

OmnipathR: utility functions to work with Omnipath in R

Alberto Valdeolivas^{*1}, Attila Gabor^{†1}, Denes Turei^{‡1}, and Julio Saez-Rodriguez^{§1,2}

¹Institute of Computational Biomedicine, Heidelberg University, Faculty of Medicine, 69120 Heidelberg, Germany

²RWTH Aachen University, Faculty of Medicine, Joint Research Centre for Computational Biomedicine (JRC-COMBINE), 52074 Aachen, Germany

*alvaldeolivas@gmail.com †attila.gabor@bioquant.uni-heidelberg.de ‡denes.turei@embl.de
§julio.saez@bioquant.uni-heidelberg.de

April 28, 2020

Abstract

This vignette describes how to use the *OmnipathR* package to retrieve information from the Omnipath database:

<http://omnipathdb.org/>

In addition, it includes some utility functions to filter, analyse and visualize the data.

Package

OmnipathR 1.2.0

Feedbacks and bugreports are always very welcomed!

Please use the Github issue page to report bugs or for questions:

<https://github.com/saezlab/OmnipathR/issues>

Many thanks for using *OmnipathR*!

Contents

1	Introduction	3
1.1	Query types	3
1.2	Mouse and rat	4
2	Installation of the <i>OmnipathR</i> package	4
3	Usage Examples	5
3.1	Interactions.	5
3.1.1	Protein-protein interaction networks	6
3.1.2	Other interaction datasets	7
3.2	Post-translational modifications (PTMs)	10
3.3	Complexes	12
3.4	Annotations	13
3.5	Intercell	15
3.6	Conclusion	17
A	Session info	20

1 Introduction

OmnipathR is an R package built to provide easy access to the data stored in the Omnipath webservice [1]:

<http://omnipathdb.org/>

The webservice implements a very simple REST style API. This package make requests by the HTTP protocol to retrieve the data. Hence, fast Internet access is required for a proper use of *OmnipathR*.

1.1 Query types

OmnipathR can retrieve five different types of data:

- **Interactions:** protein-protein interactions organized in different datasets:
 - **Omnipath:** the OmniPath data as defined in the original publication [1] and collected from different databases.
 - **Pathwayextra:** activity flow interactions without literature reference.
 - **Kinaseextra:** enzyme-substrate interactions without literature reference.
 - **Ligrecextra:** ligand-receptor interactions without literature reference.
 - **Tfregulons:** transcription factor (TF)-target interactions from DoRothEA [2, 3].
 - **Mirnatarget:** miRNA-mRNA and TF-miRNA interactions.
- **Post-translational modifications (PTMs):** It provides enzyme-substrate reactions in a very similar way to the aforementioned interactions. Some of the biological databases related to PTMs integrated in Omnipath are Phospho.ELM [4] and PhosphoSitePlus [5].
- **Complexes:** it provides access to a comprehensive database of more than 22000 protein complexes. This data comes from different resources such as: CORUM [6] or Hu.map [7].
- **Annotations:** it provides a large variety of data regarding different annotations about proteins and complexes. These data come from dozens of databases covering different topics such as: The Topology Data Bank of Transmembrane Proteins (TOPDB) [8] or ExoCarta [9], a database collecting the proteins that were identified in exosomes in multiple organisms.
- **Intercell:** it provides information on the roles in inter-cellular signaling. For instance, if a protein is a ligand, a receptor, an extracellular matrix (ECM) component, etc. The data does not come from original sources but combined from several databases by us. The source databases, such as CellPhoneDB [10] or Receptome [11], are also referred for each record.

Figure 1 shows an overview of the resources featured in OmniPath. For more detailed information about the original data sources integrated in Omnipath, please visit: <http://omnipathdb.org/> and <http://omnipathdb.org/info>.

OmnipathR: utility functions to work with Omnipath in R

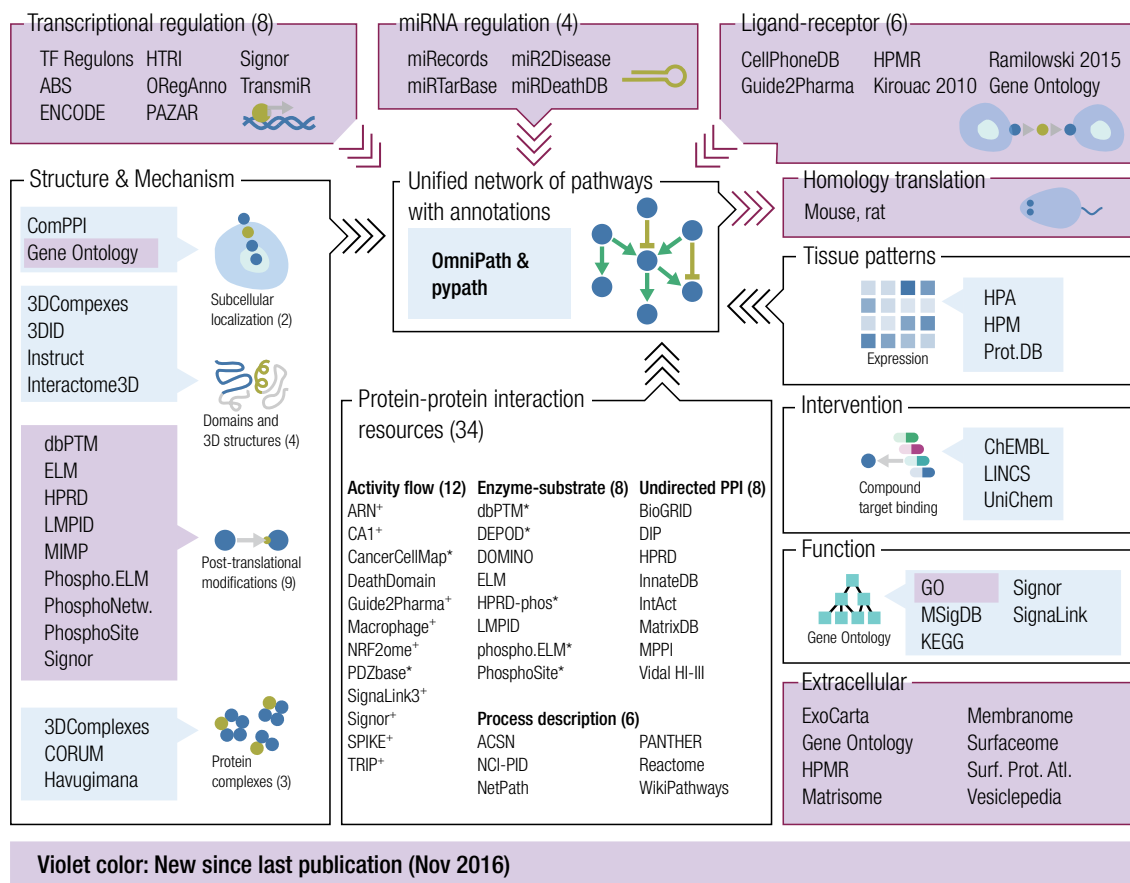


Figure 1: Overview of the resources featured in Omnipath

Causal resources (including activity-flow and enzyme-substrate resources) can provide direction (*) or sign and direction (+) of interactions.

1.2 Mouse and rat

Excluding the miRNA interactions, all interactions and PTMs are available for human, mouse and rat. The rodent data has been translated from human using the NCBI Homologene database. Many human proteins do not have known homolog in rodents hence rodent datasets are smaller than their human counterparts.

In case you work with mouse omics data you might do better to translate your dataset to human (for example using the pypath.homology module, <https://github.com/saezlab/pypath/>) and use human interaction data.

2 Installation of the OmnipathR package

First of all, you need a current version of R (www.r-project.org). OmnipathR is a freely available package deposited on <http://bioconductor.org/> and <https://github.com/saezlab/OmnipathR>. You can install it by running the following commands on an R console:

```
> if (!requireNamespace("BiocManager", quietly = TRUE))
+   install.packages("BiocManager")
> BiocManager::install("OmnipathR")
```

3 Usage Examples

In the following paragraphs, we provide some examples to describe how to use the [OmnipathR](#) package to retrieve different types of information from Omnipath webserver. In addition, we play around with the data aiming at obtaining some biological relevant information.

Noteworthy, the sections **complexes**, **annotations** and **intercell** are linked. We explore the annotations and roles in inter-cellular communications of the proteins involved in a given complex. This basic example shows the usefulness of integrating the information available in the different **Omnipath** resources.

3.1 Interactions

Proteins interact among them and with other biological molecules to perform cellular functions. Proteins also participates in pathways, linked series of reactions occurring inter/intra cells to transform products or to transmit signals inducing specific cellular responses. Protein interactions are therefore a very valuable source of information to understand cellular functioning.

We are going to download the original **Omnipath** human interactions [1]. To do so, we first check the different source databases and select some of them. Then, we print some of the downloaded interactions ("+" means activation, "-" means inhibition and "?" means undirected interactions or inconclusive data).

```
> library(OmnipathR)
> library(tidyr)
> library(dnet)
> library(gprofiler2)
> ## We check some of the different interaction databases
> head(get_interaction_databases(),10)

[1] "TRIP"      "DIP"      "Wang"      "KEGG"      "BioGRID"  "Fantom5"  "HPRD"
[8] "LRdb"      "KEA"      "MIMP"

> ## The interactions are stored into a data frame.
> interactions <-
+   import_Omnipath_Interactions(filter_databases=c("SignaLink3","PhosphoSite",
+   "Signor"))
> ## We visualize the first interactions in the data frame.
> print_interactions(head(interactions))
```

	source	interaction	target	nsources	nrefs
172	TRPM7 (Q96QT4)	==(+)==>	ANXA1 (P04083)	13	15
99	SRC (P12931)	==(+)==>	TRPV1 (Q8NER1)	9	8
45	PRKG1 (Q13976)	==(-)==>	TRPC3 (Q13507)	9	6
75	PRKG1 (Q13976)	==(-)==>	TRPC6 (Q9Y210)	8	5

```
118   LYN (P07948) == ( + ) ==> TRPV4 (Q9HBA0)      8      5
98   PRKACA (P17612) == ( + ) ==> TRPV1 (Q8NER1)     6      1
```

3.1.1 Protein-protein interaction networks

Protein-protein interactions are usually converted into networks. Describing protein interactions as networks not only provides a convenient format for visualization, but also allows applying graph theory methods to mine the biological information they contain.

We convert here our set of interactions to a network/graph (*igraph* object). Then, we apply two very common approaches to extract information from a biological network:

- **Shortest Paths:** finding a path between two nodes (proteins) going through the minimum number of edges. This can be very useful to track consecutive reactions within a given pathway. We display below the shortest path between two given proteins and all the possible shortest paths between two other proteins. It is to note that the functions *printPath_es* and *printPath_vs* display very similar results, but the first one takes as an input an edge sequence and the second one a node sequence.

```
> ## We transform the interactions data frame into a graph
> OPI_g <- interaction_graph(interactions = interactions)
> ## Find and print shortest paths on the directed network between proteins
> ## of interest:
> printPath_es(shortest_paths(OPI_g, from = "TYR03", to = "STAT3",
+   output = 'epath')$epath[[1]], OPI_g)

      source interaction      target nsources nrefs
1 TYR03 (Q06418) == ( ? ) ==> AKT1 (P31749)      2      0
2 AKT1 (P31749) == ( ? ) ==>  BTK (Q06187)      3      1
3 BTK (Q06187) == ( ? ) ==> STAT3 (P40763)      2      2

> ## Find and print all shortest paths between proteins of interest:
> printPath_vs(all_shortest_paths(OPI_g, from = "DYRK2",
+   to = "MAPKAPK2")$res, OPI_g)
```

- **Clustering:** grouping nodes (proteins) in such a way that nodes belonging to the same group (called cluster) are more connected in the network to each other than to those in other groups (clusters). Since proteins interact to perform their functions, proteins within the same cluster are likely to be implicated in similar biological tasks. Figure 2 shows the subgraph containing the proteins and interactions of a specific protein. The *igraph* R package contains functions to apply several different cluster methods on graphs (visit <https://igraph.org/r/doc/> for detailed information.)

```
> ## We apply a clustering algorithm (Louvain) to group proteins in
> ## our network. We apply here Louvain which is fast but can only run
> ## on undirected graphs. Other clustering algorithms can deal with
> ## directed networks but with longer computational times,
> ## such as cluster_edge_betweenness. These cluster methods are directly
> ## available in the igraph package.
> OPI_g_undirected <- as.undirected(OPI_g, mode=c("mutual"))
> cl_results <- cluster_louvain(OPI_g_undirected)
> ## We extract the cluster where a protein of interest is contained
```

OmnipathR: utility functions to work with Omnipath in R

```
> cluster_id <- cl_results$membership[which(cl_results$names == "CD22")]
> module_graph <- induced_subgraph(OPI_g_undirected,
+   V(OPI_g)$name[which(cl_results$membership == cluster_id)])

> ## We print that cluster with its interactions.
> par(mar=c(0.1,0.1,0.1,0.1))
> plot(module_graph, vertex.label.color="black",vertex.frame.color="#ffffff",
+   vertex.size= 15, edge.curved=.2,
+   vertex.color = ifelse(igraph::V(module_graph)$name == "CD22","yellow",
+   "#00CCFF"), edge.color="blue",edge.width=0.8)
```

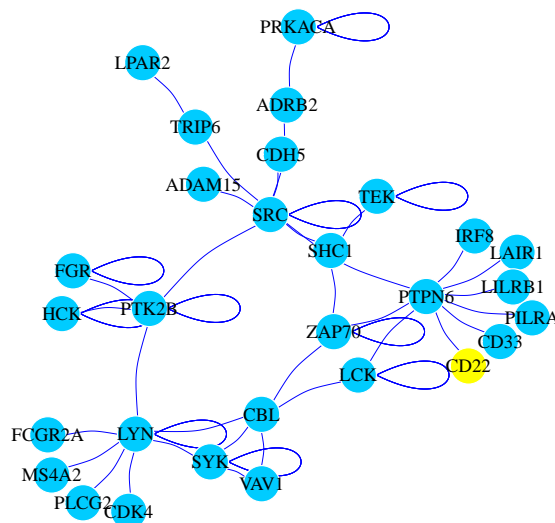


Figure 2:
Subnetwork extracted from the interactions graph representing the cluster where we can find the gene CD22 (yellow node).

3.1.2 Other interaction datasets

We used above the interactions from the dataset described in the original **Omnipath** publication [1]. In this section, we provide examples on how to retry and deal with interactions from the remaining datasets. The same functions can be applied to every interaction dataset.

In the first example, we are going to get the interactions from the **pathwayextra** dataset, which contains activity flow interactions without literature reference. We are going to focus on the mouse interactions for a given gene in this particular case.

```
> ## We query and store the interactions into a dataframe
> interactions <-
+   import_PathwayExtra_Interactions(filter_databases=c("BioGRID","IntAct"),
+   select_organism = 10090)
> ## We select all the interactions in which Amfr gene is involved
> interactions_Amfr <- dplyr::filter(interactions, source_genesymbol == "Amfr" |
+   target_genesymbol == "Amfr")
```

OmnipathR: utility functions to work with Omnipath in R

```
> ## We print these interactions:
> print_interactions(interactions_Amfr)

      source interaction      target nsources
1 Amfr (Q9R049) ==(+ )==> Vcp (Q01853)      3
```

Next, we download the interactions from the **kinaseextra** dataset, which contains enzyme-substrate interactions without literature reference. We are going to focus on rat reactions targeting a particular gene.

```
> ## We query and store the interactions into a dataframe
> interactions <-
+   import_KinaseExtra_Interactions(filter_databases=c("PhosphoPoint",
+   "PhosphoSite"), select_organism = 10116)
> ## We select the interactions in which Dpysl2 gene is a target
> interactions_TargetDpysl2 <- dplyr::filter(interactions,
+   target_genesymbol == "Dpysl2")
> ## We print these interactions:
> print_interactions(interactions_TargetDpysl2)

      source interaction      target nsources
1 Gsk3b (P18266) ==(+/-)==> Dpysl2 (P47942)      18
2 Rock2 (Q62868) ==(+ )==> Dpysl2 (P47942)      11
4 Cdk5 (Q03114) ==(+ )==> Dpysl2 (P47942)      10
6 Rock1 (Q63644) ==( ? )==> Dpysl2 (P47942)       8
5 Gsk3a (P18265) ==( ? )==> Dpysl2 (P47942)       7
3 Fer (P09760) ==( ? )==> Dpysl2 (P47942)       5
```

In the following example we are going to work with the **ligrecextra** dataset, which contains ligand-receptor interactions without literature reference. Our goal is to find the potential receptors associated to a given ligand. For a more global overview, we induce a network containing the genes involved in these interactions (Figure 3).

```
> ## We query and store the interactions into a dataframe
> interactions <- import_LigrecExtra_Interactions(filter_databases=c("HPRD",
+   "Guide2Pharma"),select_organism=9606)
> ## Receptors of the CDH1 ligand.
> interactions_CDH1 <- dplyr::filter(interactions, source_genesymbol == "CDH1")
> ## We transform the interactions data frame into a graph
> OPI_g <- interaction_graph(interactions = interactions_CDH1)
> ## We induce a network with the genes involved in the shortest path and their
> ## first neighbors to get a more general overview of the interactions
> Induced_Network <- dNetInduce(g=OPI_g,
+   nodes_query=as.character( V(OPI_g)$name), knn=0,
+   remove_loops=FALSE, largest.comp=FALSE)
>
```

```
> ## We print the induced network
> par(mar=c(0.1,0.1,0.1,0.1))
> plot(Induced_Network, vertex.label.color="black",
+   vertex.frame.color="#ffffff",vertex.size= 20, edge.curved=.2,
+   vertex.color =
```


OmnipathR: utility functions to work with Omnipath in R

```
+ ifelse(igraph::V(Induced_Network)$name %in% c("CDH1"),
+ "yellow", "#00CCFF"), edge.color="blue", edge.width=0.8)
```

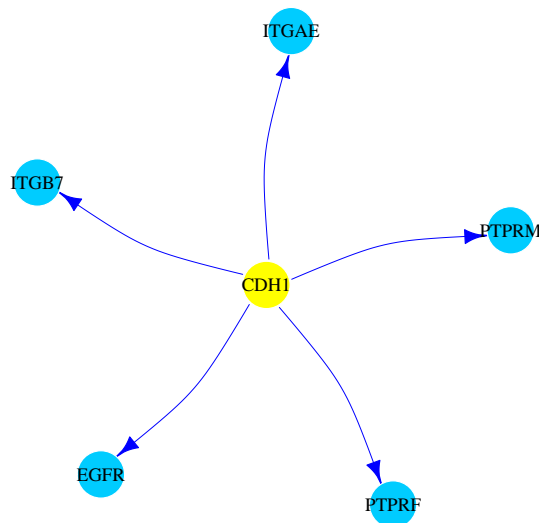


Figure 3:

Subnetwork extracted from the **kinaseextra** interactions graph containing the shortest path between *B2M* and *TFR2* (yellow nodes). The first neighbors of the genes involved in the shortest path are also shown.

Another very interesting interaction dataset also available in Omnipath are the **tfregulons** from DoRothEA [2, 3]. It contains transcription factor (TF)-target interactions with confidence score, ranging from A-E, being A the most confident interactions. In the code chunk shown below, we select and print the most confident interactions for a given TF.

```
> ## We query and store the interactions into a dataframe
> interactions <- import_TFregulons_Interactions(filter_databases=c("DoRothEA_A",
+ "ARACNe-GTEX"), select_organism=9606)
> ## We select the most confident interactions for a given TF and we print
> ## the interactions to check the way it regulates its different targets
> interactions_A_GLI1 <- dplyr::filter(interactions, tfregulons_level=="A",
+ source_genesymbol == "GLI1")
> print_interactions(interactions_A_GLI1)
```

	source	interaction	target	nsources
1	GLI1 (P08151)	==(+)==>	PTCH1 (Q13635)	5
2	GLI1 (P08151)	==(+)==>	BCL2 (P10415)	5
3	GLI1 (P08151)	==(-)==>	EGR2 (P11161)	3
4	GLI1 (P08151)	==(+)==>	IGFBP6 (P24592)	3
5	GLI1 (P08151)	==(+)==>	SFRP1 (Q8N474)	3
6	GLI1 (P08151)	==(-)==>	SLIT2 (O94813)	3

The last dataset describing interactions is **mirnatarget**. It stores miRNA-mRNA and TF-miRNA interactions. These interactions are only available for human so far. We next select the miRNA interacting with the TF selected in the previous code chunk, *GLI1*. The main function of miRNAs seems to be related with gene regulation. It is therefore interesting to see

OmnipathR: utility functions to work with Omnipath in R

how some miRNA can regulate the expression of a TF which in turn regulates the expression of other genes. Figure 4 shows a schematic network of the miRNA targeting *GLI1* and the genes regulated by this TF.

```
> ## We query and store the interactions into a dataframe
> interactions <-
+   import_miRNAtarget_Interactions(filter_databases=c("miRTarBase","miRecords"))
> ## We select the interactions where a miRNA is interacting with the TF
> ## used in the previous code chunk and we print these interactions.
> interactions_miRNA_GLI1 <-
+   dplyr::filter(interactions, target_genesymbol == "GLI1")
> print_interactions(interactions_miRNA_GLI1)
```

	source	interaction	target	nsources	nrefs
2	hsa-miR-324-5p (MIMAT0000761)	==(?)==>	GLI1 (P08151)	3	3
1	hsa-miR-125b (MIMAT0000423)	==(?)==>	GLI1 (P08151)	2	1
3	hsa-miR-326 (MIMAT0000756)	==(?)==>	GLI1 (P08151)	2	1
4	hsa-miR-202 (MIMAT0002811)	==(?)==>	GLI1 (P08151)	1	1
5	hsa-miR-133b (MIMAT0000770)	==(?)==>	GLI1 (P08151)	1	1

```
> ## We transform the previous selections to graphs (igraph objects)
> OPI_g_1 <-interaction_graph(interactions = interactions_A_GLI1)
> OPI_g_2 <-interaction_graph(interactions = interactions_miRNA_GLI1)
```

```
> ## We print the union of both previous graphs
> par(mar=c(0.1,0.1,0.1,0.1))
> plot(OPI_g_1 %u% OPI_g_2, vertex.label.color="black",
+   vertex.frame.color="#ffffff",vertex.size= 20, edge.curved=.25,
+   vertex.color = ifelse(grepl("miR",igraph::V(OPI_g_1 %u% OPI_g_2)$name),
+   "red",ifelse(igraph::V(OPI_g_1 %u% OPI_g_2)$name == "GLI1",
+   "yellow","#00CCFF")), edge.color="blue",
+   vertex.shape = ifelse(grepl("miR",igraph::V(OPI_g_1 %u% OPI_g_2)$name),
+   "vrectangle","circle"),edge.width=0.8)
```

3.2 Post-translational modifications (PTMs)

Another query type available is PTMs which provides enzyme-substrate reactions in a very similar way to the aforementioned interactions. PTMs refer generally to enzymatic modification of proteins after their synthesis in the ribosomes. PTMs can be highly context-specific and they play a main role in the activation/inhibition of biological pathways.

In the next code chunk, we download the **PTMs** for human. We first check the different available source databases, even though we do not perform any filter. Then, we select and print the reactions involving a specific enzyme-substrate pair. Those reactions lack information about activation or inhibition. To obtain that information, we match the data with **Omnipath** interactions. Finally, we show that it is also possible to build a graph using this information, and to retrieve PTMs from mouse or rat.

```
> ## We check the different PTMs databases
> get_ptms_databases()

[1] "KEA" "MIMP"
```

OmnipathR: utility functions to work with Omnipath in R

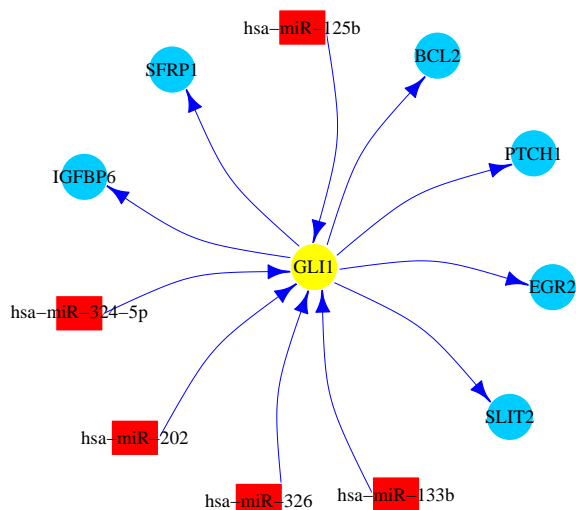


Figure 4:

Schematic network of the miRNA (red square nodes) targeting *GLI1* (yellow node) and the genes regulated by this TF (blue round nodes).

```
[3] "PhosphoNetworks"      "PhosphoSite"
[5] "PhosphoSite_MIMP"    "PhosphoSite_ProtMapper"
[7] "ProtMapper"          "phosphoELM"
[9] "phosphoELM_MIMP"     "BEL-Large-Corpus_ProtMapper"
[11] "HPRD"                "HPRD_MIMP"
[13] "RLIMS-P_ProtMapper"  "SIGNOR"
[15] "SIGNOR_ProtMapper"   "dbPTM"
[17] "Li2012"              "NCI-PID_ProtMapper"
[19] "REACH_ProtMapper"    "Sparsen_ProtMapper"
[21] "Reactome_ProtMapper" "DEPOD"
```

```
> ## We query and store the ptms into a dataframe. No filtering by
> ## databases in this case.
> ptms <- import_Omnipath_PTMS()
> ## We can select and print the reactions between a specific kinase and
> ## a specific substrate
> print_interactions(dplyr::filter(ptms, enzyme_genesymbol=="MAP2K1",
+   substrate_genesymbol=="MAPK3"))
```

	enzyme	interaction	substrate	modification	nresources
1	MAP2K1 (Q02750)	====>	MAPK3_Y204 (P27361)	phosphorylation	16
2	MAP2K1 (Q02750)	====>	MAPK3_T202 (P27361)	phosphorylation	15
3	MAP2K1 (Q02750)	====>	MAPK3_Y210 (P27361)	phosphorylation	3
4	MAP2K1 (Q02750)	====>	MAPK3_T207 (P27361)	phosphorylation	3
5	MAP2K1 (Q02750)	====>	MAPK3_T80 (P27361)	phosphorylation	1
6	MAP2K1 (Q02750)	====>	MAPK3_Y222 (P27361)	phosphorylation	1

```
> ## In the previous results, we can see that ptms does not contain sign
> ## (activation/inhibition). We can generate this information based on the
> ## protein-protein Omnipath interaction dataset.
```

```
> interactions <- import_Omnipath_Interactions()
> ptms <- get_signed_ptms(ptms, interactions)
> ## We select again the same kinase and substrate. Now we have information
> ## about inhibition or activation when we print the ptms
> print_interactions(dplyr::filter(ptms, enzyme_genesymbol=="MAP2K1",
+   substrate_genesymbol=="MAPK3"))
```

	enzyme	interaction	substrate	modification	nresources
5	MAP2K1 (Q02750)	==(+)==>	MAPK3_Y204 (P27361)	phosphorylation	16
6	MAP2K1 (Q02750)	==(+)==>	MAPK3_T202 (P27361)	phosphorylation	15
1	MAP2K1 (Q02750)	==(+)==>	MAPK3_T207 (P27361)	phosphorylation	3
2	MAP2K1 (Q02750)	==(+)==>	MAPK3_Y210 (P27361)	phosphorylation	3
3	MAP2K1 (Q02750)	==(+)==>	MAPK3_T80 (P27361)	phosphorylation	1
4	MAP2K1 (Q02750)	==(+)==>	MAPK3_Y222 (P27361)	phosphorylation	1

```
> ## We can also transform the ptms into a graph.
> ptms_g <- ptms_graph(ptms = ptms)
> ## We download PTMs for mouse
> ptms <- import_Omnipath_PTMS(filter_databases=c("PhosphoSite", "Signor"),
+   select_organism=10090)
```

3.3 Complexes

Some studies indicate that around 80% of the human proteins operate in complexes, and many proteins belong to several different complexes [12]. These complexes play critical roles in a large variety of biological processes. Some well-known examples are the proteasome and the ribosome. Thus, the description of the full set of protein complexes functioning in cells is essential to improve our understanding of biological processes.

The **complexes** query provides access to more than 20000 protein complexes. This comprehensive database has been created by integrating different resources. We now download these molecular complexes filtering by some of the source databases. We check the complexes where a couple of specific genes participate. First, we look for the complexes where any of these two genes participate. We then identify the complex where these two genes are jointly involved. Finally, we perform an enrichment analysis with the genes taking part in that complex. You should keep an eye on this complex since it will be used again in the forthcoming sections.

```
> ## We check the different complexes databases
> get_complexes_databases()
```

[1]	"Compleat"	"CORUM"	"hu.MAP"	"ComplexPortal"
[5]	"Signor"	"PDB"	"CellPhoneDB"	"Havugimana2012"
[9]	"HPMR"	"Guide2Pharma"	"CFinder"	"NetworkBlast"

```
> ## We query and store complexes from some sources into a dataframe.
> complexes <- import_Omnipath_complexes(filter_databases=c("CORUM", "hu.MAP"))
> ## We check all the molecular complexes where a set of genes participate
> query_genes <- c("WRN", "PARP1")
> ## Complexes where any of the input genes participate
> complexes_query_genes_any <- unique(get_complex_genes(complexes, query_genes,
+   total_match=FALSE))
```

OmnipathR: utility functions to work with Omnipath in R

```
> ## We print the components of the different selected components
> head(complexes_query_genes_any$components_genesymbols,6)

[1] "NCAPD2_NCAPH_PARP1_SMC2_SMC4_XRCC1"
[2] "CCNA2_CDK2_LIG1_PARP1_POLA1_POLD1_POLE_RFC1_RFC2_RPA1_RPA2_RPA3_TOP1"
[3] "CCNA2_CCNB1_CDK1_PARP1_POLA1_POLD1_POLE_RFC1_RFC2_RPA1_RPA2_RPA3_TOP1"
[4] "MRE11_PARP1_RAD50_TERF2_TERF2IP_XRCC5_XRCC6"
[5] "TERF2_WRN"
[6] "CALR_DHX30_H2AFX_HIST3H2BB_HSPA5_NPM1_PARP1"

> ## Complexes where all the input genes participate jointly
> complexes_query_genes_join <- unique(get_complex_genes(complexes,query_genes,
+   total_match=TRUE))
> ## We print the components of the different selected components
> complexes_query_genes_join$components_genesymbols

[1] "PARP1_WRN_XRCC5_XRCC6"
```

```
> genes_complex <-
+   unlist(strsplit(complexes_query_genes_join$components_genesymbols, "-"))
> ## We can perform an enrichment analyses with the genes in the complex
> EnrichmentResults <- gost(genes_complex, significant = TRUE,
+   user_threshold = 0.001, correction_method = c("fdr"),
+   sources=c("GO:BP", "GO:CC", "GO:MF"))
> ## We show the most significant results
> EnrichmentResults$result %>%
+   dplyr::select(term_id, source, term_name, p_value) %>%
+   dplyr::top_n(5, -p_value)
```

	term_id	source	term_name	p_value
1	GO:0010332	GO:BP	response to gamma radiation	5.249561e-08
2	GO:0032392	GO:BP	DNA geometric change	3.532059e-07
3	GO:0032508	GO:BP	DNA duplex unwinding	3.532059e-07
4	GO:0010212	GO:BP	response to ionizing radiation	6.490259e-07
5	GO:0000781	GO:CC	chromosome, telomeric region	3.067319e-07

3.4 Annotations

Biological annotations are statements, usually traceable and curated, about the different features of a biological entity. At the genetic level, annotations describe the biological function, the subcellular situation, the DNA location and many other related properties of a particular gene or its gene products.

The annotations query provides a large variety of data about proteins and complexes. These data come from dozens of databases and each kind of annotation record contains different fields. Because of this, here we have a `record_id` field which is unique within the records of each database. Each row contains one key value pair and you need to use the `record_id` to connect the related key-value pairs (see examples below).

OmnipathR: utility functions to work with Omnipath in R

Now, we focus in the annotations of the complex studied in the previous section. We first inspect the different available databases in the omnipath webserver. Then, we download the annotations for our complex itself as a biological entity. We find annotations related to the nucleus and transcriptional control, which is in agreement with the enrichment analysis results of its individual components.

```
> ## We check the different annotation databases
> get_annotation_databases()

[1] "Ramilowski_location" "MSigDB"          "HGNC"
[4] "GO_Intercell"        "HPMR"            "LOCATE"
[7] "Zhong2015"           "HPA_secretome"   "CPAD"
[10] "kinase.com"          "Guide2Pharma"    "Baccin2019"
[13] "KEGG"                "DGIdb"           "CompPPI"
[16] "Adhesome"            "Integrins"       "HPA_subcellular"
[19] "MatrixDB"            "DisGeNet"        "Surfaceome"
[22] "CSPA"                "TopDB"           "NetPath"
[25] "Int0Gen"             "TFcensus"        "OPM"
[28] "Matrisome"           "Kirouac2010"     "CancerSEA"
[31] "Vesiclepedia"        "Phosphatome"     "Exocarta"
[34] "Ramilowski2015"      "CancerGeneCensus" "Membranome"
[37] "HPA_tissue"          "LRdb"            "HPMR_complex"
[40] "CORUM_Funcat"        "CORUM_GO"        "Signalink3"
[43] "SIGNOR"

> ## We can further investigate the features of the complex selected
> ## in the previous section.
>
> ## We first get the annotations of the complex itself:
> annotations <- import_Omnipath_annotations(select_genes=paste0("COMPLEX:",
+ complexes_query_genes_join$components_genesymbols))
> head(dplyr::select(annotations,source,label,value),10)
```

	source	label	value
1	Ramilowski_location	location	nucleus
2	MSigDB	collection	reactome_pathways
3	MSigDB	geneset	REACTOME_DNA_DOUBLE_STRAND_BREAK_REPAIR
4	MSigDB	collection	chemical_and_genetic_perturbations
5	MSigDB	geneset	COLLIS_PRKDC_SUBSTRATES
6	MSigDB	collection	chemical_and_genetic_perturbations
7	MSigDB	geneset	PUJANA_CHEK2_PCC_NETWORK
8	MSigDB	collection	chemical_and_genetic_perturbations
9	MSigDB	geneset	PUJANA_BRCA1_PCC_NETWORK
10	MSigDB	collection	reactome_pathways

Afterwards, we explore the annotations of the individual components of the complex in some databases. We check the pathways where these proteins are involved. Once again, we also find many nucleus related annotations when checking their cellular location.

```
> ## Then, we explore some annotations of its individual components
>
> ## Pathways where the proteins belong:
> annotations <- import_Omnipath_annotations(select_genes=genes_complex,
+ filter_databases=c("NetPath"))
```

OmnipathR: utility functions to work with Omnipath in R

```
> dplyr::select(annotations, genesymbol, value)

  genesymbol                                value
1    PARP1                                Androgen receptor (AR)
2    PARP1 TNF-related weak inducer of apoptosis (TWEAK)
3    PARP1          Corticotropin-releasing hormone (CRH)
4    PARP1          Tumor necrosis factor (TNF) alpha
5    PARP1                                Oncostatin-M (OSM)
6    XRCC5                                Androgen receptor (AR)
7    XRCC6                                Androgen receptor (AR)

> ## Cellular localization of our proteins
> annotations <- import_Omnipath_annotations(select_genes=genes_complex,
+   filter_databases=c("ComPPI"))
> ## Since we have same record_id for some results of our query, we spread
> ## these records across columns
> spread(annotations, label, value) %>%
+   dplyr::arrange(desc(score)) %>%
+   dplyr::top_n(10, score)

  uniprot genesymbol entity_type source record_id location
1  P12956    XRCC6      protein ComPPI    2975   nucleus
2  P09874    PARP1      protein ComPPI   11259   nucleus
3  Q14191     WRN      protein ComPPI   16129   nucleus
4  P13010    XRCC5      protein ComPPI   13400   nucleus
5  P13010    XRCC5      protein ComPPI   13398 membrane
6  P12956    XRCC6      protein ComPPI    2976   cytosol
7  P13010    XRCC5      protein ComPPI   13399   cytosol
8  Q14191     WRN      protein ComPPI   16130   cytosol
9  P12956    XRCC6      protein ComPPI    2972 extracellular
10 P12956    XRCC6      protein ComPPI    2974   membrane
11 P13010    XRCC5      protein ComPPI   13397 extracellular

  score
1 0.99999997629184
2 0.999999887104
3 0.9999996544
4 0.99999868288
5 0.972
6 0.958
7 0.958
8 0.94
9 0.8600000000000001
10 0.8600000000000001
11 0.8600000000000001
```

3.5 Intercell

Cells perceive cues from their microenvironment and neighboring cells, and respond accordingly to ensure proper activities and coordination between them. The ensemble of these communication process is called inter-cellular signaling (**intercell**).

OmnipathR: utility functions to work with Omnipath in R

Intercell query provides information about the roles of proteins in inter-cellular signaling (e.g. if a protein is a ligand, a receptor, an extracellular matrix (ECM) component, etc.) This query type is very similar to annotations. However, **intercell** data does not come from original sources, but combined from several databases by us into categories (we also refer to the original sources).

We first inspect the different categories available in the Omnipath webserver. Then, we focus again in our previously selected complex and we check its potential roles in inter-cellular signaling. We repeat the analysis with its individual components.

```
> ## We check some of the different intercell categories
> head(get_intercell_categories(),10)

[1] "cell_surface"          "cell_surface_cspa"
[3] "cell_surface_dgidb"    "cell_surface_go"
[5] "cell_surface_hpmr"     "cell_surface_membranome"
[7] "cell_surface_surfaceome" "chemokine_ligands_hgnc"
[9] "ecm"                   "ecm_go"

> ## We import the intercell data into a dataframe
> intercell <- import_Omnipath_intercell()
> ## We check the intercell annotations for our previous complex itself
> dplyr::filter(intercell,
+   genesymbol == complexes_query_genes_join$components_genesymbols,
+   mainclass != "") %>%
+   dplyr::select(category,genesymbol, mainclass)

[1] category    genesymbol mainclass
<0 rows> (or 0-length row.names)

> ## We check the intercell annotations for the individual components of
> ## our previous complex. We filter our data to print it in a good format
> dplyr::filter(intercell,genesymbol %in% genes_complex, mainclass!="") %>%
+   dplyr::distinct(genesymbol,mainclass, .keep_all = TRUE) %>%
+   dplyr::select(category, genesymbol, mainclass) %>%
+   dplyr::arrange(genesymbol)

      category genesymbol      mainclass
1 intracellular_locate  PARP1  intracellular
2 extracellular_comppi  PARP1  extracellular
3 intracellular_locate   WRN   intracellular
4 intracellular_locate  XRCC5  intracellular
5 extracellular_matrixdb XRCC5  extracellular
6 extracellular_enzyme_go XRCC5 extracellular_enzyme
7 cell_surface_cspa     XRCC6   cell_surface
8 intracellular_locate  XRCC6  intracellular
9 extracellular_matrixdb XRCC6  extracellular
10 surface_enzyme_go    XRCC6  surface_enzyme
11 extracellular_enzyme_go XRCC6 extracellular_enzyme

> ## We close graphical connections
> while (!is.null(dev.list())) dev.off()
```


3.6 Conclusion

OmnipathR provides access to the wealth of data stored in the Omnipath webservice <http://omnipathdb.org/> from the *R* environment. In addition, it contains some utility functions for visualization, filtering and analysis. The main strength of *OmnipathR* is the straightforward transformation of the different Omnipath data into commonly used *R* objects, such as dataframes and graphs. Consequently, it allows an easy integration of the different types of data and a gateway to the vast number of *R* packages dedicated to the analysis and representation of biological data. We highlighted these abilities in some of the examples detailed in previous sections of this document.

References

- [1] Dénes Türei, Tamás Korcsmáros, and Julio Saez-Rodriguez. OmniPath: guidelines and gateway for literature-curated signaling pathway resources. *Nature Methods*, 13(12):966–967, November 2016. URL: <https://doi.org/10.1038/nmeth.4077>, doi:10.1038/nmeth.4077.
- [2] Luz Garcia-Alonso, Francesco Iorio, Angela Matchan, Nuno Fonseca, Patricia Jaaks, Gareth Peat, Miguel Pignatelli, Fiammetta Falcone, Cyril H. Benes, Ian Dunham, Graham Bignell, Simon S. McDade, Mathew J. Garnett, and Julio Saez-Rodriguez. Transcription factor activities enhance markers of drug sensitivity in cancer. *Cancer Research*, 78(3):769–780, December 2017. URL: <https://doi.org/10.1158/0008-5472.can-17-1679>, doi:10.1158/0008-5472.can-17-1679.
- [3] Luz Garcia-Alonso, Christian H. Holland, Mahmoud M. Ibrahim, Denes Turei, and Julio Saez-Rodriguez. Benchmark and integration of resources for the estimation of human transcription factor activities. *Genome Research*, 29(8):1363–1375, July 2019. URL: <https://doi.org/10.1101/gr.240663.118>, doi:10.1101/gr.240663.118.
- [4] H. Dinkel, C. Chica, A. Via, C. M. Gould, L. J. Jensen, T. J. Gibson, and F. Diella. Phospho.ELM: a database of phosphorylation sites—update 2011. *Nucleic Acids Research*, 39(Database):D261–D267, November 2010. URL: <https://doi.org/10.1093/nar/gkq1104>, doi:10.1093/nar/gkq1104.
- [5] Peter V. Hornbeck, Bin Zhang, Beth Murray, Jon M. Kornhauser, Vaughan Latham, and Elzbieta Skrzypek. PhosphoSitePlus, 2014: mutations, PTMs and recalibrations. *Nucleic Acids Research*, 43(D1):D512–D520, December 2014. URL: <https://doi.org/10.1093/nar/gku1267>, doi:10.1093/nar/gku1267.
- [6] Madalina Giurgiu, Julian Reinhard, Barbara Brauner, Irmtraud Dunger-Kaltenbach, Gisela Fobo, Goar Frishman, Corinna Montrone, and Andreas Ruepp. CORUM: the comprehensive resource of mammalian protein complexes—2019. *Nucleic Acids Research*, 47(D1):D559–D563, October 2018. URL: <https://doi.org/10.1093/nar/gky973>, doi:10.1093/nar/gky973.
- [7] Kevin Drew, Chanjae Lee, Ryan L Huizar, Fan Tu, Blake Borgeson, Claire D McWhite, Yun Ma, John B Wallingford, and Edward M Marcotte. Integration of over 9, 000 mass spectrometry experiments builds a global map of human protein complexes. *Molecular Systems Biology*, 13(6):932, June 2017. URL: <https://doi.org/10.15252/msb.20167490>, doi:10.15252/msb.20167490.
- [8] László Dobson, Tamás Langó, István Reményi, and Gábor E. Tusnády. Expediting topology data gathering for the TOPDB database. *Nucleic Acids Research*, 43(D1):D283–D289, November 2014. URL: <https://doi.org/10.1093/nar/gku1119>, doi:10.1093/nar/gku1119.
- [9] Shivakumar Keerthikumar, David Chisanga, Dinuka Ariyaratne, Haidar Al Saffar, Sushma Anand, Kening Zhao, Monisha Samuel, Mohashin Pathan, Markandeya Jois, Naveen Chilamkurti, Lahiru Gangoda, and Suresh Mathivanan. ExoCarta: A web-based compendium of exosomal cargo. *Journal of Molecular Biology*, 428(4):688–692, February 2016. URL: <https://doi.org/10.1016/j.jmb.2015.09.019>, doi:10.1016/j.jmb.2015.09.019.

- [10] Roser Vento-Tormo, Mirjana Efremova, Rachel A. Botting, Margherita Y. Turco, Miquel Vento-Tormo, Kerstin B. Meyer, Jong-Eun Park, Emily Stephenson, Krzysztof Polański, Angela Goncalves, Lucy Gardner, Staffan Holmqvist, Johan Henriksson, Angela Zou, Andrew M. Sharkey, Ben Millar, Barbara Innes, Laura Wood, Anna Wilbrey-Clark, Rebecca P. Payne, Martin A. Ivarsson, Steve Lisgo, Andrew Filby, David H. Rowitch, Judith N. Bulmer, Gavin J. Wright, Michael J. T. Stubbington, Muzlifah Haniffa, Ashley Moffett, and Sarah A. Teichmann. Single-cell reconstruction of the early maternal–fetal interface in humans. *Nature*, 563(7731):347–353, November 2018. URL: <https://doi.org/10.1038/s41586-018-0698-6>, doi:10.1038/s41586-018-0698-6.
- [11] I. Ben-Shlomo, S. Yu Hsu, R. Rauch, H. W. Kowalski, and A. J. W. Hsueh. Signaling receptome: A genomic and evolutionary perspective of plasma membrane receptors involved in signal transduction. *Science Signaling*, 2003(187):re9–re9, June 2003. URL: <https://doi.org/10.1126/stke.2003.187.re9>, doi:10.1126/stke.2003.187.re9.
- [12] Tord Berggård, Sara Linse, and Peter James. Methods for the detection and analysis of protein–protein interactions. *PROTEOMICS*, 7(16):2833–2842, August 2007. URL: <https://doi.org/10.1002/pmic.200700131>, doi:10.1002/pmic.200700131.

A Session info

- R version 4.0.0 (2020-04-24), x86_64-pc-linux-gnu
- Locale: LC_CTYPE=en_US.UTF-8, LC_NUMERIC=C, LC_TIME=en_US.UTF-8, LC_COLLATE=C, LC_MONETARY=en_US.UTF-8, LC_MESSAGES=en_US.UTF-8, LC_PAPER=en_US.UTF-8, LC_NAME=C, LC_ADDRESS=C, LC_TELEPHONE=C, LC_MEASUREMENT=en_US.UTF-8, LC_IDENTIFICATION=C
- Running under: Ubuntu 18.04.4 LTS
- Matrix products: default
- BLAS: /home/biocbuild/bbs-3.11-bioc/R/lib/libRblas.so
- LAPACK: /home/biocbuild/bbs-3.11-bioc/R/lib/libRlapack.so
- Base packages: base, datasets, grDevices, graphics, methods, stats, utils
- Other packages: BiocStyle 2.16.0, OmnipathR 1.2.0, dnet 1.1.7, dplyr 0.8.5, ggplot2 3.3.0, ggraph 2.0.2, gprofiler2 0.1.9, hexbin 1.28.1, igraph 1.2.5, supraHex 1.26.0, tidyr 1.0.2
- Loaded via a namespace (and not attached): BiocGenerics 0.34.0, BiocManager 1.30.10, MASS 7.3-51.6, Matrix 1.2-18, R6 2.4.1, RCurl 1.98-1.2, Rcpp 1.0.4.6, Rgraphviz 2.32.0, ape 5.3, assertthat 0.2.1, bitops 1.0-6, bookdown 0.18, cli 2.0.2, colorspace 1.4-1, compiler 4.0.0, crayon 1.3.4, data.table 1.12.8, digest 0.6.25, ellipsis 0.3.0, evaluate 0.14, fansi 0.4.1, farver 2.0.3, ggforce 0.3.1, ggrepel 0.8.2, glue 1.4.0, graph 1.66.0, graphlayouts 0.7.0, grid 4.0.0, gridExtra 2.3, gtable 0.3.0, htmltools 0.4.0, htmlwidgets 1.5.1, http 1.4.1, jsonlite 1.6.1, knitr 1.28, labeling 0.3, lattice 0.20-41, lazyeval 0.2.2, lifecycle 0.2.0, magick 2.3, magrittr 1.5, munsell 0.5.0, nlme 3.1-147, parallel 4.0.0, pillar 1.4.3, pkgconfig 2.0.3, plotly 4.9.2.1, polyclip 1.10-0, purrr 0.3.4, rlang 0.4.5, rmarkdown 2.1, scales 1.1.0, stats4 4.0.0, stringi 1.4.6, stringr 1.4.0, tibble 3.0.1, tidygraph 1.1.2, tidyselect 1.0.0, tools 4.0.0, tweenr 1.0.1, utf8 1.1.4, vctrs 0.2.4, viridis 0.5.1, viridisLite 0.3.0, withr 2.2.0, xfun 0.13, yaml 2.2.1