

**Supplemental Table 1** A survey of AAS prediction methods and their observations.

Groups with key observations have been highlighted in bold.

Group	Observations
Wang and Moulton (2001) (77)	<ul style="list-style-type: none"> <li>● Observed that disease mutations can be distinguished from neutral SNPs using structure.</li> <li>● Modeled amino acid change on to structure and observed that 83% of disease SNPs affect protein stability; 5% affect ligand binding.</li> </ul>
Sunyaev et al. (2000, 2001, 2002) (53, 68, 69)	<ul style="list-style-type: none"> <li>● Prediction method incorporates sequence conservation, structure, and annotation from Swiss-Prot database for prediction.</li> <li>● ~20% of nsSNPs estimated to be damaging to protein function. To resolve the large number of damaging nsSNPs with the fecundity of inbred marriages, damaging nsSNPs must have mild effects.</li> </ul>
Chasman and Adams (2001) (9)	<ul style="list-style-type: none"> <li>● Accessibility, <i>B</i>-factor and sequence conservation are the most useful for prediction according to ANOVA, PCA and correlation analysis. Implemented a prediction method based on the three features.</li> <li>● A protein structure that has <math>\geq 60\%</math> sequence identity to the input protein gives the best performance, with no increase in performance as sequence identity increases.</li> </ul>
Ng and Henikoff (2001, 2002, 2003) (45--47)	<ul style="list-style-type: none"> <li>● Prediction method uses sequence conservation and position-specific scoring matrices with Dirichlet priors to model the allowed amino acids at a position.</li> <li>● Using an AAS prediction method has better accuracy than an amino acid substitution scoring matrix.</li> <li>● The fraction of nsSNPs predicted to be damaging is comparable to the false positive error so the number of damaging nsSNPs is low and cannot be estimated.</li> </ul>
Ferrer-Costa et al. (16--18) (2002, 2004, 2005)	<ul style="list-style-type: none"> <li>● Sequence-based prediction method. Structural properties (secondary structure and accessibility) are predicted and used for prediction. A neural network provides the final prediction.</li> </ul>
Saunders and Baker (2002) (59)	<ul style="list-style-type: none"> <li>● Found that the most accurate predictions are obtained using a combination of sequence and structural features.</li> <li>● Estimated <math>C^\beta</math> density values from ab initio structure prediction which can be beneficial when there are few sequences available (<math>\leq 3</math> homologues).</li> </ul>
Terp et al. (2002) (72)	<ul style="list-style-type: none"> <li>● Studied 20 biophysical parameters, 9 of which were significant for prediction.</li> <li>● Score includes the likelihood that a mutation will be observed clinically. The disease-causing mutation is severe enough to be included in a disease database, but not preclinically lethal.</li> </ul>
Mooney et al. (2003, 2003) (41, 42)	<ul style="list-style-type: none"> <li>● Using sequence homology, calculate negative entropy to measure conservation.</li> <li>● Found degree of conservation scores for mutations are correlated with the severity of the phenotype for syndromes caused by mutations in the androgen receptor.</li> </ul>

Group	Observations
Stitzel et al. (2003, 2004) (64, 65)	<ul style="list-style-type: none"> <li>●TopoSNP uses sequence conservation and structure. Structure is used to classify if the AAS is located on a surface, in a pocket, or buried. Under this classification, only 3% of disease-associated mutations were buried.</li> <li>●Using structure, 88% of disease-associated mutations are in pockets but 68% of nondisease SNPs are also in these pockets.</li> </ul>
Krishnan and Westhead (2003) (31)	<ul style="list-style-type: none"> <li>●Implemented 2 different machine-learning methods: decision trees and support vector machines.</li> <li>●Methods use sequence and structure. Structural attributes (secondary structure and solvent accessibility) are predicted and homologous structure is not necessary.</li> </ul>
Thomas et al. (2003) (73)	<ul style="list-style-type: none"> <li>●Position-specific evolutionary conservation scores calculated from Hidden Markov Models in the PANTHER library.</li> <li>●Some disease-causing mutations predicted to be gain-of-function rather than loss-of-function.</li> </ul>
del sol Mesa et al. (2003) (14)	<ul style="list-style-type: none"> <li>●Implemented three different methods to partition a protein family into subfamilies for better prediction of positions involved in functional specificity.</li> </ul>
Fleming et al. (2003) (19)	<ul style="list-style-type: none"> <li>●Identify conserved sites through sliding window.</li> <li>●Incorporates DNA sequence to identify sites evolving under positive selection.</li> </ul>
Santib���ez-Koref et al. (2003) (58)	<ul style="list-style-type: none"> <li>●Uses sequence conservation and takes into account phylogeny.</li> <li>●Altering tree structure decreases performance; altering branch lengths does not affect performance as much.</li> </ul>
Herrgard et al. (2003) (27)	<ul style="list-style-type: none"> <li>●Using sequence and structure, this prediction method focuses on mutations in the active site of an enzyme to find residues that affect catalytic ability rather than protein stability.</li> <li>●Average prediction accuracy for deleterious mutations 85%, prediction for nondeleterious mutations is 81%.</li> </ul>
Cai et al. (2004) (7)	<ul style="list-style-type: none"> <li>●Bayesian network evaluating only positions residing inside a PFAM domain. Also uses structural annotation provided by SWISS-PROT.</li> </ul>
Lau and Chasman (2004)(33)	<ul style="list-style-type: none"> <li>●From a set of sequences, chooses the optimal sequence subalignment which gives as many tolerated amino acids but excludes proteins that are functionally divergent from the query protein.</li> <li>●Performs better than SIFT.</li> </ul>
Balasubramanian et al. (2005)(2)	<ul style="list-style-type: none"> <li>●Uses logistic regression analysis based on sequence and structural features to obtain high prediction accuracy on G-protein coupled receptors.</li> </ul>
Stone and Sidow (2005) (67)	<ul style="list-style-type: none"> <li>●Sequence-based MAPP method, which performs better than SIFT.</li> <li>●Takes a protein alignment and a tree for input.</li> <li>●Quantified the benefit of using orthologues instead of paralogues in sequence alignment.</li> </ul>

Group	Observations
Yue and Moulton (2005) (82, 83)	<ul style="list-style-type: none"><li>● Two methods based on support vector machine. One is sequence-only, one is structure-only.</li><li>● Larger change in free energy is correlated with a higher fraction of disease mutations.</li><li>● SwissProt functional annotation decreases overall accuracy.</li><li>● Protein models based on structure with <math>\geq 40\%</math> sequence identity has comparable performance to using the structure of the input protein.</li><li>● Sequence-based method has better overall performance compared to structure.</li><li>● When number of sequences is <math>&lt; 10</math>, false positive rate increases but false negative rate stays the same.</li></ul>