

Interpreting the effect of mutations to protein binding sites from large-scale genomic screens

Sara Jamshidi Parvar^a, Benjamin A Hall^b, David Shorthouse^{a,*}

^a UCL School of Pharmacy, 29-39 Brunswick Square, London WC1N 1AX, UK

^b UCL Department of Medical Physics and Biomedical Engineering, Mallet Place Engineering Building, University College London, Gower Street, London WC1E 6BT, UK

ARTICLE INFO

Keywords:

Binding calculations
FoldX
Protein structure
Missense mutation

ABSTRACT

Predicting the functionality of missense mutations is extremely difficult. Large-scale genomic screens are commonly performed to identify mutational correlates or drivers of disease and treatment resistance, but interpretation of how these mutations impact protein function is limited. One such consequence of mutations to a protein is to impact its ability to bind and interact with partners or small molecules such as ATP, thereby modulating its function. Multiple methods exist for predicting the impact of a single mutation on protein–protein binding energy, but it is difficult in the context of a genomic screen to understand if these mutations with large impacts on binding are more common than statistically expected. We present a methodology for taking mutational data from large-scale genomic screens and generating functional and statistical insights into their role in the binding of proteins both with each other and their small molecule ligands. This allows a quantitative and statistical analysis to determine whether mutations impacting protein binding or ligand interactions are occurring more or less frequently than expected by chance. We achieve this by calculating the potential impact of any possible mutation and comparing an expected distribution to the observed mutations. This method is applied to examples demonstrating its ability to interpret mutations involved in protein–protein binding, protein–DNA interactions, and the evolution of therapeutic resistance.

1. Introduction

With the increasing accessibility of genome sequencing methods, such as Illumina and NanoSeq methods, mutational data for a wide range of organisms and conditions are rapidly becoming available[1–3]. Genome sequencing screens are routinely being applied to study evolution of organisms[4], population diversity[5], and complex diseases such as cancer[6]. As the amount of this data increases, new methods are needed to enable interpretation and understanding of this high dimensional and functionally heterogeneous data. When a mutation is observed within a gene compared to a reference sequence, this mutation can be synonymous (does not change the amino acid composition of the protein encoded by the gene), or non-synonymous (somehow changes the amino acid composition of the protein encoded by the gene). Within non-synonymous mutations, missense mutations (which switch the amino acid at a specific site of the protein with another one) are some of the hardest to functionally interpret[7–9]. These mutations may induce any of a myriad of effects on a protein, including having no effect at all. Missense mutations may modify protein folding, alter its ability to

interact with protein partners or complexes, damage a functionally important region such as an enzymatic binding site, or have other impacts. Understanding how these mutations change a protein has large implications for comprehension of protein structure biophysics, treatment decisions, and design of drugs.

Calculations can be performed on a protein structure to predict the energetic effects of a missense mutation, thereby predicting quantitative changes in the stability and folding of a protein (referred to as the Gibbs free energy or $\Delta\Delta G$). We have previously applied these calculations to understand mutational function in disease[10–13]. These calculations can additionally be applied to the interaction energy between a protein and another molecule or protein, generating the $\Delta\Delta G$ of binding, a measure of how much a specific missense mutation energetically impacts the interface between the two molecules. Many methods exist for calculating the $\Delta\Delta G$ of binding, ranging from machine-learning based calculators[14,15], to “force-field” type methods that use chemical descriptions of atoms and their bonds to calculate energies[16], through to dynamical and “alchemical” methods that use computationally intensive dynamics-based methods[17,18]. In the same field of study, Brownian

* Corresponding author at: UCL School of Pharmacy, 29-39 Brunswick Square, London WC1N 1AX, UK.

E-mail address: d.shorthouse@ucl.ac.uk (D. Shorthouse).

<https://doi.org/10.1016/j.ymeth.2023.12.008>

Received 31 August 2023; Received in revised form 27 November 2023; Accepted 22 December 2023

Available online 5 January 2024

1046-2023/© 2024 The Authors. Published by Elsevier Inc. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

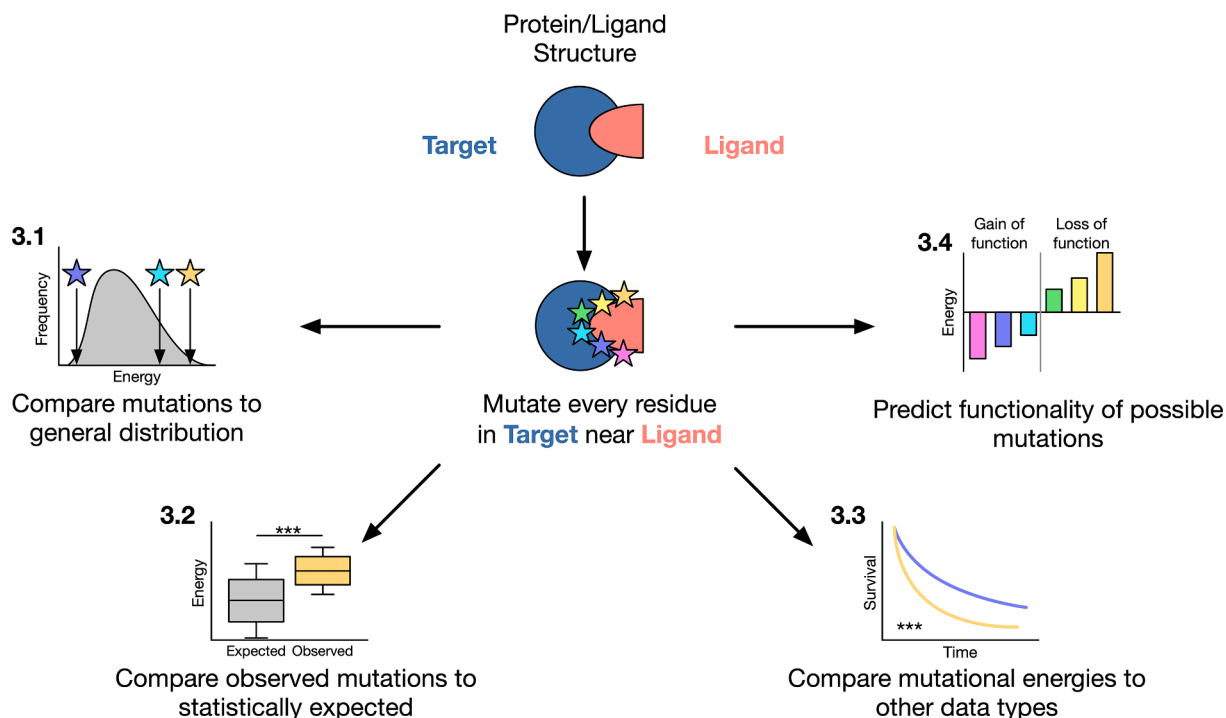


Fig. 1. Workflow for performing binding energy mutagenesis saturation screens. Results can then be used to study single mutations, groups of mutations, compare mutational energies to other data types, and predict the functionality of any potential mutation.

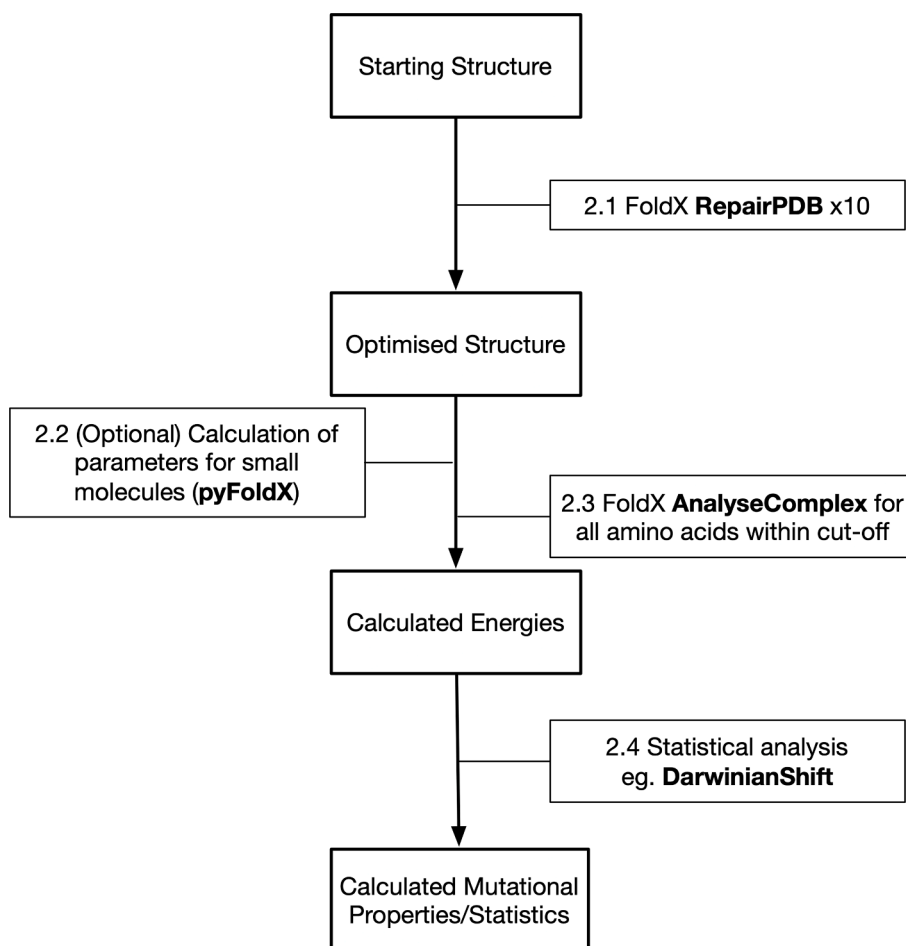


Fig. 2. Flow diagram for workflow presented in this manuscript. Numbers represent sections of this manuscript where the methods are explained.

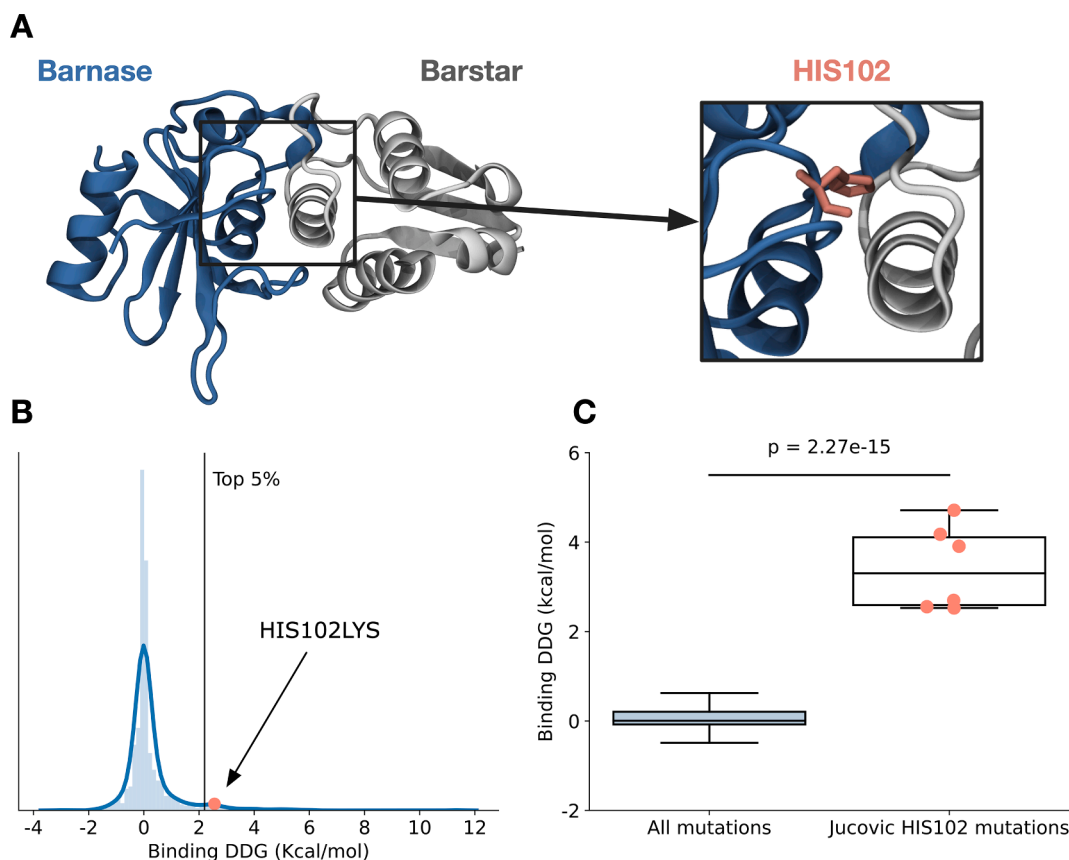


Fig. 3. Comparing single and groups of mutations. A - Structure of barnase-barstar complex with HIS102 highlighted in red (insert – right). B - Binding energy distribution showing the top 5% of mutations, and highlighting the location of HIS102LYS. C - Binding energy distributions for all mutations, and a subset of HIS102 mutations known to disrupt binding of barnase-barstar. P value represents independent *t*-test. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

dynamics calculations can also be used to estimate binding affinity of two biophysical objects, and provide the advantage of incorporating changes to environmental conditions and interactions with surfaces [19,20].

Interpretation of these energy predictions however can be challenging. Individual calculations may denote a large change in interaction energy, but knowing whether they are larger than an “expected” change due to randomness or the mutational background is necessary to evidence that they are under evolutionary selection, and are therefore highly likely to be functionally impactful. Mutation sets can be compared to all possible mutations to study whether they are enriched. Previous studies have developed statistical tests for analysing mutational data applied to protein structures, producing an “expected distribution” for a biophysical property [21]. In this work, we demonstrate how binding site saturation mutagenesis screens can be performed *in silico* using the tool FoldX [16], and how the data from these screens can be statistically analysed to understand mutational selection and functional impact. Similar results would be expected to be obtained using any method for calculating the estimated mutational $\Delta\Delta G$ such as Rosetta, as previous studies have reported high similarities between results from a range of methods [22–25]. We provide all code and analysis in the form of scripts and notebooks available at https://github.com/shorthouse-lab/binding_ddg, and apply our methods to a range of biological systems including protein–protein, protein–DNA, and protein–ligand interactions (Fig. 1).

2. Methods

We present methods for calculating binding energies – here applied

to the example case of the barnase-barstar complex 1BRS (<https://doi.org/10.2210/pdb1BRS/pdb>). All code required to replicate this workflow is provided at https://github.com/shorthouse-lab/binding_ddg. We explore analysis of the generated data for this and other complexes in the results section. This workflow involves the energetic optimisation of a target structure, the calculation of mutational energies for every possible amino acid mutation within a defined cutoff, and then statistical analysis of the results (Fig. 2). For this study, calculations were performed on a virtual machine hosted by NMRbox [26].

2.1. Structure preparation

Our workflow utilises the energy calculation methods included as part of FoldX [16]. Prior to performing calculations, structures need to undergo cleaning and energy minimisation. The structure is cleaned by removing excess molecules and complexes from the.pdb file. In the case of barnase-barstar complex, we used the.pdb file 1BRS. This file contains the structures for three barnase-barstar complexes; we opted to retain only chains A and D, which corresponds to one complex. Subsequently, this complex is energy minimised using the FoldX command ‘repairpdb’, performed ten times in sequence. As ‘repairpdb’ optimises each amino acid in the protein sequentially, conducting multiple rounds of repair result in a more optimised structure, determined by a lower overall calculated energy. We chose to repair the structure ten times to ensure the structure has reached a local minima, as for large structures multiple rounds of repair will result in a more optimal energy (An example of this for 3KMD can we found in the appendix). We used the following command to repair the pdb, a python script to do this ten times automatically is available at https://github.com/shorthouse-lab/binding_ddg:

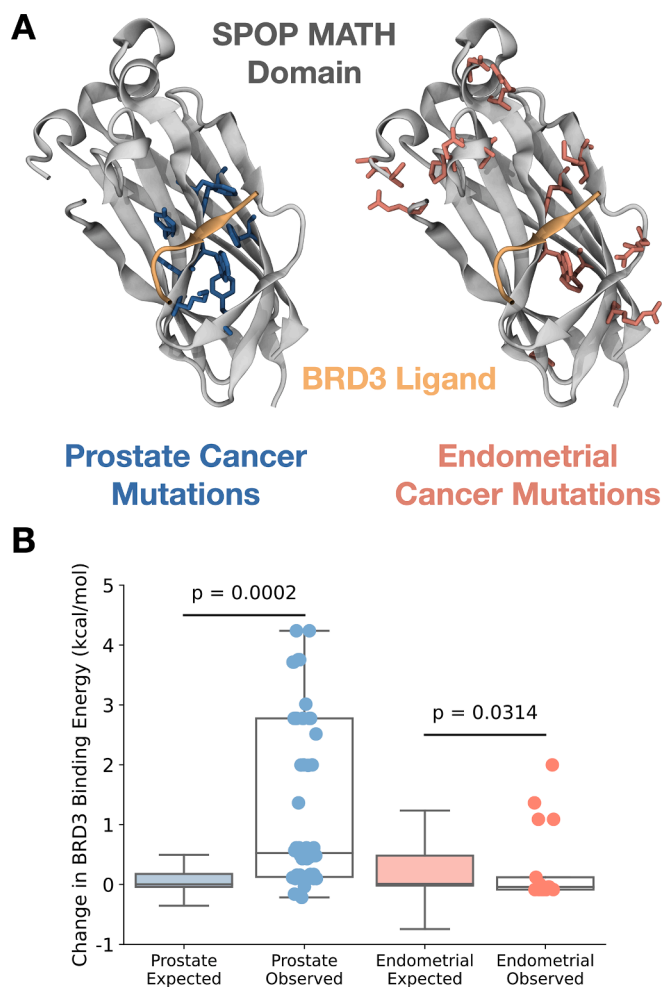


Fig. 4. Comparing mutations to an expected distribution. A - Structure of SPOP MATH domain bound to BRD3 ligand. Mutations observed in the TCGA in prostate cancer are shown in blue (left), mutations observed in the TCGA in endometrial cancer are shown in red (right). B - SPOP-BRD3 binding energy distributions for an expected mutational distribution generated from the mutational signature of the cancer, and the observed mutations for prostate and endometrial cancer in the TCGA. P value represents Monte Carlo Cumulative Distribution Function. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

```
$foldx --command=RepairPDB --pdb=1brs.pdb --water=CRYSTAL
--pH=7 --vdwDesign=2 --ionStrength=0.05
```

Inclusion of the “--water = CRYSTAL” command allows waters in the interface of the complex to be retained and included in energy calculations. “--pH” and “--ionStrength” flags enable to tool to perform calculations under physiologically relevant pH and ionic conditions.

2.2. Molecule parameterisation

FoldX recognises certain molecules, such as ATP, using prebuilt parameters. This allows the calculation of energetic effects of mutations on binding, potentially indicating mutations that increase or reduce binding affinity and consequently altering enzymatic activity. These calculations can be performed with a recognised molecule through setting the molecule as the “ligand”, and mutating all residues within a defined radius in the “target” protein. For other molecules however, parameters may not be available in FoldX, and so we must parameterise them. We used the existing method incorporated into pyFoldX[27] to perform

parameterisation, an example of which is included at https://github.com/shorthouse-lab/binding_ddg.

2.3. Binding DDG calculation

Binding DDG calculations are then performed on every residue of the “target” protein within a defined radius of the “ligand”. We chose a cut-off of 10 Å for calculations, and use Biopython[28] to calculate these, as FoldX is only expected to capture local changes in protein structures, and does not allow large conformational or backbone shifts. For each residue within the cut-off, we mutate it to each other amino acid using the FoldX command BuildModel:

```
$foldx --command=BuildModel --pdb=1brs_repair.pdb
--numberOfRuns=1 --pH=7 --vdwDesign=2 --ionStrength=0.05
--mutant-file=mutantfile.txt --water=CRYSTAL
```

Where “mutantfile.txt” is a file containing the mutation to induce as denoted by FoldX (Wildtype residue, chain, residue number, mutant residue). For example, to mutate the Threonine at position 100 of barnase (chain A) in 1brs.pdb to Alanine, the file will contain: “TA100A;”. Energy is then calculated using the ‘AnalyseComplex’ command applied to either a wildtype or mutant structure:

```
$foldx --command=AnalyseComplex --pdb=pdbfile.pdb
--analyseComplexChains=A,D --pH=7 --vdwDesign=2
--ionStrength=0.05 --water=CRYSTAL
```

By comparing the energies of the wildtype and mutant structures for each mutation, we can calculate the change in binding energy induced by each mutation (Binding $\Delta\Delta G$). Results are then aggregated to generate a single file containing the binding $\Delta\Delta G$ of every possible mutation to the “target” protein within the defined radius of the “ligand”. Scripts to automatically perform and summarise these calculations for any complex of interest are included at https://github.com/shorthouse-lab/binding_ddg.

2.4. Data analysis

Analysis is performed with python3, using the libraries pandas, numpy[29], and scipy[30]. Plotting is performed with matplotlib[31] and seaborn[32]. Survival analysis was performed using the lifelines python package[33]. To calculate evidence of mutational selection we used the library darwinian_shift (https://github.com/michaelhall28/darwinian_shift)[21]. Jupyter notebooks containing the code for all analysis presented in this paper are available at https://github.com/shorthouse-lab/binding_ddg.

3. Results

3.1. Barnase-barstar – Probing single and small groups of mutations

We applied this workflow to the barnase-barstar complex as a proof-of-principle. The barnase-barstar complex is a classically studied pair of proteins whereby barnase is an extracellularly secreted ribonuclease from *Bacillus Amyloliquefaciens*. Barnase is lethal to a bacterial cell that does not also express its inhibitor barstar, which tightly binds it and prevents activity[34,35]. The barnase-barstar structure 1BRS[36] was downloaded from the pdb, repaired, and binding $\Delta\Delta G$ calculated for every possible mutation in barnase. A conserved histidine (HIS102) in barnase is known to destabilise the protein complex when mutated to Lysine[37] (Fig. 3A). When plotting the distribution of binding $\Delta\Delta G$ for all potential mutations to barnase within 10 Å of barstar, we find that HIS102LYS, with a $\Delta\Delta G$ of 2.56 kcal/mol, is in the top 5 % (>2.2 kcal/

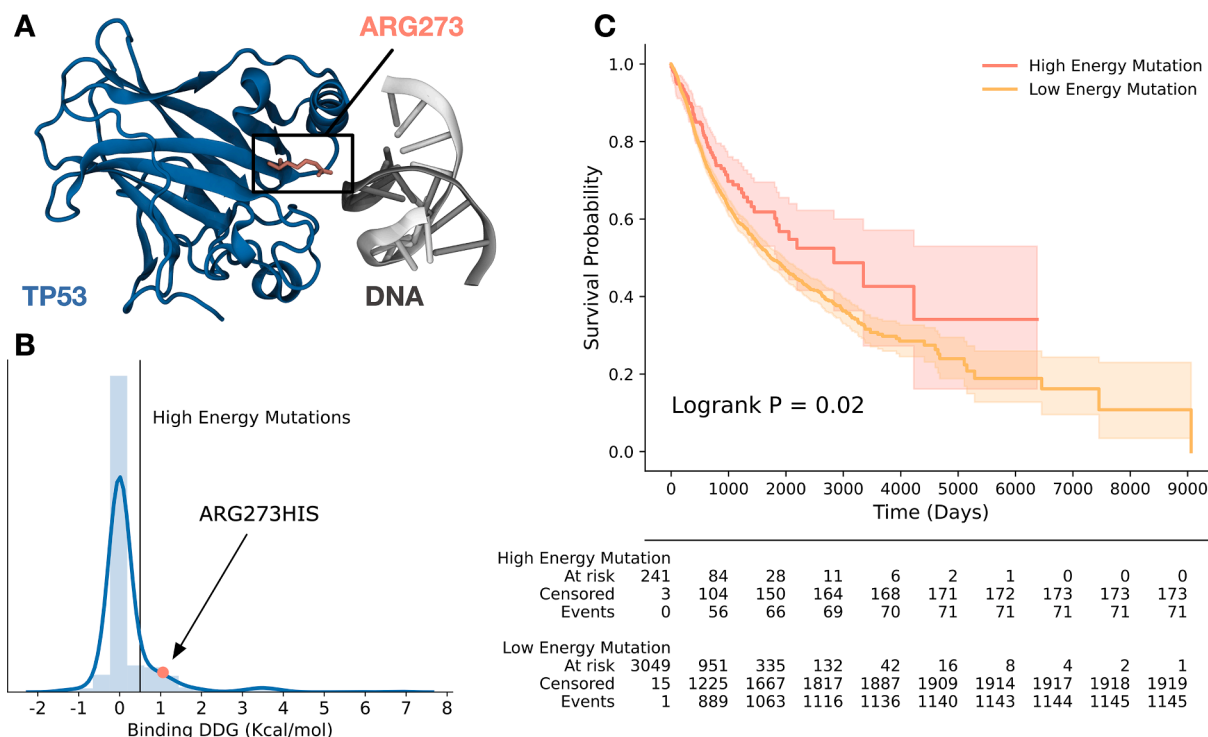


Fig. 5. Assessing mutational impact through comparison with other data. A – Structure of TP53 bound to DNA with ARG273 highlighted. B – TP53-DNA binding energy distribution for all mutations in TP53, highlighting ARG273HIS. C – Kaplan-Meier analysis of patients with high energy ($\Delta\Delta G > 0.5$ kcal/mol – red) and low energy ($\Delta\Delta G < 0.5$ kcal/mol – orange) mutations in TP53 across the pan-cancer TCGA dataset. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

mol) of all possible mutations (Fig. 3B). Furthermore, a set of other mutations to HIS102 have been experimentally determined to reduce binding affinity (HIS102ASP, HIS102GLN, HIS102GLY, HIS102LEU, HIS102ALA). By extracting the predicted binding $\Delta\Delta G$ for this group of mutations, we can compare them to the total population, and therefore calculate statistics to determine if they are outliers. We find that for these six mutations in barnase, they have a statistically significantly higher energy than the “base” distribution (independent *t*-test $p = 2.27 \times 10^{-15}$) and therefore are higher energy than would be expected (Fig. 3C). This example illustrates how our method of calculating binding $\Delta\Delta G$ for all possible mutations allows simple statistical comparison of individual or groups of observed mutations.

3.2. SPOP – Evidencing mutational selection

We next applied this workflow to an example of mutations to protein–protein binding in large-scale genomics datasets to identify evolutionary selection. SPOP is an E3 ligase involved in degradation of protein targets, and is known to be under selection in a number of cancers[38]. Notably, SPOP missense mutations occurring in prostate cancer are thought to be loss of function (LoF), decreasing SPOP activity, whilst missense mutations occurring in endometrial cancer are thought to be gain of function (GoF), increasing protein activity[39] (Fig. 4A). We used the structure of SPOP bound to its protein–ligand target BRD3 (pdb id: 6I41) to calculate every possible mutation within 10 Å of the ligand. We downloaded the mutational data for a large number of patients with prostate[40] ($n = 494$) and endometrial[40] ($n = 447$) cancers from The Cancer Genome Atlas. Both cancer cohorts have numerous observed mutations within the ligand binding domain, and we have used a previously published method known as “Darwinian shift”[21] to generate expected distributions - distributions of binding $\Delta\Delta G$ s that would be expected if mutations to the binding site were only subject to the mutational signature of the cancer. If our observed mutations are significantly different to this expected distribution, then we can

conclude that there is evidence of evolutionary selection. We find that mutations near the ligand of SPOP in prostate cancer are statistically significantly (CDF Monte Carlo $p = 0.0002$) higher energy than expected, and therefore would likely decrease the ability of the protein to bind BRD3, whilst the mutations in endometrial cancer are significantly (CDF Monte Carlo $p = 0.0314$) lower energy than expected, and therefore do not significantly increase binding energy (Fig. 4B). This demonstrates how binding $\Delta\Delta G$ calculations can be used to explore evolutionary selection in large-scale genomics datasets, with the example of understanding selection of ligand binding mutations in SPOP.

3.3. TP53 – Calculating binding to non-protein molecules

TP53 is a DNA binding protein playing a role in a majority of human cancers[41]. Most human cancers mutate TP53 to some extent, and evidence suggests that different types of mutations have different effects on cell phenotype, tumour progression, and treatment response[42]. Many mutations to TP53 impact its ability to bind DNA, whilst others damage its ability to tetramerise, and potentially have numerous other effects. To illustrate how our methodology can explore mutational effects and their relation to other data sources, we first calculated the binding $\Delta\Delta G$ for all mutations in the TP53 structure 3KMD within 10 Å of the cocrystallised DNA fragment (Fig. 5A). Energies follow a broadly normal distribution, and the known DNA contact disrupting hotspot mutation ARG273HIS is found to have a comparatively high mutational energy (1.059 kcal/mol, within the top 8 % of mutations) (Fig. 5B). We next collected the mutational and clinical data from the TCGA pan-cancer atlas[6], and extracted the clinical data for all patients with a missense mutation in TP53. We subset mutations into those expected to have any (even minor) destabilising effect on TP53-DNA binding (denoted “high energy” - those with a $\Delta\Delta G > 0.5$ kcal/mol), and those not expected to reduce the ability of TP53 to bind DNA (denoted “low energy” - those with a $\Delta\Delta G < 0.5$ kcal/mol). We find a statistically

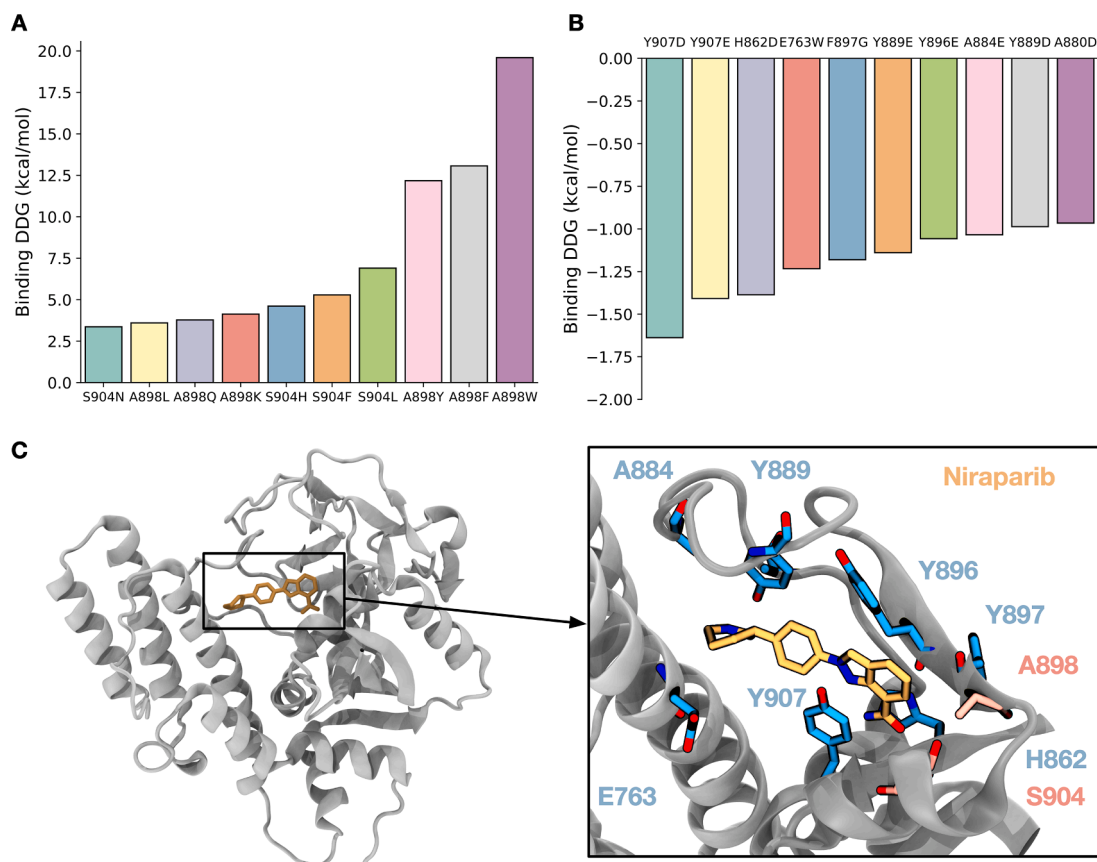


Fig. 6. Predicting mutational effects from binding $\Delta\Delta G$. A – Highest ten mutations by $\Delta\Delta G$ for the interaction between PARP1 and the inhibitor niraparib. B – Lowest ten mutations by $\Delta\Delta G$ for the interaction between PARP1 and the inhibitor niraparib. C – Structure of PARP1-niraparib (Orange), highlighting the residues with high (blue) and low (red) energy mutations. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

significant (logrank $P = 0.02$) difference in overall survival for our pan-cancer cohort ($n = 3290$), and patients with high energy TP53-DNA binding mutations have a better overall survival (Fig. 5C). This suggests that mutations that do not disrupt DNA binding lead to poorer clinical outcomes, possibly because their biophysical effects are more potent, or they are more likely to respond to treatment. With this example we demonstrate how to apply our method to correlate biophysical properties of mutations in protein structures with other information such as clinical or disease outcome data.

3.4. PARP inhibition by niraparib – Predicting resistance mutations

A final potential use of this workflow is demonstrated by prediction of mutations that might lead to resistance to drugs. Mutations in proteins that reduce the ability of a drug to bind are likely to lead to reduced toxicity to those cells, and so these mutations will be selected for in competitive conditions. This is a large problem for antibiotics, where resistance is emerging and becoming increasingly prevalent, but also for chemotherapy, whereby introduction of targeted therapies provides evolutionary pressure on the tumour to evolve resistance around them [43]. We chose to demonstrate this workflow with the example of the targeted ovarian cancer treatment, niraparib, in complex with its target protein PARP1. The niraparib-PARP1 structure has been solved (4R6E) [44]. We again repeated the presented methodology – first repairing the structure ten times before calculating all possible mutations within 10 Å of the drug. For this example, we had to generate parameters for niraparib. We chose to use pyFoldX[27] to generate parameters - each atom in the molecule is assigned to an ‘atomtype’ included in FoldX. This resulted in a parameter file for niraparib that, when included in the

directory where calculations are performed, allows FoldX to recognise small molecules and perform energy calculations on them.

We calculated every possible mutation to PARP1 in the vicinity of niraparib, surmising that mutations which increase the binding energy significantly are likely to confer resistance. We find that within the top ten mutations by binding $\Delta\Delta G$, only two amino acids are represented – residues 898 and 904, both of which have direct contact with the ligand (Fig. 6A, C). In particular, mutation ALA898TRP is predicted to have an extremely high $\Delta\Delta G$ value of 19.60 kcal/mol, suggesting that it will significantly reduce the ability of niraparib to bind to PARP1. Similarly, looking at mutations that reduce $\Delta\Delta G$, we can infer that these mutations may increase niraparib binding affinity and thus sensitise the protein to inhibition. The top two mutations that reduce binding $\Delta\Delta G$ are both to residue TYR907 (Fig. 6B, C). Interestingly, this residue is known to be phosphorylated by c-Met, which subsequently reduces the effect of PARP inhibitors[45]. This highlights that whilst these predictions can give some insights into mutational effects on protein activity, biological context is required in order to fully understand and interpret them.

4. Discussion

These methods provide a toolkit for interpreting functionality of mutations to protein binding. Whilst we demonstrate that this method can interpret single or small groups of mutations, its most powerful application is in the interpretation of large-scale missense mutational data, which allows the generation of a mutational ‘expected’ signature to compare against. We chose to use the empirical energy calculation method available in FoldX[16] as we feel that this method is a good trade-off between computational time (readily parallelisable and does

not involve dynamics calculations) and accuracy, as it has been shown to be comparable to other available methods[22–24]. However, it is possible to apply these methods to any distribution of energy calculations calculated by a different method, such as Monte Carlo sampling-based rosetta[17], or significantly more computationally intensive but higher accuracy alchemy methods[18]. Ultimately however, this method is dependent on the accuracy of the underlying energy calculations, and whilst evidence suggests that average accuracy of FoldX is acceptable[24], it is not necessarily accurate for predicting individual mutations. We also note that this methodology is limited to cases where protein structures are available, and whilst recent advances in AI based structure and protein–protein interface predictions mitigate some of this [46,47], it is unknown exactly how reliable individual predicted structures are, and therefore how much weight can be put on the outputs from their predictions.

5. Conclusions

We present a workflow for performing and analysing calculations to understanding the impact of missense mutations on protein–ligand binding, with a focus on “big data”. As increasingly large datasets of mutations are generated, methodologies such as these which enable interpretation of this complex and large information become ever more important. Our methodology involves calculating all possible mutations to an interaction site to generate a base distribution that can be used for statistical comparisons. We apply this method to single and small groups of mutations, data from large-scale genomic screens, use our results to understand clinical and other data types, and use it predictively to study potential effects of unobserved mutations. This workflow provides a toolkit for improving our understanding of notoriously hard to functionally interpret whole genome sequencing and CRISPR screens with reasonable computational cost and space requirements – for the example of 3KMD, this study generated 1 GB of data, taking just 6 h on a 40 CPU virtual machine. The workflow presented involves the use of FoldX, but can be applied to any current or future methods that predict or calculate the binding $\Delta\Delta G$, including Rosetta, which has been demonstrated to generate comparable results[10,23,24]. This methodology has applications in areas of genomics research such as biomedical study of diseases like cancer, evolutionary biology, and protein or therapeutic engineering. In an effort to make this workflow as accessible as possible, we include code to implement these analyses at https://github.com/shorthouse-lab/binding_ddg.

Funding

SJP is funded through the EPSRC & SFI Centre for Doctoral Training in Transformative Pharmaceutical Technologies, United Kingdom (grant no. EP/S023054/1). BAH acknowledges support from the Royal Society, United Kingdom (grant no. UF130039).

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

All data is included or linked in the methods or available at the associated GitHub: https://github.com/shorthouse-lab/binding_ddg

Appendix A. Supplementary material

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.ymeth.2023.12.008>.

References

- [1] E. Turro, W.J. Astle, K. Megy, S. Gräf, D. Greene, O. Shamardina, H.L. Allen, A. Sanchis-Juan, M. Frontini, C. Thys, J. Stephens, R. Mapeta, O.S. Burren, K. Downes, M. Haimel, S. Tuna, S.V.V. Deevi, T.J. Aitman, D.L. Bennett, P. Calleja, K. Carss, M.J. Caulfield, P.F. Chinnery, P.H. Dixon, D.P. Gale, R. James, A. Koziell, M.A. Laffan, A.P. Levine, E.R. Maher, H.S. Markus, J. Morales, N.W. Morrell, A. D. Mumford, E. Ormondroyd, S. Rankin, A. Rendon, S. Richardson, I. Roberts, N.B. A. Roy, M.A. Saleem, K.G.C. Smith, H. Stark, R.Y.Y. Tan, A.C. Themistocleous, A. J. Thrasher, H. Watkins, A.R. Webster, M.R. Wilkins, C. Williamson, J. Whitworth, S. Humphray, D.R. Bentley, S. Abbs, L. Abulhoul, J. Adlard, M. Ahmed, H. Alachkar, D.J. Allsup, J. Almeida-King, P. Ancliff, R. Antrobus, R. Armstrong, G. Arno, S. Ashford, A. Attwood, P. Aurora, C. Babbs, C. Bacchelli, T. Bakchoul, S. Banka, T. Bariana, J. Barwell, J. Batista, H.E. Baxendale, P.L. Beales, D. R. Bentley, A. Bierzynska, T. Biss, M.A.K. Bitner-Glindzicz, G.C. Black, M. Bleda, I. Blesneac, D. Bockenbauer, H. Bogaard, C.J. Bourne, S. Boyce, J.R. Bradley, E. Bragin, G. Breen, P. Brennan, C. Brewer, M. Brown, A.C. Browning, M. J. Browning, R.J. Buchan, M.S. Buckland, T. Bueser, C.B. Diz, J. Burn, S.O. Burns, O.S. Burren, N. Burrows, C. Campbell, G. Carr-White, K. Carss, R. Casey, J. Chambers, J. Chambers, M.M.Y. Chan, C. Cheah, F. Cheng, P.F. Chinnery, M. Chitre, M.T. Christian, C. Church, J. Clayton-Smith, M. Cleary, N.C. Brod, G. Coghlan, E. Colby, T.R.P. Cole, J. Collins, P.W. Collins, C. Colombo, C. J. Compton, R. Condliffe, S. Cook, H.T. Cook, N. Cooper, P.A.A. Corris, A. Furnell, F. Cunningham, N.S. Curry, A.J. Cutler, M.J. Daniels, M. Dattani, L.C. Daugherty, J. Davis, A. De Soya, S.V.V. Deevi, T. Dent, C. Deshpande, E.F. Dewhurst, P. H. Dixon, S. Douzgou, K. Downes, A.M. Drazzyk, E. Drewe, D. Duarte, T. Dutt, J.D. M. Edgar, K. Edwards, W. Egnor, M.N. Ekani, P. Elliott, W.N. Erber, M. Erwood, M. C. Estiu, D.G. Evans, G. Evans, T. Everington, M. Eyries, H. Fassihi, R. Favier, J. Findhammer, D. Fletcher, F.A. Flinter, R.A. Floto, T. Fowler, J. Fox, A.J. Frary, C. E. French, K. Freson, M. Frontini, D.P. Gale, H. Gall, V. Ganesan, M. Gattens, C. Geoghegan, T.S.A. Gerighty, A.G. Gharavi, S. Ghio, H.A. Ghofrani, J.S.R. Gibbs, K. Gibson, K.C. Gilmour, B. Girerd, N.S. Gleadall, S. Goddard, D.B. Goldstein, K. Gomez, P. Gordins, D. Gosal, S. Gräf, J. Graham, L. Grassi, D. Greene, L. Greenhalgh, A. Greinacher, P. Gresele, P. Griffiths, S. Grigoriadou, R.J. Grocock, D. Grozeva, M. Gurnell, S. Hackett, C. Hadinnapola, W.M. Hague, R. Hague, M. Haimel, M. Hall, H.L. Hanson, E. Hague, K. Harkness, A.R. Harper, C.L.L. Harris, D. Hart, A. Hassan, G. Hayman, A. Henderson, A. Herwadkar, J. Hoffman, S. Holden, R. Horvath, H. Houlden, A.C.C. Houweling, L.S. Howard, F. Hu, G. Hudson, J. Hughes, A.P. Huissoon, M. Humbert, S. Humphray, S. Hunter, M. Hurles, M. Irving, L. Izatt, S.A. Johnson, S. Jolles, J. Jolley, D. Josifova, N. Jurkute, T. Karten, J. Karten, M.A. Kasanicki, H. Kazkaz, R. Kazmi, P. Kelleher, A.M. Kelly, W. Kelsall, C. Kempster, D.G. Kiely, N. Kingston, R. Klima, N. Koelling, M. Kostadima, G. Kovacs, A. Koziell, R. Kreuzhuber, T.W. Kuijpers, A. Kumar, D. Kumararatne, M.A. Kurian, M.A. Laffan, F. Lalloo, M. Lambert, A. Lawrie, D. M. Layton, N. Lench, C. Lentaigne, T. Lester, A.P. Levine, R. Linger, H. Longhurst, L.E. Lorenzo, E. Louka, P.A. Lyons, R.D. Machado, R.V. MacKenzie Ross, B. Madan, E.R. Maher, J. Maimaris, S. Malka, S. Mangles, R. Mapeta, K.J. Marchbank, S. Marks, H.U. Marschall, A. Marshall, J. Martin, M. Mathias, E. Matthews, H. Maxwell, P. McAlinden, M.I. McCarthy, H. McKinney, A. McMahon, S. Meacham, A.J. Mead, I.M. Castello, K. Megy, S.G.G. Mehta, M. Michaelides, C. Millar, S.N. Mohammed, S. Moledina, D. Montani, A.T. Moore, J. Morales, N. W. Morrell, M. Mozere, K.W. Muir, A.D. Mumford, A.H. Nemeth, W.G. Newman, M. Newnham, S. Noorani, P. Nurden, J. O'Sullivan, S. Obaji, C. Odhams, S. Okoli, A. Olschewski, H. Olschewski, K.R. Ong, S.H. Oram, E. Ormondroyd, W. H. Ouwehand, C. Palles, S. Papadia, S.M. Park, D. Parry, S. Patel, J. Paterson, A. Peacock, S.H.H. Pearce, J. Peden, K. Peerlinck, C.J. Penkett, J. Pepke-Zaba, R. Petersen, C. Pilkington, K.E.S. Poole, R. Prathalingam, B. Psaila, A. Pyle, R. Quinton, S. Rahman, S. Rankin, A. Rao, F.L. Raymond, P.J. Rayner-Matthews, C. Rees, T. Renton, C.J. Rhodes, A.S.C. Rice, S. Richardson, A. Richter, L. Robert, I. Roberts, A. Rogers, S.J. Rose, R. Ross-Russell, C. Roughley, N.B.A. Roy, D. M. Ruddy, O. Sadeghi-Alavijeh, M.A. Saleem, N. Samani, C. Samaraghitana, A. Sanchis-Juan, R.B. Sargur, R.N. Sarkany, S. Satchell, S. Savic, J.A. Sayer, G. Sayer, L. Scelsi, A.M. Schaefer, S. Schulman, R. Scott, M. Scully, C. Searle, W. Seeger, A. Sen, W.A.C. Sewell, D. Seyres, N. Shah, O. Shamardina, S.E. Shapiro, A.C. Shaw, P.J. Short, K. Sibson, L. Side, I. Simeoni, M.A.A. Simpson, M.C. Sims, S. Sivapalaratnam, D. Smedley, K.R. Smith, K. Snape, N. Soranzo, F. Soubrier, L. Southgate, O. Spasic-Boskovic, S. Staines, E. Staples, H. Stark, J. Stephens, C. Steward, K.E. Stirrups, A. Stuckey, J. Suntharalingam, E.M. Swietlik, P. Syrris, R. C. Tait, K. Talks, R.Y.Y. Tan, K. Tate, J.M. Taylor, J.C. Taylor, J.E. Thaventhiran, A. C. Themistocleous, E. Thomas, D. Thomas, M.J. Thomas, P. Thomas, K. Thomson, A.J. Thrasher, G. Threadgold, C. Thys, T. Tilly, M. Tischkowitz, C. Titterton, J. A. Todd, C.H. Toh, B. Tolhuis, I.P. Tomlinson, M. Toshner, M. Traylor, C. Treacy, P. Treadaway, R. Trembath, S. Tuna, W. Turek, E. Turro, P. Twiss, T. Vale, C. Van Geet, N. van Zuydam, M. Vandekuilen, A.M. Vandersteen, M. Vazquez-Lopez, J. von Ziegenweid, A.V. Noordegraaf, A. Wagner, Q. Waisfisz, S.M. Walker, N. Walker, K. Walter, J.S. Ware, H. Watkins, C. Watt, A.R. Webster, L. Wedderburn, W. Wei, S.B. Welch, J. Wessels, S.K. Westbury, J.P. Westwood, J. Wharton, D. Whitehorn, J. Whitworth, A.O.M. Wilkie, M.R. Wilkins, C. Williamson, B. T. Willson, E.K.S. Wong, N. Wood, Y. Wood, C.G. Woods, E.R.R. Woodward, S. J. Wort, A. Worth, M. Wright, K. Yates, P.F.K. Yong, T. Young, P. Yu, P. Yu-Wai-Man, E. Zlamalova, N. Kingston, N. Walker, C.J. Penkett, K. Freson, K.E. Stirrups, F. L. Raymond, Whole-genome sequencing of patients with rare diseases in a national health system, *Nature* 583 (2020), <https://doi.org/10.1038/s41586-020-2434-2>.
- [2] D.M. Altshuler, R.A. Gibbs, L. Peltonen, S.F. Schaffner, F. Yu, E. Dermitzakis, P. E. Bonnen, P.L.W. De Bakker, P. Deloukas, S.B. Gabriel, R. Gwilliam, S. Hunt, M. Inouye, X. Jia, M. Aarno Palotie, P. Parkin, K. Whittaker, A. Chang, L.R. Hawes,

- enhances anti-tumor effects of PARP inhibitors, *Nat. Med.* 22 (2016), <https://doi.org/10.1038/nm.4032>.
- [46] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. Žídek, A. Potapenko, A. Bridgland, C. Meyer, S.A. Kohl, A.J. Ballard, A. Cowie, B. Romera-Paredes, S. Nikolov, R. Jain, J. Adler, T. Back, S. Petersen, D. Reiman, E. Clancy, M. Zielinski, M. Steinegger, M. Pacholska, T. Berghammer, S. Bodenstein, D. Silver, O. Vinyals, A.W. Senior, K. Kavukcuoglu, P. Kohli, D. Hassabis, Highly accurate protein structure prediction with AlphaFold, *Nature* 596 (2021), <https://doi.org/10.1038/s41586-021-03819-2>.
- [47] D.F. Burke, P. Bryant, I. Barrio-Hernandez, D. Memon, G. Pozzati, A. Shenoy, W. Zhu, A.S. Dunham, P. Albanese, A. Keller, R.A. Scheltema, J.E. Bruce, A. Leitner, P. Kundrotas, P. Beltrao, A. Elofsson, Towards a structurally resolved human protein interaction network, *Nat. Struct. Mol. Biol.* 30 (2023), <https://doi.org/10.1038/s41594-022-00910-8>.