

KEGGgraph: a graph approach to KEGG PATHWAY in R and Bioconductor

Jitao David Zhang and Stefan Wiemann

April 16, 2010

Abstract

We demonstrate the capabilities of the *KEGGgraph* package, an interface between KEGG pathways and graph model in R as well as a collection of tools for these graphs. Superior to preceding approaches, *KEGGgraph* maintains the pathway topology and allows further analysis or dissection of pathway graphs. It parses the regularly updated KGML (KEGG XML) files into graph models maintaining all essential pathway attributes.

1 Introduction

Since its first introduction in 1995, KEGG PATHWAY has been widely used as a reference knowledge base for understanding biological pathways and functions of cellular processes. The knowledge from KEGG has proven of great value by numerous work in a wide range of fields [Kanehisa *et al.*, 2008].

Pathways are stored and presented as graphs on the KEGG server side, where nodes are molecules (protein, compound, etc) and edges represent relation types between the nodes, e.g. activation or phosphorylation. The graph nature raised our interest to investigate them with powerful graph tools implemented in R and Bioconductor [Gentleman *et al.*, 2004], including *graph*, *RBGL* and *Rgraphviz* [Carey *et al.*, 2005]. While it is barely possible to query the graph characteristics by manual parsing, a native and straightforward client-side tool is currently missing to parse pathways and analyze them consequently with tools for graph in R.

To address this problem, we developed the open-source software package *KEGGgraph*, an interface between KEGG pathway and graph object as well as a collection of tools to analyze, dissect and visualize these graphs.

The package requires KGML (KEGG XML) files, which can be downloaded from KEGG FTP site (<ftp://ftp.genome.jp/pub/kegg/xml>) without license permission for academic purposes. To demonstrate the functionality, in 'extdata/' sub-directory of *KEGGgraph* we have pre-installed several KGML files.

2 Software features

KEGGgraph offers the following functionalities:

Parsing: It should be noted that one 'node' in KEGG pathway does not necessarily map to merely one gene product, for example the node 'ERK' in the human TGF-Beta signaling pathway contains two homologues, MAPK1 and MAPK3. Therefore, among several parsing options, user can set whether to expand these nodes topologically. Beyond facilitating the interpretation of pathways in a gene-oriented manner, the approach also entitles unique identifiers to nodes, enabling merging graphs from different pathways.

Graph operations: Two common operations on graphs are subset and merge. A sub-graph of selected nodes and the edges in between are returned when subsetting, while merging produces a new graph that contains nodes and edges of individual ones. Both are implemented in *KEGGgraph*.

Visualization: *KEGGgraph* provides functions to visualize KEGG graphs with custom style. Nevertheless users are not restricted by them, alternatively they are free to render the graph with other tools like the ones in *Rgraphviz*.

Besides the functionalities described above, *KEGGgraph* also has tools for remote KGML file retrieval, graph feature study and other related tasks. We will demonstrate them later in this vignette.

3 Case studies

We load the *KEGGgraph* by typing or pasting the following codes in R command line:

```
> library(KEGGgraph)
```

3.1 Get KGML files

There are at least two possibilities to get KGML (KEGG XML) files:

- Manual download from KEGG FTP site at `ftp://ftp.genome.jp/pub/kegg/xml/`.
- Automatic retrieval from KEGG FTP site with the function `retrieveKGML`.

To retrieve KGML file automatically from KEGG FTP site, one has to know the pathway identifier (in the form of [a-z]3[0-9]5, where the three-alphabet code represent the organism and the five digits represent pathway). One method to find the mapping between pathway name and identifier is use `KEGGPATHNAME2ID` environment in *KEGG.db*. For example, the following codes retrieve p53 signaling pathway of *C.elegans* from KEGG FTP site.

```
> library(KEGG.db)
> tmp <- tempfile()
```

```
> pName <- "p53 signaling pathway"
> pId <- mget(pName, KEGGPATHNAME2ID)[[1]]
> retrieveKGML(pId, organism="cel", destfile=tmp, method="wget", quiet=TRUE)
```

Note: `retrieveKGML` uses a *try-download* mechanism (since the *KEGGgraph* version 1.1.2) to retrieve the KGML file from remote KEGG FTP server. Since August 2009 the KGML files are stored separately in different subdirectories depending on whether they record metabolic or non-metabolic pathways. Since from the KEGG pathway accession ID alone (e.g. hsa00010) alone it is not possible to determine its content, the `retrieveKGML` first tries to download the file from the non-metabolic subdirectory, and tries the metabolic directory in case no file was found in the non-metabolic category (the *try* step). In case the corresponding file is found, it is downloaded (the *download* step). Even if the file is found in the first try round, it still needs to be downloaded in the *download* step. However, this does not actually need to the network overhead, since thanks to the common cache system the file is only downloaded once.

3.2 Parsing and graph feature query

First we read in KGML file for human MAPK signaling pathway (with KEGG ID hsa04010):

```
> mapkKGML <- system.file("extdata/hsa04010.xml",
+                           package="KEGGgraph")
```

Once the file is ready, we can either parse them into an object of *KEGGPathway* or an object of *graph*. *KEGGPathway* object maintains the information of the pathway (title, link, organism, etc), while *graph* objects are more natural approach and can be directly plugged in many other tools. To convert the pathway into graph, we use

```
> mapkG <- parseKGML2Graph(mapkKGML, expandGenes=TRUE)
> mapkG
```

```
A graphNEL graph with directed edges
Number of Nodes = 265
Number of Edges = 876
```

Alternatively we can parse the KGML file first into an object of *KEGGpathway*, which can be later converted into the graph object, as the following lines show:

```
> mapkpathway <- parseKGML(mapkKGML)
> mapkpathway
```

```
KEGG Pathway
[ Title ]: MAPK signaling pathway
```

```
[ Name ]: path:hsa04010
[ Organism ]: hsa
[ Number ] :04010
[ Image ] :http://www.genome.jp/kegg/pathway/hsa/hsa04010.gif
[ Link ] :http://www.genome.jp/dbget-bin/show_pathway?hsa04010
```

Statistics:

```
    136 node(s)
    171 edge(s)
     0 reaction(s)
```

```
> mapkG2 <- KEGGpathway2Graph(mapkpathway, expandGenes=TRUE)
> mapkG2
```

A graphNEL graph with directed edges

Number of Nodes = 265

Number of Edges = 876

There is no difference between graph objects derived from two approaches.

The option 'expandGenes' in parsing controls whether the nodes of paralogues in pathways should be expanded or not. Since one 'node' in KEGG pathway does not necessarily map to only one gene/gene product (e.g. 'ERK' maps to MAPK1 and MAPK3), the option allows expanding these nodes and takes care of copying existing edges.

Another option users may find useful is 'genesOnly', when set to TRUE, the nodes of other types than 'gene' (compounds, for example) are neglected and the result graph consists only gene products. This is especially desired when we want to query network characteristics of gene products. Its value is set to 'TRUE' by default.

The following commands extract node and edge information:

```
> mapkNodes <- nodes(mapkG)
> nodes(mapkG)[1:3]
```

```
[1] "hsa:5923" "hsa:5924" "hsa:11072"
```

```
> mapkEdges <- edges(mapkG)
> edges(mapkG)[1]
```

```
$`hsa:5923`
```

```
[1] "hsa:22800" "hsa:22808" "hsa:3265" "hsa:3845" "hsa:4893" "hsa:6237"
```

Edges in KEGG pathways are directional, that is, an edge starting at node A pointing to node B does not guarantee a reverse relation, although reciprocal edges are also allowed.

When listing edges, a list indexed with node names is returned. Each item in the list records the nodes pointed to.

We can also extract the node attributes specified by KEGG with `getKEGGnodeData`:

```
> mapkGnodedata <- getKEGGnodeData(mapkG)
> mapkGnodedata[[2]]
```

KEGG Node (Entry 'hsa:5924'):

```
-----
[ displayName ]: RASGRF1, GRF1...
[ Name ]: hsa:5924
[ Type ]: gene
[ Link ]: http://www.genome.jp/dbget-bin/www_bget?hsa+5923+5924
-----
```

An alternative to use `getKEGGnodeData` is

```
> getKEGGnodeData(mapkG, 'hsa:5924')
```

, returning identical results.

Similarly the `getKEGGedgeData` is able to extract edge information:

```
> mapkGedgeData <- getKEGGedgeData(mapkG)
> mapkGedgeData[[4]]
```

KEGG Edge (Type: PPre1):

```
-----
[ Entry 1 ID ]: hsa:5923
[ Entry 2 ID ]: hsa:3845
[ Subtype ]:
  [ Subtype name ]: activation
  [ Subtype value ]: -->
-----
```

Alternatively the query above can be written as:

```
> getKEGGedgeData(mapkG, 'hsa:627~hsa:4915')
```

For `KEGGNode` and `KEGGEdge` objects, methods are implemented to fetch their attributes, for example `getName`, `getType` and `getDisplayName`. Guides to use these methods as well as examples can be found in help pages.

This case study finishes with querying the degree attributes of the nodes in the graph. We ask the question which nodes have the highest out- or in-degrees. Roughly speaking the out-degree (number of out-going edges) reflects the regulatory role, while the in-degree (number of in-going edges) suggests the subjectability of the protein to intermolecular regulations.

```

> mapkGoutdegrees <- sapply(edges(mapkG), length)
> mapkGindegrees <- sapply(inEdges(mapkG), length)
> topouts <- sort(mapkGoutdegrees, decreasing=T)
> topins <- sort(mapkGindegrees, decreasing=T)
> topouts[1:3]

```

```

hsa:5594 hsa:5595 hsa:1432
      26      26      13

```

```

> topins[1:3]

```

```

hsa:5923 hsa:5924 hsa:10125
      26      26      26

```

3.3 Graph subset and merge

We demonstrate the subsetting of the graph with 25 randomly chosen nodes of MAPK pathway graph:

```

> library(Rgraphviz)
> set.seed(123)
> randomNodes <- sample(nodes(mapkG), 25)
> mapkGsub <- subGraph(randomNodes, mapkG)
> mapkGsub

```

A graphNEL graph with directed edges

Number of Nodes = 25

Number of Edges = 9

The subgraph is visualized in figure 1, where nodes with in-degree or out-degree in red and others in grey.¹ And in the example we also demonstrate how to convert KEGG ID into other other identifiers via the Entrez GeneID. More details on the conversion of IDs can be found on page 13.

Another common operation on graphs is merging, that is, combining different graphs together. It is inspired by the fact that many KEGG pathways embed other pathway, for example MAPK signaling pathway embeds 6 pathways including Wnt signaling pathway. **mergeGraphs** provides the possibility to merge them into one graph for further analysis. Next we merge MAPK and Wnt signaling pathway into one graph. The graphs to be merged should be organized into a list, and it is commandary to use 'expandGenes=TRUE' option when parsing to make sure the nodes are unique and indexed by KEGGID.

¹The **makeAttr** function is used to assign nodes with rendering attributes, whose code can be found in the Rnw file.


```

> wntKGML <- system.file("extdata/hsa04310.xml", package="KEGGgraph")
> wntG <- parseKGML2Graph(wntKGML)
> graphs <- list(mapk=mapkG, wnt=wntG)
> merged <- mergeGraphs(graphs)
> merged

```

```

A graphNEL graph with directed edges
Number of Nodes = 386
Number of Edges = 1628

```

We observe that the node number in the merged graph (386) is less than the sum of two graphs (265 and 148 for MAPK and Wnt pathway respectively), reflecting the crosstalk between the pathways by sharing nodes.

3.4 Using other graph tools

In R and Bioconductor there are powerful tools for graph algorithms and operations, including *graph*, *Rgraphviz* and *RBGL*. The KEGG graphs can be analyzed with their functionalities to describe characteristics of the pathway and to answer biological relevant questions.

Here we demonstrate the use of other graph tools with asking the question which nodes are of the highest importance in MAPK signalling pathway. To this end we turn to relative betweenness centrality [Aittokallio and Schwikowski, 2006, Carey *et al.*, 2005]. Betweenness is a centrality measure of a node within a graph. Nodes that occur on many shortest paths between other vertices have higher betweenness than those that do not. It is scaled by the factor of $(n-1)(n-2)/2$ to get relative betweenness centrality, where n is the number of nodes in the graph. Both measurements estimate the importance or the role of the node in the graph.

With the function implemented in *RBGL*, our aim is to identify most important nodes in MAPK signalling pathway.

```

> library(RBGL)
> bcc <- brandes.betweenness centrality(mapkG)
> rbccs <- bcc$relative.betweenness centrality.vertices[1L,]
> toprbccs <- sort(rbccs, decreasing=TRUE)[1:4]
> toprbccs

```

```

hsa:4214 hsa:2885 hsa:5605 hsa:5604
0.2685233 0.2467144 0.2366450 0.2366450

```

We identify the top 4 important nodes judged by betweenness centrality as MAP3K1 (hsa:4214), GRB2 (hsa:2885), MAP2K2 (hsa:5605) and MAP2K1 (hsa:5604) (the mapping between KEGG ID and gene symbol is done using *biomaRt*, see page 14). In figure 2 we illustrate them as well as their interacting partners in MAPK pathway.


```

> toprbccName <- names(toprbccs)
> toprin <- sapply(toprbccName, function(x) inEdges(mapkG)[x])
> toprout <- sapply(toprbccName, function(x) edges(mapkG)[x])
> toprSubnodes <- unique(unname(c(unlist(toprin), unlist(toprout), toprbccName)))
> toprSub <- subGraph(toprSubnodes, mapkG)
> nAttrs <- list()
> tops <- c("MAPK3K1", "GRB2", "MAP2K2", "MAP2K1")
> topLabels <- lapply(toprbccName, function(x) x); names(topLabels) <- tops
> nAttrs$label <- makeAttr(toprSub, "", topLabels)
> nAttrs$fillcolor <- makeAttr(toprSub, "lightblue", list(orange=toprbccName))
> nAttrs$width <- makeAttr(toprSub, "", list("0.8"=toprbccName))
> plot(toprSub, "twopi", nodeAttrs=nAttrs, attrs=list(graph=list(start=2)))

```

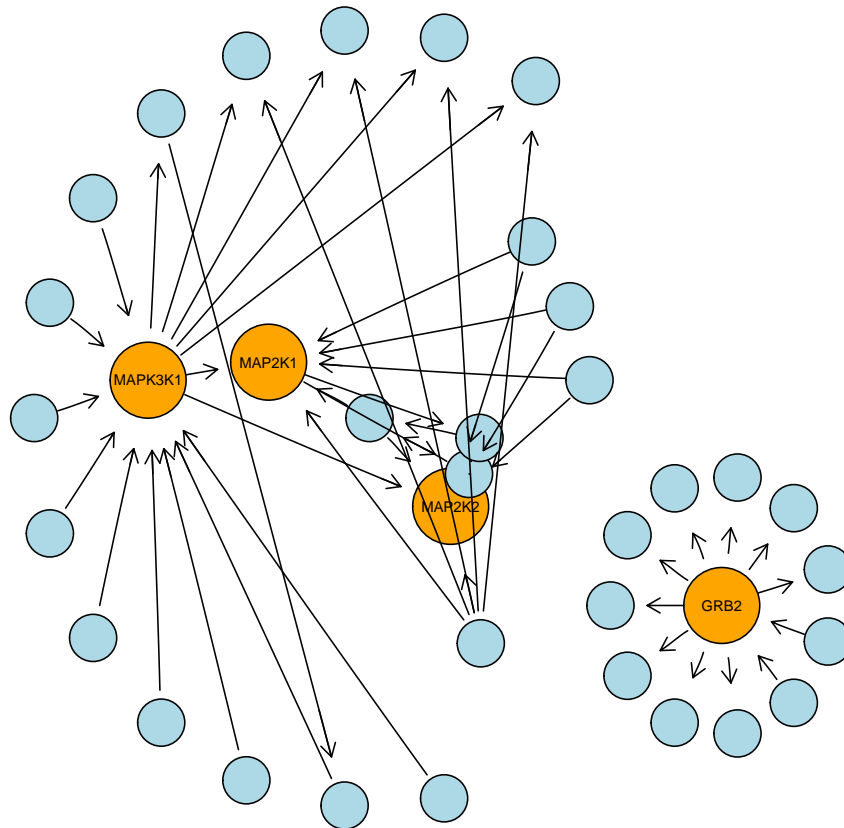


Figure 2: Nodes with the highest relative betweenness centrality in MAPK pathway (in orange) and their interacting partners (in blue).

4 Other functionalities

Besides the ability to parse and operate on KEGG PATHWAY graph objects, the *KEGG-graph* package also provides functionalities to complement tasks related to deal with KEGG pathways. We introduce some of them here, for a full list of functions please see the package help file:

```
> help(package=KEGGgraph)
```

4.1 Parsing chemical compound reaction network

KEGG PATHWAY captures two kinds of network: the protein network and the chemical network. The protein network consists *relations* (*edges*) between gene products, while the chemical network illustrates the *reactions* between chemical compounds. Since the metabolic pathway can be viewed both as a network of proteins (enzymes) and as a network of chemical compounds, metabolic pathways can be viewed as both protein networks and chemical networks, whereas regulatory pathways are always viewed as protein networks only. KEGGPathway provides methods to access this network.

We show the example of Glycine, serine and threonine metabolism pathway.

```
> mapfile <- system.file("extdata/map00260.xml", package="KEGGgraph")
> map <- parseKGML(mapfile)
> map
```

KEGG Pathway

```
[ Title ]: Glycine, serine and threonine metabolism
[ Name ]: path:map00260
[ Organism ]: map
[ Number ] :00260
[ Image ] :http://www.genome.jp/kegg/pathway/map/map00260.gif
[ Link ] :http://www.genome.jp/dbget-bin/show_pathway?map00260
```

Statistics:

```
    144 node(s)
    371 edge(s)
    68 reaction(s)
```

```
> reactions <- getReactions(map)
> reactions[[1]]
```

KEGG Reaction(rn:R08211)

```
[ Name ]: rn:R08211
[ Type ]: irreversible
[ Substrate Name ]: cpd:C00576
[ Product Name ]: cpd:C00719
```

Figure 3 shows how to extract reactions from the pathway and to build a directed graph with them.

4.2 Expand embedded pathways

Function `parseKGMLexpandMaps` is a function to handle with pathways embedding other pathways. For example, pancreatic cancer pathway embeds 9 other pathways including MAPK and ErbB signaling pathway, cell cycle and apoptosis pathway, etc. To parse them into one graph, the users only have to download the KGML file and feed the file name to `parseKGMLexpandMaps`, the function parses the file, analyze the embedded pathways, download their files from KEGG FTP site automatically (alternatively a local repository can be specified for KGML files) and merge the individual pathways into a single graph. For example, the following single line parses MAPK signaling pathway with all its embedded pathways

```
> mapkGembed <- parseKGMLexpandMaps(mapkKGML)
```

As its name suggests, function `subGraphByNodeType` subsets the graph by node type, the nodes to subset are those of the type given by the user. It is useful when the KGML file was parsed with 'genesOnly=FALSE' option and later on the user wants only certain kind of node, 'gene' for example, remained. The following example shows how to use it.

```
> mapkGall <- parseKGML2Graph(mapkKGML, genesOnly=FALSE)
> mapkGall
```

```
A graphNEL graph with directed edges
Number of Nodes = 277
Number of Edges = 891
```

```
> mapkGsub <- subGraphByNodeType(mapkGall, "gene")
> mapkGsub
```

```
A graphNEL graph with directed edges
Number of Nodes = 265
Number of Edges = 876
```

```

> chemicalGraph <- KEGGpathway2reactionGraph(map)
> outDegrees <- sapply(edges(chemicalGraph), length)
> maxout <- names(sort(outDegrees,decreasing=TRUE))[1:3]
> nAttrs <- list()
> maxoutlabel <- as.list(maxout); names(maxoutlabel) <- maxout
> nAttrs$label <- makeAttr(chemicalGraph, "", maxoutlabel)
> nAttrs$fillcolor <- makeAttr(chemicalGraph, "lightblue", list(orange=maxout))
> nAttrs$width <- makeAttr(chemicalGraph, "0.8", list("1.8"=maxout))
> plot(chemicalGraph, nodeAttrs=nAttrs)

```

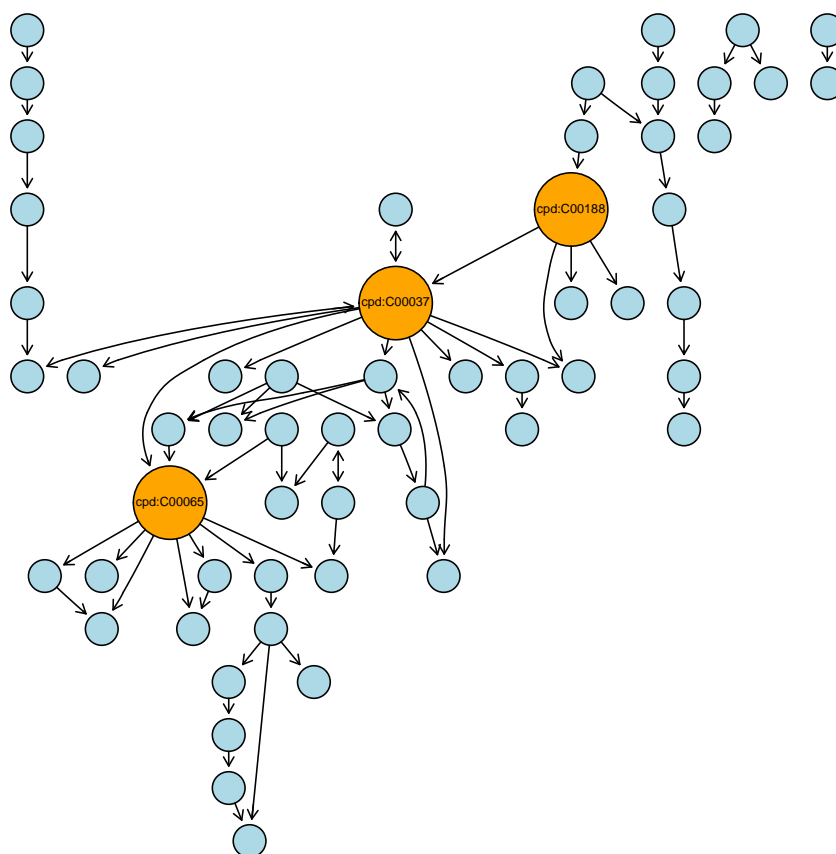


Figure 3: Reaction network built of chemical compounds: the orange nodes are the three compounds with maximum out-degree in this network.

4.3 Annotation

`translateKEGGID2GeneID` translates KEGG identifiers (KEGGID) into Entrez GeneID. For example, if we want to find the Entrez GeneID of the nodes in MAPK pathway having the highest relative betweenness centrality, the following codes do the job.

```
> toprbccKEGGID <- names(toprbccs)
> toprbccKEGGID

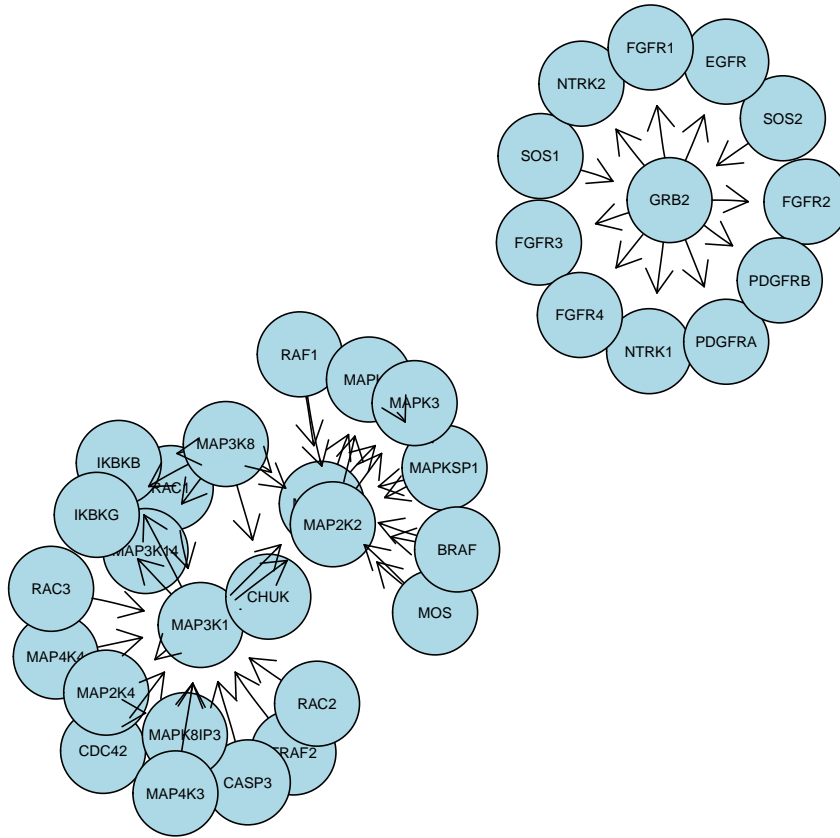
[1] "hsa:4214" "hsa:2885" "hsa:5605" "hsa:5604"

> toprbccGeneID <- translateKEGGID2GeneID(toprbccKEGGID)
> toprbccGeneID

[1] "4214" "2885" "5605" "5604"
```

To convert GeneID to other identifiers, we recommend genome wide annotation packages, for human it is *org.Hs.eg.db* and the packages for other organisms can be found at <http://www.bioconductor.org/packages/release/data/annotation/>. To demonstrate its use, we draw the sub-network in the figure 2 again, whereas nodes are now labeled with gene symbols.

```
> if(require(org.Hs.eg.db)) {
+   tnodes <- nodes(toprSub)
+   tgeneids <- translateKEGGID2GeneID(tnodes)
+   tgenesymbols <- sapply(mget(tgeneids, org.Hs.egSYMBOL, ifnotfound=NA), "[",1)
+   toprSubSymbol <- toprSub
+   nodes(toprSubSymbol) <- tgenesymbols
+   plot(toprSubSymbol, "neato",attrs=list(node=list(font=5, fillcolor="lightblue")))
+ }
```



Alternatively, users could use R package *biomaRt* [Durinck *et al.*, 2005, Durinck and Huber, 2008] for ID conversion, whereas it assumes that the user has an internet connection. The following example shows how to translate the node hits we acquired in the example above into HGNC symbols:

```
> library(biomaRt)
> hsapiens <- useMart("ensembl", "hsapiens_gene_ensembl" )
> filters <- listFilters(hsapiens)
> getBM(attributes=c("entrezgene", "hgnc_symbol"),
+       filters="entrezgene",
+       values=toprbccGeneID, mart=hsapiens)
```

5 Acknowledgement

We thank Vincent Carey, Holger Fröhlich and Wolfgang Huber for comments and suggestions on the package, and the reviewers from the Bioconductor community.

6 Conclusion

Before the release of *KEGGgraph*, several R/Bioconductor packages have been introduced and proven their usefulness in understanding biological pathways with KEGG. However, *KEGGgraph* is the first package able to parse any KEGG pathways from KGML files into graphs. In comparison, existing tools can not achieve the results we present here. They either neglects the graph topology (*KEGG.db*), do not parse pathway networks (*keggorth*), or are specialized for certain pathways (*cMAP* and *pathRender*).

With *KEGGgraph*, we contribute a direct and natural approach to KEGG pathways, and the possibilities to study them in R and Bioconductor.

7 Session Info

The script runs within the following session:

```
R version 2.11.0 RC (2010-04-16 r51754)
```

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

```
other attached packages:
```

```
[1] RBGL_1.23.0      org.Hs.eg.db_2.4.1  RSQLite_0.8-4
[4] DBI_0.2-5        AnnotationDbi_1.9.8 Biobase_2.7.6
[7] KEGGgraph_1.3.3  Rgraphviz_1.25.1    graph_1.25.25
[10] XML_2.8-1
```

References

[Gentleman *et al.*, 2004] Gentleman *et al.* (2004) Bioconductor: open software development for computational biology and bioinformatics, *Genome Biology*, **5**, R80.

- [Carey *et al.*, 2005] Carey *et al.* (2005) Network structures and algorithms in Bioconductor, *Bioinformatics*, **21**, 135-136.
- [Kanehisa *et al.*, 2008] Kanehisa *et al.* (2008) KEGG for linking genomes to life and the environment, *Nucleic Acids Research, Database issue*, **36**, 480-484.
- [Klukas and Schreiber, 2007] Klukas and Schreiber. (2007) Dynamic exploration and editing of KEGG pathway diagrams, *Bioinformatics*, **23**, 344-350.
- [Aittokallio and Schwikowski, 2006] Aittokallio and Schwikowski (2006) Graph-based methods for analysing networks in cell biology, *Briefings in Bioinformatics*, **7**, 243-255.
- [Durinck *et al.*, 2005] Durinck *et al.* (2005) BioMart and Bioconductor: a powerful link between biological databases and microarray data analysis, *Bioinformatics*, **21**, 3439-3440.
- [Durinck and Huber, 2008] Durinck and Huber (2008) R/Bioconductor package *biomaRt*, 2008