

Documentation for Software GARFIELD

Valentina Iotchkova, Matthias Geihs, Graham Ritchie

October 25, 2017

1 Introduction

Software GARFIELD (GWAS analysis of regulatory and functional information enrichment with LD correction) implements the functional enrichment analysis approach described in [Iotchkova et al, 2016] using C++ for data pre-processing, and R for enrichment estimation and significance testing and visualisation. It provides a tool for assessing the enrichment of association analysis signals in 1005 features extracted from ENCODE, GENCODE and Roadmap Epigenomics projects, including genic annotations, chromatin states, histone modifications, DNaseI hypersensitive sites and transcription factor binding sites, among others, in a number of publicly available cell lines.

1.1 Method Overview

Genome-Wide Association Studies (GWAS) have been increasingly fruitful in discovering genotype-phenotype associations. The mechanisms underlying these associations, however, are still largely unknown as only a small fraction of these SNPs are known to directly alter protein-coding genes. The interpretation of functional consequences of non-coding variants has been greatly enhanced by large-scale efforts to identify regulatory genomic regions (e.g ENCODE). However, robust methods are still lacking to systematically evaluate the contribution of these regions to genetic variation implicated in diseases or quantitative traits.

We have developed a novel approach, named GARFIELD, that leverages GWAS findings with regulatory or functional annotations to find features relevant to a phenotype of interest. It performs greedy pruning of GWAS SNPs ($LD\ r^2 > 0.1$) and then annotates them based on functional information overlap. Next, it quantifies Odds Ratio (OR) at various GWAS significance cutoffs and assesses them by employing a generalized linear model framework, while matching for minor allele frequency, distance to nearest transcription start site and number of LD proxies ($r^2 > 0.8$). Within this framework, GARFIELD accounts for major sources of confounding that current methods do not offer.

1.2 How to cite

Valentina Iotchkova, Graham R.S. Ritchie, Matthias Geihs, Sandro Morganello, Josine L. Min, Klaudia Walter, Nicholas J. Timpson, UK10K Consortium, Ian Dunham, Ewan Birney and Nicole Soranzo. *GARFIELD - GWAS Analysis of Regulatory or Functional Information Enrichment with LD correction*. doi: <https://doi.org/10.1101/085738>

1.3 Contacts

For issues or further questions about GARFIELD you can get in touch at valentina.iotchkova@gmail.com.

2 Download and Installation

We provide source code for your own compilation, which can be downloaded as a compressed tarball from

```
http://www.ebi.ac.uk/birney-srv/GARFIELD/
```

or by opening a terminal and typing

```
wget http://www.ebi.ac.uk/birney-srv/GARFIELD/package/garfield-v2.tar.gz
```

We further provide all data necessary for running GARFIELD for analysis of genome-wide association studies in European populations. Those can be downloaded from

```
http://www.ebi.ac.uk/birney-srv/GARFIELD/
```

or by typing the following in a terminal

```
wget http://www.ebi.ac.uk/birney-srv/GARFIELD/package/garfield-data.tar.gz
```

Note: You should download both files in the same directory or you will need to modify the paths for the data files in the `garfield` script.

2.1 Installation

To decompress and extract the files in the software bundle in a terminal from the location where it is downloaded type

```
tar -xvf garfield-v2.tar.gz
```

In order to compile GARFIELD a C++ compiler is required. Additionally an R distribution needs to already be pre-installed for the end figure creation. Please make sure you have both of them before proceeding any further.

To compile the code type

```
cd garfield-v2  
make
```

This would create an executables `garfield-prep-chr`. Further details of how to run them will be given in the **Running GARFIELD** section.

In order to use the data files, they must also be decompressed. To do so execute

```
cd ../  
tar -xvf garfield-data.tar.gz
```

2.2 License

Copyright (C) 2017 Genome Research Ltd / EMBL - European Bioinformatics Institute

Author : Valentina Iotchkova <valentina.iotchkova@gmail.com>

Author : Matthias Geihs <mg18@sanger.ac.uk>

This file is part of GARFIELD - GWAS analysis of regulatory or functional information enrichment with LD correction.

GARFIELD is free software; you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation; either version 2 of the License, or (at your option) any later version.

This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.

You should have received a copy of the GNU General Public License along with this program; if not, write to the Free Software Foundation, Inc., 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301, USA.

2.3 Tree structures

```
garfield-v2
|----- garfield
|----- garfield_annotate_uk10k.sh
|----- garfield_extract_variants_overlapping_enriched_annotations.sh
|----- garfield-create-input-gwas.sh
|----- garfield-Meff-Padj.R
|----- garfield-test.R
|----- garfield-plot-function.R
|----- garfield-plot.R
|----- garfield-prep-chr.cpp
|----- LICENSE
|----- makefile
|----- README

garfield-data
|----- annotation
| |----- chr1
| |----- ...
| |----- chr22
| |----- link_file.txt
|----- maftssd
```

```

| |----- chr1
| |----- ...
| |----- chr22
|----- output
| |----- cd-meta
| | |----- garfield.Meff.cd-meta.out
| | |----- garfield.prep.cd-meta.out
| | |----- garfield.test.cd-meta.out
| | |----- garfield.test.cd-meta.out.Chromatin_States.pdf
| | |----- garfield.test.cd-meta.out.FAIRE.pdf
| | |----- garfield.test.cd-meta.out.Footprints.pdf
| | |----- garfield.test.cd-meta.out.Genic.pdf
| | |----- garfield.test.cd-meta.out.Histone_Modifications.pdf
| | |----- garfield.test.cd-meta.out.Hotspots.pdf
| | |----- garfield.test.cd-meta.out.Peaks.pdf
| | |----- garfield.test.cd-meta.out.TFBS.pdf
| |----- GIANT_HEIGHT
| | |----- garfield.Meff.cd-meta.out
| | |----- garfield.prep.GIANT_HEIGHT.out
| | |----- garfield.test.GIANT_HEIGHT.out
| | |----- garfield.test.GIANT_HEIGHT.out.Chromatin_States.pdf
| | |----- garfield.test.GIANT_HEIGHT.out.FAIRE.pdf
| | |----- garfield.test.GIANT_HEIGHT.out.Footprints.pdf
| | |----- garfield.test.GIANT_HEIGHT.out.Genic.pdf
| | |----- garfield.test.GIANT_HEIGHT.out.Histone_Modifications.pdf
| | |----- garfield.test.GIANT_HEIGHT.out.Hotspots.pdf
| | |----- garfield.test.GIANT_HEIGHT.out.Peaks.pdf
| | |----- garfield.test.GIANT_HEIGHT.out.TFBS.pdf
|----- pval
| |----- cd-meta
| | |----- chr1
| | |----- ...
| | |----- chr22
| |----- GIANT_HEIGHT
| | |----- chr1
| | |----- ...
| | |----- chr22
|----- tags
| |----- r01
| | |----- chr1
| | |----- ...
| | |----- chr22
| |----- r08
| | |----- chr1
| | |----- ...
| | |----- chr22

```

3 Input files

GARFIELD requires a number of input files containing the GWAS, LD and annotation data as well as minor allele frequencies (MAF) and distances to nearest transcription start site (TSS) needed for the analysis. We have pre-processed data on all variants from the UK10K study, which can be used for the enrichment analysis of European sample GWAS. This data is contained in the `garfield-data.tar.gz` file. The size of the compressed tarball is 5.9Gb which amounts to 83Gb after unpacking spread into five sub-directories `annotation`, `pval`, `maftssd`, `tags` and `output`.

If you have downloaded the data somewhere else than the same folder as the software package you would need to change the `$DATADIR` paths set in the `./garfield-v2/garfield` file.

If you are interested in running enrichment analysis for non-European samples then allele frequencies and LD need to be re-calculated for them and put into the types of files described in the next subsections.

3.1 GWAS data

Summary statistics from a GWA study of interest need to be provided in a separate subfolder of the `./garfield-data/pval` folder, e.g. `./garfield-data/pval/GIANT_HEIGHT`, in files split by chromosome and named `chr1`, `chr2`, etc.. The files must contain no header and have genomic position (build 37) in the first column and P-value from association analysis in the second column using a space as a column delimiter. In addition, files must be numerically sorted. An example by the supplied in the package `./garfield-data/pval/GIANT_HEIGHT/chr22` file is given below

```
16381711 0.9274
16966809 0.4753
16970900 0.5239
16983606 0.4753
16989290 0.4689
17030792 0.2353
```

The GIANT Height association analysis summary statistics have been downloaded from

http://www.broadinstitute.org/collaboration/giant/index.php/GIANT_consortium_data_files

and processed by us and added here as an example together with a GWAS on Crohn's disease downloaded from

<http://www.ibdgenetics.org/downloads.html>

Tree structure

```
garfield-data
|----- pval
| |----- cd-meta
| | |----- chr1
| | |----- ...
```


Tree structure

```
garfield-data
|----- tags
|   |----- r01
|   |   |----- chr1
|   |   |----- ...
|   |   |----- chr22
|   |----- r08
|   |   |----- chr1
|   |   |----- ...
|   |   |----- chr22
```

3.4 MAF, TSS distance

Finally, for testing we need to match variants by MAF and distance to nearest TSS. This data is provided in the `./garfield-data/maftssd` folder, again split into space separated numerically sorted `chr` files. An example is given below, where the first column contains the genomic position of each variant, the second it's MAF calculated from the UK10K sequence data and the third contains it's distance to the nearest TSS.

```
16051237 0.000145677788349 -10919
16051249 0.106493504342 -10907
16051477 0.00320565708155 -10679
16051497 0.36456744798 -10659
16052080 0.104518066225 -10076
16052107 0.000285031428574 -10049
16052216 0.000143040758483 -9940
```

Tree structure

```
garfield-data
|----- maftssd
|   |----- chr1
|   |----- ...
|   |----- chr22
```

4 Running GARFIELD

The file `./garfield-v2/garfield` contains all necessary commands and shows how to run the software for the GIANT Height data example. To execute simply open a terminal, go to the `garfield-v2` software directory and type

```
./garfield
```

If you follow the same directory structure as described above, to run analysis for your own GWAS data you only need to put it in a subfolder of `./garfield-data/pval` and then change the

`$INPUTNAME` variable in the `garfield` file to the name of this subfolder.

In more detail, the enrichment analysis is performed with the use of four separate tools

1. `garfield-prep-chr` for preparing the data to be used for analysis
2. `garfield-Meff-Padj.R` for estimating the effective number of annotations and multiple testing adjusted p-value
3. `garfield-test.R` for calculating enrichment ORs and associated p-values
4. `garfield-plot.R` for visualisation of the computed results.

Additionally `garfield_extract_variants_overlapping_enriched_annotations.sh` can be used for identifying the exact variants driving the enrichment signal for a given threshold and given annotations. There are further helper scripts for annotating variants with custom annotations from bed files `garfield_annotate_uk10k.sh` and `garfield_annotate_uk10k_helper`, and for creating GARFIELD input p-value files from general GWAS summary statistics `garfield-create-input-gwas.sh`.

4.1 `garfield-prep-chr`

prepares the input data for `garfield-test.R` and `garfield-Meff-Padj.R`.

Usage:

```
./garfield-prep-chr -ptags prunetags -ctags clumptags -maftss maftssd -pval pvalue \
  -ann annot -o output [-excl exclude] [-chr chr]
```

- `prunetags` is a file with LD tags format like `./garfield-data/tags/r01/chr22` which is to be used for LD pruning
- `clumptags` is a file with LD tags format like `./garfield-data/tags/r08/chr22` which is to be used for further annotation and feature matching for the permutation testing
- `maftssd` is a file like `./garfield-data/maftssd/chr22` which contains the MAF and TSS distances for feature matching in the permutation testing
- `pvalue` is a file with format as in `./garfield-data/pval/GIANT_HEIGHT/chr22` which contains the p-values from the association analysis
- `annot` is a file with format at in `./garfield-data/annotation/chr22` which contains the annotations of all genetic variants
- `exclude` is a comma separated string of integer values of annotation indices not to be additionally inheriting the annotations of their high LD proxies. This may be useful for very broad annotations such as intergenic regions and repressed segmentation states. An example would be `1,2,3`, which corresponds to the first three annotations in the `./garfield-data/annotation/link_file.txt`. Conversely if no such annotations are to be present `-1` should be used.

- **chr** is an identified for the chromosome, which can be used for identification of variants with given annotations. If this is of no interest, it can be ignored or set to anything for this part of the script.
- **output** is the output file form the preparation step

The output of the tool is space separated and contains the remaining independent variants after greedy pruning with the following fields

```
chr snpid pvalue clump_tag_count maf tssd a1a2...
```

where **chr** is the chromosome of the variant, **snpid** is the genomic position of a variant, **pvalue** is its association analysis p-value, **clump_tag_count** is the number of high LD proxies that variant has, **maf** is its MAF, **tssd** is its distance to the nearest TSS, **a1a2...** are all annotations of that variant after additionally adding the annotations of all high LD proxies with format as in the `./garfield-data/annotation/chr22` file. A sample output is given below

[illegible]

4.2 garfield-Meff-Padj.R

computes the effective number of annotations and adjusted for multiple testing p-values from the output of the `garfield-prep-chr` tool.

Usage:

```
Rscript garfield-Meff-Padj.R -i prepfiler -o outfile [-s subset]
```

- **-i prepfile** specifies the input file, where **prepfile** is a file created by the **garfield-prep-chr** tool
- **-o outfile** specifies the output file
- **-s** specifies a subset of annotations for which to run the enrichment analysis (if not interested in all of them). The format is a comma separated list of ranges, e.g. **2-5,10-80,99** . This is an optional parameter and if left unspecified, the tool will use all annotations.

The output of the tool is space separated and contains the results from the permutation testing along with some summaries and has the following format

```
Meff 496.41
Padj 0.0001
```

where

- **Meff** shows the effective number of annotations
- **Padj** shows the enrichment p-value adjusted for multiple testing (on the effective number of annotations)

4.3 garfield-test.R

computes odds ratios and enrichment p-values, while using accounting for variant differences in number of high LD proxies (`clump_tag_count`), minor allele frequency (MAF) and transcription start site (TSS) distance. Variants for each of these three features are split into custom number of bins and are represented by categorical covariates in a logistic regression model.

Usage:

```
Rscript garfield-test.R -i prefile -o outfile -l linkfile -pt pthreshs -b binning \
-c condition [-s subset] [-ct condthresh -padj padj]
```

- `-i prefile` specifies the input file, where `prefile` is a file created by the `garfield-prep-chr` tool
- `-o outfile` specifies the output file
- `-l linkfile` specifies the annotation `./garfield-data/annotation/link_file.txt`.
- `-pt` specifies the GWAS p-value thresholds to test enrichment for, e.g. `-pt 1e-2,1e-3`
- `-b` specifies number of binning quantiles for each of the binning dimensions: number of high LD tags, minor allele frequency and nearest transcription start site distance, e.g. `-b n5,m5,t5` will create 5 LD tag bins, 5 MAF bins and 5 TSS bins, respectively
- `-c` specifies if GARFIELD is to run one annotation at a time (0) or if additionally heuristic model selection is to be performed (1). Note: to use `-c 1` GARFIELD must first be run with the `-c 0` option
- `-ct` specifies the p-value threshold for considering two annotations conditionally independent. This option is only needed when using the `-c 1` option
- `-padj` specifies the adjusted for multiple testing p-value threshold for considering an annotation to be significantly enriched/depleted for a given trait. This option is only needed when using the `-c 1` option
- `-s` specifies a subset of annotations for which to run the enrichment analysis (if not interested in all of them). The format is a comma separated list of ranges, e.g. `2-5,10-80,99`. This is an optional parameter and if left unspecified, the tool will test all annotations.

The output of the tool is space separated and contains the results from the enrichment testing along with some summaries, where the columns denote

- `ID` shows the ID of the annotation from the annotation file (`INDEX+1` in the link file)
- `PThresh` is a GWAS threshold used for enrichment analysis testing
- `OR` is the odds ratio for that annotation at that threshold
- `Pvalue` is the p-value of the significance of the observed enrichment
- `Beta` denotes the effect size (or log odds ratio) for that annotation at that threshold

- `SE` denotes the standard error of the effect size for that annotation at that threshold
- `CI95_lower` denotes the lower bound for the 95% CI of the effect size
- `CI95_upper` denotes the upper bound for the 95% CI of the effect size
- `NAnnotThresh` is the number of (independent) annotated variants with the considered annotation passing the GWAS significance threshold `PThresh` (after pruning)
- `NAnnot` is the total number of (independent) annotated variants with the given annotation (after pruning)
- `Nthresh` is the number of (independent) variants passing the GWAS significance threshold `PThresh` (after pruning)
- `N` is the total number of LD pruned variants
- `linkID` is the ID of the annotation from the link file
- `Annotation` specifies a unique name for the annotation
- `Celltype` specifies the cell type of the annotation (if applicable)
- `Tissue` specifies the tissue of the annotation (if applicable)
- `Type` specifies the subtype of the annotation
- `Category` specifies whether the annotation is a chromatin state, transcription factor binding site, genic annotation, etc...

Additionally, when using the conditional analysis step, there is an extra column `ID_c` at the beginning of the file denoting the index of the annotation from the order in which conditional analysis was performed. Conditional analysis creates three output files - one for conditionally independent annotations, one for the remaining annotations that were found to be enriched on their own but no longer show significant enrichment in the presence of other annotations and one for a summary of the best fitted model after model selection.

4.4 garfield-plot.R

produces the final figures from the table of results. It has the following usage

```
Rscript garfield-plot.R -i test.out -o output_path_prefix -t plot_title -f min \
-padj thresh -s subset [-col set_of_colour_names]
```

- `test.out` is an output file from the `garfield-test` tool
- `output_path_prefix` is a path prefix for the output files and figures
- `plot_title` is the title label for the figures. If you do not want a title " " should be used here
- `min` is the minimum number of variants at a certain threshold to be used for filtering the data before plotting

- **thresh** is the significance threshold to be used for plotting. Value of zero sets it to the default value of 0.05/498.
- **subset** specifies a subset of annotations for which to plot the enrichment analysis results (if not interested in all of them). The format is a comma separated list of ranges, e.g. 2-5,10-80,99 . This is an optional parameter and if left unspecified, the tool will plot all annotations.
- **col** is an optional parameter for changing the colours corresponding to the different GWAS thresholds used. Note this needs to be a comma separated list of length 1+Number of thresholds used and it needs to contain names recognised by R as colour names, e.g. for a set of 3 colours "red,blue,white".

The output of the function is a set of figures for different classes of functional annotations. Sample output files can be found in the `./garfield-data/output/GIANT_HEIGHT` folder, which are a result from the enrichment analysis of the GIANT Height GWAS data.

4.5 `garfield_annotate_uk10k.sh`

Creates annotations for GARFIELD from user supplied annotations in bed file format. It requires a separate bed file for each annotation and requires all annotations to reside within a single folder with no additional files being present in the same folder. The function uses internally `garfield_annotate_uk10k_helper` for creating the overlaps.

- **ANNOTATION_DIR** points to the folder containing the bed files
- **BED_FILES** takes all (bed) files in **ANNOTATION_DIR**. Note: if there are other files in that folder, you will need to change this to only point to the files of interest
- **OUTPUT_DIR** points to the directory where the output is stored
- **VARIANTS_DIR** is the folder containing all variants to annotate. This is already pre-specified, so no need to change unless you want to run this step on a different reference SNP set. Note: if that is the case then all other GARFIELD input files need to be calculated for the same SNP set.

4.6 `garfield_extract_variants_overlapping_enriched_annotations.sh`

Extracts the genetic variants driving the enrichment signal for given annotations and given GWAS significance threshold. It has the following usage

```
./garfield_extract_variants_overlapping_enriched_annotations.sh $PRUNETAGSDIR \
$CLUMPTAGSDIR $ANNOTDIR $PVALDIR $PREPFILE $TESTFILE $GARFIELD_significant_annotations \
$GARFIELD_VARS $PTHRESH $PENRICH
```

- **PRUNETAGSDIR** is the folder containing the LD tags for pruning
- **CLUMPTAGSDIR** is the folder containing the LD tags for annotation
- **ANNOTDIR** is the folder containing all annotation files

- PVALDIR is the folder containing the P-values used as input for GARFIELD
- PREPFILE is the GARFIELD output file from the `garfield-prep` step
- TESTFILE is the GARFIELD output file from the `garfield-test` step
- GARFIELD_significant_annotations is an output file containing the significant annotations found after enrichment analysis at GWAS threshold PTHRESH and enrichment significance level PENRICH. Columns denote annotation index, odds ratio of enrichment, p-value of enrichment, annotation name
- GARFIELD_VARS specifies the output, which contains for each significantly enriched annotation, each pruned variant used for the enrichment analysis along with summary information from the enrichment (OR, P-value) and whether the variants itself is annotated and what GWAS p-value it has, what LD tags it has with annotation overlap and GWAS p-value information.
 - ID annotation index
 - ANNOTATION annotation name
 - OR odds ratio for given annotation at given threshold
 - PVAL p-value of enrichment at given threshold
 - VAR_INFO chr:pos(GWAS Pvalue, Annotationoverlap)
 - INFO_TAGS_USED_TO_ANNOTATElist of chr:pos(GWAS Pvalue, Annotation overlap) for all LD tags ($r^2 > 0.8$). Separated by |.
 - INFO_TAGS_PRUNED_OUTlist of chr:pos(GWAS Pvalue, Annotation overlap) for all LD tags ($r^2 > 0.1$). Separated by |.
- PTHRESH is the GWAS threshold of interest
- PENRICH is the enrichment significance threshold for which to extract relevant annotations (recommended to set it equal to Padj).

4.7 garfield-create-input-gwas.sh

Creates GARFIELD input p-value files from general GWAS summary statistics files. After setting the required parameters internally it has the following usage

```
./garfield-create-input-gwas.sh
```

- chrcol is the column in GWAS file containing chromosome information
- poscol is the column in GWAS file containing genomic position information (build37)
- pvalcol is column in GWAS file containing GWAS p-value information
- TRAITNAME is the name of directory for GWAS trait to be created
- GWASFILENAME is the name of file containing GWAS summary statistics
- OUTDIR is the output directory to be used as input for GARFIELD analysis, by default this is `../garfield-data/pval/$TRAITNAME`.

5 Output

5.1 Tree structure

The final basic usage output is created in subfolders of the `./garfield-data/output/` directory and contains a total of 11 files: an output from the `garfield-prep` step, an output from the `garfield-Meff-Padj.R` step, a final results file from the `garfield-test.R` step and 8 figures, one for each of the different types of annotations used in our analysis, namely genic annotations, chromatin segmentation states, transcription factor binding sites, histone modifications and open chromatin data (FAIRE, DHS Hotspots, peaks and footprints).

garfield-data

```
|----- output
|----- cd-meta
| |----- garfield.Meff.cd-meta.out
| |----- garfield.prep.cd-meta.out
| |----- garfield.test.cd-meta.out
| |----- garfield.test.cd-meta.out.Chromatin_States.pdf
| |----- garfield.test.cd-meta.out.FAIRE.pdf
| |----- garfield.test.cd-meta.out.Footprints.pdf
| |----- garfield.test.cd-meta.out.Genic.pdf
| |----- garfield.test.cd-meta.out.Histone_Modifications.pdf
| |----- garfield.test.cd-meta.out.Hotspots.pdf
| |----- garfield.test.cd-meta.out.Peaks.pdf
| |----- garfield.test.cd-meta.out.TFBS.pdf
|----- GIANT_HEIGHT
| |----- garfield.Meff.cd-meta.out
| |----- garfield.prep.GIANT_HEIGHT.out
| |----- garfield.test.GIANT_HEIGHT.out
| |----- garfield.test.GIANT_HEIGHT.out.Chromatin_States.pdf
| |----- garfield.test.GIANT_HEIGHT.out.FAIRE.pdf
| |----- garfield.test.GIANT_HEIGHT.out.Footprints.pdf
| |----- garfield.test.GIANT_HEIGHT.out.Genic.pdf
| |----- garfield.test.GIANT_HEIGHT.out.Histone_Modifications.pdf
| |----- garfield.test.GIANT_HEIGHT.out.Hotspots.pdf
| |----- garfield.test.GIANT_HEIGHT.out.Peaks.pdf
| |----- garfield.test.GIANT_HEIGHT.out.TFBS.pdf
```

5.2 Results output file

An example of the `./garfield-data/output/cd-meta/garfield.test.cd-meta.out` file is shown below

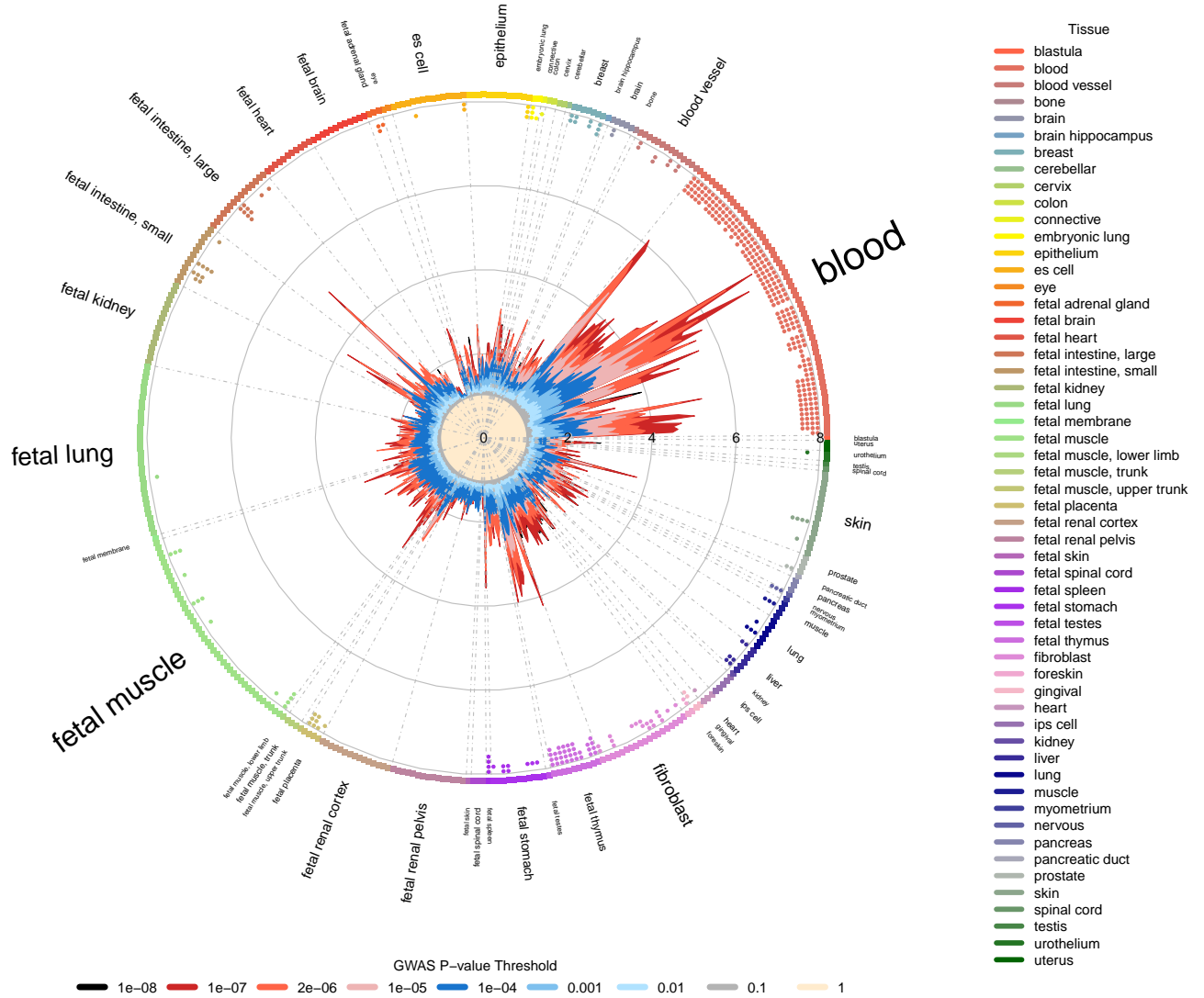
ID	PThresh	OR	Pvalue	Beta	SE	CI95_lower	CI95_upper	NAnnotThesh	NAnnot	NThresh	N	linkID	Annotation	Celltype	Tissue	Type	Category
1	1e-05	1.75	9.67e-3	0.56	0.21	0.13	0.99	30	5095	182	82981	0	AG10803_footprints.txt	AG10803	skin	footprints	Footprints
1	1e-08	1.66	0.12	0.51	0.33	-0.14	1.17	12	5095	72	82981	0	AG10803_footprints.txt	AG10803	skin	footprints	Footprints
2	1e-05	1.55	0.03	0.44	0.20	0.03	0.85	36	6670	182	82981	1	AoAF_footprints.txt	AoAF	blood_vessel	footprints	Footprints
2	1e-08	1.51	0.19	0.41	0.32	-0.21	1.04	15	6670	72	82981	1	AoAF_footprints.txt	AoAF	blood_vessel	footprints	Footprints
3	1e-05	2.26	7.02e-4	0.81	0.24	0.34	1.29	24	2599	182	82981	2	CD20+_footprints.txt	CD20+	blood	footprints	Footprints
3	1e-08	3.17	5.94e-4	1.15	0.33	0.49	1.81	14	2599	72	82981	2	CD20+_footprints.txt	CD20+	blood	footprints	Footprints
4	1e-05	2.75	7.43e-07	1.01	0.20	0.61	1.41	39	4304	182	82981	3	CD34+_Mobilized_footprints.txt	CD34+	blood	footprints	Footprints
4	1e-08	3.58	1.91e-05	1.27	0.29	0.69	1.86	21	4304	72	82981	3	CD34+_Mobilized_footprints.txt	CD34+	blood	footprints	Footprints
5	1e-05	0.94	0.83	-0.05	0.27	-0.59	0.47	17	5062	182	82981	4	fBrain_footprints.txt	fBrain	fetal_brain	footprints	Footprints

where

- **ID** shows the ID of the annotation (INDEX+1 from link file `./garfield-data/annotation/link_file.txt`)
- **PThresh** is a GWAS threshold used for enrichment analysis testing
- **OR** is the observed odds ratio for that annotation at that threshold
- **Pvalue** is the p-value of the significance of the observed enrichment.
- **Beta** is the effect size from the generalized linear model also equal to $\log(\text{OR})$
- **SE** is the standard error from the generalized linear model
- **CI95_lower** and **CI95_upper** show the 95% confidence interval for **Beta**.
- **NAnnotThresh** is the number of (independent) annotated variants with the considered **Annotation** passing the GWAS significance threshold **PThresh** (after pruning)
- **NAnnot** is the total number of (independent) annotated variants with the given annotation (after pruning)
- **Nthresh** is the number of (independent) variants passing the GWAS significance threshold **PThresh** (after pruning)
- **N** is the total number of LD pruned variants.
- **linkID** shows the ID of the annotation from the link file `./garfield-data/annotation/link_file.txt`.
- **Annotation** specifies a unique name for the annotation
- **Celltype** specifies the cell type of the annotation (if applicable)
- **Tissue** specifies the tissue of the annotation (if applicable)
- **Type** specifies the subtype of the annotation
- **Category** specifies whether the annotation is a Chromatin state, transcription factor binding site, genic annotation, etc...

5.3 Figures

./garfield-data/output/cd-meta/garfield.test.cd-meta.out.Hotspots.pdf



Enrichment of Crohn's Disease variants in DNaseI Hypersensitive sites (broad peaks) from ENCODE and Roadmap Epigenomics data. Radial plot shows the enrichment (OR) in each cell type (dots on the outside of the circle sorted by tissue) for each GWAS significance threshold between 10^{-8} and all (independent) variants (shown by inner colours and bottom legend). Furthermore, small dots on the outer side of the plot show if the observed enrichment is significant (if there is a dot present) or not (if there isn't) for thresholds 10^{-5} , 10^{-6} , 10^{-7} , 10^{-8} in direction from outside to inside.

References

Valentina Iotchkova, Graham R.S. Ritchie, Matthias Geihs, Sandro Morganello, Josine L. Min, Klaudia Walter, Nicholas J. Timpson, UK10K Consortium, Ian Dunham, Ewan Birney and Nicole Soranzo. *GARFIELD - GWAS Analysis of Regulatory or Functional Information Enrichment with LD correction*. doi: <https://doi.org/10.1101/085738>