggplot2 plot or list of plots depending on mode selected for "plots"

Examples

```
frag.type <- "DNase"

data.dir <- "./my_sasq_database/human/DNase/"

tissue <- "blood_tissue"

comp <- CompareSequences(
  sequence1="CAGTTTCATGAGG",
  sequence2="CAGTTTATGAGG",
  kl=7,
  data.dir=data.dir,
  damage.mode="exhaustive",
  tissue=tissue,
  vocab.flag=TRUE,
  frag.type="DNase",
```

```
plots="highest"  
)
```

DecodeKmer

Description

Takes a kmer and dissects all ambivalent FASTA characters to create a character vector of all matching kmers.

Usage

```
DecodeKmer(kmer)
```

Arguments

kmer character string, allowing all FASTA standard and ambivalent characters

Value

Returns a vector containing of all matching definite kmers as character strings.

Examples

```
DecodeKmer("WGATAA")
```

DissectSequence

Description

Split an input sequence into kmer windows.

Usage

```
(sequence, kl, list=FALSE)
```

Arguments

sequence character string, allowing all FASTA standard and ambivalent characters
kl length in bp to split the sequence into
list {FALSE, TRUE} select if output should be a list or a vector [default=FALSE]

Value

Returns a vector or list of all kmers.

Examples

```
DissectList("AGGGATACGTAGACGGTGTAA", kl=7, list=FALSE)
```

GetCount

Description

Wrapper function to get the count of the k-mer in DHS in the tissue of interest.

Usage

```
GetCount(kmer, tissue, data.dir, frag.type)
```

Arguments

kmer	input kmer, FASTA characters string 5 - 7 bp
tissue	name of tissue directory of interest
data.dir	path to directory containing preprocessed data
frag.type	{"DNase", "ATAC"} fragmentation type.

Value

Returns the count as single numeric value.

Examples

```
frag.type <- "DNase"

data.dir <- "./my_sasq_database/human/DNase/"

tissue <- "blood_tissue"

GetCount(kmer="CACGTG", data.dir=data.dir, tissue=tissue, frag.type=frag.type)
```

GetFootprint

Description

Wrapper to retrieve the merged & smoothed profile of kmer.

Usage

```
GetFootprint(kmer,
             tissue,
             data.dir,
             frag.type,
             smooth=TRUE,
             smooth.bandwidth=5)
```

Arguments

kmer	input kmer, FASTA characters string 5 - 7 bp
tissue	name of tissue directory of interest
data.dir	path to directory containing preprocessed data
frag.type	{"DNase", "ATAC"} fragmentation type.
smooth	{FALSE, TRUE} smooth the profile [default=TRUE]
smooth.bandwidth	bandwidth to smooth in bp [default=5]

Value

Returns list containing the profile and counts. ...\$profile, ...\$count.

Examples

```
frag.type <- "DNase"

data.dir <- "./my_sasq_database/human/DNase/"

tissue <- "blood_tissue"

GetFootprint(kmer="CACGTG",
             tissue=tissue,
             data.dir=data.dir,
             frag.type=frag.type,
             smooth=TRUE)
```

GetFootprintStrand

Description

Wrapper to retrieve strand-specific (smoothed) profiles of a kmer from tissue or background

Usage

```
GetFootprintStrand(kmer,
                  tissue,
                  data.dir,
                  frag.type,
                  smooth=TRUE,
                  smooth.bandwidth=5,
                  background.flag=FALSE)
```


Arguments

kmer	input kmer, FASTA characters string 5 - 7 bp
tissue	name of tissue directory of interest
data.dir	path to directory containing preprocessed data
frag.type	{"DNase", "ATAC"} fragmentation type.
smooth	{FALSE, TRUE} smooth the profile [default=TRUE]
smooth.bandwidth	bandwidth to smooth in bp [default=5]
background.flag	select if to retrieve the profile from the genome-wide background [default=FALSE]

Value

Returns list with strand-specific profiles and counts. ...\$profile.plus, ...\$profile.minus, ...\$count.plus, ...\$count.minus

Examples

```
frag.type <- "DNase"

data.dir <- "./my_sasq_database/human/DNase/"

tissue <- "blood_tissue"

GetFootprintStrand(kmer="CACGTG",
                   tissue=tissue,
                   data.dir=data.dir,
                   frag.type=frag.type,
                   smooth=TRUE,
                   background.flag=FALSE)
```

GetPossibleMutations

Description

Take a sequence input and split into kmers of length $kl * 2 - 1$ (e.g. for kl 7 always takes the 13 surrounding bases) and separate into reference and variance. Take the parsed position as index for the first base to mutate. List all possible mutations filling the kl 'th position in the variant column with all possible substitutions depending on the reference base. First position of the sequence to be mutated is therefore the kl 'th position of the sequence string.

Usage

```
GetPossibleMutations(sequence, kl=7, chr=".", position=1)
```

Arguments

sequence	character string, allowing all FASTA standard and ambivalent characters
kl	length in bp to split the sequence into
chr	chromosome of the sequence to print
position	bp position of the sequence to print

Value

Returns a 6 column data frame with "chr" "pos" "ref.base" "var.base" "ref.seq" "var.seq"

Examples

```
GetPossibleMutations(sequence=c("AGGGATACGTAGACGGTGTA"), kl=7, chr="chrX", position=1345990)
```

GetSFR

Description

Wrapper function to get the SFR ratio. If indicated and available, use the present vocabulary file to directly grep the SFR. Else get the average profile, estimate the borders and calculate the SFR. Note that for using the vocabulary file only nonambivalent DNA chars are allowed. For the alternative ambivalent chars are decoded.

Usage

```
GetSFR(kmer,
       tissue,
       data.dir,
       vocab.flag=FALSE,
       vocab.file=paste0(data.dir,"/",tissue,"/vocabulary_",tissue,".txt"),
       frag.type="")
```

Arguments

kmer	input kmer, FASTA characters string 5 - 7 bp
tissue	name of tissue directory of interest
data.dir	path to directory containing preprocessed data
vocab.flag	{FALSE, TRUE} flag if a precalculated vocabulary file is present and should be used [default=FALSE] Speeds up the analysis but precalculation is time consuming, therefore has to be done manually after initial processing of new tissue data. Recommended if tissue is frequently used for analysis.
vocab.file	path to precalculated vocabulary file [default=paste0(data.dir,"",tissue,"vocabulary_",tissue,".txt")]
frag.type	{"DNase", "ATAC"} fragmentation type.

Value

Returns SFR as single numeric value.

Examples

```
GetSFR(kmer="CACGTG",
       tissue="blood_tissue",
       data.dir="/my_sasq_database/human/DNase/",
       vocab.flag=TRUE,
       frag.type="DNase")
```

GrepProfile

Description

Grep strand-specific profile of 250 bp surrounding kmer and normalise for total cuts in 250 bp window.

Usage

```
GrepProfile(kmer, infile)
```

Arguments

kmer input kmer, FASTA characters string 5 - 7 bp
infile path to the input file

Value

Returns list of `..$profile` and `..$count`.

Examples

```
frag.type <- "DNase"

data.dir <- "./my_sasq_database/human/DNase/"

tissue <- "blood_tissue"

infile.plus=file.path(data.dir,
                      tissue,
                      "counts",
                      paste0("kmers_", kl, "_count_", tissue, "_pnorm_JH60_plus.txt"))

profile.list <- GrepProfile(kmer="CACGTG", infile=infile)
```

InSilicoMutation

Description

Wrapper for maximum/absolute damage insilico mutation. Takes an input sequence, splits into data frame of desired kmer length matching window sizes for running comparison and compares reference sequence against all possible mutated sequences (single base pair substitutions). Reports according to report mode ("all", "max", "maxabs").

Usage

```
InSilicoMutation(sequence,
  kl=7,
  chr=".",
  position=1,
  report="all",
  damage.mode="exhaustive",
  tissue=tissue,
  data.dir="./my_sasq_database/human/DNase/",
  vocab.flag=TRUE,
  vocab.file=paste0(data.dir,"/",tissue,"/vocabulary_",tissue,".txt"),
  frag.type=frag.type,
  progress.bar=FALSE)
```

Arguments

sequence1	input sequence, character string
kl	length of the kmer windows, to split up the sequence
chr	chromosome of the sequence to print
position	bp position of the sequence to print, first mutated base is the kl'th base. Therefore input +- 6 bp positions of sequence surrounding your sequence of interest and set position to the kl'th index.
report	{ "all", "max", "maxabs" } Select which damage per position to report. [default="all"] "all" = report all 3 possible substitutions per position. Reports three rows per bp position. "max" = only report substitution with highest positive damage. Report one row per bp position, easy to convert to wig. "maxabs" = only report substitution with highest absolute damage. Reports one row per bp position as well.
damage.mode	{ "exhaustive", "local" } mode to calculate the total damage for the comparison. "exhaustive" sum over all pair-wise comparison. "local" extract the highest pair-wise damage score. [default="exhaustive"]
tissue	name of tissue directory of interest
data.dir	path to directory containing preprocessed data
vocab.flag	{ FALSE, TRUE } flag if a precalculated vocabulary file is present and should be used [default=FALSE] Speeds up the analysis but precalculation is time consuming, therefore has to be done manually after initial processing of new tissue data. Recommended if tissue is frequently used for analysis.
vocab.file	path to precalculated vocabulary file [default=paste0(data.dir,"",tissue,"vocabulary_",tissue,".txt")]
frag.type	{ "DNase", "ATAC" } fragmentation type.
progress.bar	{ FALSE, TRUE } Select if to display a progress.bar when running. (Requires packages pbapply if set to TRUE!) [default=FALSE]

Value

Seven columns dataframe c(chr, position, ref.base, var.base, ref.sequence, var.sequence, damage).

Examples

```
# use (for example) the BSgenome package to extract reference genome sequence
library(BSgenome)
library(BSgenome.Hsapiens.UCSC.hg18)
```

```

genome <- BSgenome.Hsapiens.UCSC.hg18

#set the sequence coordinates of the desired genomic location
chr <- "chr16"
start.pos <- 145852
end.pos <- start.pos + 30

# Get the sequence sequence
seq <- as.character(getSeq(genome, "chr16", start=start.pos-6, end=end.pos+6))

#perform the i silcio mutation
df.insilico <- InSilicoMutation(sequence=seq,
                                kl=7,
                                chr="chr16",
                                position=start.pos,
                                report="all",
                                damage.mode="exhaustive",
                                tissue="blood_tissue",
                                data.dir="./my_sasq_database/human/DNase/",
                                vocab.flag=TRUE,
                                frag.type=frag.type,
                                progress.bar = TRUE
                                )

#display as rainbowplot
rp <- RainbowPlot(df.insilico, ylim=c(-4,4))

```

MakeRainbowTrackHub

Description

Wrapper to make a UCSC browser track hub from the insilico mutation data frame. Input is a data frame in the same format as the output data frame of the InSilicoMutation function. 7 columns: chr pos ref.base var.base ref.seq var.seq damage

Usage

```

MakeRainbowTrackHub(input.df,
                     id.tag = "SasQ_RP",
                     store.tracks = paste0("~/", id.tag, "_tracks"),
                     store.hub = paste0("~/", id.tag, "_track_hub"),
                     genome.build = "hg19",
                     path.chr.sizes,
                     short.label = "SasQ In silico mutation Rainbow plot",
                     long.label = "",
                     set.email = "none",
                     bedgraph.to.bigwig.path,
                     make.softlinks = FALSE)

```

Arguments

input.df	input 7 column data frame. chr pos ref.base var.base ref.seq var.seq damage
id.tag	id tag to name the hub directory and bw tracks.
store.tracks	path to directory to store the bigwig tracks.
store.hub	where to store the hub & visualization folder.
genome.build	select genome build ("hg19", "hg18", "mm9", ...) [default="hg19"]
path.chr.sizes	full.path to chr sizes file matching to the selected genome. Not provided in package.
short.label	shortLabel for track hub.
long.label	longLabel for track hub [default = short.label]
set.email	email contact address to appear in trackHub. [default="none"]
bedgraph.to.bigwig.path	full path to UCSC bedGraphToBigWig conversion tool. Not included.
make.softlinks	{FALSE, TRUE} set flag if to directly make softlinks in hub folder [default=FALSE] (e.g. set TRUE if running directly on a cluster so that the final softlinks paths are already correct [default = FALSE] if data files will be copied to a different direction afterwards (e.g. when mounted and ran locally) if FALSE: After creation copy data hub to desired location and create softlinks in the hub folder to the tracks in the track folder e.g. "ln -s store.tracks/*.bw store.hub")

Value

Writes bigWig tracks into desired directory and creates a trackHub structure to migrate to public domain and import to UCSC.

Examples

```
frag.type <- "DNase"
data.dir <- "./my_sasq_database/human/DNase/"
tissue <- "blood_tissue"

# use (for example) the BSgenome package to extract reference genome sequence
library(BSgenome)
library(BSgenome.Hsapiens.UCSC.hg18)
genome <- BSgenome.Hsapiens.UCSC.hg18

#set the sequence coordinates of the desired genomic location
chr <- "chr16"
start.pos <- 145852
end.pos <- start.pos + 30

# Get the sequence sequence
seq <- as.character(getSeq(genome, "chr16", start=start.pos-6, end=end.pos+6))

#perform the in silico mutation
df.insilico <- InSilicoMutation(sequence=seq,
                                kl=7,
                                chr="chr16",
                                position=start.pos,
```

```

        report="all",
        damage.mode="exhaustive",
        tissue="blood_tissue",
        data.dir=data.dir,
        vocab.flag=TRUE,
        frag.type=frag.type,
        progress.bar = TRUE
    )

out.dir <- "./myoutputpath/"

# write rainbowplot ucsc track hub
MakeRainbowTrackHub(
  input.df = df.insilico,
  id.tag = "SasQ_hub_blood",
  store.tracks = paste0(out.dir, "/", id.tag, "_tracks"),
  store.hub = paste0(out.dir, "/", id.tag, "_track_hub"),
  genome.build = "hg19",
  path.chr.sizes = "~/mydatabase/chrom_sizes/hg19_chrom_sizes.txt",
  short.label = "SasQ In silico mutation Rainbow plot",
  bedgraph.to.bigwig.path = "~/mytools/tools/bedGraphToBigWig",
  make.softlinks = FALSE
)

```

PlotOverlap

Description

Plot two average profiles overlapping or on top of each other given two input profiles. Plot the shoulders if plot.shoulders is set to TRUE, then use shoulders list as provided or determine.

Usage

```

PlotOverlap(profile1,
            profile2,
            kmer1,
            kmer2,
            count1 = "NA",
            count2 = "NA",
            ymode = "separate",
            ylim = c(0,0.01),
            xlim = c(-125,125),
            plot.shoulders=FALSE,
            shoulders1=FALSE,
            shoulders2=FALSE)

```

Arguments

profile1	input profile1
profile2	input profile2
kmer1	input k-mer 1 (reference)
kmer2	input k-mer 2 (variant)
count1	count of k-mer1 occurrence
count2	count of k-mer2 occurrence
ymode	mode how to plot the overlapping profiles ("merged" or as "separate" profiles above each other) [default=separate]
ylim	ylim to fix for plot [default c(0, 0.01))]
xlim	xlim to fix for plot [default c(-125, 125)]
plot.shoulders	{FALSE, TRUE} if to plot the estimated shoulders with the profiles [default=FALSE] note that it is only plotted if the separate profile option was selected to keep the plots tidy
shoulders1	list object of estimated shoulder position and ranges for profile 1 [default=FALSE]
shoulders2	list object of estimated shoulder position and ranges for profile 2 [default=FALSE]

Value

Returns ggplot2 plot object of the overlay plot.

Examples

```
frag.type <- "DNase"
data.dir <- "./my_sasq_database/human/DNase/"
tissue <- "blood_tissue"

kmer1 <- "WGATAA" #note FASTA ambiguous code is supported
kmer2 <- "WGATTA"

fp1 <- GetFootprint(kmer=kmer1,
                    tissue=tissue,
                    data.dir=data.dir,
                    frag.type=frag.type,
                    smooth=TRUE)
fp2 <- GetFootprint(kmer=kmer2,
                    tissue=tissue,
                    data.dir=data.dir,
                    frag.type=frag.type,
                    smooth=TRUE)

# make an overlap plot
PlotOverlap(
  fp1$profile,
  fp2$profile,
  kmer1,
  kmer2,
  fp1$count,
  fp2$count,
  ymode="separate"
)
```

PlotOverlapKmers

Description

Wrapper function to produce an overlay plot from two kmers and a tissue input only.

Usage

```
PlotOverlapKmers(kmer1,
                  kmer2,
                  tissue1,
                  tissue2,
                  data.dir,
                  frag.type,
                  smooth=TRUE,
                  ylim=c(0,0.01),
                  xlim=c(-125,125),
                  plot.shoulders=FALSE)
```

Arguments

kmer1	input k-mer 1 (reference)
kmer2	input k-mer 2 (variant)
tissue1	name of tissue 1 directory of interest
tissue2	name of tissue 2 directory of interest
data.dir	path to directory containing preprocessed data
frag.type	{"DNase", "ATAC"} fragmentation type.
smooth	{FALSE TRUE} if to smooth the profiles. [default=TRUE]
ymode	mode how to plot the overlapping profiles ("merged" or as "separate" profiles above each other) [default=separate]
ylim	ylim to fix for plot [default c(0, 0.01)]
xlim	xlim to fix for plot [default c(-125, 125)]
plot.shoulders	{FALSE, TRUE} if to plot the estimated shoulders with the profiles [default=FALSE] note that it is only plotted if the separate profile option was selected to keep the plots tidy

Value

Returns ggplot2 plot object of the overlay plot.

Examples

```
frag.type <- "DNase"
data.dir <- "./my_sasq_database/human/DNase/"
tissue <- "blood_tissue"

PlotOverlapKmers(
  kmer1="CACGTG",
  kmer2="CACGTT",
  tissue1=tissue,
```

```

issue2=tissue,
data.dir=data.dir,
frag.type="DNase",
smooth=TRUE,
plot.shoulders=FALSE
)

```

PlotSingle

Description

Plot the average profile plot given an input profile. Plot the shoulders if plot.shoulders is set to TRUE and use shoulders list provided or determine.

Usage

```

PlotSingle(profile,
           kl=7,
           plot.shoulders=FALSE,
           shoulders=FALSE,
           ylim=c(0,0.01),
           xlim=c(-125,125),
           color="black")

```

Arguments

profile	input profile
kl	length of the kmer windows, to split up the sequences
ylim	ylim to fix for plot [default=c(0, 0.01)]
xlim	xlim to fix for plot [default=c(-125, 125)]
plot.shoulders	{FALSE, TRUE} if to plot the estimated shoulders with the profiles [default=FALSE]
shoulders	list object of estimated shoulder position and ranges for profile 1 [default=FALSE]
color	Select a color for the profile [default="black"]

Value

Returns single ggplot2 plot object of the profile.

Examples

```

frag.type <- "DNase"
data.dir <- "./my_sasq_database/human/DNase/"
tissue <- "blood_tissue"
kmer <- "CACGTG"

# get the footprint
fp <- GetFootprint(kmer=kmer, tissue=tissue, data.dir=data.dir, frag.type=frag.type, smooth=TRUE)

# estimate the shoulders from the profile (use smoothed profile or smooth within call!)

```

```
sh <- SobelBorders(fp$profile, kl=nchar(kmer))

# make single, merged profile plot
p <- PlotSingle(profile=fp$profile,
               kl=nchar(kmer),
               plot.shoulders=TRUE,
               shoulders=sh
               )
```

PlotSingleKmer

Description

Wrapper function to produce a plot from kmer and tissue input only.

Usage

```
PlotSingleKmer(kmer,
               tissue,
               data.dir,
               frag.type,
               smooth=TRUE,
               smooth.bandwidth=5,
               plot.shoulders=FALSE,
               ylim=c(0,0.01),
               xlim=c(-125,125),
               color="black")
```

Arguments

kmer	input k-mer
tissue	name of tissue 1 directory of interest
data.dir	path to directory containing preprocessed data
frag.type	{"DNase", "ATAC"} fragmentation type.
smooth	{FALSE TRUE} if to smooth the profiles. [default=TRUE]
smooth.bandwidth	bandwidth to smooth in bp [default=5]
plot.shoulders	{FALSE, TRUE} if to plot the estimated shoulders with the profiles [default=FALSE]
ylim	ylim to fix for plot [default c(0, 0.01)]
xlim	xlim to fix for plot [default c(-125, 125)]
color	Select a color for the profile [default="black"]

Value

Returns single ggplot2 plot object of the profile.

Examples

```
PlotSingleKmer(kmer="CACGTG",
               tissue="blood_tissue",
               data.dir="./my_sasq_database/human/DNase/",
               frag.type="DNase")
```

PlotSingleStrands

Description

Wrapper function to produce strand specific plots from kmer and tissue input only.

Usage

```
PlotSingleStrands(kmer,
                  tissue,
                  data.dir,
                  frag.type,
                  smooth=TRUE,
                  smooth.bandwidth=5,
                  background.flag=FALSE,
                  ylim=c(0,0.01),
                  xlim=c(-125,125))
```

Arguments

kmer	input k-mer
tissue	name of tissue 1 directory of interest
data.dir	path to directory containing preprocessed data
frag.type	{"DNase", "ATAC"} fragmentation type.
smooth	{FALSE TRUE} if to smooth the profiles. [default=TRUE]
smooth.bandwidth	bandwidth to smooth in bp [default=5]
background.flag	FALSE TRUE Select if to visualize the genome-wide background cut profiles or from tissue data. Will adjust the required and queried data structure accordingly. Note that if set to TRUE, a different data directory, that pointing to the storage of the deproteinized, genome-wide background has to be provided. [default=FALSE]
ylim	ylim to fix for plot [default c(0, 0.01)]
xlim	xlim to fix for plot [default c(-125, 125)]

Value

Returns list of strand-specific profile plots: for plus ...**\$plot.plus** and minus strand ...**\$plot.minus**.

Examples

```
PlotSingleStrands(kmer="WGATAA",
                  tissue = "blood_tissue",
                  data.dir = "./my_sasq_database/human/DNase/",
                  frag.type = "DNase")
```

PruneProfile

Description

Prune an retrieved (250 bp) average profile equally from both directions given the profile and the desired window size around the kmer which to retrieve. Requires an even number as length to prune. Length of the output profile will always be (desired.length + kmer.length).

Usage

```
PruneProfile(profile, desired.length)
```

Arguments

profile	input profile
desired.length	length of to prune to (will be desired.length + kmer.length)

Value

Returns pruned profile as numeric vector.

Examples

```
frag.type <- "DNase"
data.dir <- "./my_sasq_database/human/DNase/"
tissue <- "blood_tissue"
kmer <- "CACGTG"

# get the footprint
fp <- GetFootprint(kmer=kmer,
                   tissue=tissue,
                   data.dir=data.dir,
                   frag.type=frag.type,
                   smooth=TRUE)

#prune
pruned.profile <- PruneProfile(fp$profile, 100)
```

QueryJaspar

Description

Take a sequence input and query it against a set of Jaspar2014 PWMs. (as provided or saved from the JASPAR2014 R package). Requires "Biostrings" and "TFBSTools" packages.

Usage

```
QueryJaspar(sequence, threshold=0.8, pwm.data)
```

Arguments

sequence input sequence
threshold relative percentage score threshold above which to report matches [default=0.8]
pwm.data a stored pwm.RData object as provided with the Sasquatch distribution or as retrieved and saved from JASPAR2014 R package

Value

Returns character string listing the PWM matches above the relative percentage score threshold.

Examples

```
#requires Biostrings and TFBSTools R packages
library(Biostrings)
library(TFBSTools)

#load human.pwm object
load("./my_sasq_database/jaspar/jaspar2014.human.9606.all.versions")

# Single JASPAR query
QueryJaspar(sequence="AGATAATAG", threshold=0.8, pwm.data=human.pwm)
```

QueryJasparBatch

Description

Take a data frame from the RefVarBatch query as input and query it against the set of selected Jaspars2014 PWMs using a selected match.threshold. Select an absolute footprinting damage above which to query jasper.

Usage

```
QueryJasparBatch <- function(df, damage.threshold=0, match.threshold=0.8, pwm.data)
```

Arguments

df 9 column data frame as output from RefVarBatch()
damage.threshold absolute predicted footprinting damage above which a sequence should be selected for the query. [default=0, query all]
match.threshold relative percentage score threshold above which to report matches [default=0.8]
pwm.data a stored pwm.RData object as provided with the Sasquatch distribution or as retrieved and saved from JASPAR2014 R package

Value

Returns data frame with additional column for jasper query results.

Examples

```
#ty example data frame
tdf <- data.frame(
  id=c("1", "2", "3"),
  ref=c("ATAGATAATCGCT", "ATAGATAATCGCT", "ATATATTCTCGCT"),
  var=c("ATAGATCATCGCT", "ATAGATTATCGCT", "ATAGATGATCGCT")
)

#make a reference vs. variant batch query
comp.df <- RefVarBatch(ref.var.df=tdf,
  kl=7,
  damage.mode="exhaustive",
  tissue="blood_tissue",
  data.dir="./my_sasq_database/human/DNase/",
  vocab.flag=TRUE,
  frag.type="DNase")

#requires Biostrings and TFBSTools R packages
library(Biostrings)
library(TFBSTools)

#load human.pwm object
load("./my_sasq_database/jaspar/jaspar2014.human.9606.all.versions")

#query refvar data frame against jaspar pwms
comp.df.jaspar <- QueryJasparBatch(df=comp.df,
  damage.threshold=0.3,
  match.threshold=0.8,
  pwm.data=human.pwm)
```

QueryLongSequence

Description

Wrapper function to split a longer sequence into kmers of length kl and return kmer, SFR and plots if specified.

Usage

```
QueryLongSequence(
  sequence,
  kl,
  tissue,
  data.dir,
  vocab.flag=FALSE,
  vocab.file=paste0(data.dir,"/",tissue,"/vocabulary_",tissue,".txt"),
  frag.type,
  plots=FALSE,
  smooth=TRUE,
  plot.shoulders=TRUE,
  ylim=c(0,0.01),
```

```
xlim=c(-125,125))
```

Arguments

sequence	input sequence, character string
kl	length of the kmer windows, to split up the sequence
tissue	name of tissue directory of interest
data.dir	path to directory containing preprocessed data
vocab.flag	{FALSE, TRUE} flag if a precalculated vocabulary file is present and should be used [default=FALSE] Speeds up the analysis but precalculation is time consuming, therefore has to be done manually after initial processing of new tissue data. Recommended if tissue is frequently used for analysis.
vocab.file	path to precalculated vocabulary file [default=paste0(data.dir,"",tissue,"vocabulary_",tissue,".txt")]
frag.type	{"DNase", "ATAC"} fragmentation type.
plots	{FALSE, TRUE} indicate if to plot a profile per kmer [default=FALSE]
smooth	{FALSE, TRUE} smooth the profile [default=TRUE]
plot.shoulders	{FALSE, TRUE} print the estimated shoulders with the plots [default=TRUE]
ylim	y-limit for plots [default c(0,0.01)]
xlim	x-limit for plots [default c(-125,125)]

Value

Returns `...$df` a data frame listing the splitted kmers with the respective SFR.
`...$plots` if specified list of profile plots with one plot per splitted k-mer.

Examples

```
QueryLongSequence(
  sequence="ATAGATAATCGCT",
  kl,
  tissue="blood_tissue",
  data.dir="./my_sasq_database/human/DNase/",
  vocab.flag=FALSE,
  vocab.file=paste0(data.dir,"/",tissue,"/vocabulary_",tissue,".txt"),
  frag.type="DNase",
  plots=FALSE)
```

RainbowPlot

Description

Make a rainbow plot from data frame as output from InSilicoMutation. Must have been run with `report="all"`

Usage

```
RainbowPlot <- function(df, ylim=c(-2,2))
```


Arguments

df data frame as output from InSilicoMutation (report="all")
ylim y-limits for plot[default=c(-2,2)]

Value

Returns ggplot2 plot object containing the rainbowplot.

Examples

```
# use (for example) the BSgenome package to extract reference genome sequence
library(BSgenome)
library(BSgenome.Hsapiens.UCSC.hg18)
genome <- BSgenome.Hsapiens.UCSC.hg18

#set the sequence coordinates of the desired genomic location
chr <- "chr16"
start.pos <- 145852
end.pos <- start.pos + 30

# Get the sequence sequence
seq <- as.character(getSeq(genome, "chr16", start=start.pos-6, end=end.pos+6))

#perform the i silcio mutation
df.insilico <- InSilicoMutation(sequence=seq,
                                kl=7,
                                chr="chr16",
                                position=start.pos,
                                report="all",
                                damage.mode="exhaustive",
                                tissue="blood_tissue",
                                data.dir="./my_sasq_database/human/DNase/",
                                vocab.flag=TRUE,
                                frag.type=frag.type,
                                progress.bar = TRUE
                                )

#display as rainbowplot
rp <- RainbowPlot(df.insilico, ylim=c(-4,4))
```

RefVarBatch

Description

Wrapper function to analyse multiple Ref-Var-Sequence pairs. Split each sequence into kmers of length kl, get their SFRs (ideally for speed from vocab.file or calculate per instance. Calculate the damage associated with each kmer pair and from that the local or the exhaustive summed up damage of entire sequence pair.

Usage

```
RefVarBatch(ref.var.df,
            kl,
            damage.mode="exhaustive",
            tissue,
            data.dir,
            vocab.flag=FALSE,
            vocab.file=paste0(data.dir,"/",tissue,"/vocabulary_",tissue,".txt"),
            frag.type)
```

Arguments

ref.var.df	three column data frame listing id reference and variance sequence (id reference variant)
kl	length of the kmer windows, to split up the sequences
damage.mode	{ "exhaustive" , "local" } mode to calculate the total damage for the comparison. "exhaustive" sum over all pair-wise comparison. "local" extract the highest pair-wise damage score. [default= "exhaustive"]
tissue	name of tissue directory of interest
data.dir	path to directory containing preprocessed data
vocab.flag	{FALSE, TRUE} flag if a precalculated vocabulary file is present and should be used [default=FALSE] Speeds up the analysis but precalculation is time consuming, therefore has to be done manually after initial processing of new tissue data. Recommended if tissue is frequently used for analysis.
vocab.file	path to precalculated vocabulary file [default=paste0(data.dir,"",tissue,"vocabulary_",tissue,".txt")]
frag.type	{ "DNase" , "ATAC" } fragmentation type.

Value

Returns data frame listing Ref and Var sequence with highest scoring kmer pair and according SFRs and calculated (exhaustive or local) total damage.

Examples

```
#ty example data frame
tdf <- data.frame(
  id=c("1", "2", "3"),
  ref=c("ATAGATAATCGCT", "ATAGATAATCGCT", "ATATATTCTCGCT"),
  var=c("ATAGATCATCGCT", "ATAGATTATCGCT", "ATAGATGATCGCT")
)

#make a reference vs. variant batch query
comp.df <- RefVarBatch(ref.var.df=tdf,
                      kl=7,
                      damage.mode="exhaustive",
                      tissue="blood_tissue",
                      data.dir="./my_sasq_database/human/DNase/",
                      vocab.flag=TRUE,
                      frag.type="DNase")
```

SmoothProfile

Description

Helper function to smooth a profile given the specified bandwidth and a gaussian, normal kernel.

Usage

```
SmoothProfile(profile, bandwidth=5)
```

Arguments

profile input profile (numeric vector)
smooth.bandwidth bandwidth to smooth in bp [default=5]

Value

Returns smoothed profile as numeric vector.

Examples

```
fp <- GetFootprint(kmer="CACGTG",  
                  tissue="blood_tissue",  
                  data.dir="./my_sasq_database/human/DNase/",  
                  frag.type="DNase",  
                  smooth=FALSE)  
  
smoothed.profile <- SmoothProfile(fp$profile, 5)
```

SobelBorders

Description

Estimate the footprint shoulders by based on zero crossings of the 1D 1st derivative approximation of a smoothed profile and get the optimal shoulder range that maximizes the SFRatio of the footprint. (Approximations for speeding up the process are: Estimate the optimal positions based on 6 bp wide shoulders. Then optimize the shoulder width/range with allowed ranges in {4,6,8,10})

Usage

```
SobelBorders(profile, kl)
```

Arguments

profile input profile (numeric vector) has to be smoothed for proper function!
kl {5,6,7} kmer length input

Value

Returns list object containing the shoulder centric positions and their ranges.

- ..**\$us** upstream shoulder center
- ..**\$ds** downstream shoulder center
- ..**\$us.range** width/range of upstream shoulder
- ..**\$ds.range** width/range of downstream shoulder
- ..**\$flag** flag {TRUE FALSE} if shoulders could be estimated

Examples

```
# get the footprint
fp <- GetFootprint(kmer="CACGTG",
                   tissue="blood_tissue",
                   data.dir="./my_sasq_database/human/DNase/",
                   frag.type="DNase",
                   smooth=TRUE)

# estimate the shoulders from the profile
# (use smoothed profile or smooth within call!)
sh <- SobelBorders(fp$profile, kl=nchar(kmer))
```

Sobeln

Description

Helper function to calculate 1D 1st derivative approximation of profile by 1D sobel filtering.

Usage

```
Sobeln(profile)
```

Arguments

profile input profile (numeric vector) should be smoothed!

Value

Returns 1D 1st derivative approximation of the profile as numeric vector.

Examples

```
fp <- GetFootprint(kmer="CACGTG",
                   tissue="blood_tissue",
                   data.dir="./my_sasq_database/human/DNase/",
                   frag.type="DNase",
                   smooth=TRUE)

# estimate the shoulders from the profile
```

```
# (use smoothed profile or smooth within call!)  
approx.1st.deriv.prof<- Sobeln(fp$profile)
```