

Using the NHLBI GRASP repository of GWAS test results with Bioconductor

Vincent J. Carey

October 9, 2014

Contents

1	Introduction	1
2	Installation	2
3	Demonstration	2
3.1	Attachment and messaging	2
3.2	Indexed p-value bins	2
3.3	Tabulations	3
3.4	Inspecting some relatively weak associations in asthma	4
4	Quick view of the basic interfaces	4
4.1	dplyr-oriented	4
4.2	RSQLite-oriented	4
5	Some QC (Consistency check): Are NHGRI GWAS catalog loci included?	5

1 Introduction

GRASP (Genome-Wide Repository of Associations Between SNPs and Phenotypes) v2.0 was released in September 2014. The [primary GRASP web resource](#) includes links to a web-based query interface. This document describes a Bioconductor package that replicates information in the v2.0 [textual release](#).

The primary reference for version 2 is: Eicher JD, Landowski C, Stackhouse B, Sloan A, Chen W, Jensen N, Lien J-P, Leslie R, Johnson AD (2014) GRASP v 2.0: an update to the genome-wide repository of associations between SNPs and phenotypes. Nucl Acids Res, published online Nov 26, 2014 PMID 25428361

From the main web page:

GRASP includes all available genetic association results from papers, their supplements and web-based content meeting the following guidelines:

- All associations with $P < 0.05$ from GWAS defined as $\geq 25,000$ markers tested for 1 or more traits.
- Study exclusion criteria: CNV-only studies, replication/follow-up studies testing $< 25K$ markers, non-human only studies, article not in English, gene-environment or gene-gene GWAS where single SNP main effects are not given, linkage only studies, aCGH/LOH only studies, heterozygosity/homozygosity (genome-wide or long run) studies, studies only presenting gene-based or pathway-based results, simulation-only studies, studies which we judge as redundant with prior studies since they do not provide significant inclusion of new samples or exposure of new results (e.g., many methodological papers on the WTCCC and FHS GWAS).
- More detailed methods and resources used in constructing the catalog are described at [Methods & Resources](#).
- Terms of Use for GRASP data: <http://apps.nhlbi.nih.gov/Grasp/Terms.aspx>
- Medical disclaimer: [<http://apps.nhlbi.nih.gov/Grasp/Overview.aspx>] (<http://apps.nhlbi.nih.gov/Grasp/Overview.aspx>)

2 Installation

Install the package in Bioconductor version 3.1 or later using `BiocInstaller::biocLite("grasp2db")`.

The first time the `grasp2` data base is referenced, via the `GRASP2()` function described below, a large (5.3Gb) file is downloaded to a local cache using [AnnotationHub](#). This may take a considerable length of time (10's of minutes, perhaps an hour or more depending on internet connections). Subsequent uses refer to the locally cached file, and are fast.

3 Demonstration

3.1 Attachment and messaging

Attach the package.

```
library(grasp2db)
```

Workflows start with a reference to the data base.

```
grasp2 <- GRASP2()
## snapshotDate(): 2017-06-29
## loading from cache '/home/biocadmin/.AnnotationHub/25509'
grasp2
## src:  sqlite 3.19.3 [/home/biocadmin/.AnnotationHub/25509]
## tbls:  count, study, variant
```

There are three tables. The 'variant' table summarizes NA variants. The 'study' table contains NA citations from which the data are extracted; the 'variant' and 'study' tables are related by PMID identifier. The 'count' table contains NA records summarizing SNPs observed in Discovery and Replication samples from 12 distinct Populations; the 'variant' and 'count' tables are related by NHLBIkey.

3.2 Indexed p-value bins

The database has been indexed on a number of fields, and an integer rounding of $-\log_{10}$ of the quantity recorded as the Pvalue of the association is available.

```
variant <- tbl(grasp2, "variant")
q1 = (variant %>% select(Pvalue, NegativeLog10PBin) %>%
      filter(NegativeLog10PBin > 8) %>%
      summarize(maxp = max(Pvalue), n=n()))
q1
## # Source:   lazy query [?? x 2]
## # Database:  sqlite 3.19.3 [/home/biocadmin/.AnnotationHub/25509]
##           maxp      n
##           <dbl> <int>
## 1 3.16181e-09 322794
```

This shows that the quantities in `NegativeLog10PBin` are upper bounds on the exponents of the p-values in the integer-labeled bins defined by this quantity.

A useful utility from `dplyr` is the query explain method:

```
explain(q1)
## <SQL>
## SELECT MAX(`Pvalue`) AS `maxp`, COUNT() AS `n`
```

```
## FROM (SELECT `Pvalue` AS `Pvalue`, `NegativeLog10PBin` AS `NegativeLog10PBin`
## FROM `variant`)
## WHERE (`NegativeLog10PBin` > 8.0)
##
## <PLAN>
##   addr      opcode p1      p2 p3      p4 p5 comment
## 1      0        Init  0      19 0        00      NA
## 2      1        Null  0      1 3        00      NA
## 3      2      OpenRead 2      2 0        10 00     NA
## 4      3      OpenRead 3 5465555 0      k(2,,) 00     NA
## 5      4        Real  0      4 0        8 00     NA
## 6      5      Affinity 4      1 0        D 00     NA
## 7      6      SeekGT  3     14 4        1 00     NA
## 8      7        Seek  3      0 2        00      NA
## 9      8      Column  2      8 5        00      NA
## 10     9 RealAffinity 5      0 0        00      NA
## 11    10      CollSeq 0      0 0 (BINARY) 00     NA
## 12    11      AggStep0 0      5 1      max(1) 01     NA
## 13    12      AggStep0 0      0 2      count(0) 00     NA
## 14    13        Next  3      7 0        00      NA
## 15    14      AggFinal 1      1 0      max(1) 00     NA
## 16    15      AggFinal 2      0 0      count(0) 00     NA
## 17    16        Copy  1      6 1        00      NA
## 18    17      ResultRow 6      2 0        00      NA
## 19    18        Halt  0      0 0        00      NA
## 20    19      Transaction 0      0 11      0 01     NA
## 21    20        Goto  0      1 0        00      NA
```

3.3 Tabulations

This query select variants with large effect from the 'variant' table, and joins them to their published phenotypic effect in the 'study' table.

```
study <- tbl(grasp2, "study")
large_effect <-
  variant %>% select(PMID, SNPId_dbSNP134, NegativeLog10PBin) %>%
  filter(NegativeLog10PBin > 5)
phenotype <-
  left_join(large_effect,
            study %>% select(PMID, PaperPhenotypeDescription))
## Joining, by = "PMID"
phenotype
## # Source:   lazy query [?? x 4]
## # Database: sqlite 3.19.3 [/home/biocadmin/.AnnotationHub/25509]
##   PMID SNPId_dbSNP134 NegativeLog10PBin PaperPhenotypeDescription
##   <chr>      <int>          <int>          <chr>
## 1 20502693      253            6 Gene expression in monocytes
## 2 19913121      255            6 Lipid level measurements
## 3 19913121      256            6 Lipid level measurements
## 4 19913121      263            6 Lipid level measurements
## 5 19913121      264            6 Lipid level measurements
## 6 19913121      271            6 Lipid level measurements
## 7 19913121      285            6 Lipid level measurements
```

```
## 8 19913121      301      6      Lipid level measurements
## 9 20502693      326      6      Gene expression in monocytes
## 10 20502693     327      6      Gene expression in monocytes
## # ... with more rows
```

3.4 Inspecting some relatively weak associations in asthma

```
lkaw <- semi_join(
  variant %>%
    filter(NegativeLog10PBin <= 4) %>%
    select(PMID, chr_hg19, SNPid_dbSNP134, PolyPhen2),
  study %>% filter(PaperPhenotypeDescription == "Asthma")
## Joining, by = "PMID"
```

We materialize the filtered table into a data.frame and check how many PolyPhen2 notations including a substring of 'Damaging':

```
lkaw %>% filter(PolyPhen2 %like% "%amaging%")
## # Source:   lazy query [?? x 4]
## # Database:  sqlite 3.19.3 [/home/biocadmin/.AnnotationHub/25509]
##   PMID chr_hg19 SNPid_dbSNP134
##   <chr> <chr> <int>
## 1 20860503 1 4762
## 2 20860503 1 880633
## 3 20860503 1 880633
## 4 20860503 3 1053338
## 5 20860503 19 1054940
## 6 20860503 1 1060622
## 7 20860503 1 1156281
## 8 20860503 13 1536207
## 9 20860503 10 1799853
## 10 20860503 11 2186797
## # ... with more rows, and 1 more variables: PolyPhen2 <chr>
```

4 Quick view of the basic interfaces

4.1 dplyr-oriented

```
grasp2
## src:  sqlite 3.19.3 [/home/biocadmin/.AnnotationHub/25509]
## tbls: count, study, variant
```

4.2 RSQLite-oriented

```
gcon = grasp2$con
library(RSQLite)
gcon
## <SQLiteConnection>
```

```
## Path: /home/biocadmin/.AnnotationHub/25509
## Extensions: TRUE
dbListTables(gcon)
## [1] "count" "study" "variant"
```

Note that the package opens the SQLite data base in 'read-only' mode, but updates are possible (e.g., directly opening a connection to the data base without restricting access). There may be implicit control if the user does not have write access to the file.

5 Some QC (Consistency check): Are NHGRI GWAS catalog loci included?

We have an image of the NHGRI GWAS catalog inheriting from GRanges.

```
library(gwascat)
## Loading required package: Homo.sapiens
## Loading required package: AnnotationDbi
## Loading required package: stats4
## Loading required package: BiocGenerics
## Loading required package: parallel
##
## Attaching package: 'BiocGenerics'
## The following objects are masked from 'package:parallel':
##
##   clusterApply, clusterApplyLB, clusterCall, clusterEvalQ,
##   clusterExport, clusterMap, parApply, parCapply, parLapply,
##   parLapplyLB, parRapply, parSapply, parSapplyLB
## The following objects are masked from 'package:dplyr':
##
##   combine, intersect, setdiff, union
## The following objects are masked from 'package:stats':
##
##   IQR, mad, sd, var, xtabs
## The following objects are masked from 'package:base':
##
##   Filter, Find, Map, Position, Reduce, anyDuplicated, append,
##   as.data.frame, cbind, colMeans, colSums, colnames, do.call,
##   duplicated, eval, evalq, get, grep, grepl, intersect,
##   is.unsorted, lapply, lengths, mapply, match, mget, order,
##   paste, pmax, pmax.int, pmin, pmin.int, rank, rbind, rowMeans,
##   rowSums, rownames, sapply, setdiff, sort, table, tapply,
##   union, unique, unsplit, which, which.max, which.min
## Loading required package: Biobase
## Welcome to Bioconductor
##
##   Vignettes contain introductory material; view with
##   'browseVignettes()'. To cite Bioconductor, see
##   'citation("Biobase")', and for packages 'citation("pkgnam")'.
## Loading required package: IRanges
## Loading required package: S4Vectors
##
## Attaching package: 'S4Vectors'
## The following objects are masked from 'package:dplyr':
```

```
##
## first, rename
## The following object is masked from 'package:base':
##
## expand.grid
##
## Attaching package: 'IRanges'
## The following objects are masked from 'package:dplyr':
##
## collapse, desc, slice
##
## Attaching package: 'AnnotationDbi'
## The following object is masked from 'package:dplyr':
##
## select
## Loading required package: OrganismDbi
## Loading required package: GenomicFeatures
## Loading required package: GenomeInfoDb
## Loading required package: GenomicRanges
## No methods found in "RSQLite" for requests: dbGetQuery
## Loading required package: GO.db
##
## Loading required package: org.Hs.eg.db
##
## Loading required package: TxDb.Hsapiens.UCSC.hg19.knownGene
## gwascat loaded. Use data(ebicat38) for hg38 coordinates;
## data(ebicat37) for hg19 coordinates.
data(gwrngs19) # hg19 addresses; NHGRI ships hg38
gwrngs19
## gwasloc instance with 17254 records and 35 attributes per record.
## Extracted: Mon Sep 8 13:08:13 2014
## Genome: hg19
## Excerpt:
## GRanges object with 5 ranges and 35 metadata columns:
##      seqnames      ranges strand | Date.Added.to.Catalog
##      <Rle>         <IRanges> <Rle> | <character>
## 1 chr19 [ 7739177, 7739177] * | 04/16/2014
## 2 chr6 [ 32626601, 32626601] * | 08/02/2014
## 3 chr4 [ 38799710, 38799710] * | 08/02/2014
## 4 chr5 [110467499, 110467499] * | 08/02/2014
## 5 chr2 [102966549, 102966549] * | 08/02/2014
##      PUBMEDID First.Author      Date      Journal
##      <integer> <character> <character> <character>
## 1 24123702 Chung CM 03/03/2014 Diabetes Metab Res Rev
## 2 24388013 Ferreira MA 12/30/2013 J Allergy Clin Immunol
## 3 24388013 Ferreira MA 12/30/2013 J Allergy Clin Immunol
## 4 24388013 Ferreira MA 12/30/2013 J Allergy Clin Immunol
## 5 24388013 Ferreira MA 12/30/2013 J Allergy Clin Immunol
##      Link
##      <character>
## 1 http://www.ncbi.nlm.nih.gov/pubmed/24123702
## 2 http://www.ncbi.nlm.nih.gov/pubmed/24388013
## 3 http://www.ncbi.nlm.nih.gov/pubmed/24388013
## 4 http://www.ncbi.nlm.nih.gov/pubmed/24388013
```

```

## 5 http://www.ncbi.nlm.nih.gov/pubmed/24388013
##
##
## 1 Common quantitative trait locus downstream of RETN gene identified by genome-wide association study
## 2 Genome-wide
## 3 Genome-wide
## 4 Genome-wide
## 5 Genome-wide
##      Disease.Trait
##      <character>
## 1      Resistin levels
## 2 Asthma and hay fever
## 3 Asthma and hay fever
## 4 Asthma and hay fever
## 5 Asthma and hay fever
##      Initial.Sample.Size
##      <character>
## 1      382 Han Chinese ancestry individuals
## 2 6,685 European ancestry cases, 14,091 European ancestry controls
## 3 6,685 European ancestry cases, 14,091 European ancestry controls
## 4 6,685 European ancestry cases, 14,091 European ancestry controls
## 5 6,685 European ancestry cases, 14,091 European ancestry controls
##      Replication.Sample.Size
##      <character>
## 1      559 Han Chinese ancestry individuals
## 2 878 European ancestry cases, 2,455 European ancestry controls
## 3 878 European ancestry cases, 2,455 European ancestry controls
## 4 878 European ancestry cases, 2,455 European ancestry controls
## 5 878 European ancestry cases, 2,455 European ancestry controls
##      Region      Chr_id Chr_pos.hg38 Reported.Gene.s.      Mapped_gene
##      <character> <character> <numeric> <character> <character>
## 1      19p13.2      19      7674291      RETN      RETN - C19orf59
## 2      6p21.32      6      32658824      HLA-DQB1      TRNAI25
## 3      4p14      4      38798089      TLR1      TLR1
## 4      5q22.1      5      111131801      WDR36      WDR36 - RPS3AP21
## 5      2q12.1      2      102350089      IL1RL1      IL1RL1
##      Upstream_gene_id Downstream_gene_id Snp_gene_ids
##      <character> <character> <character>
## 1      56729      199675
## 2      <NA>      <NA>      100189401
## 3      <NA>      <NA>      7096
## 4      134430      402287
## 5      <NA>      <NA>      9173
##      Upstream_gene_distance Downstream_gene_distance
##      <character> <character>
## 1      3.84      2.77
## 2      <NA>      <NA>
## 3      <NA>      <NA>
## 4      1.3      60.38
## 5      <NA>      <NA>
##      Strongest.SNP.Risk.Allele      SNPs      Merged      Snp_id_current
##      <character> <character> <character> <character>
## 1      rs1423096-G      rs1423096      0      1423096

```

```
## 2          rs9273373-G  rs9273373          0          9273373
## 3          rs4833095-T  rs4833095          0          4833095
## 4          rs1438673-C  rs1438673          0          1438673
## 5          rs10197862-A rs10197862          0          10197862
##      Context Intergenic Risk.Allele.Frequency  p.Value Pvalue_mlog
##      <character> <character>          <character> <numeric> <numeric>
## 1 Intergenic          1          0.78      1e-07      7.00000
## 2                   0          0.54      4e-14     13.39794
## 3 missense           0          0.74      5e-12     11.30103
## 4 Intergenic          1          0.49      3e-11     10.52288
## 5 intron              0          0.85      4e-11     10.39794
##      p.Value..text. OR.or.beta          X95..CI..text.
##      <character> <numeric>          <character>
## 1                   0.32 [0.25-0.40] ug/mL increase
## 2                   1.24          [1.17-1.30]
## 3                   1.20          [1.14-1.26]
## 4                   1.16          [1.11-1.21]
## 5                   1.24          [1.16-1.32]
##      Platform..SNPs.passing.QC.          CNV
##      <character> <character>
## 1                   Illumina [NR]          N
## 2 Illumina [up to 4,972,397] (imputed)          N
## 3 Illumina [up to 4,972,397] (imputed)          N
## 4 Illumina [up to 4,972,397] (imputed)          N
## 5 Illumina [up to 4,972,397] (imputed)          N
##      num.Risk.Allele.Frequency
##      <numeric>
## 1                   0.78
## 2                   0.54
## 3                   0.74
## 4                   0.49
## 5                   0.85
## -----
##      seqinfo: 23 sequences from hg19 genome
```

We would like to verify that the majority of the variants enumerated in the NHGRI catalog are also present in GRASP 2.0. We supply a function called `checkAnti` which obtains the anti-join between a chromosome-specific slice of the NHGRI catalogue and the slice of GRASP2 for the same chromosome. We compute for chromosome 22 the fraction of NHGRI records present in GRASP2.

```
gr22 = variant %>% filter(chr_hg19 == "22")
abs22 = checkAnti("22")
1 - (abs22 %>% nrow()) /
  (gr22 %>% count %>% collect %>% `[`("n"))
## [1] 0.9982036
```

The absent records can be seen to be relatively recent additions to the NHGRI catalog.

```
abs22
## # A tibble: 217 x 42
##   seqnames      start      end width strand Date.Added.to.Catalog PUBMEDID
##   <fctr>      <int>    <int> <int> <fctr>          <chr>    <int>
## 1          22 37545505 37545505     1      *          08/05/2014 24390342
## 2          22 39747671 39747671     1      *          08/05/2014 24390342
## 3          22 21979096 21979096     1      *          08/05/2014 24390342
```



```
## 4      22 48923459 48923459      1      *      07/28/2014 24322204
## 5      22 36125264 36125264      1      *      07/29/2014 24348519
## 6      22 21431054 21431054      1      *      07/23/2014 24324551
## 7      22 35144411 35144411      1      *      07/23/2014 24324551
## 8      22 22047969 22047969      1      *      07/23/2014 24324551
## 9      22 37586792 37586792      1      *      07/23/2014 24324551
## 10     22 34386473 34386473      1      *      04/26/2014 24165912
## # ... with 207 more rows, and 35 more variables: First.Author <chr>,
## #   Date <chr>, Journal <chr>, Link <chr>, Study <chr>,
## #   Disease.Trait <chr>, Initial.Sample.Size <chr>,
## #   Replication.Sample.Size <chr>, Region <chr>, Chr_id <chr>,
## #   Chr_pos.hg38 <dbl>, Reported.Gene.s. <chr>, Mapped_gene <chr>,
## #   Upstream_gene_id <chr>, Downstream_gene_id <chr>, Snp_gene_ids <chr>,
## #   Upstream_gene_distance <chr>, Downstream_gene_distance <chr>,
## #   Strongest.SNP.Risk.Allele <chr>, SNPs <chr>, Merged <chr>,
## #   Snp_id_current <chr>, Context <chr>, Intergenic <chr>,
## #   Risk.Allele.Frequency <chr>, p.Value <dbl>, Pvalue_mlog <dbl>,
## #   p.Value..text. <chr>, OR.or.beta <dbl>, X95..CI..text. <chr>,
## #   Platform..SNPs.passing.QC. <chr>, CNV <chr>,
## #   num.Risk.Allele.Frequency <dbl>, chr_hg19 <chr>, pos_hg19 <int>
```