

Using the PAnnBuilder Package

Hong Li^{‡*}

October 28, 2009

[‡]Key Lab of Systems Biology
Shanghai Institutes for Biological Sciences
Chinese Academy of Sciences, P. R. China

Contents

1	Overview	2
2	Getting Started	2
2.1	Requirements	2
2.2	Installation	3
2.3	Public Data Sources	3
2.4	Annotation Packages Produced by PAnnBuilder	4
2.4.1	Packages produced by PAnnBuilder	4
2.4.2	Using annotation data package	5
3	Function Description	13
3.1	Getting URL and Version	13
3.2	Parsing and Writing Data	13
3.2.1	Employing perl program to parse data	13
3.2.2	Writing data using R	14
3.3	Writing Help Documents	15
3.4	Building Data Packages	15
4	Building Annotation Data Packages	16
5	Session Information	18

*sysptm@gmail.com

1 Overview

In genomic era, genome-scale experiments and data analyses require genes to be annotated from different sources, an R Development Core Team (2008) package `AnnBuilder` was developed for this purpose Zhang et al. (2003). In post genomic era, advances in proteomics highlight the urgency of understanding protein language Jensen (2006). However, relative to genes, `AnnBuilder` is limited in protein annotation due to the complexity of proteins caused by 3-D structure, alternative splicing, modification, dynamic location and other features. The package `PAnnBuilder` focuses on assembling proteomic annotation data, which should be very useful for proteomic data interpretation. It not only inherits the good features of `AnnBuilder` such as automatically parsing protein annotation data and building R packages from selected sources, but also emphasizes specific features of proteins. `PAnnBuilder` currently supports 16 databases involving diverse aspects of proteomics, such as protein primary data, protein domain/family, subcellular location, protein interaction, post-translational modifications, body fluid proteomics, homolog/ortholog groups and so on. Additionally, `PAnnBuilder` allows annotation to unknown proteins based on sequence similarities to other well-annotated proteins. To extend the use of `PAnnBuilder`, 54 standard R annotation packages are produced from main protein databases, which are freely available for download via `biocLite`.

2 Getting Started

2.1 Requirements

`PAnnBuilder` requires the support from the following items:

1. For `PAnnBuilder` $\geq 1.5.0$, R $\geq 2.7.0$ is needed for building SQLite-based package. Dependent R packages are needed to be installed: `methods`, `utils`, `RSQLite`, `Biobase` ($\geq 1.17.0$), `AnnotationDbi` ($\geq 1.3.12$). If you use the installation script "biocLite.R" to install `PAnnBuilder` ($\geq 1.5.0$), it will automatically check these dependent packages and install missing packages from CRAN or Bioconductor.
2. `Rtools` is needed for Windows user. The matched version of `Rtools` with R can be downloaded via <http://www.murdoch-sutherland.com/Rtools/>.
3. Perl is required to parse the rather large annotation source data files. It can be downloaded from <http://www.perl.com/download.csp>.
4. Program `formatdb` and `BLASTP` is needed for function `crossBuilder`, `crossBuilder_DB`, `pSeqBuilder`, `pSeqBuilder_DB`.
`BLAST` can be downloaded from <http://www.ncbi.nlm.nih.gov/BLAST/download.shtml>.

Note: It is better to have enough space for the temporary directory. The path of the per-session temporary directory can be acquired by:

```
> tempdir()
[1] "/tmp/Rtmp4nlp3x"
```

2.2 Installation

The biocLite script is used to install PAnnBuilder from within R:

```
> source("http://bioconductor.org/biocLite.R")
> biocLite("PAnnBuilder")

> library(PAnnBuilder)
```

Note: Web Connection is needed to install PAnnBuilder and its depended packages.

2.3 Public Data Sources

PAnnBuilder parses proteomics annotation data from public sources and build R annotation packages. It also provides convenient functions to access these sources. For example, you can get all supported databases for "Homo sapiens" by:

```
> library(PAnnBuilder)
> getALLUrl("Homo sapiens")
> getALLBuilt("Homo sapiens")
```

Detail description of all supported public data sources in PAnnBuilder are as follows:

UniProt The data `ftp://ftp.ebi.ac.uk/pub/databases/uniprot/knowledgebase` will be used to map protein UniProt identifiers to diverse annotation available in UniProt database.

IPI The data `ftp://ftp.ebi.ac.uk/pub/databases/IPI/current` will be used to map protein IPI identifiers to diverse annotation available in International Protein Index database.

RefSeq The data `ftp://ftp.ncbi.nih.gov/refseq` will be used to map protein RefSeq identifiers to diverse annotation available in NCBI RefSeq database.

Entrez Gene The data `ftp://ftp.ncbi.nih.gov/gene/DATA` will be used to annotate genes after the Entrez Gene identifiers have been obtained.

Gene Ontology The data `ftp://ftp.ebi.ac.uk/pub/databases/G0/goa/proteomes` and `http://archive.geneontology.org/latest-termdb` will be used to obtain gene ontology information.

KEGG Some data at `ftp://ftp.genome.jp/pub/kegg` will be used to obtain pathway information.

HomoloGene A data file provided by <ftp://ftp.ncbi.nlm.nih.gov/pub/HomoloGene/current/> will be used to extract mappings between GeneID/ProteinGI and HomoloGene ids.

InParanoid The data <http://inparanoid.sbc.su.se> will be used to obtain ortholog protein groups between two organisms.

Gene Interaction The data <ftp://ftp.ncbi.nih.gov/gene/GeneRIF> will be used to extract protein-protein interactions between Entrez GeneID or Protein RefSeq ids.

IntAct The data <ftp://ftp.ebi.ac.uk/pub/databases/intact/current/psimitab> will be used to extract protein-protein interactions between UniProt protein accession numbers.

MPPI The data <http://mips.gsf.de/proj/ppi/data/mppi.gz> will be used to extract protein-protein interactions between UniProt protein accession numbers.

3DID The data <http://gatealoy.pcb.ub.es/3did/download> will be used to extract domain-domain interactions between Pfam domain identifiers.

DOMINE The data <http://domine.utdallas.edu> will be used to extract domain-domain interactions between Pfam domain identifiers.

DBSubLoc The data <http://www.bioinfo.tsinghua.edu.cn/~guotao> will be used to obtain subcellular localization for protein from SWISS-PROT and PIR database.

BaCelLo The data <http://gpcr2.biocomp.unibo.it/bacello> will be used to map SwissProt eukaryotes protein identifiers to subcellular localization.

SCOP The data <http://scop.mrc-lmb.cam.ac.uk/scop> will be used to map PDB structure identifiers to SCOP domain identifiers.

PeptideAtlas The data <http://www.peptideatlas.org/builds> will be used to obtain experimentally identified peptides and their coordinates on chromosomes.

SysPTM The data <http://www.biosino.org/papers/SysPTM> will be used to obtain protein post-translational modifications information.

Sys-BodyFluid The data <http://www.biosino.org/papers/Sys-BodyFluid> will be used to map IPI protein identifiers to body fluids.

2.4 Annotation Packages Produced by PAnnBuilder

2.4.1 Packages produced by PAnnBuilder

PAnnBuilder has powerful ability to build R package for assembling proteome annotation. However, the process of building new package may be time-consuming because of the downloading and parsing of large data files. To make PAnnBuilder useful for any users, we have

built many frequently used annotation packages in advance. These pre-built packages can be downloaded via `biocLite`.

`PAnnBuilder` produces SQLite-based R packages using the `"*Builder_DB"` functions (see Table 1). Each SQLite-based annotation package contains a number of `AnnDbBimap` objects.

The pre-built packages provide a quick start for R beginners. If one wants to analyze protein set in Human IPI database, the quickest way is to download and use the SQLite-based package `"org.Hs.ipi.db"`. However, if the package one wants has not been built or a new-version database is released, new package should be built using functions in `PAnnBuilder` (See Section 4 for methods of building annotation packages).

2.4.2 Using annotation data package

SQLite-based `*.db` packages are capable of flexible data queries, reverse mapping, and data filtering. Vignette of `"AnnotationDbi"` detailedly described how to use SQLite based packages (<http://www.bioconductor.org/packages/release/bioc/vignettes/AnnotationDbi/inst/doc/AnnotationDbi.pdf>). Here package `"org.Hs.ipi.db"` produced by `PAnnBuilder` is illustrated as an example in the following code chunk.

1. Install and load `*.db` annotation package.

```
> biocLite("org.Hs.ipi.db")
> library(org.Hs.ipi.db)
```

2. Browse data in the package.

```
> `?`(org.Hs.ipiGENEID)
```

3. Convert object into a `"list"` object, and get values by index or name.

```
> xx <- as.list(org.Hs.ipiGENEID)
> xx[!is.na(xx)][1:3]
```

```
$IPI00000001.2
[1] " 6780"
```

```
$IPI00000005.1
[1] " 4893"
```

```
$IPI00000006.1
[1] " 3265"
```

```
> xx[["IPI00792103.1"]]
[1] " 28957"
```

Table 1: SQLite-based Annotation Packages Produced by PAnnBuilder.

Description	R Function	Source	Organism	Package
complete and canonical annotation for all proteins of a specific organism, including protein description, Entrez gene identifier, KEGG pathway, gene ontology, domain, and so on.	pBaseBuilder_DB	IPI Swiss-Prot RefSeq	Homo sapiens Mus musculus Rattus norvegicus Homo sapiens Mus musculus Rattus norvegicus Homo sapiens Mus musculus Rattus norvegicus	org.Hs.ipi.db org.Mm.ipi.db org.Rn.ipi.db org.Hs.sp.db org.Mm.sp.db org.Rn.sp.db org.Hs.ref.db org.Mm.ref.db org.Rn.ref.db
protein identifier mapping	crossBuilder_DB	Swiss-Prot, IPI, RefSeq	Homo sapiens Mus musculus Rattus norvegicus	org.Hs.cross.db org.Mm.cross.db org.Rn.cross.db
protein-protein or domain-domain interaction data	intBuilder_DB	IntAct NCBI Gene MPPI 3did Domine		int.intact.db int.geneint.db int.mppi.db int.did.db int.domine.db
protein subcell location	subcellBuilder_DB	BaCelLo DBSubLoc		sc.bacello.db sc.dbsubloc.db
protein structure classification	scopBuilder_DB	SCOP		scop.db
protein post-translational modification	ptmBuilder_DB	SysPTM		sysptm.db
body fluid protein	bfBuilder_DB	Sys-BodyFluid	Homo sapiens	org.Hs.bf.db
homolog protein group	HomoloGeneBuilder_DB	HomoloGene		homolog.db
ortholog protein group	InParanoidBuilder_DB	InParanoid	Homo sapiens, Mus musculus	org.HsMm.ortholog.db
peptides identified by mass spectrometry	PeptideAtlasBuilder_DB	PeptideAtlas	Homo sapiens	org.Hs.pep.db
gene ontology	GOABuilder_DB	GOA	Homo sapiens	org.Hs.goa.db
identifier and name	dNameBuilder_DB			dName.db

4. Specific utilities for SQLite-based *.db packages.

```
> toTable(org.Hs.ipiPATH[1:3])
```

	ipi_id	path_id
1	IPI00000005.1	hsa04010
2	IPI00000005.1	hsa04012
3	IPI00000005.1	hsa04360
4	IPI00000005.1	hsa04370
5	IPI00000005.1	hsa04530
6	IPI00000005.1	hsa04540
7	IPI00000005.1	hsa04650
8	IPI00000005.1	hsa04660
9	IPI00000005.1	hsa04662
10	IPI00000005.1	hsa04664
11	IPI00000005.1	hsa04720
12	IPI00000005.1	hsa04730
13	IPI00000005.1	hsa04810
14	IPI00000005.1	hsa04910
15	IPI00000005.1	hsa04912
16	IPI00000005.1	hsa04916
17	IPI00000005.1	hsa05211
18	IPI00000005.1	hsa05213
19	IPI00000005.1	hsa05214
20	IPI00000005.1	hsa05215
21	IPI00000005.1	hsa05216
22	IPI00000005.1	hsa05218
23	IPI00000005.1	hsa05219
24	IPI00000005.1	hsa05220
25	IPI00000005.1	hsa05221
26	IPI00000005.1	hsa05223
27	IPI00000006.1	hsa04010
28	IPI00000006.1	hsa04012
29	IPI00000006.1	hsa04360
30	IPI00000006.1	hsa04370
31	IPI00000006.1	hsa04510
32	IPI00000006.1	hsa04530
33	IPI00000006.1	hsa04540
34	IPI00000006.1	hsa04650
35	IPI00000006.1	hsa04660
36	IPI00000006.1	hsa04662
37	IPI00000006.1	hsa04664
38	IPI00000006.1	hsa04720

```
39 IPI00000006.1 hsa04730
40 IPI00000006.1 hsa04810
41 IPI00000006.1 hsa04910
42 IPI00000006.1 hsa04912
43 IPI00000006.1 hsa04916
44 IPI00000006.1 hsa05211
45 IPI00000006.1 hsa05213
46 IPI00000006.1 hsa05214
47 IPI00000006.1 hsa05215
48 IPI00000006.1 hsa05216
49 IPI00000006.1 hsa05218
50 IPI00000006.1 hsa05219
51 IPI00000006.1 hsa05220
52 IPI00000006.1 hsa05221
53 IPI00000006.1 hsa05223
54 IPI00000023.4 hsa04080
```

```
> tmp1 <- revmap(org.Hs.ipiPATH)
> class(tmp1)
```

```
[1] "AnnDbBimap"
attr(,"package")
[1] "AnnotationDbi"
```

```
> as.list(tmp1)[1]
```

```
$hsa00010
```

```
[1] "IPI00003925.6" "IPI00005118.2" "IPI00006663.1" "IPI00009744.2"
[5] "IPI00009790.1" "IPI00015911.1" "IPI00016768.3" "IPI00018031.1"
[9] "IPI00018246.5" "IPI00019500.3" "IPI00021338.2" "IPI00024087.3"
[13] "IPI00026663.2" "IPI00027165.3" "IPI00027497.5" "IPI00028066.2"
[17] "IPI00028155.2" "IPI00030116.2" "IPI00060200.3" "IPI00073772.5"
[21] "IPI00102864.3" "IPI00103467.4" "IPI00148061.3" "IPI00166751.3"
[25] "IPI00169383.3" "IPI00184775.3" "IPI00215979.3" "IPI00216171.3"
[29] "IPI00216792.2" "IPI00216932.4" "IPI00217872.3" "IPI00217966.8"
[33] "IPI00218407.6" "IPI00218474.5" "IPI00218570.6" "IPI00218767.1"
[37] "IPI00218768.1" "IPI00218896.3" "IPI00218899.6" "IPI00219018.7"
[41] "IPI00219217.3" "IPI00219526.6" "IPI00219568.4" "IPI00219585.4"
[45] "IPI00220271.3" "IPI00220644.8" "IPI00220663.3" "IPI00220665.6"
[49] "IPI00220667.3" "IPI00221234.7" "IPI00244083.4" "IPI00248961.1"
[53] "IPI00292698.3" "IPI00292709.2" "IPI00294380.5" "IPI00294630.2"
[57] "IPI00296183.7" "IPI00299456.4" "IPI00306044.3" "IPI00306301.2"
[61] "IPI00332371.9" "IPI00333619.4" "IPI00374975.2" "IPI00382746.1"
```



```

[65] "IPI00383071.1" "IPI00384116.3" "IPI00384883.3" "IPI00384967.3"
[69] "IPI00385347.1" "IPI00386733.2" "IPI00394758.1" "IPI00413730.4"
[73] "IPI00418262.5" "IPI00451401.3" "IPI00465025.1" "IPI00465028.7"
[77] "IPI00465179.3" "IPI00465248.5" "IPI00465343.3" "IPI00465439.5"
[81] "IPI00473031.6" "IPI00479186.6" "IPI00479877.4" "IPI00513830.1"
[85] "IPI00549564.6" "IPI00549725.6" "IPI00549885.4" "IPI00550486.1"
[89] "IPI00552290.3" "IPI00552617.3" "IPI00553155.1" "IPI00554498.3"
[93] "IPI00555601.1" "IPI00555722.1" "IPI00555728.1" "IPI00556013.1"
[97] "IPI00556284.1" "IPI00604528.3" "IPI00607708.3" "IPI00639981.1"
[101] "IPI00640568.1" "IPI00640862.1" "IPI00642546.1" "IPI00642664.1"
[105] "IPI00642732.2" "IPI00642820.1" "IPI00643166.1" "IPI00643196.1"
[109] "IPI00643575.1" "IPI00644994.1" "IPI00645848.1" "IPI00646468.1"
[113] "IPI00647702.1" "IPI00654709.1" "IPI00743142.2" "IPI00743713.1"
[117] "IPI00746777.3" "IPI00759806.1" "IPI00784216.2" "IPI00787680.1"
[121] "IPI00788640.1" "IPI00788737.1" "IPI00788836.1" "IPI00788938.1"
[125] "IPI00789081.1" "IPI00789134.1" "IPI00789171.1" "IPI00789173.1"
[129] "IPI00789301.1" "IPI00790892.1" "IPI00791170.1" "IPI00791428.1"
[133] "IPI00791564.2" "IPI00791666.1" "IPI00792207.2" "IPI00792375.1"
[137] "IPI00792448.1" "IPI00792655.1" "IPI00792715.1" "IPI00793665.1"
[141] "IPI00793922.1" "IPI00794508.1" "IPI00794605.1" "IPI00794991.1"
[145] "IPI00795075.1" "IPI00795257.2" "IPI00795549.1" "IPI00795622.2"
[149] "IPI00795914.1" "IPI00796111.1" "IPI00796116.1" "IPI00796333.1"
[153] "IPI00796633.1" "IPI00796735.1" "IPI00796823.1" "IPI00796852.1"
[157] "IPI00797038.1" "IPI00797221.7" "IPI00797270.3" "IPI00797580.1"
[161] "IPI00798351.1" "IPI00815786.1" "IPI00815793.1" "IPI00815950.1"
[165] "IPI00816065.1" "IPI00830064.1" "IPI00844159.4" "IPI00847989.2"
[169] "IPI00871353.1" "IPI00872487.1" "IPI00872991.2" "IPI00873455.1"
[173] "IPI00873817.1" "IPI00902547.1" "IPI00902982.1" "IPI00903108.1"
[177] "IPI00909694.1" "IPI00910754.1" "IPI00915933.1" "IPI00916206.1"
[181] "IPI00916990.1" "IPI00916994.1" "IPI00917139.1" "IPI00917193.1"
[185] "IPI00917237.1" "IPI00917473.1" "IPI00917841.1" "IPI00922320.1"
[189] "IPI00922697.1"

```

```

> tmp2 <- reverseSplit(as.list(org.Hs.ipiPATH))
> class(tmp2)

```

```
[1] "list"
```

```
> tmp2[1]
```

```
$hsa00010
```

```

[1] "IPI00003925.6" "IPI00005118.2" "IPI00006663.1" "IPI00009744.2"
[5] "IPI00009790.1" "IPI00015911.1" "IPI00016768.3" "IPI00018031.1"

```

[9] "IPI00018246.5" "IPI00019500.3" "IPI00021338.2" "IPI00024087.3"
[13] "IPI00026663.2" "IPI00027165.3" "IPI00027497.5" "IPI00028066.2"
[17] "IPI00028155.2" "IPI00030116.2" "IPI00060200.3" "IPI00073772.5"
[21] "IPI00102864.3" "IPI00103467.4" "IPI00148061.3" "IPI00166751.3"
[25] "IPI00169383.3" "IPI00184775.3" "IPI00215979.3" "IPI00216171.3"
[29] "IPI00216792.2" "IPI00216932.4" "IPI00217872.3" "IPI00217966.8"
[33] "IPI00218407.6" "IPI00218474.5" "IPI00218570.6" "IPI00218767.1"
[37] "IPI00218768.1" "IPI00218896.3" "IPI00218899.6" "IPI00219018.7"
[41] "IPI00219217.3" "IPI00219526.6" "IPI00219568.4" "IPI00219585.4"
[45] "IPI00220271.3" "IPI00220644.8" "IPI00220663.3" "IPI00220665.6"
[49] "IPI00220667.3" "IPI00221234.7" "IPI00244083.4" "IPI00248961.1"
[53] "IPI00292698.3" "IPI00292709.2" "IPI00294380.5" "IPI00294630.2"
[57] "IPI00296183.7" "IPI00299456.4" "IPI00306044.3" "IPI00306301.2"
[61] "IPI00332371.9" "IPI00333619.4" "IPI00374975.2" "IPI00382746.1"
[65] "IPI00383071.1" "IPI00384116.3" "IPI00384883.3" "IPI00384967.3"
[69] "IPI00385347.1" "IPI00386733.2" "IPI00394758.1" "IPI00413730.4"
[73] "IPI00418262.5" "IPI00451401.3" "IPI00465025.1" "IPI00465028.7"
[77] "IPI00465179.3" "IPI00465248.5" "IPI00465343.3" "IPI00465439.5"
[81] "IPI00473031.6" "IPI00479186.6" "IPI00479877.4" "IPI00513830.1"
[85] "IPI00549564.6" "IPI00549725.6" "IPI00549885.4" "IPI00550486.1"
[89] "IPI00552290.3" "IPI00552617.3" "IPI00553155.1" "IPI00554498.3"
[93] "IPI00555601.1" "IPI00555722.1" "IPI00555728.1" "IPI00556013.1"
[97] "IPI00556284.1" "IPI00604528.3" "IPI00607708.3" "IPI00639981.1"
[101] "IPI00640568.1" "IPI00640862.1" "IPI00642546.1" "IPI00642664.1"
[105] "IPI00642732.2" "IPI00642820.1" "IPI00643166.1" "IPI00643196.1"
[109] "IPI00643575.1" "IPI00644994.1" "IPI00645848.1" "IPI00646468.1"
[113] "IPI00647702.1" "IPI00654709.1" "IPI00743142.2" "IPI00743713.1"
[117] "IPI00746777.3" "IPI00759806.1" "IPI00784216.2" "IPI00787680.1"
[121] "IPI00788640.1" "IPI00788737.1" "IPI00788836.1" "IPI00788938.1"
[125] "IPI00789081.1" "IPI00789134.1" "IPI00789171.1" "IPI00789173.1"
[129] "IPI00789301.1" "IPI00790892.1" "IPI00791170.1" "IPI00791428.1"
[133] "IPI00791564.2" "IPI00791666.1" "IPI00792207.2" "IPI00792375.1"
[137] "IPI00792448.1" "IPI00792655.1" "IPI00792715.1" "IPI00793665.1"
[141] "IPI00793922.1" "IPI00794508.1" "IPI00794605.1" "IPI00794991.1"
[145] "IPI00795075.1" "IPI00795257.2" "IPI00795549.1" "IPI00795622.2"
[149] "IPI00795914.1" "IPI00796111.1" "IPI00796116.1" "IPI00796333.1"
[153] "IPI00796633.1" "IPI00796735.1" "IPI00796823.1" "IPI00796852.1"
[157] "IPI00797038.1" "IPI00797221.7" "IPI00797270.3" "IPI00797580.1"
[161] "IPI00798351.1" "IPI00815786.1" "IPI00815793.1" "IPI00815950.1"
[165] "IPI00816065.1" "IPI00830064.1" "IPI00844159.4" "IPI00847989.2"
[169] "IPI00871353.1" "IPI00872487.1" "IPI00872991.2" "IPI00873455.1"
[173] "IPI00873817.1" "IPI00902547.1" "IPI00902982.1" "IPI00903108.1"

```

[177] "IPI00909694.1" "IPI00910754.1" "IPI00915933.1" "IPI00916206.1"
[181] "IPI00916990.1" "IPI00916994.1" "IPI00917139.1" "IPI00917193.1"
[185] "IPI00917237.1" "IPI00917473.1" "IPI00917841.1" "IPI00922320.1"
[189] "IPI00922697.1"

```

```
> Lkeys(org.Hs.ipiPATH) [1:3]
```

```
[1] "IPI00000005.1" "IPI00000006.1" "IPI00000023.4"
```

```
> Rkeys(org.Hs.ipiPATH) [1:3]
```

```
[1] "hsa00010" "hsa00020" "hsa00030"
```

```
> org.Hs.ipi_dbschema()
```

```

--
-- IPI_DB schema
-- =====
--

```

```

CREATE TABLE basic (
  ipi_id CHAR(13) NOT NULL,           -- IPI Protein Identifier
  ipi_ac VARCHAR(20) NOT NULL,       -- IPI Acession Number in current vers
  len INTEGER NOT NULL,
  mw INTEGER NOT NULL,
  de VARCHAR(255) NOT NULL,         -- Protein Description
  symbol VARCHAR(80) NOT NULL,
  sp_ac VARCHAR(20) NOT NULL,       -- Swiss-Prot Primary Acession Number
  sp_id VARCHAR(20) NOT NULL,       -- Swiss-Prot Protein Identifier
  ref_id VARCHAR(20) NOT NULL,      -- RefSeq accession number
  gi VARCHAR(10) NOT NULL,          -- NCBI Protein GI
  gene_id VARCHAR(10) NOT NULL      -- Entrez Gene ID
  unigene_id VARCHAR(10) NOT NULL,  -- UniGene ID
  kegg_id VARCHAR(20) NOT NULL,     -- KEGG gene ID
  FOREIGN KEY (ipi_id)
);

```

```

CREATE TABLE seq (
  ipi_id CHAR(13) NOT NULL,           -- IPI Protein Identifier
  seq text NOT NULL,                 -- Protein Sequence
  FOREIGN KEY (ipi_id)
);

```

```

CREATE TABLE ipiac (
  ipi_id CHAR(13) NOT NULL,           -- IPI Protein Identifier
  ipi_acs CHAR(11) NOT NULL,         -- IPI Protein Accession Number
  FOREIGN KEY (ipi_id)
);
CREATE TABLE go (
  ipi_id CHAR(13) NOT NULL,           -- IPI Protein Identifier
  go_id CHAR(10) NOT NULL,           -- GO ID
  evidence CHAR(3) NOT NULL,         -- GO evidence code
  ontology CHAR(2) NOT NULL,        -- GO ontology
  FOREIGN KEY (ipi_id)
);
CREATE TABLE path (
  ipi_id CHAR(13) NOT NULL,           -- IPI Protein Identifier
  path_id CHAR(5) NOT NULL,          -- KEGG pathway short ID
  FOREIGN KEY (ipi_id)
};
CREATE TABLE pfam (
  ipi_id CHAR(13) NOT NULL,           -- IPI Protein Identifier
  pfam_id CHAR(7) NULL,              -- Pfam ID
  FOREIGN KEY (ipi_id)
);
CREATE TABLE interpro (
  ipi_id CHAR(13) NOT NULL,           -- IPI Protein Identifier
  interpro_id CHAR(9) NOT NULL,      -- InterPro ID
  FOREIGN KEY (ipi_id)
);
CREATE TABLE prosite (
  ipi_id CHAR(13) NOT NULL,           -- IPI Protein Identifier
  prosite_id CHAR(7) NULL,           -- PROSITE ID
  FOREIGN KEY (ipi_id)
);

```

```

> selectSQL <- paste("SELECT ipi_id, de", "FROM basic", "WHERE de like '%histone%'")
> tmp3 <- dbGetQuery(org.Hs.ipi_dbconn(), selectSQL)
> tmp3[1:3, ]

```

```

      ipi_id                                     de
1 IPI00001830.1 HETEROCHROMATIN-SPECIFIC NONHISTONE PROTEIN (FRAGMENT).
2 IPI00002831.4          HISTONE DEACETYLASE COMPLEX SUBUNIT SAP30L.
3 IPI00002929.2 ISOFORM 3 OF HISTONE-LYSINE N-METHYLTRANSFERASE SUV39H2.

```

3 Function Description

3.1 Getting URL and Version

To download data file from public database, the first step is getting its URL and release/version information. URLs of supported databases are stored in "data/sourceURLs.txt". Following functions are used to get the url:

- `getSrcUrl` - return url according given database.
- `getALLUrl` - return urls for all databases used in PAnnBuilder packages.
- `getSrcBuilt` - return release/version according given database.
- `getALLBuilt` - return release/version information for all databases used in PAnnBuilder packages.

3.2 Parsing and Writing Data

Parsing is a key step to convert original data file to R object. Sometimes R is directly used to parse and write data. But for large data file or complicated data format, perl is firstly employed to quickly process data, and then R function reads the result file into R objects.

3.2.1 Employing perl program to parse data

Segment of perl program is written into file in "inst/scripts". Name and function of these parser files are as follows:

- `spParser` - parse protein data from SwissProt or TrEMBL
- `ipiParser` - parse protein data from IPI
- `refseqParser` - parse protein data from NCBI RefSeq
- `equalParser` - find protein ID mapping with equal sequences
- `mergeParser` - merge different ID mapping files
- `mppiParser` - parse protein protein interaction data from MIPS
- `paParser` - parse data from PeptideAtlas
- `dbsublocParser` - parse data from DBSubLoc
- `pfamNameParser` - parse domain id and name from Pfam
- `blastParser` - filter the results of blast

Function `fileMuncher` and `fileMuncher_DB` perl file based on given parser file and additional input data file, then perform this perl program via R.

3.2.2 Writing data using R

Besides using perl program, R functions also parse data from simple data file and store them as R environment objects or Bimap objects.

- `createEmptyDPkg` - create an empty R package at given directory.
- `writeSQ` - write sequence data into R package.
- `writeName` - parse multiple data file, and write the mapping of id and name into R package. It employs `writeGOName`, `writeKEGGName`, `writePFAMName`, `writeINTERPROName` and `writeTAXName` to respectively write data from GO, KEGG, Pfam, InterPro and TAX.
- `writeSCOPData` - parse structural classification of proteins from SCOP database.
- `writeSubCellData` - parse data from protein subcellular location databases, and write into R package. It employs `writeBACELLOData` and `writeDBSUBLOCData` to respectively write data from BaCelLo and DBSubLoc.
- `writeIntData` - parse data from protein-protein/domain-domain interaction databases, and write into R package. It employs `writeGENEINTData`, `writeINTACTData`, `writeMP-PIDData`, `write3DIDData` and `writeDOMINEDData` to respectively write data from NCBI gene interaction data file, EBI intact, MIPS interaction data, 3DID database and DOMINE database.
- `writePtmData` - parse database involving protein post-translational modifications, and write into R package. It employs `writeSYSPTMData` to write data from SysPTM database.
- `writeBfData` - write data involving body fluids proteomics into R package. It employs `writeSYSBODYFLUIDData` to write data from Sys-BodyFluid database.
- `writeGOAData` - write gene ontology terms from GOA database into R package.
- `writeHomoloGeneData` - write homolog groups from NCBI HomoloGene into R package.
- `writeInParanoidData` - write paralog groups from InParanoid into R package.
- `writePeptideAtlasData` - write peptides identified by Mass Spectrometry from PeptideAtlas database into R package.
- `writeMeta_DB` - write meta information about the annotation package into SQLite-based package.

- `writeData_DB` - parse data from databases, and write data as tables into SQLite-based R package. It employs `writeSPData_DB`, `writeIPIData_DB`, `writeREFSEQData_DB`, `writeGENEINTData_DB`, `writeINTACTData_DB`, `writeMPPIDData_DB`, `write3DIDData_DB`, `writeDOMINEDData_DB`, `writeSYSBODYFLUIDData_DB`, `writeSYSPTMData_DB`, `writeSCOPData_DB`, `writeBACELLOData_DB`, `writeDBSUBLOCDData_DB`, `writeGOADData_DB`, `writeHomoloGeneData_DB`, `writeInParanoidData_DB`, and `writePeptideAtlasData_DB` to respectively write data from Swiss-Prot, IPI, NCBI RefSeq database, and so on.
- `writeName_DB` - parse multiple data file, and write the mapping of id and name into SQLite-based R package. It employs `writeGOName_DB`, `writeKEGGName_DB`, `writePFAMName_DB`, `writeINTERPROName_DB` and `writeTAXName_DB` to respectively write data from GO, KEGG, Pfam, InterPro and TAX.
- `createSeeds` - define a list of `AnnDbBimap` objects which indicates key and value of .
- `createAnnObjs` - produce `AnnDbBimap` objects based on the definition in `createSeeds`.

3.3 Writing Help Documents

Help documents is an important part for new package. Diverse templates of help documents are stored in the "inst/templates" directory. When building new package, R functions use these templates to create "*.rd" help file in the "man" directory:

- `getRepList` - return a list which will replace the symbols in template file.
- `copyTemplates_DB` - will produce an "*.rd" file in the "man" subdirectory for SQLite-based annotation package. Related template files in the "inst/templates" subdirectory are needed for this function. .
- `writeDescription_DB` - write description file for SQLite-based annotation package.

3.4 Building Data Packages

Basic functions described above make it possible to build proteomic annotation data packages. Based on these, `PAnnBuilder` develops multiple sophisticated functions to assemble proteomic annotation data. Each function is implemented by the "*Builder_DB" R functions.

- `pBaseBuilder_DB` - build annotation data packages for primary protein database such as SwissProt, TREMBL, IPI or NCBI RefSeq protein data.
- `pSeqBuilder_DB` - build annotation data packages for query protein sequences based on sequence similarity.

- `crossBuilder_DB` - build annotation data packages for protein id mapping in SwissProt, TrEMBL, IPI and NCBI RefSeq databases.
- `subcellBuilder_DB` - build annotation data packages for protein subcellular location from BaCelLo or DBSubLoc database.
- `HomoloGeneBuilder_DB` - build annotation data packages for homolog protein group from NCBI HomoloGene database.
- `InParanoidBuilder_DB` - build annotation data packages for ortholog protein group between two given organisms from InParanoid database.
- `GOABuilder_DB` - build annotation data packages for mapping proteins of UniProt to Gene Ontology from GOA database.
- `scopBuilder_DB` - build annotation data packages for Structural Classification of Proteins.
- `intBuilder_DB` - build annotation data packages for protein-protein or domain-domain interaction from IntAct, MPPI, 3DID, DOMINE or NCBI Gene interaction database.
- `PeptideAtlasBuilder_DB` - build annotation data packages for experimentally identified peptides from PeptideAtlas database.
- `ptmBuilder_DB` - build annotation data packages for post-translational modifications from SysPTM database.
- `bfBuilder_DB` - build annotation data packages for proteins in body fluids from SysBodyFluid database.
- `dNameBuilder_DB` - build annotation data packages for mapping between entry ID and name from GO, KEGG, Pfam, InterPro and NCBI Taxonomy databases.

4 Building Annotation Data Packages

1. The first thing you need to do is setting basic parameters such as "pkgpath", "version", and "author".

```

> library(PAnnBuilder)
> pkgPath <- tempdir()
> version <- "1.0.0"
> author <- list()
> author[["authors"]] <- "Hong Li"
> author[["maintainer"]] <- "Hong Li <sysptm@gmail.com>"

```


- Then you can run diverse `"*Builder_DB"` functions to build packages by yourselves. `pBaseBuilder_DB`, `subcellBuilder_DB`, and `pSeqBuilder_DB` are taken as examples to build annotation packages.

pBaseBuilder_DB builds annotation data packages for proteins in three primary protein databases (SwissProt, IPI, RefSeq). It is a convenient way to obtain complete and canonical annotation, including protein description, Entrez gene identifier, KEGG pathway, gene ontology, domain, coordinates on chromosomes and so on. For example, if you want to build annotation package for Mouse IPI database, you can use codes as follows:

```
> pBaseBuilder_DB(baseMapType = "ipi", organism = "Mus musculus",
+   prefix = "org.Mm.ipi", pkgPath = pkgPath, version = version,
+   author = author)
```

After running, a subdirectory called `"org.Mm.ipi"` will be produced in the path given by `"pkgPath"`. This directory contains all data and files, which can be used to build R package by `"R CMD build"` command.

subcellBuilder_DB builds annotation data Package which provides protein subcellular location information.

```
> subcellBuilder_DB(src = "BaCello", prefix = "sc.bacello", pkgPath,
+   version, author)
> dir(file.path(pkgPath, "sc.bacello.db"))

[1] "DESCRIPTION" "inst"           "man"           "NAMESPACE"    "R"
```

pSeqBuilder_DB uses blast to calculate sequence similarity between query proteins and subject proteins, then assign annotation for query protein according to existing annotation of its similar proteins. `pSeqBuilder_DB` is useful for proteins which have not well annotated. Following code chunk gives an example for annotation query proteins by `pSeqBuilder_DB`. Needed R packages `org.Hs.sp.db`, `org.Hs.ipi.db`, can be downloaded via `biocLite`.

```
> tmp = system.file("data", "query.example", package = "PAnnBuilder")
> tmp = readLines(tmp)
> tag = grep("^>", tmp)
> query <- sapply(1:(length(tag) - 1), function(x) {
+   paste(tmp[(tag[x] + 1):(tag[x + 1] - 1)], collapse = "")
+ })
> query <- c(query, paste(tmp[(tag[length(tag)] + 1):length(tmp)],
+   collapse = ""))
> names(query) = sub(">", "", tmp[tag])
> blast <- c("blastp", "10.0", "BLOSUM62", "0", "-1", "-1", "T",
+   "F")
> names(blast) <- c("p", "e", "M", "W", "G", "E", "U", "F")
> match <- c(1e-05, 0.9, 0.9)
```

```

> names(match) <- c("e", "c", "i")
> if (!require("org.Hs.sp.db")) {
+   biocLite("org.Hs.sp.db")
+ }
> if (!require("org.Hs.ipi.db")) {
+   biocLite("org.Hs.ipi.db")
+ }
> annPkgs <- c("org.Hs.sp.db", "org.Hs.ipi.db")
> seqName <- c("org.Hs.spSEQ", "org.Hs.ipiSEQ")
> pSeqBuilder_DB(query, annPkgs, seqName, blast, match, prefix = "test1",
+   pkgPath, version, author)

```

3. After the running of `"*Builder"` or `"*Builder_DB"` function has been finished, a subdirectory named `"pkgName"` will be produced in given `"pkgPath"`. Then the command `"R CMD build"` can be used to build source R package, and `"R CMD --binary build"` can be used to build binary R package for Windows.

4. Note:

- Web connection is needed to download files from public databases, and Perl is needed to parse data files. Additionally, Rtools is needed for Windows user.
- Users should be aware that downloading, parsing, and saving data may take a long time, in addition to requiring enough disk space to store temporary data files.
- `"R CMD build"` and `"R CMD --binary build"` should be used in command line, not in R. Detailed document about how to create your own packages can be found in the book `"Writing R Extensions"` (<http://cran.r-project.org/doc/manuals/R-exts.pdf>). For Windows users, `"R CMD build"` needs you to have installed the files for building source packages (which is the default), as well as the Windows toolset (see the `"R Installation and Administration"` manual at <http://cran.r-project.org/doc/manuals/R-admin.pdf>).

5 Session Information

This vignette was generated using the following package versions:

```

R version 2.10.0 (2009-10-26)
x86_64-unknown-linux-gnu

```

locale:

```

[1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
[3] LC_TIME=en_US.UTF-8      LC_COLLATE=en_US.UTF-8
[5] LC_MONETARY=C            LC_MESSAGES=en_US.UTF-8

```

```
[7] LC_PAPER=en_US.UTF-8      LC_NAME=C
[9] LC_ADDRESS=C              LC_TELEPHONE=C
[11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
```

attached base packages:

```
[1] tools      stats      graphics  grDevices  utils      datasets  methods
[8] base
```

other attached packages:

```
[1] org.Hs.ipi.db_1.1.1 PAnnBuilder_1.8.0  AnnotationDbi_1.8.0
[4] Biobase_2.6.0      RSQLite_0.7-3      DBI_0.2-4
```

References

- O. N. Jensen. Interpreting the protein language using proteomics. *Nat Rev Mol Cell Biol*, 7(6):391–403, 2006.
- R Development Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2008. URL <http://www.R-project.org>. ISBN 3-900051-07-0.
- J. Zhang, V. Carey, and R. Gentleman. An extensible application for assembling annotation for genomic data. *Bioinformatics*, 19(1):155–6, 2003.