

# Package ‘PGA’

October 16, 2019

**Type** Package

**Title** An package for identification of novel peptides by customized database derived from RNA-Seq

**Description** This package provides functions for construction of customized protein databases based on RNA-Seq data with/without genome guided, database searching, post-processing and report generation. This kind of customized protein database includes both the reference database (such as Refseq or ENSEMBL) and the novel peptide sequences form RNA-Seq data.

**Version** 1.14.0

**Date** 2019-04-23

**Author** Shaohang Xu, Bo Wen

**Maintainer** Bo Wen <wenbostar@gmail.com>, Shaohang Xu <xsh.skye@gmail.com>

**Depends** R (>= 3.5.0), IRanges, GenomicRanges, Biostrings (>= 2.26.3), data.table, rTANDEM

**Imports** S4Vectors (>= 0.9.25), Rsamtools (>= 1.10.2), GenomicFeatures (>= 1.19.8), biomaRt (>= 2.17.1), stringr, RCurl, Nozzle.R1, VariantAnnotation (>= 1.7.28), rtracklayer, RSQLite, ggplot2, AnnotationDbi, customProDB (>= 1.21.5), pheatmap

**Suggests** RMariaDB, BSgenome.Hsapiens.UCSC.hg19, RUnit, BiocGenerics, BiocStyle, knitr, R.utils

**VignetteBuilder** knitr

**License** GPL-2

**biocViews** Proteomics, ImmunoOncology, MassSpectrometry, Software, ReportWriting, RNASeq, Sequencing

**git\_url** <https://git.bioconductor.org/packages/PGA>

**git\_branch** RELEASE\_3\_9

**git\_last\_commit** 48259d4

**git\_last\_commit\_date** 2019-05-02

**Date/Publication** 2019-10-15

## R topics documented:

addGeneName4Ensembl . . . . .	2
createProDB4DenovoRNASeq . . . . .	3

dbCreator . . . . .	4
easyRun . . . . .	5
parserGear . . . . .	6
PrepareAnnotationEnsembl2 . . . . .	8
PrepareAnnotationRefseq2 . . . . .	9
reportGear . . . . .	10
runTandem . . . . .	11

<b>Index</b>	<b>13</b>
--------------	-----------

---

addGeneName4Ensembl    *Add gene names*

---

### Description

Add "gene name" column to the attachments (eg. snv.tsv and idl.tsv) of report if there is a demand (optional function). This function is useful for "Ensembl" IDs.

### Usage

```
addGeneName4Ensembl(mart, report = "report")
```

### Arguments

mart                    See detail in function [PrepareAnnotationEnsembl](#).

report                  The report directory for adding gene names.

### Value

none

### Examples

```
## Not run:
library(PGA)
# set the biomaRt parameters just as "PrepareAnnotationEnsembl2" did

mart <- biomaRt::useMart("ENSEMBL_MART_ENSEMBL",
  dataset="hsapiens_gene_ensembl",
  host="grch37.ensembl.org",
  path="/biomart/martservice",
  archive=FALSE)
addGeneName4Ensembl(mart=mart,report="report")

## End(Not run)
```

---

`createProDB4DenovoRNASeq`

*Create protein database based on the transcripts from de novo reconstruction of transcriptome from RNA-seq data*

---

## Description

Create protein database based on the transcripts from de novo reconstruction of transcriptome from RNA-seq data.

## Usage

```
createProDB4DenovoRNASeq(infa = "./trinity.fasta", bool_use_3frame = FALSE,
  outmtab = "novel_transcripts_ntx.tab",
  outfa = "./novel_transcripts_ntx.fasta", bool_get_longest = TRUE,
  make_decoy = TRUE, decoy_tag = "#REV#", outfile_name = "test")
```

## Arguments

<code>infa</code>	A FASTA format file containing transcript sequences (such as derived from de novo transcriptome assembly by Trinity)
<code>bool_use_3frame</code>	A logical variable indicating whether to translate the raw sequences with 3-frame (forward). Default is 6-frame translation (FALSE).
<code>outmtab</code>	A txt format file containing the novel transcripts information
<code>outfa</code>	The output fasta format protein sequence file
<code>bool_get_longest</code>	When it's set as TRUE, only the longest protein sequences will be retained after the raw sequences are translated with 3-frame or 6-frame. Otherwise, all the protein sequences with longer than 30 aa will be retained.
<code>make_decoy</code>	A logical variable indicating whether to add the decoy sequences.
<code>decoy_tag</code>	The prefix of decoy sequence IDs.
<code>outfile_name</code>	Output file name

## Value

The database file(s)

## Examples

```
transcript_seq_file <- system.file("extdata/input", "Trinity.fasta", package="PGA")
createProDB4DenovoRNASeq(infa=transcript_seq_file)
```

dbCreator

*Create customized protein database from RNA-Seq data***Description**

The main function to create customized protein database from RNA-Seq data

**Usage**

```
dbCreator(gtffFile = NULL, vcfFile = NULL, bedFile = NULL,
         tabFile = NULL, annotation_path = NULL, outdir, outfile_name,
         lablersid = FALSE, COSMIC = FALSE, bool_get_longest = TRUE,
         organism = "Homo sapiens", make_decoy = TRUE, genome = NULL,
         debug = FALSE, ...)
```

**Arguments**

gtffFile	A GTF format file containing novel transcripts information
vcfFile	A VCF format file containing SNV and INDEL information
bedFile	A BED format file from Tophat2 containing junction information.
tabFile	A Tab format file from Hisat2 containing junction information (e.g. splicesites.tab). As HISAT2 is a successor to both HISAT and TopHat2, we recommend that users switch to this option.
annotation_path	This directory contains numerous pieces of genome annotation information which can be downloaded by <a href="#">PrepareAnnotationEnsembl2</a> or <a href="#">PrepareAnnotationRefseq2</a> .
outdir	Output directory
outfile_name	Output file name
lablersid	A logical indicating whether to do the SNV annotation(dbSNP)
COSMIC	A logical indicating whether to do the SNV annotation(COSMIC)
bool_get_longest	When it's set as TRUE, the longest sequences will be retained after the DNA sequences are six-frame translated into protein sequences. Otherwise, the protein sequences more than 30 aa are retained.
organism	What is the Genus and species of this organism. Please use proper scientific nomenclature for example: "Homo sapiens" and not "human", default is "Homo sapiens".
make_decoy	A logical indicating whether to add the decoy sequences
genome	Genome information. This is a BSgenome object(e.g. Hsapiens). Default is NULL.
debug	Output the debug info (Default is <False>).
...	Additional arguments

**Value**

The database file

**Examples**

```
vcffile <- system.file("extdata/input", "PGA.vcf", package="PGA")
bedfile <- system.file("extdata/input", "junctions.bed", package="PGA")
gtffile <- system.file("extdata/input", "transcripts.gtf", package="PGA")
annotation <- system.file("extdata", "annotation", package="PGA")
outfile_path<-"db/"
outfile_name<-"test"
library(BSgenome.Hsapiens.UCSC.hg19)
dbCreator(gtffile=gtffile,vcffile=vcffile,bedfile=bedfile,
          annotation_path=annotation,outfile_name=outfile_name,
          genome=Hsapiens,outdir=outfile_path)
```

easyRun

*easyRun***Description**

This function is used to automate the peptide identification based on searching the customized database derived from RNA-Seq data.

**Usage**

```
easyRun(gtffile = NULL, vcffile = NULL, bedfile = NULL, spectra = NULL,
        annotation_path = NULL, outdir = "pga_dir", outPrefix = "pga",
        lablersid = FALSE, COSMIC = FALSE, bool_get_longest = TRUE,
        organism = "Homo sapiens", genome = NULL, enzyme = "[KR][X]",
        tol = 10, tolu = "ppm", itol = 0.6, itolu = "Daltons",
        varmod = NULL, fixmod = NULL, miss = 2, maxCharge = 8, ti = FALSE,
        cpu = 0, alignment = 1, xmx = 2, ...)
```

**Arguments**

gtffile	A GTF format file containing novel transcripts information
vcffile	A VCF format file containing SNV and INDEL information
bedfile	A BED format file containing junction information
spectra	MS/MS peak list file
annotation_path	This directory contains numerous pieces of genome annotation information which can be downloaded by <a href="#">PrepareAnnotationEnsembl2</a> or <a href="#">PrepareAnnotationRefseq2</a> .
outdir	Output directory.
outPrefix	The prefix of output file.
lablersid	A logical indicating whether to do the SNV annotation(dbSNP)
COSMIC	A logical indicating whether to do the SNV annotation(COSMIC)
bool_get_longest	When it's set as TRUE, the longest sequences will be retained after the DNA sequences are six-frame translated into protein sequences. Otherwise, the protein sequences more than 30 aa are retained.

organism	What is the Genus and species of this organism. Please use proper scientific nomenclature for example: "Homo sapiens" and not "human", default is "Homo sapiens".
genome	Genome information. This is a BSgenome object (e.g. Hsapiens).
enzyme	Specification of specific protein cleavage sites. Default is "[KR][X]".
tol	Parent ion mass tolerance (monoisotopic mass).
tolu	Parent ion M+H mass tolerance window units.
itol	Fragment ion mass tolerance (monoisotopic mass).
itolu	Unit for fragment ion mass tolerance (monoisotopic mass).
varmod	Specification of potential modifications of residues.
fixmod	Specification of modifications of residues.
miss	The number of missed cleavage sites. Default is 2.
maxCharge	The Maximum parent charge, default is 8
ti	anticipate carbon isotope parent ion assignment errors. Default is false.
cpu	The number of CPU used for X!Tandem search. Default is 1.
alignment	0 or 1 to determine if peptide should be alignment or not. Default is 0.
xmx	The maximum Java heap size. The unit is "G".
...	Additional arguments

### Value

none

### Examples

```
vcffile <- system.file("extdata/input", "PGA.vcf", package="PGA")
bedfile <- system.file("extdata/input", "junctions.bed", package="PGA")
gtffile <- system.file("extdata/input", "transcripts.gtf", package="PGA")
annotation <- system.file("extdata", "annotation", package="PGA")
library(BSgenome.Hsapiens.UCSC.hg19)
msfile <- system.file("extdata/input", "pga.mgf", package="PGA")
easyRun(gtffile=gtffile,vcffile=vcffile,bedfile=bedfile,spectra=msfile,
        annotation_path=annotation,genome=Hsapiens,cpu = 6,
        enzyme = "[KR][X]", varmod = "15.994915@M",itol = 0.05,
        fixmod = "57.021464@C", tol = 10, tolu = "ppm", itolu = "Daltons",
        miss = 2, maxCharge = 8, ti = FALSE,xmx=1)
```

---

parserGear

*Post-processing for the identification result*

---

### Description

This function is mainly for q-value calculation, protein inference and novel peptides spectra annotation.

**Usage**

```
parserGear(file = NULL, db = NULL, outdir = "parser_outdir",
  prefix = "pga", novelPrefix = "VAR", decoyPrefix = "###REV###",
  alignment = 1, xmx = NULL, thread = 1, msfile = NULL)
```

**Arguments**

file	MS/MS search file. Currently, only XML format file of X!Tandem, DAT result of Mascot and mzIdentML result file (MS-GF+, MyriMatch, IPeak, OMSSA, ...) are supported.
db	A FASTA format database file used for MS/MS searching. Usually, it is from the output of the function dbCreator.
outdir	Output directory.
prefix	The prefix of output file.
novelPrefix	The prefix of novel protein ID. Default is "VAR". "VAR" is the prefix which used by function dbCreator. This value should be left to the default when your database is constructed by the function getTrinityDB.
decoyPrefix	The prefix of decoy sequences ID. Default is "###REV###". "###REV###" is the prefix which used by function dbCreator.
alignment	0 or 1 to determine if peptide should be alignment or not. Default is 1.
xmx	The maximum Java heap size. The unit is "G".
thread	This parameter is used to specify the number of threads. "0" represents that all of the available threads are used; "1" represents one thread is used; "2" represents two threads are used, and so on. Default is 1.
msfile	The MS/MS data (mgf format) used.

**Value**

none

**Examples**

```
vcffile <- system.file("extdata/input", "PGA.vcf", package="PGA")
bedfile <- system.file("extdata/input", "junctions.bed", package="PGA")
gtffile <- system.file("extdata/input", "transcripts.gtf", package="PGA")
annotation <- system.file("extdata", "annotation", package="PGA")
outfile_path <- "db/"
outfile_name <- "test"
library(BSgenome.Hsapiens.UCSC.hg19)
dbfile <- dbCreator(gtffile=gtffile,vcffile=vcffile,bedfile=bedfile,
  annotation_path=annotation,outfile_name=outfile_name,
  genome=Hsapiens,outdir=outfile_path)

msfile <- system.file("extdata/input", "pga.mgf", package="PGA")

## X!Tandem as the peptide identification software
idfile <- runTandem(spectra = msfile, fasta = dbfile, outdir = "./", cpu = 6,
  enzyme = "[KR]|[X]", varmod = "15.994915@M", itol = 0.05,
  fixmod = "57.021464@C", tol = 10, tolu = "ppm",
  itolu = "Daltons", miss = 2, maxCharge = 8, ti = FALSE)
parserGear(file = idfile, db = dbfile, decoyPrefix="###REV###",xmx=1,thread=8,
```

```

        outdir = "parser_outdir")

## Mascot as the peptide identification software
dat_file <- system.file("extdata/input", "mascot.dat", package="PGA")
parserGear(file = dat_file, db = dbfile, decoyPrefix="#REV#", xmx=1, thread=8,
           outdir = "parser_outdir_mascot")

```

---

```
PrepareAnnotationEnsembl2
```

*Prepare annotation from ENSEMBL*

---

## Description

Prepare the annotation from ENSEMBL through biomaRt. This function is modified from the function [PrepareAnnotationEnsembl](#) in **customProDB**.

## Usage

```
PrepareAnnotationEnsembl2(mart, annotation_path, splice_matrix = FALSE,
                          dbsnp = NULL, transcript_ids = NULL, COSMIC = FALSE, ...)
```

## Arguments

mart	See detail in function <a href="#">PrepareAnnotationEnsembl</a> .
annotation_path	See detail in function <a href="#">PrepareAnnotationEnsembl</a> .
splice_matrix	See detail in function <a href="#">PrepareAnnotationEnsembl</a> .
dbsnp	See detail in function <a href="#">PrepareAnnotationEnsembl</a> .
transcript_ids	See detail in function <a href="#">PrepareAnnotationEnsembl</a> .
COSMIC	See detail in function <a href="#">PrepareAnnotationEnsembl</a> .
...	Additional arguments

## Value

Several .RData file containing annotations needed for following analysis.

## See Also

[PrepareAnnotationRefseq2](#).

## Examples

```

ensembl <- biomaRt::useMart("ENSEMBL_MART_ENSEMBL", dataset="hsapiens_gene_ensembl",
                           host="grch37.ensembl.org", path="/biomart/martservice",
                           archive=FALSE)

annotation_path <- tempdir()
transcript_ids <- c("ENST00000234420", "ENST00000269305", "ENST00000445888")

PrepareAnnotationEnsembl2(mart=ensembl, annotation_path=annotation_path,
                          splice_matrix=FALSE, dbsnp=NULL, transcript_ids=transcript_ids,
                          COSMIC=FALSE)

```



---

PrepareAnnotationRefseq2

*Prepare annotation from Refseq*

---

## Description

Prepare the annotation for Refseq through UCSC table browser. This function is modified from the function [PrepareAnnotationRefseq](#) in **customProDB**.

## Usage

```
PrepareAnnotationRefseq2(genome = "hg19", CDSfasta, pepfasta, annotation_path,  
  dbsnp = NULL, transcript_ids = NULL, splice_matrix = FALSE,  
  COSMIC = FALSE, ...)
```

## Arguments

genome	See detail in function <a href="#">PrepareAnnotationRefseq</a> .
CDSfasta	See detail in function <a href="#">PrepareAnnotationRefseq</a> .
pepfasta	See detail in function <a href="#">PrepareAnnotationRefseq</a> .
annotation_path	See detail in function <a href="#">PrepareAnnotationRefseq</a> .
dbsnp	See detail in function <a href="#">PrepareAnnotationRefseq</a> .
transcript_ids	See detail in function <a href="#">PrepareAnnotationRefseq</a> .
splice_matrix	See detail in function <a href="#">PrepareAnnotationRefseq</a> .
COSMIC	See detail in function <a href="#">PrepareAnnotationRefseq</a> .
...	Additional arguments

## Value

Several .RData file containing annotations needed for following analysis.

## See Also

[PrepareAnnotationEnsembl2](#).

## Examples

```
## Not run:  
transcript_ids <- c("NM_001126112", "NM_033360", "NR_073499")  
pepfasta <- system.file("extdata", "refseq_pro_seq.fasta",  
  package="customProDB")  
CDSfasta <- system.file("extdata", "refseq_coding_seq.fasta",  
  package="customProDB")  
annotation_path <- tempdir()  
PrepareAnnotationRefseq2(genome='hg19', CDSfasta, pepfasta, annotation_path,  
  dbsnp=NULL, transcript_ids=transcript_ids,  
  splice_matrix=FALSE, COSMIC=FALSE)  
  
## End(Not run)
```

---

reportGear	<i>The main function for report generation</i>
------------	--

---

## Description

The main function for report generation

## Usage

```
reportGear(parser_dir, tab_dir, report_dir)
```

## Arguments

parser_dir	The directory which contains the peptide identification results
tab_dir	The directory which contains the database annotation files
report_dir	The report output directory

## Value

none

## Examples

```
vcffile <- system.file("extdata/input", "PGA.vcf", package="PGA")
bedfile <- system.file("extdata/input", "junctions.bed", package="PGA")
gtffile <- system.file("extdata/input", "transcripts.gtf", package="PGA")
annotation <- system.file("extdata", "annotation", package="PGA")
outfile_path<-"db/"
outfile_name<-"test"
library(BSgenome.Hsapiens.UCSC.hg19)
dbfile <- dbCreator(gtffile=gtffile,vcffile=vcffile,bedfile=bedfile,
                  annotation_path=annotation,outfile_name=outfile_name,
                  genome=Hsapiens,outdir=outfile_path)

msfile <- system.file("extdata/input", "pga.mgf", package="PGA")
idfile <- runTandem(spectra = msfile, fasta = dbfile, outdir = "./", cpu = 6,
                  enzyme = "[KR]|[X]", varmod = "15.994915@M", itol = 0.05,
                  fixmod = "57.021464@C", tol = 10, tolu = "ppm",
                  itolu = "Daltons", miss = 2, maxCharge = 8, ti = FALSE)
parserGear(file = idfile, db = dbfile, decoyPrefix="#REV#",xmx=1,thread=8,
           outdir = "parser_outdir")
reportGear(parser_dir = "parser_outdir", tab_dir = outfile_path,
           report_dir = "report")
```

---

runTandem	<i>run X!Tandem</i>
-----------	---------------------

---

## Description

run X!Tandem

## Usage

```
runTandem(spectra = "", fasta = "", outdir = ".", cpu = 1,
  enzyme = "[KR]|[X]", tol = 10, tolu = "ppm", itol = 0.6,
  itolu = "Daltons", varmod = NULL, fixmod = NULL, miss = 2,
  maxCharge = 8, ti = FALSE)
```

## Arguments

spectra	MS/MS peak list file
fasta	Protein database file for searching.
outdir	The output directory.
cpu	The number of CPU used for X!Tandem search. Default is 1.
enzyme	Specification of specific protein cleavage sites. Default is "[KR] [X]".
tol	Parent ion mass tolerance (monoisotopic mass).
tolu	Parent ion M+H mass tolerance window units.
itol	Fragment ion mass tolerance (monoisotopic mass).
itolu	Unit for fragment ion mass tolerance (monoisotopic mass).
varmod	Specification of potential modifications of residues.
fixmod	Specification of modifications of residues.
miss	The number of missed cleavage sites. Default is 2.
maxCharge	The Maximum parent charge, default is 8
ti	anticipate carbon isotope parent ion assignment errors. Default is false.

## Value

The search result file path

## Examples

```
vcffile <- system.file("extdata/input", "PGA.vcf", package="PGA")
bedfile <- system.file("extdata/input", "junctions.bed", package="PGA")
gtffile <- system.file("extdata/input", "transcripts.gtf", package="PGA")
annotation <- system.file("extdata", "annotation", package="PGA")
outfile_path <- "db/"
outfile_name <- "test"
library(BSgenome.Hsapiens.UCSC.hg19)
dbfile <- dbCreator(gtffile=gtffile,vcffile=vcffile,bedfile=bedfile,
  annotation_path=annotation,outfile_name=outfile_name,
  genome=Hsapiens,outdir=outfile_path)
```

```
msfile <- system.file("extdata/input", "pga.mgf", package="PGA")
runTandem(spectra = msfile, fasta = dbfile, outdir = "./", cpu = 6,
          enzyme = "[KR][X]", varmod = "15.994915eM",
          fixmod = "57.021464eC", tol = 10, tolu = "ppm", itol = 0.05,
          itolu = "Daltons", miss = 2, maxCharge = 8, ti = FALSE)
```

# Index

addGeneName4Ensembl, [2](#)

createProDB4DenovoRNASeq, [3](#)

dbCreator, [4](#)

easyRun, [5](#)

parserGear, [6](#)

PrepareAnnotationEnsembl, [2](#), [8](#)

PrepareAnnotationEnsembl2, [4](#), [5](#), [8](#), [9](#)

PrepareAnnotationRefseq, [9](#)

PrepareAnnotationRefseq2, [4](#), [5](#), [8](#), [9](#)

reportGear, [10](#)

runTandem, [11](#)