Corresponding author(s): Geoffrey J. Barton

Last updated by author(s): Feb 2, 2024

# Reporting Summary

Nature Portfolio wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Portfolio policies, see our Editorial Policies and the Editorial Policy Checklist.

## Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

| n/a | Confirmed | |
|---|---|---|
| ☐ | ☒ | The exact sample size (*n*) for each experimental group/condition, given as a discrete number and unit of measurement |
| ☐ | ☒ | A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly |
| ☐ | ☒ | The statistical test(s) used AND whether they are one- or two-sided *Only common tests should be described solely by name; describe more complex techniques in the Methods section.* |
| ☐ | ☒ | A description of all covariates tested |
| ☐ | ☒ | A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons |
| ☐ | ☒ | A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals) |
| ☐ | ☒ | For null hypothesis testing, the test statistic (e.g. *F*, *t*, *r*) with confidence intervals, effect sizes, degrees of freedom and *P* value noted *Give P values as exact values whenever suitable.* |
| ☒ | ☐ | For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings |
| ☒ | ☐ | For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes |
| ☐ | ☒ | Estimates of effect sizes (e.g. Cohen's *d*, Pearson's *r*), indicating how they were calculated |

*Our web collection on statistics for biologists contains articles on many of the points above.*

## Software and code

Policy information about availability of computer code

| Data collection | Python 3.6.7; Arpeggio; DSSP 3.0.0; Hmmer 3.3.2; OC 2.1b |
|---|---|
| Data analysis | STAMP 4.4.2; Python 3.11; ProIntVar 0.1.0-patched; ProteoFAV 0.2.3; VarAlign; Biopython 1.74; Keras 2.10.0; Matplotlib 3.6.3; Numpy 1.24.1; Pandas 1.5.3; Scipy 1.10.0; Seaborn 0.12.2; Scikit-learn 1.2.2; Tensorflow 2.10.0 |

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Portfolio guidelines for submitting code & software for further information.

## Data

Policy information about availability of data

All manuscripts must include a data availability statement. This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A description of any restrictions on data availability
- For clinical datasets or third party data, please ensure that the statement adheres to our policy

Data is available at https://github.com/bartongroup/FRAGSYS (DOI: 10.5281/zenodo.10606595).

# Research involving human participants, their data, or biological material

Policy information about studies with human participants or human data. See also policy information about sex, gender (identity/presentation), and sexual orientation and race, ethnicity and racism.

| | |
|---|---|
| Reporting on sex and gender | No stratifications on sex employed. |
| Reporting on race, ethnicity, or other socially relevant groupings | No relevant stratifications were employed. |
| Population characteristics | Population variants were obtained from the gnomAD v2.1 dataset. |
| Recruitment | No-one was recruited for this study. |
| Ethics oversight | Ethical oversight was not applied for this study as no aspect required such oversight. |

Note that full information on the approval of the study protocol must also be provided in the manuscript.

# Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

☒ Life sciences  ☐ Behavioural & social sciences  ☐ Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see nature.com/documents/nr-reporting-summary-flat.pdf

# Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

| | |
|---|---|
| Sample size | Sample size of protein dataset was determined by the number of fragment screening experiments for which density maps were obtained using the PanDDA algorithm. Homologous sequences, and variants within them were determined by the size of the source databases (SwissProt &PDB). |
| Data exclusions | Fragment screening experiments of multi-protein complexes, as well as those proteins with no human homologues were removed from the dataset. |
| Replication | Our results originate from a global analysis of the structural and sequence properties of 1,500 crystallographically resolved proteins depicting the interaction of 35 proteins with hundreds of ligands. |
| Randomization | Samples were not allocated in any way in this work. "Groups" (e.g., clusters 1-4) were defined based on calculated properties. |
| Blinding | Blinding was not relevant as this was a computational, post-hoc analysis. |

# Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

## Materials & experimental systems

| n/a | Involved in the study |
|---|---|
| ☒ | ☐ Antibodies |
| ☒ | ☐ Eukaryotic cell lines |
| ☒ | ☐ Palaeontology and archaeology |
| ☒ | ☐ Animals and other organisms |
| ☒ | ☐ Clinical data |
| ☒ | ☐ Dual use research of concern |
| ☒ | ☐ Plants |

## Methods

| n/a | Involved in the study |
|---|---|
| ☒ | ☐ ChIP-seq |
| ☒ | ☐ Flow cytometry |
| ☒ | ☐ MRI-based neuroimaging |

## Plants

Seed stocks

-

Novel plant genotypes

-

Authentication

-