

1 **Title:**

2 **aRgus: multilevel visualization of non-synonymous single nucleotide variants &**
3 **advanced pathogenicity score modeling for genetic vulnerability assessment**

4 Authors:

5 Julian Schröter^{a,*}, Tal Dattner^{b,*}, Jennifer Hüllein^{c,*}, Alejandra Jayme^d, Vincent Heuveline^d,
6 Georg F. Hoffmann^b, Stefan Kölker^b, Dominic Lenz^b, Thomas Opladen^b, Bernt Popp^e,
7 Christian P. Schaaf^f, Christian Staufner^b, Steffen Syrbe^a, Sebastian Uhrig^c, Daniel
8 Hübschmann^{c,g,h,†}, Heiko Brennenstuhl^{b,f,†,§}

9 ^a *Division of Pediatric Epileptology, Center for Pediatrics and Adolescent Medicine,*
10 *University Hospital Heidelberg, Im Neuenheimer Feld 430, D-69120 Heidelberg, Germany.*

11 ^b *Division of Neuropediatrics and Metabolic Medicine, Center for Pediatrics and Adolescent*
12 *Medicine, University Hospital Heidelberg, Im Neuenheimer Feld 430, D-69120 Heidelberg,*
13 *Germany.*

14 ^c *Computational Oncology, Molecular Precision Oncology Program, National Center for*
15 *Tumor Diseases (NCT), German Cancer Research Center (DKFZ), Im Neuenheimer Feld*
16 *460, D-69120 Heidelberg, Germany.*

17 ^d *Engineering Mathematics and Computing Lab (EMCL), Interdisciplinary Center for*
18 *Scientific Computing (IWR), University of Heidelberg, Im Neuenheimer Feld 205, D-69120*
19 *Heidelberg, Germany.*

20 ^e *Institute of Human Genetics, University Medical Center Leipzig, Philipp-Rosenthal-Str. 55*
21 *(Haus W), D-04103 Leipzig, Germany.*

22 ^f *Institute of Human Genetics, University Hospital Heidelberg, Im Neuenheimer Feld 440, D-*
23 *69120 Heidelberg, Germany.*

24 ^g *German Cancer Consortium (DKTK), Im Neuenheimer Feld 280, D-69120 Heidelberg,*
25 *Germany.*

26 ^h *Heidelberg Institute for Stem Cell Technology and Experimental Medicine (HI-STEM), Im*
27 *Neuenheimer Feld 280, D-69120 Heidelberg, Germany.*

28

29 § To whom correspondence should be addressed.

30 * Equal contributors (first authors).

31 † Equal contributors (senior authors).

32

33 **Corresponding author:**

34 Heiko Brennenstuhl, MD, MBA
35 Institute of Human Genetics
36 University Hospital Heidelberg
37 Im Neuenheimer Feld 440
38 D-69120 Heidelberg
39 Germany

40

41 Email: heiko.brennenstuhl@med.uni-heidelberg.de

42 Phone: 06221 56-5081

43

1 **Abstract**

2 The widespread use of high-throughput sequencing techniques is leading to a rapidly
3 increasing number of disease-associated variants of unknown significance and candidate
4 genes. Integration of knowledge concerning their genetic, protein as well as functional and
5 conservational aspects is necessary for an exhaustive assessment of their relevance and for
6 prioritization of further clinical and functional studies investigating their role in human
7 disease. In order to collect the necessary information, a multitude of different databases has to
8 be accessed and data extraction from the original sources commonly is not user-friendly and
9 requires advanced bioinformatics skills. This leads to a decreased data accessibility for a
10 relevant number of potential users such as clinicians, geneticist, and clinical researchers. Here,
11 we present aRgus (<https://argus.urz.uni-heidelberg.de/>), a standalone webtool for simple
12 extraction and intuitive visualization of multi-layered gene, protein, variant, and variant effect
13 prediction data. aRgus provides interactive exploitation of these data within seconds for any
14 known gene of the human genome. In contrast to existing online platforms for compilation of
15 variant data, aRgus complements visualization of chromosomal exon-intron structure and
16 protein domain annotation with ClinVar and gnomAD variant distributions as well as
17 position-specific variant effect prediction score modeling. aRgus thereby enables timely
18 assessment of protein regions vulnerable to variation with single amino acid resolution and
19 provides numerous applications in variant and protein domain interpretation as well as in the
20 design of *in vitro* experiments.

21 **Keywords**

22 Pathogenicity scores; variant effect prediction; variant assessment; computational genetics

1 1. Introduction

2 In recent years, high-throughput sequencing methods have led to a tremendous increase in the
3 extent of genetic and variant data related to human disease (1, 2). Upon identification of
4 disease-associated genetic variants of unknown significance or in novel candidate genes, an
5 investigator may need to integrate of multi-layered information concerning exon-intron
6 structure, protein domain annotation, mutational constraint, as well as known variants present
7 in patients and healthy individuals including their allele frequency. Additionally, the potential
8 biological impact of variants on protein structure and function can be predicted using *in silico*
9 pathogenicity scores that assign a numerical value to each amino acid substitution. This is
10 particularly helpful for estimation of damaging variant effects when no functional *in vitro*
11 data is available. This information has to be taken into consideration for variant interpretation
12 according to the ACMG guidelines (3). Although the majority of the above-mentioned data
13 are publicly accessible, they are only available in abstract, tabular form, stored in a multitude
14 of different databases that have to be accessed individually, and their extraction, formatting,
15 and analysis often require extensive bioinformatic capabilities. User-friendly platforms have
16 previously been developed in order to facilitate access to genetic data from several resources
17 but lack detailed integration and visualization of different pathogenicity scoring models (4-7).

18 Therefore, we developed aRgus (<https://argus.urz.uni-heidelberg.de/>) as a standalone
19 webtool for user-friendly and intuitive compilation and visualization of complex data on
20 genetic variants and *in silico* pathogenicity scores from the extensive databases Ensembl,
21 Simple ClinVar, the Universal Protein Resource (UniProt), the Genome Aggregation
22 Database (gnomAD), and dbNSFP (4, 5, 7-9). The Ensembl database contains comprehensive
23 genomic information including chromosomal gene and transcript localization (4). Simple
24 ClinVar is an interactive webtool using a custom algorithm to retrieve simplified summary
25 statistics on variant and phenotype information from ClinVar, the largest archive of genetic
26 variants associated with human disease (5, 10). UniProt represents the largest database for
27 protein sequence and domain annotation data (7). The gnomAD database contains variant
28 data from nearly 150,000 healthy individuals identified in exome and genome sequencing
29 studies (8). The dbNSFP database represents a rich resource containing values of numerous *in*
30 *silico* pathogenicity scores precalculated for all biologically possible non-synonymous single-
31 nucleotide variants (nsSNVs) and related information, such as their gnomAD allele
32 frequencies, that can be used for variant annotation (9). dbNSFP is implemented in several
33 annotation tools such as ANNOVAR, VarSome, the UCSC Genome Browser, and the
34 Ensembl Variant Effect Predictor and also offers an own application but can only be used for
35 single queries or short lists of SNVs (6, 11-13).

36 In contrast, aRgus provides the synopsis of both variant and pathogenicity score data using an
37 intuitive graphical user interface. aRgus allows display of exon-intron structure and protein
38 domain annotation together with ClinVar and gnomAD variant distributions, a vivid
39 visualization of pathogenicity score values and their statistical comparison in different variant
40 groups, as well as an interactive table comprising ClinVar- and dbNSFP-derived variants.
41 The use of aRgus enables identification of protein regions susceptible to missense variation
42 up to single amino acid (AA) resolution and represents a powerful tool for enhanced
43 inference-based variant interpretation.

1 2. Methods

2 2.1. Implementation

3 aRgus is implemented as a standalone application using the RStudio shiny framework
4 (<https://shiny.rstudio.com/>) that allows translation of remote user operations into HTML code.
5 Chromosomal coordinates and the UniProt ID of the transcript are retrieved through Ensembl
6 (14) using the R package *AnnotationHub*. In order to achieve user-friendliness and to
7 maximize the quality of data retrieval, the canonical transcript is automatically determined
8 via query of the MANE transcript (15) or the highest quality APPRIS isoform (16). ClinVar
9 variant and phenotype annotation is retrieved using a monthly updated dataset generated via
10 the Simple ClinVar filter (5). Domain and region annotations of the corresponding protein are
11 directly retrieved from UniProt using the R package *drawProteins* (7, 17). We use a tabix-
12 indexed dbNSFP (v.4.3a) file to access up to 43 *in silico* pathogenicity scores for all possible
13 nsSNVs and their gnomAD (exomes v.2.1, genomes v.3.0) allele counts (9). All databases are
14 updated in regular intervals according to their respective release cycle. All visualizations are
15 realized using the R library *ggplot2* v3.2.1 (18). Each plot (.svg/.png) and table (.csv/.xlsx)
16 can be exported separately for offline data processing. The aRgus web server is compatible
17 with all common web browser applications including versions for mobile devices. The source
18 code is available at <https://github.com/huellejn/argus>. The application can be deployed locally
19 using a Docker image.

20 2.2. Visualization of tabular pathogenicity score data

21 Theoretically, a gene transcript can mutate at any base position into three alternate bases
22 leading to nsSNVs on the gene level as well as amino acid substitutions or truncations on the
23 protein level, depending on the position within the base triplet. The damaging effect on
24 protein function can be predicted *in silico* by an individual value of different pathogenicity
25 scores assigned to each amino acid substitution (Fig. S1). Thus, all biologically possible
26 nsSNVs can be simulated and result in several datapoints per amino acid position. In order to
27 visualize these data intuitively and vividly, a dual approach was conducted: First, the
28 *geom_smooth()* function of the R package *ggplot2* was used to generate a polynomial
29 regression of smoothed conditional means displayed by an approximation curve with 95%
30 confidence interval. Local Polynomial Regression Fitting (*loess*, formula = $y \sim x$) and a
31 generalized additive model (*GAM*, formula = $y \sim s(x, bs = "cs")$) are used for $<$ and \geq 1,000
32 datapoints, respectively. Second, the arithmetic means of multiple pathogenicity score values
33 at one amino acid position were calculated and visualized as a heat-strip color-coded by the
34 predicted degree of the damaging effect on protein function (Fig. S1).

35 2.3. Statistics

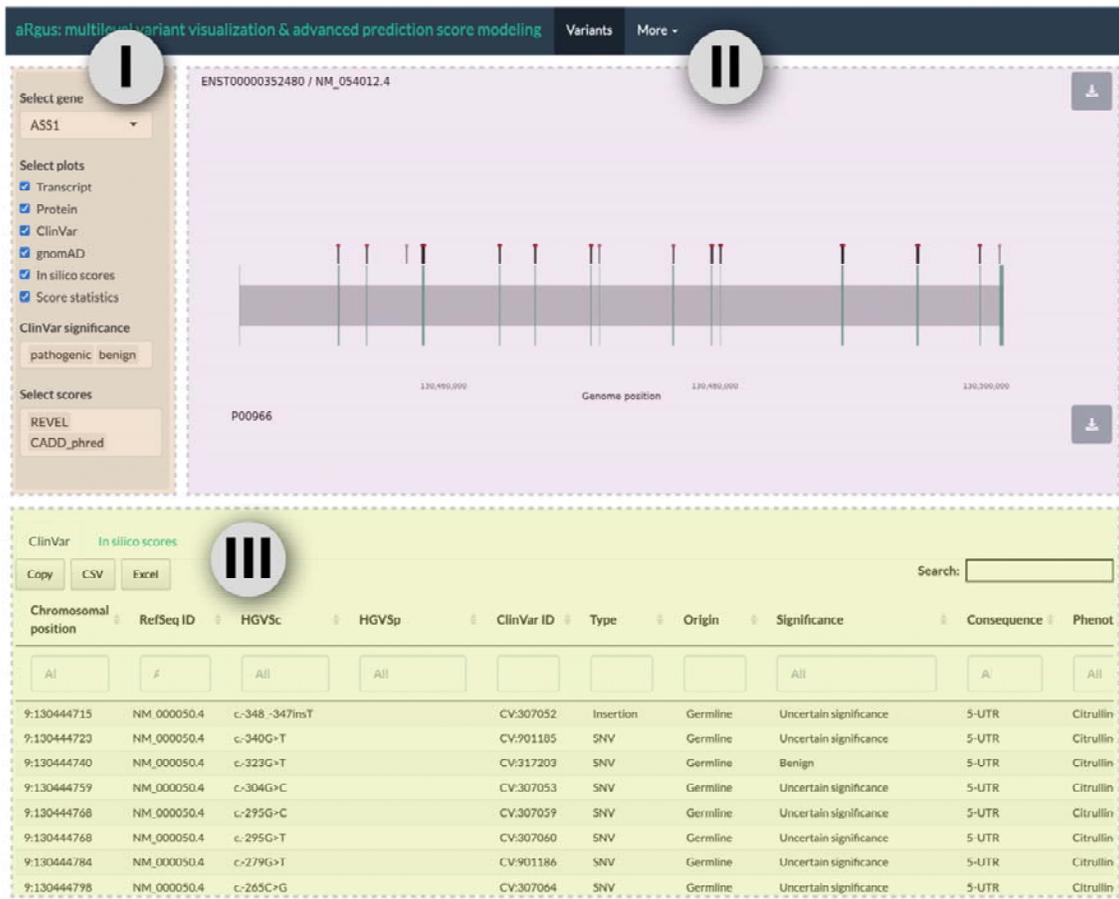
36 All pathogenicity scores can be subjected to t-test comparisons between four pre-defined
37 groups: 1.) variants stored in ClinVar and classified as pathogenic/likely pathogenic
38 (*ClinVar_pathogenic*), 2.) variants stored in ClinVar and classified as benign/likely benign
39 (*ClinVar_benign*), 3.) variants stored in gnomAD (*gnomAD*), and 4.) all biologically possible
40 variants stored in dbNSFP (*InSilico*). Score value distributions within these groups are
41 displayed as violin plots with integrated quartiles. The level of significance is shown as
42 asterisks as follows: * = $p < 0.05$, ** = $p < 0.01$, *** = $p < 0.001$.

43

1 3. Results

2 3.1. Main user interface

3 aRgus provides intuitive use and accessibility. It can be accessed via all common browsers
 4 and operating systems including mobile devices. On the aRgus main page, the user can enter
 5 a gene of interest via its HGNC symbol (Fig. 1) The tool subsequently provides the MANE-
 6 and APPRIS-curated canonical transcripts. The user can then choose from a panel of six plots
 7 that can be displayed in a modular way in order to allow an individual compilation: 1.)
 8 Unspliced transcript plot; 2.) protein plot; 3.) mutational constraint plots of disease-associated
 9 and putatively benign ClinVar as well as 4.) tolerated gnomAD variants; 5.) a combined
 10 pathogenicity score model including a polynomial fit and heat-strip with position-coded
 11 annotation of score mean values; and 6.) a statistical comparison of score values of different
 12 variant groups. Additionally, two interactive tables are available including a tab for all
 13 ClinVar variants (*ClinVar*) and all biologically possible nsSNVs together with corresponding
 14 score values derived from the dbNSFP database (*In silico scores*), respectively. Plots can be
 15 exported in two file formats (.png/.svg) with user-specified aspect ratios. Tables can be
 16 exported as .csv or .xlsx files for individual data storage and further offline data manipulation.



17 **Fig. 1: aRgus user interface.** I) Interactive input mask with control elements, II) dynamic
 18 results area, III) tables from which variants can be selected for display with label.

19

1 **3.2. Applications**

2 **3.2.1. Unspliced transcript plot**

3 The unspliced transcript plot (UTP) displays the gene's scaled exon-intron structure from left
4 to right starting with the first exon for improved readability regardless of the genomic
5 localization on the forward or reverse strand. By default, pathogenic and likely pathogenic
6 (P/LP) ClinVar variants are shown as lollipops which allows convenient visualization of
7 intronic variants. In order to display the variant description, ClinVar and simulated dbNSFP
8 variants can be manually selected in the respective tables. Figure 2A shows the UTP for the
9 gene *ASS1*, encoding the enzyme argininosuccinate synthase (ASS), with selected P/LP
10 ClinVar variants (red) and variants from the *In silico scores* table (gray), containing the
11 dbNSFP-derived variants.

12 **3.2.2. Protein plot**

13 The primary structure of the resulting protein is visualized by the protein plot showing a
14 linearized representation together with annotated domains retrieved from UniProt. As in the
15 UTP, variants can be manually selected from the provided tables. Thereby, distribution of
16 known and novel variants and their relation to protein domains/regions can easily be assessed.
17 This versatile visualization provides useful insights for assessment of the pathophysiological
18 relevance of potentially functionally relevant domains, given a gene scarcely associated with
19 pathogenic variants. Figure 2B shows respective amino acid changes and protein domains of
20 ASS.

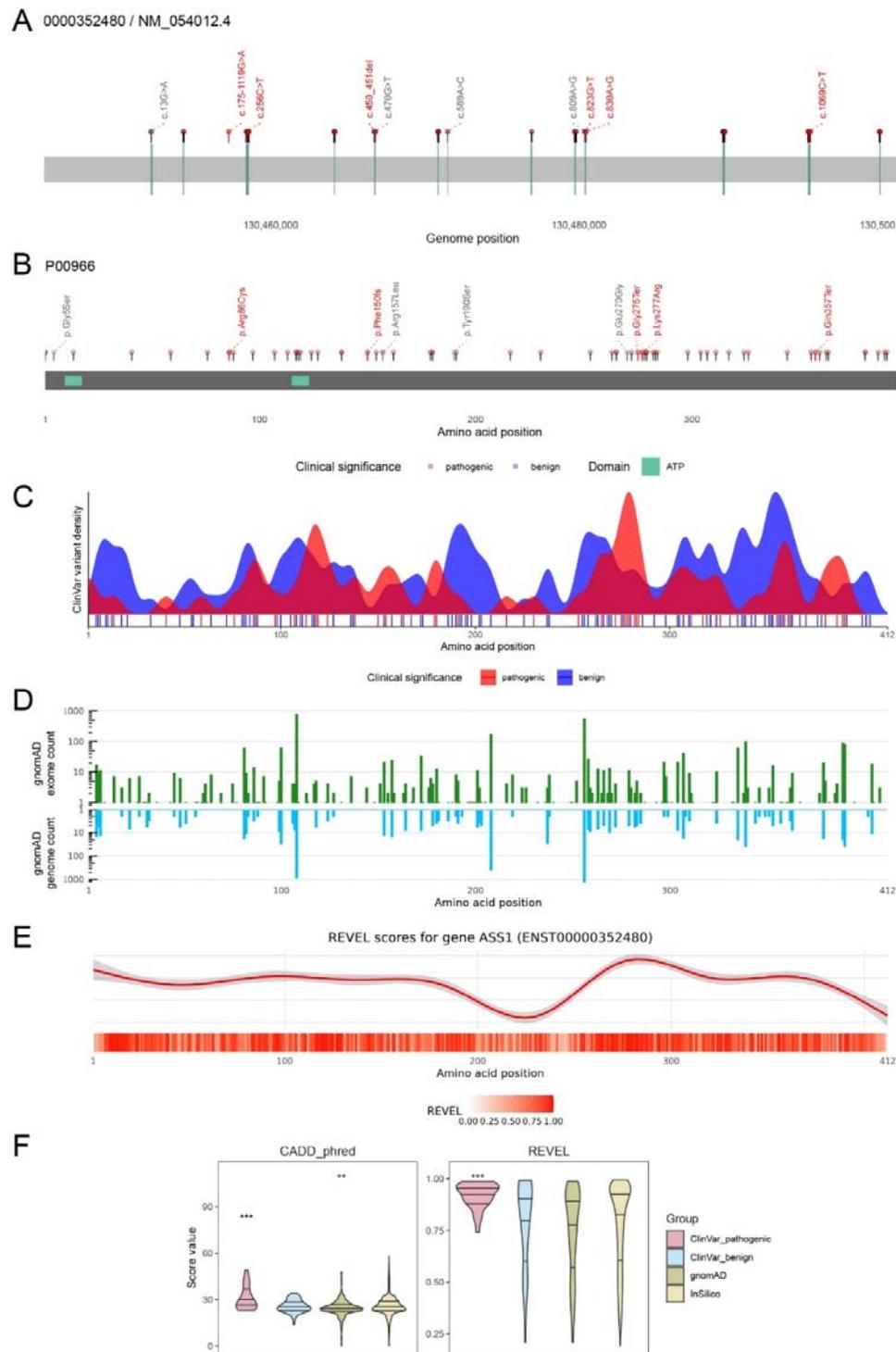
21 **3.2.3. ClinVar and gnomAD mutational constraint plots**

22 Distributions of ClinVar and gnomAD variants with respect to their protein position and
23 allele frequency are visualized by density and bar plots, respectively, facilitating assessment
24 of a protein's mutational constraint. This includes sections of mutational hotspots, recurrent
25 pathogenic and benign variants as well as the position-specific degree of tolerance towards
26 missense variation. For more precise localization, ClinVar variants are additionally shown as
27 vertical lines underneath the density curves (Fig. 2C). gnomAD variants are displayed in two
28 separate logarithmic bar plots depending on their origin from the exomes (green) or genomes
29 (blue) dataset (Fig. 2D). For ASS, ClinVar density curves reveal an accumulation of
30 pathogenic variants in the region of AA 260-280 whereas gnomAD variants from both
31 exomes and genomes show low population allele frequencies or are completely absent from
32 the dataset (Figure 2E).

33 **3.2.4. In silico pathogenicity score model**

34 Pre-calculated pathogenicity score values of all biologically possible nsSNVs are retrieved
35 from the dbNSFP database. In order to improve data accessibility, the resulting multiple data
36 points per protein position are simplified and visualized using a polynomial regression model
37 combined with a heat-strip scaled to the linear protein representation. Depending on the
38 user's research question, the desired pathogenicity scoring model can immediately be
39 selected from a list of up to 43 different scores. Plots for three different scores can be
40 displayed simultaneously. This enables assessment of the predicted, position-specific impact
41 of amino acid substitutions within the context of known protein domains and facilitates
42 detection of regions of increased or decreased susceptibility to missense variation. Thereby,
43 the functional impact of novel variants can be estimated and investigation of unknown

- 1 sections of predicted damaging variant effects can be addressed in order to formulate future
- 2 research hypotheses.



1 **Fig. 2: aRgus plots.** A) UTP of the gene *ASS1*. Labels show P/LP variants (red) and selected
2 variants from the *in silico* tab (gray). B) Protein plot with AA exchanges corresponding to
3 variants shown in A). C) Density plot of P/LP (red) and benign/likely benign (blue) Simple
4 ClinVar variants. D) Logarithmic histogram of gnomAD exomes (green) and genomes (blue)
5 variant allele frequencies. E) Polynomial regression of REVEL score (top) and heat-strip of
6 mean score values (bottom). F) t-test group comparisons shown as violin plots with quartiles,
7 * (p -value < 0.05), ** (p -value < 0.01), and *** (p -value < 0.001).

1 In our practical example, regions with low (AA 200-250) and high (AA 270–300) values of
2 the pathogenicity score *REVEL* correspond to local minima and maxima of the curve. The
3 heat-strip representation displays mean score values allowing a more fine-granular resolution
4 (Fig. 2E).

5 **3.2.5. Statistical comparisons**

6 Pathogenicity score values within the four variant groups *ClinVar_pathogenic*,
7 *ClinVar_benign*, *gnomAD*, and *InSilico* are shown as violin plots with integrated quartiles
8 (for definitions see Methods section 2.3). Additionally, score value distributions are
9 statistically compared in order to assess the capability of the specific score to discriminate
10 between variants of the different categories and hence its possible suitability for variant
11 classification. For example, *ASSI* variants, that were annotated as P/LP, yield significantly
12 higher *CADD* and *REVEL* score values than variants in the other three groups (Fig. 2F).

13 **3.2.6. Interactive table**

14 On the bottom side of the user interface, an interactive table, that remains sticky during
15 scrolling, is available (Fig. 1). It comprises two tabs with all ClinVar variants as well as all
16 simulated nsSNVs and corresponding pathogenicity score values. In order to provide
17 interactivity to the user, selected variants are displayed in the UTP and protein plot. Both
18 tables can be filtered, e.g., by variant type. Individual cells with score values in the *in silico*
19 table are color-coded according to the predicted variant effect using score-specific cut-offs.

20

1 **4. Discussion**

2 The availability of databases with clinical and genetic information has never been greater
3 than it is today. Scientific and medical advances, particularly in terms of sequencing and
4 storage capabilities, will lead to an exponential growth of information in the coming decades.
5 However, database queries often require bioinformatic tools, which ultimately limit the yield
6 and usability of such. To enable clinicians, scientists, and other users without prior
7 bioinformatic knowledge to explore rich yet complex datasets, user-friendly tools with an
8 intuitive interface and the possibility to easily export data for further processing are needed.
9 Web server applications allow users to make such queries regardless of the device and
10 operating system. aRgus is therefore designed as a lightweight, multidimensional R/Shiny
11 application to enable fast database queries.

12 aRgus uses minimal user input in the form of the gene name according to HUGO Gene
13 Nomenclature Committee (HGNC) standard. aRgus can thus retrieve information of variable
14 complexity on the localization and distribution of pathogenic variants at the chromosomal
15 and protein levels, which can be used to explore biological and biochemical properties, such
16 as mutational hotspots of pathogenic and benign variance within proteins. Visual linkages of
17 pathogenic variation can be generated by annotating functionally important regions and
18 domains from the UniProt database. aRgus provides simple means of displaying complex
19 distributional information using complexity-reduced density representation that is quick and
20 easy for the human eye to comprehend. The user is offered a wide range of possibilities to
21 select relevant information to answer respective research questions.

22 By allowing simultaneous display of variants stored in gnomAD, the issue of survivorship
23 bias, as a form of selection bias, can be overcome. Survivorship bias occurs in all clinical
24 genetic databases and potentially leads to oversight of variants, that did not pass biological
25 selection, by sole assessment of pathogenic variants from clinical databases such as ClinVar.
26 This often results in misconceptions in the interpretation of mutational hotspots. The
27 gnomAD database v2.1 contains over 125,000 exomes and 15,000 genomes from different
28 populations. A comparison of benign variants derived from gnomAD and pathogenic variants
29 listed in ClinVar and other genetic databases thus enables an improved assessment of putative
30 pathogenic hotspots on the gene and protein level.

31 Beyond pure visualization of information on known pathogenic variants, a polynomial
32 regression model and heatmap visualization offer an additional way of data exploitation
33 which can be particularly advantageous for proteins that have previously been described to
34 only a limited extent. These models overcome inaccessible, tabular data on pathogenicity
35 scores and simplify the comprehensibility of visualized predicted variant effects up to single
36 amino acid resolution. By annotation of all biologically possible missense variants using 36
37 different pathogenicity scores, statements can be made about protein regions with high
38 impact of amino acid exchanges without existing *in vitro* studies. Alternatively, resulting
39 information can be used to plan functional *in vitro* studies, e.g., in order to investigate the
40 functional relevance of regions in scarcely described proteins or with only limited data on
41 pathogenic variants.

42 **4.1. Limitations**

43 Despite of its scientific value, aRgus is subject to some limitations. The quality of the
44 visualizations and analyses produced by aRgus heavily depends on the quality of data
45 available. According to our use cases, ClinVar data does not represent the entirety of all
46 previously reported pathogenic variants. This is largely due to the lack of obligation of
47 genetic laboratories to enter newly discovered disease-causing variants in centralized
48 repositories. Extensive literature reviews are therefore necessary to obtain a comprehensive

1 picture of mutational distribution. This could be significantly improved by the addition of
2 further, commercial databases such as HGMD or LOVD (19, 20). To enable users to
3 visualize variants identified through their own literature research or genetic studies, variants
4 can be selected from the dbNSFP-derived table of pathogenicity score values and are
5 automatically highlighted in all plots.

6 **4.2. Conclusion**

7 Combining accessible and interactive visualizations of genetic and variant data with
8 pathogenicity analysis in a synoptic, standalone tool, aRgus outstands existing applications
9 for genetic data exploitation regarding output versatility and flexibility (5, 21). With each
10 update of the databases connected to aRgus, the diversity and analysis capabilities of its
11 visualizations and datasets will also improve. Thus, aRgus will provide useful and previously
12 mostly inaccessible information to a broad usership with limited bioinformatics skills such as
13 practicing clinicians, basic scientists, and geneticists, and thus be helpful to answer scientific
14 questions.
15

1 **5. Acknowledgments**

2 **5.1. Funding**

3 This work was supported by the NCT Molecular Precision Oncology Program and the
4 Physician Scientist Program of the Medical Faculty of the University of Heidelberg (JS, HB).
5 JS and SS received funding by the Dietmar Hopp Stiftung (grant 1DH1813319 to SS). BP
6 was supported by the Deutsche Forschungsgemeinschaft (grant PO2366/2-1).

7 **5.2. Conflict of interest**

8 None declared.

9 **6. Author contributions**

10 JS, TD, and HB devised the project and main conceptual ideas and designed the study. JS,
11 HB, JH, AJ, SU, and DH have designed and delivered the technical realization and
12 implementation of aRgus. All authors were involved in the further development of aRgus
13 during the development period through their intellectual input and the execution of targeted
14 analyses. All authors provided critical feedback and helped shape the research, analysis, and
15 manuscript.
16

1 7. References

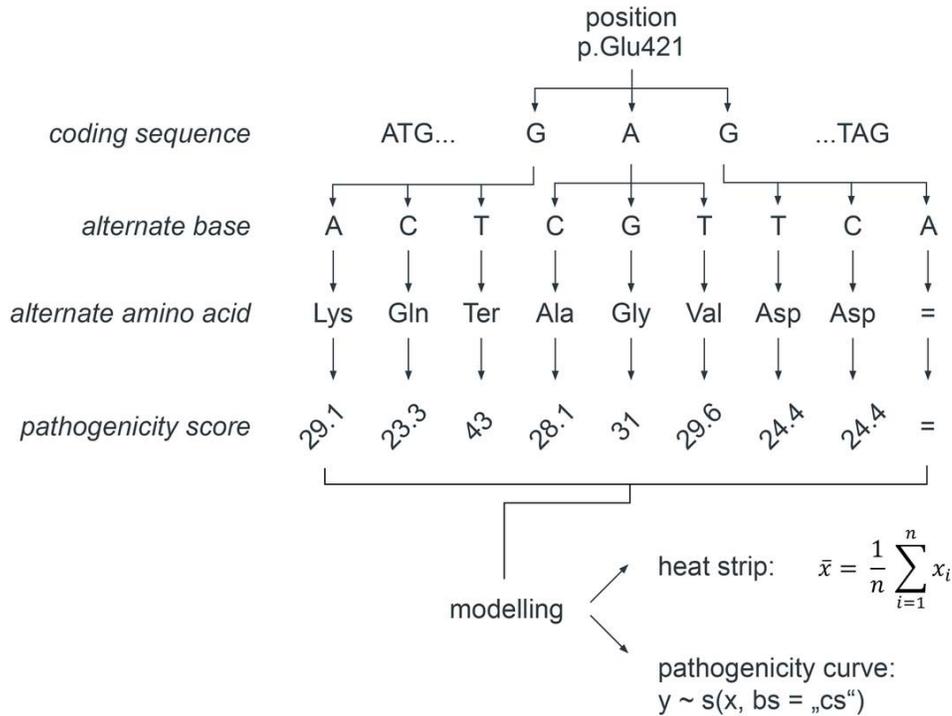
- 2 1. Goodwin S, McPherson JD, McCombie WR. Coming of age: ten years of next-
3 generation sequencing technologies. *Nat Rev Genet.* 2016;17(6):333-51.
- 4 2. Kingsmore SF. 2022: a pivotal year for diagnosis and treatment of rare genetic
5 diseases. *Cold Spring Harb Mol Case Stud.* 2022;8(2).
- 6 3. Richards S, Aziz N, Bale S, Bick D, Das S, Gastier-Foster J, et al. Standards and
7 guidelines for the interpretation of sequence variants: a joint consensus recommendation of
8 the American College of Medical Genetics and Genomics and the Association for Molecular
9 Pathology. *Genet Med.* 2015;17(5):405-24.
- 10 4. Howe KL, Achuthan P, Allen J, Allen J, Alvarez-Jarreta J, Amode MR, et al.
11 Ensembl 2021. *Nucleic Acids Res.* 2021;49(D1):D884-D91.
- 12 5. Perez-Palma E, Gramm M, Nurnberg P, May P, Lal D. Simple ClinVar: an interactive
13 web server to explore and retrieve gene and disease variants aggregated in ClinVar database.
14 *Nucleic Acids Res.* 2019;47(W1):W99-W105.
- 15 6. Kopanos C, Tsiolkas V, Kouris A, Chapple CE, Albarca Aguilera M, Meyer R, et al.
16 VarSome: the human genomic variant search engine. *Bioinformatics.* 2019;35(11):1978-80.
- 17 7. UniProt C. UniProt: the universal protein knowledgebase in 2021. *Nucleic Acids Res.*
18 2021;49(D1):D480-D9.
- 19 8. Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alfoldi J, Wang Q, et al. The
20 mutational constraint spectrum quantified from variation in 141,456 humans. *Nature.*
21 2020;581(7809):434-43.
- 22 9. Liu X, Li C, Mou C, Dong Y, Tu Y. dbNSFP v4: a comprehensive database of
23 transcript-specific functional predictions and annotations for human nonsynonymous and
24 splice-site SNVs. *Genome Med.* 2020;12(1):103.
- 25 10. Landrum MJ, Lee JM, Riley GR, Jang W, Rubinstein WS, Church DM, et al. ClinVar:
26 public archive of relationships among sequence variation and human phenotype. *Nucleic*
27 *Acids Res.* 2014;42(Database issue):D980-5.
- 28 11. Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM, et al. The
29 human genome browser at UCSC. *Genome Res.* 2002;12(6):996-1006.
- 30 12. McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GR, Thormann A, et al. The Ensembl
31 Variant Effect Predictor. *Genome Biol.* 2016;17(1):122.
- 32 13. Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants
33 from high-throughput sequencing data. *Nucleic Acids Res.* 2010;38(16):e164.
- 34 14. Morgan M, Shepherd L. AnnotationHub: Client to access AnnotationHub resources.
35 R package version 2.22.0. 2020.
- 36 15. Morales J, Pujar S, Loveland JE, Astashyn A, Bennett R, Berry A, et al. A joint NCBI
37 and EMBL-EBI transcript set for clinical genomics and research. *Nature.*
38 2022;604(7905):310-5.
- 39 16. Rodriguez JM, Rodriguez-Rivas J, Di Domenico T, Vazquez J, Valencia A, Tress ML.
40 APPRIS 2017: principal isoforms for multiple gene sets. *Nucleic Acids Res.*
41 2018;46(D1):D213-D7.
- 42 17. Brennan P. drawProteins: a Bioconductor/R package for reproducible and
43 programmatic generation of protein schematics. *F1000Res.* 2018;7:1105.
- 44 18. Wickham H. *ggplot2 : Elegant Graphics for Data Analysis.* Cham: Springer
45 International Publishing : Imprint: Springer,; 2016.
- 46 19. Fokkema IF, Taschner PE, Schaafsma GC, Celli J, Laros JF, den Dunnen JT. LOVD
47 v.2.0: the next generation in gene variant databases. *Hum Mutat.* 2011;32(5):557-63.

- 1 20. Stenson PD, Mort M, Ball EV, Chapman M, Evans K, Azevedo L, et al. The Human
2 Gene Mutation Database (HGMD((R))): optimizing its use in a clinical diagnostic or research
3 setting. *Hum Genet.* 2020;139(10):1197-207.
- 4 21. Perez-Palma E, May P, Iqbal S, Niestroj LM, Du J, Heyne HO, et al. Identification of
5 pathogenic variant enriched regions across genes and gene families. *Genome Res.*
6 2020;30(1):62-71.

7
8

1 **8. Supplementary files**

2 **8.1. Figure S1**



3

4 **Fig. S1: Schematic illustration of dbNSFP-derived variant simulation and aRgus-**
 5 **mediated visualization.** Starting from the coding sequence of a gene transcript, any base at
 6 any position is exchanged with its three non-synonymous alternate bases (top). Individual
 7 pathogenicity score values (bottom) are assigned to the corresponding amino acid
 8 substitutions (middle). In aRgus, the resulting tabular data is modelled and visualized using a
 9 dual approach with a polynomial regression curve and a heat strip.

10

1 **8.2. Table S1. Pathogenicity scores available on aRgus.**

Score	Version	Source
REVEL	Release May 3, 2021	https://sites.google.com/site/revelgenomics/
CADD_phred	v1.6	http://cadd.gs.washington.edu/
SIFT	ensembl 66	https://sift.bii.a-star.edu.sg/www/history.html
SIFT4G	v2.4	http://sift.bii.a-star.edu.sg/sift4g/public//Homo_sapiens/
Polyphen HDIV	v2.2.2	http://genetics.bwh.harvard.edu/pph2/
Polyphen HVAR	v2.2.2	http://genetics.bwh.harvard.edu/pph2/
PROVEAN	v1.1 ensembl 66	http://provean.jcvi.org/index.php
M-CAP	v1.3	http://bejerano.stanford.edu/MCAP/
VEST4	v4.0	http://karchinlab.org/apps/appVest.html
FATHMM	v2.3	http://fathmm.biocompute.org.uk
MetaSVM	n/a	doi: 10.1093/hmg/ddu733
MetaLR	n/a	doi: 10.1093/hmg/ddu733
ClinPred	n/a	https://sites.google.com/site/clinpred/home
MutationTaster	v2	http://www.mutationtaster.org/
MutationAssessor	Release 3	http://mutationassessor.org/
DANN	n/a	https://cbel.ics.uci.edu/public_data/DANN/
MutPred	v1.2	http://mutpred.mutdb.org/
MVP	v1.0	https://github.com/ShenLab/missense
MPC	Release 1	ftp://ftp.broadinstitute.org/pub/ExAC_release/release1/regional_missense_constraint/
LRT	Release 11/2009	http://www.genetics.wustl.edu/jflab/lrt_query.html
Primate AI	n/a	https://github.com/Illumina/PrimateAI
DEOGEN2	n/a	https://deogen2.mutaframe.com/
BayesDel_addAF	v1	http://fengbj-laboratory.org/BayesDel/BayesDel.html
BayesDel_noAF	v1	http://fengbj-laboratory.org/BayesDel/BayesDel.html
fathmm.MKL_coding	v2.3	http://fathmm.biocompute.org.uk/fathmmMKL.htm
fathmm.XF_coding	v2.3	http://fathmm.biocompute.org.uk/fathmm-xf/
Eigen.raw	v1.1	http://www.columbia.edu/~ii2135/eigen.html

Eigen.PC.raw	v1.1	http://www.columbia.edu/~ji2135/eigen.html
GenoCanyon	v1.0.3	http://genocanyon.med.yale.edu/index.html
integrated_fitCons	v1.01	http://compgen.bscb.cornell.edu/fitCons/
GM12878_fitCons	v1.01	http://compgen.bscb.cornell.edu/fitCons/
H1.hESC_fitCons	v1.01	http://compgen.bscb.cornell.edu/fitCons/
HUVEC_fitCons	v1.01	http://compgen.bscb.cornell.edu/fitCons/
LINSIGHT	n/a	http://compgen.cshl.edu/~yihuang/LINSIGHT/
GERP++_RS	n/a	http://mendel.stanford.edu/SidowLab/downloads/gerp/
phyloP100way_vertibrate	n/a	http://hgdownload.soe.ucsc.edu/goldenPath/hg38/phyloP100way/
phyloP30way_mammalian	n/a	http://hgdownload.soe.ucsc.edu/goldenPath/hg38/phyloP30way/
phyloP17way_primate	n/a	http://hgdownload.soe.ucsc.edu/goldenPath/hg38/phyloP17way/
phastCons100way_vertibrate	n/a	http://hgdownload.soe.ucsc.edu/goldenPath/hg38/phastCons100way/
phastCons30way_mammalian	n/a	http://hgdownload.soe.ucsc.edu/goldenPath/hg38/phastCons30way/
phastCons17way_primate	n/a	http://hgdownload.soe.ucsc.edu/goldenPath/hg38/phastCons17way/
SiPhy_29way_logOdds	n/a	https://www.broadinstitute.org/mammals-models/29-mammals-project-supplementary-info
LIST.S2_score	Release: 2019_10	https://precomputed.list-s2.msl.ubc.ca/

1 *n/a: not applicable.*