

Blind prediction of deleterious amino acid variations with SNPs&GO

Emidio Capriotti¹ | Pier Luigi Martelli¹ | Piero Fariselli² | Rita Casadio¹

¹Biocomputing Group, BiGeA / Giorgio Prodi Interdepartmental Center for Cancer Research, University of Bologna University of Bologna, Bologna, Italy

²Department of Comparative Biomedicine and Food Science, University of Padova, Legnaro, Padova, Italy

Correspondence

Emidio Capriotti, BioFold Unit, Department of Biological, Geological, and Environmental Sciences (BiGeA), University of Bologna, Via F. Selmi 3, Bologna, 40126, Italy.
Email: emidio.capriotti@unibo.it

Availability: SNPs&GO is accessible at <http://snps.biofold.org/snps-and-go> or <http://snps-and-go.biocomp.unibo.it>

Contract grant sponsors: Italian MIUR (PRIN 2010-2011, project 20108XYHJS); European Union RTD Framework Program (COST BMBS Action TD1101, Action BM1405); Italian MIUR (PON projects PON01_02249, PAN Lab PONa3_00166); FARB UNIBO 2012.

For the CAGI Special Issue

Abstract

SNPs&GO is a machine learning method for predicting the association of single amino acid variations (SAVs) to disease, considering protein functional annotation. The method is a binary classifier that implements a support vector machine algorithm to discriminate between disease-related and neutral SAVs. SNPs&GO combines information from protein sequence with functional annotation encoded by gene ontology (GO) terms. Tested in sequence mode on more than 38,000 SAVs from the SwissVar dataset, our method reached 81% overall accuracy and an area under the receiving operating characteristic curve of 0.88 with low false-positive rate. In almost all the editions of the Critical Assessment of Genome Interpretation (CAGI) experiments, SNPs&GO ranked among the most accurate algorithms for predicting the effect of SAVs. In this paper, we summarize the best results obtained by SNPs&GO on disease-related variations of four CAGI challenges relative to the following genes: *CHEK2* (CAGI 2010), *RAD50* (CAGI 2011), *p16-INK* (CAGI 2013), and *NAGLU* (CAGI 2016). Result evaluation provides insights about the accuracy of our algorithm and the relevance of GO terms in annotating the effect of the variants. It also helps to define good practices for the detection of deleterious SAVs.

KEYWORDS

disease-related variation, gene ontology, genome interpretation, machine learning, protein function, single amino acid variation, variant annotation

1 | INTRODUCTION

Large-scale genomic experiments are generating a huge amount of genetic variants whose effect is still unknown (Capriotti, Nehrt, Kann, & Bromberg, 2012). Among all possible genetic alterations, single-nucleotide variants (SNVs) are the most frequent type of variations between individual genomes (Durbin et al., 2010) and nonsynonymous SNVs (inducing single amino acid variations in the encoded protein) are the variant class most frequently associated with disease. Despite the improvements in the characterization of the human genome, the relationship between genotype and phenotype is still an open problem. In this context, the development of more accurate methods for the detection and annotation of SNVs becomes one of the key challenges for personalized medicine (Fernald, Capriotti, Daneshjoui, Karczewski, & Altman, 2011). During the last few years, several initiatives have been established to promote, disseminate, and evaluate research in the field of disease-associated phenomics. International consortiums have collected data from thousands of individuals for defining functional regions of the human genome (Durbin et al., 2010; ENCODE Project Consortium, 2012) and for characteriz-

ing the landscape of genetic alterations associated to human pathologies (Cancer Genome Atlas Research Network, et al., 2013; International Cancer Genome Consortium, et al., 2010). At the same time, many meetings contributed to the dissemination of the increasing number of computational methods (Niroula & Vihinen, 2016) for the identification and annotation of the genetic variants (Bromberg, Capriotti, & Carter, 2016; Oetting, 2011). Finally, in silico experiments with different computational challenges were organized to evaluate the available tools for predicting the impact of genetic variants and/or the association between genotype and phenotype (Brownstein et al., 2014; Saez-Rodriguez et al., 2016). Among the computational experiments, the Critical Assessment for Genome Interpretation (CAGI) provided several blind datasets for testing the accuracy of the predictive algorithms (<https://genomeinterpretation.org/>). The Bologna Biocomputing Group and the BioFold Unit, as active members of this community, participated in all the CAGI editions since 2010 submitting predictions for many challenges adopting SNPs&GO (Calabrese, Capriotti, Fariselli, Martelli, & Casadio, 2009; Capriotti et al., 2013). SNPs&GO is a support vector machine-based approach to predict the impact of single amino acid variations (SAVs). Our method

takes in input information extracted from the protein sequence profile and functional information encoded through the gene ontology (GO) terms. In a previous independent evaluation, SNPs&GO was scored among the most accurate methods for predicting the impact of SAVs (Thusberg, Olatubosun, & Vihinen, 2011). In this work, we analyze the best predictions submitted using two versions of SNPs&GO, trained on data sets of different size and performing among the state-of-the-art predictors (Calabrese et al., 2009; Capriotti et al., 2013). The assessment of the results of the four challenges of the CAGI experiments confirmed that SNPs&GO consistently scores among the best methods for predicting the impact of SAVs.

2 | MATERIAL AND METHODS

2.1 | SNPs&GO predictions

SNPs&GO is a support vector machine-based approach that takes in input information from protein sequence and function. SNPs&GO internally runs a BLAST (Altschul et al., 1997) search against the UniRef90 database (Suzek, Huang, McGarvey, Mazumder, & Wu, 2007) to build the protein sequence profile. Functional information encoded by GO terms are extracted from UniProt database (Magrane & UniProt Consortium, 2011). For each GO term, all the human proteins reported in SwissVar database (Mottaz, David, Veuthey, & Yip, 2010) are collected and a log-odd score (LGO) is calculated as the logarithm of the fraction of disease and neutral SAVs. Thus, the functional score of each protein is obtained by summing the LGO values of the associated GO terms and their parents in the GO-rooted graph. The SNPs&GO functional score contributes to the performance of our method providing an empirical estimation of the probability of having a deleterious SAV in a protein, given the associated GO terms.

The prediction output of SNPs&GO is a score ranging between 0 and 1 that represents the probability of a SAV to be pathogenic. By construction, a threshold (t) of 0.5 is selected to discriminate between benign ($t \leq 0.5$) and pathogenic ($t > 0.5$) SAVs. Depending on the score, a reliability index (RI) ranging from 0 to 10 is defined to estimate the level of confidence of the prediction. In this paper, we considered two versions of SNPs&GO: the first version (SNPs&GO⁰⁹) implemented before 2009 (Calabrese et al., 2009) used by the Biocomputing Group and the updated version (SNPs&GO¹³) used and maintained by the Bio-Fold Unit (Capriotti et al., 2013). With respect to the older version of SNPs&GO, the new one has been trained on an updated version of the SwissVar database (Mottaz et al., 2010), including ~4,700 more SAVs (~14%). Furthermore, the conservation and functional scores are calculated using the updated versions of the UniRef90 database and GO that correspond to ~8,900 more sequences with at least one associated GO term (32%).

2.2 | CHEK2 challenge (CAGI 2010)

For the *CHEK2* challenge, predictors were asked to classify variants as occurring in breast cancer cases or controls and to provide an estimation of the probability of a given variant to be in the case set (f_{case}).

We focused our analysis on the subset of 32 SAVs (MUT-CHEK2). We predicted the probability f_{case} with SNPs&GO⁰⁹ (f_{case}^p), considering both the binary prediction (disease/neutral) and the RI; predictions were transformed into probability with a linear function so that $f_{\text{case}}^p = 1$ corresponds to disease predictions with RI = 10, and $f_{\text{case}}^p = 0$ corresponds to neutral predictions with RI = 10. The list of MUT-CHEK2 variants with the experimental values of f_{case} (f_{case}^e) was released (Le Calvez-Kelm et al., 2011), and it is reported in Supp. Table S1, along with predictions performed with SNPs&GO⁰⁹, SIFT (Ng & Henikoff, 2003), and AlignGVGD (Mathe et al., 2006). To evaluate the quality of the predictions, we transformed the experimental f_{case} (f_{case}^e) in a binary classification (pathogenic/benign), by applying a threshold equal to 0.7 (which represents the median of the optimal f_{case}^e using the default prediction thresholds). If $f_{\text{case}}^e > 0.7$, the variation is classified as pathogenic, otherwise benign (see Supp. Materials). For the predicted f_{case} (f_{case}^p), the thresholds were selected by maximizing the performance of each method (see Supp. Materials). With this assumption, the MUT-CHEK2 dataset is divided, on the basis of f_{case}^e , in 21 pathogenic and 11 benign SAVs, and the performance of the algorithms was calculated using the standard evaluation measures for binary classifiers (see Supp. Materials). For the *CHEK2* challenge, we compared the performance of SNPs&GO⁰⁹ (Calabrese et al., 2009) with SIFT (Ng & Henikoff, 2003) and AlignGVGD (Mathe et al., 2006), which have been used by the assessors as baseline methods. More information about the *CHEK2* challenge is available in Supp. Materials and at <http://goo.gl/2Wlr6M>.

2.3 | RAD50 dataset (CAGI 2011)

As in the case of *CHEK2*, for this challenge, too, SNPs&GO⁰⁹ was used to predict the probability of each variant to be in the case set. With SNPs&GO, we scored the pathogenicity of 35 SAVs (MUT-RAD50) carried by up to 20 individuals. The MUT-RAD50 list of variations and the associated predictions are reported in Supp. Table S2. This list of variants has been released in a recent publication (Damiola et al., 2014). As we did for the *CHEK2* challenge, we classified each variant according to the fraction of carriers in the case set (f_{case}^e) defined in Supp. Eq. S3. Using a threshold of 0.7, the MUT-RAD50 set splits in 17 pathogenic and 18 benign missense SNVs. More information about the *RAD50* challenge is available in Supp. Materials and at <http://goo.gl/y4nwl1>.

2.4 | p16INK4A challenge (CAGI 2013)

For the *p16* challenge in CAGI 2013, predictors were asked to estimate the proliferation rates (p) of mutation-like cells. Considering experimental results, a score of 0.50 was assigned to samples with same proliferation rate as the control; variations leading to an increase or decrease of the proliferation rate are labeled with a score higher (up to 1) or lower (down to 0) than 0.5, respectively. We predicted the proliferation rates with SNPs&GO¹³, using the raw output of the method, which represents the probability of a variant to be related to disease. The list of variations and the associated predictions are reported in Supp. Table S3. The data providers also included a set 19 proliferation rates from mutation-like cells as possible training set

(TRAIN-P16). For the *p16* challenge, we compared the prediction submitted by the BioFolD Unit using SNPs&GO¹³ and DrCancer (Capriotti & Altman, 2011) with the most accurate prediction in the CAGI assessment, developed by the SPARKS-Lab (<http://sparks-lab.org/>), and implementing a method specifically optimized on the TRAIN-P16 dataset. More information about the *p16* challenge is available in Supp. Materials and at <http://goo.gl/51hGuZ>.

2.5 | NAGLU challenge (CAGI 2016)

For the NAGLU challenge, CAGI 2016 participants were asked to predict the relative change in enzymatic activity (RelAct) associated to each SAV. In this paper, we perform the a posteriori comparison of the submitted predictions obtained with SNPs&GO⁰⁹ (Calabrese et al., 2009) with the most accurate predictions in the CAGI assessment, performed with MutPred (Li et al., 2009). In this analysis, we include the new predictions from the last version of SNPs&GO¹³ (Capriotti et al., 2013), which were not submitted to the CAGI. The list of NAGLU amino acid variations and the associated predictions are reported in Supp. Table S4. More information about the NAGLU challenge is available in the Supp. Materials and at <http://goo.gl/wp17aB>.

2.6 | Comparison with other methods

In this study, we compared two versions of SNPs&GO (SNPs&GO⁰⁹ and SNPs&GO¹³) with other computational methods. In detail, for CHEK2 and RAD50 challenges, we compared SNPs&GO⁰⁹ predictions submitted by the Biocomputing Group with AlignGVGD (Mathe et al., 2006) and SIFT (Ng & Henikoff, 2003). Align-GVGD, which has been used by the assessor as baseline method, is a program that combines the biophysical characteristics of amino acids and protein multiple sequence alignments. It is based on the calculation of Grantham score (Grantham, 1974) on a multiple sequence alignment. AlignGVGD classifies SAVs in seven classes from C0 to C65, which correspond to the lowest and highest level of enrichment for pathogenic variants. For the AlignGVGD predictions, we used the precalculated multiple sequence alignments including all the sequences from *Homo sapiens* to *Sea urchin* (see <http://agvgd.hci.utah.edu/>).

SIFT is one of the most popular tools for scoring the impact of genetic variants based on sequence homology. The algorithm is based on the assumption that important amino acids will be conserved in the protein family, and changes at well-conserved positions tend to be predicted as deleterious. SIFT returns a probabilistic score ranging from 0 to 1, which represents the normalized probability that an amino acid change is tolerated. In standard predictions, variations with score below 0.05 are classified as pathogenic. The predictions from SIFT algorithm were calculated using the Web server <http://sift.bii.a-star.edu.sg/> with default parameters.

Although AlignGVGD and SIFT are not among the most updated tools currently available for predicting the impact of the genetic variations, we included them in our analysis as baseline methods to compare with SNPs&GO. This is in agreement with the procedure followed by the assessor of CHEK2 and RAD50 challenges, who selected AlignGVGD as reference for benchmarking different predictors.

For the *p16INKA4* challenge, we compared the predictions of SNPs&GO¹³ and DrCancer (Capriotti & Altman, 2011) submitted by the BioFolD Unit with those from an ad hoc method implemented by the SPARK-LAB. DrCancer is a modification of the SNPs&GO algorithm that is based on the slim version of the GO (<http://geneontology.org/page/go-slim-and-subset-guide>). The disease-specific method has been trained and tested on a set of more than 3,000 cancer-causing variants. Similar to SNPs&GO, DrCancer returns in output a score from 0 to 1 representing the probability of SAVs of being cancer causing. The SPARK-LAB method used SVM with linear kernel trained on the TRAIN-P16 dataset. The input features of the algorithm include a combination of the position-specific scoring matrix values for wild-type and mutant residues and the predicted free energy change upon SAV computed by ROSETTA3 (Leaver-Fay et al., 2011) and dMutant (Zhou & Zhou, 2002).

For the NAGLU challenge, only the binary predictions derived from SNPs&GO⁰⁹ were officially submitted by the Bologna Biocomputing Group. To better evaluate the accuracy of our algorithm, we compared the predictions from SNPs&GO⁰⁹ with those from the latest version of SNPs&GO (SNPs&GO¹³) maintained by the BioFolD Unit and two versions of MutPred2 algorithm (Li et al., 2009). In detail, for MutPred2, we considered the predictions of the algorithm running in default mode (MutPred2) and the predictions without gene-level homology count features (MutPred2*). MutPred2 is a machine-learning approach based on an ensemble of neural networks trained on a combination of features including the SIFT output, conservation scores, and predicted structural and functional residue properties. Similar to SNPs&GO, MutPred2 output represents the probability that the amino acid substitution is deleterious.

For the NAGLU challenge, SNPs&GO¹³ and MutPred2 predictions were obtained subtracting the raw outputs to one.

2.7 | Prediction evaluation

The evaluation of the accuracy of computational methods for variant annotation is a difficult task whose solution depends on the complexity of the prediction. For the CAGI challenges here discussed, we use two evaluation systems. The first evaluation is based on the regression between the experimental and predicted values (r_{Pearson}) and their ranking (r_{Spearman} , $r_{\text{KendallTau}}$). For this test, the root mean square error (RMSE) after linear fitting is also calculated. The second evaluation is based on the standard evaluation measures for binary classifiers suggested in recent papers (Vihinen, 2012, 2013). They are: true-positive rate (TPR) and true-negative rate (also referred as sensitivity and specificity), positive and negative predicted values (PPV, NPV), overall accuracy (Q_2), Matthews correlation coefficient (MC), and area under the receiver operating characteristic curve (AUC). The thresholds for the classification of the experimental and predicted data were optimized for each challenge. More details about the evaluation measures and classification thresholds used for the evaluation of the CHEK2, RAD50, *p16*, and NAGLU challenges are described in Supp. Materials.

TABLE 1 Performance of the predictors for the *CHEK2* challenge (CAGI 2010)

Method	Q ₂	TPR	PPV	TNR	NPV	AUC	MC	RMSE	r _{Pearson}	r _{Spearman}	r _{KendallTau}
SNPs&GO ⁰⁹	0.72	0.81	0.77	0.55	0.60	0.73	0.36	0.46	0.29	0.32	0.25
SIFT	0.69	0.95	0.69	0.18	0.67	0.53	0.22	0.43	0.19	0.10	0.08
AlignGVGD	0.66	0.67	0.78	0.64	0.50	0.70	0.29	0.67	0.32	0.26	0.25

Notes: The overall accuracy (Q₂), true-positive/negative rates (TPR/TNR), positive/negative predicted values (PPV/NPV), area under the ROC curve (AUC), and Matthews correlation coefficient (MC) are calculated using an f_{case}^e threshold of 0.70 for dividing cases from controls. The positive and negative classes refer to pathogenic and benign variations, respectively. For SNPs&GO⁰⁹, SIFT, and AlignGVGD, the f_{case}^p thresholds are 0.35, 0.60, and C0, respectively. All the binary (Q₂, TPR, TNR, PPV, NPV, AUC, and MC) and regression (RMSE, r_{Pearson}, r_{Spearman}, and r_{KendallTau}) evaluation measures are described in Supp. Material. The confusion matrices for calculating the performance of the binary classifiers are reported in Supp. Table S5. Bold correlation coefficients correspond to *P* values < 0.05. The SNPs&GO⁰⁹ *P* value for the Spearman test is 0.07.

TABLE 2 Performance of the predictors for the *RAD50* challenge (CAGI 2011)

Method	Q ₂	TPR	PPV	TNR	NPV	AUC	MC	RMSE	r _{Pearson}	r _{Spearman}	r _{KendallTau}
SNPs&GO ⁰⁹	0.66	0.41	0.78	0.89	0.62	0.64	0.34	0.64	0.27	0.32	0.28
SNPs&GO ^{09*}	0.73	0.57	1.00	1.00	0.57	0.82	0.57	0.66	0.43	0.62	0.56
SIFT	0.63	0.65	0.61	0.61	0.65	0.57	0.26	0.49	0.12	0.23	0.19
AlignGVGD	0.57	0.12	1.00	1.00	0.55	0.55	0.25	0.68	0.08	0.08	0.07

Notes: SNPs&GO^{09*} refers to the performance on the subset of 11 amino acid variations in the Zn hook and P-loop hydrolase domains. The overall accuracy (Q₂), true-positive/negative rates (TPR/TNR), positive/negative predicted values (PPV/NPV), area under the ROC curve (AUC), and Matthews correlation coefficient (MC) are calculated using an f_{case}^e threshold of 0.70 for dividing cases from controls. The positive and negative classes refer to pathogenic and benign variations, respectively. For SNPs&GO⁰⁹, SIFT, and AlignGVGD, the f_{case}^p thresholds are 0.10, 0.15, and C45, respectively. All the binary (Q₂, TPR, TNR, PPV, NPV, AUC, and MC) and regression (RMSE, r_{Pearson}, r_{Spearman}, and r_{KendallTau}) evaluation measures are described in Supp. Material. The confusion matrices for calculating the performance of the binary classifiers are reported in Supp. Table S5. Bold correlation coefficients correspond to *P* values < 0.05. The SNPs&GO⁰⁹ *P* value for the Spearman test is 0.06.

3 | RESULTS

3.1 | *CHEK2* and *RAD50* challenges

The *CHEK2* and *RAD50* challenges run in the first two editions of the CAGI experiments. For these challenges, the predictors were asked to estimate the probability of the carrier of a specific SAV to be in the case set (f_{case}). The predictions were evaluated by Sean Tavtigian (University of Utah), who also provided the experimental data for both challenges. According to his assessment, we compared the prediction performed with SNPs&GO⁰⁹ with those performed with AlignGVGD and SIFT, by estimating the evaluation measures for binary classification (Q₂, FPR, TPR, NPV, PPV, AUC, MC) and regression (RMSE, r_{Pearson}, r_{Spearman}, r_{KendallTau}) described in Supp. Material. The performances of the three predictors for the *CHEK2* and *RAD50* challenges are summarized in Tables 1 and 2. SNPs&GO⁰⁹ resulted in better performance than SIFT and AlignGVGD in the regression tests (RMSE, r_{Pearson}, r_{Spearman}, r_{KendallTau}). Although all the predictors achieved relatively low correlation coefficient values, SNPs&GO is the only one scoring with a consistently significant r_{Kendall/Tau} (*P* value < 0.05). It must be noted that the experimental values of f_{case}^e are biased toward the extreme values: SAVs with f_{case}^e either equal to 0 or 1 correspond to 78% and 74% of the *CHEK2* and *RAD50* datasets, respectively. This bias can hamper the estimation of the correlation coefficients.

In a second test, we evaluated the performances of SNPs&GO⁰⁹, SIFT, and AlignGVGD as binary classifiers. For each method, we transformed the probability predictions into classes by optimizing the separating threshold. For each method and challenge, the threshold is the value maximizing the product among Q₂, AUC, and MC, as described in

Supp. Materials. With this procedure, SNPs&GO⁰⁹ reaches a good performance on the *CHEK2* dataset showing an Q₂ of 72%, a MC of 0.36 and an AUC of 0.73 when the output threshold is set to 0.35.

For the *RAD50* challenge, SNPs&GO⁰⁹ shows better performance than the other methods, and the performance becomes significantly better when we focus on the variations in the Zn hook and P-loop hydrolase domains. On this subset of 11 SAVs, SNPs&GO⁰⁹ achieves good performances both in the binary classification and regression tests. For the *RAD50* challenge, SIFT resulted in better performance than AlignGVGD in terms of Q₂, but both methods showed AUCs close to those of the random predictors.

3.2 | *p16* challenge

For the *p16* challenge, predictors were asked to estimate the proliferation rate of mutation-like cells with respect to wild-type cells (RelPro). In this experiment, a prediction near 0.5 indicates a proliferation rate similar to wild-type cell, whereas values close to 1 are associated to the highest proliferation rates in mutated cells. Here, we compared the predictions of SNPs&GO¹³ and DrCancer submitted by the BioFold Unit with the most successful predictions submitted by the SPARK-LAB. With this comparison, we show that the automatic methods (SNPs&GO¹³ and DrCancer) can achieve similar level of accuracy with respect to the SPARK-LAB algorithm, which has been specifically developed for the *p16* challenge. Our comparison, based on a regression test (Table 3), reveals that SPARK-LAB predictions achieved better correlation coefficients. In detail, SPARK-LAB results in 0.16 better r_{Pearson} and r_{Spearman}, with respect to SNPs&GO¹³. The

TABLE 3 Performance of the predictors for the *p16* challenge (CAGI 2013)

Method	Q_2	TPR	PPV	TNR	NPV	AUC	MC	RMSE	r_{Pearson}	r_{Spearman}	$r_{\text{KendallTau}}$
SPARK-LAB	0.90	0.80	1.00	1.00	0.83	0.92	0.82	76	0.83	0.87	0.69
SNPs&GO ¹³	0.70	1.00	0.63	0.40	1.00	0.88	0.50	76	0.66	0.81	0.60
DrCancer	0.60	1.00	0.56	0.20	1.00	0.84	0.33	76	0.58	0.67	0.47

Notes: The overall accuracy (Q_2), true-positive/negative rates (TPR/TNR), positive/negative predicted values (PPV/NPV), area under the ROC curve (AUC), and Matthews correlation coefficient (MC) are calculated using an experimental relative proliferation (RelPro) rate threshold of 75 and a predicted threshold of 0.75. The positive and negative classes refer to pathogenic and benign variations, respectively. All the binary (Q_2 , TPR, TNR, PPV, NPV, AUC, and MC) and regression (RMSE, r_{Pearson} , r_{Spearman} , and $r_{\text{KendallTau}}$) evaluation measures are described in Supp. Material. The confusion matrices for calculating the performance of the binary classifiers are reported in Supp. Table S5. Bold correlation coefficients correspond to P values < 0.05.

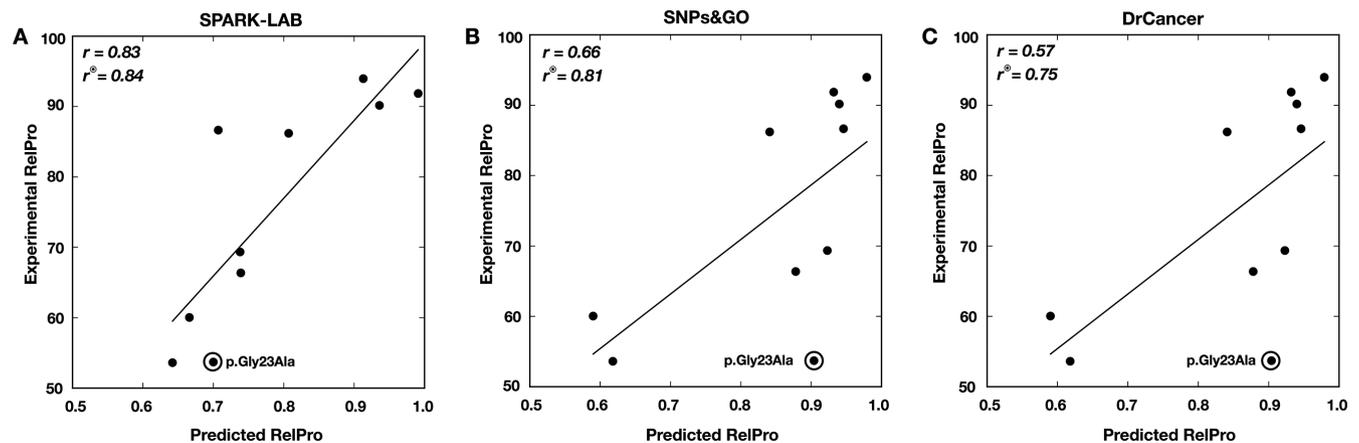


FIGURE 1 Comparison between predicted and experimental relative proliferation (RelPro) rates for the *p16* challenge. Linear regression for SPARK-LAB (A), SNPs&GO¹³ (B), and DrCancer (C) predictions. r and r^0 are the Pearson's correlation coefficients with and without the amino acid variation p.Gly23Ala, respectively

difference in the value of $r_{\text{KendallTau}}$ is ~ 0.09 . After plotting the linear regression curves between predicted and experimental values (Figure 1), we noticed that the difference in the performances is mainly due to the wrong prediction of the amino acid variation p.Gly23Ala. As shown in Figure 1, removing prediction of the amino acid variation p.Gly23Ala in the calculation, the r_{Pearson} values, the SPARK-LAB method, and SNPs&GO¹³ differ by 0.02. According to the suggestion of CAGI assessors, the predictors were also evaluated as binary classifiers (Carraro et al. 2017). In Table 3, we report the performance considering all predictions with score higher than 0.75 as deleterious variants. With this assumption, we observed a decreasing level of accuracy going from SPARK-LAB to DrCancer predictions. Despite the differences in the scores, it is still remarkable that a general method like SNPs&GO resulted in a good level of performance with respect to the problem-specific method developed by the SPARK-LAB. The analysis of the assessors showed SNPs&GO and DrCancer score among the best predictors for this challenge.

3.3 | NAGLU challenge

For the NAGLU challenge, participants were asked to predict the value of the RelAct of the mutated NAGLU with respect to the wild type. In this experiment, predictions close to one correspond to SAV with similar enzymatic activity with respect to the wild type. RelAct equal to zero is associated to the variants with no enzymatic activity. We used

SNPs&GO by setting the RelAct equal to 1 minus the probability for the variant to be related to disease.

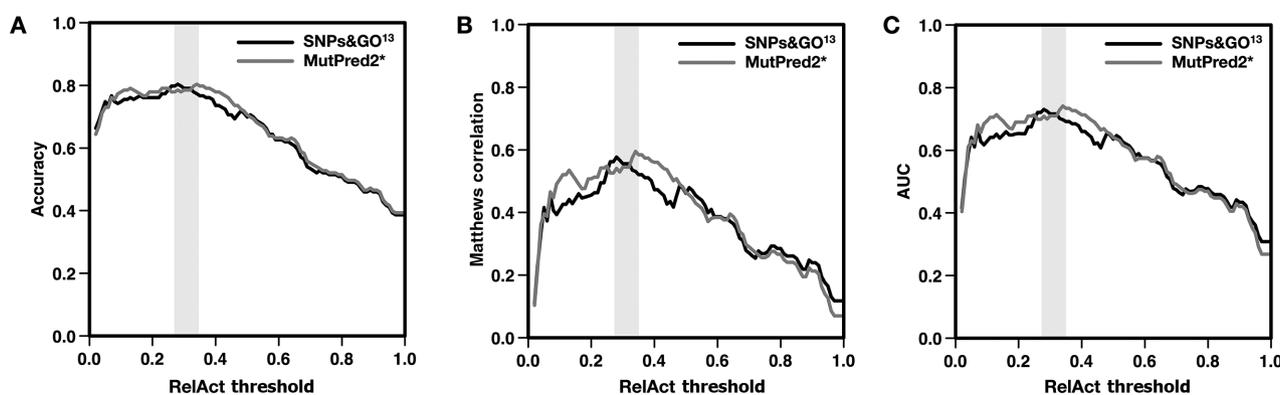
In our analysis, we compared the performance of two versions of MutPred2 with the two versions of SNPs&GO (SNPs&GO⁰⁹ and SNPs&GO¹³). The MutPred2 predictions were performed in default mode (MutPred2) and without gene-level homology count features (MutPred2*).

For SNPs&GO, the first set of predictions have been submitted by the Bologna Biocomputing Group using SNPs&GO⁰⁹. The second set of predictions, which were not submitted to the CAGI experiments, have been directly derived from the raw output of the last version of SNPs&GO (SNPs&GO¹³), maintained by the BioFold Unit. For the NAGLU challenge, we report the results of the regression and binary classification tests in Table 4. Our analysis shows that the accuracy of SNPs&GO¹³ is comparable to MutPred2*, which is the best method for this challenge. The average difference in the correlation coefficients between SNPs&GO¹³ and MutPred2* is ~ 0.02 . The results of the binary classification test, performed by optimizing the RelAct thresholds for all the methods and by considering the same output classification threshold equal to 0.5, confirm the similarity between the performance of SNPs&GO¹³ and MutPred2*. Indeed, SNPs&GO¹³ and MutPred2* achieve the same Q_2 and AUC (with RelAct thresholds equal to 0.28 and 0.34, respectively). In Figure 2, we show that the performance of SNPs&GO¹³ and MutPred2* in terms of Q_2 , AUC, and MC are consistently similar at different RelAct threshold.

TABLE 4 Performance of the predictors for the NAGLU challenge (CAGI 2016)

Method	Q_2	TPR	PPV	TNR	NPV	AUC	MC	RMSE	r_{Pearson}	r_{Spearman}	$r_{\text{KendallTau}}$
MutPred2*	0.80	0.88	0.80	0.70	0.81	0.85	0.60	0.30	0.60	0.61	0.43
MutPred2	0.68	0.38	0.69	0.89	0.68	0.79	0.31	0.30	0.51	0.54	0.37
SNPs&GO ¹³	0.80	0.87	0.82	0.69	0.77	0.84	0.58	0.32	0.56	0.58	0.42
SNPs&GO ⁰⁹	0.72	0.70	0.82	0.74	0.61	0.72	0.43	0.48	0.42	0.43	0.35

Notes: The overall accuracy (Q_2), true-positive/negative rates (TPR/TNR), positive/negative predicted values (PPV/NPV), area under the ROC curve (AUC), and Matthews correlation coefficient (MC) are computed by choosing the threshold maximizing their product. The best performance for MutPred2*, MutPred2, SNPs&GO¹³, and SNPs&GO⁰⁹ are obtained using an experimental relative activity (RelAct) threshold equal to 0.34, 0.55, 0.28, and 0.28, respectively. A threshold on the prediction equal to 0.5 is considered for all the methods. The positive and negative classes refer to pathogenic and benign variations, respectively. All the binary (Q_2 , TPR, TNR, PPV, NPV, AUC, and MC) and regression (RMSE, r_{Pearson} , r_{Spearman} , and $r_{\text{KendallTau}}$) evaluation measures are described in Supp. Material. The confusion matrices for calculating the performance of the binary classifiers are reported in Supp. Table S5. Bold correlation coefficients correspond to P values < 0.05.

**FIGURE 2** Comparison between the binary classification performance of SNPs&GO¹³ (black) and MutPred2* (gray) on the NAGLU dataset

4 | DISCUSSION

In this work, we analyzed the performance of SNPs&GO algorithm in predicting the impact of SAVs. From 2010, the Bologna Biocomputing Group and the BioFold Unit participated in all the editions of the CAGI experiments with two different versions of SNPs&GO, namely, SNPs&GO⁰⁹ and SNPs&GO¹³. The first version of SNPs&GO (SNPs&GO⁰⁹), used by the Bologna Biocomputing Group, resulted among the best algorithm for predicting the impact on SAVs in *CHEK2* and *RAD50* challenges. The last version of SNPs&GO (SNPs&GO¹³), maintained by BioFold unit, was successful in scoring the impact of genetic variants in the latest CAGI challenges (*p16* and *NAGLU*). In particular, the predictions submitted by the BioFold Unit were among the most accurate in the prediction of the impact of *p16INK4A* variants. In our a posteriori evaluation of nonsubmitted predictions for the *NAGLU* challenge, SNPs&GO¹³ resulted in performance similar to the best version of MutPred2 algorithm.

Our analysis shows that the automatic annotation of SAVs with our tools scores better when predicting the functional impact of the variants (*p16* and *NAGLU* challenges in Tables 3 and 4) than the frequency of disease variant carriers (f_{case}) (*CHEK2* and *RAD50* challenges in Tables 1 and 2). This observation derives from the comparison of the correlation coefficients for the *p16* and *NAGLU* challenges (in almost all the cases above 0.5) with those of the *CHEK2* and *RAD50* challenges (around 0.29).

The better performance of the last version of SNPs&GO¹³ with respect to the oldest SNPs&GO⁰⁹ is likely due to the more informative training set, in terms of the number of sequences available for alignments in the newer version of UniRef90 and variations in the training set as collected from SwissVar. In particular, for the *NAGLU* challenge, the release of SwissVar used for the training of SNPs&GO⁰⁹ contained only 25 disease-related SAVs, which is significantly lower than the 67 disease-related amino acid variants present in the more recent version of SwissVar used for training SNPs&GO¹³.

In general, it is difficult to evaluate the gain in the performance associated with the improvement of the GO annotations. Nevertheless, comparing SNPs&GO with AlignGVGD and SIFT in the *CHEK2* and *RAD50* challenges, we learnt that the functional contribution to the predictions is particularly helpful when evolutionary information is not discriminative enough.

Finally, we would like to point out that the improvement in the performance obtained by SNPs&GO⁰⁹ in the *RAD50* challenge on the subset of variants falling in specific protein domains (Table 2) supports the notion that evolution information is important for the quality of the prediction. Indeed, conserved regions, such as protein domains, result in more informative sequence alignments.

In the case of multiple SAVs in the same position, evolutionary information may not be sufficient for discrimination, and other features (such as physicochemical characteristics, steric hindrance, solvent accessibility, specific position in the protein structure) may be relevant

for discriminating disease related to neutral variations. SNPs&GO is based on sequence and function.

5 | CONCLUSIONS

The analysis of the results of four CAGI challenges (*CHEK2*, *RAD50*, and *p16*, *NAGLU*) shows that SNPs&GO was consistently among the best algorithms for predicting the effect of the SAVs. Although the prediction of the real value of the functional impact is still a difficult task, SNPs&GO has shown a good level of generalization reaching good performance as a binary classifier when the predictions are directly generated from the raw output without any gene/problem-specific customization.

ACKNOWLEDGMENTS

E.C. acknowledge the computational resources accessed at the GeneFix - Genome Informatics Service at the University of Alabama at Birmingham, AL and the Institute for Mathematical Modeling of Biological Systems at the University of Düsseldorf (Germany). The CAGI experiment coordination is supported by NIH U41 HG007346 and the CAGI conference by NIH R13 HG006650.

DISCLOSURE STATEMENT

The authors declare no conflict of interest.

REFERENCES

- Altschul, S. F., Madden, T. L., Schaffer, A. A., Zhang, J., Zhang, Z., Miller, W., ... Lipman, D. J. (1997). Gapped BLAST and PSI-BLAST: A new generation of protein database search programs. *Nucleic Acids Research*, *25*(17), 3389–3402.
- Bromberg, Y., Capriotti, E., & Carter, H. (2016). VarI-SIG 2015: Methods for personalized medicine—The role of variant interpretation in research and diagnostics. *BMC Genomics*, *17*(Suppl 2), 425.
- Brownstein, C. A., Beggs, A. H., Homer, N., Merriman, B., Yu, T. W., Flannery, K. C., ... Margulies, D. M. (2014). An international effort towards developing standards for best practices in analysis, interpretation and reporting of clinical genome sequencing results in the CLARITY Challenge. *Genome Biology*, *15*(3), R53.
- Calabrese, R., Capriotti, E., Fariselli, P., Martelli, P. L., & Casadio, R. (2009). Functional annotations improve the predictive score of human disease-related mutations in proteins. *Human Mutation*, *30*(8), 1237–1244.
- Cancer Genome Atlas Research Network, Weinstein, J. N., Collisson, E. A., Mills, G. B., Shaw, K. R., Ozenberger, B. A., ... Stuart, J. M. (2013). The Cancer Genome Atlas Pan-Cancer analysis project. *Nature Genetics*, *45*(10), 1113–1120.
- Capriotti, E., & Altman, R. B. (2011). A new disease-specific machine learning approach for the prediction of cancer-causing missense variants. *Genomics*, *98*(4), 310–317.
- Capriotti, E., Calabrese, R., Fariselli, P., Martelli, P. L., Altman, R. B., & Casadio, R. (2013). WS-SNPs&GO: A web server for predicting the deleterious effect of human protein variants using functional annotation. *BMC Genomics*, *14*(Suppl 3), S6.
- Capriotti, E., Nehrt, N. L., Kann, M. G., & Bromberg, Y. (2012). Bioinformatics for personal genome interpretation. *Briefings in Bioinformatics*, *13*(4), 495–512.
- Carraro, M., Minervini, G., Giollo, M., Bromberg, Y., Capriotti, E., Casadio, R., ... Tosatto, S. (2017). Performance of in silico tools for the evaluation of p16INK4a (*CDKN2A*) variants in CAGI. *Human Mutation*. (Minor Revisions).
- Damiola, F., Pertesi, M., Oliver, J., Le Calvez-Kelm, F., Voegelé, C., Young, E. L., ... Tavtigian S. V. (2014). Rare key functional domain missense substitutions in *MRE11A*, *RAD50*, and *NBN* contribute to breast cancer susceptibility: Results from a Breast Cancer Family Registry case-control mutation-screening study. *Breast Cancer Research: BCR*, *16*(3), R58.
- Durbin, R. M., Abecasis, G. R., Altshuler, D. L., Auton, A., Brooks, L. D., Gibbs, R. A., ... McVean, G. A. (2010). A map of human genome variation from population-scale sequencing. *Nature*, *467*(7319), 1061–1073.
- ENCODE Project Consortium. (2012). An integrated encyclopedia of DNA elements in the human genome. *Nature*, *489*(7414), 57–74.
- Fernald, G. H., Capriotti, E., Daneshjou, R., Karczewski, K. J., & Altman, R. B. (2011). Bioinformatics challenges for personalized medicine. *Bioinformatics*, *27*(13), 1741–1748.
- Grantham, R. (1974). Amino acid difference formula to help explain protein evolution. *Science*, *185*(4154), 862–864.
- International Cancer Genome Consortium, Hudson, T. J., Anderson, W., Artez, A., Barker, A. D., Bell, C., ... Yang, H. (2010). International network of cancer genome projects. *Nature*, *464*(7291), 993–998.
- Leaver-Fay, A., Tyka, M., Lewis, S. M., Lange, O. F., Thompson, J., Jacak, R., ... Bradley, P. (2011). ROSETTA3: An object-oriented software suite for the simulation and design of macromolecules. *Methods in Enzymology*, *487*, 545–574.
- Le Calvez-Kelm, F., Lesueur, F., Damiola, F., Vallee, M., Voegelé, C., Babikyan, D., ... Tavtigian, S. V. (2011). Rare, evolutionarily unlikely missense substitutions in *CHEK2* contribute to breast cancer susceptibility: Results from a breast cancer family registry case-control mutation-screening study. *Breast Cancer Research: BCR*, *13*(1), R6.
- Li, B., Krishnan, V. G., Mort, M. E., Xin, F., Kamati, K. K., Cooper, D. N., ... Radivojac, P. (2009). Automated inference of molecular mechanisms of disease from amino acid substitutions. *Bioinformatics*, *25*(21), 2744–2750.
- Magrane, M., & UniProt Consortium. (2011). UniProt Knowledgebase: A hub of integrated protein data. *Database (Oxford)*, *2011*, bar009.
- Mathe, E., Olivier, M., Kato, S., Ishioka, C., Hainaut, P., & Tavtigian, S. V. (2006). Computational approaches for predicting the biological effect of p53 missense mutations: A comparison of three sequence analysis based methods. *Nucleic Acids Research*, *34*(5), 1317–1325.
- Mottaz, A., David, F. P., Veuthey, A. L., & Yip, Y. L. (2010). Easy retrieval of single amino-acid polymorphisms and phenotype information using Swiss-Var. *Bioinformatics*, *26*(6), 851–852.
- Ng, P. C., & Henikoff, S. (2003). SIFT: Predicting amino acid changes that affect protein function. *Nucleic Acids Research*, *31*(13), 3812–3814.
- Niroula, A., & Vihinen, M. (2016). Variation Interpretation predictors: Principles, types, performance, and choice. *Human Mutation*, *37*(6), 579–597.
- Oetting, W. S. (2011). Exploring the functional consequences of genomic variation: The 2010 Human Genome Variation Society Scientific Meeting. *Human Mutation*, *32*(4), 486–490.
- Saez-Rodriguez, J., Costello, J. C., Friend, S. H., Kellen, M. R., Mangravite, L., Meyer, P., ... Stolovitzky, G. (2016). Crowdsourcing biomedical research: Leveraging communities as innovation engines. *Nature Reviews Genetics*, *17*(8), 470–486.
- Suzek, B. E., Huang, H., McGarvey, P., Mazumder, R., & Wu, C. H. (2007). UniRef: Comprehensive and non-redundant UniProt reference clusters. *Bioinformatics*, *23*(10), 1282–1288.

- Thusberg, J., Olatubosun, A., & Vihinen, M. (2011). Performance of mutation pathogenicity prediction methods on missense variants. *Human Mutation*, 32(4), 358–368.
- Vihinen, M. (2012). How to evaluate performance of prediction methods? Measures and their interpretation in variation effect analysis. *BMC Genomics*, 13Suppl 4, S2.
- Vihinen, M. (2013). Guidelines for reporting and using prediction tools for genetic variation analysis. *Human Mutation*, 34(2), 275–282.
- Zhou, H., & Zhou, Y. (2002). Distance-scaled, finite ideal-gas reference state improves structure-derived potentials of mean force for structure selection and stability prediction. *Protein Science*, 11(11), 2714–2726.

SUPPORTING INFORMATION

Additional Supporting Information may be found online in the supporting information tab for this article.

How to cite this article: Capriotti E, Martelli PL, Fariselli P, Casadio R. Blind prediction of deleterious amino acid variations with SNPs&GO. *Human Mutation*. 2017;00:1–8. doi: <https://doi.org/10.1002/humu.23179>