

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

The Landscape of Mutations in Fumarate Hydratase

David Shorthouse¹, Michael W J Hall^{1,2}, Benjamin A Hall^{1}*

*Correspondence to be addressed to B.A.H (bh418@mrc-cu.cam.ac.uk)

¹ MRC Cancer Unit, University of Cambridge, Cambridge, CB2 0XZ, UK

² Wellcome Trust Sanger Institute, Hinxton, CB10 1SA, UK

26 **ABSTRACT**

27 Fumarate Hydratase (FH) is an enzyme of the citric acid (TCA) cycle that is
28 responsible for reversibly catalysing the conversion between fumarate and malate.
29 FH loss and subsequent buildup of the oncometabolite fumarate causes hereditary
30 leiomyomatosis and renal cell carcinoma.

31 We sought to explore the mutational landscape of FH in silico, to predict the
32 functional effects of many detected mutations, and categorise detected but un-
33 characterised mutations in human populations. Using mutational energy predicting
34 tools such as Rosetta and FoldX we can accurately predict mutations and mutational
35 hotspots with high disruptive capability. Furthermore, through performing molecular
36 dynamics simulations we show that hinge regions of the protein can be stabilized or
37 destabilized by mutations, with new mechanistic implications of the consequences
38 on the binding affinity of the enzyme for its substrates.

39 We can additionally categorise a large majority of mutations and potential mutations
40 into functional groups. This allows us to predict which detected mutations in the
41 human population are likely to be loss-of-function, and therefore predispose patients
42 to papillary renal carcinoma through considering only mutations to the protein
43 binding site, hinges, and those that are buried deep within the protein. We
44 additionally link mutation data to publicly available metabolomics data, and show that
45 we can accurately predict which mutations in cancer cell lines are functionally
46 relevant.

47

48

49

50

51

52

53

54

55 INTRODUCTION

56 Fumarate hydratase (FH) is a member of the tricarboxylic acid (TCA) cycle occurring
57 in the mitochondria, and enzymatically metabolises fumarate within the cytosol. FH
58 activity in the cell is responsible for the reversible conversion of the metabolite
59 fumarate into malate, and the knockout or mutational inactivation of FH in kidneys is
60 linked to an oncogenically-associated buildup of fumarate^{1,2}. As a result the enzyme
61 FH has been described as a tumor suppressor, and fumarate as part of a novel
62 classification of molecules named “oncometabolites”. Precisely how the buildup of
63 fumarate can be oncogenic is unknown, but recent work points towards suppression
64 of DNA repair responses, EMT, and promotion of mitotic entry upon fumarate
65 buildup³⁻⁵.

66 Understanding the effects of mutations on the activity and assembly of FH is of
67 importance for the understanding and stratification of germline mutations in FH,
68 which can predispose patients with a single mutated or deleted allele to hereditary
69 leiomyomatosis and renal cell cancer (HLRCC) upon mutational inactivation of their
70 remaining wild-type copy^{6,7}. Previous work has identified mutants linked with
71 inherited and de-novo FH-related conditions, including cancer⁸ – most notably, the
72 FH mutation database represents a comprehensive list of mutations and their
73 effects, if known, on FH activity⁹.

74 In recent years numerous methods have been developed for estimating the effects of
75 single point mutations (SNPs) on the stability of a protein structure in silico. Notable
76 methods include FoldX^{10,11}, which uses an empirical force field to predict the
77 alterations in a protein induced by mutation, and methods included as part of the
78 Rosetta suite^{12,13}, which uses Monte-Carlo based dynamics to predict energetic
79 effects of mutations. Additionally, molecular dynamics can be used to more
80 comprehensively investigate mutant protein structure, though at significantly higher
81 computational cost. With the advent of high-throughput methods such as CRISPR
82 screening, and larger projects being undertaken to screen populations for mutations
83 and disease, coupled with large-scale disease-focussed data generating projects
84 such as The Cancer Genome Atlas (TCGA)¹⁴ and the International Cancer Genome
85 Consortium (ICGC)¹⁵, the number and diversity of mutations being implicated in
86 disease is rapidly expanding. Whilst methods to attempt to sift functionally relevant

87 mutations from synonymous to detect highly mutated genes exist in the form of
88 statistical tests such as DN/DS¹⁶, mutsig¹⁷, and oncodrive¹⁸, including some methods
89 that take into account structure of the protein such as Rhapsody¹⁹, there is scope for
90 detailed, structure-informed, chemically aware methods to classify mutations,
91 including those not yet observed, into Loss-of-Function (LOF) and benign categories.

92 Here we computationally screen and classify every potential mutation in the available
93 fumarate hydratase structure to study the landscape of potential mutations. We
94 consider the structural and biological implications of each mutation, and thus can
95 predict mechanistic details of every potential mutant. We confirm that our method
96 predicts known functionally relevant mutations, and predict from existing databases
97 of mutations which have an unknown effect, which of them will be damaging to the
98 activity of FH. Overall we predict that 66% of all mutations to FH influence activity or
99 assembly. We further validate our predictions through studying the Cancer Cell Line
100 Encyclopaedia (CCLE)^{20,21} and show that previously unstudied mutations that we
101 predict to be damaging to the function of FH result in altered metabolite levels
102 expected from disruption to the activity of FH.

103

104

105

106

107

108

109

110

111

112

113

114

115

116 RESULTS

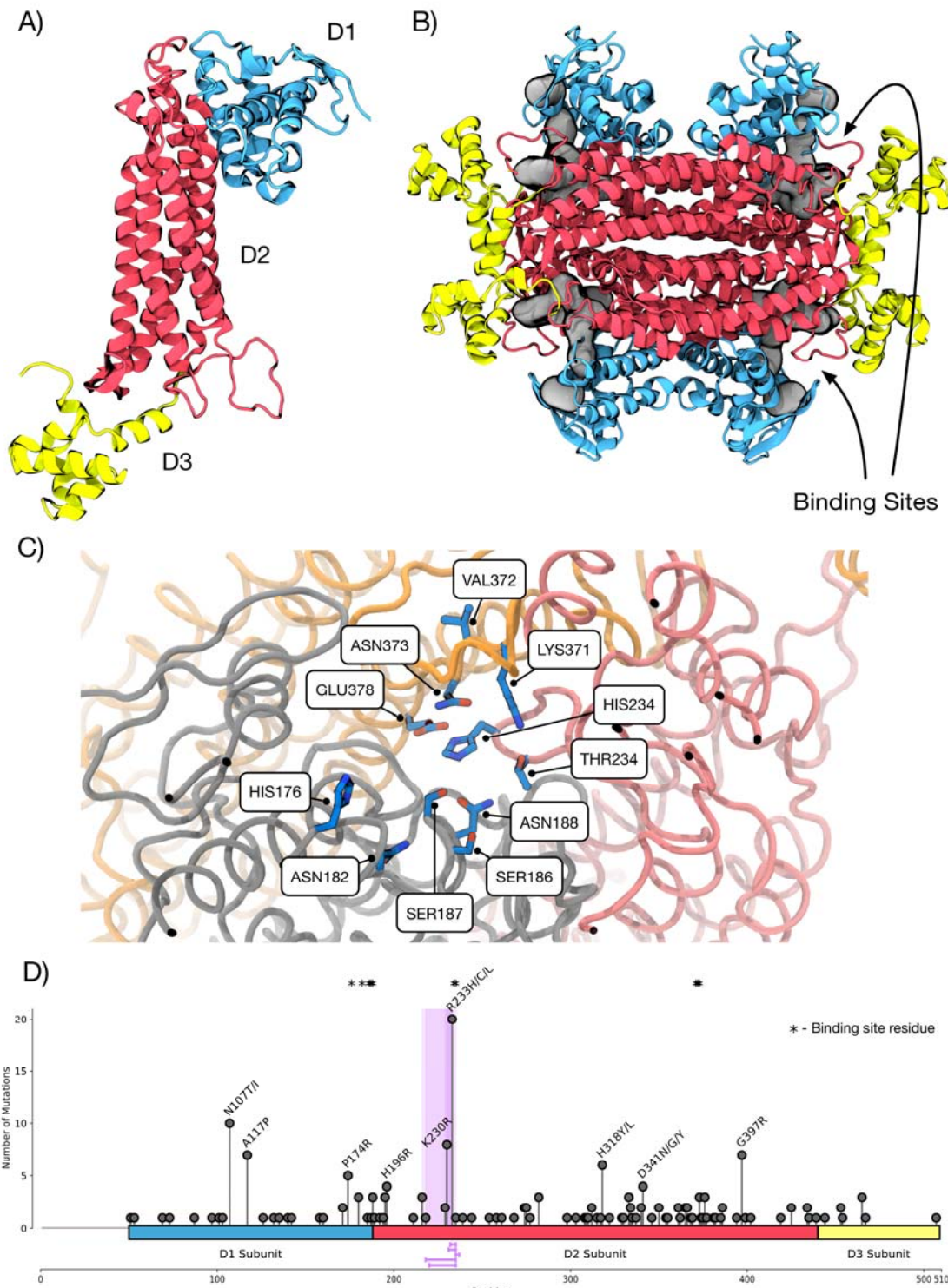
117 Evidence of mutational clustering in FH

118 Human FH is formed as a homotetramer of subunits generated from the *FH* gene.
119 Each subunit contains 3 domains, Domain 1, Domain 2, and Domain 3 (D1, D2, and
120 D3 respectively) (**Figure 1 A**). D1 is formed from residues in the range 49-188, D2 is
121 formed from residues in the range 189-439, and D3 from residues in the range 440-
122 510. The full functional protein is an assembly of 4 subunits and contains 4 identical
123 binding pockets made of interactions between 3 subunits (**Figure 1 B**). There are
124 two proposed regions of importance for catalysis of the fumarate/malate conversion;
125 Site A, the known active site (hereafter referred to as the binding site), and Site B, a
126 region of proposed but unknown functional importance^{22,23}. For this study we chose
127 to only include the known catalytic site, Site A, defined as residues HIS176, ASN182,
128 SER186, SER187, ASN188, THR234, HIS235, LYS371, VAL372, ASN373, and
129 GLU378 (**Figure 1 C**). We do not consider Site B due to the unknown and conflicting
130 evidence surrounding its importance. For this study we chose to focus on the crystal
131 structure 5upp²⁴, which covers residues 49-510 of the 510 residue protein
132 assembled into a homotetramer.

133 To study mutations known or suspected to have roles in human disease, we
134 investigated the Fumarate Hydratase Mutation Database⁹, which contains 378
135 mutations, including 113 that are distinct missense, at the time of this study. The
136 Fumarate Hydratase Mutation Database attempts to pool all observed mutations in
137 FH, including those that are benign, and a large number of mutations have no clinical
138 or functional annotation. Mutations that are not known to be benign (i.e those either
139 labelled as loss-of-function, or those which are uncharacterised) are shown in **Figure**
140 **1 D**. In particular, mutations at amino acids 107, 117, 230, and 233 are reported at a
141 higher frequency than other mutations and may indicate regions of mutational
142 vulnerability in the structure.

143 We applied the NMC clustering method to look for clustering of mutations across the
144 1D sequence of the protein²⁵. We chose to include the top 5 predicted clusters,
145 ranked by significance, and with a size less than 50 residues long. We find the most
146 significant clusters are all within the region of the more prevalent mutations in

147 residues 230 and 233, indicating that this region is statistically highly over mutated,
 148 and potentially a mutationally vulnerable site.



149 **Figure 1:** Structure and current mutations in Fumarate Hydratase. A) Structure of a
 150 single subunit of FH showing the D1, D2, and D3 regions. B) Structure of an
 151

152 assembled homotetramer of FH. Binding sites are highlighted and made up on an
153 interface between 3 subunits. C) Close up of the binding site of FH showing the
154 residues involved in catalytic activity. D) Mutational spectrum of non-benign SNPs in
155 FH. D1, D2, and D3 regions are highlighted in blue, red, and yellow respectively.
156 Stars indicate residues involved in catalytic activity that make up the binding site of
157 FH. Purple highlight and lines represent the top 5 mutational clusters as calculated
158 by the NMC algorithm.

159

160 **Classification of mutations by proximity to the binding site and protein hinges**

161 Residues of the catalytic site in FH have been previously identified as essential for
162 the conversion of fumarate to malate. We chose to define binding site-associated
163 residues as those with alpha-carbons (CA) within 6 Å of the CA of any binding site
164 residue. Plotting the resultant distances for each residue in the FH structure shows
165 that specific clusters of residues in the vicinity of the binding site are also significantly
166 mutated (**Figure 2 A**). In particular, there is a high frequency of mutations between
167 residues 172-189, 232-237, 277-278, and 369-381 that correspond to mutations
168 likely to alter the binding site via proximity by our definition. Additionally, generating
169 the Rhapsody scores¹⁹ for each residue results in regions of predicted pathogenicity
170 that also align with the binding site regions – reinforcing that mutations neighbouring
171 binding sites are likely to be pathogenic purely via proximity and disruption of the
172 precise conformation of sidechains necessary for catalysis. Whilst Rhapsody
173 represents a potentially useful single metric for assessing mutational disruption,
174 incorporates evidence from sequence and structure alone, without biological context.

175 Due to the three-domain structure of FH we surmised that regions involved in the
176 “hinging” of these domains may influence the binding site assembly, due to the
177 proximity and reliance of the quaternary structure of multiple domains to make up the
178 binding pocket. To calculate predicted hinges within the structure we used Gaussian
179 Network Modelling (GNM) within prody^{26,27} to calculate the major normal modes for
180 an individual subunit of FH. We find that the second mode best represents the
181 hinging mode expected around the three domains of the protein. Calculating the
182 hinge residues from the second normal mode results in residues 196, 198, 232, 242,
183 270, 317, 401, 411, and 448 being the most likely “hinge points” in the structure

184 **(Figure 2 B)**, these residues are shown on a single subunit of FH, coloured by
185 eigenvector direction in **Figure 2 C**. In order to assess whether mutation of these
186 domains was sufficient to disrupt the quaternary structure of the protein and thus the
187 binding site we chose to simulate a known mutation within a hinge region that is
188 found at a high frequency in the FH mutation database using molecular dynamics
189 simulations. We chose to simulate the R233H mutant, and the wild type (WT)
190 tetrameric assemblies for 200ns each. Measuring the angles between CA atoms of
191 two residues in the centre of the D2 and D3 regions with respect to the hinge reveals
192 that the R233H mutant reduces the angle of the domains by an average of 8
193 degrees, and so leads to a partial occlusion of the catalytic site of FH **(Figure 2 D)**.
194 From this evidence we conclude that disruption of these hinges are likely to alter the
195 binding site and assembly of FH – and are likely pathogenic. We chose to treat all
196 mutations with CA atoms within 6 Å of any hinge residue as potentially LOF through
197 disruption of the protein quaternary structure.

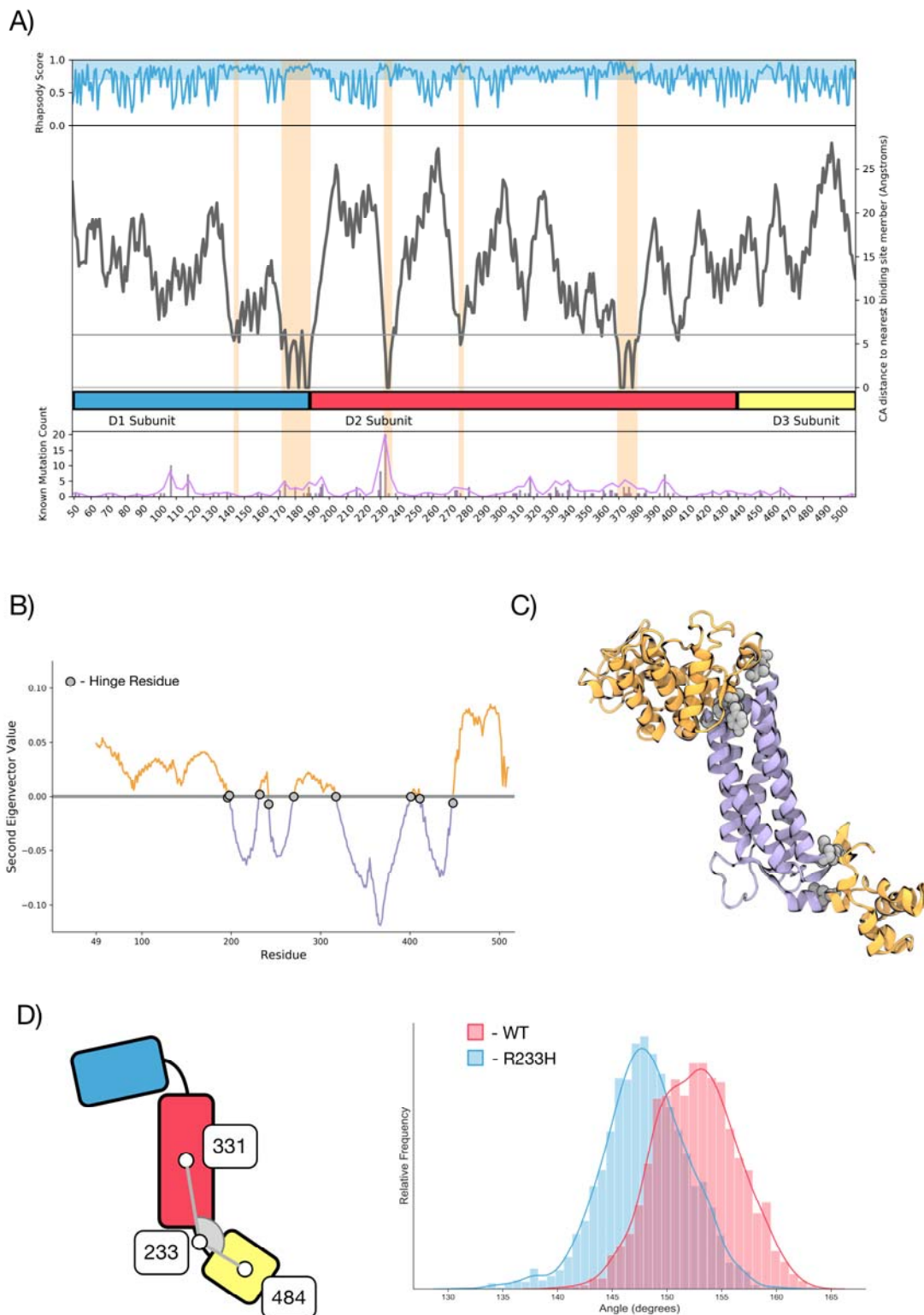
198 Overall, we find that mutations near to either the binding site, or a hinge region of the
199 protein are likely to disrupt or alter the protein function. We find that, from the FH
200 mutation database, a significant proportion of mutations can be classed as either
201 binding site-associated, or hinge-associated, including a number of known loss-of-
202 function (LOF) variants. Whilst 42 residues in the 461 amino acid protein structure
203 (9%) are classified as being “binding site-associated”, we find that 11 of the 30 (36%)
204 known LOF mutations are within these residues, showing a clear bias towards
205 binding site-associated mutations. Similarly, 55 of the 461 (12%) amino acids in the
206 protein structure are classified as “hinge-associated”, and we find 7 of the 30 (23%)
207 within the FH mutation database fulfil this classification, showing a lesser, but still
208 large occurrence bias. Distance calculations for all potential mutations are included
209 in **Supplementary Table 1**.

210

211

212

213



214

215 **Figure 2:** Mutations can be categorised on proximity to functional regions of FH. A)

216 Alpha carbon (CA) distance from a binding site residue. Shown is: Top: average

217 Rhapsody score for each residue, Middle: distance of each residue from a binding
218 site residue by CA distance, Bottom: mutational frequency for each residue. Orange
219 highlights show some regions have high Rhapsody scores, low binding site distance,
220 and high mutational frequency. B) Second normal mode eigenvectors per residue for
221 a single subunit of FH. Residues with an eigenvector above the line are moving
222 generally opposed to those with an eigenvector below the line. Predicted hinge
223 residues are shown in grey. C) Single subunit of FH coloured according to
224 eigenvector direction (positive as orange and negative as purple). Hinge residues
225 are highlighted as grey. D) Molecular Dynamics simulations of hinge mutations
226 shows altered hinge flexibility. Left: Schematic of the angle measured in each
227 simulation, Right: Angle of WT (red), and R233H mutant FH (blue) over a 200 ns
228 equilibrium molecular dynamics simulation.

229

230 **High-Throughput mutational stability screen of FH *in silico***

231 To study how mutations that are not near the binding site or hinge regions may have
232 effects on the structure of the protein, we sought to generate predicted mutational
233 energy changes ($\Delta\Delta G$) for every potential amino acid substitution in the FH
234 structure. We chose to use two conceptually different methods and use an average
235 between the two methods to study each potential mutant. We chose to utilize the
236 FoldX method^{10,11}, and the Rosetta cartesian_ddg method^{12,13}, (hereafter described
237 as the Rosetta method) to perform mutational energy calculations. Both methods
238 have been shown to generate accurate predictions on the CAGI5 blind challenge
239 datasets, but overlaps between the two methods on the same dataset indicate that
240 they generally predict different mutations correctly, the combined overlap between
241 the two being a good indicator of mutational $\Delta\Delta G$ ¹³.

242 To perform mutant calculations, the pdb structure 5upp was first relaxed using the
243 FoldX RelaxPDB method, before each mutation and its resultant $\Delta\Delta G$ was
244 calculated. We additionally calculate the Relative Solvent Accessible Surface Area
245 (RSA) for each wild-type (WT) residue. Mutations on the surface of the protein are
246 unlikely to dramatically alter the folding of the protein, so we chose to only consider a
247 mutation potentially destabilizing if it is buried, defined as having an RSA ≤ 0.2 .

248

249 We find a good agreement between the FoldX and Rosetta methods, with an r of
250 0.67 ($p < 0.0001$) for all mutational energies (**Figure 3 A**). Notably however, both
251 methods appear to agree on predictions of mutations with extremely high energy, but
252 there is a significant portion of the distribution that shows a reasonably poor
253 correlation, particularly mutations that have a predicted $\Delta\Delta G$ between 1 and -1
254 kcal/mol. We chose to study the average predicted energy of each mutation by
255 taking the average $\Delta\Delta G$ from the two methods. Ranking the average $\Delta\Delta G$ over all
256 ~9000 mutations results in a distribution of all mutational energies across the
257 mutational landscape (**Figure 3 B**). We find that mutations known to be LOF, and
258 that are not within 6 Å of either the binding site or hinge regions tend to cluster near
259 the upper end of the distribution, indicating that they affect the stability of the protein.
260 Mutations that are known benign tend to fall near the lower end of the distribution.
261 We chose a cutoff of 2.5 kcal/mol for classifying mutations as potentially
262 destabilizing, and any mutation over this value for average energy, and with a RSA $<$
263 0.2 was classified as destabilizing. Across all potential mutations we find that ~45%
264 (3968 out of 8778) meet this criterion (**Figure 3 C**). This fits roughly with historical
265 data of mutational stability in T4 lysozyme, which found that 45% of mutational sites
266 lead to structural inactivation of enzymatic function²⁸.

267

268

269

270

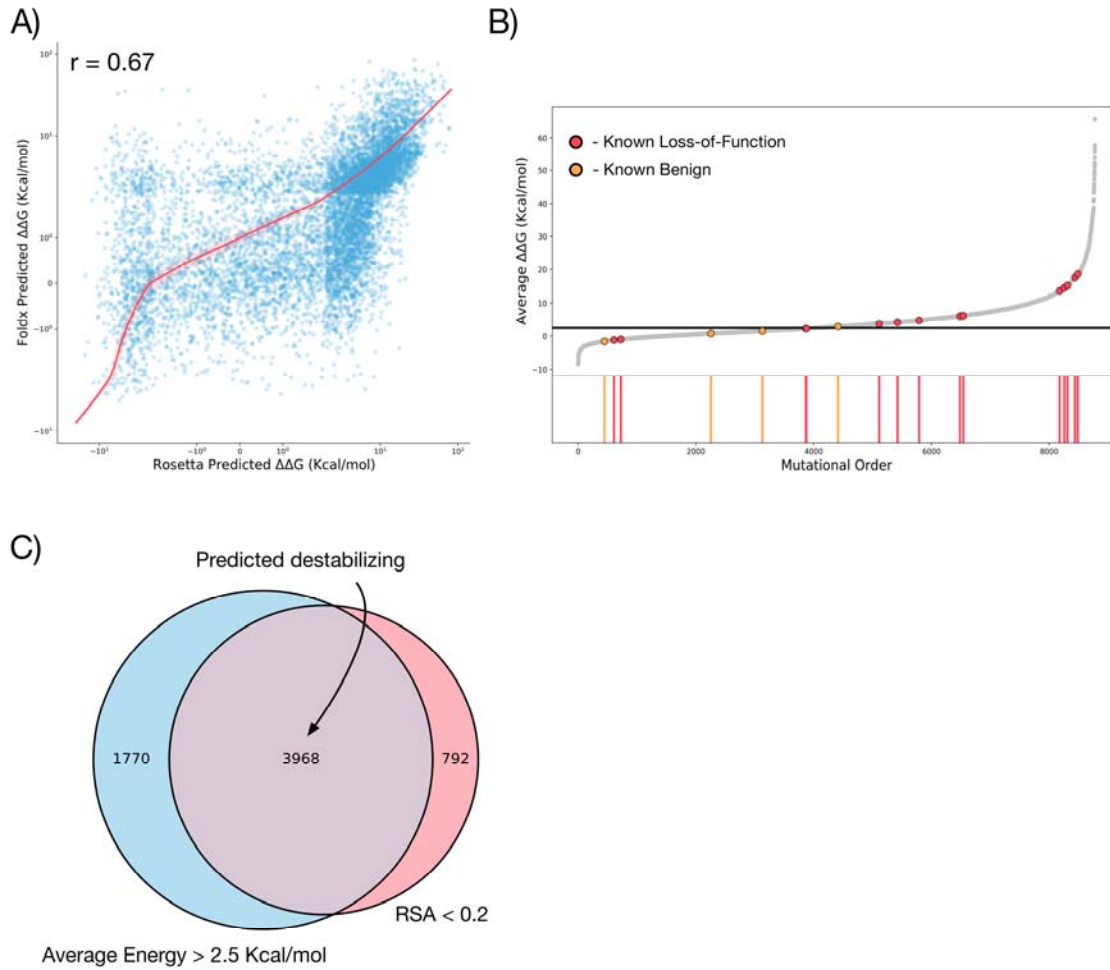
271

272

273

274

275



276

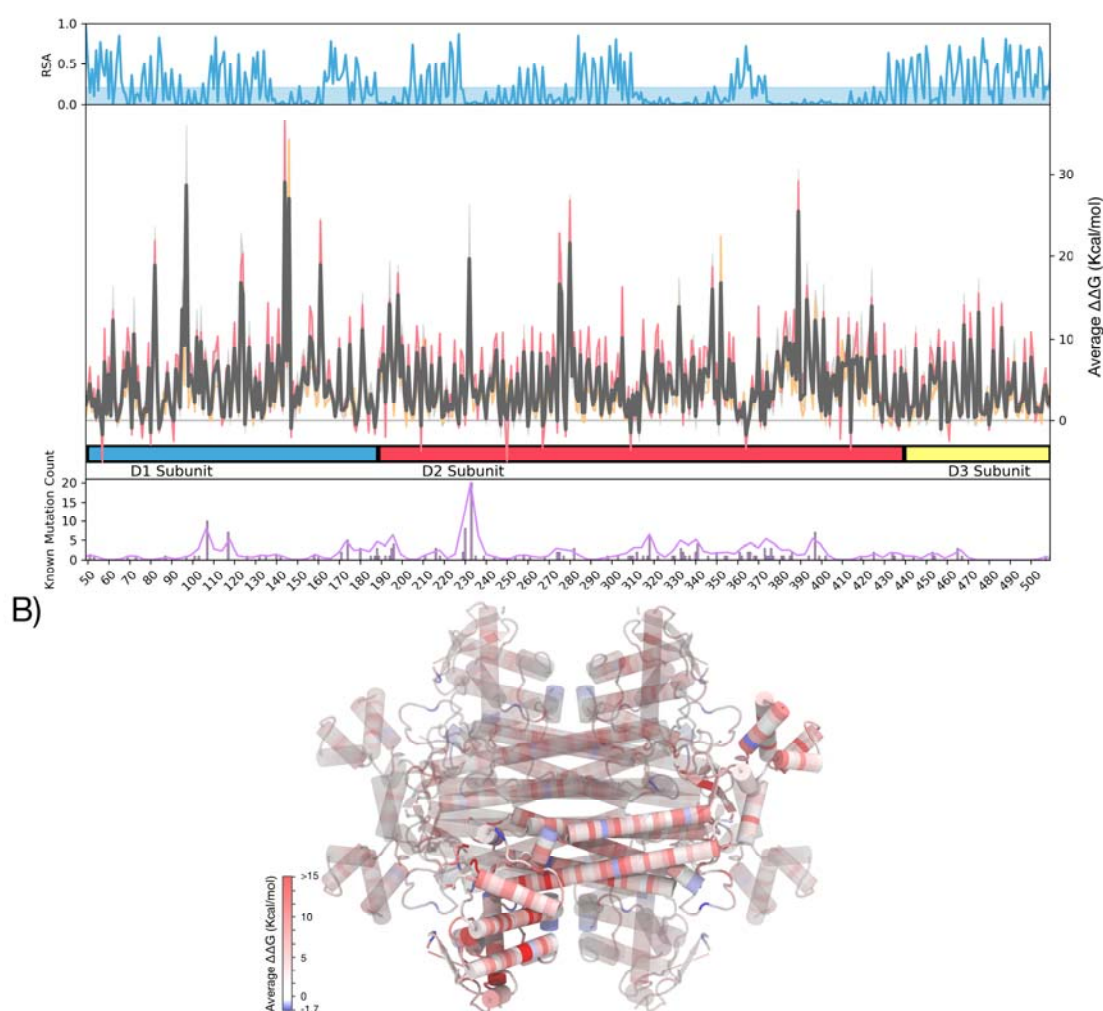
277 **Figure 3:** Prediction of Destabilizing Mutations. A) Comparison of $\Delta\Delta G$ calculations
278 from FoldX and Rosetta. Correlation r is spearman's rank. B) Position of known loss-
279 of-function (red) and benign mutations (orange) on the $\Delta\Delta G$ spectrum. Mutations are
280 ordered in ascending mutational $\Delta\Delta G$. Black line represents 2.5 Kcal/mol cutoff. C)
281 Overlap between residues with a high predicted mutational energy (Defined as those
282 with average $\Delta\Delta G > 2.5$ Kcal/mol) and buried residues (RSA < 0.2). In total 3968
283 mutations are classified as destabilizing by taking the overlap between these two
284 criterion.

285

286

287

288 Plotting mutational frequency for both methods, and their average for each residue
289 **(Figure 4 A)** reveals that the most destabilizing mutations predicted by either
290 method are in regions with a large number of buried amino acids, as expected.
291 When plotting these mutations on the structure of the protein **(Figure 4 B)**, we find
292 the most significantly destabilizing mutations are those packed within the centre of
293 D1, and on the interface between D1 and D2. This location suggests mutational
294 disruption will alter the position of the D1/D2 interface, and thus will affect the binding
295 site conformation, whereas mutations within the core D2 region are likely to influence
296 the stability of the fully assembled tetramer.



297

298 **Figure 4:** Predicted $\Delta\Delta G$ for every mutation in Fumarate Hydratase. A) Average
299 mutational energy per residue in the FH structure. Top: Relative solvent accessible
300 surface area (RSA) for every residue. Blue highlight indicates RSA < 0.2, classified

301 as buried. Middle: average $\Delta\Delta G$ for each residue (Grey). Red and Orange lines
302 represent average Rosetta and FoldX calculations respectively. Bottom: Mutational
303 frequency from the FH mutation database for each residue in FH. B) Mutational $\Delta\Delta G$
304 applied to the structure 5upp of FH. Red indicates high average $\Delta\Delta G$, and so
305 represents areas where mutations are likely to disrupt the structure. Blue represents
306 regions of generally stabilizing mutations.

307

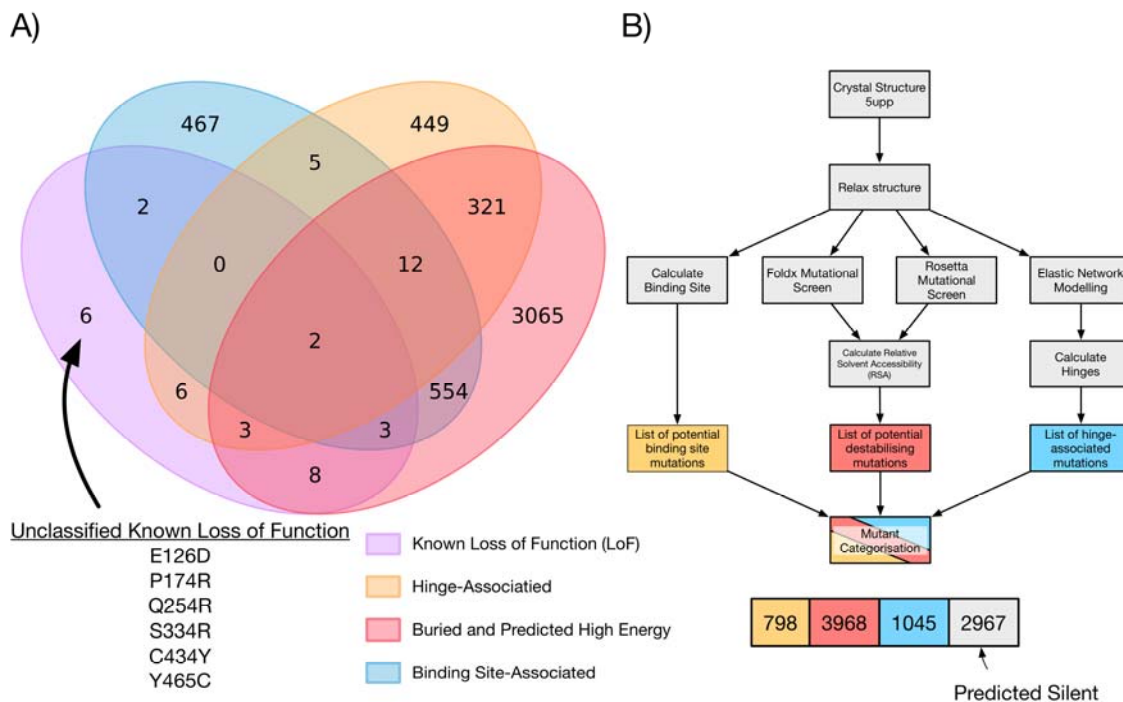
308 **Existing mutations are accurately categorised based on known phenotypic** 309 **effects**

310 Overall, we define a scheme for classifying mutations into different categories of
311 potentially disrupting, predicted LOF substitutions (**Figure 5 A**). The initial structure
312 is relaxed using FoldX, before the binding site and hinge regions are calculated and
313 classified, additionally mutations that are potentially destabilizing are defined based
314 on average energy from the Rosetta and FoldX mutation methods, plus screened for
315 buried mutations through calculating the RSA. This results in a categorisation for
316 every mutation, where each is classified as predicted silent, binding site, hinge site,
317 or destabilizing (including combinations of disruptive mutation types)
318 (**Supplementary Table 1**). Overall we classify 5811 out of 8778 (66%) mutations as
319 potentially functionally disruptive, similar to a study of mutational effects on TP53,
320 which found that roughly 50-60% of all possible mutations were functionally
321 disruptive²⁹.

322

323 To study the accuracy of our classification we chose to interrogate all mutations
324 within the FH mutation database. We classed all mutations within the database as
325 either loss of function (LOF), benign, or of unknown functional effect. In total 34
326 mutations had a known (or implied) functional effect, whilst 73 were classified as
327 unknown (**Supplementary Table 2**). We sought to validate our functional
328 classifications (binding site associated, hinge associated, or destabilizing) against
329 the mutations that are known to be LOF (**Figure 5 B**). We find that 24 out of 30
330 (80%) mutations are correctly classified as LOF, and 3 out of 4 (75%) are correctly
331 classed as benign. Of the mutations incorrectly classified as benign when they are

332 known to be LOF, two mutations involve cysteine (C434Y, Y465C), which is known
 333 to be modelled poorly by Rosetta cartesian_ddg, and in fact, is classified as a
 334 stabilizing mutations by Rosetta, with a predicted $\Delta\Delta G$ of -5.2 kcal/mol, though FoldX
 335 classifies it as destabilizing. The mutation incorrectly classified as deleterious when it
 336 is listed as benign within the FH mutation database is R268G. We predict the R268G
 337 mutation to be both destabilizing ($\Delta\Delta G > 2.5$ kcal/mol, RSA < 0.2) and hinge-
 338 associated. Whilst the mutation is listed as benign, no experimental information is
 339 cited, and PolyPhen-2³⁰, and Rhapsody also classify this particular mutation as
 340 damaging, indicating that the benign classification for this mutation may be
 341 questionable. To further explore this mutation we ran a molecular dynamics
 342 simulation of the R268G mutant. Simulations predict that mutant R268G reduces the
 343 hinge angle of the D1/D2 domains by ~5 degrees (**Supplementary Figure 1**), and
 344 supports previous evidence from the R233H mutant, that hinges within the protein
 345 can effect binding site assembly. Of the 73 unknown mutations, we predict that 28
 346 are functionally benign, and 45 are potential LOF mutations.



347

348 **Figure 5:** Prediction of known Loss-of-Function (LOF) mutations. Venn diagram
 349 showing overlap between Hinge-associated (orange), destabilizing (red) and Binding
 350 site-associated (blue) mutations. 30 known LOF mutations are included (purple) 24

351 mutations are correctly categorized as LOF, whilst 6 are incorrectly categorized as
352 benign mutations. B) Schema for categorization of mutations in FH. The structure is
353 initially relaxed using FoldX RelaxPDB, residues within 6 Å of the binding site are
354 calculated resulting in a list of 798 binding site associated mutations (orange). FoldX
355 and Rosetta are used to calculate the $\Delta\Delta G$ for every mutation and this is subset by
356 the relative solvent accessible surface area resulting in 3968 potentially destabilizing
357 mutations (red). Elastic network modelling is performed to generate hinge regions of
358 the protein, and residues within 6 Å of hinges are calculated, resulting in 1045 hinge
359 associated mutations (blue). 2967 mutations are predicted to be silent.

360

361 **Mutations with unknown properties can be accurately predicted to be**
362 **functional or neutral**

363 To study all potential mutations in FH we chose to plot all mutations using umap³¹.
364 We ran umap on the 4 major axis involved in the classification in this study for every
365 mutations – minimum distance to a binding site residue, minimum distance to a
366 hinge residue, average $\Delta\Delta G$ of mutation, and RSA for each residue (**Figure 6 A**).
367 We find that distinct regions of the plot cluster into functionally different mutations
368 when coloured by classification. There is a region specifically for hinge-associated
369 mutations, binding site-associated, and unknown (not predicted damaging)
370 mutations. In particular, the region of “unknown” (not classified as damaging)
371 mutations overlaps significantly with a number of predicted destabilizing mutations,
372 indicating that discrimination between these mutations is difficult, and perhaps not
373 accurate with currently available data. Also shown are the known benign and loss of
374 function mutations. We find that most of the benign mutations, aside from R268G are
375 found clearly within the regions of predicted benign mutations. R268G clusters with
376 the hinge mutation region as expected from our previous classification. For the
377 known LOF mutations, we find they mostly cluster within the well defined regions for
378 binding site, hinge, and destabilizing mutations. There are some mutations,
379 particularly those which were misclassified, that fall within ambiguous regions of
380 state space in the mutational landscape, and so are hard to classify using our
381 defined criterion. Overall, we find that mutations broadly separate as expected using
382 umap, and that unclassified mutations can be plotted on the resultant distribution –

383 confidence of the classification of any individual mutation can be inferred from where
384 it fits within the landscape.

385

386 To test the predictive power of our classification we used the Cancer Cell Line
387 Encyclopedia to look for changes in metabolite levels associated with mutations in
388 FH^{21,32}. We find 42 mutations (35 unique) in FH within 34 individual cell lines
389 (**Supplementary Table 3**). Selecting only for missense mutations yielded 25
390 mutations (20 unique) within 23 unique cell lines. We classified the mutations
391 according to our criterion as either predicted LOF, or predicted benign. We find that
392 by analysis of metabolomics data included in the CCLE database, mutations that we
393 predict to be LOF have a higher average level of fumarate/mateate/alpha-
394 ketoisovalerate detected in media than cells with predicted benign mutations ($p =$
395 0.035) – indicating that these cell lines may have an accumulation of fumarate as a
396 result of inactive levels of FH (**Figure 6 B,C**).

397

398

399

400

401

402

403

404

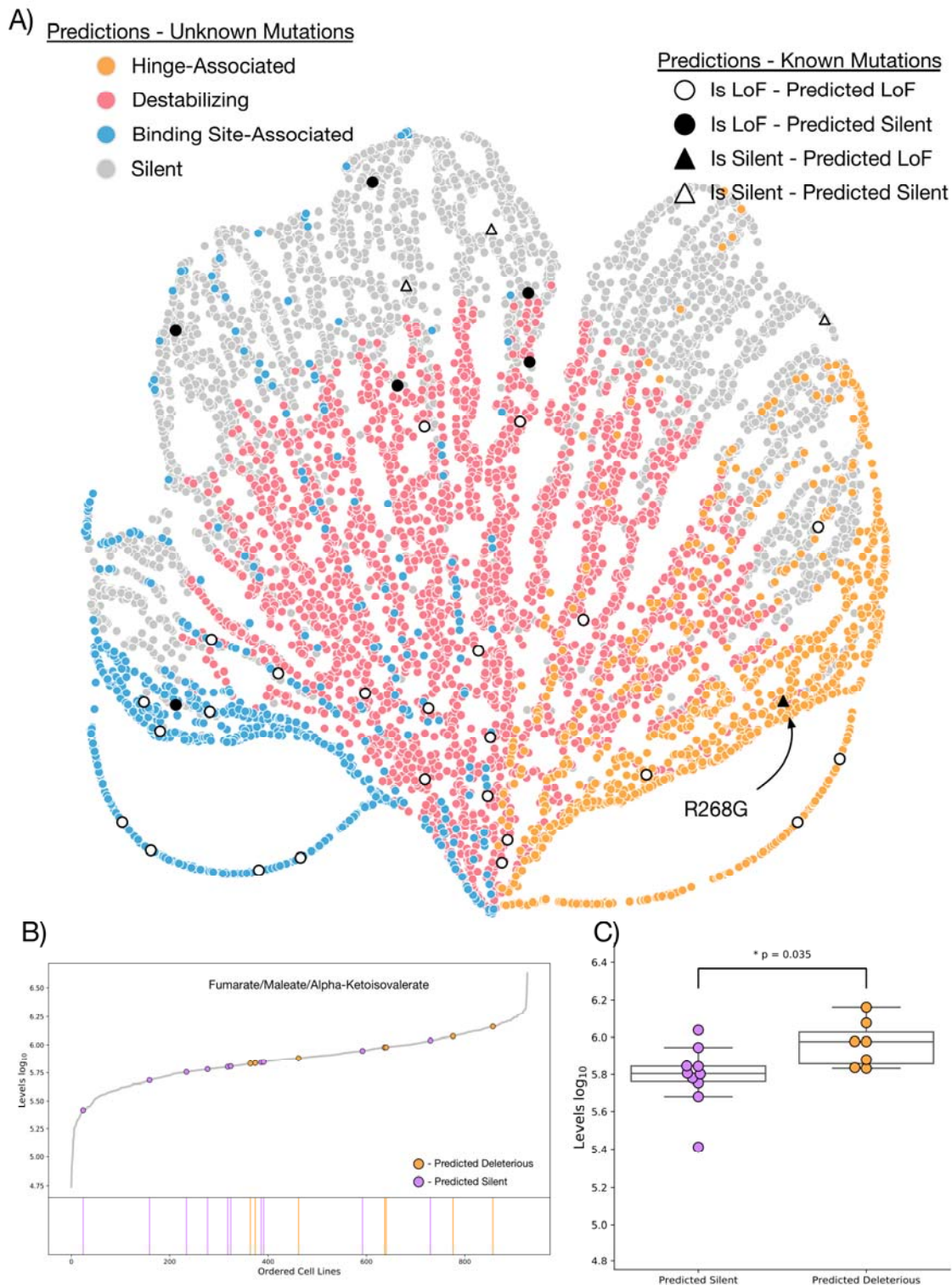
405

406

407

408

409



410

411 **Figure 6:** Mutational Landscape of Fumarate Hydratase. A) Umap for all mutations
 412 in FH. Mutations are coloured by classification. Hinge-associated (orange),
 413 Destabilizing (red), and Binding site-associated (blue) are shown clustered into
 414 groups. Predicted silent mutations (grey) are also shown. Overlaid are our

415 predictions for characterized mutations in the FH mutation database. Mutations that
416 are known Loss-of-Function (LOF) are circular and coloured according to whether
417 we predict them to be LOF (black) or silent (white). Known benign mutations are in
418 triangles, and also coloured according to whether we predict them to be LOF (black)
419 or silent (white). The questionable known benign mutation R268G is labelled B)
420 Mutations in the Cancer Cell Line Encyclopedia (CCLE) metabolomics data. All cell
421 lines are ranked according to their detected levels of Fumarate/Maleate/Alpha-
422 Ketoisovalerate. Coloured are cell lines with mutations in FH that we predict to be
423 LOF (orange), or silent (purple). C) Swarmplot for levels of Fumarate/Maleate/Alpha-
424 Ketoisovalerate in mutant FH cell lines. Mutations predicted to be silent are
425 significantly lower than mutations predicted to be LOF (p value represents
426 independent T test).

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445 **DISCUSSION**

446 In conclusion, we have shown, using a comprehensive combination of techniques,
447 that we can categorise accurately the functional effects of any potential missense
448 mutation in FH. Beyond FH, we present an integrated series of methods that can be
449 adapted for mutationally screening any protein for functionally relevant mutations in a
450 reasonably small amount of computational time. Our workflow predicts the functional
451 effects of all mutations that can be compared to existing methods based on machine-
452 learning principles such as Rhapsody and PolyPhen, at significantly lower time and
453 effort expenditure than experimental characterization. Whilst some other methods
454 incorporate some manner of structural analysis in their predictions, ours
455 demonstrates a new perspective, as it explicitly models every potential mutation in a
456 structure, allowing it to interface directly with other computational techniques in the
457 field such as molecular dynamics simulations to further study mutations of interest.

458

459 Biologically we propose three ways in which mutations can potentially disrupt the
460 catalytic activity of FH. In particular we find that addition of hinge altering mutations
461 are necessary for classification of many known LoF mutations, indicating that there is
462 a biological relevance, and hinting at a mechanism for, mutations that change the
463 flexibility and stiffness of protein hinges in this case. Additionally, we chose to
464 exclude site B from our analysis of mutation disruption and find that we are able to
465 classify almost all known mutations without its inclusion. This implies that mutations
466 in site B may not have functional or disease-related relevance, despite some
467 evidence that site B can alter catalytic activity of the enzyme³³. This is reinforced by
468 the fact that 27 of the 461 residues within the protein structure are classified as near
469 site B (6%), and only 3 of 30 residues in the FH mutation database (10%) are near to
470 site B, showing a poor to negligible enrichment of mutations in site B when
471 compared to similar calculations for site A.

472

473 Fumarate hydratase represents a good first-use case for high-throughput mutational
474 screen due to the need to understand mutations in their functional context, but as

475 mutational detection techniques, and high-throughput mutational studies increase
476 the need to be able to classify mutations confidently as benign and LOF is more
477 important. Here we show that our method accurately classifies known LOF and
478 benign mutations with a high degree of accuracy, and predict which mutations
479 discovered in the human population are likely to have functional relevance, and
480 therefore predispose patients to particular metabolic diseases.

481

482 Whilst the accuracy of our method with the current data is high, there are clear
483 regions where the analysis is not able to discriminate between mutations on the
484 borderline between destabilizing and benign, this results from the lack of accuracy in
485 the mutational $\Delta\Delta G$ calculations, despite using the best available methods at time of
486 study¹³. As better methods become available it will be of interest to improve upon
487 this work to attempt a more accurate classification.

488 Finally, whilst the work here focusses on a single molecule within the TCA cycle, FH,
489 structural data has existed for a large number of enzymes within the cycle for some
490 time³⁴⁻³⁷, and it would be of great interest to look into mutations across entire
491 metabolic pathways. With this study laying the groundwork, it will be of future interest
492 to model all mutations in all enzymes, and attempt to further link these with genomic
493 and metabolomic data that is already available.

494

495

496

497

498

499

500

501 **AUTHOR CONTRIBUTIONS**

502 DS and BAH conceived the study and wrote the manuscript. DS generated all data
503 and performed all analysis. All authors were responsible for editing of the
504 manuscript.

505

506 **ACKNOWLEDGEMENTS**

507 We thank the Frezza group, in particular Christian Frezza for support and
508 constructive feedback during the generation of this manuscript.

509

510 **COMPETING FINANCIAL INTEREST**

511 The authors declare no competing financial interest.

512

513 **DATA AVAILABILITY**

514 All data used in this study, including the code used in generating all figures from raw
515 data is available publicly at: https://github.com/shorthouse-mrc/Fumarate_Hydratase

516

517

518

519

520

521

522

523 **METHODS**

524 **FH mutation database**

525 The FH mutation database was downloaded from the Leiden Open Variation
526 Database⁹ (<https://databases.lovd.nl/shared/variants/FH/unique>). Missense
527 mutations were manually curated into categories (Loss of Function, Benign, and
528 Unknown) based on their implied clinical classification, and variant remarks, which
529 contained information regarding FH enzymatic activity.

530 **Mutational Clustering**

531 Mutational clustering was performed with the NMC clustering algorithm, which
532 attempts to discern the likelihood of a mutation spectrum occurring by random
533 chance. NMC returns clusters of mutations that are statistically significant. We chose
534 to run the NMC algorithm using the R library iPAC²⁵, using an alpha cutoff value of
535 0.05, and the Bonferroni multiple test correction method (see supplementary code).

536 **Gaussian Network Modelling**

537 GNM was implemented using the Prody package in python³⁸.

538 **Molecular Dynamics Simulations**

539 Molecular dynamics was performed using Gromacs version 2018.1³⁹. We chose to
540 simulate proteins using the GROMOS 54a7 forcefield⁴⁰.

541 The protein structure was first repaired using FoldX¹⁰ “RepairPDB” with the following
542 command:

```
$foldx --command=RepairPDB --pdb=5upp.pdb --ionStrength=0.05 --pH=7 --  
vdwDesign=2
```

543 The protein was then placed in a cubic box and solvated with spc water. Counterions
544 were introduced to a neutral charge, and to a concentration of 0.05 mol/litre. The
545 system was energy minimized using the steepest descents algorithm until the
546 maximum force, F_{\max} , of the system reached below 1000 kJ/mol/nm.
547

548 Equilibration was performed using the NVT, followed by the NPT ensembles for 100
549 ps each. We chose to use the verlet cutoff scheme and PME electrostatics, and
550 utilized periodic boundary conditions in the x,y, and z planes.

551 Molecular dynamics was performed for 200 ns retaining velocities from the NPT
552 equilibration. We used the V-rescale temperature coupling scheme, and Parrinello-
553 Rahman isotropic pressure coupling.

554 **FoldX $\Delta\Delta G$ Calculations**

555 FoldX predicted $\Delta\Delta G$ was calculated using the PositionScan command within
556 FoldX4. Positionscan was run on each residue in the protein structure sequentially
557 using the following command:

```
$foldx --command=PositionScan --pdb=5upp.pdb --ionStrength=0.05 --pH=7 --  
vdwDesign=2 --pdbHydrogens=false --positions=49
```

558
559

560 For positionscan on the 49th residue.

561 **Rosetta $\Delta\Delta G$ Calculations**

562 Rosetta predicted $\Delta\Delta G$ was calculated using the cartesian_ddg method as described
563 in Kellogg et al:

```
$path/to/source/bin/cartesian_ddg.static.linuxgccrelease -in:file:s 5upp.pdb -  
in:file::fullatom -database /path/to/database/ -ignore_unrecognized_res true -  
ignore_zero_occupancy false -fa_max_dis 9.0 -ddgccartesian -ddg::mut_file  
mutfile.txt -ddg::iterations 3 -ddg::dump_pdbs true -ddg::suppress_checkpointing  
true -ddg::mean true -ddg::min true -ddg::output_silent true -bbnbr 1 -  
beta_nov16_cart > logfile.log
```

564
565

566 $\Delta\Delta G$ was calculated by averaging the energy of 3 models of each mutation and
567 comparing it to the WT calculation.

568 **Umap**

569 We used Umap³¹ based on the github repository at
570 www.github.com/lmcinnes/unmap

571 **Cancer Cell Line Encyclopedia Data**

572 Cancer Cell Line Encyclopedia (CCLE) mutation data was downloaded from the
573 Broad Institute at: <https://portals.broadinstitute.org/ccle/data> . Metabolomics data
574 was obtained from the supplementary data of Li et al²¹ .

575 **Data Analysis**

576 Both MDanalysis⁴¹ and Biopython⁴² were used for analysis of structural data.

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592 **BIBLIOGRAPHY**

- 593 1. Ewbank, C., Kerrigan, J. F. & Aleck, K. *Fumarate Hydratase Deficiency*.
594 *GeneReviews*® (University of Washington, Seattle, 1993).
- 595 2. Yang, M., Soga, T., Pollard, P. J. & Adam, J. The emerging role of fumarate as
596 an oncometabolite. *Front. Oncol.* **2**, 85 (2012).
- 597 3. Leshets, M., Ramamurthy, D., Lisby, M., Lehming, N. & Pines, O. Fumarase is
598 involved in DNA double-strand break resection through a functional interaction
599 with Sae2. *Curr. Genet.* **64**, 697–712 (2018).
- 600 4. Yogev, O. *et al.* Fumarase: A Mitochondrial Metabolic Enzyme and a
601 Cytosolic/Nuclear Component of the DNA Damage Response. *PLoS Biol.* **8**,
602 e1000328 (2010).
- 603 5. Schmidt, C., Sciacovelli, M. & Frezza, C. Fumarate hydratase in cancer: A
604 multifaceted tumour suppressor. *Semin. Cell Dev. Biol.* (2019).
- 605 6. Skala, S. L., Dhanasekaran, S. M. & Mehra, R. Hereditary Leiomyomatosis
606 and Renal Cell Carcinoma Syndrome (HLRCC): A Contemporary Review and
607 Practical Discussion of the Differential Diagnosis for HLRCC-Associated Renal
608 Cell Carcinoma. *Arch. Pathol. Lab. Med.* **142**, 1202–1215 (2018).
- 609 7. Alam, N. A., Olpin, S. & Leigh, I. M. Fumarate hydratase mutations and
610 predisposition to cutaneous leiomyomas, uterine leiomyomas and renal
611 cancer. *Br. J. Dermatol.* **153**, 11–17 (2005).
- 612 8. Clark, G. R. *et al.* Germline *FH* Mutations Presenting With
613 Pheochromocytoma. *J. Clin. Endocrinol. Metab.* **99**, E2046–E2050 (2014).
- 614 9. Bayley, J.-P., Launonen, V. & Tomlinson, I. P. The FH mutation database: an
615 online database of fumarate hydratase mutations involved in the MCUL
616 (HLRCC) tumor syndrome and congenital fumarase deficiency. *BMC Med.*
617 *Genet.* **9**, 20 (2008).
- 618 10. Delgado, J., Radusky, L. G., Cianferoni, D. & Serrano, L. FoldX 5.0: working
619 with RNA, small molecules and a new graphical interface. *Bioinformatics* **35**,

- 620 4168–4169 (2019).
- 621 11. Schymkowitz, J. *et al.* The FoldX web server: an online force field. *Nucleic*
622 *Acids Res.* **33**, W382–W388 (2005).
- 623 12. Alford, R. F. *et al.* The Rosetta All-Atom Energy Function for Macromolecular
624 Modeling and Design. *J. Chem. Theory Comput.* **13**, 3031–3048 (2017).
- 625 13. Strokach, A., Corbi-Verge, C. & Kim, P. M. Predicting changes in protein
626 stability caused by mutation using sequence- and structure-based methods in a
627 CAGI5 blind challenge. *Hum. Mutat.* **40**, 1414–1423 (2019).
- 628 14. Tomczak, K., Czerwińska, P. & Wiznerowicz, M. The Cancer Genome Atlas
629 (TCGA): an immeasurable source of knowledge. *Contemp. Oncol. (Poznan,*
630 *Poland)* **19**, A68-77 (2015).
- 631 15. International Cancer Genome Consortium, T. I. C. G. *et al.* International
632 network of cancer genome projects. *Nature* **464**, 993–8 (2010).
- 633 16. Martincorena, I. *et al.* Universal Patterns of Selection in Cancer and Somatic
634 Tissues. *Cell* **171**, 1029-1041.e21 (2017).
- 635 17. Lawrence, M. S. *et al.* Mutational heterogeneity in cancer and the search for
636 new cancer-associated genes. *Nature* **499**, 214–218 (2013).
- 637 18. Tamborero, D., Lopez-Bigas, N. & Gonzalez-Perez, A. Oncodrive-CIS: A
638 Method to Reveal Likely Driver Genes Based on the Impact of Their Copy
639 Number Changes on Expression. *PLoS One* **8**, e55489 (2013).
- 640 19. Ponzoni, L., Oltvai, Z. N. & Bahar, I. Rhapsody: Pathogenicity prediction of
641 human missense variants based on protein sequence, structure and dynamics.
642 *bioRxiv* 737429 (2019).
- 643 20. Ghandi, M. *et al.* Next-generation characterization of the Cancer Cell Line
644 Encyclopedia. *Nature* **569**, 503–508 (2019).
- 645 21. Li, H. *et al.* The landscape of cancer cell line metabolism. *Nat. Med.* **25**, 850–
646 860 (2019).

- 647 22. Picaud, S. *et al.* Structural basis of fumarate hydratase deficiency. *J. Inherit.*
648 *Metab. Dis.* **34**, 671–6 (2011).
- 649 23. Rose, I. A. & Weaver, T. M. The role of the allosteric B site in the fumarase
650 reaction. *Proc. Natl. Acad. Sci. U. S. A.* **101**, 3393 (2004).
- 651 24. Ajalla Aleixo, M. A., Rangel, V. L., Rustiguel, J. K., de Pádua, R. A. P. &
652 Nonato, M. C. Structural, biochemical and biophysical characterization of
653 recombinant human fumarate hydratase. *FEBS J.* **286**, 1925–1940 (2019).
- 654 25. Ye, J., Pavlicek, A., Lunney, E. A., Rejto, P. A. & Teng, C. H. Statistical
655 method on nonrandom clustering with application to somatic mutations in
656 cancer. *BMC Bioinformatics* (2010).
- 657 26. Haliloglu, T., Bahar, I. & Erman, B. Gaussian Dynamics of Folded Proteins.
658 *Phys. Rev. Lett.* **79**, 3090–3093 (1997).
- 659 27. Bahar, I., Atilgan, A. R. & Erman, B. Direct evaluation of thermal fluctuations in
660 proteins using a single-parameter harmonic potential. *Fold. Des.* **2**, 173–181
661 (1997).
- 662 28. Rennell, D., Bouvier, S. E., Hardy, L. W. & Poteete, A. R. Systematic mutation
663 of bacteriophage T4 lysozyme. *J. Mol. Biol.* **222**, 67–88 (1991).
- 664 29. Kotler, E. *et al.* A Systematic p53 Mutation Library Links Differential Functional
665 Impact to Cancer Mutation Pattern and Evolutionary Conservation. *Mol. Cell*
666 **71**, 178-190.e8 (2018).
- 667 30. Adzhubei, I. A. *et al.* A method and server for predicting damaging missense
668 mutations. *Nat. Methods* **7**, 248–249 (2010).
- 669 31. McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation
670 and Projection for Dimension Reduction. (2018).
- 671 32. Barretina, J. *et al.* The Cancer Cell Line Encyclopedia enables predictive
672 modelling of anticancer drug sensitivity. *Nature* **483**, 603–607 (2012).
- 673 33. Alam, N. A. *et al.* Missense Mutations in Fumarate Hydratase in Multiple

- 674 Cutaneous and Uterine Leiomyomatosis and Renal Cell Cancer. *J. Mol.*
675 *Diagnostics* **7**, 437–443 (2005).
- 676 34. Remington, S., Wiegand, G. & Huber, R. Crystallographic refinement and
677 atomic models of two different forms of citrate synthase at 2.7 and 1.7 Å
678 resolution. *J. Mol. Biol.* **158**, 111–152 (1982).
- 679 35. Lauble, H., Kennedy, M. C., Beinert, H. & Stout, C. D. Crystal structures of
680 aconitase with isocitrate and nitroisocitrate bound. *Biochemistry* **31**, 2735–
681 2748 (1992).
- 682 36. Chapman, A. D. ., Cortés, A., Dafforn, T. ., Clarke, A. . & Brady, R. . Structural
683 basis of substrate specificity in malate dehydrogenases: crystal structure of a
684 ternary complex of porcine cytoplasmic malate dehydrogenase, α -
685 Ketomalonate and TetrahydroNAD 1 1 Edited by R. Huber. *J. Mol. Biol.* **285**,
686 703–712 (1999).
- 687 37. Yankovskaya, V. *et al.* Architecture of Succinate Dehydrogenase and Reactive
688 Oxygen Species Generation. *Science (80-.)*. **299**, 700–704 (2003).
- 689 38. Bakan, A., Meireles, L. M. & Bahar, I. ProDy: Protein Dynamics Inferred from
690 Theory and Experiments. *Bioinformatics* **27**, 1575–1577 (2011).
- 691 39. Abraham, M. J. *et al.* GROMACS: High performance molecular simulations
692 through multi-level parallelism from laptops to supercomputers. *SoftwareX* **1–**
693 **2**, 19–25 (2015).
- 694 40. Schmid, N. *et al.* Definition and testing of the GROMOS force-field versions
695 54A7 and 54B7. *Eur. Biophys. J.* **40**, 843–856 (2011).
- 696 41. Michaud-Agrawal, N., Denning, E. J., Woolf, T. B. & Beckstein, O. MDAAnalysis:
697 A toolkit for the analysis of molecular dynamics simulations. *J. Comput. Chem.*
698 **32**, 2319–2327 (2011).
- 699 42. Cock, P. J. A. *et al.* Biopython: freely available Python tools for computational
700 molecular biology and bioinformatics. *Bioinformatics* **25**, 1422–1423 (2009).

701