1    **A systematic genotype-phenotype map for missense variants in the human**

2    **intellectual disability-associated gene *GDI1***

3    Rachel A. Silverstein[1,2,3,5], Song Sun[1,2,3,4,6], Marta Verby[1,2,3], Jochen Weile[1,2,3,4], Yingzhou

4    Wu[1,2,3,4], Marinella Gebbia[1,2,3], Iosifina Fotiadou[1,2,3,4], Julia Kitaygorodsky[1, 2, 3,4], Frederick P.

5    Roth[1,2,3,4,*]

6    **Author affiliations**

7    1.  Lunenfeld-Tanenbaum Research Institute, Sinai Health, Toronto, ON M5G 1X5, Canada

8    2.  The Donnelly Centre, University of Toronto, Toronto, ON M5S 3E1, Canada

9    3.  Department of Molecular Genetics, University of Toronto, Toronto, ON M5S 3E1,

10      Canada

11   4.  Department of Computer Science, University of Toronto, Toronto, ON M5S 2E4, Canada

12   5.  Present address: Division of Medical Sciences, Harvard Medical School, 260 Longwood

13      Ave, Boston, MA 02115, USA

14   6.  Present address: Analytical Sciences, Sanofi Pasteur, Toronto, ON M2R 3T4, Canada

15   •  Correspondence should be addressed to F.P.R. via fritz.roth@utoronto.ca

16

17

18

19

20

21

22

23

24

**Abstract**

26  Next generation sequencing has become a common tool in the diagnosis of genetic diseases.

27  However, for the vast majority of genetic variants that are discovered, a clinical interpretation is

28  not available. Variant effect mapping allows the functional effects of many single amino acid

29  variants to be characterized in parallel. Here, we combine multiplexed functional assays with

30  machine learning to assess the effects of amino acid substitutions in the human intellectual

31  disability-associated gene, *GDI1*. We show that the resulting variant effect map can be used to

32  discriminate pathogenic from benign variants. Our variant effect map recovers known

33  biochemical and structural features of *GDI1* and reveals additional aspects of *GDI1* function. We

34  explore how our functional assays can aid in the interpretation of novel *GDI1* variants as they are

35  discovered, and to re-classify previously observed variants of unknown significance.

36

**Background**

38       Next-generation sequencing is now routinely practiced in the diagnosis of genetic

39  conditions. However, the usefulness of these methods is limited by our ability to interpret the

40  genetic variants that are discovered. The Genome Aggregation Database (gnomAD) (1), has

41  amassed over 4.6 million unique missense variants present in the human population. Of these

42  missense variants, 99% are rare (minor allele frequency $< 0.5\%$) (2) and only 13% have a

43  definitive clinical interpretation available on ClinVar (3). Therefore, methods to close the gap

44  between variant identification and interpretation are needed.

45       Several approaches to variant interpretation are available, including genome wide

46  association studies (GWAS), family segregation analysis, functional assays, and computational

2

47    prediction of variant effects. Of these, GWAS and computational prediction can both be used to

48    interpret data at a scale commensurate with the numbers of human genetic variants. However,

49    GWAS is of limited value for the interpretation of rare variants due to limited statistical power

50    and error in associations that is increased due to small sample sizes (4). Current computational

51    prediction approaches are considered at best weak evidence for clinical variant interpretation (5).

52    Functional assays have traditionally been used to test variants on an individual basis, but these

53    experiments are resource-intensive and this evidence is unlikely to be available at the time a

54    newly-discovered variant is first classified. However, it has become possible to perform

55    multiplexed assays of variant effect (MAVE), enabling the testing of functional effects for large

56    numbers of missense variants in parallel (2,6–8). For example, a framework for variant effect

57    mapping of human genes by complementation in *S. cerevisiae* has been previously described and

58    applied to multiple genes (8–10). This framework has been shown to identify, at stringent

59    confidence thresholds (90% precision), two to three times more pathogenic variants than are

60    identified by computational prediction alone (8–10). Here, we apply this framework to carry out

61    large-scale testing of missense variants of human *GDI1*, one of multiple genes on the X

62    chromosome that have been found to contain mutations causing X-linked non-syndromic

63    intellectual disability (11).

64        The *GDI1* gene encodes the protein GDI1 (Rab GDP dissociation inhibitor alpha). In

65    mammals, GDI1 is expressed primarily in the brain and is necessary for the control of endocytic

66    and exocytic pathways in neurons and astrocytes through the spatial and temporal control of

67    numerous Rab proteins (12,13). GDI1 functions to extract inactive GDP-bound Rab from

68    membranes by binding and solubilizing the genranylgeranyl anchor (a post-translational

69    modification at C-terminal cysteine residues which anchors Rabs to membranes) (14). *GDI1*-null

3

70    mouse models show deficits in short- and long-term synaptic plasticity and behavioral

71    phenotypes including alteration of hippocampus-dependent forms of short-term memory, spatial

72    working memory and associative fear-related memory (12). In humans, *GDI1* loss-of-function

73    variants can cause non-syndromic intellectual disability (ID), characterized by cognitive

74    impairment in the absence of other symptoms or physical anomalies (11). The form of ID caused

75    by *GDI1* variants follows an X-linked semi-dominant pattern of inheritance, with hemizygous

76    males being most severely affected and female carriers showing milder or no symptoms (15,16).

77          As a common condition which has been estimated to affect up to 3% of the general

78    population (11), ID presents a diagnostic challenge due to its many potential causes. Alterations

79    in over 700 genes have been associated with ID, few of which are frequently-occurring (17,18).

80    Separating causal from benign genetic variation in ID patients is therefore a significant clinical

81    challenge. Indeed, although an etiological diagnosis brings substantial benefits for patients and

82    their families (19), including more accurate prognosis, genetic counselling on recurrence risk,

83    and earlier access to resources within the community and specialized education programs, only

84    ~30% of ID patients receive an etiological diagnosis (20,21). Proactive functional testing for

85    variants in genes associated with ID could aid in the identification of causal variants and

86    facilitate earlier etiological diagnosis.

87          Here, we present large-scale measurements of the functional effects of missense variation

88    in *GDI1*.  Variant assay results are consistent with our knowledge of GDI1 function. A

89    comparison of variant scores with ClinVar annotations suggests that the map will prove useful in

90    assigning pathogenicity to genetic variation.

91

92    **Results**

4

93 **Multiplexed yeast complementation efficiently identifies damaging *GDI1* variants**

94       To efficiently test the deleteriousness of *GDI1* missense variants, we used a previously-

95 validated humanized yeast model system(22). In this system, the *Homo sapiens GDI1* (*HsGDI1*)

96 can complement a temperature sensitive allele of the orthologous *Saccharomyces cerevisiae* gene

97 Gdi1 (*Sc*Gdi1 (Ts)) and thereby restore yeast growth at restrictive temperatures. Importantly,

98 pathogenic variants of *HsGDI1* (L92P and R423P) showed a reduced ability to complement

99 *Sc*Gdi1(Ts) (22). This supported the possibility of a yeast-based functional assay of *HsGDI1*

100 variants, which we scaled up in order to test large numbers of missense variants in parallel (fig.

101 1a).

102       Mutagenesis of the *HsGDI1* open reading frame (ORF) was performed using a

103 previously-described pooled mutagenesis approach, Precision Oligo-Pool based Code Alteration

104 or "POPCode" (8), which uses oligonucleotide-directed codon randomization to yield a library of

105 single-codon *GDI1* variants. Following mutagenesis, the variant library was cloned into yeast

106 expression vectors and transformed *en masse* into a *S. cerevisiae* strain carrying the temperature

107 sensitive *Sc*Gdi1(Ts) allele. The yeast library was then grown competitively at restrictive

108 temperatures to induce selection for cells containing functional *HsGDI1* variants.

109       The library of *HsGDI1* ORFs was extracted from both pre- and post-selection yeast

110 populations, and sequenced deeply (with each position being observed in ~2 million reads). The

111 deep sequencing approach used was TileSeq (8), involving amplification and paired-end

112 sequencing of 12 "tiles", each ~100 nucleotides in length, that together cover the length of the

113 *GDI1* ORF. In order to decrease the rate of variants called erroneously due to sequencing error,

114 only variants that were detected in both forward and reverse reads were accepted. In total, 5534

115 unique amino acid changes were detected. To understand the rate at which missense variants are

5

139   synonymous variants was 1 and the median log ratio of variants containing a premature stop

140   codon was 0 (medians shown in fig. 1c). When calculating median log ratios, we included only

141   high confidence measurements (SD < 0.3) and, because nonsense mutations near the C-terminus

142   result in less severe loss of complementation, we only considered nonsense mutations within the

143   first 400 amino acids of the *GDI1* ORF (fig. S3).  In order to estimate fitness scores for the

144   remaining 80% of amino acid changes and refine scores of variants that were less well measured,

145   we applied a previously-described imputation pipeline (24). This pipeline uses the Gradient

146   Boosted Tree method to impute missing values based on intrinsic features of the data set

147   including average fitness of nearby variants, amino acid substitution matrix scores

148   (BLOSUM100 (25)), and variant effect scores predicted by computational methods including

149   PolyPhen-2 (26), and PROVEAN (23). To avoid low-confidence predictions based on limited

150   experimental data, imputation was not performed for amino acid positions with fewer than 3

151   well-measured variants. The result was a 'variant effect map' encompassing the majority of all

152   possible amino acid substitutions in *GDI1* (fig. 2). The most important features for predicting

153   fitness scores in this data set were the average fitness scores of the three most similar variants at

154   the same amino acid position, followed by BLOSUM100, PolyPhen2, and PROVEAN scores

155   (fig. S4).

156

157   **Our variant effect map is consistent with known biochemical features of GDI1**

158        The GDI1 protein contains four sequence conserved regions (SCRs), SCR1, SCR2,

159   SCR3A and SCR3B, common to all members of the Rab-GDI/CHM superfamily (27). Together,

160   SCR1 and SCR3B form a Rab-binding platform at the apex of the GDI1 structure (27,28) (fig.

161   3a). SCR3A contains a mobile effector loop (MEL) which constitutes a membrane receptor

7

162    binding site as well as a helix flanking the lipid binding pocket (29,30). At its N-terminal end,

163    SCR2 contains the C-terminus–binding region (CBR), which forms an essential interaction with

164    the C-terminus of Rab (28).

165        To determine overall patterns of variant deleteriousness within GDI1, we took the

166    average fitness score of all variants at a given amino acid position resulting in a "positional

167    fitness score" (fig. 3b). As expected, average fitness was significantly lower in the sequence-

168    conserved regions than in other parts of the protein (fig. 3c, 3d), supporting the notion that these

169    regions are important for biological function. We modeled the sequence of *H. sapiens* GDI1 on

170    the crystal structure of *S. cerevisiae* RabGDP-dissociation inhibitor in complex with prenylated

171    YPT1 GTPase (28) (the yeast homolog of human Rab-1A). The conserved face of GDI1

172    constituting the Rab binding platform contains the majority of residues with low positional

173    fitness scores (fig. 3a). Mutations in the SCR1 and SCR3B segments exhibited the lowest

174    positional fitness on average (fig 3d), consistent with previous mutational analysis showing that

175    disrupting these regions leads to decreased Rab binding and inability of GDI1 to extract Rab

176    from membranes (27). Since the C-terminal non-conserved region showed a striking increase in

177    average fitness scores around residue 425 (fig. 3b), we divided this region into two separate

178    sections, "linker 3", consisting of residues 460-424, and "C-terminus", consisting of residues

179    425-447. Mutations in the 22 "C-terminus" residues were significantly less deleterious than those

180    in linker 3 (Wilcoxon p<0.01).  The non-conserved region between SCR1 and SCR2 (termed

181    "linker 1") also exhibited high fitness scores, suggesting that variation here is also well tolerated

182    (fig. 3d).

183        Compared to SCR1 and SCR3B, variants in SCR2 were significantly less deleterious

184    (Wilcoxon p<0.01, and p<0.001 respectively).  On average, fitness scores of variants in SCR2

8

185 were comparable to those in the non-conserved region between SCR2 and SCR3A (termed

186 "linker 2") and the non-conserved linker 3 region (fig. 3d). Within SCR2, variants with the most

187 severe fitness effects tend towards the N-terminal CBR segment (fig. 3b). However, altering any

188 one of several hydrophobic residues within the helices flanking the lipid binding pocket,

189 especially Leu216 and Leu144, also yielded low positional fitness (fig. 3e). The location of these

190 residues, coupled with their average positional fitness scores, suggests that they may play an

191 important role in geranylgeranyl binding.

192        Deleterious mutations within the SCR3A region were observed predominantly towards

193 the C-terminus. Residues within the MEL region had moderate average positional fitness scores

194 between 0.5 to 0.75. It was previously reported that when MEL mutations Arg218Ala,

195 Tyr219Ala, and Ser222Ala are introduced into the corresponding positions of the yeast protein

196 *Sc*Gdi1, they do not cause visible growth defects in yeast.  However, when any one of these is

197 introduced in combination, they can exacerbate the effects of partial loss-of-function variants

198 elsewhere in *GDI1* (29). Our results show that single mutants Arg218Ala, Tyr219Ala, and

199 Ser222Ala each result in modest loss of function with fitness scores of 0.75 +/- 0.18, 0.67 +/-

200 0.22, and 0.66 +/- 0.13 respectively (regularized standard error for fitness scores was calculated

201 as described in materials and methods). It is possible that our competition-based assay was more

202 sensitive to minor growth changes and thus able to detect growth defects not detected by spotting

203 assays. While the study by Luan et al. only tested mutations in residues 218 - 222, we observed

204 some variants just outside of this region to be extremely deleterious, especially a short β strand

205 (termed β-strand e3 in Luan *et al.*) from residues Ser222 to Pro227 (fig. 3e). Despite the

206 importance of this segment indicated by our map, a biological function for this strand segment

207 has not been described.

9

208

**Relating fitness score to severity of intellectual disability**

210    Severity of ID is highly variable with cases ranging from mild to profound (31).

211 Although the severity of ID has been reported for only three *GDI1* missense variants have been

212 reported to date, we explored whether there was potential for variant fitness scores to predict the

213 severity of the associated ID phenotype. Males from a family with the Leu92Pro variant, were

214 reported to suffer from mild to moderate ID (11,32). For this variant, we obtained a fitness score

215 of 0.74 +/- 0.03.  Individuals in a family carrying the Gly237Val variant were reported to have

216 moderate ID (33), and we observed a corresponding lower fitness score of 0.55 +/- 0.07 for

217 Gly237Val.  Thus, the order of the fitness scores for these two variants agreed with the reported

218 order of ID severity. We note however that, like 20 (80%) of the 25 variants listed in the ClinVar

219 database, both Leu92Pro and Gly237Val are currently annotated as a variants of uncertain

220 significance, highlighting the need for better tools for interpretation.  Finally, a family carrying

221 the Arg423Pro variant suffered moderate to severe ID (15).  We did not observe Arg423Pro in

222 our assay, and were only able to impute a score with necessarily higher estimated uncertainty

223 (0.64 +/- 0.24). Although fitness scores may be predictive of ID severity; it is currently

224 insufficient to draw this conclusion from only reported ID severity data.

225

***GDI1* variant effect map predicts pathogenic variants with higher precision than**

**computational methods alone**

228    In order to test whether fitness scores from the *GDI1* map can provide useful evidence for

229 determining variant pathogenicity, we wished to determine whether our variant effect map can be

230 used to separate known benign from damaging alleles. Our set of presumed-damaging variants

10

231   included the only variant currently annotated as pathogenic (Arg423Pro (15)) and the additional

232   missense *GDI1* variants discussed above: Leu92Pro (11,32) and Gly237Val (33) based on

233   evidence from clinical reports. Because the number of currently known human pathogenic variants

234   is small, we also included four missense variants in highly-conserved regions which have been

235   previously shown to inhibit the ability of GDI1 to extract Rab3A from membranes in rat

236   synaptosomes, Tyr39Val, Glu233Ser, Met250Tyr, and Thr248Pro (27). To establish a reference

237   set of presumed-tolerated variants, we extracted all variants in gnomAD that had been observed in

238   male subjects (who are hemizygous for *GDI1* and less likely to be ID given that gnomAD excludes

239   subjects with early-childhood disease). Although we cannot rule out the possibility that our set of

240   presumed-damaging variants contains some tolerated variants, nor that our set of tolerated variants

241   contains some damaging variants, we reasoned these sets would enable a conservative estimate of

242   the ability of our scores to distinguish damaging from tolerated variation.

243       We observed that our sets of presumed tolerated and damaging alleles were well-separated

244   based on fitness score (fig 4a). Although fitness scores for presumed-damaging variants showed a

245   strong tendency to have lower scores, the lowest score amongst these was 0.5 and none were null-

246   like. We next calculated a precision-recall curve (fig. 4b) showing, as we change the fitness score

247   threshold below which a variant is deemed "damaging," the trade-off between precision (fraction

248   of below-threshold variants that are damaging) and recall (fraction of damaging variants that are

249   below the threshold). For comparison we also provide precision-recall curves for commonly used

250   computational predictors of variant effect including PolyPhen-2 (34), PROVEAN (23), and

251   VARITY (35) (fig. 4b). Our variant effect mapping framework was able to identify 6 out of 7

252   damaging variants (87% recall) with 100% precision using a fitness score threshold of <0.68. We

253   identified all damaging variants (100% recall) with 88% precision when a threshold fitness score

254    of <0.74 was used. The most accurate computational predictors were VARITY and PolyPhen-2,

255    which were each able to identify ~75% of damaging variants with ~75% precision.

256         Because single computational predictors are rarely used in isolation, we wondered

257    whether a combined computational prediction score, encompassing data from PolyPhen-2,

258    PROVEAN, and VARITY could separate damaging and tolerated variants with accuracy similar

259    to our variant effect mapping framework. We rationalized that agreement between multiple

260    computational prediction methods might be interpreted as stronger evidence for variant effect

261    than a single computational method alone. We therefore scaled the PolyPhen-2, PROVEAN, and

262    VARITY scores, and our fitness scores for the tolerated and damaging variant sets described

263    above such that scores ranged from 0 to 1 (with 0 representing most damaging and 1 representing

264    most tolerated). This  allowed us to make comparisons between the different score types (fig.

265    4c). Notably, all the prediction methods were able to accurately identify all of the damaging

266    variants (fig 4c). (Note that VARITY only generates predictions for single nucleotide variants, so

267    scores were not generated for 3 out of 7 damaging variants). However, it can be seen that the

268    increased accuracy of our VE mapping framework is due to the lack of false positives (prediction

269    of a damaging variant/low fitness score when the variant is in fact tolerated). Moreover, the three

270    computational methods tended to agree on many of the false positives, each assigning them low

271    scores when the variant was in fact tolerated. For instance, Gly114Cys, Arg141Leu, Arg218Gln,

272    and Arg292Trp were four particularly prominent false positives where each of the computational

273    methods predicted a low score, however, the fitness score generated by our VE mapping

274    framework correctly indicated that the variant was tolerated. Thus, we conclude that agreement

275    between computational predictors does not necessarily appear to be an indicator of accuracy. To

276    further illustrate this, wished to generate a "combined computational predictor score" which

12

277    takes into account the agreement between different computational predictors. We therefore took

278    the median of the scaled scores from the three computational prediction tools, reasoning that this

279    would eliminate outliers where an individual prediction tool did not agree with the other two.

280    When we plotted the precision recall curve for this "combined computational predictor" we

281    found that it did not perform better than the individual predictors (fig. 4d). This is consistent with

282    our notion that the shortcomings of the computational prediction methods are not due to

283    individual outlying predictions.

284

285    **Using our VE map to interpret clinically-relevant missense variants**

286        To facilitate the use of fitness scores as evidence to classify variants, we wished to

287    calculate likelihood ratios that convey the extent to which one should raise or lower the

288    probability that a variant is damaging, based on the fitness score. To this end, we estimated

289    probability density functions that describe the distributions of scores from our presumed-

290    damaging and -tolerated variant sets (see Methods). Then, the ratios of probability density for

291    damaging and tolerated variants can be used to obtain a damaging:tolerated likelihood ratio for

292    variants with any given fitness score. By this method, we determined that variants with fitness

293    scores below 0.72 were over 10 times more likely to be damaging than tolerated and variants

294    with fitness scores above 0.81 were over 10 times more likely to be tolerated than damaging.

295        We wondered whether our map could aid in the interpretation of GDI1 variants of

296    unknown significance which have been observed in the clinic (fig. 5). The ClinVar database lists

297    25 missense variants in GDI1, only four of which currently has a definitive clinical

298    interpretation. For 15 out of the 21 variants without a definitive interpretation, we were able to

299    generate an interpretation of either "deleterious" or "tolerated" with odds ratios greater than 1:10

13

300    based on our VE map (fig. 5). In order to be conservative with our interpretations, any variants

301    which had intermediate fitness scores leading to odds ratios less than 1:10 were labeled as

302    "unknown". Of note, we discovered four additional variants (in addition to those previously

303    included in our "likely damaging" variant set) that, were found by our assay to be highly

304    deleterious: R35W, G40V, R290S, and V381E. The latter three of these variants had almost null-

305    like scores. This highlights the possibility that ID due to *GDI1* mutations is under-diagnosed due

306    to current limitations in clinical variant interpretation.

307            Interestingly, Phe158Ser, annotated on ClinVar as "likely pathogenic" based on the

308    amino acid change being located within in a conserved region (SCR2), was non-conservative

309    with respect to amino acid properties, and was not observed as a common variant in the NHLBI

310    Exome Sequencing Project (37). However, our map score for Phe158Ser ([0.875 +/- 0.03]

311    originally, [0.870 +/- 0.03] post-refinement) does not provide strong evidence that this variant is

312    damaging. Using our current model based on the distributions of known pathogenic and benign

313    variants, and using no prior assumptions about the pathogenicity of Phe158Ser, a fitness score of

314    0.87 indicates the odds that the variant is damaging is less than 1:100. If our likelihood ratio

315    calibration is accurate, then even given a very strong prior belief (P = 0.99) that this variant is

316    damaging, the posterior odds would be less than 1:10.

317

318

319    **Discussion**

320            Towards clinical variant interpretation, the likelihood ratios that we derived for each

321    variant from our map could be discretized as strong, moderate, or supporting evidence for the

322    functional impact of a variant, and combined with other evidence using American College of

14

323  Medical Genetics and Genomics/Association for Molecular Pathology (ACMG/AMP) guidelines

324  (5). Alternatively, a Bayesian framework consistent with ACMG/AMP guidelines has been

325  proposed (36), in which the likelihood ratios we provide could be used directly and

326  quantitatively to infer variant pathogenicity (in the context of other evidence such as family

327  history, co-segregation, etc.).

328  A major drawback of our likelihood estimation approach is the limited number of known

329  damaging and tolerated variants currently available. Due to small sample sizes, our current

330  estimate of the score distributions of known damaging and tolerated variants is only an

331  approximation. As more variants in *GDI1* are discovered, assigned clinical significance and

332  added to databases such as ClinVar, this information should be incorporated to more confidently

333  estimate likelihood ratios.

334  In addition to variant interpretation, variant effect maps can also provide insights into the

335  function of a protein's structural components. In previous studies, structure-function analysis of

336  GDI1 has been largely focused on the conserved regions common to all members of the

337  *GDI/CHM* superfamily. Our results confirm that variants within the conserved regions forming

338  the Rab binding platform do tend to be the most deleterious. However, certain residues within

339  non-conserved regions exhibited fitness scores that suggested damaging substitutions. These

340  positions may be important for protein folding or stability, or contribute to functional roles of

341  GDI1 not shared by other members of the GDI family. While the MEL region has been the focus

342  for mutational analysis within SCR3A, we found that variants flanking the MEL region,

343  especially within β-strand e3, appeared markedly more deleterious. These findings can be used

344  to guide further mutational analysis of *GDI1*, aimed at discovering the specific functional roles

345  of each of these regions.

15

346          Due to the length of the *GDI1* ORF, the coverage of well-measured amino acid

347     substitutions for *GDI1* (19%) was somewhat lower than had been achieved for previous genes

348     studied using this approach (8,9). Nonetheless, precision-recall analysis revealed that, after

349     imputation by machine learning, the variant effect map was able to predict pathogenic variants

350     with greater accuracy than current computational methods alone, and with precision similar to

351     previously studied genes. Thus, using experimental data for a minority of substitutions, we could

352     accurately score variant effects for the majority of amino acid changes.

353          As genetic testing and exome sequencing continue to be used as diagnostic tools for

354     genetic disorders, it is expected that more patients with novel *GDI1* mutations will be

355     discovered. This map can be used to assist the interpretation of variants immediately upon their

356     discovery, thus accelerating the diagnostic process which is often costly, time-consuming, and

357     stressful for patients and their families. Due to the highly heterogeneous etiology of ID, it is

358     reasonable to expect that response to therapeutic and pharmacological interventions may also

359     vary in accordance with the cause of ID. Unfortunately, therapeutic guidelines rarely

360     differentiate between different forms of ID. Increased rates of etiological diagnoses could

361     improve our understanding of rare forms of ID and aid in the development of more personalized

362     guidelines for management and treatment.

363

364     **Conclusions**

365     Here we have presented the first variant effect map for single amino acid substitutions in *GDI1*,

366     and showed that map scores could distinguish presumed-damaging from presumed-tolerated

367     variants with better precision than current computational approaches (including Polyphen2,

368     VARITY, and PROVEAN) at all recall thresholds. Furthermore, our variant effect map recovers

369    known biochemical and structural features of GDI1 and provides insights into structural regions

370    which may be important for GDI1 function.

371

372

373

374    **Methods**

375

376    **Strains and Plasmids**

377    The *S. cerevisiae* strain carrying the temperature sensitive Gdi1 allele, TSA64 (*gdi1-1::KanR;*

378    *his3Δ1 leu2Δ0 ura3Δ0 met15Δ0*) (gift from G. Tan, C. Boone and B. Andrews) was used as a

379    host for the *GDI1* variant library. The Gateway destination vector used to express Hs*GDI1*,

380    pHYC-NatMX (CEN/ARS-based, ADH1 promoter, and NatMX marker), was constructed

381    previously (22). The Hs*GDI1* ORF clone (pDONR223-*GDI1*) was obtained from the Human

382    ORFeome v8.1 library (38).

383

384    **Construction of *GDI1* variant library by POPcode mutagenesis**

385    POPcode mutagenesis was performed on the *GDI1* ORF as described previously (9):

386    Oligonucleotides of 28-38 bases were designed to target each codon in the open reading frame of

387    *GDI1*, such that the targeted codon is replaced with a NNK-degenerate codon (a mixture of all

388    four nucleotides in the 1st and 2nd codon positions, and a mixture of G and T in the 3rd

389    position). Oligos were annealed to uracilated *GDI1* template, gaps between annealed

390    oligonucleotides were filled using KAPA HiFi Uracil+ DNA polymerase, nicks were sealed

391    using T4 DNA ligase, and the wild type template was degraded using Uracil-DNA-Glycosylase.

17

392    The variant library was transferred to the yeast expression vector, pHYC-NatMX, by *en masse*

393    Gateway LR reaction (8) followed by transformation into NEB5a competent E. coli cells (New

394    England Biolabs) and selection for ampicillin resistance. Plasmids extracted from a pool of

395    ~100,000 clones were transformed into the *S. cerevisiae* temperature-sensitive strain TSA64 *en*

396    *masse* using EZ Kit Yeast Transformation kit (Zymo Research). The entire transformed library

397    was grown in selective media (YPD + clonNAT) for two overnights. All yeast growth was

398    carried out at permissive temperature (25C).

399

400    **High-throughput yeast-based complementation**

401    For the pre-selection condition, plasmids were extracted from two 9 ODU samples of yeast

402    culture carrying the variant library (to be used for downstream tiling PCR). For the selective

403    condition, two replicates of 20 ODU of cells were inoculated into 200ml of YPD + clonNAT and

404    grown to full density at restrictive temperature (38°C) with shaking. Plasmids for tiling PCR

405    were extracted from 9 ODU of each culture following competitive growth. In parallel, 2 ODU of

406    TSA64 expressing wild type *GDI1* was inoculated into 20ml of YPD + clonNAT. Wild type

407    pools were grown under the same conditions as the POPcode library and plasmid was extracted

408    from 9 ODU samples to be used as a control for sequencing error during TileSeq.

409

410    **Measurement of allele frequencies in pre- and post-selective pools by TileSeq**

411    TileSeq was performed on the plasmids extracted from pre-selective, post-selective, and wild

412    type pools as described previously (8): (i) The *GDI1* ORF was amplified with primers carrying a

413    binding site for Illumina sequencing adaptors; (ii) each amplicon was indexed with an Illumina

414    sequencing adaptor; (iii) paired end sequencing was performed on the tiled amplicons to an

415    average sequencing depth of ~ 2 million reads. Raw sequencing reads were mapped to the *GDI1*

416    ORF using Bowtie2 (39). A custom Perl script (40) was used to parse the alignment files to count

417    the number of co-occurrences of a codon change in both paired reads. Mutational counts for each

418    tiled region were subsequently normalized by the corresponding sequencing depth, generating a

419    "raw data" file (table S1) where mutational counts are expressed in "reads per million", i.e. the

420    number of reads normalized to a depth of 1M reads (indicated as "reads/million" below).

421

422    **Data processing and fitness score calculation**

423    Processing of raw read count data (available in table S1) was carried out using the "legacy2.R"

424    script (41). This script is derived from the "legacy.R" script from the tileseqMave R package

425    described previously (10), with several modifications to improve filtering and fitness score

426    calculation for variants detected at low frequencies. Read counts for each variant in the wild type

427    control were subtracted from the corresponding read count for variants in each condition in order

428    to account for the detection of variants due to sequencing error. An enrichment ratio ($\phi$) was

429    calculated for each variant as the ratio of the normalized read counts after selection to before

430    selection. Since there was less agreement between replicate read counts for variants present at

431    lower frequencies in the pre-selection pool, a pre-filter was applied to remove all variants present

432    in fewer than 200 reads/million in either replicate. The cut-off value of 200 reads/million was

433    chosen in order to maximize the *t*-statistic measure of separation of mean synonymous and

434    pathogenic log ratios (fig. S2a). Additionally, any variants with read counts within 3 standard

435    deviations of zero in the post-selective condition were removed from the data set due to the

436    possibility that they were lost due to a bottleneck effect when sampled from the pre-selective

437    pool. As described previously (8), standard deviation estimates were regularized according to a

19

438     method for Bayesian regularization described by Baldi and Long (42), which improves

439     confidence estimates for measurements for which few replicates are available (in this case, two).

440     A fitness score ($FS_{MUT}$) was calculated for each variant as $\ln(\phi_{MUT}/\phi_{STOP})/\ln(\phi_{SYN}/\phi_{STOP})$, where

441     $\phi_{MUT}$ is the enrichment ratio calculated for each variant, $\phi_{STOP}$ is the median enrichment ratio of

442     all well-measured nonsense variants and $\phi_{SYN}$ is the median enrichment ratio of all well-

443     measured synonymous variants, such that $FS_{MUT}$ equals zero when $\phi_{MUT}$ equals $\phi_{STOP}$ and $FS_{MUT}$

444     equals one when $\phi_{MUT}$ equals $\phi_{SYN}$. Well-measured variants included in the calculation of the

445     medians $\phi_{STOP}$ and $\phi_{SYN}$ were those for which enrichment ratios between replicates agreed highly

446     with regularized standard deviation less than 0.3. Because nonsense mutations after residue 400

447     did not result in complete loss of function (fig. S3), nonsense mutations at amino acid positions

448     greater than 400 were excluded from the $\phi_{STOP}$ calculation. Fitness scores generated through this

449     pipeline are available in table S2.

450

451     **Imputation for missing variant effect map positions and fitness score refinement**

452     Imputation was performed using the variant effect imputation web server (24). The imputation

453     machine learning model was trained on the fitness scores of the experimentally measured

454     variants using the Gradient Boosted Tree (GBT) method. Features of the measured variants used

455     in the model include mean fitness scores of up to three nearest neighbor variants, standard fitness

456     score error of up to three most similar neighbor variants at the same position, number of

457     neighbors used, PolyPhen-2 score, PROVEAN score, and Blosum100 score. Fitness scores for

458     missing variants were not imputed for positions with fewer than three well-measured variants

459     due to insufficient functional data. Fitness scores of experimentally measured variants were also

460     refined using the weighted average of imputed and measured values (weighting by the inverse-

461    square of estimated standard error in each input value). Output of the imputation pipeline is

462    available in table S3.

463

464    **Construction of GDI1 homology model**

465    Human GDI1 (RefSeq: NP_001484.1) residues 1 - 426 were modeled on the crystal structure of

466    RabGDP-dissociation inhibitor in complex with prenylated YPT1 GTPase (PDB: 1UKV) using

467    Swiss-Model ProMod3 Version 1.3.0 (43). The poorly-aligned 21 C-terminal residues were not

468    included in the model.

469

470    **Likelihood ratio calculations**

471    Our set of presumed damaging human variants contained Leu92Pro (11), Arg423Pro (15), and

472    Gly237Val (33). Arg423pro is currently annotated as "pathogenic" on ClinVar. Leu92Pro was

473    previously annotated as pathogenic but is currently annotated as having "uncertain significance",

474    however we believe that d'Adamo et. al (11) provide strong evidence for the deleteriousness of

475    this mutation. Gly237Val was added to ClinVar more recently and is also annotated as having

476    "uncertain significance", however this variant seemed likely to be deleterious based on familial

477    segregation analysis by Duan et. al (33). We included four additional variants, Tyr39Val,

478    Glu233Ser, Met250Tyr, and Thr248Pro (27), which have not been observed in humans, but

479    which were shown to inhibit GDI1 function in functional assays.  The set of presumed tolerated

480    variants consisted of the 46 gnomAD variants from male subjects (hemizygous at the *GDI1*

481    locus), who were presumed to be healthy given that gnomAD excludes subjects with early

482    childhood disease. Normal distributions were fitted to the histograms of the fitness scores of

483    presumed damaging and tolerated variants by maximum likelihood parameter estimation in order

484    to obtain estimated probability density functions for pathogenic/disease and benign variants ($p_D$

485    and $p_B$ respectively). The normal distributions used are shown in fig. 4a (but scaled such that the

486    area under each curve equals 1 for likelihood ratio calculations). The damaging:tolerated

487    likelihood ratio for a variant with fitness score, $f$, was calculated as the ratio of the estimated

488    probability density functions evaluated at $f$: $\Lambda(D{:}T \mid f) = p_D(f)/p_T(f)$. This likelihood ratio

489    can be used together with prior beliefs about a variants' pathogenicity to calculate the odds that a

490    variant is damaging, $O(D{:}T \mid f)$, using the Odds form of Bayes' rule:

491
$$O(D{:}T \mid f) = \Lambda(D{:}T \mid f) \times \frac{P(D)}{P(T)}$$

492    where, $\Lambda(D{:}T \mid f)$ is the likelihood ratio, $P(D)$ is the prior probability that the variant is

493    damaging, and $P(T)$ is the prior probability that the variant is tolerated such that $P(T) = 1 -$

494    $P(D)$.

495

496

497    **Declarations**

498    **Ethics approval and consent to participate**

499    Not applicable

500    **Consent for publication**

501    Not applicable

502    **Availability of data and materials**

503    All data generated or analyzed during this study are included in this published article and its

504    supplementary information files.

505    **Competing interests**

22

506    F.P.R.is a scientific advisor and shareholder for SeqWell, Constantiam Biosciences and

507    BioSymetrics, and a Ranomics shareholder. S.S. is an employee of Sanofi Pasteur and a

508    Ranomics shareholder. M.V. is an employee and shareholder of Deep Genomics, Inc.  The

509    authors declare no other competing interests.

510    **Funding**

518    **Authors' contributions**

519    SS and FPR conceived the project. SS established the assay. MV,  SS, IF, and JK performed

520    mutagenesis and selection. MG performed sequencing. RAS performed primary data analyses

521    with contributions from SS and JW. RAS performed all downstream analyses of map scores,

522    including analysis of sequence-structure-function relationships. YW developed the machine

523    learning imputation pipeline with contributions from JW. RAS and FPR wrote the manuscript.

524    All authors read and approved the manuscript.

525    **Acknowledgments**

528

## References

1. Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alföldi J, Wang Q, et al. The mutational constraint spectrum quantified from variation in 141,456 humans. Nature. 2020;

2. Starita LM, Ahituv N, Dunham MJ, Kitzman JO, Roth FP, Seelig G, et al. Variant Interpretation: Functional Assays to the Rescue. American Journal of Human Genetics. 2017.

3. Landrum MJ, Lee JM, Benson M, Brown G, Chao C, Chitipiralla S, et al. ClinVar: Public archive of interpretations of clinically relevant variants. Nucleic Acids Res. 2016;

4. Wright CF, West B, Tuke M, Jones SE, Patel K, Laver TW, et al. Assessing the Pathogenicity, Penetrance, and Expressivity of Putative Disease-Causing Variants in a Population Setting. Am J Hum Genet. 2019;

5. Richards S, Aziz N, Bale S, Bick D, Das S, Gastier-Foster J, et al. Standards and guidelines for the interpretation of sequence variants: A joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology. Genet Med. 2015;

6. Fowler DM, Araya CL, Fleishman SJ, Kellogg EH, Stephany JJ, Baker D, et al. High-resolution mapping of protein sequence-function relationships. Nat Methods. 2010;

7. Fowler DM, Fields S. Deep mutational scanning: A new style of protein science. Nature Methods. 2014.

8. Weile J, Sun S, Cote AG, Knapp J, Verby M, Mellor JC, et al. A framework for exhaustively mapping functional missense variants. Mol Syst Biol. 2017;

9. Sun S, Weile J, Verby M, Wu Y, Wang Y, Cote AG, et al. A proactive genotype-to-

552        patient-phenotype map for cystathionine beta-synthase. Genome Med. 2020;

553   10.   Weile J, Kishore N, Sun S, Maaieh R, Verby M, Li R, et al. Shifting landscapes of human

554        MTHFR missense-variant effects. Am J Hum Genet. 2021;

555   11.   D'Adamo P, Menegon A, Nigro C Lo, Grasso M, Gulisano M, Tamanini F, et al.

556        Mutations in GDI1 are responsible for X-linked non-specific mental retardation. Nat

557        Genet. 1998;

558   12.   D'Adamo P, Welzl H, Papadimitriou S, Di Barletta MR, Tiveron C, Tatangelo L, et al.

559        Deletion of the mental retardation gene Gdi1 impairs associative memory and alters social

560        behavior in mice. Hum Mol Genet. 2002;

561   13.   Potokar M, Jorgačevski J, Lacovich V, Kreft M, Vardjan N, Bianchi V, et al. Impaired

562        αGDI Function in the X-Linked Intellectual Disability: The Impact on Astroglia Vesicle

563        Dynamics. Mol Neurobiol. 2017;

564   14.   Goody RS, Rak A, Alexandrov K. The structural and mechanistic basis for recycling of

565        Rab proteins between membrane compartments. Cellular and Molecular Life Sciences.

566        2005.

567   15.   Bienvenu T, Des Portes V, Saint Martin A, McDonell N, Billuart P, Carrié A, et al. Non-

568        specific X-linked semidominant mental retardation by mutations in a Rab GDP-

569        dissociation inhibitor. Hum Mol Genet. 1998;

570   16.   Strobl-Wildemann G, Kalscheuer VM, Hu H, Wrogemann K, Ropers HH, Tzschach A.

571        Novel GDI1 mutation in a large family with nonsyndromic X-linked intellectual disability.

572        Am J Med Genet Part A. 2011;

573   17.   Stessman HAF, Xiong B, Coe BP, Wang T, Hoekzema K, Fenckova M, et al. Targeted

574        sequencing identifies 91 neurodevelopmental-disorder risk genes with autism and

575  developmental-disability biases. Nat Genet. 2017;

576 18. Kvarnung M, Nordgren A. Intellectual disability & rare disorders: A diagnostic challenge.

577  In: Advances in Experimental Medicine and Biology. 2017.

578 19. Bélanger SA, Caron J. Evaluation of the child with global developmental delay and

579  intellectual disability. Paediatr Child Heal. 2018;

580 20. Monroe GR, Frederix GW, Savelberg SMC, De Vries TI, Duran KJ, Van Der Smagt JJ, et

581  al. Effectiveness of whole-exome sequencing and costs of the traditional diagnostic

582  trajectory in children with intellectual disability. Genet Med. 2016;

583 21. Thevenon J, Duffourd Y, Masurel-Paulet A, Lefebvre M, Feillet F, El Chehadeh-Djebbar

584  S, et al. Diagnostic odyssey in severe neurodevelopmental disorders: Toward clinical

585  whole-exome sequencing as a first-line diagnostic test. Clin Genet. 2016;

586 22. Sun S, Yang F, Tan G, Costanzo M, Oughtred R, Hirschman J, et al. An extended set of

587  yeast-based functional assays accurately identifies human disease mutations. Genome Res.

588  2016;

589 23. Choi Y, Chan AP. PROVEAN web server: A tool to predict the functional effect of amino

590  acid substitutions and indels. Bioinformatics. 2015;

591 24. Wu Y, Weile J, Cote AG, Sun S, Knapp J, Verby M, et al. A web application and service

592  for imputing and visualizing missense variant effect maps. Bioinformatics. 2019;

593 25. Henikoff S, Henikoff JG. Amino acid substitution matrices from protein blocks. Proc Natl

594  Acad Sci U S A. 1992;

595 26. Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P, et al. A

596  method and server for predicting damaging missense mutations. Nature Methods. 2010.

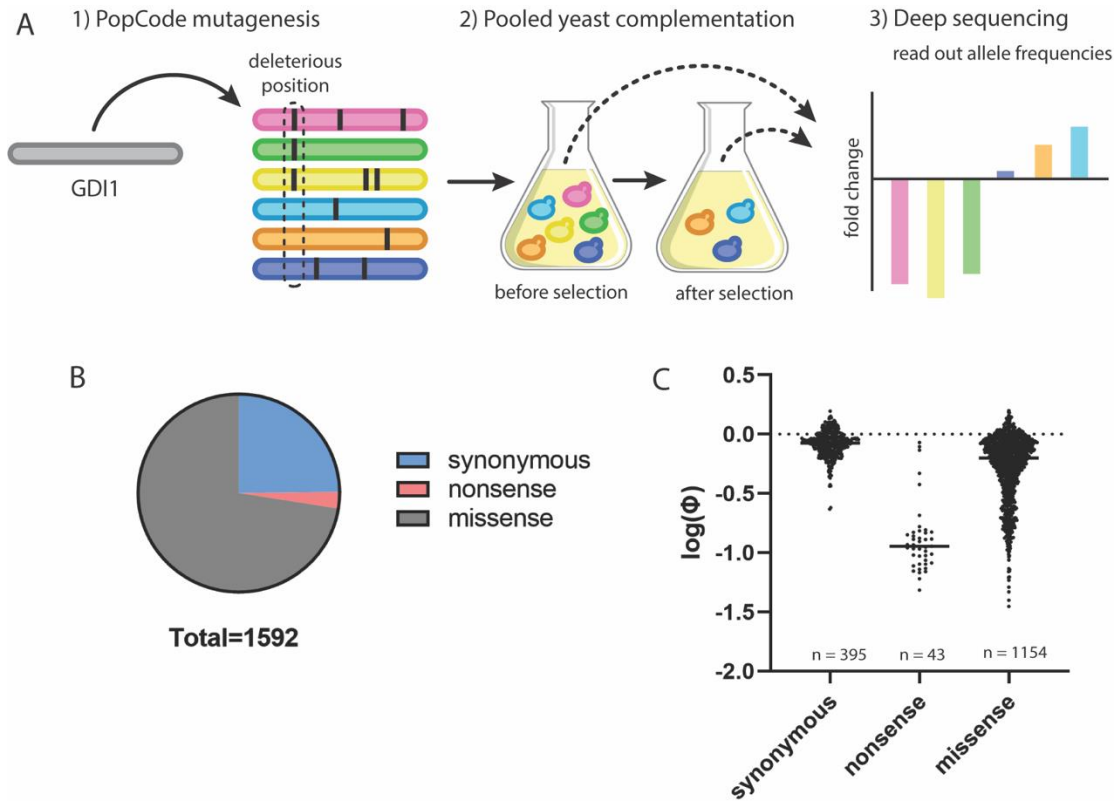597 27. Schalk I, Zeng K, Wu SK, Stura EA, Matteson J, Huang M, et al. Structure and mutational

26

598   analysis of Rab GDP-dissociation inhibitor. Nature. 1996;

599   28.   Rak A, Pylypenko O, Durek T, Watzke A, Kushnir S, Brunsveld L, et al. Structure of Rab

600         GDP-Dissociation Inhibitor in Complex with Prenylated YPT1 GTPase. Science (80- ).

601         2003;

602   29.   Luan P, Heine A, Zeng K, Moyer B, Greasely SE, Kuhn P, et al. A new functional domain

603         of guanine nucleotide dissociation inhibitor (alpha-GDI) involved in Rab recycling.

604         Traffic. 2000;

605   30.   An Y, Shao Y, Alory C, Matteson J, Sakisaka T, Chen W, et al. Geranylgeranyl switching

606         regulates GDI-Rab GTPase recycling. Structure. 2003;

607   31.   Patel DR, Apple R, Kanungo S, Akkal A. Intellectual disability: Definitions, evaluation

608         and principles of treatment. Pediatric Medicine. 2018.

609   32.   Hamel BCJ, Kremer H, Wesby-van Swaay E, Van Den Helm B, Smits APT, Oostra BA,

610         et al. A gene for nonspecific X-linked mental retardation (MRX41) is located in the distal

611         segment of Xq28. Am J Med Genet. 1996;

612   33.   Duan Y, Lin S, Xie L, Zheng K, Chen S, Song H, et al. Exome sequencing identifies a

613         novel mutation of the GDI1 gene in a Chinese non-syndromic X-linked intellectual

614         disability family. Genet Mol Biol. 2017;

615   34.   Adzhubei I, Jordan DM, Sunyaev SR. Predicting functional effect of human missense

616         mutations using PolyPhen-2. Curr Protoc Hum Genet. 2013;

617   35.   Wu Y, Li R, Sun S, Weile J, Roth FP. Improved pathogenicity prediction for rare human

618         missense variants. Am J Hum Genet. 2021;

619   36.   Tavtigian S V., Greenblatt MS, Harrison SM, Nussbaum RL, Prabhu SA, Boucher KM, et

620         al. Modeling the ACMG/AMP variant classification guidelines as a Bayesian

621        classification framework. Genet Med. 2018;

622   37.   NHLBI GO Exome Sequencing Project (ESP). Exome Variant Server. NHLBI. 2018.

623   38.   Yang X, Boehm JS, Yang X, Salehi-Ashtiani K, Hao T, Shen Y, et al. A public genome-

624        scale lentiviral expression library of human ORFs. Nat Methods. 2011;

625   39.   Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. Nat Methods.

626        2012;

627   40.   TileSeq package [Internet]. Available from: https://github.com/rothlab/tileseq_package

628   41.   Weile J, Silverstein R. Github/tileseqMave/legacy2.R [Internet]. Available from:

629        https://github.com/RachelSilverstein/tileseqMave/blob/master/R/legacy2.R

630   42.   Baldi P, Long AD. A Bayesian framework for the analysis of microarray expression data:

631        Regularized t-test and statistical inferences of gene changes. Bioinformatics. 2001;

632   43.   Studer G, Tauriello G, Bienert S, Biasini M, Johner N, Schwede T. ProMod3 - A versatile

633        homology modelling toolbox. PLoS Comput Biol. 2021;

634   44.   McGuffin LJ, Bryson K, Jones DT. The PSIPRED protein structure prediction server.

635        Bioinformatics. 2000;
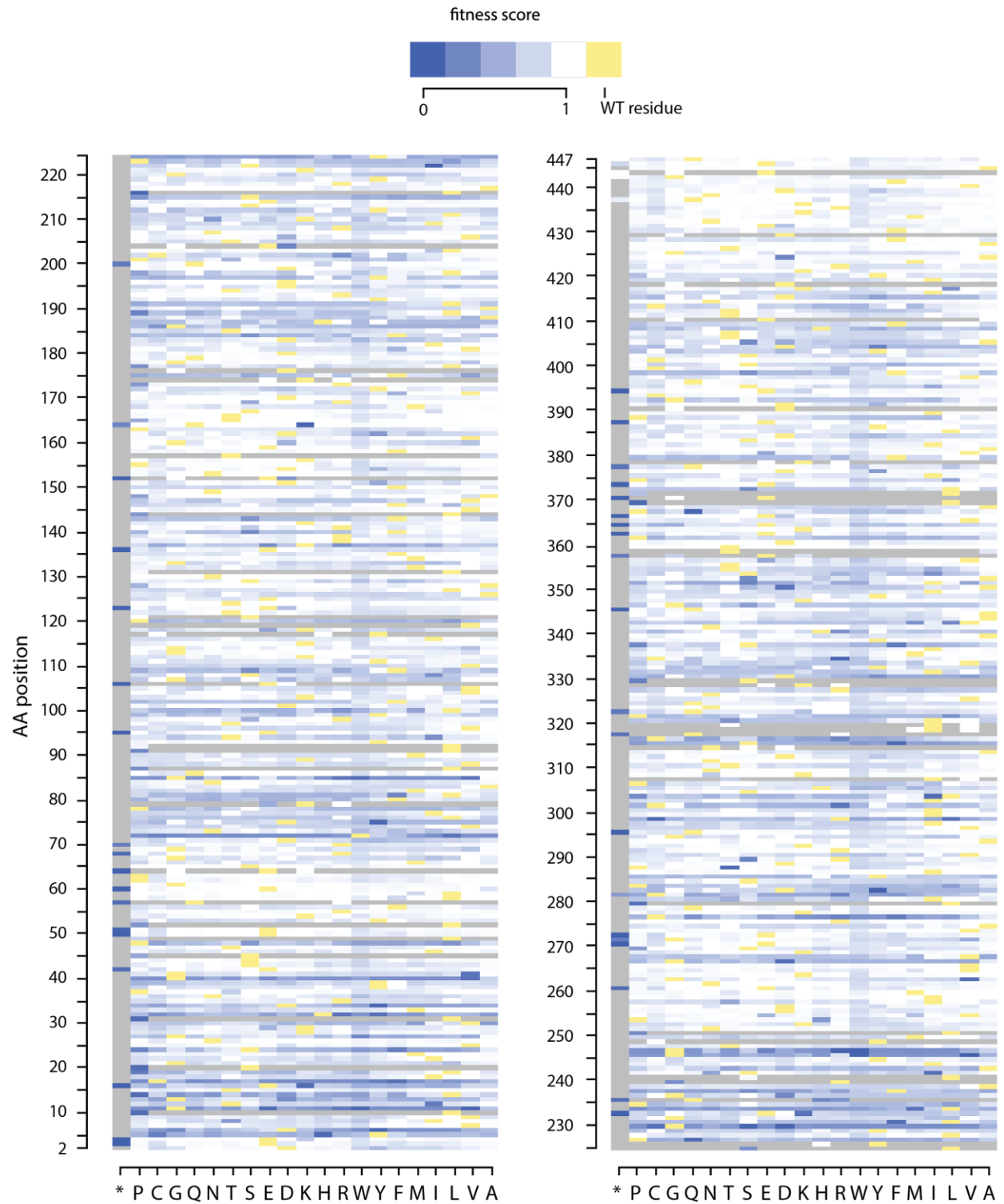
636

637   **Figures and Legends**

638

639

**Figure 1: High throughput yeast complementation screen separates synonymous and**

**nonsense *GDI1* variants**

a) Graphical overview of the variant effect mapping framework.

b) Number of well-measured variants recovered from the complementation screen.

c) Log($\phi$) values comparing pre- and post-selection variant frequencies for all well measured
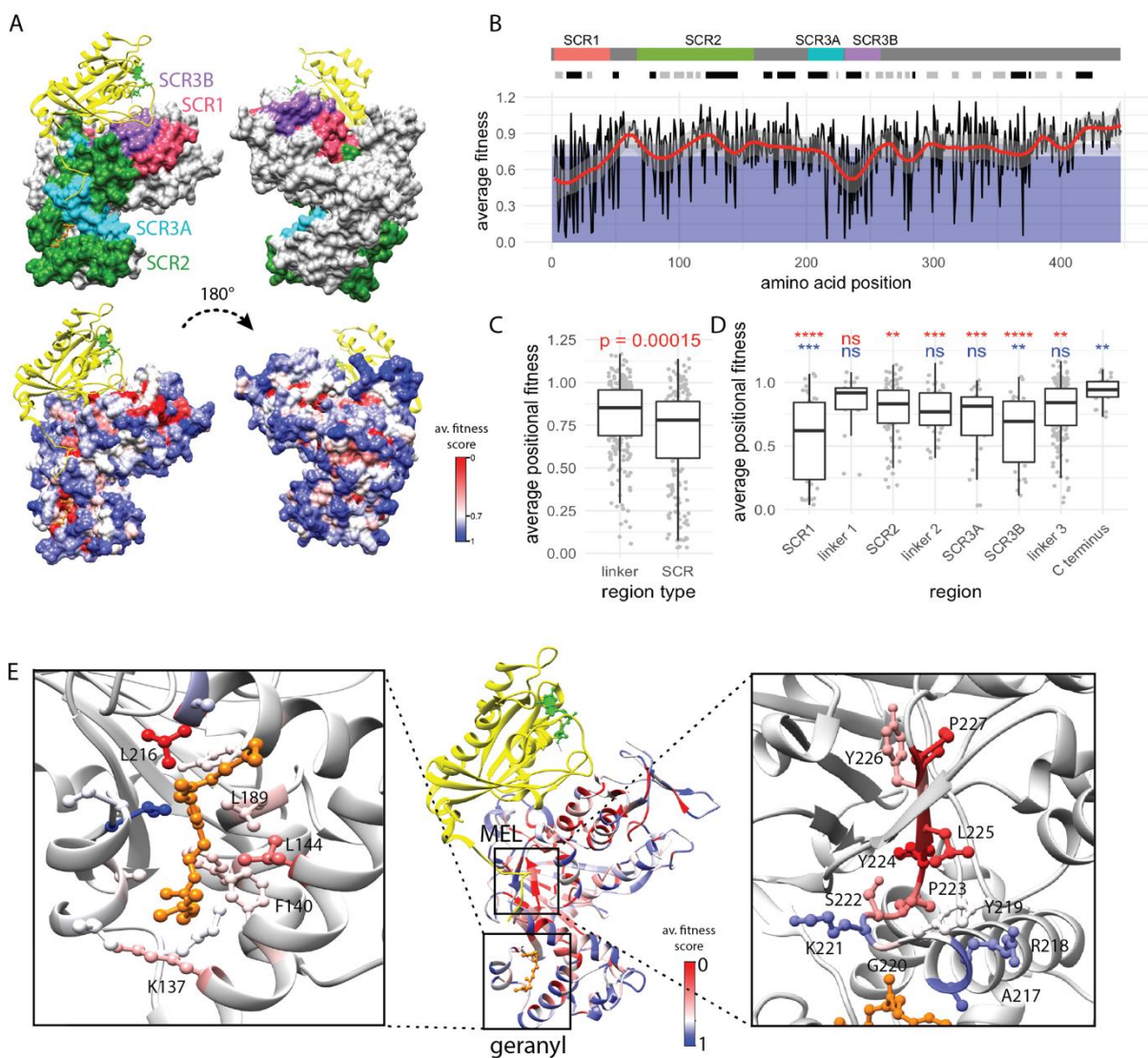
synonymous, nonsense and missense *GDI1* variants.

646

647

**Figure 2: *GDI1* variant effect map**

30

649    A *GDI1* missense variant effect map resulting from the complementation screen coupled with

650    imputation and refinement by machine learning. Fitness scores of 0 (blue) represent the median

651    behavior of complete loss of function variants (based on observed fitness of nonsense variants)

652    and fitness scores of 1 (white) represent wild type-like function (based on observed fitness of

653    synonymous variants). Yellow tiles represent the wild type amino acid at that position. Gray tiles

654    represent substitutions for which scores were not imputed due to insufficient data for

655    substitutions at that amino acid position.

656



657

658     **Figure 3: Fitness scores enable structure-function analysis of *GDI1***

659     (a) Homology model of human GDI1 (colored surface) modeled on the structure of *S. cerevisiae*

660     RabGDP-dissociation inhibitor in complex with prenylated YPT1 GTPase (yellow ribbon). In the

661     bottom panel, residues are colored according to their average positional fitness scores with 0

662     representing null-like scores (red) and 1 representing wild type-like scores (blue).
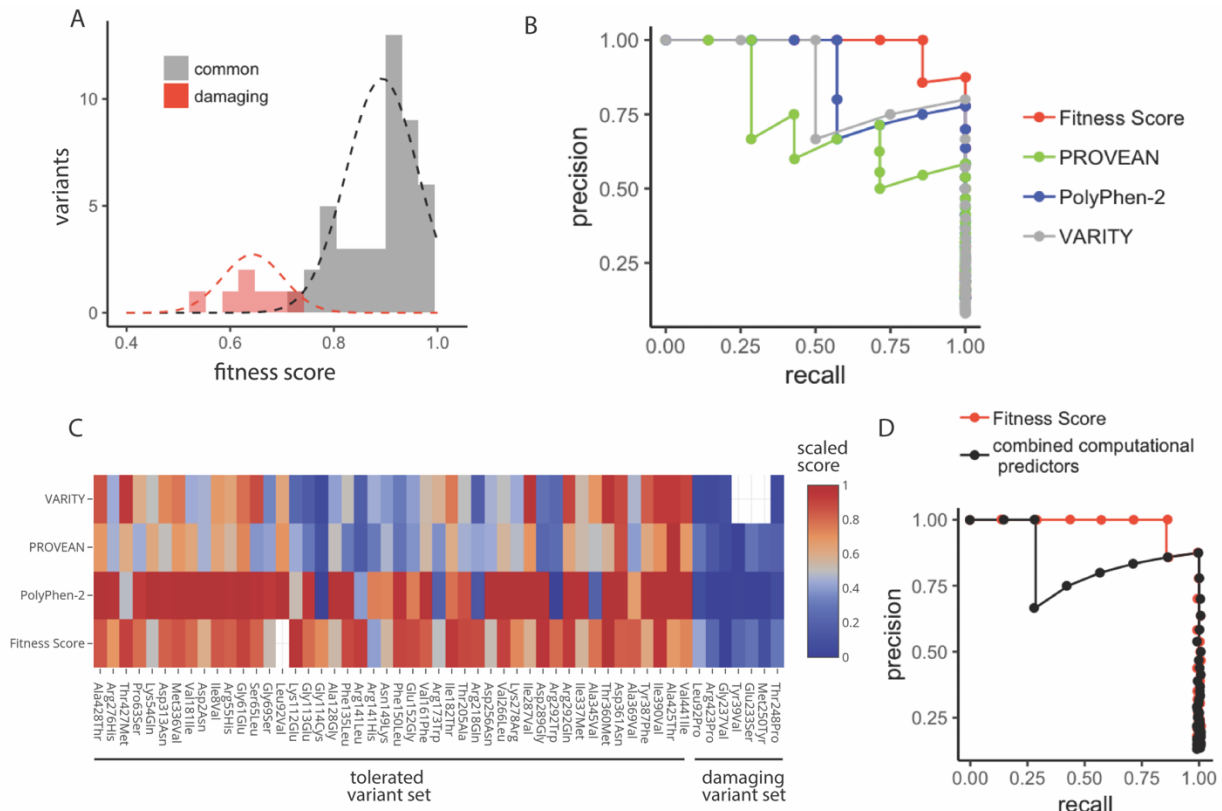
663     (b) Average fitness score of all variants at each amino acid position (black line) overlaid with a

664     smoothed summary curve (red). The dark blue region of the plot represents fitness scores less

665     than 0.72 (over 10 times more likely to be damaging than tolerated) and the white region

666     represents fitness scores over 0.81 (over 10 times more likely to be tolerated than damaging).

667     [Tolerated:damaging odds ratios were calculated as described in methods]. The tracks above the

668     plot represent: depiction of GDI1 with sequence conserved regions (SCRs) common to all

669     members of the GDI/CHM superfamily (top track) and; the secondary structures of human GDI1

670     as predicted by PSIPRED 4.0 (44) (bottom track; black = helix, gray = strand).

671     (c) Average fitness scores of amino acid positions within non-conserved or "linker" regions

672     versus sequence conserved regions. Significance level was determined using Wilcoxon signed-

673     rank test.

674     (d) Region-wise comparison of average positional fitness scores. Wilcoxon signed-rank tests

675     were performed comparing each region to the "C-terminus" region (red asterisks) and to SCR2

676     (blue asterisks). Significance levels are denoted by: * (p<0.05), ** (p<0.01), *** (p<0.001), and

677     **** (p<0.0001).

678     e) Center: Ribbon representation of human GDI1 modeled on the structure of *S. cerevisiae*

679     RabGDP-dissociation inhibitor in complex with prenylated YPT1 GTPase (yellow).

680     GDI1☐residues are colored by average positional fitness score. Left: side chains of all

681    hydrophobic residues within 5A of the geranylgeranyl group (orange). Right: side chains of

682    residues comprising the mobile effector loop and proximal beta strand.

683



684

**Figure 4: *GDI1* variant effect map separates damaging and common variants with higher**

**precision than current computational methods**

(a) Distribution of fitness scores for known damaging and known common (presumed tolerated)

*GDI1* missense variants. Common variants are comprised of 46 missense variants listed in

gnomAD which have been observed in at least one hemizygous individual.

(b) Precision-recall curve for our fitness scores compared to various computational methods for

variant interpretation. A sliding threshold was used for each score type starting at the lowest

score; variants below this threshold were called as damaging. For each threshold value, the

number of true damaging variants identified (true positives) and the number of benign variants

33

694    identified in error (false positives) was evaluated. The precision [true positives/(true positives +

695    false positives)] versus the recall [true positives/(true positives + false negatives)] is shown for

696    each threshold value.

697    c) Scaled fitness scores and computational predictor scores for all variants from our tolerated and

698    damaging variant sets. Scores were scaled such that all score types range from 0 to 1 with 0

699    representing most damaging and 1 representing most tolerated.

700    d) Precision recall curves for our fitness scores and for a "combined computational predictor

701    score" which is the median of scaled PolyPhen-2, PROVEAN, and VARITY scores (scaling was

702    performed as described in c).

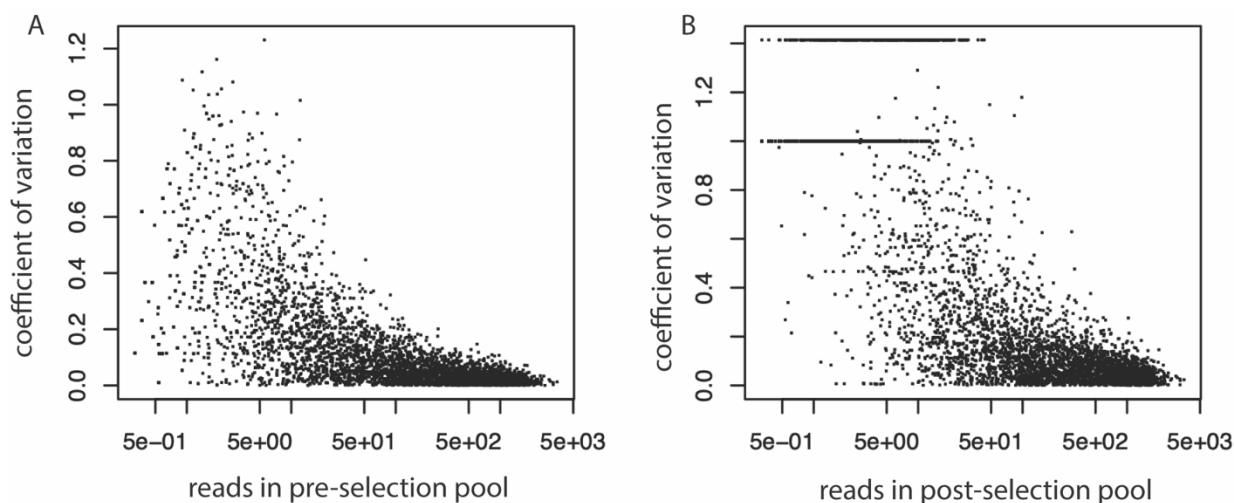| Variant | ClinVar annotation | fitness score | imputed score | standard error | conclusion |
|---------|--------------------|---------------|---------------|----------------|------------|
| R35W | Uncertain significance | 0.63 | 0.63 | 0.03 | deleterious |
| G40V | Uncertain significance | 0.10 | 0.11 | 0.05 | deleterious |
| S65T | Uncertain significance |  | 0.78 | 0.15 | unknown |
| S65L | Uncertain significance | 1.07 | 0.91 | 0.05 | tolerated |
| L76V | Uncertain significance |  | 0.71 | 0.20 | unknown |
| L92P | Uncertain significance | 0.74 | 0.74 | 0.03 | unknown |
| G113E | Uncertain significance | 0.91 | 0.90 | 0.03 | tolerated |
| A128G | Uncertain significance |  | 0.79 | 0.14 | unknown |
| R138W | Uncertain significance | 0.87 | 0.87 | 0.04 | unknown |
| F158S | Likely pathogenic | 0.88 | 0.87 | 0.03 | tolerated |
| Y192C | Uncertain significance | 0.99 | 0.99 | 0.02 | tolerated |
| R193H | Uncertain significance | 1.06 | 0.94 | 0.02 | tolerated |
| R193L | Uncertain significance | 0.92 | 0.91 | 0.04 | tolerated |
| G237V | Uncertain significance | 0.55 | 0.55 | 0.07 | deleterious |
| D289H | Uncertain significance |  | 0.77 | 0.19 | unknown |
| R290S | Uncertain significance | 0.17 | 0.20 | 0.05 | deleterious |
| R290H | Uncertain significance | 0.95 | 0.95 | 0.03 | tolerated |
| V381E | Uncertain significance | 0.19 | 0.36 | 0.17 | deleterious |
| T413A | Likely benign | 1.15 | 0.87 | 0.01 | tolerated |
| R423P | Pathogenic |  | 0.64 | 0.24 | unknown |
| T427M | Conflicting interpretations | 0.97 | 0.96 | 0.04 | tolerated |
| A428V | Uncertain significance | 1.16 | 0.86 | 0.03 | tolerated |

703

704 **Figure 5: Interpretation of clinically-relevant GDI1 missense variants from the ClinVar**

705 **database**

706 Fitness scores (experimentally measured where available, and computationally imputed for all

707 variants) are listed for all GDI1 missense variants listed on ClinVar. We concluded that a variant

708 is "deleterious" where the damaging:tolerated odds ratio was greater than 1:10 and vice versa for

709 "tolerated" variants.

710

711 **Supplemental Information**

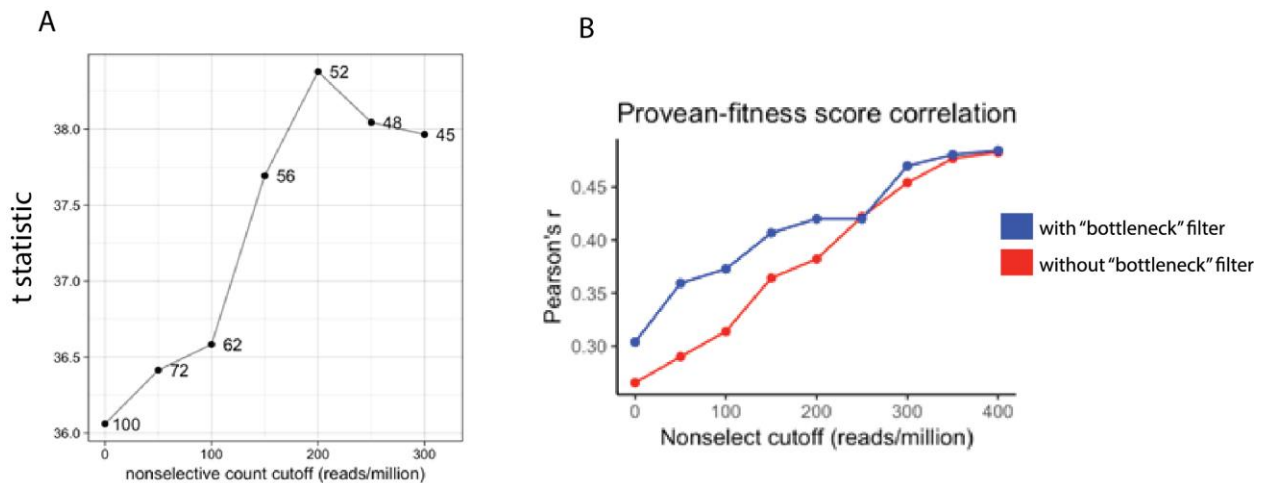712 **Supplemental Figures and Legends**

713 

714

715 **Figure S1: Variants present at low frequencies in complementation screen show poorer**

716 **agreement between replicates**

717 a) Coefficient of variation between two read count replicates for all detected variants in the pre-

718 selection pool versus frequency in the pre-selection pool (as measured by mean read count of the

719 two replicates).

35

720 b) Coefficient of variation between two read count replicates for all detected variants in the post-

721 selection pool versus frequency in the post-selection pool (as measured by mean read count of
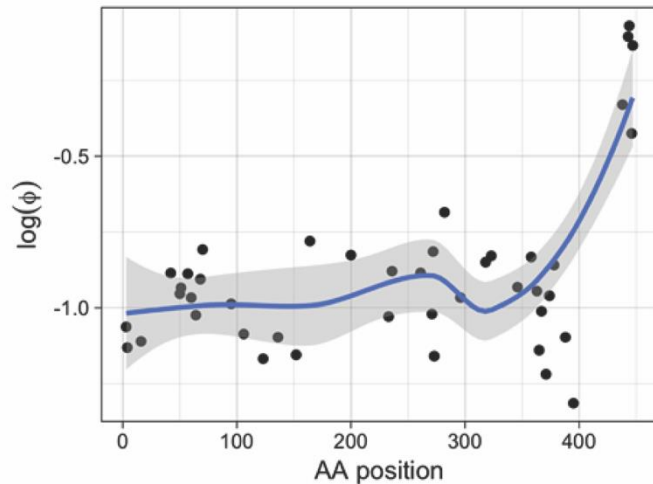
722 the two replicates).

723



724

**Figure S2: Filtering out variants present at low frequencies in the pre-selection pool**

**improves metrics of fitness measurement accuracy**

727 a) Multiple read count cut-offs were tested wherein read counts in the pre-selection pool were

728 filtered to include only high-frequency variants (present at frequencies greater than the cut-off

729 value). For each cut-off value tested, a two-sample t-statistic was calculated to evaluate the

730 separation of fold changes between nonsense variants and synonymous variants. A cut-off value

731 of 200 reads/million maximized the separation of synonymous and nonsense variants.

732 b) Correlation between PROVEAN scores and our fitness scores increase as variants are filtered

733 for higher frequency in the pre-selection variant pool. For each read count cutoff, the correlation

734 (Pearson's R) between our calculated fitness score (prior to imputation) and PROVEAN scores

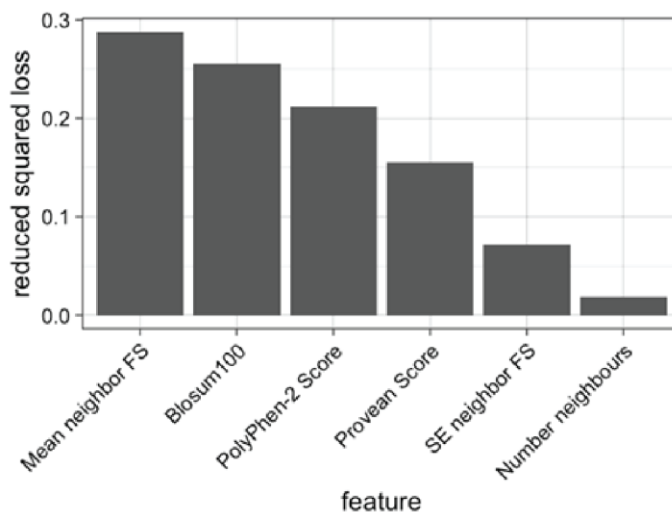735 for all missense variants was calculated.

736



737 **Figure S3: Fold changes (pre-selection/post-selection) of all measured nonsense mutations**

738 **in *GDI1***

739 Nonsense mutations after amino acid position 400 lead to less severe loss of complementation.



740

741 **Figure S4: Feature importance for gradient boosted trees imputation model**

742 Mean neighbor FS: the mean fitness scores of the 3 most similar amino acids at the same residue

743 position. SE neighbor FS: Standard error of the fitness scores of the 3 most similar amino acids

744 at the same residue position. Number neighbors: Number of variants measured at the same amino

745 acid position

746

37

747    **Descriptions of Supplementary Tables**

748    **Table S1: Table of raw yeast complementation data**

749    Table of unfiltered *GDI1* variant frequencies in pre- and post- selection deep sequencing pools.

750    Variants counts are presented in reads/million.

751    **Table S2: Fitness score table**

752    Table of fitness score data calculated for all well-measured *GDI1* variants.

753    **Table S3: Imputed scores**

754    Table of all fitness scores including computationally imputed scores for amino acid substitutions

755    not measured experimentally.