

**AnnotationHub 获取kegg org数据库：除了公开的19个之外，其他的也都可以获取、下载**

笔记本： R

创建时间： 2020/7/14 16:46

更新时间： 2020/12/7 11:06

作者： 干冰

URL: <https://www.bioinfo-scrounger.com/archives/512/>

## GO KEGG 富集分析数据库准备

注意：请使用最新的R下载

### (1) 加载R包、创建链接

```
library(AnnotationHub)
library(AnnotationDbi)
ah <- AnnotationHub()
```

### (2) 搜索org数据库

```
# 获取所有orgdb
org <- ah[ah$rdataclass == "OrgDb",]

# 搜索 物种
hm <- query(org, "Homo sapiens") # 人
hm <- query(org, "mellifera") # 蜜蜂
hm # 查看搜索结果
# 结果见下图，得到了很多个，第一个就是目标

# query的完整写法
# query(x, pattern, ignore.case=TRUE) pattern 是正则匹配
# 这里我找一个特殊的物种蜜蜂做示例 https://www.ncbi.nlm.nih.gov/genome/?
term=txid7460[orgn] Apis mellifera (honey bee)
```

```
> hm <- query(org, "mellifera")
> hm
AnnotationHub with 4 records
# snapshotDate(): 2018-10-24
# $dataprovder: ftp://ftp.ncbi.nlm.nih.gov/gene/DATA/
# $species: Apis mellifera, Apis mellifera cerana, Apis mellifera dorsata, A.
# $rdataclass: OrgDb
# additional mcalls(): taxonomyid, genome, description,
# coordinate_1_based, maintainer, rdatadateadded, preparerclass, tags,
# rdatapath, sourceurl, sourcetype
# retrieve records with, e.g., 'object[["AH67105"]]'

      title
AH67105 | org.Apis_mellifera.eg.sqlite
AH67153 | org.Apis_mellifera_cerana.eg.sqlite
AH67162 | org.Apis_mellifera_florea.eg.sqlite
AH67213 | org.Apis_mellifera_dorsata.eg.sqlite
>
```

### (3) 下载数据库

```
org_db <- ah[["AH67105"]] # 这一步会使用网络下载数据，缓存文件通常保存在个人
~/.AnnotationHub/id 文件中。见下方截图示例
```

```
> library(AnnotationDbi)
> ah <- AnnotationHub()
snapshotDate(): 2018-10-24
> org <- ah[ah$rdataclass == "OrgDb",]
> hm <- query(org, "Homo sapiens") # 人
> hm_org <- hm[[1]]
downloading 1 resources
retrieving 1 resource
|=====100%
loading from cache
'/home/ganb//.AnnotationHub/72902'
```

这个文件就是当前物种的数据库文件  
可以直接拷贝出来，重命名，用  
loadDb加载

## (4) 数据库保存、加载 (saveDb不能用了，报错，原因未知，此时请采用4.1备选方案)

```
# 保存到文件，下次直接加载即可
saveDb(org_db, file = "mellifera.orgdb")

# 加载
org_db = loadDb(file = "mellifera.orgdb")
```

### (4.1) 备选数据保存方案

```
cp ~/.AnnotationHub/id abc.orgdb # 直接拷贝数据库缓存文件
org_db = loadDb(file = "abc.orgdb")
```

文件名建议统一命名为：物种名称.kegg简称.orgdb

escherichia\_coli.eco.orgdb 大肠杆菌  
honey\_bee\_apis\_mellifera.ame.orgdb 蜜蜂  
human.hsa.orgdb 人  
mouse.mmu.orgdb 小鼠  
rabbit\_oryctolagus\_cuniculus.ocu.orgdb 兔子  
rat.rno.orgdb 大鼠  
zea\_mays.zma.orgdb 玉米  
zebrafish\_danio\_rerio.dre.orgdb 斑马鱼

### 数据库相关操作

- columns(org\_db)
  - 查看数据库包含哪些信息
- head(keys(org\_db, keytype = "SYMBOL"))
  - 获取所有SYMBOL信息

```
> columns(org_db)
[1] "ACCNUM" "ALIAS" "CHR" "ENTREZID" "EVIDENCE"
[6] "EVIDENCEALL" "GENENAME" "GID" "GO" "GOALL"
[11] "ONTOLOGY" "ONTOLOGYALL" "PMID" "REFSEQ" "SYMBOL"
[16] "UNIGENE"
```

```
> head(keys(org_db, keytype = "SYMBOL"), 20)
[1] "14-3-3zeta" "18-w" "18S rRNA" "28S rRNA" "5-HT1"
[6] "5-HT2alpha" "5-HT2beta" "5-ht7" "A4" "ACSF2"
[11] "AChE-2" "AGLU2" "AQP" "ATP5G2" "Abcam"
[16] "Ac3" "Acph-1" "Ada2b" "Adar" "Adk1"
```

**注意：要完成GO/KEGG 分析，org数据库要包含  
SYMBOL/ENTREZID/GO 这三个信息  
SYMBOL/ENTREZID**

- 这两个信息主要用来把输入基因SYMBOL转化为ENTREZID (ENTREZID编号唯一，且GO/KEGG富集分析用的都是用这个编号，而不是SYMBOL)
  - **注意：有的物种不支持ENTREZID转换，下面做详细说明**
- SYMBOL严格区分大小写，一定要保证与NCBI一致

**GO:**

- GO富集分析要使用
- 查看kegg对应物种的ncbi-geneid转换api能否打开，例如：

<http://rest.kegg.jp/conv/hsa/ncbi-geneid>

- 能打开，一切正常
- 不能打开，则需要换一种映射方式了，请仔细阅读后面的资料

常见的几个数据库  
人

```
> hm <- query(org, "Homo sapiens") # 人
> hm
AnnotationHub with 1 record
# snapshotDate(): 2018-10-24
# names(): AH66156
# $dataProvider: ftp://ftp.ncbi.nlm.nih.gov/gene/DATA/
# $species: Homo sapiens
# $rdataclass: OrgDb
# $rdataadded: 2018-10-22
# $title: org.Hs.eg.db.sqlite
# $description: NCBI gene ID based annotations about Homo sapiens
# $taxonomyid: 9606
# $genome: NCBI genomes
# $sourcetype: NCBI/ensembl
# $sourceurl: ftp://ftp.ncbi.nlm.nih.gov/gene/DATA/, ftp://ftp.ensembl.org/p...
# $sourcesize: NA
# $tags: c("NCBI", "Gene", "Annotation")
# retrieve record with 'object[["AH66156"]]'
>
```

小鼠

```
# retrieve records with, e.g., 'object[["AH66157"]]'

title
AH66157 | org.Mm.eg.db.sqlite
AH66327 | org.Musa_AA_Group.eg.sqlite
AH66328 | org.Musa_acuminata.eg.sqlite
AH66329 | org.Musa_acuminata_AA_Group.eg.sqlite
AH66330 | org.Musa_papa.eg.sqlite
```

大鼠

```
> hm <- query(org, "rattus") # 人
> hm
AnnotationHub with 1 record
# snapshotDate(): 2018-10-24
# names(): AH66159
# $dataProvider: ftp://ftp.ncbi.nlm.nih.gov/gene/DATA/
# $species: Rattus norvegicus
# $rdataclass: OrgDb
# $rdataadded: 2018-10-22
# $title: org.Rn.eg.db.sqlite
# $description: NCBI gene ID based annotations about Rattus norvegicus
# $taxonomyid: 10116
# $genome: NCBI genomes
# $sourcetype: NCBI/ensembl
# $sourceurl: ftp://ftp.ncbi.nlm.nih.gov/gene/DATA/, ftp://ftp.ensembl.org/p...
# $sourcesize: NA
# $tags: c("NCBI", "Gene", "Annotation")
# retrieve record with 'object[["AH66159"]]'
> org_db = ah[["AH66159"]]
downloading 1 resources
retrieving 1 resource
=====| 100%

loading from cache
'/home/ganb//.AnnotationHub/72905'
> org_db
org_db
```

## 关于KEGG富集分析支持的输入ID说明

kegg支持的物种列表: [https://www.genome.jp/kegg/catalog/org\\_list.html](https://www.genome.jp/kegg/catalog/org_list.html)

kegg物种基因id使用: <https://www.genome.jp/kegg/genes.html>

### Data Source of KEGG GENES

The following table shows the data source of the KEGG GENES database.

Category	Original DB <sup>1</sup>	Content <sup>2</sup>	Genome identifier	Gene identifier
Eukaryotes	RefSeq	RefSeq release (complete)	T0 numbers (three or four letter organism codes)	GeneID
Prokaryotes	RefSeq	NCBI reference genomes		Locus_tag
	GenBank	Other complete genomes		Locus_tag
Viruses	RefSeq	Refseq release (viral)	T40000 (vg) T4 numbers	GeneID
Addendum	KEGG	Functionally characterized proteins	T10000 (ag)	ProteinID

<sup>1</sup> Original DB name is shown in the definition field of each GENES entry.

<sup>2</sup> RefSeq bimonthly releases are used to update eukaryotes and viruses.

Prokaryoteic genomes are selected from [ftp://ftp.ncbi.nlm.nih.gov/genomes/ASSEMBLY\\_REPORTS/](ftp://ftp.ncbi.nlm.nih.gov/genomes/ASSEMBLY_REPORTS/).

(1) 通常情况下, 富集分析时, 我们输入的都是基因symbol

(2) GO/KEGG 基本不支持symbol的识别, 因为symbol存在别名问题

- 那么就得换一个id作为唯一识别码

(3) 我们需要org\_db数据库, 把symbol映射为GO/KEGG可识别的标签, 通常情况下, 就是

**ENTREZID**

- 查看物种在kegg中使用的基因名称（可以看到，都是一个数字编号，但是也有特殊情况，例如下面的bmyo）

- 人 <http://rest.kegg.jp/list/hsa>

```
hsa:65249      ZSWIM4: zinc finger SWIM-type containing 4
hsa:100462981 MIRNR2L2, HN2; MT-RNR2 like 2
hsa:100526771 SMIM35, TMFRS54-A51; small integral membrane protein 35
hsa:3441      IFNA4, IFN-alpha4a, INFA4; interferon alpha 4
hsa:317754    POTE, A26B3, ANKRD21, CT104.1, POTE, POTE-21, POTE21; POTE ankyrin domain family member D
hsa:149685    ADIG, SMAF1; adipogenin
hsa:134864    TAAR1, TAI, TARI, TRAR1; trace amine associated receptor 1
hsa:9720      CCDC144A; coiled-coil domain containing 144A
hsa:3442      IFNA5, IFN-alpha-5, IFN-alphaG, INA5, INFA5, leIF.G; interferon alpha 5
hsa:84210     ANKRD20A1, ANKRD20A; ankyrin repeat domain 20 family member A1
hsa:101928917 HSF3; heat shock transcription factor family, X-linked member 3
hsa:139741    ACTRT1, AIP1, ARIP1, ARPT1, HSD27; actin related protein T1
hsa:100302174 MIR1307, MIRN1307, hsa-mir-1307, mir-1307; microRNA 1307
```

- 小鼠 <http://rest.kegg.jp/list/mmu>

- 大鼠 <http://rest.kegg.jp/list/rno>

- 蕈状芽孢杆菌 <http://rest.kegg.jp/list/bmyo>

```
bmyo:BG05_1    glutamine amidotransferases class-II family protein
bmyo:BG05_2    aluminum activated malate transporter family protein
bmyo:BG05_3    ABC transporter family protein
bmyo:BG05_4    hypothetical protein
bmyo:BG05_5    ygaB-like family protein
bmyo:BG05_6    small, acid-soluble spore, gamma-type family protein
bmyo:BG05_7    putative yfhS
bmyo:BG05_8    mutY; A/G-specific adenine glycosylase
bmyo:BG05_9    hypothetical protein
bmyo:BG05_10   WVLL family protein
bmyo:BG05_11   sspK; small, acid-soluble spore protein K
bmyo:BG05_12   ypzG-like family protein
bmyo:BG05_13   hypothetical protein
bmyo:BG05_14   recX family protein
bmyo:BG05_15   NAD dependent epimerase/dehydratase family protein
bmyo:BG05_16   amidohydrolase family protein
bmyo:BG05_17   acetyltransferase domain protein
bmyo:BG05_18   yfhE-like family protein
bmyo:BG05_19   ugtP; processive diacylglycerol glucosyltransferase
bmyo:BG05_20   pflA; pyruvate formate-lyase 1-activating enzyme
bmyo:BG05_21   pflB; formate acetyltransferase
bmyo:BG05_22   hypothetical protein
bmyo:BG05_23   glycosyltransferase like 2 family protein
bmyo:BG05_24   3-beta hydroxysteroid dehydrogenase/isomerase family protein
bmyo:BG05_25   nucleotide sugar dehydrogenase family protein
bmyo:BG05_26   glycosyltransferase like 2 family protein
```

- 通常情况下，KEGG使用的唯一基因标签是ENTREZID，这个ID就是ncbi-geneid，这个ID是唯一的，不会因为别名改变而改变

See also 134 discontinued or replaced items.

Name/Gene ID	Description	Location	Aliases	MIM
<input type="checkbox"/> TP53 ID: 7157	tumor protein p53 [Homo sapiens (human)]	Chromosome 17, NC_000017.11 (7668421..7687490, complement)	BCC7, BMFS5, LFS1, P53, TRP53	191170
<input type="checkbox"/> TP53 ID: 24842	tumor protein p53 [Rattus norvegicus (Norway rat)]	Chromosome 10, NC_005109.4 (56186299..56198449)	Trp53, p53	
<input type="checkbox"/> tp53 ID: 30590	tumor protein p53 [Danio rerio (zebrafish)]	Chromosome 5, NC_007116.7 (24086227..24097807)	brp, brp53, drp, drp53, etlD22686, etlD22686.5, fb40d06, p5, p53, wu.fb40d06, zgc:111919	
<input type="checkbox"/> TP53 ID: 403869	tumor protein p53 [Canis lupus familiaris (dog)]	Chromosome 5, NC_006587.3 (32561406..32565149, complement)	P53	

- 使用 <http://rest.kegg.jp/conv/hsa/ncbi-geneid> 可以查询物种在kegg中，keggid与ncbi-geneid的映射关系，可以发现，数字编号id是完全一样的，那么，我们就可以直接使用ENTREZID做富集分了

```
← → × rest.kegg.jp/conv/hsa/ncbi-geneid
应用 NCBI NCBI GEO UniProt 天星 天星软件 GitHub Linux命令大全 linux

ncbi-geneid: 1 hsa: 1
ncbi-geneid: 10 hsa: 10
ncbi-geneid: 100 hsa: 100
ncbi-geneid: 1000 hsa: 1000
ncbi-geneid: 10000 hsa: 10000
ncbi-geneid: 100008586 hsa: 100008586
ncbi-geneid: 100008587 hsa: 100008587
ncbi-geneid: 100008588 hsa: 100008588
ncbi-geneid: 100008589 hsa: 100008589
ncbi-geneid: 10001 hsa: 10001
ncbi-geneid: 10002 hsa: 10002
ncbi-geneid: 10003 hsa: 10003
ncbi-geneid: 100033411 hsa: 100033411
ncbi-geneid: 100033413 hsa: 100033413
ncbi-geneid: 100033414 hsa: 100033414
ncbi-geneid: 100033415 hsa: 100033415
ncbi-geneid: 100033416 hsa: 100033416
ncbi-geneid: 100033417 hsa: 100033417
ncbi-geneid: 100033418 hsa: 100033418
```

- 但是，有的物种，在kegg中是不存在ncbi-geneid映射关系的，例如前面说的细菌 bmyo, <http://rest.kegg.jp/conv/bmyo/ncbi-geneid> 你会发现无法访问。

- 面对非ncbi-geneid的情况，你就不能用ENTREZID了，需要换一种ID。这就需要你自己找一下org\_db，哪一个库中是keggid

- 这里以bmyo为例

1. 查看包含的数据库 columns(org\_db)

```
> columns(org_db)
[1] "ACCNUM" "ALIAS" "ENTREZID" "EVIDENCE" "EVIDENCEALL"
[6] "GENENAME" "GID" "GO" "GOALL" "ONTOLOGY"
[11] "ONTOLOGYALL" "PMID" "REFSEQ" "SYMBOL"
```

2. 查看数据库的中内容 head(keys(org\_db, keytype = "ALIAS"))

```
> head(keys(org_db, keytype = "ALIAS"))
[1] "BG05_10" "BG05_100" "BG05_1000" "BG05_1001" "BG05_1002" "BG05_1003"
```

3. 发现，上面的ALIAS就是KEGGID，那么，富集分析时，上面的id就是转换目标

## KEGG目前的总结：

1. 查看物种在kegg里的kegg id <http://rest.kegg.jp/list/hsa>
2. 查看物种在kegg中是否存在ncbi-geneid 映射关系 <http://rest.kegg.jp/conv/hsa/ncbi-geneid>
  1. 如果有，那么基本上id映射方式就是 ENTREZID
  2. 如果没有，那么，就需要打开org\_db，查一下哪个数据库能与 list/物种 列表里的keggid 对应上
    1. 目前来看，通常是ALIAS数据库对应kegg 编号
3. 注意：并不是所有kegg列出来的物种，都有org 数据库，只是某一个菌群存在数据库，例如 Rahnella只有 raa有orgdb

Chania	sfo	Chania multitudinisentens	2014	GenBank
	rah	Rahnella sp. Y9602	2011	GenBank
	raq	Rahnella aquatilis CIP 78.65 = ATCC 33071	2012	GenBank
Rahnella	raa	Rahnella aquatilis HX2	2012	GenBank
	rox	Rahnella sp. ERM1:05	2018	GenBank

1. 所以，在创建物种数据库前，先看一下kegg 支持哪一个菌群

end