

# Package ‘PhyloProfile’

October 17, 2020

**Version** 1.2.8

**Date** 2020-07-21

**Title** PhyloProfile

## Description

PhyloProfile is a tool for exploring complex phylogenetic profiles. Phylogenetic profiles, presence/absence patterns of genes over a set of species, are commonly used to trace the functional and evolutionary history of genes across species and time. With PhyloProfile we can enrich regular phylogenetic profiles with further data like sequence/structure similarity, to make phylogenetic profiling more meaningful. Besides the interactive visualisation powered by R-Shiny, the package offers a set of further analysis features to gain insights like the gene age estimation or core gene identification.

**URL** <https://github.com/BIONF/PhyloProfile/>

**BugReports** <https://github.com/BIONF/PhyloProfile/issues>

**License** MIT + file LICENSE

**Depends** R (>= 4.0.0)

**Encoding** UTF-8

**biocViews** Software, Visualization, DataRepresentation,  
MultipleComparison, FunctionalPrediction

**Imports** ape, bioDist, BiocStyle, Biostrings, colourpicker, data.table,  
DT, energy, ExperimentHub, ggplot2, gridExtra, pbapply,  
RColorBrewer, shiny, shinyBS, shinyjs, OmaDB, plyr, xml2, zoo

**RoxygenNote** 7.1.1

**Suggests** knitr, rmarkdown, testthat

**VignetteBuilder** knitr

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

**git\_branch** RELEASE\_3\_11

**git\_last\_commit** 9e60604

**git\_last\_commit\_date** 2020-07-21

**Date/Publication** 2020-10-16

**Author** Vinh Tran [aut, cre],  
Bastian Greshake Tzovaras [aut],  
Ingo Ebersberger [aut],  
Carla Mölbert [ctb]

**Maintainer** Vinh Tran <[tran@bio.uni-frankfurt.de](mailto:tran@bio.uni-frankfurt.de)>

**R topics documented:**

calcPresSpec	3
checkInputValidity	4
checkNewick	5
checkOmaID	5
clusterDataDend	6
compareMedianTaxonGroups	7
compareTaxonGroups	8
createArchiPlot	9
createGeneAgePlot	10
createLongMatrix	10
createPercentageDistributionData	11
createProfileFromOma	12
createRootedTree	13
createVarDistPlot	13
createVariableDistributionData	15
createVariableDistributionDataSubset	15
dataCustomizedPlot	16
dataFeatureTaxGroup	17
dataMainPlot	18
dataVarDistTaxGroup	19
distributionTest	20
estimateGeneAge	20
fastaParser	21
featureDistTaxPlot	22
filterProfileData	23
finalProcessedProfile	24
fromInputToProfile	25
fullProcessedProfile	27
geneAgePlotDf	28
generateSinglePlot	28
getAllDomainsOma	29
getAllFastaOma	30
getCommonAncestor	31
getCoreGene	31
getDataClustering	32
getDataForOneOma	33
getDendrogram	34
getDistanceMatrix	35
getDomainFolder	36
getFastaFromFasInput	36
getFastaFromFile	37
getFastaFromFolder	38
getIDsRank	39
getInputTaxaID	39
getInputTaxaName	40
getNameList	41
getOmaDataForOneOrtholog	41
getOmaDomainFromURL	42
getOmaMembers	42
getQualColForVector	43

getSelectedFastaOma . . . . .	44
getSelectedTaxonNames . . . . .	44
getTaxonomyInfo . . . . .	45
getTaxonomyMatrix . . . . .	46
getTaxonomyRanks . . . . .	47
gridArrangeSharedLegend . . . . .	47
heatmapPlotting . . . . .	48
highlightProfilePlot . . . . .	49
idList . . . . .	51
mainLongRaw . . . . .	51
mainTaxonomyRank . . . . .	52
pairDomainPlotting . . . . .	52
parseDomainInput . . . . .	53
parseInfoProfile . . . . .	54
ppTaxonomyMatrix . . . . .	55
ppTree . . . . .	55
processNcbiTaxonomy . . . . .	56
profileWithTaxonomy . . . . .	56
qualitativeColours . . . . .	57
rankIndexing . . . . .	58
rankList . . . . .	58
reduceProfile . . . . .	59
runPhyloProfile . . . . .	59
singleDomainPlotting . . . . .	60
sortDomains . . . . .	61
sortInputTaxa . . . . .	62
sortTaxaFromTree . . . . .	63
taxa2dist . . . . .	64
taxonNamesReduced . . . . .	64
taxonomyMatrix . . . . .	65
taxonomyTableCreator . . . . .	65
varDistTaxPlot . . . . .	66
wideToLong . . . . .	67
xmlParser . . . . .	68

## Index 69

---

calcPresSpec	<i>Calculate percentage of present species in each super taxon</i>
--------------	--

---

### Description

Calculate percentage of present species in each super taxon

### Usage

```
calcPresSpec(profileWithTax, taxaCount)
```

### Arguments

profileWithTax	data frame of main PhyloProfile input together with their taxonomy info (see ?profileWithTaxonomy)
taxaCount	number of species occur in each supertaxon (e.g. phylum or kingdom)

**Value**

A data frame with

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[profileWithTaxonomy](#) for a demo input data

**Examples**

```
# NOTE: for internal testing only - not recommended for outside using
data("profileWithTaxonomy", package="PhyloProfile")
taxaCount <- plyr::count(profileWithTaxonomy, "supertaxon")
taxaCount$freq <- 1
calcPresSpec(profileWithTaxonomy, taxaCount)
```

---

checkInputValidity      *Check the validity of the input phylogenetic profile file*

---

**Description**

Check if input file has one of the following format: orthoXML, multiple FASTA, tab-delimited matrix (wide or long), or list of OMA IDs.

**Usage**

```
checkInputValidity(filein)
```

**Arguments**

filein            input file

**Value**

The format of the input file format, or type of error

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[checkOmaID](#)

**Examples**

```
filein <- system.file(
  "extdata", "test.main.wide", package = "PhyloProfile", mustWork = TRUE
)
checkInputValidity(filein)
```

---

checkNewick	<i>Check the validity of input newick tree</i>
-------------	--

---

**Description**

Check the validity of input newick tree

**Usage**

```
checkNewick(tree, inputTaxonID = NULL)
```

**Arguments**

tree	input newick tree
inputTaxonID	list of all input taxon IDs for the phylogenetic profiles

**Value**

Possible formatting error of input tree. 0 = suitable tree for using with PhyloProfile, 1 = missing parenthesis; 2 = missing comma; 3 = tree has singleton; or a list of taxa that do not exist in the input phylogenetic profile.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[getInputTaxaID](#) for getting input taxon IDs, [ppTree](#) for an example of input tree

**Examples**

```
data("ppTree", package="PhyloProfile")
checkNewick(ppTree, c("ncbi3702", "ncbi3711", "ncbi7029"))
```

---

checkOmaID	<i>Check the validity of input OMA IDs</i>
------------	--

---

**Description**

Check if input IDs are valid OMA IDs for OMA Browser

**Usage**

```
checkOmaID(ids)
```

**Arguments**

ids	list of ids needs to be checked
-----	---------------------------------

**Value**

List of invalid IDs (not readable for OMA)

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**Examples**

```
print("Uncomment the following line to run the function")
# checkOmaID("HUMAN29398")
```

---

clusterDataDend	<i>Create a hclust object from the distance matrix</i>
-----------------	--

---

**Description**

Create a hclust object from the distance matrix

**Usage**

```
clusterDataDend(distanceMatrix = NULL, clusterMethod = "complete")
```

**Arguments**

`distanceMatrix` calculated distance matrix (see `?getDistanceMatrix`)  
`clusterMethod` clustering method ("single", "complete", "average" for UPGMA, "mcquitty" for WPGMA, "median" for WPGMC, or "centroid" for UPGMC). Default = "complete".

**Value**

An object class hclust generated based on input distance matrix and a selected clustering method.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[getDataClustering](#), [getDistanceMatrix](#), [hclust](#)

**Examples**

```
data("finalProcessedProfile", package="PhyloProfile")
data <- finalProcessedProfile
profileType <- "binary"
profiles <- getDataClustering(
  data, profileType, var1AggregateBy, var2AggregateBy)
distMethod <- "mutualInformation"
distanceMatrix <- getDistanceMatrix(profiles, distMethod)
clusterMethod <- "complete"
clusterDataDend(distanceMatrix, clusterMethod)
```

---

`compareMedianTaxonGroups`*Compare the median values of a variable between 2 taxon groups*

---

## Description

Given the phylogenetic profiles that contains up to 2 additional variables besides the presence/absence information of the orthologous proteins. This function will compare the median scores of those variables between 2 different taxon groups (e.g. parasitic species vs non-parasitic species), which are defined as in-group and out-group. In-group is identified by the user. Out-group contains all taxa in the input phylogenetic profiles that are not part of the in-group.

## Usage

```
compareMedianTaxonGroups(data, inGroup, useCommonAncestor, variable)
```

## Arguments

<code>data</code>	input phylogenetic profile in long format (see <code>?mainLongRaw</code> and <code>?createLongMatrix</code> )
<code>inGroup</code>	ID list of in-group taxa (e.g. "ncbi1234")
<code>useCommonAncestor</code>	TRUE/FALSE if using all taxa that share the same common ancestor with the pre-selected in-group as the in-group taxa. Default = TRUE.
<code>variable</code>	name of the variable that need to be compared

## Value

List of genes that have a difference in the variable's median scores between the in-group and out-group taxa and their corresponding delta-median.

## Author(s)

Vinh Tran (tran@bio.uni-frankfurt.de)

## Examples

```
data("mainLongRaw", package="PhyloProfile")
data <- mainLongRaw
inGroup <- c("ncbi9606", "ncbi10116")
variable <- colnames(data)[4]
compareMedianTaxonGroups(data, inGroup, TRUE, variable)
```

---

compareTaxonGroups      *Compare the score distributions between 2 taxon groups*

---

### Description

Given the phylogenetic profiles that contains up to 2 additional variables besides the presence/absence information of the orthologous proteins. This function will compare the distribution of those variables between 2 different taxon groups (e.g. parasitic species vs non-parasitic species), which are defined as in-group and out-group. In-group is identified by the user. Out-group contains all taxa in the input phylogenetic profiles that are not part of the in-group.

### Usage

```
compareTaxonGroups(data, inGroup, useCommonAncestor, variable,
  significanceLevel)
```

### Arguments

data	input phylogenetic profile in long format (see ?mainLongRaw and ?createLongMatrix)
inGroup	ID list of in-group taxa (e.g. "ncbi1234")
useCommonAncestor	TRUE/FALSE if using all taxa that share the same common ancestor with the pre-selected in-group as the in-group taxa. Default = TRUE.
variable	name of the variable that need to be compared
significanceLevel	significant cutoff for the statistic test (between 0 and 1). Default = 0.05.

### Value

list of genes that have a significant difference in the variable distributions between the in-group and out-group taxa and their corresponding p-values.

### Author(s)

Vinh Tran (tran@bio.uni-frankfurt.de)

### Examples

```
data("mainLongRaw", package="PhyloProfile")
data <- mainLongRaw
inGroup <- c("ncbi9606", "ncbi10116")
variable <- colnames(data)[4]
compareTaxonGroups(data, inGroup, TRUE, variable, 0.05)
```



---

createArchiPlot      *Create protein's domain architecture plot*

---

## Description

Create architecture plot for both seed and orthologous protein. If domains of ortholog are missing, only architecture of seed protein will be plotted. NOTE: seed protein ID is the one being shown in the profile plot, which normally is also the orthologous group ID.

## Usage

```
createArchiPlot(info = NULL, domainDf = NULL, labelArchiSize = 12,  
               titleArchiSize = 12)
```

## Arguments

`info`                    a list contains seed and ortholog's IDs

`domainDf`                dataframe contains domain info for the seed and ortholog. This including the seed ID, orthologs IDs, sequence lengths, feature names, start and end positions, feature weights (optional) and the status to determine if that feature is important for comparison the architecture between 2 proteins\* (e.g. seed protein vs ortholog) (optional).

`labelArchiSize`        lable size (in px). Default = 12.

`titleArchiSize`        title size (in px). Default = 12.

## Value

A domain plot as `arrangeGrob` object. Use `grid::grid.draw(plot)` to render.

## Author(s)

Vinh Tran tran@bio.uni-frankfurt.de

## See Also

[singleDomainPlotting](#), [sortDomains](#), [parseDomainInput](#), [getQualColForVector](#)

## Examples

```
seedID <- "101621at6656"  
orthoID <- "101621at6656|AGRPL@224129@0|224129_0:001955|1"  
info <- c(seedID, orthoID)  
domainFile <- system.file(  
  "extdata", "domainFiles/101621at6656.domains",  
  package = "PhyloProfile", mustWork = TRUE  
)  
domainDf <- parseDomainInput(seedID, domainFile, "file")  
plot <- createArchiPlot(info, domainDf, 9, 9)  
grid::grid.draw(plot)
```

---

createGeneAgePlot      *Create gene age plot*

---

**Description**

Create gene age plot

**Usage**

```
createGeneAgePlot(geneAgePlotDf, textFactor = 1)
```

**Arguments**

geneAgePlotDf    data frame required for plotting gene age (see ?geneAgePlotDf)  
textFactor        increase factor of text size

**Value**

A gene age distribution plot as a ggplot2 object

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[estimateGeneAge](#) and [geneAgePlotDf](#)

**Examples**

```
geneAgePlotDf <- data.frame(  
  name = c("Streptophyta (Phylum)", "Bikonta", "Eukaryota (Superkingdom)"),  
  count = c(7, 1, 30),  
  percentage = c(18, 3, 79)  
)  
createGeneAgePlot(geneAgePlotDf)
```

---

createLongMatrix      *Create a long matrix format for all kinds of input phylogenetic profiles*

---

**Description**

Create a long matrix format for all kinds of input phylogenetic profiles

**Usage**

```
createLongMatrix(inputFile = NULL)
```

**Arguments**

inputFile        input profile file in orthoXML, multiple FASTA, tab-delimited matrix format (wide or long).

**Value**

A data frame of input data in long-format containing seed gene IDs ( or orthologous group IDs), their orthologous proteins together with the corresponding taxonomy IDs and values of (up to) two additional variables.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[xmlParser](#), [fastaParser](#), [wideToLong](#)

**Examples**

```
inputFile <- system.file(
  "extdata", "test.main.wide", package = "PhyloProfile", mustWork = TRUE
)
createLongMatrix(inputFile)
```

---

createPercentageDistributionData

*Create data for percentage present taxa distribution*

---

**Description**

Create data for percentage present taxa distribution

**Usage**

```
createPercentageDistributionData(inputData, rankName = NULL)
```

**Arguments**

inputData        dataframe contains raw input data in long format (see ?mainLongRaw)  
rankName         name of the working taxonomy rank (e.g. "species", "family")

**Value**

A dataframe for analysing the distribution of the percentage of species in the selected supertaxa, containing the seed protein IDs, percentage of their orthologs in each supertaxon and the corresponding supertaxon names.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[mainLongRaw](#)

**Examples**

```
data("mainLongRaw", package="PhyloProfile")
createPercentageDistributionData(mainLongRaw, "class")
```

---

`createProfileFromOma` *Create a phylogenetic profile from a raw OMA dataframe*

---

**Description**

Create a phylogenetic profile from a raw OMA dataframe

**Usage**

```
createProfileFromOma(finalOmaDf = NULL)
```

**Arguments**

`finalOmaDf` raw OMA data for a list of proteins (see `?getDataForOneOma`)

**Value**

Dataframe of the phylogenetic profiles in long format, which contains the seed protein IDs, their orthologous proteins and the corresponding taxonmy IDs of the orthologs.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[getDataForOneOma](#)

**Examples**

```
print("Uncomment the following lines to run the function")
# omaData <- getDataForOneOma("HUMAN29397", "OG")
# createProfileFromOma(omaData)
```

---

createRootedTree      *Create rooted tree from a taxonomy matrix*

---

### Description

Create rooted tree from a taxonomy matrix

### Usage

```
createRootedTree(df, rootTaxon = NULL)
```

### Arguments

df	data frame contains taxonomy matrix used for generating tree (see distDf in example)
rootTaxon	taxon used for rooting the taxonomy tree

### Value

A rooted taxonomy tree as an object of class "phylo".

### Author(s)

Vinh Tran tran@bio.uni-frankfurt.de

### See Also

[taxa2dist](#) for distance matrix generation from a taxonomy matrix, [getTaxonomyMatrix](#) for getting taxonomy matrix, [ppTaxonomyMatrix](#) for a demo taxonomy matrix data

### Examples

```
data("ppTaxonomyMatrix", package = "PhyloProfile")
# prepare matrix for calculating distances
distDf <- subset(ppTaxonomyMatrix, select = -c(ncbiID, fullName))
row.names(distDf) <- distDf$abbrName
distDf <- distDf[, -1]
# create taxonomy tree rooted by ncbi10090
createRootedTree(distDf, "ncbi10090")
```

---

createVarDistPlot      *Create distribution plot*

---

### Description

Create distribution plot for one of the additional variable or the percentage of the species present in the supertaxa.

## Usage

```
createVarDistPlot(data, varName = "var", varType = "var1",  
  percent = c(0, 1), textSize = 12)
```

## Arguments

data	dataframe contains data for plotting (see <code>?createVariableDistributionData</code> , <code>?createVariableDistributionDataSubset</code> or <code>?createPercentageDistributionData</code> )
varName	name of the variable that need to be analyzed (either name of variable 1 or variable 2 or "percentage of present taxa"). Default = "var".
varType	type of variable (either "var1", "var2" or "presSpec"). Default = "var1".
percent	range of percentage cutoff (between 0 and 1). Default = c(0,1)
textSize	text size of the distribution plot (in px). Default = 12.

## Value

A distribution plot for the selected variable as a ggplot object

## Author(s)

Vinh Tran tran@bio.uni-frankfurt.de

## See Also

[mainLongRaw](#), [createVariableDistributionData](#), [createVariableDistributionDataSubset](#), [createPercentageDistributionData](#)

## Examples

```
data("mainLongRaw", package="PhyloProfile")  
data <- createVariableDistributionData(  
  mainLongRaw, c(0, 1), c(0.5, 1)  
)  
varName <- "Variable abc"  
varType <- "var1"  
percent <- c(0,1)  
textSize <- 12  
createVarDistPlot(  
  data,  
  varName,  
  varType,  
  percent,  
  textSize  
)
```

---

createVariableDistributionData  
*Create data for additional variable distribution*

---

**Description**

Create data for additional variable distribution

**Usage**

```
createVariableDistributionData(inputData, var1Cutoff = c(0, 1),  
                             var2Cutoff = c(0, 1))
```

**Arguments**

inputData	dataframe contains raw input data in long format (see ?mainLongRaw)
var1Cutoff	min and max cutoff for var1. Default = c(0, 1).
var2Cutoff	min and max cutoff for var2. Default = c(0, 1).

**Value**

A dataframe for analysing the distribution of the additional variable(s) containing the protein (ortholog) IDs and the values of their variables (var1 and var2).

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[mainLongRaw](#)

**Examples**

```
data("mainLongRaw", package="PhyloProfile")  
createVariableDistributionData(  
  mainLongRaw, c(0, 1), c(0.5, 1)  
)
```

---

createVariableDistributionDataSubset  
*Create data for additional variable distribution (for a subset data)*

---

**Description**

Create data for additional variable distribution (for a subset data)

**Usage**

```
createVariableDistributionDataSubset(fullProfileData,  
                                   distributionData, selectedGenes, selectedTaxa)
```

**Arguments**

`fullProfileData` dataframe contains the full processed profiles (see `?fullProcessedProfile`, `?filterProfileData` or `?fromInputToProfile`)

`distributionData` dataframe contains the full distribution data (see `?createVariableDistributionData`)

`selectedGenes` list of genes of interest. Default = "all".

`selectedTaxa` list of taxa of interest Default = "all".

**Value**

A dataframe for analysing the distribution of the additional variable(s) for a subset of genes and/or taxa containing the protein (ortholog) IDs and the values of their variables (`var1` and `var2`).

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[parseInfoProfile](#), [createVariableDistributionData](#), [fullProcessedProfile](#), [mainLongRaw](#)

**Examples**

```
data("fullProcessedProfile", package="PhyloProfile")
data("mainLongRaw", package="PhyloProfile")
distributionData <- createVariableDistributionData(
  mainLongRaw, c(0, 1), c(0.5, 1)
)
selectedGenes <- "100136at6656"
selectedTaxa <- c("Mammalia", "Saccharomycetes", "Insecta")
createVariableDistributionDataSubset(
  fullProcessedProfile,
  distributionData,
  selectedGenes,
  selectedTaxa
)
```

---

`dataCustomizedPlot`      *Create data for customized profile plot*

---

**Description**

Create data for customized profile plot based on a selected list of genes and/or taxa, containing seed protein IDs (`geneID`), ortholog IDs (`orthoID`) together with their ncbi taxonomy IDs (`ncbiID` and `abbrName`), full names (`fullName`), indexed supertaxa (`supertaxon`), values for additional variables (`var1`, `var2`) and the aggregated values of those additional variables for each supertaxon (`mVar1`, `mVar2`), number of original and filtered co-orthologs in each supertaxon (`paralog` and `paralogNew`), number of species in each supertaxon (`numberSpec`) and the each supertaxon (`presSpec`).



**Usage**

```
dataCustomizedPlot(dataHeat = NULL, selectedTaxa = "all",
  selectedSeq = "all")
```

**Arguments**

dataHeat            a data frame contains processed profiles (see `?fullProcessedProfile`, `?filterProfileData`)

selectedTaxa        selected subset of taxa. Default = "all".

selectedSeq         selected subset of genes. Default = "all".

**Value**

A dataframe contains data for plotting the customized profile.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[filterProfileData](#)

**Examples**

```
data("fullProcessedProfile", package="PhyloProfile")
selectedTaxa <- c("Mammalia", "Saccharomycetes", "Insecta")
selectedSeq <- "all"
dataCustomizedPlot(fullProcessedProfile, selectedTaxa, selectedSeq)
```

---

dataFeatureTaxGroup    *Create data for feature distribution comparison plot*

---

**Description**

Create data for plotting the distribution of the protein domain features between 2 group of taxa for a selected gene (average number of feature occurrence per protein/ortholog).

**Usage**

```
dataFeatureTaxGroup(mainDf, domainDf, inGroup, gene)
```

**Arguments**

mainDf              input phylogenetic profile in long format (see `?mainLongRaw` and `?createLongMatrix`)

domainDf            dataframe contains domain info for the seed and ortholog. This including the seed ID, orthologs IDs, sequence lengths, feature names, start and end positions, feature weights (optional) and the status to determine if that feature is important for comparison the architecture between 2 proteins\* (e.g. seed protein vs ortholog) (optional). (see `?parseDomainInput`)

inGroup ID list of in-group taxa (e.g. "ncbi1234")  
 gene ID of gene that need to be plotted the feature distribution comparison between in- and out-group taxa.

**Value**

Dataframe containing all feature names, their frequencies (absolute count and the average instances per protein - IPP) in each taxon group and the corresponding taxa group type (in- or out-group).

**Author(s)**

Vinh Tran (tran@bio.uni-frankfurt.de)

**See Also**

[createLongMatrix](#), [parseDomainInput](#)

**Examples**

```
data("mainLongRaw", package="PhyloProfile")
mainDf <- mainLongRaw
gene <- "101621at6656"
inputFile <- system.file(
  "extdata", "domainFiles/101621at6656.domains",
  package = "PhyloProfile", mustWork = TRUE
)
type <- "file"
domainDf <- parseDomainInput(gene, inputFile, type)
inGroup <- c("ncbi9606", "ncbi10116")
dataFeatureTaxGroup(mainDf, domainDf, inGroup, gene)
```

---

dataMainPlot

*Create data for main profile plot*

---

**Description**

Create data for main profile plot

**Usage**

```
dataMainPlot(dataHeat = NULL)
```

**Arguments**

dataHeat a data frame contains processed profiles (see [?fullProcessedProfile](#), [?filterProfileData](#))

**Value**

A dataframe for plotting the phylogenetic profile, containing seed protein IDs (geneID), ortholog IDs (orthoID) together with their ncbi taxonomy IDs (ncbiID and abbrName), full names (fullName), indexed supertaxa (supertaxon), values for additional variables (var1, var2) and the aggregated values of those additional variables for each supertaxon (mVar1, mVar2), number of original and filtered co-orthologs in each supertaxon (paralog and paralogNew), number of species in each supertaxon (numberSpec) and the species that have orthologs in each supertaxon (presSpec).

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[filterProfileData](#)

**Examples**

```
data("fullProcessedProfile", package="PhyloProfile")
dataMainPlot(fullProcessedProfile)
```

---

dataVarDistTaxGroup     *Create data for variable distribution comparison plot*

---

**Description**

Create data for plotting the distribution comparison between 2 groups of taxa for a selected gene.

**Usage**

```
dataVarDistTaxGroup(data, inGroup, gene, variable)
```

**Arguments**

data	input phylogenetic profile in long format (see ?mainLongRaw and ?createLongMatrix)
inGroup	ID list of in-group taxa (e.g. "ncbi1234")
gene	ID of gene that need to be plotted the distribution comparison between in- and out-group taxa.
variable	var1 or c(var1, var2)

**Value**

Dataframe containing list of values for all available variables for the selected genes in in-group and out-group taxa (max. 3 columns).

**Author(s)**

Vinh Tran (tran@bio.uni-frankfurt.de)

**See Also**

[createLongMatrix](#)

**Examples**

```
data("mainLongRaw", package="PhyloProfile")
data <- mainLongRaw
inGroup <- c("ncbi9606", "ncbi10116")
variable <- colnames(data)[c(4, 5)]
dataVarDistTaxGroup(data, inGroup, "101621at6656", variable)
```

---

distributionTest      *Compare the distribution of 2 numeric vectors*

---

### Description

This function tests the difference between the distributions of two input numeric samples using the statistical test. First the Kolmogorov-Smirnov is used to check if 2 samples have the same distribution. If yes, Wilcoxon-Mann-Whitney will be used to compare the distribution difference.

### Usage

```
distributionTest(varIn, varOut, significanceLevel)
```

### Arguments

varIn                  first numeric vector  
varOut                 second numeric vector  
significanceLevel      significant cutoff of the Kolmogorov-Smirnov test. Default = 0.05.

### Value

p-value of the comparison test.

### Author(s)

Carla Mölbert (carla.moelbert@gmx.de)

---

estimateGeneAge      *Calculate the phylogenetic gene age from the phylogenetic profiles*

---

### Description

Calculate the phylogenetic gene age from the phylogenetic profiles

### Usage

```
estimateGeneAge(processedProfileData, rankName, refTaxon,  
                  var1C0, var2C0, percentC0)
```

### Arguments

processedProfileData      dataframe contains the full processed phylogenetic profiles (see ?fullProcessed-Profile or ?parseInfoProfile)  
rankName                  working taxonomy rank (e.g. "species", "genus", "family")  
refTaxon                  reference taxon name (e.g. "Homo sapiens", "Homo" or "Hominidae")  
var1C0                    cutoff for var1. Default: c(0, 1)  
var2C0                    cutoff for var2. Default: c(0, 1)  
percentC0                 cutoff for percentage of species present in each supertaxon. Default: c(0, 1)

**Value**

A dataframe contains estimated gene ages for the seed proteins.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[parseInfoProfile](#) for creating a full processed profile dataframe; [getNameList](#) and [getTaxonomyMatrix](#) for getting taxonomy info, [fullProcessedProfile](#) for a demo input dataframe

**Examples**

```
data("fullProcessedProfile", package="PhyloProfile")
rankName <- "class"
refTaxon <- "Mammalia"
processedProfileData <- fullProcessedProfile
var1Cutoff <- c(0, 1)
var2Cutoff <- c(0, 1)
percentCutoff <- c(0, 1)
estimateGeneAge(
  processedProfileData,
  rankName,
  refTaxon,
  var1Cutoff, var2Cutoff, percentCutoff
)
```

---

fastaParser

*Parse multi-fasta input file*

---

**Description**

Parse multi-fasta input file

**Usage**

```
fastaParser(inputFile = NULL)
```

**Arguments**

**inputFile** input multiple fasta file. Check `extdata/test.main.fasta` or <https://github.com/BIONF/PhyloProfile/wiki/Data#multi-fasta-format> for the supported FASTA header.

**Value**

A data frame of input data in long-format containing seed gene IDs ( or orthologous group IDs), their orthologous proteins together with the corresponding taxonomy IDs and values of (up to) two additional variables.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**Examples**

```
inputFile <- system.file(
  "extdata", "test.main.fasta", package = "PhyloProfile", mustWork = TRUE
)
fastaParser(inputFile)
```

---

featureDistTaxPlot     *Create feature distribution comparison plot*

---

**Description**

Create protein feature distribution plots between 2 groups of taxa for a selected gene.

**Usage**

```
featureDistTaxPlot(data, plotParameters)
```

**Arguments**

data	dataframe for plotting (see ?dataFeatureTaxGroup)
plotParameters	plot parameters, including size of x-axis, y-axis, legend and title; position of legend ("right", "bottom" or "none"); names of in-group and out-group; flip the plot coordinate ("Yes" or "No"). NOTE: Leave blank or NULL to use default values.

**Value**

Distribution plots as a ggplot2 object.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[dataFeatureTaxGroup](#)

**Examples**

```
data("mainLongRaw", package="PhyloProfile")
data <- mainLongRaw
gene <- "101621at6656"
inputFile <- system.file(
  "extdata", "domainFiles/101621at6656.domains",
  package = "PhyloProfile", mustWork = TRUE
)
type <- "file"
domainDf <- parseDomainInput(gene, inputFile, type)
inGroup <- c("ncbi9606", "ncbi10116")
plotDf <- dataFeatureTaxGroup(data, domainDf, inGroup, gene)
plotParameters <- list(
  "xSize" = 12,
  "ySize" = 12,
```

```

    "angle" = 15,
    "legendSize" = 12,
    "inGroupName" = "In-group",
    "outGroupName" = "Out-group",
    "flipPlot" = "No"
  )
  featureDistTaxPlot(plotDf, plotParameters)

```

---

filterProfileData      *Filter phylogenetic profiles*

---

### Description

Create a filtered data needed for plotting or clustering phylogenetic profiles. NOTE: this function require some intermediate steps using the results from other functions. If you would like to get a full processed data from the raw input, please use the function from `InputToProfile()` instead!

### Usage

```

filterProfileData(DF, refTaxon = NULL,
  percentCO = c(0, 1), coorthoCOMax = 9999,
  var1CO = c(0, 1), var2CO = c(0, 1), var1Rel = "protein",
  var2Rel = "protein", groupByCat = FALSE, catDt = NULL)

```

### Arguments

DF	a reduced dataframe contains info for all phylogenetic profiles in the selected taxonomy rank.
refTaxon	selected reference taxon. NOTE: This taxon will not be affected by the filtering. If you want to filter all, set <code>refTaxon &lt;- NULL</code> . Default = NULL.
percentCO	min and max cutoffs for percentage of species present in a supertaxon. Default = <code>c(0, 1)</code> .
coorthoCOMax	maximum number of co-orthologs allowed. Default = 9999.
var1CO	min and max cutoffs for var1. Default = <code>c(0, 1)</code> .
var2CO	min and max cutoffs for var2. Default = <code>c(0, 1)</code> .
var1Rel	relation of var1 ("protein" for protein-protein or "species" for protein-species). Default = "protein".
var2Rel	relation of var2 ("protein" for protein-protein or "species" for protein-species). Default = "protein".
groupByCat	group genes by their categories (TRUE or FALSE). Default = FALSE.
catDt	dataframe contains gene categories (optional, NULL if <code>groupByCat = FALSE</code> or no info provided). Default = NULL.

### Value

A filtered dataframe for generating profile plot including seed gene IDs (or orthologous group IDs), their ortholog IDs and the corresponding (super)taxa, (super)taxon IDs, number of co-orthologs in each (super)taxon, values for two additional variables var1, var2, supertaxon, and the categories of seed genes (or ortholog groups).

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[parseInfoProfile](#) and [reduceProfile](#) for generating input dataframe, [fullProcessedProfile](#) for a demo full processed profile dataframe, [fromInputToProfile](#) for generating fully processed data from raw input.

**Examples**

```
# NOTE: this function require some intermediate steps using the results from
# other functions. If you would like to get a full processed data from the
# raw input, please use the function fromInputToProfile() instead!
data("fullProcessedProfile", package="PhyloProfile")
superTaxonDf <- reduceProfile(fullProcessedProfile)
refTaxon <- "Mammalia"
percentCutoff <- c(0.0, 1.0)
coorthologCutoffMax <- 10
var1Cutoff <- c(0.75, 1.0)
var2Cutoff <- c(0.5, 1.0)
var1Relation <- "protein"
var2Relation <- "species"
groupByCat <- FALSE
catDt <- NULL
filterProfileData(
  superTaxonDf,
  refTaxon,
  percentCutoff,
  coorthologCutoffMax,
  var1Cutoff,
  var2Cutoff,
  var1Relation,
  var2Relation,
  groupByCat,
  catDt
)
```

---

finalProcessedProfile *An example of a final processed & filtered phylogenetic profile.*

---

**Description**

An example of a final processed & filtered phylogenetic profile.

**Usage**

```
data(finalProcessedProfile)
```



**Format**

A data frame with 91 rows and 10 variables:

- geneID Seed or ortholog group ID, e.g. "100136at6656"
- supertaxon Supertaxon name together with its ordered index, e.g. "1001\_Mammalia"
- supertaxonID Supertaxon ID (only different than ncbiID in case working with higher taxonomy rank than input's). e.g. "40674"
- var1 First additional variable
- presSpec The percentage of species presenting in each supertaxon
- category "cat"
- orthoID Ortholog ID, e.g. "100136at6656|RAT@10116@1|G3V7R8|1"
- var2 Second additional variable
- paralog Number of co-orthologs in the current taxon
- taxonMon Name of supertaxon but without index, e.g. "Mammalia"

---

fromInputToProfile      *Complete processing of raw input phylogenetic profiles*

---

**Description**

Create a processed and filtered data for plotting or analysing phylogenetic profiles from raw input file (from raw input to final filtered dataframe)

**Usage**

```
fromInputToProfile(rawInput, rankName, refTaxon = NULL,
  taxaTree = NULL, var1AggregateBy = "max", var2AggregateBy = "max",
  percentCutoff = c(0, 1), coorthologCutoffMax = 9999,
  var1Cutoff = c(0, 1), var2Cutoff = c(0, 1), var1Relation = "protein",
  var2Relation = "protein", groupByCat = FALSE, catDt = NULL)
```

**Arguments**

rawInput	input file (in long, wide, multi-fasta or orthoxml format)
rankName	taxonomy rank (e.g. "species", "phylum", ...)
refTaxon	selected reference taxon name (used for sorting and will be protected from filtering). Default = NULL.
taxaTree	input taxonomy tree for taxa in input profiles (optional). Default = NULL.
var1AggregateBy	aggregate method for var1 (min, max, mean or median). Default = "max".
var2AggregateBy	aggregate method for VAR2 (min, max, mean or median). Default = "max".
percentCutoff	min and max cutoffs for percentage of species present in a supertaxon. Default = c(0, 1).
coorthologCutoffMax	maximum number of co-orthologs allowed. Default = 9999.

var1Cutoff	min and max cutoffs for var1. Default = c(0, 1).
var2Cutoff	min and max cutoffs for var2. Default = c(0, 1).
var1Relation	relation of var1 ("protein" for protein-protein or "species" for protein-species). Default = "protein".
var2Relation	relation of var2 ("protein" for protein-protein or "species" for protein-species). Default = "protein".
groupByCat	group genes by their categories (TRUE or FALSE). Default = FALSE.
catDt	dataframe contains gene categories. Default = NULL.

### Value

Dataframe required for generating phylogenetic profile plot or clustering analysis. It contains seed gene IDs (or orthologous group IDs), their ortholog IDs and the corresponding (super)taxa, (super)taxon IDs, number of co-orthologs in each (super)taxon, values for two additional variables var1, var2, categories of seed genes (or ortholog groups).

### Author(s)

Vinh Tran tran@bio.uni-frankfurt.de

### See Also

[createLongMatrix](#), [getInputTaxaID](#), [getInputTaxaName](#), [sortInputTaxa](#), [parseInfoProfile](#), [reduceProfile](#), [filterProfileData](#)

### Examples

```
rawInput <- system.file(
  "extdata", "test.main.long", package = "PhyloProfile", mustWork = TRUE
)
rankName <- "class"
refTaxon <- "Mammalia"
taxaTree <- NULL
var1AggregateBy <- "max"
var2AggregateBy <- "mean"
percentCutoff <- c(0.0, 1.0)
coorthologCutoffMax <- 10
var1Cutoff <- c(0.75, 1.0)
var2Cutoff <- c(0.5, 1.0)
var1Relation <- "protein"
var2Relation <- "species"
groupByCat <- FALSE
catDt <- NULL
fromInputToProfile(
  rawInput,
  rankName,
  refTaxon,
  taxaTree,
  var1AggregateBy,
  var2AggregateBy,
  percentCutoff,
  coorthologCutoffMax,
  var1Cutoff,
  var2Cutoff,
```

```

    var1Relation,
    var2Relation,
    groupByCat,
    catDt
  )

```

---

fullProcessedProfile *An example of a fully processed phylogenetic profile.*

---

### Description

An example of a fully processed phylogenetic profile.

### Usage

```
data(fullProcessedProfile)
```

### Format

A data frame with 168 rows and 17 variables:

- supertaxon Supertaxon name together with its ordered index, e.g. "1001\_Mammalia"
- geneID Seed or ortholog group ID, e.g. "100136at6656"
- ncbiID Taxon ID, e.g. "ncbi10116"
- orthoID Ortholog ID, e.g. "100136at6656|HUMAN@9606@1|Q9UNQ2|1"
- var1 First additional variable
- var2 Second additional variable
- paralog Number of co-orthologs in the current taxon
- abbrName NCBI ID of the ortholog, e.g. "ncbi9606"
- taxonID Taxon ID of the ortholog, in this case: "0"
- fullName Full taxon name of the ortholog, e.g. "Homo sapiens"
- supertaxonID Supertaxon ID (only different than ncbiID in case working with higher taxonomy rank than input's). e.g. "40674"
- rank Rank of the supertaxon, e.g. "class"
- category "cat"
- presSpec The percentage of species presenting in each supertaxon
- mVar1 Value of the 1. variable after grouping into supertaxon
- mVar2 Value of the 2. variable after grouping into supertaxon
- numberSpec Total number of species in each supertaxon

---

geneAgePlotDf      *Create data for plotting gene ages*

---

**Description**

Create data for plotting gene ages

**Usage**

```
geneAgePlotDf(geneAgeDf)
```

**Arguments**

geneAgeDf      data frame containing estimated gene ages for seed proteins

**Value**

A dataframe for plotting gene age plot containing the absolute number and percentage of genes for each calculated evolutionary ages and the corresponding position for writing those number on the plot.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[estimateGeneAge](#)

**Examples**

```
geneAgeDf <- data.frame(  
  geneID = c("100136at6656", "100265at6656", "101621at6656", "103479at6656"),  
  cat = c("0000001", "0000011", "0000001", "0000011"),  
  age = c("07_LUCA", "06_Eukaryota", "07_LUCA", "06_Eukaryota")  
)  
geneAgePlotDf(geneAgeDf)
```

---

generateSinglePlot      *Create a single violin distribution plot*

---

**Description**

Create a single violin distribution plot

**Usage**

```
generateSinglePlot(plotDf, parameters, variable)
```

**Arguments**

plotDf	dataframe for plotting containing values for each variable in in-group and out-group.
parameters	plot parameters, including size of x-axis, y-axis, legend and title; position of legend ("right", "bottom" or "none"); mean/median point; names of in-group and out-group; and plot title. NOTE: Leave blank or NULL to use default values.
variable	name of variable that need to be plotted (one of the column names of input dataframe plotDf).

**Value**

A violin plot as a ggplot object.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**Examples**

```
data("mainLongRaw", package="PhyloProfile")
data <- mainLongRaw
inGroup <- c("ncbi9606", "ncbi10116")
varNames <- colnames(data)[c(4, 5)]
plotDf <- dataVarDistTaxGroup(data, inGroup, "101621at6656", varNames)
plotParameters <- list(
  "xSize" = 12,
  "ySize" = 12,
  "titleSize" = 15,
  "legendSize" = 12,
  "legendPosition" = "right",
  "mValue" = "mean",
  "inGroupName" = "In-group",
  "outGroupName" = "Out-group",
  "title" = "101621at6656"
)
generateSinglePlot(plotDf, plotParameters, colnames(plotDf)[1])
```

---

getAlldomainsOma

*Create domain annotation dataframe from a raw OMA dataframe*

---

**Description**

Create domain annotation dataframe from a raw OMA dataframe

**Usage**

```
getAlldomainsOma(finalOmaDf = NULL)
```

**Arguments**

finalOmaDf      raw OMA data for a list of proteins (see ?getDataForOneOma)

**Value**

Dataframe of the domain annotation used for PhyloProfile, which contains seed IDs, ortholog IDs, ortholog lengths, annotated features, start and end positions of those features.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[getDataForOneOma](#)

**Examples**

```
print("Uncomment the following line to run the function")
# omaData <- getDataForOneOma("HUMAN29397", "OG")
# getAllDomainsOma(omaData)
```

---

getAllFastaOma

*Get all fasta sequences from a raw OMA dataframe*

---

**Description**

Get all fasta sequences from a raw OMA dataframe

**Usage**

```
getAllFastaOma(finalOmaDf = NULL)
```

**Arguments**

finalOmaDf      raw OMA data for a list of proteins (see ?getDataForOneOma)

**Value**

A list contains all protein sequences in fasta format.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[getDataForOneOma](#)

**Examples**

```
print("Uncomment the following line to run the function")
# omaData <- getDataForOneOma("HUMAN29397", "OG")
# getAllFastaOma(omaData)
```

---

getCommonAncestor	<i>Get all taxa that share a common ancestor</i>
-------------------	--

---

### Description

Identify the common ancestor for a selected taxa and return a list of all taxa that have that common ancestor from an large input taxa set.

### Usage

```
getCommonAncestor(inputTaxa = NULL, inGroup = NULL)
```

### Arguments

inputTaxa	ID list of all input taxa (e.g. "ncbi12345")
inGroup	ID list of selected taxa used for identify the common ancestor (e.g.: "ncbi55555")

### Value

A list containing the taxonomy rank and name of the common ancestor, together with a dataframe storing the full taxonomy info of all taxa that share that corresponding common ancestor.

### Author(s)

Vinh Tran (tran@bio.uni-frankfurt.de)

### Examples

```
inputTaxa <- c("ncbi34740", "ncbi9606", "ncbi374847", "ncbi123851",
              "ncbi5664", "ncbi189518", "ncbi418459", "ncbi10116", "ncbi284812",
              "ncbi35128", "ncbi7070")
inGroup <- c("ncbi9606", "ncbi10116")
getCommonAncestor(inputTaxa, inGroup)
```

---

getCoreGene	<i>Identify core genes for a list of selected taxa</i>
-------------	--

---

### Description

Identify core genes for a list of selected (super)taxa. The identified core genes must be present in at least a certain proportion of species in each selected (super)taxon (identified via percentCutoff) and that criteria must be fulfilled for a certain percentage of selected taxa or all of them (determined via coreCoverage).

### Usage

```
getCoreGene(rankName, taxaCore = c("none"), profileDt,
            var1Cutoff = c(0, 1), var2Cutoff = c(0, 1), percentCutoff = c(0, 1),
            coreCoverage = 1)
```

**Arguments**

rankName	working taxonomy rank (e.g. "species", "genus", "family")
taxaCore	list of selected taxon names
profileDt	dataframe contains the full processed phylogenetic profiles (see ?fullProcessedProfile or ?parseInfoProfile)
var1Cutoff	cutoff for var1. Default = c(0, 1).
var2Cutoff	cutoff for var2. Default = c(0, 1).
percentCutoff	cutoff for percentage of species present in each supertaxon. Default = c(0, 1).
coreCoverage	the least percentage of selected taxa should be considered. Default = 1.

**Value**

A list of identified core genes.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[parseInfoProfile](#) for creating a full processed profile dataframe

**Examples**

```
data("fullProcessedProfile", package="PhyloProfile")
rankName <- "class"
taxaCore <- c("Mammalia", "Saccharomycetes", "Insecta")
profileDt <- fullProcessedProfile
var1Cutoff <- c(0.75, 1.0)
var2Cutoff <- c(0.75, 1.0)
percentCutoff <- c(0.0, 1.0)
coreCoverage <- 1
getCoreGene(
  rankName,
  taxaCore,
  profileDt,
  var1Cutoff, var2Cutoff,
  percentCutoff, coreCoverage
)
```

---

getDataClustering      *Get data for calculating distance matrix from phylogenetic profiles*

---

**Description**

Get data for calculating distance matrix from phylogenetic profiles

**Usage**

```
getDataClustering(data, profileType = "binary", var1AggBy = "max",
  var2AggBy = "max")
```



**Arguments**

data	a data frame contains processed and filtered profiles (see ?fullProcessedProfile and ?filterProfileData, ?fromInputToProfile)
profileType	type of data used for calculating the distance matrix. Either "binary" (consider only the presence/absence status of orthologs), or "var1"/"var2" for taking values of the additional variables into account. Default = "binary".
var1AggBy	aggregate method for VAR1 (min, max, mean or median). Default = "max".
var2AggBy	aggregate method for VAR2 (min, max, mean or median). Default = "max".

**Value**

A wide dataframe contains values for calculating distance matrix.

**Author(s)**

Carla Mölbert (carla.moelbert@gmx.de), Vinh Tran (tran@bio.uni-frankfurt.de)

**See Also**

[fromInputToProfile](#)

**Examples**

```
data("fullProcessedProfile", package="PhyloProfile")
data <- fullProcessedProfile
profileType <- "binary"
var1AggregateBy <- "max"
var2AggregateBy <- "mean"
getDataClustering(data, profileType, var1AggregateBy, var2AggregateBy)
```

---

getDataForOneOma	<i>Get OMA info for a query protein and its orthologs</i>
------------------	---

---

**Description**

Get taxonomy IDs, sequences, length and annotations for an OMA orthologous group (or OMA HOG).

**Usage**

```
getDataForOneOma(seedID = NULL, orthoType = "OG")
```

**Arguments**

seedID	OMA protein ID
orthoType	type of OMA orthologs ("OG" or "HOG"). Default = "OG".

**Value**

Data frame contains info for all sequences of the input OMA group (or HOG). That info contains the protein IDs, taxonomy IDs, sequences, lengths, domain annotations (tab delimited) and the corresponding seed ID.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**Examples**

```
print("Uncomment the following line to run the function")
# getDataForOneOma("HUMAN29397", "OG")
```

---

getDendrogram	<i>Plot dendrogram tree</i>
---------------	-----------------------------

---

**Description**

Plot dendrogram tree

**Usage**

```
getDendrogram(dd = NULL)
```

**Arguments**

dd                      dendrogram object (see ?clusterDataDend)

**Value**

A dendrogram plot for the genes in the input phylogenetic profiles.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[clusterDataDend](#)

**Examples**

```
data("finalProcessedProfile", package="PhyloProfile")
data <- finalProcessedProfile
profileType <- "binary"
profiles <- getDataClustering(
  data, profileType, var1AggregateBy, var2AggregateBy)
distMethod <- "mutualInformation"
distanceMatrix <- getDistanceMatrix(profiles, distMethod)
clusterMethod <- "complete"
dd <- clusterDataDend(distanceMatrix, clusterMethod)
getDendrogram(dd)
```

---

getDistanceMatrix	<i>Calculate the distance matrix</i>
-------------------	--------------------------------------

---

### Description

Calculate the distance matrix

### Usage

```
getDistanceMatrix(profiles = NULL, method = "mutualInformation")
```

### Arguments

profiles	dataframe contains profile data for distance calculating (see ?getDataClustering)
method	distance calculation method ("euclidean", "maximum", "manhattan", "canberra", "binary", "distanceCorrelation", "mutualInformation" or "pearson" for binary data; "distanceCorrelation" or "mutualInformation" for non-binary data). Default = "mutualInformation".

### Value

A calculated distance matrix for input phylogenetic profiles.

### Author(s)

Carla Mölbert (carla.moelbert@gmx.de), Vinh Tran (tran@bio.uni-frankfurt.de)

### See Also

[getDataClustering](#)

### Examples

```
data("finalProcessedProfile", package="PhyloProfile")
data <- finalProcessedProfile
profileType <- "binary"
profiles <- getDataClustering(
  data, profileType, var1AggregateBy, var2AggregateBy)
method <- "mutualInformation"
getDistanceMatrix(profiles, method)
```

---

getDomainFolder      *Get domain file from a folder for a seed protein*

---

**Description**

Get domain file from a folder for a seed protein

**Usage**

```
getDomainFolder(seed, domainPath)
```

**Arguments**

seed	seed protein ID
domainPath	path to domain folder

**Value**

Domain file and its complete directory path for the selected protein.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**Examples**

```
## Not run:
domainPath <- paste0(
  path.package("PhyloProfile", quiet = FALSE), "/extdata/domainFiles"
)
getDomainFolder("OG_1009", domainPath)

## End(Not run)
```

---

getFastaFromFasInput      *Get fasta sequences from main input file in multi-fasta format*

---

**Description**

Get fasta sequences from main input file in multi-fasta format

**Usage**

```
getFastaFromFasInput(seqIDs = NULL, file = NULL)
```

**Arguments**

seqIDs	list of sequences IDs. Set seqIDs = "all" if you want to get all fasta sequences from the input file.
file	raw phylogenetic profile input file in multi-fasta format.

**Value**

A dataframe with one column contains sequences in fasta format.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**Examples**

```
file <- system.file(
  "extdata", "test.main.fasta",
  package = "PhyloProfile", mustWork = TRUE
)
getFastaFromFasInput("all", file)
```

---

getFastaFromFile	<i>Get fasta sequences from main input file in multi-fasta format</i>
------------------	---

---

**Description**

Get fasta sequences from main input file in multi-fasta format

**Usage**

```
getFastaFromFile(seqIDs = NULL, concatFasta = NULL)
```

**Arguments**

seqIDs	list of sequences IDs. Set seqIDs = "all" if you want to get all fasta sequences from the concatenated input fasta file.
concatFasta	input concatenated fasta file.

**Value**

A dataframe with one column contains sequences in fasta format.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**Examples**

```
concatFasta <- system.file(
  "extdata", "fastaFiles/concatenatedFile.fa",
  package = "PhyloProfile", mustWork = TRUE
)
getFastaFromFasInput("all", concatFasta)
```

---

getFastaFromFolder      *Get fasta sequences*

---

### Description

Get fasta sequences for the input phylogenetic profiles.

### Usage

```
getFastaFromFolder(seqIDs = NULL, path = NULL, dirFormat = NULL,  
  fileExt = NULL, idFormat = NULL)
```

### Arguments

seqIDs	list of sequences IDs.
path	path to fasta folder.
dirFormat	directory format (either 1 for "path/speciesID.fa*" or 2 for "path/speciesID/speciesID.fa*")
fileExt	fasta file extension ("fa", "fasta", "fas" or "txt")
idFormat	fasta header format (1 for ">speciesID:seqID", 2 for ">speciesID@seqID", 3 for ">speciesID seqID" or 4 for "seqID")

### Value

A dataframe with one column contains sequences in fasta format.

### Author(s)

Vinh Tran tran@bio.uni-frankfurt.de

### See Also

[mainLongRaw](#)

### Examples

```
seqIDs <- "RAT@10116@1|D3ZUE4"  
path <- system.file(  
  "extdata", "fastaFiles", package = "PhyloProfile", mustWork = TRUE  
)  
dirFormat <- 1  
fileExt <- "fa"  
idFormat <- 3  
getFastaFromFolder(seqIDs, path, dirFormat, fileExt, idFormat)
```

---

getIDsRank	<i>Get taxonomy info for a list of taxa</i>
------------	---

---

**Description**

Get NCBI taxonomy IDs, ranks and names for an input taxon list.

**Usage**

```
getIDsRank(inputTaxa = NULL, currentNCBIinfo = NULL)
```

**Arguments**

inputTaxa            NCBI ID list of input taxa.

currentNCBIinfo

table/dataframe of the pre-processed NCBI taxonomy data (/PhyloProfile/data/preProcessedTaxonom

**Value**

A list of 3 dataframes: idList, rankList and reducedInfoList. The "rankList" contains taxon names and all taxonomy ranks of the input taxa including also the noranks from the input rank to the taxonomy root. The "idList" contains input taxon IDs, taxon names, all the ranks from current rank to the taxonomy root together with their IDs (with the format "id#rank"). The reducedInfoList is a subset of preProcessedTaxonomy.txt file, containing the NCBI IDs, taxon fullnames, their current rank and their direct parent ID.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**Examples**

```
inputTaxa <- c("272557", "176299")
ncbiFilein <- system.file(
  "extdata", "data/preProcessedTaxonomy.txt",
  package = "PhyloProfile", mustWork = TRUE
)
currentNCBIinfo <- as.data.frame(data.table::fread(ncbiFilein))
getIDsRank(inputTaxa, currentNCBIinfo)
```

---

getInputTaxaID	<i>Get ID list of input taxa from the main input</i>
----------------	--

---

**Description**

Get ID list of input taxa from the main input

**Usage**

```
getInputTaxaID(rawProfile = NULL)
```

**Arguments**

rawProfile      A dataframe of input phylogenetic profile in long format

**Value**

List of all input taxon IDs (e.g. ncbi1234). Default = NULL.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[createLongMatrix](#), [mainLongRaw](#)

**Examples**

```
data("mainLongRaw", package="PhyloProfile")
getInputTaxaID(mainLongRaw)
```

---

<code>getInputTaxaName</code>	<i>Get NCBI taxon names for a selected list of taxa</i>
-------------------------------	---

---

**Description**

Get NCBI taxon names from "PhyloProfile/data/taxonNamesReduced.txt" for a list of input taxa

**Usage**

```
getInputTaxaName(rankName, taxonIDs = NULL)
```

**Arguments**

rankName      taxonomy rank (e.g. "species", "phylum",...)  
 taxonIDs      list of taxon IDs (e.g. ncbi1234). Default = NULL

**Value**

Data frame contains a list of full names, taxonomy ranks and parent IDs for the input taxa.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[getInputTaxaID](#) for getting input taxon IDs, [getNameList](#) for getting the full taxon name list

**Examples**

```
taxonIDs <- c("ncbi9606", "ncbi10116")
getInputTaxaName("species", taxonIDs)
```



---

getNameList	<i>Get list of pre-installed NCBI taxon names</i>
-------------	---

---

**Description**

Get all NCBI taxon names from "PhyloProfile/data/taxonNamesReduced.txt"

**Usage**

```
getNameList()
```

**Value**

List of taxon IDs, their full names, taxonomy ranks and parent IDs obtained from "PhyloProfile/data/taxonNamesReduced.txt"

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**Examples**

```
getNameList()
```

---

getOmaDataForOneOrtholog	<i>Get taxonomy ID, sequence and annotation for one OMA protein</i>
--------------------------	---

---

**Description**

Get taxonomy ID, sequence and annotation for one OMA protein

**Usage**

```
getOmaDataForOneOrtholog(id = NULL)
```

**Arguments**

id	oma ID of one protein
----	-----------------------

**Value**

Data frame contains the input protein ID with its taxonomy ID, sequence, length and domain annotations (tab delimited) for input OMA protein

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**Examples**

```
print("Uncomment the following line to run the function")  
# getOmaDataForOneOrtholog("HUMAN29397")
```

---

getOmaDomainFromURL     *Get domain annotation from OMA Browser*

---

### Description

Get domain annotation from OMA Browser based on a URL or a raw data frame contains annotation info from OMA

### Usage

```
getOmaDomainFromURL(domainURL = NULL)
```

### Arguments

domainURL     URL address for domain annotation of ONE OMA id or a raw data frame contains annotation info from OMA

### Value

Data frame contains feature names with their start and end positions

### Author(s)

Vinh Tran tran@bio.uni-frankfurt.de

### Examples

```
print("Uncomment the following line to run the function")
# getOmaDomainFromURL("https://omabrowser.org/api/protein/7916808/domains/")
```

---

getOmaMembers     *Get OMA members*

---

### Description

Get OMA ortholog group, OMA HOG or OMA pair's members for a seed protein from OMA Browser.

### Usage

```
getOmaMembers(id = NULL, orthoType = "OG")
```

### Arguments

id     ID of the seed protein (OMA or UniProt ID)  
orthoType     type of OMA orthologs: either "HOG", "OG" (orthologous group) or "PAIR" (orthologous pair - CURRENTLY NOT WORKING). Default = "OG".

### Value

List of OMA orthologs for an input seed protein.

**Author(s)**

Carla Mölbert carla.moelbert@gmx.de

**Examples**

```
print("Uncomment the following line to run the function")
# getOmaMembers("HUMAN29397", "OG")
```

---

getQualColForVector     *Get color for a list of items*

---

**Description**

Get color for a list of items

**Usage**

```
getQualColForVector(x = NULL)
```

**Arguments**

x                    input list

**Value**

list of colors for each element (same elements will have the same color)

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[qualitativeColours](#)

**Examples**

```
items <- c("a", "b", "c")
getQualColForVector(items)
```

getSelectedFastaOma *Get selected fasta sequences from a raw OMA dataframe*

---

**Description**

Get selected fasta sequences from a raw OMA dataframe

**Usage**

```
getSelectedFastaOma(finalOmaDf = NULL, seqID = NULL)
```

**Arguments**

finalOmaDf      raw OMA data for a list of proteins (see ?getDataForOneOma)  
seqID            OMA ID of selected protein

**Value**

Required protein sequence in fasta format.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[getDataForOneOma](#)

**Examples**

```
print("Uncomment the following line to run the function")  
# omaData <- getDataForOneOma("HUMAN29397", "OG")  
# getSelectedFastaOma(omaData, "HUMAN29397")
```

---

getSelectedTaxonNames *Get a subset of input taxa based on a selected taxonomy rank*

---

**Description**

Get a subset of taxon ncbi IDs and names from an input list of taxa based on a selected supertaxon (identified by its taxonomy rank and supertaxon name or supertaxon ID).

**Usage**

```
getSelectedTaxonNames(inputTaxonIDs, rank, higherRank, higherID,  
                      higherName)
```

**Arguments**

inputTaxonIDs	list of input taxon IDs (e.g. c("10116", "122586"))
rank	taxonomy rank of input taxa (e.g. "species")
higherRank	selected taxonomy rank (e.g. "phylum")
higherID	supertaxon ID (e.g. 7711). NOTE: either supertaxon ID or name is required, not necessary to give both.
higherName	supertaxon name (e.g. "Chordata"). NOTE: either supertaxon ID or name is required, not necessary to give both.

**Value**

A data frame contains ncbi IDs and names of taxa from the input taxon list that belong to the selected supertaxon.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**Examples**

```
inputTaxonIDs <- c("10116", "122586", "123851", "13616", "188937", "189518",
"208964", "224129", "224324", "237631", "243230")
rank <- "species"
higherRank <- "phylum"
higherID <- 7711
getSelectedTaxonNames(inputTaxonIDs, rank, higherRank, higherID, NULL)
higherName <- "Chordata"
getSelectedTaxonNames(inputTaxonIDs, rank, higherRank, NULL, higherName)
```

---

getTaxonomyInfo	<i>Get taxonomy info for a list of input taxa</i>
-----------------	---

---

**Description**

Get taxonomy info for a list of input taxa

**Usage**

```
getTaxonomyInfo(inputTaxa = NULL, currentNCBIinfo = NULL)
```

**Arguments**

inputTaxa	NCBI taxonomy IDs of input taxa.
currentNCBIinfo	table/dataframe of the pre-processed NCBI taxonomy data (/PhyloProfile/data/preProcessedTaxonom

**Value**

A list of NCBI taxonomy info for input taxa, including the taxonomy IDs, full scientific names, taxonomy ranks and the parent IDs.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**Examples**

```
inputTaxa <- c("272557", "176299")
ncbiFilein <- system.file(
  "extdata", "data/preProcessedTaxonomy.txt",
  package = "PhyloProfile", mustWork = TRUE
)
currentNCBIinfo <- as.data.frame(data.table::fread(ncbiFilein))
getTaxonomyInfo(inputTaxa, currentNCBIinfo)
```

---

getTaxonomyMatrix      *Get taxonomy matrix*

---

**Description**

Get the (full or subset) taxonomy matrix from "data/taxonomyMatrix.txt" based on an input taxon list

**Usage**

```
getTaxonomyMatrix(subsetTaxaCheck = FALSE, taxonIDs = NULL)
```

**Arguments**

subsetTaxaCheck      TRUE/FALSE subset taxonomy matrix based on input taxon IDs. Default = FALSE.

taxonIDs              list of input taxon IDs (e.g. ncbi1234). Default = NULL.

**Value**

Data frame contains the (subset of) taxonomy matrix for list of input taxa.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**Examples**

```
# get full pre-installed taxonomy matrix
getTaxonomyMatrix(FALSE, NULL)
# get taxonomy matrix for a list of taxon IDs
taxonIDs <- c("ncbi9606", "ncbi10116")
getTaxonomyMatrix(TRUE, taxonIDs)
```

---

getTaxonomyRanks      *Create a list containing all main taxonomy ranks*

---

**Description**

Create a list containing all main taxonomy ranks

**Usage**

```
getTaxonomyRanks()
```

**Value**

A list of all main ranks (from strain to superkingdom)

**Author(s)**

Carla Mölbert (carla.moelbert@gmx.de)

**Examples**

```
getTaxonomyRanks()
```

---

gridArrangeSharedLegend  
*Plot Multiple Graphs with Shared Legend in a Grid*

---

**Description**

Plot Multiple Graphs with Shared Legend in a Grid

**Usage**

```
gridArrangeSharedLegend(..., ncol = length(list(...)), nrow = 1,  
  position = c("bottom", "right"), title = NA, titleSize = 12)
```

**Arguments**

...	Plots to be arranged in grid
ncol	Number of columns in grid
nrow	Number of rows in grid
position	Grid position (bottom or right)
title	Title of grid
titleSize	Size of grid title

**Value**

Grid of plots with common legend

**Note**

adapted from <https://rdr.io/github/PhilBoileau/CLSAR/src/R/gridArrangeSharedLegend.R>

**Author(s)**

Phil Boileau, <philippe.boileau (at) rimuhc.ca>

**Examples**

```
data("mainLongRaw", package="PhyloProfile")
data <- mainLongRaw
inGroup <- c("ncbi9606", "ncbi10116")
varNames <- colnames(data)[c(4, 5)]
plotDf <- dataVarDistTaxGroup(data, inGroup, "101621at6656", varNames)
plotParameters <- list(
  "xSize" = 12,
  "ySize" = 12,
  "titleSize" = 15,
  "legendSize" = 12,
  "legendPosition" = "right",
  "mValue" = "mean",
  "inGroupName" = "In-group",
  "outGroupName" = "Out-group",
  "title" = "101621at6656"
)
plotVar1 <- generateSinglePlot(plotDf, plotParameters, colnames(plotDf)[1])
plotVar2 <- generateSinglePlot(plotDf, plotParameters, colnames(plotDf)[2])
g <- gridArrangeSharedLegend(
  plotVar1, plotVar2,
  position = plotParameters$legendPosition,
  title = plotParameters$title,
  size = plotParameters$titleSize
)
```

---

heatmapPlotting

*Create profile heatmap plot*

---

**Description**

Create profile heatmap plot

**Usage**

```
heatmapPlotting(data = NULL, parm = NULL)
```

**Arguments**

data	dataframe for plotting the heatmap phylogenetic profile (either full or subset profiles)
parm	plot parameters, including (1) type of x-axis "taxa" or "genes" - default = "taxa"; (2+3) names of 2 variables var1ID and var2ID - default = "var1" & "var2"; (4) color for lowest var1 - default = "#FF8C00"; (5) color for highest var1 - default = "#4682B4"; (6) color for lowest var2 - default = "#FFFFFF", (7) color



for highest var2 - default = "#F0E68C", (8) color of co-orthologs - default = "#07D000"; (9+10+11) text sizes for x, y axis and legend - default = 9 for each; (12) legend position "top", "bottom", "right", "left" or "none" - default = "top"; (13) zoom ratio of the co-ortholog dots from -1 to 3 - default = 0; (14) angle of x-axis from 0 to 90 - default = 60; (14) show/hide separate line for reference taxon 1/0 - default = 0; (15) enable/disable coloring gene categories TRUE/FALSE - default = FALSE). NOTE: Leave blank or NULL to use default values.

### Value

A profile heatmap plot as a ggplot object.

### Author(s)

Vinh Tran tran@bio.uni-frankfurt.de

### See Also

[dataMainPlot](#), [dataCustomizedPlot](#)

### Examples

```
data("fullProcessedProfile", package="PhyloProfile")
plotDf <- dataMainPlot(fullProcessedProfile)
plotParameter <- list(
  "xAxis" = "taxa",
  "var1ID" = "FAS_FW",
  "var2ID" = "FAS_BW",
  "lowColorVar1" = "#FF8C00",
  "highColorVar1" = "#4682B4",
  "lowColorVar2" = "#FFFFFF",
  "highColorVar2" = "#F0E68C",
  "paraColor" = "#07D000",
  "xSize" = 8,
  "ySize" = 8,
  "legendSize" = 8,
  "mainLegend" = "top",
  "dotZoom" = 0,
  "xAngle" = 60,
  "guideline" = 0,
  "colorByGroup" = FALSE
)

heatmapPlotting(plotDf, plotParameter)
```

---

highlightProfilePlot *Highlight gene and/or taxon of interest on the phylogenetic profile plot*

---

### Description

Highlight gene and/or taxon of interest on the phylogenetic profile plot

**Usage**

```
highlightProfilePlot(data, plotParameter = NULL, taxonHighlight =
  "none", rankName = "none", geneHighlight = "none")
```

**Arguments**

<code>data</code>	dataframe for plotting the heatmap phylogenetic profile (either full or subset profiles)
<code>plotParameter</code>	plot parameters, including (1) type of x-axis "taxa" or "genes" - default = "taxa"; (2+3) names of 2 variables var1ID and var2ID - default = "var1" & "var2"; (4) color for lowest var1 - default = "#FF8C00"; (5) color for highest var1 - default = "#4682B4"; (6) color for lowest var2 - default = "#FFFFFF", (7) color for highest var2 - default = "#F0E68C", (8) color of co-orthologs - default = "#07D000"; (9+10+11) text sizes for x, y axis and legend - default = 9 for each; (12) legend position "top", "bottom", "right", "left" or "none" - default = "top"; (13) zoom ratio of the co-ortholog dots from -1 to 3 - default = 0; (14) angle of x-axis from 0 to 90 - default = 60; (15) show/hide separate line for reference taxon 1/0 - default = 0; (16) enable/disable coloring gene categories TRUE/FALSE - default = FALSE). NOTE: Leave blank or NULL to use default values.
<code>taxonHighlight</code>	taxon of interest. Default = "none".
<code>rankName</code>	working taxonomy rank (needed only for highlight taxon).
<code>geneHighlight</code>	gene of interest. Default = "none".

**Value**

A profile heatmap plot with highlighted gene and/or taxon of interest as ggplot object.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[dataMainPlot](#), [dataCustomizedPlot](#)

**Examples**

```
data("fullProcessedProfile", package="PhyloProfile")
plotDf <- dataMainPlot(fullProcessedProfile)
plotParameter <- list(
  "xAxis" = "taxa",
  "var1ID" = "FAS_FW",
  "var2ID" = "FAS_BW",
  "lowColorVar1" = "#FF8C00",
  "highColorVar1" = "#4682B4",
  "lowColorVar2" = "#FFFFFF",
  "highColorVar2" = "#F0E68C",
  "paraColor" = "#07D000",
  "xSize" = 8,
  "ySize" = 8,
  "legendSize" = 8,
  "mainLegend" = "top",
  "dotZoom" = 0,
```

```

    "xAngle" = 60,
    "guideline" = 0,
    "colorByGroup" = FALSE
  )
  taxonHighlight <- "none"
  rankName <- "class"
  geneHighlight <- "100265at6656"
  highlightProfilePlot(
    plotDf, plotParameter, taxonHighlight, rankName, geneHighlight
  )

```

---

idList	<i>NCBI ID list for experimental data sets</i>
--------	--

---

### Description

Data frame, in which each row contains the complete taxonomy ranks from the lowest systematic level (strain/species) upto the taxonomy root and the corresponding IDs for one taxon in the experimental data sets.

### Usage

```
data(idList)
```

### Format

A data frame with up to 41 columns and 95 rows corresponding to 95 taxa in the 2 experimental data sets

---

mainLongRaw	<i>An example of a raw long input file.</i>
-------------	---

---

### Description

An example of a raw long input file.

### Usage

```
data(mainLongRaw)
```

### Format

A data frame with 168 rows and 5 variables:

- geneID Seed or ortholog group ID, e.g. "100136at6656"
- ncbiID Taxon ID, e.g. "ncbi36329"
- orthoID Ortholog ID, e.g. "100136at6656|PLAF7@36329@1|Q8ILT8|1"
- FAS\_F First additional variable
- FAS\_B Second additional variable

---

mainTaxonomyRank	<i>Get all NCBI taxonomy rank names</i>
------------------	---

---

**Description**

Get all NCBI taxonomy rank names

**Usage**

```
mainTaxonomyRank()
```

**Value**

A list of all available NCBI taxonomy rank names.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**Examples**

```
mainTaxonomyRank()
```

---

pairDomainPlotting	<i>Create architecture plot for a pair of seed and ortholog protein</i>
--------------------	---

---

**Description**

Create architecture plot for a pair of seed and ortholog protein

**Usage**

```
pairDomainPlotting(seed, ortho, seedDf, orthoDf, minStart, maxEnd,
  labelSize, titleSize)
```

**Arguments**

seed	Seed ID
ortho	Ortho ID
seedDf	domain dataframe for seed domains containing the seed ID, ortholog ID, sequence length, feature names, start and end positions, feature weights (optional) and the status to determine if that feature is important for comparison the architecture between 2 proteins* (e.g. seed protein vs ortholog) (optional).
orthoDf	domain dataframe for ortholog domains (same format as seedDf).
minStart	the smallest start position of all domains
maxEnd	the highest stop position of all domains
labelSize	lable size. Default = 12.
titleSize	title size. Default = 12.

**Value**

Domain plot of a pair proteins as a arrangeGrob object.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**Examples**

```
## Not run:
seed <- "101621at6656"
ortho <- "101621at6656|AGRPL@224129@0|224129_0:001955|1"
ortho <- gsub("\\|", ":", ortho)
grepID <- paste(seed, "#", ortho, sep = "")
domainFile <- system.file(
  "extdata", "domainFiles/101621at6656.domains",
  package = "PhyloProfile", mustWork = TRUE
)
domainDf <- parseDomainInput(seed, domainFile, "file")
subdomainDf <- domainDf[grepl(grepID, domainDf$seedID), ]
subdomainDf$feature <- as.character(subdomainDf$feature)
orthoDf <- subdomainDf[subdomainDf$orthoID == ortho,]
seedDf <- subdomainDf[subdomainDf$orthoID != ortho,]
minStart <- min(subdomainDf$start)
maxEnd <- max(c(subdomainDf$end, subdomainDf$length))
g <- pairDomainPlotting(seed, ortho, seedDf, orthoDf, minStart, maxEnd, 9, 9)
grid::grid.draw(g)

## End(Not run)
```

---

parseDomainInput	<i>Parse domain input file</i>
------------------	--------------------------------

---

**Description**

Get all domain annotations for one seed protein IDs.

**Usage**

```
parseDomainInput(seed = NULL, inputFile = NULL, type = "file")
```

**Arguments**

seed	seed protein ID
inputFile	name of input file (file name or path to folder contains individual domain files)
type	type of data (file" or "folder"). Default = "file".

**Value**

A dataframe for protein domains including seed ID, its orthologs IDs, sequence lengths, feature names, start and end positions, feature weights (optional) and the status to determine if that feature is important for comparison the architecture between 2 proteins\* (e.g. seed protein vs ortholog) (optional).

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[getDomainFolder](#)

**Examples**

```
seed <- "101621at6656"  
inputFile <- system.file(  
  "extdata", "domainFiles/101621at6656.domains",  
  package = "PhyloProfile", mustWork = TRUE  
)  
type <- "file"  
parseDomainInput(seed, inputFile, type)
```

---

parseInfoProfile

*Parsing info for phylogenetic profiles*

---

**Description**

Creating main dataframe for the input phylogenetic profiles based on selected input taxonomy level (e.g. strain, species) and reference taxon. The output contains the number of paralogs, percentage of species presence in each supertaxon, and the max/min/mean/median of VAR1 and VAR2.

**Usage**

```
parseInfoProfile(inputDf, sortedInputTaxa, var1AggregateBy = "max",  
  var2AggregateBy = "max")
```

**Arguments**

inputDf	input profiles in long format
sortedInputTaxa	sorted taxonomy data for the input taxa (check sortInputTaxa())
var1AggregateBy	aggregate method for VAR1 (max, min, mean or median), applied for calculating var1 of supertaxa. Default = "max".
var2AggregateBy	aggregate method for VAR2 (max, min, mean or median), applied for calculating var2 of supertaxa. Default = "max".

**Value**

A dataframe contains all info for the input phylogenetic profiles. This full processed profile that is required for several profiling analyses e.g. estimation of gene age (?estimateGeneAge) or identification of core gene (?getCoreGene).

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[createLongMatrix](#), [sortInputTaxa](#), [calcPresSpec](#), [mainLongRaw](#)

**Examples**

```
data("mainLongRaw", package="PhyloProfile")
taxonIDs <- getInputTaxaID(mainLongRaw)
sortedInputTaxa <- sortInputTaxa(
  taxonIDs, "class", "Mammalia", NULL
)
var1AggregateBy <- "max"
var2AggregateBy <- "mean"
parseInfoProfile(
  mainLongRaw, sortedInputTaxa, var1AggregateBy, var2AggregateBy
)
```

---

ppTaxonomyMatrix      *An example of a taxonomy matrix.*

---

**Description**

An example of a taxonomy matrix.

**Usage**

```
data(ppTaxonomyMatrix)
```

**Format**

A data frame with 10 rows and 162 variables:

- abbrName e.g. "ncbi10090"
- ncbiID e.g. "10090"
- fullName e.g. "Mus musculus"
- strain e.g. "10090" ...

---

ppTree      *An example of a taxonomy tree in newick format.*

---

**Description**

An example of a taxonomy tree in newick format.

**Usage**

```
data(ppTree)
```

**Format**

A data frame with only one entry

**V1** tree in newick format

---

processNcbiTaxonomy    *Pre-processing NCBI taxonomy data*

---

### Description

Download NCBI taxonomy database and parse information that are needed for PhyloProfile, including taxon IDs, their scientific names, systematic ranks, and parent (next higher) rank IDs.

### Usage

```
processNcbiTaxonomy()
```

### Value

A dataframe contains NCBI taxon IDs, taxon names, taxon ranks and the next higher taxon IDs (parent's IDs) of all taxa in the NCBI taxonomy database.

### Author(s)

Vinh Tran tran@bio.uni-frankfurt.de

### Examples

```
?processNcbiTaxonomy
## Not run:
preProcessedTaxonomy <- processNcbiTaxonomy()
# save to text (tab-delimited) file
write.table(
  preProcessedTaxonomy,
  file = "preProcessedTaxonomy.txt",
  col.names = TRUE,
  row.names = FALSE,
  quote = FALSE,
  sep = "\t"
)
# save to rdata file
save(
  preProcessedTaxonomy, file = "preProcessedTaxonomy.RData", compress='xz'
)

## End(Not run)
```

---

profileWithTaxonomy    *An example of a raw long input file together with the taxonomy info.*

---

### Description

An example of a raw long input file together with the taxonomy info.

### Usage

```
data(profileWithTaxonomy)
```



**Format**

A data frame with 20 rows and 12 variables:

- geneID Seed or ortholog group ID, e.g. "OG\_1017"
- ncbiID Taxon ID, e.g. "ncbi176299"
- orthoID Ortholog ID, e.g. "A.fabrum@176299@1582"
- var1 First additional variable
- var2 Second additional variable
- paralog Number of co-orthologs in the current taxon
- abbrName e.g. "ncbi176299"
- taxonID Taxon ID, e.g. "176299"
- fullName Full taxon name, e.g. "Agrobacterium fabrum str. C58"
- supertaxonID Supertaxon ID (only different than ncbiID in case working with higher taxonomy rank than input's)
- supertaxon Name of the corresponding supertaxon
- rank Rank of the supertaxon

---

qualitativeColours      *Create qualitative colours*

---

**Description**

Create qualitative colours

**Usage**

```
qualitativeColours(n, light = FALSE)
```

**Arguments**

n	number of colors
light	light colors TRUE or FALSE

**Value**

list of n different colors

**Source**

Modified based on <https://gist.github.com/peterk87/6011397>

**Examples**

```
## Not run:  
qualitativeColours(5)  
  
## End(Not run)
```

---

rankIndexing	<i>Indexing all available ranks (including norank)</i>
--------------	--

---

**Description**

Indexing all available ranks (including norank)

**Usage**

```
rankIndexing(rankListFile = NULL)
```

**Arguments**

rankListFile    Input file, where each row is a rank list of a taxon (see rankListFile in example)

**Value**

A dataframe containing a list of all possible ranks and their indexed values.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**Examples**

```
## Not run:
rankListFile <- system.file(
  "extdata", "data/rankList.txt", package = "PhyloProfile", mustWork = TRUE
)
rankIndexing(rankListFile)

## End(Not run)
```

---

rankList	<i>NCBI rank list for experimental data sets</i>
----------	--

---

**Description**

Data frame, in which each row contains the complete taxonomy ranks from the lowest systematic level (strain/species) upto the taxonomy root for one taxon in the experimental data sets.

**Usage**

```
data(rankList)
```

**Format**

A data frame with up to 41 columns and 95 rows corresponding to 95 taxa in the 2 experimental data sets

---

reduceProfile	<i>Reduce the full processed profile data into supertaxon level</i>
---------------	---

---

**Description**

Reduce data of the processed phylogenetic profiles from input taxonomy rank into supertaxon level (e.g. from species to phylum)

**Usage**

```
reduceProfile(fullProfile)
```

**Arguments**

fullProfile	dataframe contains the full processed profiles (see ?parseInfoProfile and ?full-ProcessedProfile)
-------------	---

**Value**

A reduced dataframe contains only profile data for the selected supertaxon rank. This dataframe contains only supertaxa and their value (

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[parseInfoProfile](#) for creating a full processed profile dataframe, [fullProcessedProfile](#) for a demo full processed profile dataframe

**Examples**

```
data("fullProcessedProfile", package="PhyloProfile")
reduceProfile(fullProcessedProfile)
```

---

runPhyloProfile	<i>Run PhyloProfile app</i>
-----------------	-----------------------------

---

**Description**

Run PhyloProfile app

**Usage**

```
runPhyloProfile()
```

**Value**

A shiny application - GUI version of PhyloProfile

**Examples**

```
?runPhyloProfile
## Not run:
runPhyloProfile()

## End(Not run)
```

---

singleDomainPlotting *Create architecture plot for a single protein*

---

**Description**

Create architecture plot for a single protein

**Usage**

```
singleDomainPlotting(df, geneID = "GeneID", sep = "|", labelSize = 12,
  titleSize = 12, minStart = NULL, maxEnd = NULL, colorScheme)
```

**Arguments**

df	domain dataframe for plotting containing the seed ID, ortholog ID, ortholog sequence length, feature names, start and end positions, feature weights (optional) and the status to determine if that feature is important for comparison the architecture between 2 proteins* (e.g. seed protein vs ortholog) (optional).
geneID	ID of seed or orthologous protein
sep	separate indicator for title. Default = " ".
labelSize	label size. Default = 12.
titleSize	title size. Default = 12.
minStart	the smallest start position of all domains
maxEnd	the highest stop position of all domains
colorScheme	color scheme for all domain types

**Value**

Domain plot of a single protein as a ggplot object.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[getQualColForVector](#), [parseDomainInput](#)

**Examples**

```

## Not run:
# get domain data
domainFile <- system.file(
  "extdata", "domainFiles/101621at6656.domains",
  package = "PhyloProfile", mustWork = TRUE
)
domainDf <- parseDomainInput(seedID, domainFile, "file")
df <- domainDf[
  domainDf$orthoID == "101621at6656|AGRPL@224129@0|224129_0:001955|1",]
# create color scheme for all domain types
allFeatures <- levels(as.factor(df$feature))
allColors <- getQualColForVector(allFeatures)
colorScheme <- structure(
  allColors,
  .Names = allFeatures
)
# other parameters
geneID <- "AGRPL@224129@0|224129_0:001955|1"
sep <- "|"
labelSize <- 9
titleSize <- 9
minStart <- min(df$start)
maxEnd <- max(df$end)
# do plotting
singleDomainPlotting(
  df,
  geneID,
  sep,
  labelSize, titleSize,
  minStart, maxEnd,
  colorScheme
)
## End(Not run)

```

---

sortDomains

*Sort one domain dataframe based on the other domain dataframe*


---

**Description**

Sort domain dataframe of one protein (either seed or ortholog) based on the dataframe of the its paired protein, in order to bring the common domain feature in the same order which make it easy for comparing.

**Usage**

```
sortDomains(seedDf, orthoDf)
```

**Arguments**

seedDf	data of seed protein
orthoDf	data of ortholog protein

**Value**

Dataframe contains sorted domain list.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**Examples**

```
## Not run:
# get domain data
domainFile <- system.file(
  "extdata", "domainFiles/101621at6656.domains",
  package = "PhyloProfile", mustWork = TRUE
)
domainDf <- parseDomainInput(seedID, domainFile, "file")
# get seedDf and orthoDf
subDf <- domainDf[
  domainDf$seedID ==
  "101621at6656#101621at6656:AGRPL0224129@0:224129_0:001955:1",]
orthoDf <- subDf[subDf$orthoID == "101621at6656:DROME@7227@1:Q9VG04",]
seedDf <- subDf[subDf$orthoID != "101621at6656:DROME@7227@1:Q9VG04",]
# sort
sortDomains(seedDf, orthoDf)

## End(Not run)
```

---

sortInputTaxa

*Sort list of (super)taxa based on a selected reference (super)taxon*

---

**Description**

Sort list of (super)taxa based on a selected reference (super)taxon

**Usage**

```
sortInputTaxa(taxonIDs = NULL, rankName, refTaxon = NULL,
  taxaTree = NULL)
```

**Arguments**

taxonIDs	list of taxon IDs (e.g.: ncbi1234, ncbi9999, ...). Default = NULL.
rankName	working taxonomy rank (e.g. "species", "phylum",...)
refTaxon	selected reference taxon. Default = NULL.
taxaTree	taxonomy tree for the input taxa (optional). Default = NULL.

**Value**

A taxonomy matrix for the input taxa ordered by the selected reference taxon. This matrix is sorted either based on the NCBI taxonomy info, or based on an user-defined taxonomy tree (if provided).

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[getNamelist](#), [getTaxonomyMatrix](#), [createRootedTree](#), [sortTaxaFromTree](#), [getInputTaxaName](#), [getInputTaxaID](#), [createLongMatrix](#)

**Examples**

```
taxonIDs <- c(
  "ncbi10116", "ncbi123851", "ncbi3702", "ncbi13616", "ncbi9606"
)
sortInputTaxa(taxonIDs, "species", "Homo sapiens", NULL)
```

---

sortTaxaFromTree

*Get sorted supertaxon list based on a rooted taxonomy tree*

---

**Description**

Get sorted supertaxon list based on a rooted taxonomy tree

**Usage**

```
sortTaxaFromTree(tree)
```

**Arguments**

tree            an "phylo" object for a rooted taxonomy tree

**Value**

A list of sorted taxa obtained the input taxonomy tree.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[ppTaxonomyMatrix](#) for a demo taxonomy matrix data

**Examples**

```
data("ppTaxonomyMatrix", package = "PhyloProfile")
# prepare matrix for calculating distances
distDf <- subset(ppTaxonomyMatrix, select = -c(ncbiID, fullName))
row.names(distDf) <- distDf$abbrName
distDf <- distDf[, -1]
# create taxonomy tree rooted by ncbi10090
rootedTree <- createRootedTree(distDf, "ncbi10090")
# get taxon list sorted from tree
sortTaxaFromTree(rootedTree)
```

---

taxa2dist	<i>taxa2dist</i>
-----------	------------------

---

**Description**

taxa2dist

**Usage**

```
taxa2dist(x, varstep = FALSE, check = TRUE, labels)
```

**Arguments**

x	taxa matrix
varstep	var-step
check	check
labels	labels

**Value**

a distance matrix

**Author(s)**

function from taxize library

---

taxonNamesReduced	<i>NCBI Taxonomy reduced data set</i>
-------------------	---------------------------------------

---

**Description**

A list of NCBI taxonomy info (including taxon IDs, taxon names, their systematic taxonomy rank and IDs of their next rank - parent IDs) for 95 taxa in two experimental sets included in PhyloProfileData package.

**Usage**

```
data(taxonNamesReduced)
```

**Format**

A data frame with 4 columns:

- ncbiID e.g. "10090"
- fullName e.g. "Mus musculus"
- rank e.g. "species"
- parentID e.g. "862507"



---

taxonomyMatrix	<i>Taxonomy matrix for experimental data sets</i>
----------------	---

---

**Description**

Data frame containing the fully aligned taxonomy IDs of 95 taxa in the experimental data sets. By taking into account both the defined ranks (e.g. strain, This data is used for clustering and then creating a taxon tree. It is used also for cross-linking between different taxonomy ranks within a taxon.

**Usage**

```
data(taxonomyMatrix)
```

**Format**

A data frame with up to 147 columns and 95 rows corresponding to 95 taxa in the 2 experimental data sets

---

taxonomyTableCreator	<i>Align NCBI taxonomy IDs of list of taxa into a sorted rank list.</i>
----------------------	---

---

**Description**

Align NCBI taxonomy IDs of list of taxa into a sorted rank list.

**Usage**

```
taxonomyTableCreator(idListFile = NULL, rankListFile = NULL)
```

**Arguments**

idListFile      a text file whose each row is a rank+ID list of a taxon (see idListFile in example)  
rankListFile    a text file whose each row is a rank list of a taxon (see rankListFile in example)

**Value**

An aligned taxonomy dataframe which contains all the available taxonomy ranks from the id and rank list file. This dataframe can be used for creating a well resolved taxonomy tree (see ?create-RootedTree) and sorting taxa based on a selected reference taxon (see ?sortInputTaxa).

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[rankIndexing](#), [createRootedTree](#), [sortInputTaxa](#)

**Examples**

```
idListFile <- system.file(
  "extdata", "data/idList.txt", package = "PhyloProfile", mustWork = TRUE
)
rankListFile <- system.file(
  "extdata", "data/rankList.txt", package = "PhyloProfile", mustWork = TRUE
)
taxonomyTableCreator(idListFile, rankListFile)
```

---

<code>varDistTaxPlot</code>	<i>Create variable distribution comparison plot</i>
-----------------------------	---

---

**Description**

Create variable distribution plots between 2 groups of taxa for a selected gene.

**Usage**

```
varDistTaxPlot(data, plotParameters)
```

**Arguments**

<code>data</code>	dataframe for plotting. Last column indicates what type of taxon group (in- or out-group). The first (or first 2) column contains values of the variables. See <code>?dataVarDistTaxGroup</code>
<code>plotParameters</code>	plot parameters, including size of x-axis, y-axis, legend and title; position of legend ("right", "bottom" or "none"); mean/median point; names of in-group and out-group; and plot title. NOTE: Leave blank or NULL to use default values.

**Value**

Distribution plots as a grob (gtable) object. Use `grid.draw` to plot.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**See Also**

[dataVarDistTaxGroup](#)

**Examples**

```
data("mainLongRaw", package="PhyloProfile")
data <- mainLongRaw
inGroup <- c("ncbi9606", "ncbi10116")
variable <- colnames(data)[c(4, 5)]
plotDf <- dataVarDistTaxGroup(data, inGroup, "101621at6656", variable)
plotParameters <- list(
  "xSize" = 12,
  "ySize" = 12,
  "titleSize" = 15,
  "legendSize" = 12,
```

```
    "legendPosition" = "right",
    "mValue" = "mean",
    "inGroupName" = "In-group",
    "outGroupName" = "Out-group",
    "title" = "101621at6656"
  )
  g <- varDistTaxPlot(plotDf, plotParameters)
  grid::grid.draw(g)
```

---

wideToLong

*Transform input file in wide matrix into long matrix format*

---

## Description

Transform input file in wide matrix into long matrix format

## Usage

```
wideToLong(inputFile = NULL)
```

## Arguments

inputFile      input file in wide matrix format

## Value

A data frame of input data in long-format containing seed gene IDs ( or orthologous group IDs), their orthologous proteins together with the corresponding taxonomy IDs and values of (up to) two additional variables.

## Author(s)

Vinh Tran tran@bio.uni-frankfurt.de

## Examples

```
inputFile <- system.file(
  "extdata", "test.main.wide", package = "PhyloProfile", mustWork = TRUE
)
wideToLong(inputFile)
```

xmlParser

*Parse orthoXML input file*

---

**Description**

Parse orthoXML input file

**Usage**

```
xmlParser(inputFile = NULL)
```

**Arguments**

inputFile      input file in xml format

**Value**

A data frame of input data in long-format containing seed gene IDs ( or orthologous group IDs), their orthologous proteins together with the corresponding taxonomy IDs and values of (up to) two additional variables.

**Author(s)**

Vinh Tran tran@bio.uni-frankfurt.de

**Examples**

```
inputFile <- system.file(  
  "extdata", "test.main.xml", package = "PhyloProfile", mustWork = TRUE  
)  
xmlParser(inputFile)
```

# Index

calcPresSpec, [3](#), [55](#)  
checkInputValidity, [4](#)  
checkNewick, [5](#)  
checkOmaID, [4](#), [5](#)  
clusterDataDend, [6](#), [34](#)  
compareMedianTaxonGroups, [7](#)  
compareTaxonGroups, [8](#)  
createArchiPlot, [9](#)  
createGeneAgePlot, [10](#)  
createLongMatrix, [10](#), [18](#), [19](#), [26](#), [40](#), [55](#), [63](#)  
createPercentageDistributionData, [11](#),  
[14](#)  
createProfileFromOma, [12](#)  
createRootedTree, [13](#), [63](#), [65](#)  
createVarDistPlot, [13](#)  
createVariableDistributionData, [14](#), [15](#),  
[16](#)  
createVariableDistributionDataSubset,  
[14](#), [15](#)

dataCustomizedPlot, [16](#), [49](#), [50](#)  
dataFeatureTaxGroup, [17](#), [22](#)  
dataMainPlot, [18](#), [49](#), [50](#)  
dataVarDistTaxGroup, [19](#), [66](#)  
distributionTest, [20](#)

estimateGeneAge, [10](#), [20](#), [28](#)

fastaParser, [11](#), [21](#)  
featureDistTaxPlot, [22](#)  
filterProfileData, [17](#), [19](#), [23](#), [26](#)  
finalProcessedProfile, [24](#)  
fromInputToProfile, [24](#), [25](#), [33](#)  
fullProcessedProfile, [16](#), [21](#), [24](#), [27](#), [59](#)

geneAgePlotDf, [10](#), [28](#)  
generateSinglePlot, [28](#)  
getAllDomainsOma, [29](#)  
getAllFastaOma, [30](#)  
getCommonAncestor, [31](#)  
getCoreGene, [31](#)  
getDataClustering, [6](#), [32](#), [35](#)  
getDataForOneOma, [12](#), [30](#), [33](#), [44](#)  
getDendrogram, [34](#)  
getDistanceMatrix, [6](#), [35](#)  
getDomainFolder, [36](#), [54](#)  
getFastaFromFasInput, [36](#)  
getFastaFromFile, [37](#)  
getFastaFromFolder, [38](#)  
getIDsRank, [39](#)  
getInputTaxaID, [5](#), [26](#), [39](#), [40](#), [63](#)  
getInputTaxaName, [26](#), [40](#), [63](#)  
getNameList, [21](#), [40](#), [41](#), [63](#)  
getOmaDataForOneOrtholog, [41](#)  
getOmaDomainFromURL, [42](#)  
getOmaMembers, [42](#)  
getQualColForVector, [9](#), [43](#), [60](#)  
getSelectedFastaOma, [44](#)  
getSelectedTaxonNames, [44](#)  
getTaxonomyInfo, [45](#)  
getTaxonomyMatrix, [13](#), [21](#), [46](#), [63](#)  
getTaxonomyRanks, [47](#)  
gridArrangeSharedLegend, [47](#)

hclust, [6](#)  
heatmapPlotting, [48](#)  
highlightProfilePlot, [49](#)

idList, [51](#)

mainLongRaw, [12](#), [14–16](#), [38](#), [40](#), [51](#), [55](#)  
mainTaxonomyRank, [52](#)

pairDomainPlotting, [52](#)  
parseDomainInput, [9](#), [18](#), [53](#), [60](#)  
parseInfoProfile, [16](#), [21](#), [24](#), [26](#), [32](#), [54](#), [59](#)  
ppTaxonomyMatrix, [13](#), [55](#), [63](#)  
ppTree, [5](#), [55](#)  
processNcbiTaxonomy, [56](#)  
profileWithTaxonomy, [4](#), [56](#)

qualitativeColours, [43](#), [57](#)

rankIndexing, [58](#), [65](#)  
rankList, [58](#)  
reduceProfile, [24](#), [26](#), [59](#)  
runPhyloProfile, [59](#)

singleDomainPlotting, [9](#), [60](#)

sortDomains, [9](#), [61](#)  
sortInputTaxa, [26](#), [55](#), [62](#), [65](#)  
sortTaxaFromTree, [63](#), [63](#)

taxa2dist, [13](#), [64](#)  
taxonNamesReduced, [64](#)  
taxonomyMatrix, [65](#)  
taxonomyTableCreator, [65](#)

varDistTaxPlot, [66](#)

wideToLong, [11](#), [67](#)

xmlParser, [11](#), [68](#)