


# Understanding the Origins of Loss of Protein Function by Analyzing the Effects of Thousands of Variants on Activity and Abundance

Matteo Cagiada,<sup>1</sup> Kristoffer E. Johansson,<sup>1</sup> Audrone Valanciute,<sup>1</sup> Sofie V. Nielsen,<sup>1</sup> Rasmus Hartmann-Petersen,<sup>1</sup> Jun J. Yang,<sup>2,3</sup> Douglas M. Fowler,<sup>4,5</sup> Amelie Stein,<sup>1</sup> and Kresten Lindorff-Larsen <sup>\*,1</sup>

<sup>1</sup>Linderstrøm-Lang Centre for Protein Science, Department of Biology, University of Copenhagen, Copenhagen, Denmark

<sup>2</sup>Department of Pharmaceutical Sciences, St. Jude Children's Research Hospital, Memphis, TN, USA

<sup>3</sup>Department of Oncology, St. Jude Children's Research Hospital, Memphis, TN, USA

<sup>4</sup>Department of Genome Sciences, University of Washington, Seattle, WA, USA

<sup>5</sup>Department of Bioengineering, University of Washington, Seattle, WA, USA

\*Corresponding author: E-mail: lindorff@bio.ku.dk.

Associate editor: Banu Ozkan

## Abstract

Understanding and predicting how amino acid substitutions affect proteins are keys to our basic understanding of protein function and evolution. Amino acid changes may affect protein function in a number of ways including direct perturbations of activity or indirect effects on protein folding and stability. We have analyzed 6,749 experimentally determined variant effects from multiplexed assays on abundance and activity in two proteins (NUDT15 and PTEN) to quantify these effects and find that a third of the variants cause loss of function, and about half of loss-of-function variants also have low cellular abundance. We analyze the structural and mechanistic origins of loss of function and use the experimental data to find residues important for enzymatic activity. We performed computational analyses of protein stability and evolutionary conservation and show how we may predict positions where variants cause loss of activity or abundance. In this way, our results link thermodynamic stability and evolutionary conservation to experimental studies of different properties of protein fitness landscapes.

**Key words:** protein variants, multiplexed assays of variant effects, deep mutational scanning, protein stability, disease variants, genomics, protein structure–function.

## Introduction

Mutational analysis of proteins has provided us with a wealth of information about the molecular interactions that stabilize proteins and govern their functions (Fersht 1999). This information has in turn enabled us to engineer proteins with improved activities and stability (Goldenzweig and Fleishman 2018), to better understand how mutations cause disease (Stein et al. 2019), and help elucidate the role of protein stability in evolution (DePristo et al. 2005; Echave et al. 2016).

Computational analyses of missense variants in genetic diseases have suggested that loss of function via loss of protein stability is a major cause of disease (Wang and Moulton 2001; Ferrer-Costa et al. 2002; Steward et al. 2003; Yue et al. 2005; Casadio et al. 2011; Gao et al. 2015; Stein et al. 2019) because unstable proteins either aggregate or become targets for the cell's protein quality control apparatus and are degraded (Nielsen et al. 2020). Indeed, cellular studies of disease-

causing variants in a number of genes have shown that many variants are degraded in the cell (Meacham et al. 2001; Olzmann et al. 2004; Yaguchi et al. 2004; Ron and Horowitz 2005; Yang et al. 2011 2013; Arlow et al. 2013; Chen et al. 2017; Nielsen et al. 2017; Matreyek et al. 2018; Abildgaard et al. 2019; Scheller et al. 2019; Suiter et al. 2020). For this reason, several methods for predicting and understanding disease-causing variants include predictions of changes in protein stability (Yue et al. 2005; Casadio et al. 2011; De Baets et al. 2012; Ancien et al. 2018; Wagih et al. 2018; Gerasimavicius et al. 2020).

Although stability-based predictions can be relatively successful and may provide mechanistic insight into the origins of disease, it is also clear that variants can cause disease via other mechanisms such as removing key residues in an active site or perturbing interactions or regulatory mechanisms. Thus, methods used to predict the pathogenicity of missense variants often combine analysis of sequence conservation with information on protein structure and stability and other

© The Author(s) 2021. Published by Oxford University Press on behalf of the Society for Molecular Biology and Evolution.

This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact [journals.permissions@oup.com](mailto:journals.permissions@oup.com)

Open Access

sources of information (Kumar et al. 2009; Adzhubei et al. 2010; De Baets et al. 2012; Kircher et al. 2014; Choi and Chan 2015; Ioannidis et al. 2016).

In the same way that perturbed protein stability may cause disease, protein structure, stability, and folding also puts restraints of how amino acid sequences evolve (Mirny and Shakhnovich 1999; DePristo et al. 2005; Liberles et al. 2012). Obviously, however, other considerations such as intrinsic activity and interactions with other molecules also play important roles in determining how sequences evolve, and site-to-site variation in evolutionary rates appear to be determined by a complex relationship between effects on stability and other functional constraints (Echave et al. 2016; Jimenez et al. 2018; Echave 2019).

In order to understand better the relationship between protein stability, abundance, and function, we here asked the question of what fraction of single amino acid changes in a protein causes loss of function via loss of stability and cellular abundance of the proteins. Until recently, mutational analyses of proteins have mostly relied on a one-by-one approach in which individual amino acid changes are introduced and effects on various properties of a protein are tested—often using in vitro experiments on purified proteins (Shoichet et al. 1995; Fersht 1999). Such experiments can now be complemented by experiments that simultaneously probe the effects of thousands of variants in a single assay. Such multiplexed assays of variant effects (MAVEs, also often termed deep mutational scans) are based on developments in high-throughput DNA synthesis, functional assays, and sequencing techniques (Kinney and McCandlish 2019). Briefly, a selection procedure (e.g., for growth rate or a fluorescent reporter of a protein property) is applied to a large library of variants, each expressed in individual cells. Variants change in frequency depending on how they perform under the conditions of the selection, and the frequency of each variant before and after the selection is determined using next-generation DNA sequencing. Changes in variant frequency are used to compute a score that describes each variant's effect on the property under selection. Such data can be used as an input to protein engineering (Araya et al. 2012; Shin and Cho 2015), to map local regions of fitness landscapes and help elucidate genotype–phenotype relationships (Hietpas et al. 2011; Sarkisyan et al. 2016; Fernandez-de Cossio-Diaz et al. 2020), and to understand which and how mutations may cause disease (Starita et al. 2015; Weile and Roth 2018; Stein et al. 2019).

Now, for the first time, we have available measurements of thousands of variant effects on two key protein properties, activity and abundance, measured in multiple proteins. Here, we take advantage of these data to examine more broadly how substitutions affect activity and stability. We examine how variants may affect abundance and activity differently to find functionally important positions in proteins (Chiasson et al. 2020), and to understand whether different types of effects are found in different regions of a protein's structure.

To do so, we here analyze two different types of MAVEs that probe different aspects of protein function. As subjects of our study, we have chosen two medically relevant human

proteins, PTEN (phosphatase and tensin homolog) and NUDT15 (nucleoside diphosphate-linked to x hydrolase 15), because for both of these proteins multiplexed functional data exist from two different assays: One measuring the effect of variants on the activity of the protein via a growth rate (Mighell et al. 2018) or drug sensitivity (Suiter et al. 2020) phenotype, and an assay that probes the effects of amino acid changes on cellular abundance (Matreyek et al. 2018; Suiter et al. 2020). We will sometimes refer to the abundance data as reporting on “stability” and the growth-based activity data as “activity,” or “function,” recognizing that the experiments report on a complex interplay of effects during the experimental assays. Notably, low scores in the activity-based assays might occur both due to loss of intrinsic enzymatic function, but also, for example, due to decreased protein abundance. Indeed, we use the complementary information on protein abundance to disentangle effects on abundance and intrinsic activity.

PTEN is a 403 amino-acid residue long lipid phosphatase expressed throughout the human body, and mutations in the PTEN gene have been associated with cancer and autism spectrum disorders (Yehia et al. 2019). In mice, PTEN has been shown to suppress tumor development via dephosphorylation of phosphatidylinositol lipids, although in vitro PTEN has been shown to have a broader range of substrates including proteins. PTEN is composed of two domains: a catalytic tensin-like domain (residues 14–185) and a C2 domain (residues 190–350) that mediates membrane recruitment (Lee et al. 1999). The C-terminal region of PTEN is disordered with a PDZ-domain binding region (residue 401–403) (Valiente et al. 2005). Our analysis of PTEN includes a MAVE that probes the effects of most single amino acid substitutions when assayed for lipid phosphatase activity in yeast (Mighell et al. 2018), whose growth had been made dependent on the ability of PTEN to catalyze the formation of essential phosphatidylinositol bisphosphate (PIP<sub>2</sub>) from its triphosphate (PIP<sub>3</sub>). Although these experiments only probe one function of PTEN and might be affected also, for example, by expression levels, it has been shown that the resulting data accurately classifies the pathogenicity of PTEN variants (Mighell et al. 2018; Jepsen et al. 2020). We complement these data with results from a different MAVE in which variant effects on cellular abundance are determined in an experiment termed “variant abundance by massively parallel sequencing” (VAMP-seq) (Matreyek et al. 2018). In VAMP-seq, the steady-state abundance of protein variants in cultured mammalian cells is detected by fusion to a fluorescent protein, and cells are sorted using fluorescent activated cell sorting. The outcome of the VAMP-seq experiment is not substantially affected by the fusion with full-length GFP and correlates with in vitro measurements of thermal stability (Matreyek et al. 2018), but importantly also captures other effects that might affect protein abundance and which could be relevant for function, evolution, and disease. Our analysis here covers the 56% of all possible single amino acid variants in PTEN for which we have measurements for both the activity and abundance, and thus complements our recent

analysis of a small number of disease variants in PTEN (Jepsen et al. 2020).

NUDT15 is a nucleotide triphosphate diphosphatase that consists of 164 amino acids in a nudix hydrolase domain featuring a conserved nudix box that coordinates the catalytic  $Mg^{2+}$ . The biologically relevant assembly is reported to be a homodimer although the monomer also has catalytic activity (Carter et al. 2015). NUDT15 deficiency is associated with intolerance to thiopurine drugs (Yang et al. 2014; Moriyama et al. 2016, 2017; Nishii et al. 2018), which are widely used in the treatments of leukemia and autoimmune diseases (Karran and Attard 2008). Thiopurines are a class of antimetabolite drugs that form the active metabolite, thio-dGTP, which competes with dGTP and causes apoptosis when incorporated extensively into DNA. NUDT15 hydrolyzes thio-dGTP and thus negatively regulates the levels and cytotoxic effects of thiopurine metabolites. Therefore, NUDT15 variants that decrease function are a major cause of toxicity during thiopurine therapy, and thus the dose of the drug may be personalized to match the metabolism of these compounds (Relling et al. 2019). The high drug sensitivity of cells with compromised NUDT15 function has been used in a MAVE to assay 95% of all single amino acid variants for causing intolerance toward thiopurine drugs (Suiter et al. 2020). The same library and cells were also used in a VAMP-seq experiment to probe variant effects on cellular abundance, and like in the case of PTEN, the results were shown to correlate with in vitro measurements of thermal stability. As in the case of PTEN, the outcome of the MAVE might depend on the exact conditions and, for example, drug concentration used, but was shown to capture the effects of several known pharmacogenetic variants (Suiter et al. 2020).

Here, we have analyzed the effect of variants on activity and cellular abundance in both PTEN and NUDT15 to provide a global view of what fraction of variants cause substantial loss of activity in the cell, and what fraction of these variants do so via loss of protein abundance. We find that approximately one-third of all variants cause loss of protein activity, and that about half of these do so most likely because of loss of protein abundance. Variants that cause loss of abundance are often found inside the protein core, whereas variants that cause loss of activity without affecting abundance are often found in functionally important positions including those involved in catalysis or that interact with substrates. We also find that we can predict rather accurately the positions where substitutions generally give rise to decreased abundance and activity, whereas it remains difficult to quantitatively predict the effects of individual variants. Together, our results provide further insight into the link between thermodynamic stability and evolutionary conservation and experimental studies of different properties of fitness landscapes.

## Results and Discussion

### Global Analysis of Variant Effects

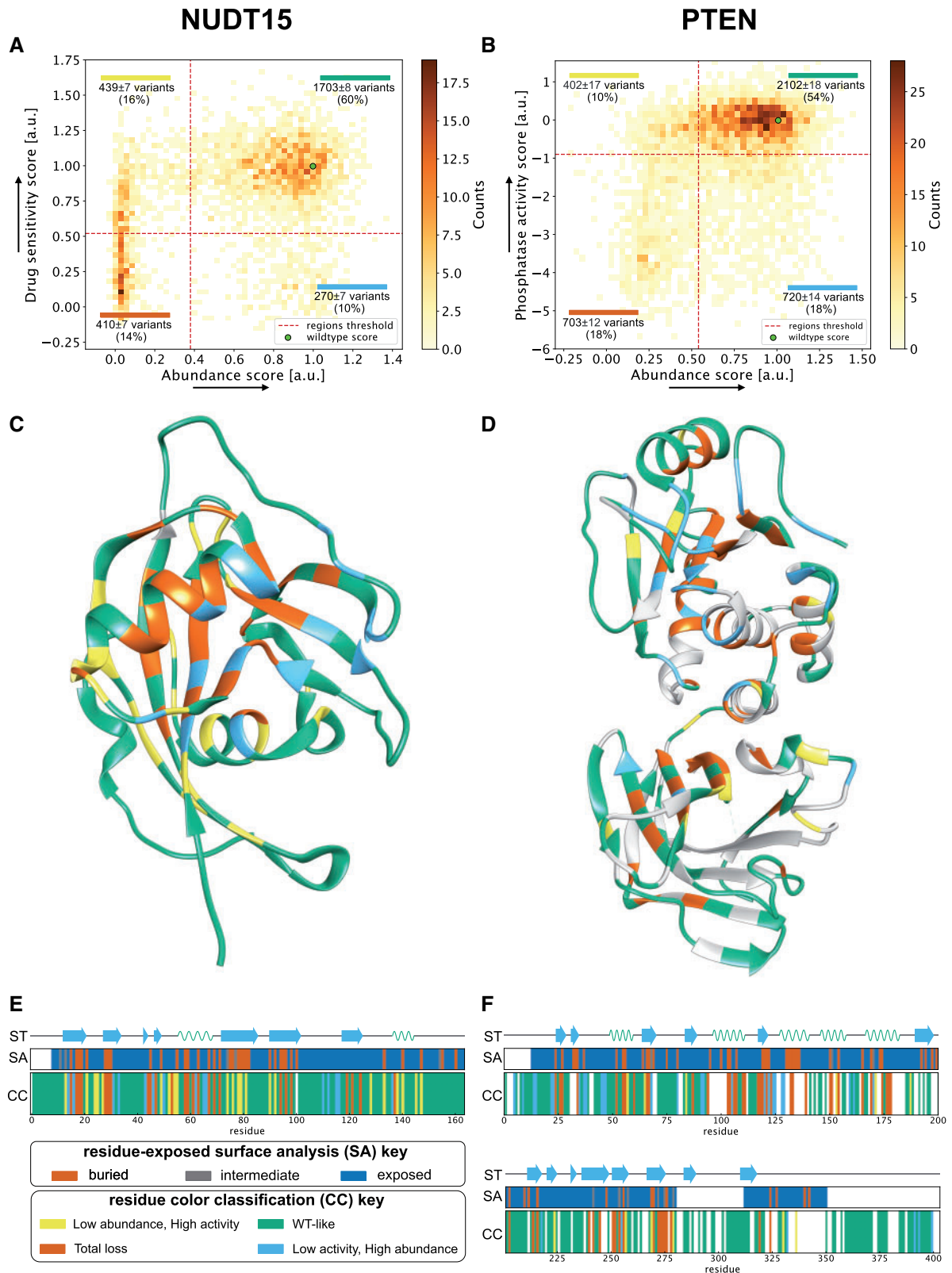
We collected data from multiplexed assays reporting on both the activity and abundance of a total of 2,822 variants in NUDT15 (Suiter et al. 2020) and 3,927 variants in PTEN

(Matreyek et al. 2018; Mighell et al. 2018) (supplementary fig. S1, Supplementary Material online). Scripts to repeat our analyses are available online at <https://github.com/KULL-Centre/papers/tree/master/2020/mave-analysis-cagiada-et-al> (last accessed April 6, 2021). Two-dimensional histograms reveal that most variants have high scores in both assays, indicating wild-type like abundance and activity under the conditions of the cellular assays (fig. 1A and B).

In order to separate wild-type like variants from those with decreased activity and/or abundance, we define a threshold value for all scores (supplementary fig. S2, Supplementary Material online). These thresholds define four classes of variants according to whether the variant showed high or low scores in the activity-based and abundance MAVEs. For simplicity, each class is also associated with a color. “WT-like” variants had wild-type like activity and abundance and are shown in green. “Low-activity, high-abundance” variants had WT-like abundance but low activity in the assays, and are shown in blue. “Low-abundance, high-activity” variants had WT-like activity but low abundance in the assays and are shown in yellow. “Total loss” variants had low activity and low abundance and are shown in red.

For both proteins, the majority of variants are wild-type like (60% for NUDT15 and 54% for PTEN; fig. 1A and B; green). The total-loss category represents variants that both show loss of activity and low cellular abundance (14% for NUDT15 and 18% for PTEN; fig. 1A and B; red), and as discussed further below, we expect that most of these variants lose activity because of their low abundance. Of the total of 680 and 1,403 variants with low activity in NUDT15 and PTEN, respectively, 60% and 50% lose activity together with loss of abundance. The low-activity, high-abundance variants are still abundant in the cell but inactivated by other means, for example, by changes in amino acids in the active site (fig. 1A and B; blue). The low-abundance, high-activity class, which contains 16% of NUDT15 and 10% of PTEN variants (fig. 1A and B; yellow), show low abundance levels, but high levels of activity in the activity-based assay and are not as easily explained by a single mechanism.

To focus our analysis on different types of variant effects in different parts of the protein structure, and to decrease uncertainty coming from examining individual variants, we converted the variant data into positional categories that represent the most frequent class among the variants at that position. We performed this classification procedure at all positions with at least five tested variants (99% for NUDT15 and 88% for PTEN), which also helped average out noise from examining individual variants with intermediate scores, and represent the classes using the same names and coloring scheme as for the variants. This results in 62% and 60% positions classified as WT-like for NUDT15 and PTEN, respectively (fig. 1C and D; green). On the other hand, at 15% and 22% of the positions most variants cause loss of activity together with loss of abundance (fig. 1C and D; red), whereas loss of activity without loss of abundance is the most common outcome at 9% and 12% of the positions (fig. 1C and D; blue). Finally, at 14% of the positions in



**Fig. 1.** Overview of the NUDT15 and PTEN multiplexed data analyzed in this work. (A) and (B) show 2D histograms that combine the data from the activity-based MAVE on the y axis with the results from the VAMP-seq experiment on the x axis. Variants are categorized based on the region of the 2D histogram (dashed lines) they belong to. The fractions of variants falling in each of the four quadrants are indicated, with errors of the mean estimated by bootstrapping using the uncertainties of the experimental scores. The two green points indicate the wild type. Arrows on the axes indicate directions of greater abundance or activity; for detailed definitions of the scores and their uncertainties, we refer the reader to the original publications (Matreyek et al. 2018; Mighell et al. 2018; Suiter et al. 2020). Panels (C) and (D) show a per-position consensus category (CC) colored onto the structure of the proteins (PDB entry 5LPG for NUDT15 and 1D5R for PTEN). Panels (E) and (F) show the positional color categories together with the secondary structure (ST) and solvent accessibility (SA). The four classes of variants/positions are represented by a color: “WT-like” (green), “Low activity, high abundance” (blue), “Low abundance, high activity” (yellow), and “Total loss” (red).



NUDT15 and 6% in PTEN the variants most often have low abundance, but high levels of activity (fig. 1C and D; yellow).

We validated the classifications using a clustering method that does not depend on defining cutoffs for the experimental scores. We grouped together positions with similar variant profiles in the two MAVEs (see Materials and Methods), and find overall very good agreement with the cutoff-based method in particular for the WT-like, total-loss and loss of activity, high-abundance categories (supplementary figs. S3 and S4, Supplementary Material online). For NUDT15, we find that 133/163 positions are classified in the same way using the two different methods, with the most variable results occurring in the category with low abundance but sufficient activity to sustain growth (supplementary fig. S3, Supplementary Material online). For PTEN, we analyzed the data using either three or four clusters, with the former appearing to be the more natural classification. In that case, 246/310 positions are classified in the same way using the two methods, with the 12 positions in the low-abundance, high-activity (yellow) category ending either as WT-like or total-loss. This indicates that three of the four categories of position effects are identified more robustly, corresponding to substitutions generally resulting in 1) WT-like activity in both assays, 2) loss of activity and abundance, or 3) loss of activity, while retaining WT-like abundance. The low-abundance, high-activity positions are, however, less robustly classified and we do not analyze them further.

As expected, amino acids at buried positions are in general sensitive to mutations. In NUDT15, 35 out of the 163 amino acids are fully buried, and half of these (49%) are classified as sensitive to mutations in both the activity- and abundance-based assays (red label) with the remaining buried positions mainly classified as low abundance, high activity (34%; fig. 1E and F). Because the variant coverage is lower in PTEN, only 355 of 403 positions can be classified in this way, and only 34 of these 355 are fully buried. Among these 34, 80% are classified as “unstable” positions (low-abundance, high-activity, and total-loss categories). Thus, loss of abundance is the typical reason for loss of activity for variants at buried positions.

Low-activity, high-abundance positions are defined as having the majority of the tested variants that have lost activity, but are still abundant in the cell. Previously such positions have been found to map to functionally important sites in the membrane protein VKOR (Chiasson et al. 2020). We find that in PTEN, these variants and positions are mainly found in the catalytic phosphatase domain (supplementary fig. S5, Supplementary Material online) and include the active site (fig. 2A and B). In NUDT15, we find the low-activity, high-abundance positions in several different regions. One group is located in proximity of the substrate-binding site and includes previously discussed Arg34 and Gly47 (Suiter et al. 2020). Another group includes the residues that coordinate a magnesium ion (Suiter et al. 2020). Finally, we find a group of residues that stretches from the substrate-interacting Arg34 and Gln44 (Carter et al. 2015) to Asn117 and Asn111 more distal from the substrate-binding pocket and connected via a hydrogen bond network (fig. 2C and D). Asn111 and Asn117 appear to help position a loop (residues 111–117) that

includes the magnesium-coordinating Glu113, and although these residues do not directly contact the substrate, many substitutions lead to loss of function without loss of abundance.

Having found that many low-activity, high-abundance positions play functional roles, we asked the question whether they are generally found near the active sites in NUDT15 and PTEN. Using Gly47 in NUDT15 and Arg130 in PTEN as reference points in the active sites in these two proteins, we find that the low-activity, high-abundance positions, where variants typically show loss of activity, but not loss of abundance, are clustered around the active sites. Specifically, we find all of these positions in NUDT15 are within 14 Å of Gly47 ( $C_{\alpha}$ -distances). The average distance between low-activity, high-abundance positions, and Gly47 is 9 Å, a value that can be compared with the average (15 Å) over all positions in NUDT15. In PTEN, we find that 29 of 32 low-activity, high-abundance positions are found in the catalytic domain. All of these 29 positions are within 22 Å of Arg130, with the average distance to Arg130 being 14 Å (compared with 21 Å over all positions).

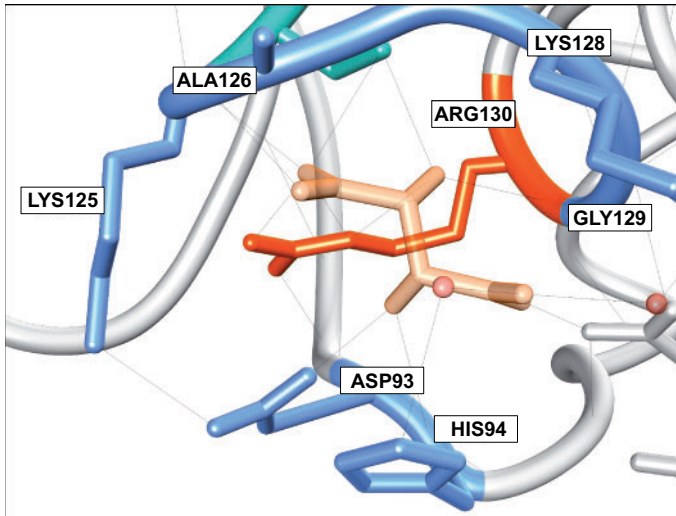
Although the typical outcome at low-activity, high-abundance positions is loss of activity, not all substitutions have equally large effects, and some amino acid substitutions are more likely to be detrimental to function than others. We thus examined the individual substitutions at the low-activity, high-abundance positions to ask whether particular types of substitutions preserve function better than others. Although such an analysis is made difficult by the small numbers of substitutions when broken down to start and end amino acid, we do find evidence suggesting differences depending on amino acid chemistry (supplementary fig. S6, Supplementary Material online). For example, at low-activity, high-abundance positions with Asn as the wild-type residue, it appears that substitutions to other small and polar amino acids (Asp, Ser, Thr) preserve function better than, for example, substitutions to hydrophobic amino acids (supplementary fig. S6, Supplementary Material online). More generally, it has previously been observed that there is a substantial effect of amino acid type on the outcome of a MAVE experiment (Gray et al. 2017; Dunham and Beltrao 2020), and we here find similar results (supplementary fig. S7, Supplementary Material online). As expected, we find that many total loss variants are substitutions of hydrophobic amino acids with charged and polar amino acids, whereas substitutions of hydrophobic residues with other hydrophobics are more common in the wild-type like category. More detailed analyses of these effects, however, are hampered by the low number of many of the types of substitutions.

### Computational Predictions of Multiplexed Data from MAVES

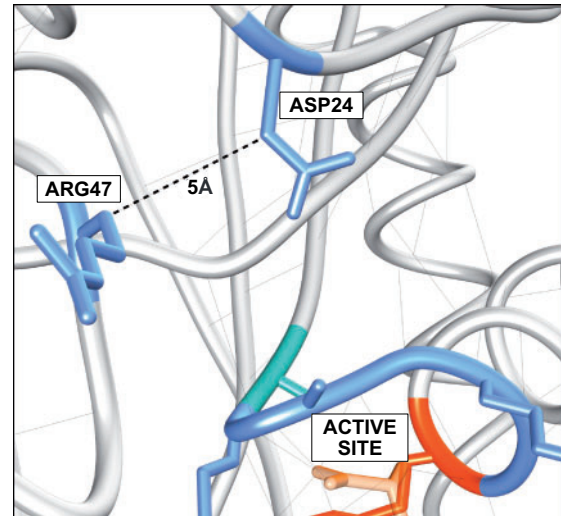
As described previously and demonstrated above, MAVES provide a wealth of data not only for use in medical applications (Weile and Roth 2018; Stein et al. 2019) but also for understanding basic properties of proteins (Dunham and Beltrao 2020). Despite recent advances in proteome-wide experiments (Després et al. 2020), it is still not possible to

## PTEN

A

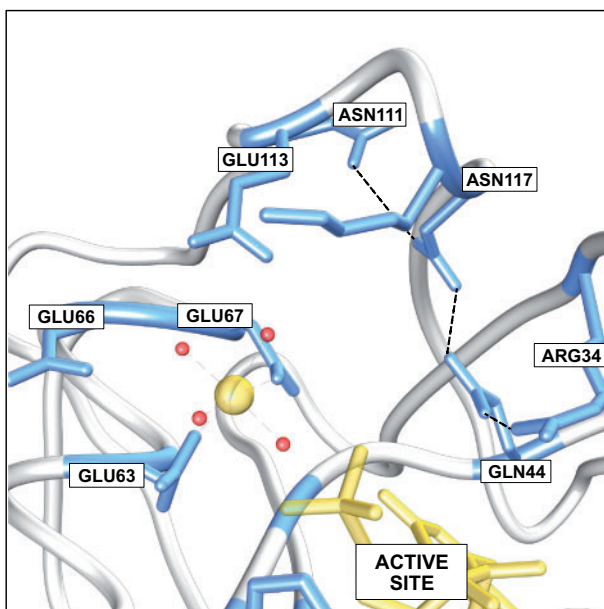


B

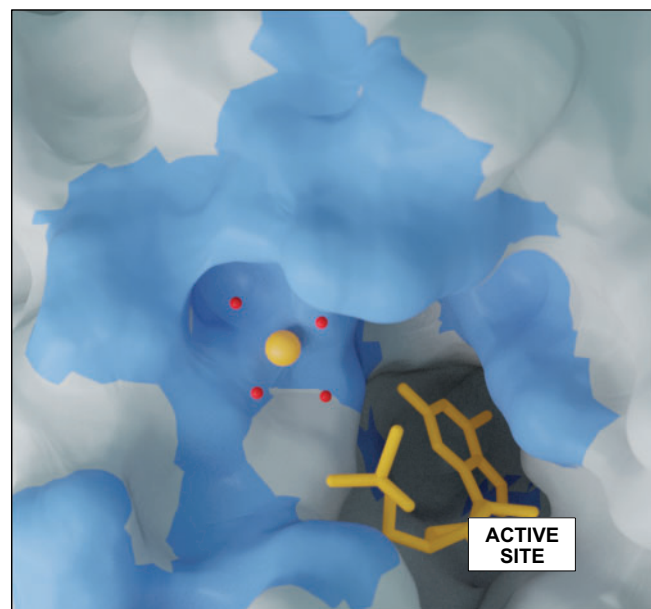


## NUDT15

C



D



**FIG. 2.** Examples of “low-activity, high-abundance” positions. (A) Residues in PTEN in the low-activity, high-abundance category (blue) include residues in and surrounding the catalytic phosphatase site including some that directly interact with the substrate (here mimicked by the inhibitor tartrate; Lee et al. 1999). (B) Other residues that are more distant to the active site also fall in this category, and variants in this region could perturb the integrity of the active site. (C and D) Examples of functionally important residues in NUDT15 that are close to, but outside of the active site. In particular, we identified four conserved residues (Asn111, Asn117, Gln44, Arg44) that appear to connect via a hydrogen bond network, and whose perturbation could affect the hydrolysis of the thiopterins.

probe all possible variants in all proteins experimentally, and thus computational methods remain an important supplement to predict and understand variant effects. Experimental data from MAVEs are thus increasingly used to benchmark prediction methods, as they provide a broad view of the effect of amino acid substitutions in proteins (Hopf et al. 2017; Jepsen et al. 2020; Livesey and Marsh 2020; Reeb et al. 2020).

Recently, we exploited the two different MAVEs for PTEN to analyze a small number of pathogenic variants together

with variants that have been observed in a broader analysis of the human population (Jepsen et al. 2020). Specifically, we compared the abundance-based (VAMP-seq) and activity-based multiplexed data with two computational methods aimed at capturing either 1) specifically protein stability or 2) function more broadly. Here, we build on this work, by 1) applying computational modeling to predict changes in thermodynamic protein stability using Rosetta (Park et al. 2016) and 2) using evolutionary conservation as a more general

view of which amino acid changes would be tolerated while maintaining function (Ekeberg et al. 2014). The former uses as input the structure of NUDT15 or PTEN to predict the change in protein stability ( $\Delta\Delta G$ ), whereas the latter uses a sequence alignment of homologous proteins as input to a computational assessment of conservation, taking both site and pair-conservation (coevolution) into account, quantified by a score (which we by analogy to  $\Delta\Delta G$  term  $\Delta\Delta E$ ) that estimates how likely a substitution would be. We note that the same kind of model can be used to predict contacts in protein structure, but in line with previous work (Lapedes et al. 2012; Lui and Tiana 2013; Hopf et al. 2017; Nielsen et al. 2017) is here used to estimate the effects of amino acid substitutions. We also note that it has previously been shown that the “pair terms” in these models, that capture effects of (apparent) coevolution between pairs of sites, improve accuracy in these predictions (Hopf et al. 2017). As previously argued (Jepsen et al. 2020), the  $\Delta\Delta G$  calculations are more akin to the results of an abundance-based MAVE (both capturing aspects of protein stability), whereas the  $\Delta\Delta E$  values capture a broader range of effects as would also be expected from an activity-based MAVE.

We thus compared the computational predictions of  $\Delta\Delta G$  and  $\Delta\Delta E$  with each of the two multiplexed assays for NUDT15 and PTEN (supplementary fig. S8, Supplementary Material online). As expected, we find that stability predictions correlate better with the abundance-based MAVE than with the activity-based MAVE, whereas for the evolutionary analysis, the situation is reversed (supplementary fig. S9, Supplementary Material online). In the case of NUDT15, for example, the data from the abundance-based MAVE correlate more strongly with the  $\Delta\Delta G$  calculations ( $r_p = 0.57$ ) than with  $\Delta\Delta E$  ( $r_p = 0.42$ ), whereas the activity-based MAVE is more poorly correlated with stability predictions ( $r_p = 0.35$ ) than with the conservation-based scores ( $r_p = 0.52$ ). Although the difference is smaller (but still present) for PTEN, the results support the notion that analysis of conservation is a better predictor of general aspects of protein function, whereas the Rosetta calculations support the expected relationship between cellular protein abundance and thermodynamic stability (Matreyek et al. 2018; Abildgaard et al. 2019; Jepsen et al. 2020). In addition, we note that whereas the correlation coefficients are not very high, the results are in line with previous analyses of similar data (Hopf et al. 2017; Jepsen et al. 2020; Livesey and Marsh 2020).

We define threshold values for the computational scores (supplementary figs. S10 and S11, Supplementary Material online) to separate wild-type like from deleterious variants and construct four categories that we label with colors as above. Using a threshold of 2 kcal/mol for the  $\Delta\Delta G$  for both proteins results in 69% (NUDT15) and 65% (PTEN) of the variants being predicted stable. Similarly, from the evolutionary conservation analysis 78% and 58% of all variants for NUDT15 and PTEN, respectively, have scores that indicate that the substitutions are tolerated. Note that, by convention, positive  $\Delta\Delta G$  and  $\Delta\Delta E$  scores indicate loss of stability or

sequence tolerance, respectively, and hence the scales are inverted compared with the scores from the MAVEs.

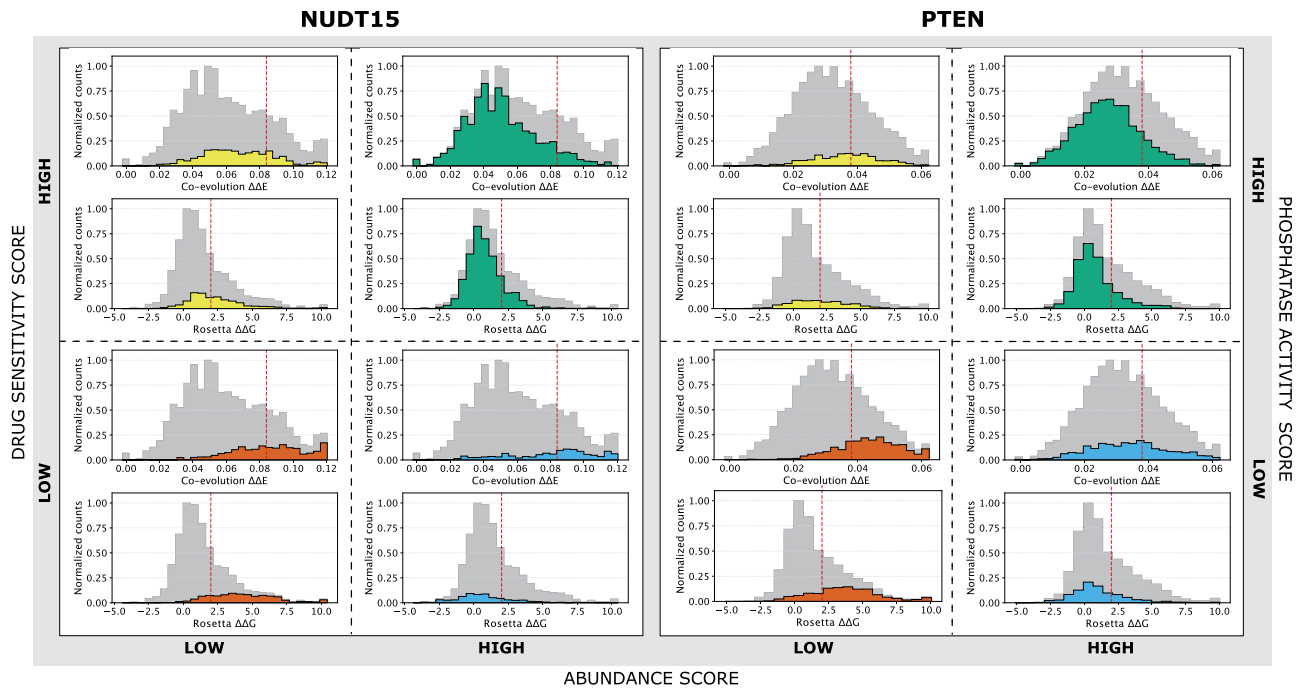
To enable a more direct comparison between the experimental and computational scores, we show histograms of the two computational scores ( $\Delta\Delta G$  and  $\Delta\Delta E$ ) for each of the four classes based on the experimental scores (fig. 3). We find that the variants that experimentally were classified as WT-like (stable and active) generally have low computational values; thus the computational predictions suggest that these substitutions have a mild effect on stability (low  $\Delta\Delta G$ ) and are compatible with substitutions observed in homologous proteins (low  $\Delta\Delta E$ ). We make similar observations for the total loss category, where the computational scores are generally above the cutoff, and for the low-activity, high-abundance category where the computational analysis finds low values of  $\Delta\Delta G$  but higher values of  $\Delta\Delta E$ . Despite these general trends, we find variable agreement in the classification of individual variants by experiments and computation (supplementary fig. S12, Supplementary Material online), with the best agreement in the WT-like and total-loss categories. To examine whether the results from the conservation analyses were specific to using IbsDCA, we also used the Evolutionary Trace Analysis algorithm (Lichtarge et al. 1996; Lua et al. 2016) to analyze the multiple sequence alignments, and found similar results (supplementary fig. S13, Supplementary Material online).

We proceeded by generating and examining the structure–function relationships that we extracted from the computational analyses (supplementary fig. S14, Supplementary Material online). We used the computational results to group the positions into four categories and found a substantial overlap with those found in experiments (supplementary fig. S15, Supplementary Material online), in particular for the WT-like and total-loss categories, with approximately 70% of the positions classified in the same way. This result suggests that the computational analyses better capture general effects at positions compared with individual variants as discussed above (fig. 3). We again used a clustering procedure as an alternative approach to classify positions and find good agreement both with the cutoff-based classification of the computational data as well as with the experiment-based classifications (supplementary fig. S16, Supplementary Material online). Thus, together these results show that a joint computational analysis of stability and conservation can be used to find positions in the protein where substitutions are likely to disrupt thermodynamic stability, and other positions where they will cause loss of activity via removing functionally important residues.

## Conclusions

Large-scale analysis of proteins using multiplexed assays provides opportunities to obtain a global view of variant effects (Gray et al. 2017; Dunham and Beltrao 2020). By combining different assays to read out different properties of a protein it becomes possible to dissect which positions contribute most to which property (Jepsen et al. 2020). Most proteins need to be folded to be active, and thus amino acid substitutions that





**Fig. 3.** Histograms of the two computational scores ( $\Delta\Delta G$  and  $\Delta\Delta E$ ) in NUDT15 and PTEN.  $\Delta\Delta G$  aims to capture effects purely on the thermodynamic stability, with high values indicating destabilized variants.  $\Delta\Delta E$  captures evolutionary conservation, as calculated by a model that takes both site and pairwise coevolution into account, and with high values indicating nonconservative substitutions. Thus, for both  $\Delta\Delta G$  and  $\Delta\Delta E$  positive values indicate detrimental substitutions, whereas in the experiments low values indicate substitutions that cause loss of activity or abundance. For both proteins, we split the histograms up according to the four categories of variants determined from the experiments, as indicated by the axes with high and low experimental scores for abundance and activity. Thus, for example, the two green histograms for NUDT15 indicate the distributions of  $\Delta\Delta G$  and  $\Delta\Delta E$  values for those variants that are classified as stable and active by the MAVEs, and indeed it is clear that most of these variants have scores that are below the cutoff (red dashed lines). In addition to the colored histograms, we also show the full histogram of all analyzed variants (gray) to ease comparison between the subsets and the full set of variants.

lead to loss of stability will often lead to loss of function. Loss of stability thus appears to be an important driver for disease (Yue et al. 2005; Stein et al. 2019) and determinant of evolutionary rates (Echave et al. 2016), and vice versa it has been shown that residues in active sites may be suboptimal for stability (Shoichet et al. 1995).

We have here exploited the availability of data generated by MAVEs for two proteins, with one experiment probing general effects on protein activity and another directly assessing cellular abundance. We show that a global analysis of these experiments can provide insight into how proteins function and how activity may be perturbed. With the assays considered here, we find that most variants have at most a modest effect on protein activity. Of the approximately 30% of the variants that cause substantial loss of activity, we find that approximately 50% also cause loss of abundance. Thus, although it is not surprising that there in many cases is a correlation between loss of function and loss of abundance, we here provide quantitative estimates on the relative importance of these effects across a wide range of substitutions in two unrelated proteins. The relative amounts of “low-activity, high-abundance” and “total-loss” variants that we find can be compared with our previous analysis of 42 disease-causing variants in PTEN, where we found a comparable fraction (~60%) of the disease-causing variants appears to cause loss of function via loss of stability and thereby cellular protein

abundance (Jepsen et al. 2020). Indeed, most (but not all) pathogenic variants in PTEN (Mighell et al. 2018; Jepsen et al. 2020) score low in the activity-based MAVE, with a substantial fraction of these also having low abundance (supplementary fig. S17, Supplementary Material online), whereas the situation for pharmacogenetic variants in NUDT15 (Suiter et al. 2020) is more complex (supplementary fig. S17, Supplementary Material online). Similarly, in our studies of pathogenic missense variants in the MLH1 gene, we found low steady-state protein levels (<50% of wild type) in seven out of 16 pathogenic variants (Abildgaard et al. 2019). Thus, at least in these cases, it appears that the fraction of variants that cause disease via this mechanism reflects the overall fraction of “total loss” variants in the protein. An interesting question for future experiments is how many of these variants would be active if protein levels could be restored for example by chemical chaperones or modulating the protein quality control apparatus (Arlow et al. 2013; Kampmeyer et al. 2017). Indeed, chaperones are known to help buffer against destabilizing variants during evolution (Rutherford and Lindquist 1998; Tokuriki and Tawfik 2009).

Building on previous work (Cheng et al. 2005; Chiasson et al. 2020), we also show how we can use variant effects on protein activity and abundance/stability to find functionally important residues both by experiments and computation. For several surface-exposed residues, many variants cause loss



of activity, but without substantial loss of abundance. We find that these include the active sites in NUDT15 and PTEN, but also discover functionally important sites adjacent to these active sites. The importance of second shell positions for modulating the structure or dynamics of active site residues has for example also emerged in studies of ligand binding (Tinberg et al. 2013) and enzyme evolution and design (Campbell et al. 2016; Broom et al. 2020). In our analysis of functional residues, we mostly focused on general effects at each position, rather than specific effects of individual substitutions. We did this to average out noise from individual measurements and to find general patterns, but with more data, it would be interesting to perform such structure–sequence–function analysis at the level of individual substitutions.

The relatively tight confinement of these low-activity/high-abundance positions may also explain why predictions of changes in protein stability can be used to predict a substantial number of disease variants: At least in NUDT15 and PTEN the number of positions where substitutions typically cause loss of abundance (and thereby activity) is greater than the number of positions where substitutions cause loss of activity while retaining protein abundance. Indeed, although functional sites induce substantial constraints on amino acid variation during evolution, the strongest effects are those closest to the active sites (Jack et al. 2016; Mayorov et al. 2019). Our ability to predict these sites by combining evolutionary analysis and stability calculations also suggests an approach for discovering new functionally important sites using a combined analysis of protein structure and sequences. We find that approximately 12% of variants in NUDT15 and PTEN appear to be able to support wild-type like growth in the cellular assays even at substantially reduced protein levels. Clearly, there can be a nonlinear relationship between a growth phenotype and protein abundance (Jiang et al. 2013), and this may help explain some of these variants. Future experiments that probe the relationship between expression levels and variant effects in NUDT15 and PTEN may shed further light on these variants. Further, the abundance-based MAVE for PTEN was performed in a cultured mammalian cell line (Matreyek et al. 2018) and the activity-based MAVE was performed in yeast (Mighell et al. 2018), leading to potential differences due to the differences in the quality control and proteostasis machinery in these cells.

In summary, we demonstrate how multiplexed assays and computational analyses are beginning to provide a coherent and comprehensive view of the global effects of variants in proteins. The results highlight that many effects are correctly predicted and thus computation can be used not only to predict whether a variant will cause loss of activity or not, but also provide some mechanistic insight. Clearly, there is room for improvement, and additional experiments on more proteins and covering more aspects of the complicated relationship between protein sequence and functions will help further our ability to predict these effects computationally (Cheng et al. 2005).

## Materials and Methods

### Conservation Analysis of Variant Effects

We used a statistical analysis of multiple sequence alignments (MSAs) of the two proteins to estimate the tolerance toward specific substitutions. In line with previous work, we use a method that includes both site and pairwise conservation (coevolution). We used the WT sequences from UniProt (P60484 and Q9NV35) as input to HHblits (Remmert et al. 2011) to build initial MSAs, which we filtered before calculating the variant effects. The first filter removes sequences (rows) in the MSA with more than 50% gaps. The second filter keeps only positions (columns) that are present in the human target sequences of NUDT15 or PTEN. Finally, we apply a similarity filter (Ekeberg et al. 2013) to remove redundant sequences (more than 80% identical). We use a modified version of the lbsDCA algorithm (Ekeberg et al. 2014), based on  $l_2$ -regularized maximization with pseudocounts to predict the likelihood of every variant of the protein. We use the energy potential generated by the algorithm to evaluate the log-likelihood difference between the wild type and the variant sequences ( $\Delta\Delta E$ ). We verified that the outcome of these analyses did not depend substantially on the parameters used to construct the MSA or to filter the alignments (supplementary fig. S18, Supplementary Material online). We performed Evolutionary Trace Analysis (Lichtarge et al. 1996; Lua et al. 2016) calculations using the webserver available at evolution.lichtargelab.org.

### Structural Analysis

We used Rosetta (GitHub SHA1 99d33ec59ce9fccc5e4f3800-c778a54afdf8504) to predict changes in thermodynamic stability ( $\Delta\Delta G$ ) from the structure of NUDT15 and PTEN using the Cartesian ddG protocol (Park et al. 2016). As starting points, we used the crystal structures of NUDT15 (Valerie et al. 2016) (PDB ID: 5LPG) and PTEN (Lee et al. 1999) (PDB ID: 1D5R). The values obtained from Rosetta were divided by 2.9 to bring them from Rosetta energy units onto a scale corresponding to kcal/mol (Frank DiMaio, University of Washington; personal correspondence) (Jepsen et al. 2020). We used DSSP-2.28 (Kabsch and Sander 1983; Touw et al. 2015) and the same crystal structures as above to classify the burial with a three-state model (Rost and Sander 1994) (buried, intermediate, or exposed).

### Defining Thresholds for Classifying Variants

We defined thresholds for the scores from both MAVEs (supplementary fig. S2, Supplementary Material online), by fitting the variant score distributions using the minimal number of Gaussians (three) needed to obtain a reasonable fit. We then used the intersection of the first and last Gaussian as cutoff for our classifications. We use a cutoff of 2 kcal/mol (similar to the value used in our previous study; Jepsen et al. 2020) for  $\Delta\Delta G$  and varied the cutoff for  $\Delta\Delta E$  to maximize the overlap in the classification of positions (supplementary fig. S11, Supplementary Material online).

To examine the threshold-based classifications, we used a hierarchical clustering algorithm (Ward 1963; Virtanen et al.

2020) to group positions with similar responses to amino acid substitutions. Each position was represented by a 40D vector that contains the scores for each of the 20 possible amino acids in the two MAVEs. Missing values were replaced by the average score over that position. We use the Euclidean distance between these vectors as similarity score in the hierarchical clustering (Ward 1963). To compare with the threshold-based classification, we analyzed this using four clusters, though, in the case of PTEN, we also show the results using only three clusters.

### Residue Classification

We assigned a category to residues for which data for at least five variants are available in both MAVEs. We used the mode (the most common class of the variants at that position) to assign the residue category.

### Supplementary Material

Supplementary data are available at *Molecular Biology and Evolution* online.

### Acknowledgments

This work is a contribution from the PRISM (Protein Interactions and Stability in Medicine and Genomics) centre funded by the Novo Nordisk Foundation (to R.H.-P., D.M.F., A.S., and K.L.-L.; NNF18OC0033950). A.S. is funded by the Lundbeck Foundation (R272-2017-4528). This research was partly supported by NIH grant (R01GM118578 to J.J.Y.). The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

### Data Availability

Code and data to repeat our analyses are available online at <https://github.com/KULL-Centre/papers/tree/master/2020/mave-analysis-cagiada-et-al> (last accessed April 6, 2021).

### References

- Abildgaard AB, Stein A, Nielsen SV, Schultz-Knudsen K, Papaleo E, Shrikhande A, Hoffmann ER, Bernstein I, Gerdes A-M, Takahashi M, et al. 2019. Computational and cellular studies reveal structural destabilization and degradation of MLH1 variants in lynch syndrome. *Elife* 8:e49138.
- 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(4):248–249.
- Ancien F, Pucci F, Godfroid M, Rooman M. 2018. Prediction and interpretation of deleterious coding variants in terms of protein structural stability. *Sci Rep* 8(1):1–11.
- Araya CL, Fowler DM, Chen W, Muniez I, Kelly JW, Fields S. 2012. A fundamental protein property, thermodynamic stability, revealed solely from large-scale measurements of protein function. *Proc Natl Acad Sci U S A* 109(42):16858–16863.
- Arlow T, Scott K, Wagenseller A, Gammie A. 2013. Proteasome inhibition rescues clinically significant unstable variants of the mismatch repair protein MSH2. *Proc Natl Acad Sci U S A* 110(1):246–251.
- Broom A, Rakotoharisoa RV, Thompson MC, Zarifi N, Nguyen E, Mukhametzhano N, Liu L, Fraser JS, Chica RA. 2020. Ensemble-based enzyme design can recapitulate the effects of laboratory directed evolution in silico. *Nat Commun* 11(1):4808.
- Campbell E, Kaltenbach M, Correy GJ, Carr PD, Porebski BT, Livingstone EK, Afriat-Jurnou L, Buckle AM, Weik M, Hollfelder F, et al. 2016. The role of protein dynamics in the evolution of new enzyme function. *Nat Chem Biol* 12(11):944–950.
- Carter M, Jemth A-S, Hagenkort A, Page BDG, Gustafsson R, Griese JJ, Gad H, Valerie NCK, Desroses M, Boström J, et al. 2015. Crystal structure, biochemical and cellular activities demonstrate separate functions of MTH1 and MTH2. *Nat Commun* 6(1):7871.
- Casadio R, Vassura M, Tiwari S, Fariselli P, Luigi Martelli P. 2011. Correlating disease-related mutations to their effect on protein stability: a large-scale analysis of the human proteome. *Hum Mutat* 32(10):1161–1170.
- Chen L, Brewer MD, Guo L, Wang R, Jiang P, Yang X. 2017. Enhanced degradation of misfolded proteins promotes tumorigenesis. *Cell Rep* 18(13):3143–3154.
- Cheng G, Qian B, Samudrala R, Baker D. 2005. Improvement in protein functional site prediction by distinguishing structural and functional constraints on protein family evolution using computational design. *Nucleic Acids Res* 33(18):5861–5867.
- Chiasson MA, Rollins NJ, Stephany JJ, Sitko KA, Matreyek KA, Verby M, Sun S, Roth F, DeSloover D, Marks DS, et al. 2020. Multiplexed measurement of variant abundance and activity reveals VKOR topology, active site and human variant impact. *eLife* 9:e58026.
- Choi Y, Chan AP. 2015. Proven web server: a tool to predict the functional effect of amino acid substitutions and indels. *Bioinformatics* 31(16):2745–2747.
- De Baets G, Van Durme J, Reumers J, Maurer-Stroh S, Vanhee P, Dopazo J, Schymkowitz J, Rousseau F. 2012. SNPeff 4.0: on-line prediction of molecular and structural effects of protein-coding variants. *Nucleic Acids Res* 40(Database issue):D935–D939.
- DePristo MA, Weinreich DM, Hartl DL. 2005. Missense meanderings in sequence space: a biophysical view of protein evolution. *Nat Rev Genet* 6(9):678–687.
- Després PC, Dubé AK, Seki M, Yachie N, Landry CR. 2020. Perturbing proteomes at single residue resolution using base editing. *Nat Commun* 11(1):1–13.
- Dunham A, Beltrao P. 2020. Exploring amino acid functions in a deep mutational landscape. *BioRxiv*. page 2020.05.26.116756.
- Echave J. 2019. Beyond stability constraints: a biophysical model of enzyme evolution with selection on stability and activity. *Mol Biol Evol* 36(3):613–620.
- Echave J, Spielman SJ, Wilke CO. 2016. Causes of evolutionary rate variation among protein sites. *Nat Rev Genet* 17(2):109–121.
- Ekeberg M, Hartonen T, Aurell E. 2014. Fast pseudolikelihood maximization for direct-coupling analysis of protein structure from many homologous amino-acid sequences. *J Comput Phys* 276:341–356.
- Ekeberg M, Lökvist C, Lan Y, Weigt M, Aurell E. 2013. Improved contact prediction in proteins: using pseudolikelihoods to infer Potts models. *Phys Rev E Stat Nonlin Soft Matter Phys* 87(1):012707.
- Fernandez-de Cossio-Diaz J, Uguzzoni G, Pagnani A. 2020. Unsupervised inference of protein fitness landscape from deep mutational scan. *Mol Biol Evol* 38(1):318–328.
- Ferrer-Costa C, Orozco M, de la Cruz X. 2002. Characterization of disease-associated single amino acid polymorphisms in terms of sequence and structure properties. *J Mol Biol* 315(4):771–786.
- Fersht A. 1999. Structure and mechanism in protein science: a guide to enzyme catalysis and protein folding. New York: W. H. Freeman.
- Gao M, Zhou H, Skolnick J. 2015. Insights into disease-associated mutations in the human proteome through protein structural analysis. *Structure* 23(7):1362–1369.
- Gerasimavicius L, Liu X, Marsh JA. 2020. Identification of pathogenic missense mutations using protein stability predictors. *Sci Rep* 10(1):15387.
- Goldenzweig A, Fleishman SJ. 2018. Principles of protein stability and their application in computational design. *Annu Rev Biochem* 87:105–129.
- Gray VE, Hause RJ, Fowler DM. 2017. Analysis of large-scale mutagenesis data to assess the impact of single amino acid substitutions. *Genetics* 207(1):53–61.

- Hietpas RT, Jensen JD, Bolon DN. 2011. Experimental illumination of a fitness landscape. *Proc Natl Acad Sci U S A*. 108(19):7896–7901.
- Hopf TA, Ingraham JB, Poelwijk FJ, Schärfe CP, Springer M, Sander C, Marks DS. 2017. Mutation effects predicted from sequence co-variation. *Nat Biotechnol*. 35(2):128–135.
- Ioannidis NM, Rothstein JH, Pejaver V, Middha S, McDonnell SK, Baheti S, Musolf A, Li Q, Holzinger E, Karyadi D, et al. 2016. Revel: an ensemble method for predicting the pathogenicity of rare missense variants. *Am J Hum Genet*. 99(4):877–885.
- Jack BR, Meyer AG, Echave J, Wilke CO. 2016. Functional sites induce long-range evolutionary constraints in enzymes. *PLoS Biol*. 14(5):e1002452.
- Jepsen MM, Fowler DM, Hartmann-Petersen R, Stein A, Lindorff-Larsen K. 2020. Classifying disease-associated variants using measures of protein activity and stability. In: Pey AL, editor. Protein homeostasis diseases. London, United Kingdom: Academic Press. p. 91–107.
- Jiang L, Mishra P, Hietpas RT, Zeldovich KB, Bolon DN. 2013. Latent effects of Hsp90 mutants revealed at reduced expression levels. *PLoS Genet*. 9(6):e1003600.
- Jimenez MJ, Arenas M, Bastolla U. 2018. Substitution rates predicted by stability-constrained models of protein evolution are not consistent with empirical data. *Mol Biol Evol*. 35(3):743–755.
- Kabsch W, Sander C. 1983. Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features. *Biopolymers* 22(12):2577–2637.
- Kampmeyer C, Nielsen SV, Clausen L, Stein A, Gerdes A-M, Lindorff-Larsen K, Hartmann-Petersen R. 2017. Blocking protein quality control to counter hereditary cancers. *Genes Chromosomes Cancer*. 56(12):823–831.
- Karran P, Attard N. 2008. Thiopurines in current medical practice: molecular mechanisms and contributions to therapy-related cancer. *Nat Rev Cancer*. 8(1):24–36.
- Kinney JB, McCandlish DM. 2019. Massively parallel assays and quantitative sequence–function relationships. *Annu Rev Genomics Hum Genet*. 20(1):99–127.
- Kircher M, Witten DM, Jain P, O’Roak BJ, Cooper GM, Shendure J. 2014. A general framework for estimating the relative pathogenicity of human genetic variants. *Nat Genet*. 46(3):310–315.
- Kumar P, Henikoff S, Ng PC. 2009. Predicting the effects of coding non-synonymous variants on protein function using the sift algorithm. *Nat Protoc*. 4(7):1073–1081.
- Lapedes A, Giraud B, Jarzynski C. 2012. Using sequence alignments to predict protein structure and stability with high accuracy. arXiv Preprint arXiv:1207.2484.
- Lee J-O, Yang H, Georgescu M-M, Di Cristofano A, Maehama T, Shi Y, Dixon JE, Pandolfi P, Pavletich NP. 1999. Crystal structure of the PTEN tumor suppressor: implications for its phosphoinositide phosphatase activity and membrane association. *Cell* 99(3):323–334.
- Liberles DA, Teichmann SA, Bahar I, Bastolla U, Bloom J, Bornberg-Bauer E, Colwell LJ, de Koning APJ, Dokholyan NV, Echave J, et al. 2012. The interface of protein structure, protein biophysics, and molecular evolution. *Protein Sci*. 21(6):769–785.
- Lichtarge O, Bourne HR, Cohen FE. 1996. An evolutionary trace method defines binding surfaces common to protein families. *J Mol Biol*. 257(2):342–358.
- Livesey BJ, Marsh JA. 2020. Using deep mutational scanning to benchmark variant effect predictors and identify disease mutations. *Mol Syst Biol*. 16:e9380.
- Lua RC, Wilson SJ, Konecki DM, Wilkins AD, Venner E, Morgan DH, Lichtarge O. 2016. UET: a database of evolutionarily-predicted functional determinants of protein sequences that cluster as functional sites in protein structures. *Nucleic Acids Res*. 44(D1):D308–D312.
- Lui S, Tiana G. 2013. The network of stabilizing contacts in proteins studied by coevolutionary data. *J Chem Phys*. 139(15):155103.
- Matreyek KA, Starita LM, Stephany JJ, Martin B, Chiasson MA, Gray VE, Kircher M, Khechaduri A, Dines JN, Hause RJ, et al. 2018. Multiplex assessment of protein variant abundance by massively parallel sequencing. *Nat Genet*. 50(6):874–882.
- Mayorov A, Dal Peraro M, Abriata LA. 2019. Active site-induced evolutionary constraints follow fold polarity principles in soluble globular enzymes. *Mol Biol Evol*. 36(8):1728–1733.
- Meacham GC, Patterson C, Zhang W, Younger JM, Cyr DM. 2001. The hsc70 co-chaperone chip targets immature cfr for proteasomal degradation. *Nat Cell Biol*. 3(1):100–105.
- Mighell TL, Evans-Dutson S, O’Roak BJ. 2018. A saturation mutagenesis approach to understanding PTEN lipid phosphatase activity and genotype-phenotype relationships. *Am J Hum Genet*. 102(5):943–955.
- Mirny LA, Shakhnovich EI. 1999. Universally conserved positions in protein folds: reading evolutionary signals about stability, folding kinetics and function. *J Mol Biol*. 291(1):177–196.
- Moriyama T, Nishii R, Lin T-N, Kihira K, Toyoda H, Nersting J, Kato M, Koh K, Inaba H, Manabe A, et al. 2017. The effects of inherited NUDT15 polymorphisms on thiopurine active metabolites in Japanese children with acute lymphoblastic leukemia. *Pharmacogenet Genomics*. 27(6):236–239.
- Moriyama T, Nishii R, Perez-Andreu V, Yang W, Klussmann FA, Zhao X, Lin T-N, Hoshitsuki K, Nersting J, Kihira K, et al. 2016. NUDT15 polymorphisms alter thiopurine metabolism and hematopoietic toxicity. *Nat Genet*. 48(4):367–373.
- Nielsen SV, Schenstrøm SM, Christensen CE, Stein A, Lindorff-Larsen K, Hartmann-Petersen R. 2020. Protein destabilization and degradation as a mechanism for hereditary disease. In: Protein homeostasis diseases. Elsevier. p. 111–125.
- Nielsen SV, Stein A, Dinitzen AB, Papaleo E, Tatham MH, Poulsen EG, Kasse M, Rasmussen LJ, Lindorff-Larsen K, Hartmann-Petersen R. 2017. Predicting the impact of lynch syndrome-causing missense mutations from structural calculations. *PLoS Genet*. 13(4):e1006739.
- Nishii R, Moriyama T, Janke LJ, Yang W, Suiter CC, Lin T-N, Li L, Kihira K, Toyoda H, Hofmann U, et al. 2018. Preclinical evaluation of NUDT15-guided thiopurine therapy and its effects on toxicity and antileukemic efficacy. *Blood* 131(22):2466–2474.
- Olzmann JA, Brown K, Wilkinson KD, Rees HD, Huai Q, Ke H, Levey AI, Li L, Chin L-S. 2004. Familial Parkinson’s disease-associated I166p mutation disrupts dj-1 protein folding and function. *J Biol Chem*. 279(9):8506–8515.
- Park H, Bradley P, Greisen P Jr, Liu Y, Mulligan VK, Kim DE, Baker D, DiMaio F. 2016. Simultaneous optimization of biomolecular energy functions on features from small molecules and macromolecules. *J Chem Theory Comput*. 12(12):6201–6212.
- Reeb J, Wirth T, Rost B. 2020. Variant effect predictions capture some aspects of deep mutational scanning experiments. *BMC Bioinformatics* 21(1):1–12.
- Relling MV, Schwab M, Whirl-Carrillo M, Suarez-Kurtz G, Pui C-H, Stein CM, Moyer AM, Evans WE, Klein TE, Antillon-Klussmann FG, et al. 2019. Clinical pharmacogenetics implementation consortium guideline for thiopurine dosing based on TPMT and NUDT 15 genotypes: 2018 update. *Clin Pharmacol Ther*. 105(5):1095–1105.
- Remmert M, Biegert A, Hauser A, Söding J. 2011. HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment. *Nat Methods*. 9(2):173–175.
- Ron I, Horowitz M. 2005. Er retention and degradation as the molecular basis underlying gaucher disease heterogeneity. *Hum Mol Genet*. 14(16):2387–2398.
- Rost B, Sander C. 1994. Conservation and prediction of solvent accessibility in protein families. *Proteins* 20(3):216–226.
- Rutherford SL, Lindquist S. 1998. Hsp90 as a capacitor for morphological evolution. *Nature* 396(6709):336–342.
- Sarkisyan KS, Bolotin DA, Meer MV, Usmanova DR, Mishin AS, Sharonov GV, Ivankov DN, Bozhanova NG, Baranov MS, Soylemez O, et al. 2016. Local fitness landscape of the green fluorescent protein. *Nature* 533(7603):397–401.
- Scheller R, Stein A, Nielsen SV, Marin FI, Gerdes A-M, Di Marco M, Papaleo E, Lindorff-Larsen K, Hartmann-Petersen R. 2019. Toward mechanistic models for genotype–phenotype correlations in phenylketonuria using protein stability calculations. *Hum Mutat*. 40(4):444–457.



- Shin H, Cho B-K. 2015. Rational protein engineering guided by deep mutational scanning. *Int J Mol Sci.* 16(9):23094–23110.
- Shoichet BK, Baase WA, Kuroki R, Matthews BW. 1995. A relationship between protein stability and protein function. *Proc Natl Acad Sci U S A.* 92(2):452–456.
- 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* 200(2):413–422.
- Stein A, Fowler DM, Hartmann-Petersen R, Lindorff-Larsen K. 2019. Biophysical and mechanistic models for disease-causing protein variants. *Trends Biochem Sci.* 44(7):575–588.
- Steward RE, MacArthur MW, Laskowski RA, Thornton JM. 2003. Molecular basis of inherited diseases: a structural perspective. *Trends Genet.* 19(9):505–513.
- Suiter CC, Moriyama T, Matreyek KA, Yang W, Scaletti ER, Nishii R, Yang W, Hoshitsuki K, Singh M, Trehan A, et al. 2020. Massively parallel variant characterization identifies NUDT15 alleles associated with thiopurine toxicity. *Proc Natl Acad Sci U S A.* 117:201915680.
- Tinberg CE, Khare SD, Dou J, Doyle L, Nelson JW, Schena A, Jankowski W, Kalodimos CG, Johnsson K, Stoddard BL, et al. 2013. Computational design of ligand-binding proteins with high affinity and selectivity. *Nature* 501(7466):212–216.
- Tokuriki N, Tawfik DS. 2009. Chaperonin overexpression promotes genetic variation and enzyme evolution. *Nature* 459(7247):668–673.
- Touw WG, Baakman C, Black J, Te Beek TA, Krieger E, Joosten RP, Vriend G. 2015. A series of PDB-related databanks for everyday needs. *Nucleic Acids Res.* 43(Database issue):D364–D368.
- Valerie NCK, Hagenkort A, Page BDG, Masuyer G, Rehling D, Carter M, Bevc L, Herr P, Homan E, Sheppard NG, et al. 2016. NUDT15 hydrolyzes 6-thio-deoxyGTP to mediate the anticancer efficacy of 6-thioguanine. *Cancer Res.* 76(18):5501–5511.
- Valiente M, Andrés-Pons A, Gomar B, Torres J, Gil A, Tapparel C, Antonarakis SE, Pulido R. 2005. Binding of PTEN to specific PDZ domains contributes to PTEN protein stability and phosphorylation by microtubule-associated serine/threonine kinases. *J Biol Chem.* 280(32):28936–28943.
- Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, Burovski E, Peterson P, Weckesser W, Bright J, et al. 2020. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat Methods.* 17(3):261–272.
- Wagih O, Galardini M, Busby BP, Memon D, Typas A, Beltrao P. 2018. A resource of variant effect predictions of single nucleotide variants in model organisms. *Mol Syst Biol.* 14(12):e8430.
- Wang Z, Moulton J. 2001. Snps, protein structure, and disease. *Hum Mutat.* 17(4):263–270.
- Ward JH Jr. 1963. Hierarchical grouping to optimize an objective function. *J Am Stat Assoc.* 58(301):236–244.
- Weile J, Roth FP. 2018. Multiplexed assays of variant effects contribute to a growing genotype–phenotype atlas. *Hum Genet.* 137(9):665–678.
- Yaguchi H, Ohkura N, Takahashi M, Nagamura Y, Kitabayashi I, Tsukada T. 2004. Menin missense mutants associated with multiple endocrine neoplasia type 1 are rapidly degraded via the ubiquitin-proteasome pathway. *Mol Cell Biol.* 24(15):6569–6580.
- Yang C, Asthagiri AR, Iyer RR, Lu J, Xu DS, Ksendzovsky A, Brady RO, Zhuang Z, Lonser RR. 2011. Missense mutations in the NF2 gene result in the quantitative loss of merlin protein and minimally affect protein intrinsic function. *Proc Natl Acad Sci U S A.* 108(12):4980–4985.
- Yang C, Huntoon K, Ksendzovsky A, Zhuang Z, Lonser RR. 2013. Proteostasis modulators prolong missense VHL protein activity and halt tumor progression. *Cell Rep.* 3(1):52–59.
- Yang S-K, Hong M, Baek J, Choi H, Zhao W, Jung Y, Haritunians T, Ye BD, Kim K-J, Park SH, et al. 2014. A common missense variant in NUDT15 confers susceptibility to thiopurine-induced leukopenia. *Nat Genet.* 46(9):1017–1020.
- Yehia L, Ngeow J, Eng C. 2019. PTEN-opathies: from biological insights to evidence-based precision medicine. *J Clin Invest.* 129(2):452–464.
- Yue P, Li Z, Moulton J. 2005. Loss of protein structure stability as a major causative factor in monogenic disease. *J Mol Biol.* 353(2):459–473.