

Expanding the Atlas of Functional Missense Variation for Human Genes

Jochen Weile^{1,2,3,4,*}, Song Sun^{1,2,3,4,5,*}, Atina G. Cote^{1,2,3}, Jennifer Knapp^{1,2,3}, Marta Verby^{1,2,3}, Joseph Mellor^{2,6}, Yingzhou Wu^{1,2,3,4}, Carles Pons⁷, Cassandra Wong^{1,2}, Natascha van Lieshout¹, Fan Yang^{1,2,3,4}, Murat Tasan^{1,2,3,4}, Guihong Tan^{2,3}, Shan Yang⁸, Douglas M. Fowler⁹, Robert Nussbaum⁸, Jesse D. Bloom¹⁰, Marc Vidal^{11,12}, David E Hill¹¹, Patrick Aloy^{7,13}, Frederick P. Roth^{1,2,3,4,14,†}

¹Lunenfeld-Tanenbaum Research Institute, Mount Sinai Hospital, Toronto Ontario M5G 1X5, Canada. ²The Donnelly Centre and Departments of ³Molecular Genetics and ⁴Computer Science University of Toronto, Toronto, Ontario M5S 3E1, Canada. ⁵Department of Medical Biochemistry and Microbiology, Uppsala University, SE-75123 Uppsala, Sweden. ⁶SeqWell Inc, Boston, MA. ⁷Institute for Research in Biomedicine (IRB Barcelona), The Barcelona Institute for Science and Technology, Barcelona, Catalonia, Spain. ⁸Invitae Corp., San Francisco, CA. ⁹Department of Genome Sciences, University of Washington, Seattle, WA. ¹⁰Fred Hutchinson Research Center, Seattle, WA. ¹¹Center for Cancer Systems Biology (CCSB), Dana-Farber Cancer Institute, Boston, MA 02215, USA. ¹²Department of Genetics, Harvard Medical School, Boston, MA 02115, USA. ¹³Institució Catalana de Recerca i Estudis Avançats (ICREA), Barcelona, Catalonia, Spain. ¹⁴Canadian Institute for Advanced Research, Toronto, ON M5G 1Z8, Canada.

*Contributed Equally; †Corresponding Author

Abstract

Although we now routinely sequence human genomes, we can confidently identify only a fraction of the sequence variants that have a functional impact. Here we developed a deep mutational scanning framework that produces exhaustive maps for human missense variants by combining random codon-mutagenesis and multiplexed functional variation assays with computational imputation and refinement. We applied this framework to four proteins corresponding to six human genes: UBE2I (encoding SUMO E2 conjugase), SUMO1 (small ubiquitin-like modifier), TPK1 (thiamin pyrophosphokinase), and CALM1/2/3 (three genes encoding the protein calmodulin). The resulting maps recapitulate known protein features, and confidently identify pathogenic variation. Assays potentially amenable to deep mutational scanning are already available for 57% of human disease genes, suggesting that DMS could ultimately map functional variation for all human disease genes.

Keywords: Genotype-Phenotype/Deep Mutational Scanning/VUS/Complementation

Running title: Atlas of functional missense variation

Character count (excluding title page and references): 66,887

Introduction

Millions of people will soon have their genomes sequenced. Unfortunately, we have only a limited ability to interpret personal genomes, each carrying 100-400 rare missense variants (The 1000 Genomes Project Consortium, 2015) of which many must currently be classified as “variants of uncertain significance” (VUS). For example, gene panel sequencing aimed at identifying germline cancer risk variants in families yielded VUS for the majority of missense variants (Maxwell *et al*, 2016). Functional variants can be predicted, but when high precision is required, computational tools (Adzhubei *et al*, 2010; Choi *et al*, 2012) detect only one third as many pathogenic variants as experimental assays (Sun *et al*, 2016). Unfortunately, validated experimental assays enabling rapid clinical interpretation of variants are not available for the vast majority of human disease genes.

Deep Mutational Scanning (DMS) (Fowler *et al*, 2010; Fowler & Fields, 2014), a strategy for large-scale functional testing of variants, can functionally annotate a large fraction of amino acid substitutions for a substantial subset of residue positions. Recent DMS studies, for example, covered the critical RING domain of BRCA1 (Starita *et al*, 2015) associated with breast cancer risk, and the PPARG protein associated with Mendelian lipodystrophy and increased risk of type 2 diabetes (Majithia *et al*, 2016). Such maps can accurately identify functionality of a clinical variant in advance of that variant’s first clinical presentation. Diverse assays can be used for DMS (see Expanded View Table EV1). Functional complementation assays test the variant gene’s ability to rescue the phenotype caused by reduced activity of the wild type gene (or its ortholog in the case of trans-species complementation) (Lee & Nurse, 1987; Osborn & Miller, 2007). Cell-based functional complementation assays can accurately identify disease variants across a diverse set of human disease genes (Sun *et al*, 2016)

Challenges to the DMS strategy include the need to establish robust assays measuring each variant’s impact on the disease-relevant functions of a gene, and to generate maps that cover all possible amino acid changes. Also, published DMS maps have not typically controlled the overall quality of measurements nor estimated the quality of individual measurements. Thus, the use of DMS maps to confidently evaluate specific variants has been limited.

Here we describe a modular DMS framework to generate complete, high-fidelity maps of variant function based on functional complementation. This framework combines elements of previous DMS studies, uses machine learning to impute and improve the map with surprisingly high accuracy, and yields a confidence measure for each reported measurement. In the following sections, we give an overview of the overall framework for DMS, describe its initial application to the SUMO E2 conjugase UBE2L, present complete high-fidelity maps for three new disease-associated proteins and explore the potential for clinical relevance. Finally, we assemble information on functional assays for known human disease genes and conclude that DMS is already potentially extensible to the majority of human disease genes, suggesting the possibility of exhaustive maps of functional variation covering all human genes.

Results

We describe a framework for comprehensively mapping functional missense variation, organized into six stages (see Figure 1A): 1) mutagenesis; 2) generation of a variant library; 3) selection of functional variants; 4) read-out of the selection results and analysis to produce an initial sequence-function map; 5) computational analysis to impute missing values; and 6) computational analysis to refine measured values via machine learning. We describe and contrast two versions of this framework: DMS-BarSeq and DMS-TileSeq.

A barcode-based deep mutational scanning strategy

We first describe DMS-BarSeq and its application to map functional missense variation for the SUMO E2 conjugase UBE2I. In DMS-BarSeq, a heterogeneous pool of cells bearing a library of different barcoded expression plasmid is quantified via barcode-sequencing before and after selection. For Stage 1 of the DMS framework—mutagenesis—we sought a relatively even representation of all possible single amino acid substitutions. We wished to allow multiple mutations per clone, both because this allowed for greater mutational coverage for any given library size, and offered an opportunity to discover intragenic epistatic relationships. To this end, we scaled up a previous mutagenesis protocol (Seyfang & Huaqian Jin, 2004) to develop Precision Oligo-Pool based Code Alteration (POPCode), which yields random codon replacements (see Materials and Methods).

For Stage 2 of the framework—generation of a variant library—we employed *en masse* recombinational cloning of mutagenized UBE2I amplicons into a pool of randomly-barcoded plasmids (see Materials and Methods). The full-length UBE2I sequence and barcode of each plasmid was established using a novel sequencing method called KiloSeq which combines plate-position-specific index sequences with Illumina sequencing to carry out full-length sequencing for thousands of samples (see Materials and Methods). We retained clones that carried at least one amino acid substitution to generate a final library comprised of 6,553 UBE2I variants, covering different combinations of 1,848 (61% of all possible) unique amino acid changes. Variant plasmids were pooled, together with empty vector and wild type control plasmids (see Materials and Methods).

For Stage 3—selection for clones encoding a functional protein—we employed a *S. cerevisiae* functional complementation assay (Sun *et al*, 2016; Jiang & Koltin, 1996), based on human *UBE2I*'s ability to rescue growth at an otherwise-lethal temperature in a yeast strain carrying a temperature sensitive (ts) allele of the *UBE2I* orthologue *UBC9*. Despite a billion years of divergence, yeast functional complementation assays can accurately discriminate pathogenic from non-pathogenic human variants (Sun *et al*, 2016). The plasmid library from Stage 3 was transformed *en masse* into the appropriate ts strain. Pools were grown for 48 hours at the permissive (25°C) and selective (37°C) temperatures, respectively (see Materials and Methods).

To assess variant functions (Stage 4), barcodes were sequenced at multiple timepoints of the selection, enabling reconstruction of individual growth curves and normalized fitness quantification for each of the 6,553 barcoded strains. Functional complementation scores

were calibrated so that 0 corresponds to the fitness of the null-allele and 1 to wild type complementation (see Materials and Methods). Using replicate agreement and extent of library representation, we estimated our uncertainty in each fitness value (see Materials and Methods).

Before further refinement in Stages 5 and 6, we wished to assess the quality of the DMS-BarSeq complementation scores. Based on both technical (Figure 1B, top) and biological replicates (different clones carrying the same mutation; Figure 1B, bottom), we found scores to be reproducible (Pearson's R of 0.97 and 0.78, respectively). Semi-quantitative manual complementation assays for a subset of mutants that spanned the range of fitness scores (see Materials and Methods) correlated well with DMS scores. Indeed, agreement between large-scale and manual scores was on par with agreement between internal replicates of the large-scale scores (Figure 1B,C).

We also examined evolutionary conservation and computational predictors of deleteriousness, such as PolyPhen-2 (Adzhubei *et al*, 2010) and PROVEAN (Choi *et al*, 2012). Although each is an imperfect measure of the functionality of amino acid changes (Sun *et al*, 2016), each should and did correlate with DMS results (Figure 1D top panel, Appendix Figure S1). Finally, we confirmed that, as expected, amino acid residues on the protein surface are more tolerant to mutation than those in the protein core or within interaction interfaces (Figure 1D, bottom panel). Taken together, these observations support the biological relevance of the DMS-BarSeq approach.

A tiled-region strategy for mapping functional variation

While DMS-BarSeq has several advantages (see Discussion), its performance comes at the cost of producing an arrayed library of clones, each with known coding and barcode sequence. We therefore also evaluated an alternative approach, DMS-TileSeq in which each functional variant is detected via the effect of selection on the abundance of clones carrying that variant. The frequency of each variant in the pool is determined, before and after selection, by deep sequencing of short amplicons that tile the complete coding region.

In terms of mutagenesis (Stage 1), DMS-TileSeq is identical to DMS-BarSeq. Given the mutagenized amplicon library, the cloning step (Stage 2) was carried out by *en masse* recombinational subcloning into expression vectors (thereby skipping the step of arraying and sequencing individual clones). This plasmid pool was next transformed *en masse* into the *ubc9-ts* strain. Deep sequencing detected 97% of all possible missense variants in our expression library, and 100% of the amino acid substitutions that can be achieved via single-nucleotide mutation. As with DMS-BarSeq, DMS-TileSeq employs pooled strains grown competitively (Stage 3) at the permissive and selective temperatures. In Stage 4, like some previous DMS efforts (Doud & Bloom, 2016), we directly sequenced the coding region from the clone population to determine variant frequency before and after selection. Use of tiled amplicons enables individual template molecules to be sequenced on both strands, allowing elimination of most base-calling errors (Fowler *et al*, 2010; Whitehead *et al*, 2012; Zhang *et al*, 2016) (see Materials and Methods for details). This reduction in base-calling error allows us to more accurately measure lower allele frequencies in mutagenized libraries.

To further assess the reliability of DMS-TileSeq, we compared results with DMS-BarSeq for UBE2I. DMS-TileSeq and DMS-BarSeq correlation was similar to that observed between biological DMS-BarSeq replicates (Pearson's $R = 0.75$, Appendix Figure S2). DMS-TileSeq and DMS-BarSeq also behaved similarly in their agreement with manual complementation assays (Appendix Figure S3). Thus, DMS-TileSeq avoids the substantial cost of arraying and sequencing thousands of individual clones, while performing on par with DMS-BarSeq in terms of reliability of functional complementation scores.

After using regression to transform the DMS-TileSeq scores to the more intuitive scale of DMS-BarSeq (where 0 corresponds to the median score of null mutant controls and 1 corresponds to the median score of wildtype controls), we combined scores from the two methods, giving greater weight to more confident measurements (see Materials and Methods). scores emerging from this procedure are referred to as 'joint scores' below.

Machine learning to complete and refine maps

Although nearly all missense variants can be detected in our UBE2I TileSeq libraries, we only considered those variants present with 'allele frequency' sufficient to allow confident detection of allele-frequency reduction post-selection (see Materials and Methods). After filtering, 2563 of 3012 possible amino acid changes (85%) were well-measured. To complete missing entries in the map (Stage 5 in the framework), we trained a random forest (Breiman, 2001) regression model using the existing joint scores in the map. The model used four types of predictive feature: intrinsic (derived from other measurements in our map); conservation-based; chemico-physical; and structural. Particularly predictive features (Figure 2D) included the average score of observed substitutions at a given position, as weighted by measurement confidence. Conservation-based features included BLOSUM62 (Henikoff & Henikoff, 1992), SIFT (Ng & Henikoff, 2001) and PROVEAN (Choi *et al*, 2012) scores, and position-specific AMAS (Livingstone & Barton, 1993) conservation. Chemico-physical features included mass and hydrophobicity of the original and wildtype amino acids, and the difference between them. Structural features included solvent accessibility and burial in interaction interfaces. Where DMS-BarSeq scores for multi-mutant clones were available, we also used the confidence-weighted average score of all clones containing a particular substitution, and variant fitness expected from a multiplicative model (St Onge *et al*, 2007) (see Materials and Methods).

We assessed imputation performance using cross-validation. Surprisingly, the error (root-mean squared deviation or RMSD) of imputed values (0.33) was on par with that of experimentally measured data (Figure 2A). As an additional validation step, we performed manual complementation assays for a set of UBE2I variants that were not present in the machine learning training data set and compared the results against imputed values (Figure 2C), again finding strong agreement. Predictions showed the least error in positions with high mutation density and the most error for hypercomplementing variants, i.e. those yielding above-WT fitness levels in yeast (Figure 2B). Although hypercomplementation may indicate that a variant is adaptive in yeast, imputation generally predicted these variants to be deleterious, a hypothesis we explore further below.

In order to examine the impact of training set size on the imputation performance, we

performed sub-sampling analysis. Performance was poor below ~5% map completeness, increased dramatically at ~10% map completeness and then improved gradually (and approximately linearly) as completeness levels rose beyond 10% (Appendix Figure S4). We also computationally examined the expected impact of changing the mutagenesis method: When training only on SNP-accessible variants (for example, if one were using libraries generated by error-prone PCR), imputation RMSD was significantly worse ($P=2.4 \cdot 10^{-5}$) compared to a training set of equivalent sample size that can, as POPCode allows, contain all possible amino acid substitutions (Appendix Figure S4).

To refine less-confident experimental measurements (Stage 6 of the framework), we combined joint experimental scores with Random Forest-predicted scores from the imputation procedure, weighting by confidence level. Scores resulting from this combination are referred to as 'refined scores' below. Overall, most values were only adjusted minimally through refinement, with 90% of values being altered by less than 2.5% of the score difference between null and wt controls (Appendix Figure S5A). This reflects the fact that most values were already of high quality. To evaluate the effect on the minority of variants that required stronger refinement, we looked for cases that were of low quality in the DMS-TileSeq dataset, but well measured in the DMS-BarSeq experiment. These cases would allow us to treat the DMS-BarSeq values as an independent reference for comparison when performing the refinement procedure only on the DMS-TileSeq dataset. We identified six cases that fulfilled these criteria. In all six cases, refinement of DMS-TileSeq resulted in improvement, i.e. adjusted the corresponding values such that they more closely resembled the gold standard (Appendix Figure S5B). However, all changes were small, suggesting that our refinement procedure was overly conservative and that alternative weighting schemes should be explored as more 'ground truth' data becomes available.

Manual complementation assays, applied to a set of variants that represented the full range of refined scores (Appendix Figure S3), served to validate the reliability of the complete, refined functional map of UBE2I after imputation and refinement. The map, as seen in Figure 3A, fulfills biochemical expectations, with the hydrophobic core, the active site and protein interaction interfaces being most strongly impacted by mutations (Figure 3B). Detailed observations with respect to structure, biochemistry and epistatic behaviour of double mutants can be found in the Appendix texts.

Hypercomplementing variants are likely to be deleterious in humans

We further investigated UBE2I variants exhibiting hypercomplementation (Figure 3A). Manual assays confirmed that complementation with these mutants allows greater yeast growth than does the wild type human protein (Appendix Figure S6A). These hypercomplementing substitutions did not reliably correspond to 'reversion' substitutions that inserted the corresponding *S. cerevisiae* residue (Appendix Figure S6B). Some substitutions could be adaptive by improving compatibility with yeast interaction partners. Indeed, a comparison with co-crystal structure data (Gareau *et al*, 2012) shows that many of the hypercomplementing residues are on the surface proximal to the substrate, with some directly contacting the substrate's sumoylation motif (Figure 2C). *In vitro* sumoylation assays performed previously for a small number of UBE2I mutants revealed increased sumoylation

for some substrates (Bernier-Villamor *et al*, 2002). Comparing our map with these sumoylation assay results, we saw that cases of hypercomplementation were enriched for substrate specificity shift (Appendix Figure S6C). However, other cases of hypercomplementation hinted at different modes of adaptation (see Appendix text).

To explore whether variants exhibiting hypercomplementation are more likely beneficial or deleterious in a human context, we used a quantitative phylogenetic approach (Bloom, 2014, 2017) to compare three models relating the (refined) complementation scores to evolutionary preference for an amino acid variant: (a) evolutionary preference is directly proportional to complementation score; (b) preference has a ceiling at the wildtype complementation score (values >1 were set to 1); or (c) preference is set to the reciprocal of complementation score for mutations with greater-than-wildtype scores, corresponding to a deleterious effect of hypercomplementing mutations. We used the phydms software (Bloom, 2017) to test which of these three approaches best described the evolutionary constraint on a set of naturally occurring UBE2I homologs, using re-calculated refined scores that excluded conservation features from the imputation and refinement process, to avoid circularity when using natural sequence data to impute or refine scores. The best fit is achieved by treating variants with greater-than-wildtype complementation in yeast as deleterious in humans (Appendix Table S1). We therefore reinterpreted cases of hyperactive complementation in our map as deleterious and repeated the machine learning training, imputation and refinement procedure. Repeated cross-validation revealed the new imputed values based on the reinterpreted score matrix to be more reliable (i.e. reducing cross-validation RMSD from 0.33 to 0.24).

Variant impact maps for five additional disease-implicated genes

Having validated the framework, we sought to map functional variation for disease-relevant genes. We applied the higher-throughput TileSeq approach, coupled with yeast complementation, to a diverse set of genes: SUMO1, for which heterozygous null variants are associated with cleft palate (Andreou *et al*, 2007); Thiamine Pyrophosphokinase 1 (TPK1), associated with vitamin B1 metabolism dysfunction (Mayr *et al*, 2011); and CALM1, CALM2 and CALM3, associated with cardiac arrhythmias (long-QT syndrome (Crotti *et al*, 2013) and catecholaminergic polymorphic ventricular tachycardia (Nyegaard *et al*, 2012)). Because the three calmodulin genes encode the same polypeptide sequence, performing DMS for CALM1 also provided maps for CALM2 and CALM3.

As no corresponding DMS-BarSeq data was available to facilitate TileSeq score rescaling for these genes, we rescaled scores such that a score of 0 corresponded to the median of nonsense variants while a score of 1 corresponded to the median of synonymous variants. The Random Forest model underlying imputation and refinement was trained anew for each map. (Differences in the stringency of each selection have the potential to introduce non-linear changes in scale that will differ between maps.) Supporting the quality of the resulting four maps, each map showed clear differences in TileSeq-score distributions between likely-neutral (synonymous) and likely-deleterious (nonsense) variants (Appendix Figure S7).

To assess the impact of the machine learning imputation and refinement on the different maps, we measured the completeness of each map before and after imputation, the cross-

validation RMSD of the imputation, as well as the maximum standard error value for each map before and after refinement (Table 1). On average, 24.6% of scores were obtained purely by imputation, and 3.96% of scores were appreciably changed by >5% of the difference between null and wt controls as a result of refinement. Proteins for which overall map quality was initially lower were improved most by refinement, while others, like SUMO1, improved only modestly. Inspection of the maps yielded a number of interesting biochemical and structural observations (see Appendix texts).

Phylogenetic analysis of SUMO1, as for UBE2I, showed that variants that complement yeast better than wild-type are best modeled as being deleterious in humans (Appendix Table S1). As was done for the first map, we transformed above-wild-type fitness scores to be deleterious before the imputation and refinement step (see Materials and Methods). Because hypercomplementing substitutions may nonetheless provide interesting clues about differences between yeast and human cellular contexts, we provide both transformed (Figure 4) and untransformed (Appendix Figure S8) map versions.

DMS functional maps reflect clinical phenotypes.

To validate the utility of our maps in the context of human disease, we extracted known disease-associated variants from ClinVar (Landrum *et al*, 2016), as well as rare and common polymorphisms observed independent of disease from GnomAD (Lek *et al*, 2016), and somatic variants previously observed in tumors from COSMIC (Forbes *et al*, 2001).

While no germline disease-associated missense variants are known for UBE2I and SUMO1 in ClinVar, somatic cancer variants have been observed for both genes according to COSMIC. Somatic variants in these three genes exhibited higher functional impact in DMS maps than germline variants (Wilcoxon $P=2.6 \times 10^{-5}$) (Figure 5A). This does not necessarily suggest that either of these genes are cancer drivers, as even passenger somatic variants should subject to less purifying selection than germline variants, but it does lend further credence to the biological relevance of our maps.

For TPK1, many very rare variants (minor allele frequency or MAF < 10^{-6}) are seen in GnomAD. The majority of these variants scored as deleterious (Appendix Figure S9A). Thiamine Metabolism Dysfunction Syndrome, reported to be caused by variants in TPK1, is a severe disease to which patients succumb in childhood (Mayr *et al*, 2011). Although GnomAD attempted to exclude subjects with severe pediatric disease, the abundance of rare predicted-deleterious variants may be understood by the disease's recessive inheritance pattern. Using phased sequence data from the 1000 Genomes Project (The 1000 Genomes Project Consortium, 2015) to determine diploid genotypes in TPK1, we assigned each subject a diploid score corresponding to the maximum (refined) score across each pair of alleles. This improved prediction performance markedly, leading to complete separation between disease and non-disease genotypes using DMS, PROVEAN or PolyPhen-2 scores (Appendix Figure S9B). However, additional compound heterozygotes with known disease status will be required to compare DMS with computational methods in the task of identifying TPK1 disease variants.

Because the inheritance pattern of calmodulin disorders is typically dominant (Crotti *et al*,

2013), we did not consider diploid genotypes but simply evaluated the ability of the (refined) DMS scores to distinguish disease from non-disease variants (Figure 5B). DMS scores performed well according to precision-recall analysis, with an area under the precision-recall curve (AUC) of 0.72, exceeding both PROVEAN (AUC=0.48) and PolyPhen-2 (AUC=0.47) (Figure 5C). At a stringent precision threshold of 90%, DMS exceeded twice the sensitivity of PROVEAN and PolyPhen-2.

We further wished to explore how classification based on these observations would perform on variants of uncertain significance (VUS). We therefore examined missense VUS substitutions seen by Invitae, a clinical genetic testing company. Ten rare calmodulin variants had been encountered, of which half were from tests ordered due to a cancer indication, the other half from tests ordered for a cardiac disease indication. Blinded to indication, we ranked and classified the 10 Invitae VUS variants by DMS score (Table 2). We classified variants as 'damaging' if they were below both the highest score of known pathogenic variants and the lowest score of GnomAD variants, and classified variants as 'benign' if they were both above the highest-scoring known-pathogenic variant and the lowest-scoring GnomAD variant. All others were classified as 'uncertain'. Using these criteria, two Invitae variants were classified as damaging, two as uncertain, and six as benign. Based on the patient test indications subsequently revealed by Invitae, five out of the six variants we classified as benign were ordered due to a non-cardiac indication, while both variants with damaging predictions and both with VUS predictions corresponded to cardiac indications. Overall, DMS scores (which do not depend on the somewhat arbitrary classification system described above) showed a significant association with cardiac indications ($P=0.008$, $U=24.5$; Mann-Whitney-U test).

Potential for applying deep mutational scanning more widely

DMS mapping requires an *en masse* functional assay that can be applied at the scale of 10^4 - 10^5 variant clones. Among ~4000 disease genes, examination of four systematic screens and curated literature suggests that ~5% of human disease genes currently have a yeast complementation assay (Sun *et al*, 2016; Kachroo *et al*, 2015; Hamza *et al*, 2015). This number could grow dramatically via systematic complementation testing under different environments and genetic backgrounds. Moreover, complementation assays can also be carried out in other model systems including human cells (Hart *et al*, 2015). Based on only three large-scale CRISPR studies (Hart *et al*, 2015; Wang *et al*, 2014; Blomen *et al*, 2015), cellular growth phenotypes (which might serve as the basis for an *en masse* selection) have already been observed in at least one cell line for 29% of human disease genes. Beyond complementation, assays of protein interaction can, in addition to identifying variants directly impacting interaction, can detect variants ablating overall function through effects on protein folding or stability. In a recent study, approximately two thirds of disease-causing variants were found to impact at least one protein interaction (Sahni *et al*, 2015). Although only a minority of human protein interactions have been mapped (Rolland *et al*, 2014), already 40% of human genes have at least one interaction partner detectable by yeast two-hybrid assay in a recent screen (Rolland *et al*, 2014). Taking the union of available assays, we estimate that 57% of known disease-associated genes (Expanded View Dataset EV2) already have an assay that is potentially amenable to DMS.

Discussion

The framework for systematically mapping functional missense variation we describe here combines elements of previous DMS studies and introduces a new mutagenesis strategy and a machine learning-based imputation and refinement strategy. This framework enables DMS maps that are ‘complete’ in the sense that high-quality functional impact scores are provided for all missense variants to full-length proteins. Application to four proteins highlighted complex relationships between the biochemical functions of these proteins with phenotypes in the yeast model system. Analysis of pathogenic variation, especially for calmodulin, supported the potential clinical utility of DMS maps from this framework.

The two described versions of DMS, DMS-BarSeq and DMS-TileSeq, each have advantages and limitations (Appendix Table S2). DMS-BarSeq permits study of the combined effects of variants at any distance along the clone, and therefore can reveal intramolecular genetic interactions. For DMS-BarSeq fully-sequenced variant clones are arrayed, enabling further investigation of individual variants. DMS-BarSeq can directly compare growth of any clone to null and wild type controls, resulting in an intuitive scoring scheme. However, despite the efficient KiloSeq strategy for sequencing arrayed clone sets we report for the first time here, DMS-BarSeq is more resource-intensive. Although the regional sequencing strategy of DMS-TileSeq can only analyze fitness of double mutant combinations falling within the same ~150bp tile, it is far less resource-intensive than DMS-BarSeq.

In the TileSeq libraries, many clones will contain multiple amino acid substitutions (for example our UBE2I library averages 2.1 amino acid changes per clone, so in this case the Poisson distribution predicts that 62% of clones will have more than one amino acid substitution). This raises the concern that the presence of multiple mutations in these clones could obscure the functional effect of any single mutation. However, the DMS-TileSeq scores in our UBE2I map follow a distribution for synonymous variants that is unimodal and distinct from the (also unimodal) distribution for nonsense codon variants, indicating that, despite the presence of multiple variants in many clones, we are able to clearly separate neutral variants from null variants (Appendix Figure S7). Indeed, despite the fact that DMS-TileSeq libraries often have multiple clones, and DMS-BarSeq analysis was based on single-mutant clones, our evaluations of DMS-TileSeq and DMS-BarSeq maps of UBE2I indicated that they are of similar quality. This may be understood by considering the large clone population analyzed (typically >100,000 clones), which means that the impact of each given query mutation will be an average effect over many genetic backgrounds. This is analogous to detecting a shift in a single SNP’s allele frequency between case and control populations in a genome-wide association study despite variation at other loci. We cannot exclude the possibility that libraries that have a higher average number of mutations per clone, or are from genes for which a higher fraction of missense variants are deleterious, may perform less well in TileSeq. In such cases it may be desirable to reduce the mutation rate in the PopCode protocol. This can be achieved by using lower concentrations of mutagenic oligos, or using “DMS by parts”, in which multiple mutagenized libraries are generated, each focusing mutagenesis on a different segment of the protein.

Given that most missense variants in individual human genes are single-nucleotide variants (Lek *et al*, 2016), and given that only ~30% of all possible amino acid substitutions are

accessible by single nucleotide mutation, one might wonder why codon mutagenesis should be preferred over single-nucleotide mutagenesis. We see four arguments for codon-level mutagenesis: 1) knowing the functional impact of all 19 possible substitutions at each positions enables clearer understanding of the biochemical properties that are required at each residue position; 2) an analysis of >60,000 unphased human exomes (Lek *et al*, 2016) found that each individual human harbors ~23 codons containing multiple nucleotide variants that together could encode an amino acid not encoded by either single variant; 3) it is not straightforward to generate balanced libraries in which every single-nucleotide variant has roughly equal representation, given that error-prone amplification methods strongly favor transition mutations over transversion mutations, while still avoiding frequent introduction of new stop codons; 4) the major cost of DMS will likely continue to be development and validation of the functional assay, so using codon-level mutagenesis instead of (or in addition to) nucleotide-level mutagenesis has a relatively small impact on overall cost.

This study yielded four DMS maps measuring functional impact of ~16,000 missense variants. The maps generated for sumoylation pathway members UBE2I and SUMO1, and disease-implicated genes CALM1/2/3 and TPK1 using our framework were consistent with biochemical expectations while providing new hypotheses. DMS maps based on functional complementation were highly predictive of disease-causing variants, outperforming popular computational prediction methods such as PolyPhen-2 or PROVEAN, confirming previous observations (Sun *et al*, 2016). Given sufficient experimental data for training, our results show that imputation can ‘fill the gaps’ with scores that are nearly as reliable as experimental measurements, and that computational refinement can improve upon experimental measures.

Currently, the machine learning model underlying our imputation and refinement is re-trained for every new map. Future imputation procedures may benefit by aggregating data from many maps to train a more general imputation model. One challenge will be ensuring that each map is measured on the same scale. For example, the score distributions for missense variants in TPK1 showed a strong bias towards deleteriousness while missense variants in Calmodulin were biased towards neutrality. However, it is unclear whether this reflects intrinsic properties of these genes as opposed to differences in the stringency of the two selection experiments.

As the community carrying out DMS experiments proceeds towards a (perhaps-distant) common goal of generating a DMS map for all human disease-associated proteins, there will be serious challenges that both TileSeq and our computational methods help to address. Previous DMS maps have assessed variation in relatively short polypeptide regions (typically less than 200 amino acids in length). As we approach the median human protein length of ~500 residues, the constraint that we have only ~1-2 missense variants per clones necessarily reduces the mutational density in each clone and thus the allele frequency in the mutagenized population. This will require a substantial increase in the scale of the library of independently transformed/transfected cells and correspondingly-increased sequencing depth to accurately quantify low allele frequencies. Through the reduction in base-calling errors permitted by sequencing both strands, Tile-Seq allows analysis of lower-allele-frequency variants. Similarly, the need to cover all single amino-acid substitutions is

substantially ameliorated by our finding that, given a critical mass of DMS data, the quality of imputed scores is nearly as good as experimental measurements. Each of these improvements could have a major impact on the field of deep mutational scanning.

Genome sequencing is likely to become common in clinical practice. Current estimates suggest that every human carries an average of 100-400 rare variants that have never before been seen in the clinic. DMS meets a critical need for fast, reliable interpretation of variant effects. Instead of generating clones and functionally testing variants of unknown significance after they are first observed, DMS offers exhaustive maps of functional variation that enable interpretation immediately upon clinical presentation, even for rare and personal variation. Our survey of assays revealed that the majority (57%) of human disease genes are potentially already accessible to DMS analysis, so that we may begin to imagine an atlas of DMS maps that reveals pathogenic missense variation for all human disease proteins.

Although our current implementation of the complementation assay uses a temperature-sensitive (ts) allele to provide a control, genes for which no ts allele is yet known are still amenable to DMS by using a null background in combination with an inducibly expressible covering allele. There is even a potential for increasing the number of genes with complementation assays by systematically screening for sensitized backgrounds (exploiting known synthetic lethality relationships or growth conditions). However, growth-based complementation assays have limitations in that they may have limited ability to detect gain- or change-of-function variants. Yeast is also limited as a platform in which to study splicing regulatory or splicing variants. While adaptation of DMS technology to human cell lines will be challenging, recent advances remove some previous hurdles. In addition to the availability of CRISPR-Cas9 to generate homozygous disruptions in target genes, recent advances in 'landing pad' technology (Matreyek *et al*, 2017) now allow transfection and integration of a specific sequence into 1-8% of a population of Hek293T cells. Thus, in theory, DMS could be done using on the order of 1M human cells.

Materials and Methods

POPCode Mutagenesis

The Precision Oligo-Pool based Code Alteration (POPCode) scales up a previous method (Seyfang & Huaqian Jin, 2004). to achieve coverage over the complete spectrum of possible amino acid changes at all protein positions. POPCode requires design of an oligonucleotide centered on each codon in the Open Reading Frame (ORF) of interest, such that the target codon is replaced with an NNK degenerate codon. This has been previously demonstrated to allow all amino acid changes while reducing the chance of generating stop codons (Pal & Fellouse, 2005). Within each mutagenic oligonucleotide, the arm flanking the target codon is varied to achieve a predicted melting temperature that is as uniform as possible to facilitate an even mutation rate across the ORF sequence. We developed a web tool that automates this design step, available online at <http://llama.mshri.on.ca/cgi/popcodeSuite/main>. (See also: Code Availability section).

The POPCode mutagenesis experiment was performed via the following steps: (i) the uracil-containing wild type template was generated by PCR-amplifying the ORF with dNTP/dUTP

mix and HotTaq DNA polymerase, (ii) the mixture of phosphorylated oligonucleotide pool and uracil-containing template was denatured by heating it to 95°C for 3 minutes and then cooled down to 4 degrees to allow the oligos hybridize to the template, (iii) gaps between hybridized oligonucleotides were filled with the non-strand-displacing *Sulpholobus* Polymerase IV (*NEB*) and sealed with T4 DNA ligase (*NEB*), (iv) after degradation of the uracil-doped wild-type strand using Uracil-DNA-Glycosylase (UDG) (*NEB*), the mutant strand was amplified with attB-sites-containing primers and subsequently transferred *en masse* to a donor vector by Gateway BP reaction to generate a library of entry clones.

Synthesis of uracil-containing template. A 50µl PCR reaction contained the following: 1ng template DNA, 1X Taq buffer, 0.2mM dNTPs-dTTP, 0.2mM dUTP, 0.4uM forward and reverse oligos, and 1U Hot Taq Polymerase. Thermal cyclers conditions are as follows: 98°C for 30s, 25 cycles of 98°C for 15s, 60°C for 30s, and 72°C for 1min. A final extension was performed at 72°C for 5 min. Uracilated amplicon was gel-purified using the Minelute gel purification kit (Qiagen).

Phosphorylation of mutagenic oligos. Desalted oligos were purchased from Eurofins or Thermo Scientific. The phosphorylation reaction is as follows: a 50µl reaction containing 1X PNK buffer, 300 pmoles oligos, 1mM ATP, and 10U Polynucleotide Kinase (*NEB*) was incubated at 37°C for 2 hours. The reaction was used directly in the subsequent POPCode reaction.

POPCode oligo annealing and fill-in. A 20µl reaction containing 20ng uracilated DNA, 0.15uM phosphorylated oligo pool, and 1.5uM 5'-oligo was incubated at 95°C for 3 minutes followed by immediate cooling to 4°C. A 30µl reaction containing 1X Taq DNA Ligase buffer, 0.2mM dNTPs, 2U *Sulfolobus* DNA Polymerase IV (*NEB*), and 40U Taq DNA Ligase (*NEB*) was added to the DNA and was incubated at 37°C for 2 hours.

Degradation of wild-type template. 1µl fill-in reaction was added to a 20µl reaction containing 1X UDG buffer and 5U Uracil DNA Glycosylase (*NEB*) and incubated at 37°C for 2 hours.

Amplification of mutagenized DNA. 1µl UDG reaction was added to a 50µl reaction containing 1X Taq buffer, 0.2mM dNTPs, 0.4uM forward and reverse oligos, and 1U Hot Taq Polymerase. Thermal cyclers conditions are as follows: 98°C for 30s, 25 cycles of 98°C for 15s, 60°C for 30s, and 72°C for 1min. A final extension was performed at 72°C for 5 min.

Single-nucleotide mutagenesis

Oxidized nucleotide PCR was performed as previously described by (Mohan *et al*, 2011). Primers were designed to attach attB sites to the product in preparation for Gateway cloning.

Preparation of oxidized nucleotides. A 100µM dNTP mixture was incubated at 37°C with 5mM FeSO₄ for 10 minutes. Addition of 0.5M Mannitol was used to stop the reaction. Oxidized nucleotides were prepared fresh for every PCR reaction.

PCR in presence of oxidized nucleotides. PCR reaction containing: 1-5ng template DNA, 1X Thermopol Buffer (Invitrogen), 1.5mM MgCl₂, 0.2mM dNTP, 0.33µM forward and reverse primers containing attB sites, 1U Taq polymerase was set up during the nucleotide oxidation reaction. Oxidized nucleotides were the last component added to the PCR reaction at a

concentration of 0.1mM (half the amount of regular dNTP). Thermal cycler program: 95°C for 10 min, 30 cycles of 95°C for 1 min, 50°C for 1 min, 72°C for 1 min, final extension at 72°C for 10 min. Mutagenized PCR product was visualised on a 1% agarose gel, and gel-extracted using a gel extraction kit (Qiagen). The gel extracted PCR product is the pooled mutagenesis product carrying attB sites that is carried through to the KiloSeq stage.

Library generation

Generation of mutagenised pool of Entries. An en masse Gateway BP reaction containing 150ng of pooled mutagenesis PCR product carrying attB sites, 150ng of pDONR223, 1μL Gateway BP Clonase II Enzyme Mix (Invitrogen), 1X TE Buffer is prepared. This reaction is incubated overnight at room temperature and then transformed into *E. coli* aiming for the maximum number of transformants (at least 100,000 CFUs) to keep complexity high. Several colonies are picked at this stage for a quality control check by Sanger sequencing, and the rest are put through a pooled DNA extraction. The result is a pool of mutagenised PCR product inserted into the entry vector pDONR223.

Generation of Barcoded Destination Pools. Barcoded destination plasmids were generated as previously reported (Yachie *et al*, 2016), but instead of being arrayed were maintained as pools with high complexity. Briefly, a linear PCR product containing two random 25 nucleotide “barcode” regions flanked by loxP and lox2272 sites along with common linker sequences for priming was combined with a gateway compatible vector at a SacI restriction site through *in vitro* DNA assembly (Gibson *et al*, 2009). This barcoded destination vector pool was transformed into One Shot ccdB Survival T1R Competent Cells (Invitrogen). The transformations were spread onto large round LB+ampicillin petri plates for increased selection capacity and pool complexity was estimated from CFU counts. The plates were combined into a single pool for plasmid DNA extraction by maxiprep.

En masse Gateway LR reaction. An en masse Gateway LR reaction was used to transfer the mutagenised pool of entries into the barcoded destination pool. This reaction takes place over five days. On Day 1 a 5μL reaction containing 150ng of mutagenised ORF pool in pDONR223 backbone, 150ng barcoded pHYC expression vector pool, 1μL LR Clonase II Enzyme Mix, 1X TE buffer is prepared. The reaction is incubated at room temperature overnight. On each of days 2-5 add in a 5μL volume consisting of 150ng barcoded pHYC expression vector, 1μL LR Clonase II Enzyme Mix, 1X TE Buffer, incubating at room temperature overnight each day. On day 5 the final volume is 25μL.

Transformations and colony picking. LR reactions were transformed into *E. coli* and plated to achieve a density of 400-600 individual colonies per plate. A Biomatrix robot (Biomatrix BM5-BC robot, S&P Robotics) was then used to automatically pick and array 384 colonies per plate for a total of ~20,000 clones in ~52 plates per ORF of interest. Each colony at this stage should contain a pHYC expression vector harbouring a variant of the ORF of interest and a unique barcode.

KiloSeq

For the BarSeq method, to establish the identity of each plasmid barcode and its associated

set of mutations in the target ORF we used KiloSeq (Appendix Figure S10) (either carried out in our laboratory or as a service from SeqWell Inc., Beverly, MA). The first step is to PCR-amplify a segment of the plasmid containing both ORF and barcode locus. PCRs were carried out using the Hydrocycler 16 (LGC Group, Ltd.), using primers with well-specific index sequences. Amplicons from each plate were pooled, and subjected to Nextera 'tagmentation' using Tn5 transposase to generate a library of amplicons with random breaks to which the adapters have been ligated. We then re-amplify those fragments to generate a library of amplicons such that one end of each amplicon bears the well-specific tag and the other 'ladder' end bears the Nextera adapter. These libraries can be re-amplified to introduce Illumina TruSeq adaptors, allowing multiple plates of amplicons to be sequenced together. Paired-end sequencing was carried out using Illumina NextSEQ 500. In each pair of reads, one read will reveal the well tag and the barcode locus, whereas the other will contain a fragment of the mutant ORF, and these fragments can be assembled into a contiguous sequence.

To perform demultiplexing, barcode identification and insert resequencing, we developed a sequence analysis pipeline (see Code Availability section). In the first step Illumina bcl2fastq is used to demultiplex the reads at the plate level using the custom Nextera indices. The resulting FASTQ files are then further demultiplexed using the well-tags in a highly parallel fashion. This results in a folder structure containing tens of thousands of individual fastq files sorted by plate and well location. These are then further processed in parallel to identify barcodes. Wells can sometimes contain more than one clone (e.g., due to incomplete washing in the robotic pinning process). Thus barcode sequences are extracted from each read and then clustered by edit distance to determine the set of barcodes in each well. The associated paired reads for each barcodes are then further split by barcode. Each barcode-specific set of ORF reads can then be analyzed with respect to mutations. Bowtie2 software (Langmead & Salzberg, 2012) is used to align reads to the ORF template, PCR duplicates are removed and nucleotide variants called using samtools pileup (Li *et al*, 2009). Given limited read lengths, identification of longer indels is not straightforward. A solution was found by extracting depth of coverage tracks for each clone and normalizing them with respect to average positional coverage across each 384-well plate, applying an edge-detection algorithm to find sudden increases or decreases within normalized coverage, indicating the presence under-covered regions that can arise as a result of insertions or deletions.

After successful genotyping with KiloSeq, we determined the subset of clones that (i) contained a minimum of one missense mutation, (ii) did not contain any insertions or deletions, (iii) did not contain mutations outside of the ORF, (iii) had unique barcodes, (iv) had sufficient read coverage during KiloSeq to allow for confident genotyping. We re-arrayed this filtered subset of clones (Biomatrix BM5-BC robot, S&P Robotics) into a condensed final library of 40 plates containing 6,548 clones.

High-throughput yeast based complementation screen

The yeast based functional assays were established and validated in our previous study (Sun *et al*, 2016). The mutant alleles of the yeast temperature sensitive strains used in this

study are *ubc9-2*, *smt3-331*, *thi80-ph*, and *cmd1-1*. The high-throughput screen was performed as follows: the POPCode generated mutant library was transferred to the expression vector PHYCDEST (Sun *et al*, 2016) by en masse Gateway LR reactions followed by transformation into NEB5 α competent *E. coli* cells (New England Biolabs) and selection for ampicillin resistance.

For the DMS-BarSeq approach, plasmids extracted from a pool of 6,548 barcoded and KiloSeq-validated mutant clones, together with barcoded null and wildtype controls, were transformed into a *S. cerevisiae* strain carrying a temperature-sensitive (ts) allele which can be functionally complemented by the corresponding wild-type human gene (Sun *et al*, 2016). Complexity for this transformation was ~100,000 CFU. For the time series BarSeq screen, the pools were grown separately at both non-selective (25°C) and selective (38°C) temperatures in triplicates to be examined at 5 different timepoints (0h, 6h, 12h, 24h, 48h) yielding 30 samples. For each sample, plasmids were extracted from 10 ODU of cells and used as templates for the downstream barcode PCR amplification. The barcode loci were amplified for each library of plasmids with primers carrying sample-specific tags and then sequenced on an Illumina NextSeq 500.

For the DMS-TileSeq approach, plasmids extracted from a pool of ~100,000 clones were transformed into the corresponding *S. cerevisiae* temperature sensitive strain yielding around 1,000,000 total transformants. Plasmids were prepared from two of 10 ODU of cells and used as templates for the downstream tiling PCR (two replicates of non-selective condition). Two of 40 ODU of cells were inoculated into 200ml medium and grown to full density with shaking at 36°C and plasmids extracted from 10 ODU of each culture were used as templates for the downstream tiling PCR (two replicates of selective condition). In parallel, plasmid expressing the wild-type ORF was transformed to the corresponding *S. cerevisiae* ts strain and grown to full density under the selection. Plasmids were extracted from two of 10 ODU of cells and used as templates for the downstream tiling PCR (two replicates of wild-type control). For each plasmid library, the tiling PCR was performed in two steps: (i) the targeted region of the ORF was amplified with primers carrying a binding site for illumina sequencing adaptors, (ii) each first-step amplicon was indexed with an illumina sequencing adaptor in the second-step PCR. We perform paired end sequencing on the tiled regions across the ORF.

Fitness scoring and refinement

For DMS-BarSeq, a computational pipeline was implemented to identify and count individual sample tags and barcode combinations within each read (see Code Availability section). We then calculate how much better (or worse) each clone grows compared to the pool average, cumulatively across timepoints. To this end, we first calculated the relative population size by dividing each clone's barcode count by the total number of barcodes in each condition. We then calculated the estimated absolute population size for each clone at each timepoint by multiplying the relative population size with the estimated total number of cells on the respective plate at the corresponding time point (obtained from OD measurements). We then treat the hourly growth rate between each individual time point compared to the pool average as an individual estimate of fitness, all of which act cumulatively. Formally, this

corresponds to the following:

Let c_{i,t_k}^τ be the barcode count for clone i , timepoint t_k at temperature τ , then $\forall i \in \{1 \leq i \leq N | i \in \mathbb{N}\}, \forall k \in \{1 \leq k \leq 5 | k \in \mathbb{N}\}, \forall \tau \in \{25^\circ, 37^\circ\}$

$$\begin{aligned} r_{i,t_k}^{(\tau)} &= \frac{c_{i,t_k}^{(\tau)}}{\sum_j c_{j,t_k}^{(\tau)}} \\ P_{i,t_k}^{(\tau)} &= r_{i,t_k}^{(\tau)} \cdot P_{*,t_k}^{(\tau)} \\ \rho_{i,t_k}^{(\tau)} &= \sqrt{(t_k - t_{k-1}) \frac{P_{i,t_k}^{(\tau)}}{P_{i,t_{k-1}}^{(\tau)}}} \\ \phi_{i,t_k}^{(\tau)} &= \frac{\rho_{i,t_k}^{(\tau)}}{\rho_{*,t_k}^{(\tau)}} \\ \phi'_{i,t_k} &= \frac{\phi_{i,t_k}^{(37^\circ)}}{\phi_{*,t_k}^{(25^\circ)}} \\ s_i &= \prod_k \phi'_{i,t_k} \\ s'_i &= \frac{s_i - s_{\text{null}}}{s_{\text{wt}} - s_{\text{null}}}, \end{aligned}$$

Where $r_{i,t_k}^{(\tau)}$ is the relative population size for clone i , timepoint t_k at temperature τ , $P_{i,t_k}^{(\tau)}$ is the absolute population size for clone i , timepoint t_k at temperature τ , $\rho_{i,t_k}^{(\tau)}$ is the measured hourly growth rate for clone i , timepoint t_k at temperature τ , $\phi_{i,t_k}^{(\tau)}$ is the fitness advantage relative to the pool growth for clone i , timepoint t_k at temperature τ , ϕ'_{i,t_k} is the normalized relative fitness advantage for clone i , timepoint t_k , and s_i is the cumulative normalized relative fitness advantage for clone i . Finally, s'_i is the fitness score relative to the internal null and wildtype controls, this results in null-like mutants receiving a score of zero and wildtype-like mutants receiving a score of one.

Given limited amounts of replicates, the empirical standard deviations calculated for each clone or variant can be expected to be imprecise. (Baldi & Long, 2001) have previously described a method for Bayesian regularization or refinement of the standard deviations which yield more robust estimates, leading to less classification error in statistical tests. Briefly, a prior estimate of the standard deviation is computed by linear regression based on the number of barcodes in the permissive condition and the fitness score. The prior is then combined with the empirical value using Baldi and Long's original formula

$$\sigma^2 = \frac{v_n \sigma_n^2}{v_n - 2} = \frac{v_0 \sigma_0^2 + (n - 1) s^2}{v_0 + n - 2}$$

where v_0 represents the degrees of freedom assigned to the prior estimate, σ_0 is the prior estimate resulting according to the regression, n represents the degrees of freedom for the empirical data (i.e. the number of replicates) and s is the empirical standard deviation. The methods were implemented as part of a larger DMS analysis package (see Code Availability)

For DMS-TileSeq, raw sequencing reads were aligned to the reference ORF cDNA sequences using Bowtie-2 (Langmead & Salzberg, 2012) and a custom Perl script was used to parse and compare the forward and reverse read alignment files to count the number of co-occurrences of a codon change in both paired reads. Mutational counts in each condition were normalized to sequencing depth at the respective position. Variants for which the number of reads in the non-permissive condition was within three standard deviations of the read count in the wildtype control were considered poorly measured and removed. Then, the normalized mutational counts from the wild type control libraries (control for sequencing errors) were subtracted from the normalized mutational counts from the non-selective and selective conditions respectively. Finally, the enrichment ratio was calculated for each variant based on the adjusted mutational counts before and after selection.

Re-scaling of fitness metrics

The results from the barcoded and regional sequencing screens do not scale linearly to each other. We used regression to find a monotonic transformation function $f(x) = a \cdot e^x + b \cdot x + c$ between the two screens' respective scales. The standard deviation is transformed accordingly using a Taylor series-based approximation:

$$\sigma' = \sigma \cdot (a \cdot e^\mu + b)$$

. After both datasets have been brought to the same scale we can join corresponding data points using weighted means, where the weight is inversely proportional to the Bayesian regularized standard error. Output standard error was adjusted to account for differences in input fitness values and increased sample size:

$$\begin{aligned} w_0 &= \frac{1}{1 + \frac{\sigma_{\bar{x}}^{(0)}}{\sigma_{\bar{x}}^{(1)}}}; \quad w_1 = \frac{1}{1 + \frac{\sigma_{\bar{x}}^{(1)}}{\sigma_{\bar{x}}^{(0)}}} \\ \mu_{\text{joint}} &= w_0 \cdot \mu_0 + w_1 \cdot \mu_1 \\ \sigma_{\text{joint}}^2 &= w_0 \cdot (\sigma_0^2 + \mu_0^2) + w_1 \cdot (\sigma_1^2 + \mu_1^2) - \mu_{\text{joint}}^2 \\ \sigma_{\bar{x}}^{(\text{joint})} &= \frac{\sigma_{\text{joint}}}{\sqrt{df_0 + df_1}} \end{aligned}$$

where μ_0 is the DMS-BarSeq value, σ_0 the associated standard deviation, $\sigma_{\bar{x}}^{(0)}$ the associated standard error, df_0 the associated degrees of freedom, μ_1 is the DMS-TileSeq value, σ_1 the associated standard deviation, $\sigma_{\bar{x}}^{(1)}$ the associated standard error, and df_1 the associated degrees of freedom. These steps were implemented as part of a larger DMS analysis package (see Code Availability).

Imputation of missing data

Next we aimed to find a machine learning method that would allow us to input the missing parts of the map. The first step towards this was to gather suitable features. We first evaluated the most promising features using linear regression and then applied a random forest model using all the available features.

The most important features were intrinsic, i.e. directly derived from unused information in the screen. These are: The average fitness across variants at the same position; The average fitness of multi-mutant clones that contain the variant of interest; the estimated fitness according to a multiplicative model to infer mutant fitness A using a double mutant AB and single mutant B. Another set of features was computed from differences between various chemical properties of the wildtype and mutant amino acids. These properties include size, volume, polarity, charge, hydropathy. A third set of features is derived from the structural context of each amino acid position. This includes secondary structure, solvent accessibility, burial in interfaces with different interaction partners and involvement in hydrogen bonds or salt bridges with interaction partners. Secondary structures were calculated using Stride (Frishman & Argos, 1995). Solvent accessibility and interface burial were calculated using the GETAREA tool (Fraczkiewicz & Braun, 1998) on the following PDB entries: For UBE2I: 3UIP (Gareau *et al*, 2012); 4W5V (Boucher *et al*. unpublished) ; 3KYD (Olsen *et al*, 2010); 2UYZ (Knipscheer *et al*, 2007); 4Y1L (Alontaga *et al*, 2015). For SUMO1: 2G4D (Xu *et al*, 2006); 2IO2 (Reverter & Lima, 2006); 3KYD (Olsen *et al*, 2010); 3UIP (Gareau *et al*, 2012); 2ASQ (Song *et al*, 2005); 4WJO (Cappadocia *et al*, 2015); 4WJQ (Cappadocia *et al*, 2015); 1WYW (Baba *et al*, 2005). For calmodulin: 3G43 (Fallon *et al*, 2009); 4DJC (Sarhan *et al*, 2012). And for TPK1: 3S4Y (Baker *et al*, 2001)

Hydrogen bond and salt bridge candidates were predicted using OpenPyMol and evaluated for validity by manual inspection. Additional features used are the BLOSUM score for a given amino acid change, the PROVEAN score, and the evolutionary conservation of the amino acid position. Conservation was obtained by generating a multiple alignment of direct functional orthologues across many eukaryotic species using CLUSTAL (Sievers & Higgins, 2014), which was used as input for AMAS (Livingstone & Barton, 1993). We then applied the complete set of features in a random forest model using the R package randomForest (Breiman, 2001) version 4.6-12 with the default settings for all hyperparameters ($n_{\text{tree}}=500$, $m_{\text{try}}=n_{\text{feat}}/3$, $\text{replace}=\text{TRUE}$, $\text{sampsize}=n_{\text{obs}}$, $\text{nodesize}=5$, $\text{maxnodes} = \text{NULL}$, $n_{\text{Perm}}=1$). These procedures were implemented as part of a larger DMS analysis package (see Code Availability section).

Refinement of low-confidence measurements

The machine-learning predictions resulting generated above can also be used to refine experimental measurements of lower confidence. To this end, the corrected standard error associated with each datapoint can be used to determine the weight of assigned to the measurement.

$$w_0 = \frac{1}{1 + \frac{\sigma_{\bar{x}}^{(0)}}{\sigma_{\bar{x}}^{(1)}}}; w_1 = \frac{1}{1 + \frac{\sigma_{\bar{x}}^{(1)}}{\sigma_{\bar{x}}^{(0)}}}$$

$$\mu_{\text{joint}} = w_0 \cdot \mu_0 + w_1 \cdot \mu_1$$

$$\sigma_{\text{joint}}^2 = w_0 \cdot (\sigma_0^2 + \mu_0^2) + w_1 \cdot (\sigma_1^2 + \mu_1^2) - \mu_{\text{joint}}^2$$

$$\sigma_{\bar{x}}^{(\text{joint})} = \frac{\sigma_{\text{joint}}}{\sqrt{df_0 + df_1}}$$

Where μ_0 is the measured value, σ_0 the associated standard deviation, $\sigma_{\bar{x}}^{(0)}$ the associated standard error, df_0 the associated degrees of freedom, μ_1 is the Random Forest-predicted value, σ_1 the associated standard deviation as approximated by cross-validation RMSD, $\sigma_{\bar{x}}^{(1)}$ the associated standard error and df_1 the associated virtual degrees of freedom. The methods were implemented as part of a larger DMS analysis package (see Code Availability section)

Experimental validation by yeast spotting assays

To validate the reliability of the fitness scores obtained during the screen, we selected three subsets of clones from our original UBE2I variant library: (1) A set of clones carrying variants with functional scores representing the full spectrum in the screen; (2) A set of clones carrying hypercomplementing variants in the screen; and (3) A set of clones carrying variants not present in the imputation training data set. After genotype verification using Sanger sequencing, each variant was transferred to the yeast expression plasmid pHYCDEST by Gateway technology and individually transformed into the yeast ts mutant strain. Cells were grown to saturation in 96-well cell culture plates at room temperature. Each culture was then adjusted to an OD600 of 1.0 and serially diluted to 5^{-1} , 5^{-2} , 5^{-3} , 5^{-4} , 5^{-5} , and 5^{-6} . These cultures (5 μ l of each) were then spotted on SC-LEU plates as appropriate to maintain the plasmid and incubated at either the permissive or non-permissive temperatures for two days. Each variant was assayed alongside negative and positive controls for loss of complementation (expression of either the wild type human protein or a GFP control). Results were interpreted by comparing the growth difference between the yeast strains expressing human genes and the corresponding control strain expressing the GFP gene.

Quantification of spotting assay images (for Appendix Figure S3) was performed as follows: Using blinded manual inspection, the following scores were assigned: 0 – No colonies visible; 0.25 – Colonies visible up to the first dilution; 0.5 – Colonies visible up to the second dilution; 0.75 – Colonies visible up to the third dilution; 1 – Colonies visible up to the fourth dilution (a value of 1 was chosen as this corresponds to growth in the wild-type control); 1.25 – Colonies visible up to the fifth dilution; 1.5 – Colonies visible up to the sixth dilution.

Assessing relationship of hyperactive complementation to reversion

To examine whether changing amino acid residues into those residues naturally occur in yeast were more likely to show hyperactive complementation we compared these cases to

changes into residues occurring in other species. The UBE2I amino acid sequence was aligned to that of its orthologues in *S. cerevisiae*, *D. discoideum* and *D. melanogaster* using CLUSTAL (Sievers & Higgins, 2014). A custom script was used to extract inter-species amino acid changes and lookup the corresponding complementation fitness values in the UBE2I map. Distributions were plotted using the R package beeswarm (Eklund, 2016). The methods were implemented as part of a larger DMS analysis package (see Code Availability section).

In vitro sumoylation comparison

Images from in vitro sumoylation assays performed for UBE2I variants by (Bernier-Villamor *et al*, 2002) were scored by visual inspection while blinded to the underlying variant information. Scores were then represented as a heatmap and compared complementation scores from the UBE2I map. The methods were implemented as part of a larger DMS analysis package provided and also available online at <https://bitbucket.org/rothlabto/dmspipeline>.

Phylogenetic comparison of different models for hypercomplementation

We used the phydms software package (Bloom, 2017) to test three different models relating the effect of complementation-enhancing substitutions in SUMO1 and UBE2I to actual preference for the substituted amino acid in a real biological context. Specifically, using the substitution models described in (Bloom, 2017), we tested three different ways of relating the evolutionary preference $\pi_{r,a}$ for amino-acid *a* at site *r* to the fitness score $f_{r,a}$ for this variant. In the first model, $\pi_{r,a} = f_{r,a}$. In the second model, $\pi_{r,a} = \min(f_{r,a}, f_{r,wt})$ where $f_{r,wt}$ is the fitness score for the wildtype amino-acid at site *r*. In the third model, $\pi_{r,a} = f_{r,a}$, if $f_{r,a} \leq f_{r,wt}$ and $1/f_{r,a}$ otherwise. We fit each of these models to the set of Ensembl homologs with at least 75% sequence identity to the human protein. As shown in Appendix Table S1, in all cases the last model (which assigns low preference to variants that strongly enhance activity) best fits the sequences. The computer code that performs this analysis is available on GitHub at https://github.com/jbloomlab/AtlasPaper_SUMO1_UBE2I_ExpCM

Statistical details

Figure 1C: Error bars show Bayesian regularized standard error based on three technical replicates and a prior based on pre-selection counts and final score (see subsection on score calculation for details).

Figure 1D: As normality cannot be assumed for the distributions of fitness scores, one-sided two-sample Wilcoxon-Mann-Whitney tests were used. Low conservation (n=60 clones) vs Medium Conservation (n=105 clones) $W = 3789$, $P = 0.015$; Medium Conservation (n=105 clones) vs High Conservation (n=404 clones) $W = 28043$, $P = 1.8 \times 10^{-7}$; Core (n=208 clones) vs Surface (n=42 clones) $W = 1649$, $P = 1.01 \times 10^{-10}$; Interface (n=215 clones) vs Surface (n=42 clones) $W = 2461$, $P = 1.58 \times 10^{-6}$.

Figure 5A: As normality cannot be assumed for the distributions of fitness scores, a one-

sided two-sample Wilcoxon-Mann-Whitney test was used: $n=\{26,31\}$ variants, $W=570.5$, $P=3.73 \times 10^{-3}$.

Code and data availability

All code associated with this work can be checked out using mercurial from the following repositories: (1) For the KiloSeq analysis pipeline: <https://bitbucket.org/rothlabto/kiloseq>; (2) for the popcode oligo design tool: <https://bitbucket.org/rothlabto/popcodesuite>; (3) For the BarSeq sequence analysis pipeline: <https://bitbucket.org/rothlabto/screenpipeline>; (4) For the TileSeq sequence analysis pipeline: https://bitbucket.org/rothlabto/tileseq_package For all raw data and downstream analyses: <https://bitbucket.org/rothlabto/dmspipeline>. Raw sequencing data can be obtained from the NCBI Short Read Archive, accession numbers SRP109101 (KiloSeq) and SRP109119 (DMS screens). All final variant maps and associated data tables can be downloaded at <http://dalai.mshri.on.ca/~jweile/projects/dmsData/>. Original data for the in vitro sumoylation analysis can be found in (Bernier-Villamor *et al*, 2002).

Acknowledgements

The authors thank Amy Caudy, Lincoln Stein, Igor Stagljar, Chris Lima and Brian Raught for their advice, and thank Brenda Andrews and Charles Boone for kindly providing temperature-sensitive yeast mutant strains.

The authors gratefully acknowledge funding by the National Human Genome Research Institute of the National Institutes of Health (NIH/NHGRI) Center of Excellence in Genomic Science (CEGS) Initiative (HG004233), the Canadian Excellence Research Chairs (CERC) Program, and the Ontario Ministry of Research and Innovation (MRI).

Author contributions

FPR, JW, SS and AGC conceived the project; SS, AGC, JK, MV and CW performed the DMS experiments and manual validations; JM, MT and FR conceived the KiloSeq method, AGC, JK and JM performed KiloSeq, JW, SS and NL developed the analysis pipeline, YW and JW developed the machine learning imputation and refinement method with advice from DF; JW, CP and PA performed structural and epistasis analyses; SS and FY curated the list of assayable genes; JB performed the evolution analysis; SY and BN helped conduct the blind test with Invitae variant data; GT constructed ts strains; DEH and MV provided human clones; and JW, FPR and SS wrote the manuscript. FPR supervised the project.

Conflicts of interest

FPR is a shareholder and scientific advisory board member of SeqWell Inc. and of Ranomics, Inc. RN and SY are employees of Invitae, Inc.

References

- Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P, Kondrashov AS & Sunyaev SR (2010) A method and server for predicting damaging missense mutations. *Nat. Methods* **7**: 248–249
- Alontaga AY, Ambaye ND, Li Y-J, Vega R, Chen C-H, Bzymek KP, Williams JC, Hu W & Chen Y (2015) RWD Domain as an E2 (Ubc9)-Interaction Module. *J. Biol. Chem.* **290**: 16550–16559
- Andreou AM, Pauws E, Jones MC, Singh MK, Bussen M, Doudney K, Moore GE, Kispert A, Brosens JJ & Stanier P (2007) TBX22 missense mutations found in patients with X-linked cleft palate affect DNA binding, sumoylation, and transcriptional repression. *Am. J. Hum. Genet.* **81**: 700–712
- Baba D, Maita N, Jee J-G, Uchimura Y, Saitoh H, Sugasawa K, Hanaoka F, Tochio H, Hiroaki H & Shirakawa M (2005) Crystal structure of thymine DNA glycosylase conjugated to SUMO-1. *Nature* **435**: 979–982
- Baker LJ, Dorocke JA, Harris RA & Timm DE (2001) The crystal structure of yeast thiamin pyrophosphokinase. *Structure* **9**: 539–546
- Baldi P & Long AD (2001) A Bayesian framework for the analysis of microarray expression data: regularized t -test and statistical inferences of gene changes. *Bioinformatics* **17**: 509–519
- Bernier-Villamor V, Sampson DA, Matunis MJ & Lima CD (2002) Structural Basis for E2-Mediated SUMO Conjugation Revealed by a Complex between Ubiquitin-Conjugating Enzyme Ubc9 and RanGAP1. *Cell* **108**: 345–356
- Blomen VA, Májek P, Jae LT, Bigenzahn JW, Nieuwenhuis J, Staring J, Sacco R, Diemen FR van, Olk N, Stukalov A, Marceau C, Janssen H, Carette JE, Bennett KL, Colinge J, Superti-Furga G & Brummelkamp TR (2015) Gene essentiality and synthetic lethality in haploid human cells. *Science*: aac7557
- Bloom JD (2014) An Experimentally Determined Evolutionary Model Dramatically Improves Phylogenetic Fit. *Mol. Biol. Evol.* **31**: 1956–1978
- Bloom JD (2017) Identification of positive selection in genes is greatly improved by using experimentally informed site-specific models. *Biol. Direct* **12**: 1
- Breiman L (2001) Random Forests. *Mach. Learn.* **45**: 5–32
- Cappadocia L, Mascle XH, Bourdeau V, Tremblay-Belzile S, Chaker-Margot M, Lussier-Price M, Wada J, Sakaguchi K, Aubry M, Ferbeyre G & Omichinski JG (2015) Structural and Functional Characterization of the Phosphorylation-Dependent Interaction between PML and SUMO1. *Structure* **23**: 126–138
- Choi Y, Sims GE, Murphy S, Miller JR & Chan AP (2012) Predicting the Functional Effect of Amino Acid Substitutions and Indels. *PLOS ONE* **7**: e46688
- Crotti L, Johnson CN, Graf E, Ferrari GMD, Cuneo BF, Ovadia M, Papagiannis J, Feldkamp MD, Rathi SG, Kunic JD, Pedrazzini M, Wieland T, Lichtner P, Beckmann B-M, Clark T, Shaffer C, Benson DW, Käåb S, Meitinger T, Strom TM, et al (2013) Calmodulin Mutations Associated with Recurrent Cardiac Arrest in Infants. *Circulation* **127**: 1009–1017
- Doud MB & Bloom JD (2016) Accurate Measurement of the Effects of All Amino-Acid Mutations on Influenza Hemagglutinin. *Viruses* **8**: 155
- Eklund A (2016) The Bee Swarm Plot, an Alternative to Stripchart Available at: <https://CRAN.R-project.org/package=beeswarm>
- Fallon JL, Baker MR, Xiong L, Loy RE, Yang G, Dirksen RT, Hamilton SL & Quiocho FA (2009) Crystal structure of dimeric cardiac L-type calcium channel regulatory domains bridged by Ca²⁺-calmodulins. *Proc. Natl. Acad. Sci.* **106**: 5135–5140
- Forbes S a., Beare D, Bindal N, Bamford S, Ward S, Cole C g., Jia M, Kok C, Boutselakis H, De T, Sondka Z, Ponting L, Stefancsik R, Harsha B, Tate J, Dawson E, Thompson S, Jubb H & Campbell P j. (2001) COSMIC: High-Resolution Cancer Genetics Using the

- Catalogue of Somatic Mutations in Cancer. In *Current Protocols in Human Genetics* John Wiley & Sons, Inc. Available at: <http://onlinelibrary.wiley.com/doi/10.1002/cphg.21/abstract> [Accessed February 5, 2017]
- Fowler DM, Araya CL, Fleishman SJ, Kellogg EH, Stephany JJ, Baker D & Fields S (2010) High-resolution mapping of protein sequence-function relationships. *Nat. Methods* **7**: 741–746
- Fowler DM & Fields S (2014) Deep mutational scanning: a new style of protein science. *Nat. Methods* **11**: 801–807
- Fraczkiewicz R & Braun W (1998) Exact and efficient analytical calculation of the accessible surface areas and their gradients for macromolecules. *J. Comput. Chem.* **19**: 319–333
- Frishman D & Argos P (1995) Knowledge-based protein secondary structure assignment. *Proteins Struct. Funct. Bioinforma.* **23**: 566–579
- Gareau JR, Reverter D & Lima CD (2012) Determinants of Small Ubiquitin-like Modifier 1 (SUMO1) Protein Specificity, E3 Ligase, and SUMO-RanGAP1 Binding Activities of Nucleoporin RanBP2. *J. Biol. Chem.* **287**: 4740–4751
- Gibson DG, Young L, Chuang R-Y, Venter JC, Hutchison CA & Smith HO (2009) Enzymatic assembly of DNA molecules up to several hundred kilobases. *Nat. Methods* **6**: 343–345
- Hamza A, Tammperre E, Kofoed M, Keong C, Chiang J, Giaever G, Nislow C & Hieter P (2015) Complementation of Yeast Genes with Human Genes as an Experimental Platform for Functional Testing of Human Genetic Variants. *Genetics* **201**: 1263–1274
- Hart T, Chandrashekar M, Aregger M, Steinhart Z, Brown KR, MacLeod G, Mis M, Zimmermann M, Fradet-Turcotte A, Sun S, Mero P, Dirks P, Sidhu S, Roth FP, Rissland OS, Durocher D, Angers S & Moffat J (2015) High-Resolution CRISPR Screens Reveal Fitness Genes and Genotype-Specific Cancer Liabilities. *Cell* **163**: 1515–1526
- Henikoff S & Henikoff JG (1992) Amino acid substitution matrices from protein blocks. *Proc. Natl. Acad. Sci.* **89**: 10915–10919
- Jiang W & Koltin Y (1996) Two-hybrid interaction of a human UBC9 homolog with centromere proteins of *Saccharomyces cerevisiae*. *Mol. Gen. Genet. MGG* **251**: 153–160
- Kachroo AH, Laurent JM, Yellman CM, Meyer AG, Wilke CO & Marcotte EM (2015) Systematic humanization of yeast genes reveals conserved functions and genetic modularity. *Science* **348**: 921–925
- Knipscheer P, Dijk WJ van, Olsen JV, Mann M & Sixma TK (2007) Noncovalent interaction between Ubc9 and SUMO promotes SUMO chain formation. *EMBO J.* **26**: 2797–2807
- Landrum MJ, Lee JM, Benson M, Brown G, Chao C, Chitipiralla S, Gu B, Hart J, Hoffman D, Hoover J, Jang W, Katz K, Ovetsky M, Riley G, Sethi A, Tully R, Villamarin-Salomon R, Rubinstein W & Maglott DR (2016) ClinVar: public archive of interpretations of clinically relevant variants. *Nucleic Acids Res.* **44**: D862–D868
- Langmead B & Salzberg SL (2012) Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**: 357–359
- Lee MG & Nurse P (1987) Complementation used to clone a human homologue of the fission yeast cell cycle control gene *cdc2*. *Nature* **327**: 31–35
- Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, O'Donnell-Luria AH, Ware JS, Hill AJ, Cummings BB, Tukiainen T, Birnbaum DP, Kosmicki JA, Duncan LE, Estrada K, Zhao F, Zou J, Pierce-Hoffman E, Berghout J, Cooper DN, et al (2016) Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**: 285–291
- Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G & Durbin

- R (2009) The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**: 2078–2079
- Livingstone CD & Barton GJ (1993) Protein sequence alignments: a strategy for the hierarchical analysis of residue conservation. *Bioinformatics* **9**: 745–756
- Majithia AR, Tsuda B, Agostini M, Gnanapradeepan K, Rice R, Peloso G, Patel KA, Zhang X, Broekema MF, Patterson N, Duby M, Sharpe T, Kalkhoven E, Rosen ED, Barroso I, Ellard S, UK Monogenic Diabetes Consortium, Kathiresan S, Myocardial Infarction Genetics Consortium, O'Rahilly S, et al (2016) Prospective functional classification of all possible missense variants in PPARG. *Nat. Genet.* **48**: 1570–1575
- Matreyek KA, Stephany JJ & Fowler DM (2017) A platform for functional assessment of large variant libraries in mammalian cells. *Nucleic Acids Res.* **45**: e102–e102
- Maxwell KN, Hart SN, Vijai J, Schrader KA, Slavin TP, Thomas T, Wubbenhorst B, Ravichandran V, Moore RM, Hu C, Guidugli L, Wenz B, Domchek SM, Robson ME, Szabo C, Neuhausen SL, Weitzel JN, Offit K, Couch FJ & Nathanson KL (2016) Evaluation of ACMG-Guideline-Based Variant Classification of Cancer Susceptibility and Non-Cancer-Associated Genes in Families Affected by Breast Cancer. *Am. J. Hum. Genet.* **98**: 801–817
- Mayr JA, Freisinger P, Schlachter K, Rolinski B, Zimmermann FA, Scheffner T, Haack TB, Koch J, Ahting U, Prokisch H & Sperl W (2011) Thiamine Pyrophosphokinase Deficiency in Encephalopathic Children with Defects in the Pyruvate Oxidation Pathway. *Am. J. Hum. Genet.* **89**: 806–812
- Mohan U, Kaushik S & Banerjee UC (2011) PCR Based Random Mutagenesis Approach for a Defined DNA Sequence Using the Mutagenic Potential of Oxidized Nucleotide Products. *Open Biotechnol. J.* **5**: 21–27
- Ng PC & Henikoff S (2001) Predicting Deleterious Amino Acid Substitutions. *Genome Res.* **11**: 863–874
- Nyegaard M, Overgaard MT, Søndergaard MT, Vranas M, Behr ER, Hildebrandt LL, Lund J, Hedley PL, Camm AJ, Wettrell G, Fosdal I, Christiansen M & Børghlum AD (2012) Mutations in calmodulin cause ventricular tachycardia and sudden cardiac death. *Am. J. Hum. Genet.* **91**: 703–712
- Olsen SK, Capili AD, Lu X, Tan DS & Lima CD (2010) Active site remodelling accompanies thioester bond formation in the SUMO E1. *Nature* **463**: 906–912
- Osborn MJ & Miller JR (2007) Rescuing yeast mutants with human genes. *Brief. Funct. Genomics* **6**: 104–111
- Pal G & Fellouse FA (2005) Methods for the Construction of Phage-Displayed Libraries. In *Phage Display In Biotechnology and Drug Discovery* pp 111–142. CRC Press Available at: <http://www.crcnetbase.com/doi/abs/10.1201/9780849359125.ch3> [Accessed February 5, 2017]
- Reverter D & Lima CD (2006) Structural basis for SENP2 protease interactions with SUMO precursors and conjugated substrates. *Nat. Struct. Mol. Biol.* **13**: 1060–1068
- Rolland T, Taşan M, Charlotiaux B, Pevzner SJ, Zhong Q, Sahni N, Yi S, Lemmens I, Fontanillo C, Mosca R, Kamburov A, Ghiassian SD, Yang X, Ghamsari L, Balcha D, Begg BE, Braun P, Brehme M, Broly MP, Carvunis A-R, et al (2014) A Proteome-Scale Map of the Human Interactome Network. *Cell* **159**: 1212–1226
- Sahni N, Yi S, Taipale M, Fuxman Bass JL, Coulombe-Huntington J, Yang F, Peng J, Weile J, Karras GI, Wang Y, Kovács IA, Kamburov A, Krykbaeva I, Lam MH, Tucker G, Khurana V, Sharma A, Liu Y-Y, Yachie N, Zhong Q, et al (2015) Widespread Macromolecular Interaction Perturbations in Human Genetic Disorders. *Cell* **161**: 647–660
- Sarhan MF, Tung C-C, Petegem FV & Ahern CA (2012) Crystallographic basis for calcium regulation of sodium channels. *Proc. Natl. Acad. Sci.* **109**: 3558–3563
- Seyfang A & Huaqian Jin J (2004) Multiple site-directed mutagenesis of more than 10 sites simultaneously and in a single round. *Anal. Biochem.* **324**: 285–291
- Sievers F & Higgins D (2014) Clustal Omega, Accurate Alignment of Very Large Numbers of

- Sequences. In *Multiple Sequence Alignment Methods*, Russell DJ (ed) pp 105–116. Humana Press Available at: http://dx.doi.org/10.1007/978-1-62703-646-7_6 [Accessed February 5, 2017]
- Song J, Zhang Z, Hu W & Chen Y (2005) Small Ubiquitin-like Modifier (SUMO) Recognition of a SUMO Binding Motif --- A reversal of the bound orientation. *J. Biol. Chem.* **280**: 40122–40129
- St Onge RP, Mani R, Oh J, Proctor M, Fung E, Davis RW, Nislow C, Roth FP & Giaever G (2007) Systematic pathway analysis using high-resolution fitness profiling of combinatorial gene deletions. *Nat. Genet.* **39**: 199–206
- Starita LM, Young DL, Islam M, Kitzman JO, Gullingsrud J, Hause RJ, Fowler DM, Parvin JD, Shendure J & Fields S (2015) Massively Parallel Functional Analysis of BRCA1 RING Domain Variants. *Genetics*: genetics.115.175802
- Sun S, Yang F, Tan G, Costanzo M, Oughtred R, Hirschman J, Theesfeld CL, Bansal P, Sahni N, Yi S, Yu A, Tyagi T, Tie C, Hill DE, Vidal M, Andrews BJ, Boone C, Dolinski K & Roth FP (2016) An extended set of yeast-based functional assays accurately identifies human disease mutations. *Genome Res.* **26**: 670–680
- The 1000 Genomes Project Consortium (2015) A global reference for human genetic variation. *Nature* **526**: 68–74
- Wang T, Wei JJ, Sabatini DM & Lander ES (2014) Genetic Screens in Human Cells Using the CRISPR-Cas9 System. *Science* **343**: 80–84
- Whitehead TA, Chevalier A, Song Y, Dreyfus C, Fleishman SJ, De Mattos C, Myers CA, Kamisetty H, Blair P, Wilson IA & Baker D (2012) Optimization of affinity, specificity and function of designed influenza inhibitors using deep sequencing. *Nat. Biotechnol.* **30**: 543–548
- Xu Z, Chau SF, Lam KH, Chan HY, Ng TB & Au SWN (2006) Crystal structure of the SENP1 mutant C603S–SUMO complex reveals the hydrolytic mechanism of SUMO-specific protease. *Biochem. J.* **398**: 345–352
- Yachie N, Petsalaki E, Mellor JC, Weile J, Jacob Y, Verby M, Ozturk SB, Li S, Cote AG, Mosca R, Knapp JJ, Ko M, Yu A, Gebbia M, Sahni N, Yi S, Tyagi T, Sheykhkarimli D, Roth JF, Wong C, et al (2016) Pooled-matrix protein interaction screens using Barcode Fusion Genetics. *Mol. Syst. Biol.* **12**: 863
- Zhang T-H, Wu NC & Sun R (2016) A benchmark study on error-correction by read-pairing and tag-clustering in amplicon-based deep sequencing. *BMC Genomics* **17**: Available at: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4751728/> Available at: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4751728/>

Figure legends

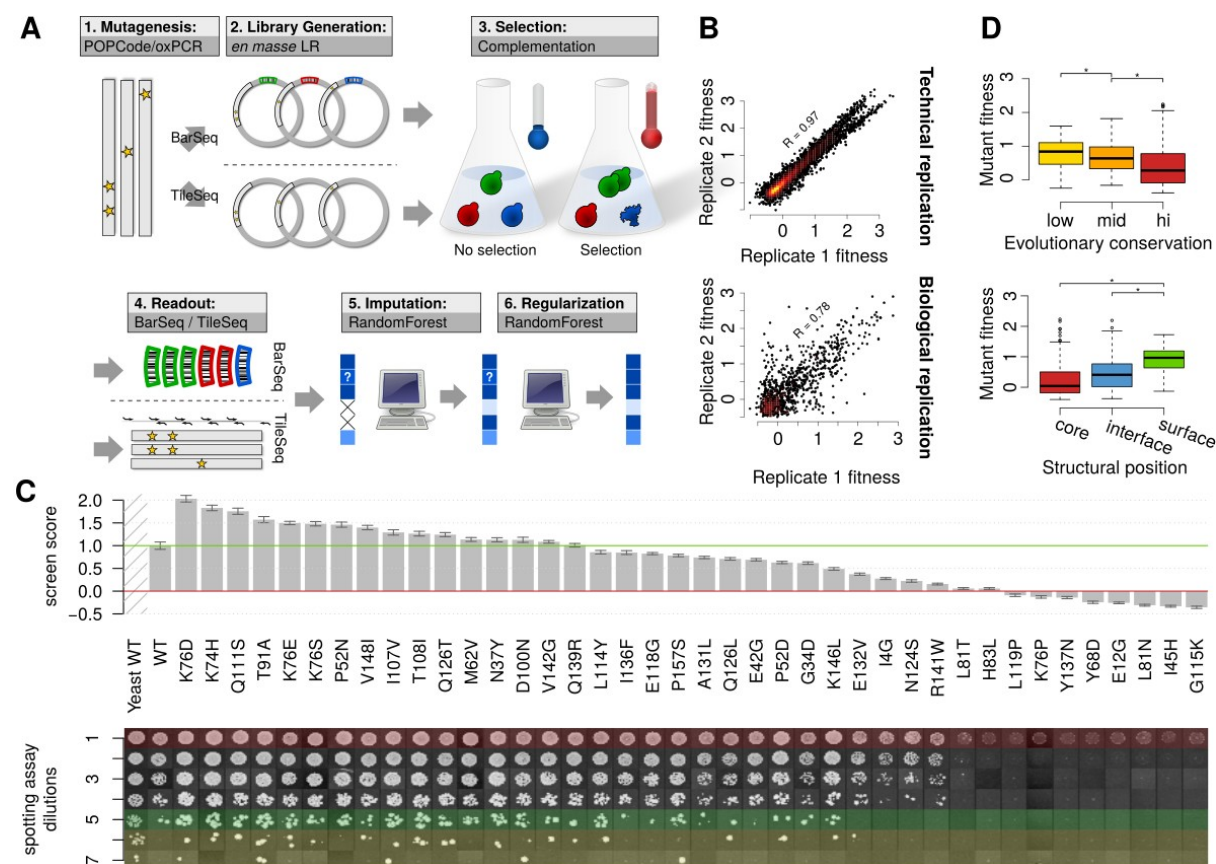
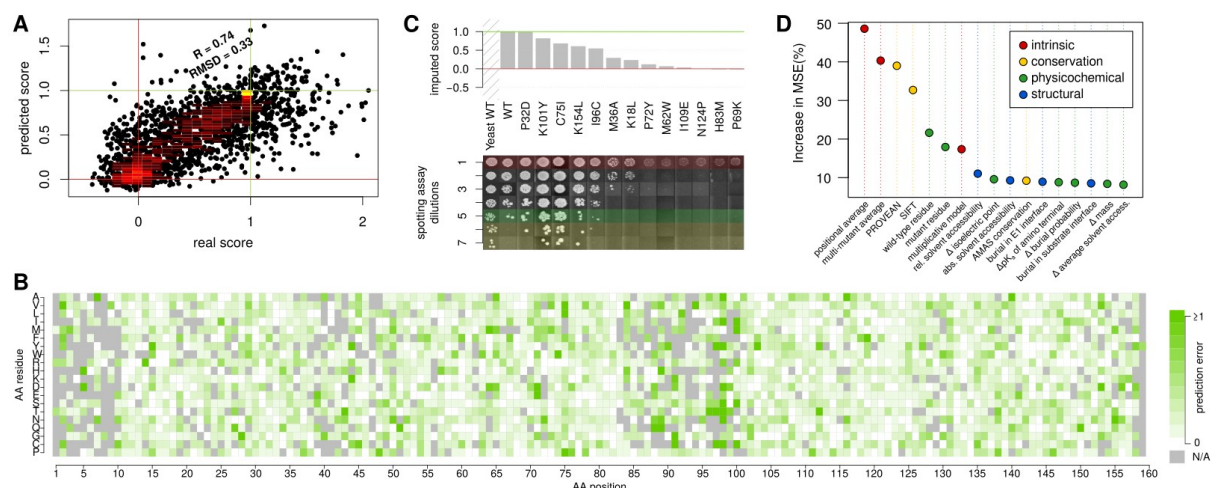


Figure 1: UBE2I screening and validation. (A) Modular structure of the screening framework. (B) Raw DMS-BarSeq fitness scores in technical replicates (separately plated assays of the same pool) and biological replicates (separate sub-strains in the pool carrying the same variants). (C) Manual spotting assay validation of a representative set of variants. Each row represents a consecutive 5-fold dilution. Marked in red: Maximal dilution visible in empty vector control. Marked in green: Maximal dilution with visible human wt control. Marked in yellow: Dilution steps exceeding visible human wt control. Bar heights represent summary screen scores. Error bars indicate bayesian refined s.e.m. (D) Variants grouped by evolutionary conservation (AMAS score) of their respective sites (top) and grouped by structural context within the protein core, within protein-protein interaction interfaces or on remaining protein surface (bottom). See Materials and Methods for statistical details.



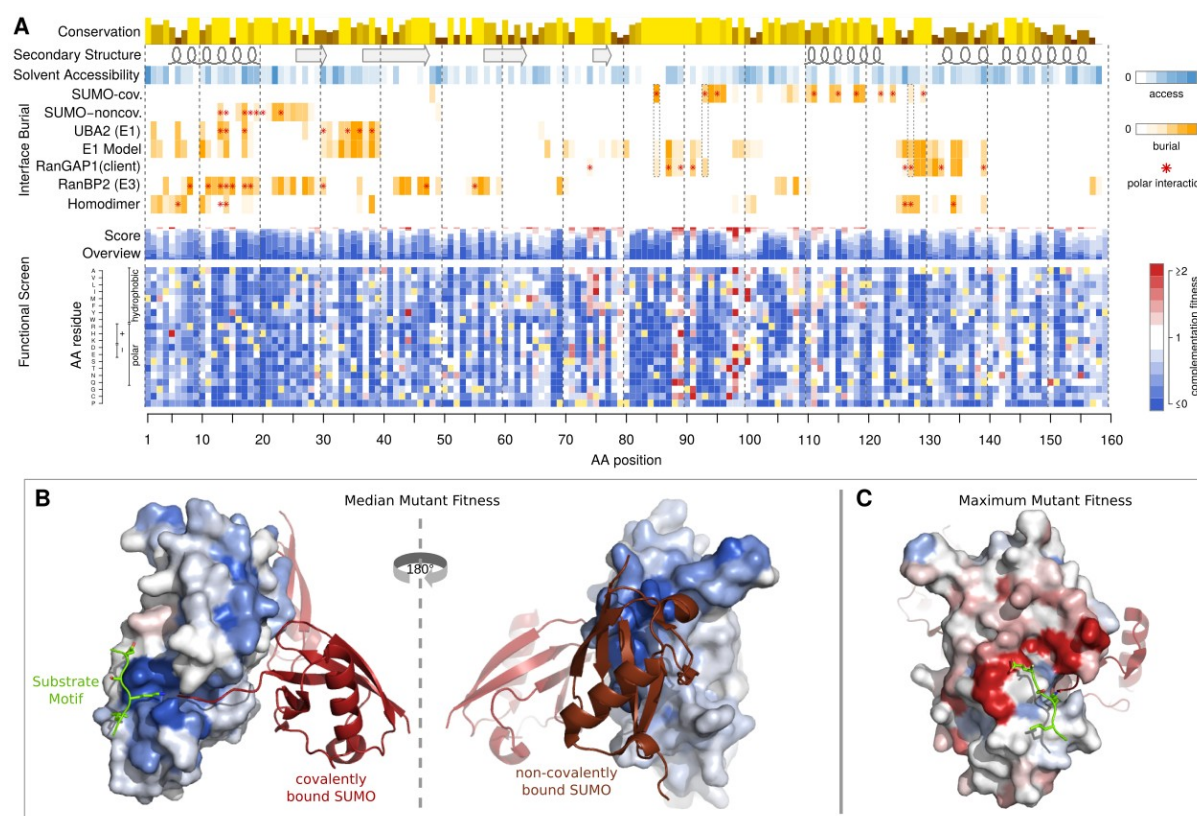


Figure 3: (A) A complete functional map of UBE2I as resulting from the combination of the complementation screen and machine learning imputation and refinement. An impact score of 0 (blue) corresponds to a fitness equivalent to the empty vector control. A score of 1 (white) corresponds to a fitness equivalent to the wildtype control. A score greater than 1 (red) corresponds to fitness above wildtype levels. Shown above, for comparison are sequence conservation, secondary structure, solvent accessibility, and burial of the respective amino acid in protein-protein interaction interfaces with covalently and non-covalently bound SUMO, the E1 UBA2, the sumoylation target RanGAP1, the E3 RanBP2 and UBE2I itself. Hydrogen bonds or salt bridges between residues and the respective interaction partner are marked with red asterisks. Residues buried in both the covalent SUMO and client interfaces are framed with dotted lines, marking the core members of the active site. (B) UBE2I crystal structure with residues colored according to the median mutant fitness. Colors as in A. The interacting substrate's ΨKxE motif is shown in green stick model; Covalently bound SUMO is shown as a red cartoon model; and non-covalently bound SUMO is shown in brown cartoon model. The structures shown were obtained by alignment of PDB entries 3UIP and 2PE6. (C) UBE2I crystal structure as in B, with residues colored according to maximum mutant fitness.

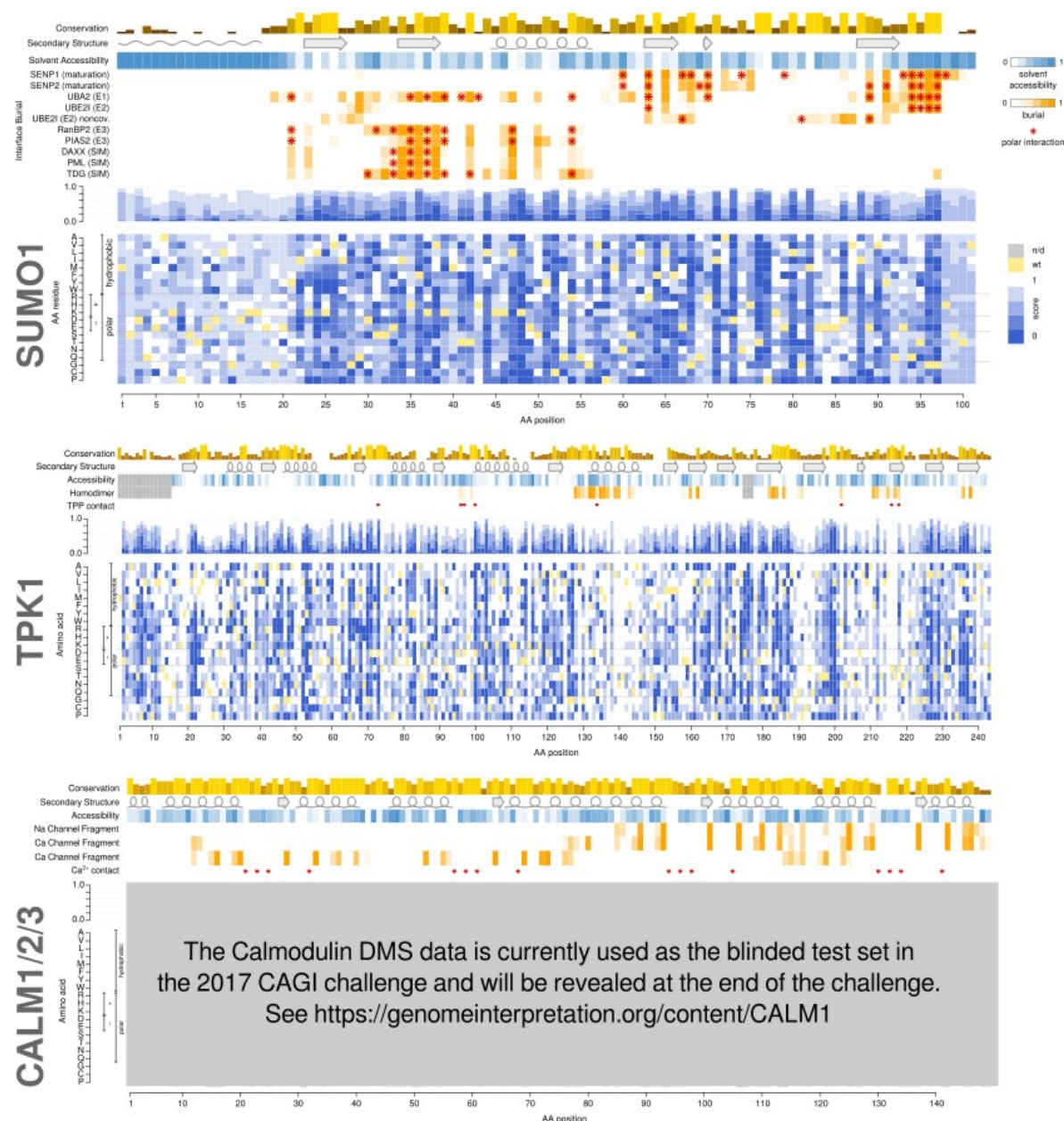


Figure 4: Functional maps of SUMO1, TPK1 and calmodulin (CALM1/2/3). Layout and colors as in Figure 3.

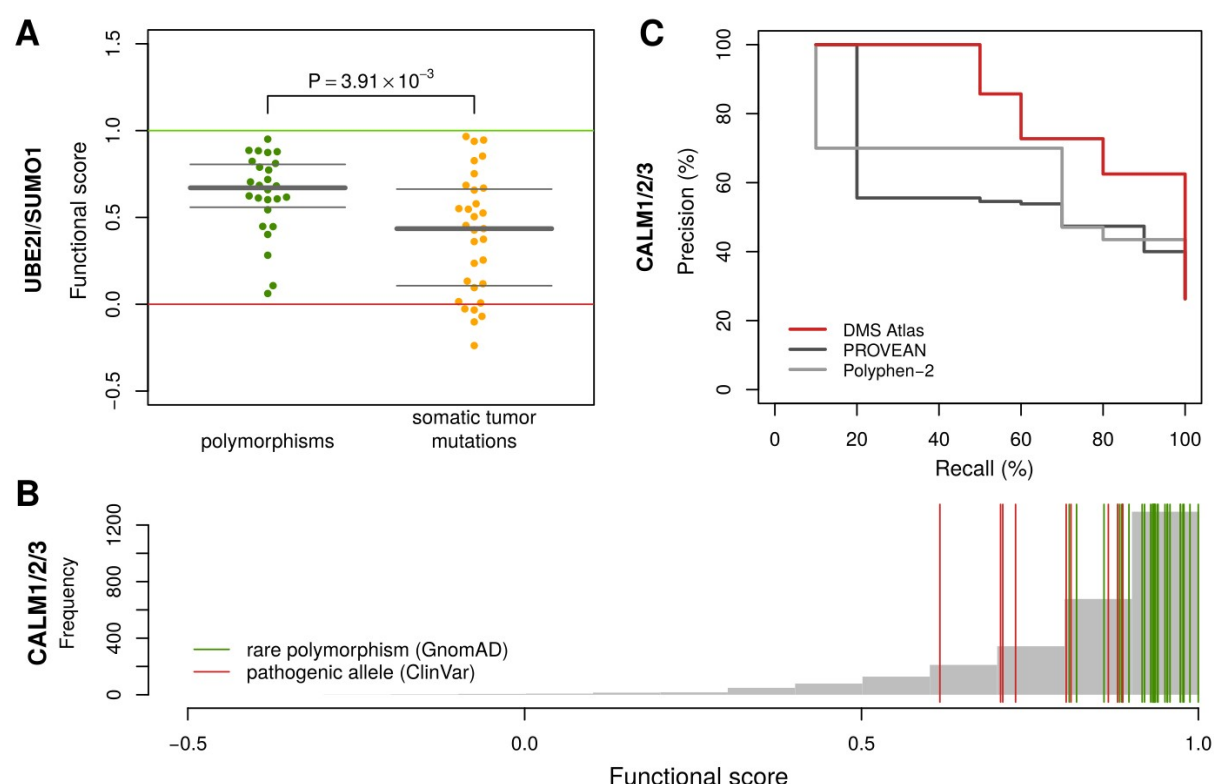


Figure 5: (A) Comparison of (refined) functional scores between rare polymorphisms (GnomAD) and somatic tumor mutations (COSMIC) in UBE2I and SUMO1. Bars show median and quartiles. One-sided Wilcoxon test, $n=\{26,31\}$ (unit:variants), $W=570.5$, $P=3.73 \times 10^{-3}$. (B) Impact score distributions in calmodulin overlaid with previously observed alleles in CALM1, CALM2 and CALM3: Rare alleles from GnomAD are shown in green; ClinVar alleles classified as pathogenic are shown in red. (C) Precision-Recall Curves for our DMS atlas, PROVEAN, and PolyPhen-2 with respect to distinguishing Gnomad variants from pathogenic alleles from ClinVar.

Tables

Table 1: Map quality comparison. Experimental max(s.e.m.): The largest standard error associated with any experimentally measured score in the given dataset; Refined max(s.e.m.): The largest standard error associated with any refined score in the given dataset. Refinement > 0.05: The percentage of variants whose scores were changed by more than 0.05 as a result of refinement.

| Gene | Possible AA changes | Achieved AA changes | Imputation RMSD | Experimental max(s.e.m.) | Refined max(s.e.m.) | Refinement > 0.05 |
|-------|---------------------------|------------------------|--------------------|-----------------------------|------------------------|----------------------|
| UBE2I | 3021 | 2563 (85%) | 0.24 | 0.36 | 0.25 | 2.46% |
| SUMO1 | 1919 | 1700 (89%) | 0.25 | 0.19 | 0.17 | 1.06% |
| TPK1 | 4617 | 3181 (69%) | 0.34 | 0.49 | 0.37 | 5.51% |
| CALM1 | 2831 | 1813 (64%) | 0.29 | 0.28 | 0.22 | 6.84% |

Table 2: Invitae VUS classification. Abbreviations: sd/rmsd = standard error (for measured values) / root-mean-squared-deviation (for imputed values); imp/ref = imputation/refinement; mild ref. = mild refinement

| Variant | MAF | sd/rmsd | imp/ref | unrefined | DMS | DMS call | indication |
|--------------|--------------------|---------|-----------|-----------|------|-----------------|------------|
| D94A | NA | 0.26 | imputed | NA | 0.46 | likely damaging | Cardio |
| D96H | NA | 0.26 | imputed | NA | 0.72 | likely damaging | Cardio |
| I28V | 10 ⁻⁵ | 0.05 | mild ref. | 0.88 | 0.88 | uncertain | Cardio |
| N98S | NA | 0.05 | mild ref. | 0.89 | 0.89 | uncertain | Cardio |
| T35I | 4x10 ⁻⁶ | 0.04 | mild ref. | 0.93 | 0.93 | likely benign | Non-Cardio |
| E48G | NA | 0.05 | mild ref. | 0.93 | 0.93 | likely benign | Cardio |
| G26D | NA | 0.06 | mild ref. | 0.94 | 0.94 | likely benign | Non-Cardio |
| T27S | 3x10 ⁻⁵ | 0.05 | mild ref. | 0.96 | 0.96 | likely benign | Non-Cardio |
| V122A | NA | 0.05 | mild ref. | 0.98 | 0.98 | likely benign | Non-Cardio |
| A104G | NA | 0.08 | mild ref. | 1.00 | 1.00 | likely benign | Non-Cardio |