# quincunx: an R package to query, download and wrangle PGS Catalog data

Ramiro Magno[1,2,¶] Isabel Duarte[1,2,¶] Ana -Teresa Maia[1,2,3,*]

[1] Centre for Biomedical Research (CBMR), Universidade do Algarve, Gambelas
Campus, 8005-139 Faro, Portugal
[2] Algarve Biomedical Center (ABC), Universidade do Algarve, Gambelas Campus,
8005-139 Faro, Portugal
[3] Faculty of Medicine and Biomedical Sciences (FMCB), Universidade do Algarve,
Gambelas Campus, 8005-139 Faro, Portugal

¶ These authors contributed equally to this work.
* To whom correspondence should be addressed. Email: atmaia@ualg.pt

## Abstract

**Motivation:** The Polygenic Score (PGS) Catalog is a recently established open
database of published polygenic scores that, to date, has collected, curated, and made
available 721 polygenic scores from over 133 publications. The PGS Catalog REST API
is the only method allowing programmatic access to this resource.
**Results:** Here, we describe *quincunx*, an R package that provides the first client
interface to the PGS Catalog REST API. *quincunx* enables users to query and quickly
retrieve, filter and integrate metadata associated with polygenic scores, as well as
polygenic scoring files in tidy table format.
**Availability:** *quincunx* is freely available under an MIT License, and can be accessed
from https://github.com/maialab/quincunx.

## Keywords

Polygenic Scores, PGS, disease susceptibility, Genomics, PGS Catalog, R, SNP, trait,
Software, human, REST client

## Introduction                                                                                      1

For two decades, GWAS identified individual variants associated with risk for complex     2
diseases, raising the hopes of a polygenic approach for disease prevention. However,        3
until recently, integration of these results was challenging delaying its prompt            4
application to the clinical setting. In 2020 alone, over 1,400 publications on polygenic    5
risk scores (PGS) appeared in PubMed, raising the need for a standardised distribution      6
of studies' key data, assuring their wide evaluation and accurate use.                      7

The Polygenic Score (PGS) Catalog, created in 2019, is a publicly available,              8
manually curated, open database of PGS and relevant metadata, that responds to this        9
need [1]. Its current release [date 2021-02-03] includes manually curated data from 133   10
publications and 721 PGS associated with 194 traits. Currently, there are three           11
alternative ways to access the data: (i) the web graphical user interface (GUI); (ii) by  12

downloading database dumps; and (iii) the recently implemented PGS Catalog    13
representational state transfer (REST) application programming interface (API),    14
released in [date 2020-06-03], which provides direct programmatic access to the    15
database, being this the preferred method for batch analyses.    16

 We developed *quincunx*, the first R package [2] to programmatically access the PGS    17
Catalog REST API. This package provides a simple user-friendly interface for querying    18
the most updated Catalog data, retrieve and map it to in-memory relational databases    19
of tidy data tables, allowing its prompt integration with tidyverse packages for    20
subsequent data transformation, visualisation and modelling [3, 4].    21

# Results    22

## Retrieving data from the PGS Catalog REST API    23

The PGS Catalog REST API is an EBI service hosted at    24
`https://www.pgscatalog.org/rest/`. The REST API uses hypermedia with resource    25
responses following the OpenAPI Specification    26
(`https://swagger.io/docs/specification/about/`). Response data is provided as    27
hierarchical data in JSON format and can be paginated, i.e., split into multiple    28
responses (`https://www.pgscatalog.org/rest/`).    29

 To ease the conversion from the hierarchical to the relational tabular format — the    30
preferred format for data analysis in R [4] — we developed a set of retrieval functions    31
(Fig. 1A). Since the REST API data is organised around five core data entities —    32
*Polygenic Scores*, *PGS Publications*, *PGS Sample Sets*, *PGS Performance Metrics* and    33
*EFO traits*— we implemented five corresponding retrieval functions that encapsulate    34
the technical aspects of resource querying and format conversion: `get_scores()`,    35
`get_publications()`, `get_sample_sets()`, `get_performance_metrics` and `get_traits()`    36
(Fig. 1A). These functions simplify the querying of PGS Catalog entities, by providing a    37
complete and consistent interface to the Catalog. For example, to query for *scores*, the    38
user needs only to know the function `get_scores()`, whereas the REST API itself    39
exposes three separate resource URL endpoints for *scores* with different querying    40
parameters. Moreover, the user can choose directly the arguments of the retrieval    41
functions from any number of available search criteria exposed by the REST API    42
(Fig. 1B). All arguments are vectorised, meaning that multiple queries are promptly    43
available from a single function call. Results obtained from multiple queries can be    44
combined with the logical operators `OR` or `AND` using the `set_operation` parameter. If    45
`set_operation` is set to `OR` (default behaviour), results are collated while removing    46
duplicates, if any. If `set_operation` is set to `AND`, only entities that concomitantly    47
match all search criteria are returned. If finer control is needed on combining results,    48
the following functions can be used: `bind()`, `union()`, `intersect()`, `setdiff()`, and    49
`setequal()`. These are S4 methods that work with the S4 classes created in *quincunx*.    50
An example of a case study (in tutorial style) can be found in Additional file 1.    51

## Representation of PGS Catalog entities    52

All S4 classes share the same design principles that make them relational databases: (i)    53
each slot corresponds to a table (dataframe in R); (ii) the first slot corresponds to the    54
main table that lists observations of the respective PGS Catalog entity, e.g., *scores*; and    55
(iii) all tables have a primary key, the identifier of the respective PGS Catalog entity:    56
`pgs_id`, `pgp_id`, `pss_id`, `ppm_id` or `efo_id`. For easy consultation of the variables present    57
in the retrieved data tables, we provide a cheatsheet (Additional file 2); and for detailed    58
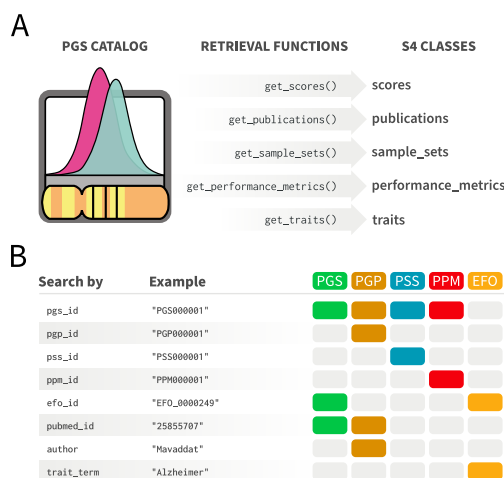descriptions, the user can issue the following help commands to open the respective help    59

**Figure 1.** *quincunx* retrieval functions. **(A)** Functions for retrieving data from the PGS Catalog: `get_scores()`, `get_publications()`, `get_sample_sets()`, `get_performance_metrics()` and `get_traits()`. **(B)** *quincunx* search criteria (function parameters) to be used with retrieval functions. Coloured boxes indicate which entities can be retrieved by each search criteria.

pages for each class: `class?scores`, `class?publications`, `class?sample_sets`, `class?performance_metrics` or `class?traits`.

## Improvements & Limitations

Compared to the exposed REST API, we have improved data accessibility in *quincunx* in several ways. Firstly, we harmonised the nomenclature of the variables in tidy tables with the nomenclature used by the GWAS Catalog [5], namely for the variables that are also used by the R package *gwasrapidd* [6]: an analogous R package that provides access to the GWAS Catalog REST API. This permits a frictionless wrangling of variables between the two R packages, allowing crosstalk between the data from the two Catalogs. Secondly, by recognising that in some cases the values of a variable are provided in its name and not in its value (a case of untidy data), we decided to perform the required refactoring to make those variables explicit columns in the relational tables, thus making the data more analysis friendly. For example, the *stage* of a sample comes implicitly coded in the JSON keys `samples_variants` and `samples_training` and are mapped in *quincunx* to the variable `stage`, with values ``discovery`` and ``training``, respectively. Additionally, the PGS Catalog REST API does not offer specific endpoints allowing direct mapping between the PGS entities, as this information is deeply nested in the hierarchical structure of the JSON responses. *quincunx* facilitates the retrieval of relationships between entities, by providing a set of mapping functions based on the entities' identifiers, e.g., `pgs_to_pgp`, `pgp_to_ppm()`, `ppm_to_pss()`, including mapping (when applicable) from PGS scores to GWAS studies: `pgs_to_study()` and `study_to_pgs()` (see online documentation for the complete list). Finally, *quincunx* provides a set of helper functions to easily browse linked web resources, such as PubMed (`open_in_pubmed()`), dbSNP (`open_in_dbsnp()`), and the PGS Catalog Web interface itself (`open_in_pgs_catalog()`).

Despite the availability of some R software packages, e.g. *bigsnpr* [7], *RápidoPGS* [8], or *SummaryLasso* [9], that allow the application of polygenic scores to particular datasets, i.e. for analyses downstream of *quincunx* data retrieval, the most popular software tools for these calculations, e.g. PRSice [10,11], *LDpred* [12], *PRS-CS* [13],

*JAMPred* [14], *lassosum* [15], *PLINK* [16,17], do not run in R. This could present an obstacle to pursuing a full PGS analysis within the same R framework (using the PGS scoring files), therefore delaying the process of polygenic score application (for an overview of the currently available methods, please see [18]).

## Conclusion

We have developed the first R client for the PGS Catalog REST API, thus greatly facilitating the programmatic access to the database. The main advantages of *quincunx* are: (i) providing a simple user-friendly interface to the REST API, allowing the programmatic querying of the most updated data from within R; (ii) the retrieval of the data in an analysis friendly format, with tidy data representations of the PGS entities, i.e., of *scores*, *publications*, *sample sets*, *performance metrics* and *traits* in the form of in-memory relational databases; (iii) allowing the automatic retrieval of polygenic scoring files from the PGS Catalog FTP server, making the data immediately available for analysis in R (as an extra feature not available via the REST API); and (iv) dedicated functions to export the retrieved objects to Excel (.xlsx) format for data inspection and sharing outside of R. *quincunx* is a package that will greatly improve the research community's ability for data mining within R, therefore accelerating the evaluation and subsequent application of published and manually curated polygenic scores.

## Availability and requirements

- **Project name**: quincunx.
- **Project home page**: https://github.com/maialab/quincunx.
- **Operating system(s)**: Platform independent.
- **Programming language**: R.
- **Other requirements**: None.
- **License**: MIT.
- **Any restrictions to use by non-academics**: None.

## Competing interests

The authors declare that they have no competing interests.

## Funding

## Author's contributions

RM & ID devised and wrote the package, and wrote the manuscript. ATM supervised the project and wrote the manuscript. All authors read and approved the final manuscript.

# Acknowledgements 126

# References 131

1. Lambert SA, Gil L, Jupp S, Ritchie SC, Xu Y, Buniello A, et al. The Polygenic 132
   Score Catalog: an open database for reproducibility and systematic evaluation. 133
   medrxiv. 2020 may. Available from: 134
   `https://www.medrxiv.org/content/10.1101/2020.05.20.20108217v1`. 135

2. R Core Team. R: A Language and Environment for Statistical Computing. 136
   Vienna, Austria; 2017. Available from: `https://www.R-project.org/`. 137

3. Wickham H, et al. Tidy data. Journal of Statistical Software. 2014;59(10):1–23. 138

4. Wickham H, Grolemund G. R for Data Science: Import, Tidy, Transform, 139
   Visualize, and Model Data. 1st ed. O'Reilly Media; 2017. Available from: 140
   `http://r4ds.had.co.nz/`. 141

5. Welter D, MacArthur J, Morales J, Burdett T, Hall P, Junkins H, et al. The 142
   NHGRI GWAS Catalog, a curated resource of SNP-trait associations. Nucleic 143
   Acids Research. 2013;42(D1):D1001–D1006. 144

6. Magno R, Maia AT. gwasrapidd: an R package to query, download and wrangle 145
   GWAS catalog data. Bioinformatics. 2019 aug. 146

7. Privé F, Aschard H, Ziyatdinov A, Blum MGB. Efficient analysis of large-scale 147
   genome-wide data with two R packages: bigstatsr and bigsnpr. Bioinformatics. 148
   2018 mar;34(16):2781–2787. 149

8. Reales G, Vigorito E, Kelemen M, Wallace C. RápidoPGS: A rapid polygenic 150
   score calculator for summary GWAS data without validation dataset. biorxiv. 151
   2020 jul. 152

9. Chen TH. SummaryLasso: Building Polygenic Risk Score Using GWAS Summary 153
   Statistics; 2019. R package version 1.2.1. Available from: 154
   `https://CRAN.R-project.org/package=SummaryLasso`. 155

10. Euesden J, Lewis CM, O'Reilly PF. PRSice: Polygenic Risk Score software. 156
    Bioinformatics. 2014 dec;31(9):1466–1468. 157

11. Choi SW, O'Reilly PF. PRSice-2: Polygenic Risk Score software for 158
    biobank-scale data. GigaScience. 2019 jul;8(7). 159

12. Vilhjálmsson BJ, Yang J, Finucane HK, Gusev A, Lindström S, Ripke S, et al. 160
    Modeling Linkage Disequilibrium Increases Accuracy of Polygenic Risk Scores. 161
    The American Journal of Human Genetics. 2015 oct;97(4):576–592. 162

13. Ge T, Chen CY, Ni Y, Feng YCA, Smoller JW. Polygenic prediction via 163
    Bayesian regression and continuous shrinkage priors. Nature Communications. 164
    2019 apr;10(1). 165

14. Newcombe PJ, Nelson CP, Samani NJ, Dudbridge F. A flexible and parallelizable approach to genome-wide polygenic risk scores. Genetic Epidemiology. 2019 jul;43(7):730–741. 166 167 168

15. Mak TSH, Porsch RM, Choi SW, Zhou X, Sham PC. Polygenic scores via penalized regression on summary statistics. Genetic Epidemiology. 2017 may;41(6):469–480. 169 170 171

16. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK: A Tool Set for Whole-Genome Association and Population-Based Linkage Analyses. The American Journal of Human Genetics. 2007 sep;81(3):559–575. 172 173 174

17. Chang CC, Chow CC, Tellier LCAM, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK: rising to the challenge of larger and richer datasets. GigaScience. 2015 feb;4(1). 175 176 177

18. Choi SW, Mak TSH, O'Reilly PF. Tutorial: a guide to performing polygenic risk score analyses. Nature Protocols. 2020 jul;15(9):2759–2772. 178 179

# Additional Files    180

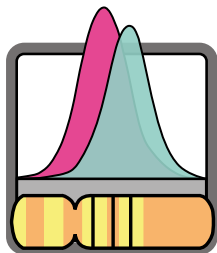## Additional file 1 — Example of a case study.    181

- **File name**: additional_file_1.pdf.    182
- **File format**: Portable Document Format (PDF).    183
- **Title**: Example Study Case.    184
- **Description**: Example of a study case exploring the PGS scores by Mavaddat et al. (2018).    185 186

## Additional file 2 — quincunx cheatsheet.    187

- **File name**: additional_file_2.pdf    188
- **File format**: Portable Document Format (PDF)    189
- **Title**: quincunx cheatsheet    190
- **Description**: Additional file 2 contains an infographics: quincunx cheatsheet.    191

**A**

| PGS CATALOG | RETRIEVAL FUNCTIONS | S4 CLASSES |
|---|---|---|
| | `get_scores()` | scores |
| | `get_publications()` | publications |
| | `get_sample_sets()` | sample_sets |
| | `get_performance_metrics()` | performance_metrics |
| | `get_traits()` | traits |

**B**

| Search by | Example | PGS | PGP | PSS | PPM | EFO |
|---|---|---|---|---|---|---|
| `pgs_id` | "PGS000001" | ■ | ■ | ■ | ■ | |
| `pgp_id` | "PGP000001" | | ■ | | | |
| `pss_id` | "PSS000001" | | | ■ | | |
| `ppm_id` | "PPM000001" | | | | ■ | |
| `efo_id` | "EFO_0000249" | ■ | | | | ■ |
| `pubmed_id` | "25855707" | ■ | ■ | | | ■ |
| `author` | "Mavaddat" | | ■ | | | |
| `trait_term` | "Alzheimer" | | | | | ■ |