

Laurens Wiel ORCID iD: 0000-0003-3410-760X

Title:

MetaDome: Pathogenicity analysis of genetic variants through aggregation of homologous human protein domains

Authors and affiliations:

Laurens Wiel^{1,2}, Coos Baakman², Daan Gilissen^{1,3}, Joris A. Veltman^{4,5}, Gerrit Vriend² and Christian Gilissen¹

1. Department of Human Genetics, Radboud Institute for Molecular Life Sciences, Radboud University Medical Center, Nijmegen, 6525 GA, the Netherlands
2. Centre for Molecular and Biomolecular Informatics, Radboud Institute for Molecular Life Sciences, Radboud University Medical Center, Nijmegen, 6525 GA, the Netherlands
3. Bio-informatica, HAN University of Applied Sciences, Nijmegen, 6525 EN, the Netherlands
4. Department of Human Genetics, Donders Institute for Brain, Cognition and Behaviour, Radboud University Medical Center, Nijmegen, 6525 GA, the Netherlands
5. Institute of Genetic Medicine, International Centre for Life, Newcastle University, Newcastle upon Tyne, NE1 3BZ, United Kingdom

Corresponding author:

Christian Gilissen

Grant numbers:

The Netherlands Organization for Scientific Research (916-14-043 to C.G. and 918-15-667 to J.A.V.) & Radboud Institute for Molecular Life Sciences, Radboud university medical center (R0002793 to G.V.).

This article has been accepted for publication and undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the Version of Record. Please cite this article as doi: 10.1002/humu.23798.

This article is protected by copyright. All rights reserved.

Accepted Article

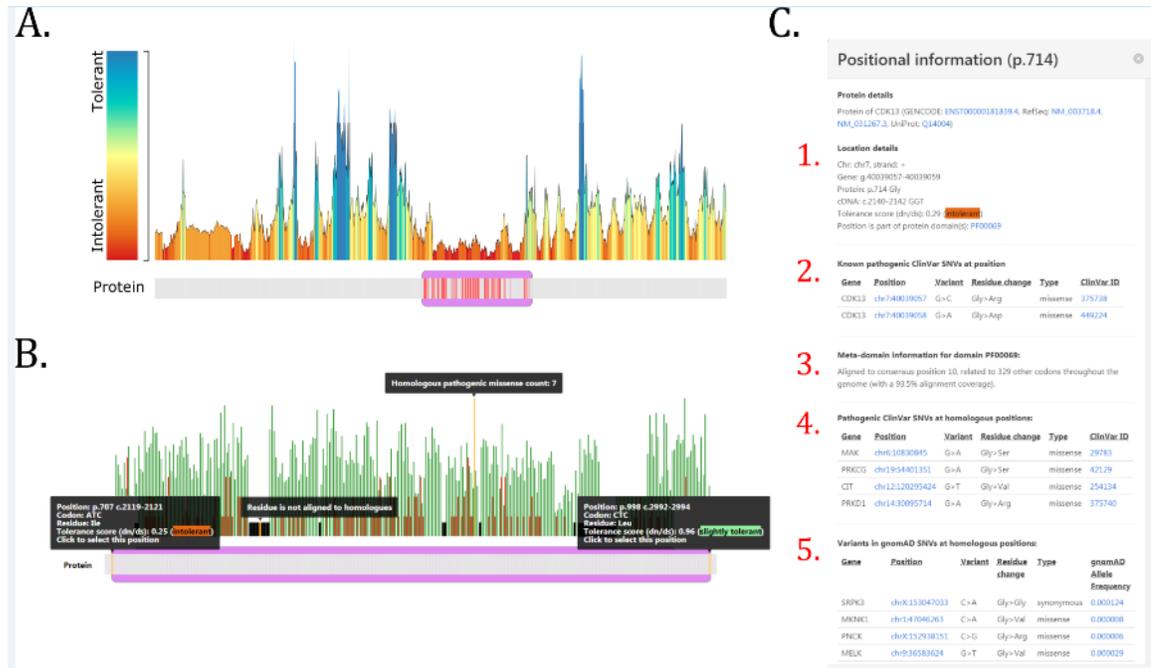
Abstract (198/200)

The growing availability of human genetic variation has given rise to novel methods of measuring genetic tolerance that better interpret variants of unknown significance. We recently developed a concept based on protein domain homology in the human genome to improve variant interpretation. For this purpose we mapped population variation from the Exome Aggregation Consortium (ExAC) and pathogenic mutations from the Human Gene Mutation Database (HGMD) onto Pfam protein domains. The aggregation of these variation data across homologous domains into meta-domains allowed us to generate amino-acid resolution of genetic intolerance profiles for human protein domains.

Here we developed MetaDome, a fast and easy-to-use web server that visualizes meta-domain information and gene-wide profiles of genetic tolerance. We updated the underlying data of MetaDome to contain information from 56,319 human transcripts, 71,419 protein domains, 12,164,292 genetic variants from gnomAD, and 34,076 pathogenic mutations from ClinVar. MetaDome allows researchers to easily investigate their variants of interest for the presence or absence of variation at corresponding positions within homologous domains. We illustrate the added value of MetaDome by an example that highlights how it may help in the interpretation of variants of unknown significance. The MetaDome web server is freely accessible at <https://stuart.radboudumc.nl/metadome>.

Graphical Abstract

We developed MetaDome, a fast and easy-to-use web server that visualizes meta-domain information and gene-wide profiles of genetic tolerance. MetaDome allows researchers to easily investigate their variants of interest for the presence or absence of variation at corresponding positions within homologous domains. The MetaDome web server is freely accessible at <https://stuart.radboudumc.nl/metadome>.



Key Words

Genetic variation; pathogenicity; web server; protein domain homology; genetic tolerance; meta-domains; gnomAD; ClinVar; Pfam

Introduction

The continuous accumulation of human genomic data has spurred the development of new methods to interpret genetic variants. There are many freely available web servers and services that facilitate the use of these data by non-bioinformaticians. For example, the ESP Exome Variant Server (Fu et al., 2012; NHLBI GO Exome Sequencing Project (ESP), 2011) and the Genome Aggregation Database (gnomAD) browser (Karczewski et al., 2017; Lek et al., 2016) help locate variants that occur frequently in the general population. These services are used for the interpretation of unknown variants based on the assumption that variants occurring frequently in the general population are unlikely to be relevant for patients with Mendelian disorders (Amr et al., 2016). There are also methods that derive information from these large human genetic databases. For example

genetic intolerance, which is commonly used to interpret variants of unknown significance by assessing whether variants stand out because they occur in regions that are genetically invariable in the general population (Ge et al., 2016; Gussow, Petrovski, Wang, Allen, & Goldstein, 2016). Examples of such methods are RVIS (Petrovski, Wang, Heinzen, Allen, & Goldstein, 2013) and subRVIS (Gussow et al., 2016). The strongest evidence for the pathogenicity of a genomic variant comes from the presence of that variant in any of the clinically relevant genetic variant databases such as the Human Gene Mutation Database (HGMD) (Stenson et al., 2017) or the public archive of clinically relevant variants (ClinVar) (Landrum et al., 2016). These databases are gradually growing in the amount of validated pathogenic information.

Another way to provide evidence for the pathogenicity of a genomic variant is to observe the effect of that variant in homologous proteins across different species. Mutations at corresponding locations in homologous proteins are found to result in similar effects on protein stability (Ashenberg, Gong, & Bloom, 2013) and can facilitate variant interpretation between disease genes and their paralogues (Lal et al., 2017). Finding homologous proteins is one the key applications of BLAST (Altschul, Gish, Miller, Myers, & Lipman, 1990). Transferring information between homologous proteins is one of the oldest concepts in bioinformatics, and can be achieved by performing a multiple sequence alignment (MSA) and locating equivalent positions between the protein sequences. We have previously used this concept and showed that it also holds for homologous Pfam protein domain relationships within the human genome. We found that ~71-72% of all disease-causing missense variants from HGMD and ClinVar occur in regions translating to a Pfam protein domain and observed that pathogenic missense variants at equivalent domain positions are often paired with the absence of population-based variation and *vice versa* (Wiel, Venselaar, Veltman, Vriend, & Gilissen, 2017). By aggregating variant information over homologous protein domains, the resolution of genetic tolerance per position is increased to the number of aligned positions. Similarly, the annotation of pathogenic variants found at equivalent domain positions also assists the interpretation of variants of unknown significance. This use of variant information from homologous protein domains was dubbed ‘meta-domains’. We realized that this type of information could be of great benefit to the genetics community and therefore developed ‘MetaDome’.

MetaDome is a freely available web server that uses our concept of meta-domains to optimally use the information from population-based and pathogenic variation datasets without the need of a bioinformatics intermediate. MetaDome is easy to use and utilizes the latest population datasets by incorporating the gnomAD and ClinVar datasets.

Results

Accessibility

The MetaDome web server is freely accessible at <https://stuart.radboudumc.nl/metadome>. MetaDome features a user-friendly web interface and features a fully interactive tour to get familiar with all parts of the analysis and visualizations.

All source code and detailed configuration instructions are available in our GitHub repository: <https://github.com/cmbi/metadome>.

The underlying database: a mapping between genes and proteins

The MetaDome web server queries genomic datasets in order to annotate positions in a protein or a protein domain. Therefore, the server needs access to genomic positional information as well as protein sequence and protein domain information. The database maps GENCODE gene translations to entries in the UniProtKB/Swiss-Prot databank in a per-position manner and corresponding protein domains or genomic variation. With respect to our criteria to map gene translations to proteins (**Methods; creating the mapping database**), 42,116 of the 56,319 full-length protein-coding GENCODE Basic transcripts for 19,728 human genes are linked to 33,492 of the 42,130 Swiss-Prot human canonical or isoform sequences. Of the total 591,556 canonical and isoform sequences present in Swiss-Prot, 42,130 result from the Human species. The resulting mappings contain 32,595,355 unique genomic positions that are linked to 19,226,961 residues in Swiss-Prot protein sequences.

71,419 Pfam domains are linked to 30,406 of the Swiss-Prot sequences in our database. Of these Pfam domain instances, 5,948 are from a unique Pfam domain family and 3,334 of these families have two or more homologues and are therefore suitable for meta-domain construction. Thus, by incorporating every protein-coding transcript, instead of only the longest ones, we increase the previously 2,750 (Wiel et al., 2017) meta-domains to 3,334. These meta-domains, on average, consist of 16 human protein domain homologues with a protein sequence length of 158 residues. **Table 1** summarizes the counting statistics for sequences, domains, etc.

How to use the MetaDome web server

At the welcome page users are offered the option to start an interactive tour or start with the analysis. The navigation bar at the top is available throughout all web pages in MetaDome and allow for further navigation to the 'About', 'Method', 'Contact' page (**Supp. Figure S1**). The user can fill in a gene symbol in the 'gene of interest' field and is aided by an auto-completion to help you find your gene of interest more easily (**Supp. Figure S2**). Clicking the 'Get transcripts' fills all GENCODE transcripts for that gene in

the dropdown box. Only the transcripts that are mapped to a Swiss-Prot protein can be used in the analysis, the others are displayed in grey (**Supp. Figure S3**).

Clicking the ‘Start Analysis’ button starts an extensive query to the back-end of the web server for the selected transcript. Firstly, all the mappings are retrieved for the transcript of interest. Secondly, the entire transcript is annotated with ClinVar and gnomAD single nucleotide variants (SNVs) and Pfam domains. Thirdly, if there are any Pfam domains suitable for meta-domain relations then all mappings for those regions are gathered and annotated with ClinVar and gnomAD variation (**methods; Composing a meta-domain**).

The web-page provided to the user as a result of the ‘Analyse Protein’ can best be explained using an example. Therefore, we have generated this result for gene *CDK13* for transcript ‘ENST00000181839.4’ (**Figure 1**). The result page features four main components that we will describe from top to bottom. Located at the top is the graph control field. Directly below the graph control is the landscape view of the protein. Below the landscape view, a schematic and interactive representation of the protein and an additional representation of the protein which controls the zooming option. Lastly, at the bottom of the page there is the list of selected positions. All of these components are interactive and the various functionalities are described in **Table 2**.

Another way to use population-based variation in the context of the entire protein is via the tolerance landscape representation in MetaDome that can be selected in the graph control component (**Figure 1.2**). The tolerance landscape depicts a missense over synonymous ratio (also known as K_a/K_s or d_N/d_S) over a sliding window of 21 residues over the entirety of the protein of interest (e.g. calculated for ten residues left and right of each residue) based on the gnomAD dataset (**methods; Computing genetic tolerance and generating a tolerance landscape; Figure 2A**). Previously, the d_N/d_S metric has been used by others and us to measure genetic tolerance and predict disease genes (Ge, Kwok, & Shieh, 2015; Gilissen et al., 2014; Lelieveld et al., 2017), and it is suitable for measuring tolerance in regions within genes (Ge et al., 2016).

An example of using the MetaDome web server for variant interpretation

The MetaDome analysis result for *CDK13* (**Figure 1**) is the longest protein coding transcript for *CDK13* with a protein sequence length of 1,512 amino acids. In the resulting schematic protein representation we can observe the Pkinase Pfam protein domain (PF00069) between positions 707 and 998 as the only protein domain in this gene (**Figure 2B**). The Pkinase domain is highly prevalent throughout the human genome with as many as 779 homologous occurrences in human proteins, of which 353 are unique genomic regions. It is the 8th most occurring domain in our mapping database. The meta-domain landscape is the default view mode and shows any missense variation found in homologous domain occurrences throughout the human genome. Population-based

(gnomAD) missense variation is displayed in green and pathogenic (ClinVar) missense variation is annotated in red bars, with the height of the bars depicting the number of variants found at each position (**Figure 2B**).

At the ‘Display ClinVar variants’ the user is provided two options; to highlight all known pathogenic information known for the current protein and/or highlight any ClinVar variants that are present at homologous positions (**Figure 2A**). All ClinVar variants highlighted are displayed in red. In total six known disease-causing SNVs are present in the *CDK13* gene itself according to ClinVar, and these all fall within the Pkinase protein domain. All of these are missense variants. If we add variants found in homologous domains there are 64 positions with one or more reported pathogenic variants (**Supp. Data S1**). Four of these positions overlap with the positions on which ClinVar variants were found in the gene itself and on position p.883 (**Supp. Figure S4**) we can observe a peak of eight missense variants annotated from other protein domains.

MetaDome helps to look in more detail to a position of interest. If we do this for protein position 714 (**Figure 2C**) in *CDK13* we find that it corresponds to consensus position 10 in the Pkinase domain (PF00069). At this position in *CDK13* there are two variants reported in ClinVar: p.Gly714Arg (ClinVar ID: 375738) submitted by (Sifrim et al., 2016), and p.Gly714Asp (ClinVar ID: 449224) submitted by GeneDX. The first is reported as a *de novo* variant and is associated to Congenital Heart Defects, Dysmorphic Facial Features, and Intellectual Developmental Disorder. For the second there is no associated phenotype provided. As MetaDome annotates variants reported at homologous positions, we can find even more information for this particular position. At the homologues aligned to this position we find a variant of identical change in *PRKDI*: p.Gly600Arg (ClinVar ID: 375740) reported as pathogenic and *de novo* in the same study (Sifrim et al., 2016). It is also associated to Congenital Heart Defects as well as associated to Ectodermal Dysplasia. There are three more reported pathogenic variants aligned to this position: *MAK*:p.Gly13Ser (ClinVar ID: 29783) associated to Retinitis Pigmentosa 62 (Özgül et al., 2011), *PRKCG*:p.Gly360Ser (ClinVar ID: 42129) associated to Spinocerebellar Ataxia Type14 (Klebe et al., 2005), and *CIT*:p.Gly106Val (ClinVar ID: 254134) associated to Microcephaly 17, primary, autosomal recessive (Özgül et al., 2011). These homologously related pathogenic variants and the severity of the associated phenotypes contributes to the evidence that this particular residue may be important at this position. Further evidence can be found from the fact that in human homologue domains this residue is extremely conserved. There are 330 unique genomic regions encoding for a codon aligned to this position (**Supp. Data S2**). Only in the gene *PIK3R4* (ENST00000356763.3) does this codon encode for another residue than Glycine, namely a Threonine at position p.Thr35.

In the same way that we explored pathogenic ClinVar variation we can also explore the variation reported in gnomAD. In *CDK13* at protein position 714 there is no reported variant in gnomAD, but there are homologously related variations. There are 65 missense variants with average allele frequency of 1.24E-05 and 76 synonymous with average allele frequency 8.71E-03 and there is no reported nonsense variation (**Supp. Data S1**).

When we inspect the tolerance landscape for *CDK13* (**Figure 2A**) we can see that all of the ClinVar variants (either annotated in *CDK13* or related via homologues) fall within the Pkinase Pfam protein domain (PF00069). In addition, the protein domain can clearly be seen as more intolerant to missense variation as compared to other parts of this protein, thereby supporting the ClinVar variants likely pathogenic role.

Conclusion

The MetaDome web server combines resources and information from different fields of expertise (e.g. genomics and proteomics) in order to increase the power in analysing population and pathogenic variation by transposing this variation to homologous protein domains. Such a transfer of information is achieved by a per-position mapping between the GENCODE and Swiss-Prot databases. 79.4% of the Human Swiss-Prot protein sequences are of identical match to one or more of 42,116 GENCODE transcripts. This means that 25.7% of the GENCODE transcriptions differ in mRNA but translate to the same Swiss-Prot protein sequence. GENCODE previously reported that this is due to alternative splicing, of which a substantial proportion only affect untranslated regions (UTRs) and thus have no impact on the protein-coding part of the gene (Harrow et al., 2006).

MetaDome is especially informative if a variant of interest falls within a protein domain that has homologues. This is highly likely as 43.6% of the positions in the MetaDome mapping database are part of a homologous protein domain. Pathogenic missense variation is also highly likely to fall within a protein domain as we previously observed for 71% of HGMD and 72% of ClinVar pathogenic missense variants (Wiel et al., 2017). By aggregating variation over protein domain homologues via MetaDome, the resolution of genetic tolerance at a single amino-acid is increased. Furthermore, we can obtain variation that could disrupt the functionality of a protein domain, as annotated throughout the entire human genome, which may potentially be disease-causing. It should be noted, that by aggregating genetic variation in this way the specific context such as haplotype information or interactions with other proteins may be lost. Aggregation via meta-domains only encapsulates general biological or molecular functions attributed to the domain. Nonetheless, we believe MetaDome can be used to better interpret variants of unknown significance through the use of meta-domains and tolerance landscapes as we have shown in our example.

As more genetic data accumulates in the years to come, MetaDome will become more and more accurate in predictions of intolerance at the base-pair level and the meta-domain landscapes will become even more populated with variation found in homologue protein domains. We can imagine many other ways of integrating this type of information to be helpful for variant interpretation. Future directions for the MetaDome web server could lead to machine learning empowered variant effect prediction, or visualization of the meta-domain information in a protein 3D structure.

Methods

Software architecture of MetaDome

MetaDome is developed in Python v3.5.1 (Rossum & Drake, 2010) and makes use of the Flask framework v0.12.4 (Ronacher, 2010) for the web server part which communicates between the front-end, the back-end, and the database. The software architecture (**Supp. Figure S5**) follows the Domain-driven design paradigm (Evans, 2004). The entities in the domain part of this software architecture are rich data representations that are based on the internal database (**Creating the mapping database**) and annotations from external resources. These entities are stored after their first creation and afterward directly used for data retrieval to make the lookup in MetaDome as efficient as possible. The code is open source and can be found at our GitHub repository:

<https://github.com/cmbi/metadome>. Detailed instructions on how to deploy the MetaDome web server can be found there too.

To ensure MetaDome can be deployed to any environment and provide a high degree of modularity, we have containerized the application via Docker v17.12.1 (Hykes, 2013). We use docker-compose v1.17.1 to ensure that different containerized aspects of the MetaDome server can work together. The following aspects are containerized to this purpose: 1.) The Flask application, 2.) a PostgreSQL v10 database wherein the mapping database is stored, 3.) a Celery v4.2.0 task queue management system to facilitate the larger tasks of the MetaDome web-based user requests, 4.) a Redis v4.0.11 for task result storage, and 5.) RabbitMQ v3.7 to mediate as a task broker between client and workers. For a full overview of the docker-compose architecture we refer to **Supp. Figure S6**.

The visualization medium of the MetaDome web server is a fully interactive and responsive HTML web page. This page is generated by the Flask framework and the navigation aesthetics are made using the CSS framework Bulma v0.7.1 (Thomas, Potiekhin, Lauhakari, Shah, & Berning, 2018). The visualizations of the various landscapes and the schematic protein are created with JavaScript, JQuery v3.3.1, and the D3 Framework v4.13.0 (Bostock, Ogievetsky, & Heer, 2011). As the visualization by the D3 Framework is highly dependent on the user's cpu power, so are the visualizations of MetaDome.

Datasets of population and disease-causing genetic variation

MetaDome makes use of single nucleotide variants (SNVs) from population and clinically relevant genetic variation databases. Population variation was obtained from the gnomAD r2.0.2 VCF file by selecting all synonymous, nonsense, and missense variants that meet the PASS filter criteria. Variants meeting the PASS criteria are considered to be true variants (Lek et al., 2016). The variants in the VCF file from ClinVar release 2018 05 03 with disease-causing (Pathogenic) status are used as the disease-causing SNVs in MetaDome.

Creating the mapping database

MetaDome stores a complete mapping between genomic, protein positions, and all domain annotations (**Supp. Figure S7**) in a PostgreSQL relational database (PostgreSQL Global Development Group, 1996). This mapping is auto-generated and stored in the PostgreSQL database by the MetaDome web server upon the first run. The genomic positions consist of each chromosomal position in the protein-coding transcripts of the GENCODE release 19 GRCh37.p13 Basic set (Harrow et al., 2012). The protein positions correspond to protein sequence positions in the UniProtKB/Swiss-Prot Release 2016_09 databank entries for the human species (Boutet et al., 2016). These mappings are created with Protein-Protein BLAST v2.2.31+ (Camacho et al., 2009) for each protein-coding translation in the GENCODE Basic set to human canonical and isoform Swiss-Prot protein sequences. We exclude sequences that do not start with a start codon (i.e. ATG encoding for methionine), or end with a stop codon. We checked if the cDNA sequence of the transcripts match the GENCODE translation via Biopython's translate function (Cock et al., 2009), if they are not identical then these are excluded too. The global information on the transcript (e.g. identifiers, sequence length) is registered in the database in the table 'genes' and, for each Swiss-Prot entry with an identical sequence match, the global information is stored in the table 'proteins'. All tables are indexed by the fields that are used in the lookups.

Next, for each identical match between translation and Swiss-Prot sequence a ClustalW2 v2.1 (Larkin et al., 2007) alignment is made between these two sequences. Each nucleotide's genomic position is mapped to the protein position and stored in the 'mappings' table. Each entry in mapping represents a single nucleotide of a codon and is linked to the corresponding entry in the 'genes' and 'proteins' table (i.e. the corresponding GENCODE translation, transcription and Swiss-Prot sequence).

Each Swiss-Prot sequence in the database is annotated via InterProScan v5.20-59.0 (Finn et al., 2017) for Pfam-A v30.0 protein domains (Finn et al., 2016) and the results are stored in the 'interpro_domains' table. After the construction of the database is finished, all meta-domain alignments can be constructed.

Composing a meta-domain

Meta-domains consist of homologous Pfam protein domain instances that are annotated using InterProScan. Meta-domains consist of domains that have at least two homologues within the human genome. MSAs are made using a three step process. 1.) Retrieve all sequences for the domain instances, 2.) Retrieve the Pfam HMM corresponding to the Pfam identifier annotated by InterProScan, and 3.) Use HMMER 3.1b2 (Finn et al., 2015) to align the sequences from the first step. The resulting Stockholm format MSA files can be inspected with alignment visualization software like Jalview (Waterhouse, Procter, Martin, Clamp, & Barton, 2009). In this Stockholm formatted file, all columns that correspond to the domain consensus represent the same homologous positions.

These Stockholm files are retrieved by the MetaDome web server when a user request meta-domain information for a position of their interest. Upon retrieval of this Stockholm file, the mapping database is used to obtain the corresponding genomic positions for each residue. These genomic positions are subsequently used to retrieve corresponding gnomAD or ClinVar variation.

Computing genetic tolerance and generating a tolerance landscape

The non-synonymous over synonymous ratio, or d_N/d_S score, is used to quantify genetic tolerance. This score is based on the observed (obs) missense and synonymous variation in gnomAD ($missense_{obs}$ and $synonymous_{obs}$). This score is corrected for the sequence composition by taking into account the background (bg) of possible missense and synonymous variants based on the codon table ($missense_{bg}$ and $synonymous_{bg}$):

$$d_N/d_S = \frac{missense_{obs}/missense_{bg}}{synonymous_{obs}/synonymous_{bg}}$$

The tolerance landscape computes this ratio as a sliding window of size 21 (i.e. ten residues before and ten after the residue of interest) over the entirety of the gene's protein, similar to the Missense Tolerance Ratio (MTR) presented by (Traynelis et al., 2017). The edges (e.g. start and end) are therefore a bit noisy as they are not the result of averaging over a full length window.

Acknowledgements

This work was in part financially supported by grants from the Netherlands Organization for Scientific Research (916-14-043 to C.G. and 918-15-667 to J.A.V.), and from the Radboud Institute for Molecular Life Sciences, Radboud university medical center (R0002793 to G.V.). We thank Hanka Venselaar for her critical reading of the manuscript.

References

- Altschul, S. F., Gish, W., Miller, W., Myers, E. W., & Lipman, D. J. (1990). Basic local alignment search tool. *Journal of Molecular Biology*, *215*(3), 403–410. [https://doi.org/10.1016/S0022-2836\(05\)80360-2](https://doi.org/10.1016/S0022-2836(05)80360-2)
- Amr, S. S., Al Turki, S. H., Lebo, M., Sarmady, M., Rehm, H. L., & Abou Tayoun, A. N. (2016). Using large sequencing data sets to refine intragenic disease regions and prioritize clinical variant interpretation. *Genetics in Medicine : Official Journal of the American College of Medical Genetics*, (July), 1–9. <https://doi.org/10.1038/gim.2016.134>
- Ashenberg, O., Gong, L. I., & Bloom, J. D. (2013). Mutational effects on stability are largely conserved during protein evolution. *Proceedings of the National Academy of Sciences of the United States of America*, *110*(52), 21071–6. <https://doi.org/10.1073/pnas.1314781111>
- Bostock, M., Ogievetsky, V., & Heer, J. (2011). D³: Data-Driven Documents. *IEEE Transactions on Visualization and Computer Graphics*, *17*(12), 2301–9. <https://doi.org/10.1109/TVCG.2011.185>
- Boutet, E., Lieberherr, D., Tognolli, M., Schneider, M., Bansal, P., Bridge, A. J., ... Xenarios, I. (2016). UniProtKB/Swiss-Prot, the Manually Annotated Section of the UniProt KnowledgeBase: How to Use the Entry View. *Methods in Molecular Biology (Clifton, N.J.)*, *1374*, 23–54. https://doi.org/10.1007/978-1-4939-3167-5_2
- Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., & Madden, T. L. (2009). BLAST+: architecture and applications. *BMC Bioinformatics*, *10*(1), 421. <https://doi.org/10.1186/1471-2105-10-421>
- Cock, P. J. A., Antao, T., Chang, J. T., Chapman, B. A., Cox, C. J., Dalke, A., ... de Hoon, M. J. L. (2009). Biopython: freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics (Oxford, England)*, *25*(11), 1422–3. <https://doi.org/10.1093/bioinformatics/btp163>
- Evans, E. (2004). *Domain-driven design: tackling complexity in the heart of software*. Addison-Wesley Professional.
- Finn, R. D., Attwood, T. K., Babbitt, P. C., Bateman, A., Bork, P., Bridge, A. J., ... Mitchell, A. L. (2017). InterPro in 2017—beyond protein family and domain annotations. *Nucleic Acids Research*, *45*(D1), D190–D199. <https://doi.org/10.1093/nar/gkw1107>

- Finn, R. D., Clements, J., Arndt, W., Miller, B. L., Wheeler, T. J., Schreiber, F., ... Eddy, S. R. (2015). HMMER web server: 2015 update. *Nucleic Acids Research*, *43*(W1), W30–W38. <https://doi.org/10.1093/nar/gkv397>
- Finn, R. D., Coghill, P., Eberhardt, R. Y., Eddy, S. R., Mistry, J., Mitchell, A. L., ... Bateman, A. (2016). The Pfam protein families database: towards a more sustainable future. *Nucleic Acids Research*, *44*(D1), D279–D285. <https://doi.org/10.1093/nar/gkv1344>
- Fu, W., O'Connor, T. D., Jun, G., Kang, H. M., Abecasis, G., Leal, S. M., ... Akey, J. M. (2012). Analysis of 6,515 exomes reveals the recent origin of most human protein-coding variants. *Nature*, *493*(7431), 216–220. <https://doi.org/10.1038/nature11690>
- Ge, X., Gong, H., Dumas, K., Litwin, J., Phillips, J. J., Waisfisz, Q., ... Shieh, J. T. C. (2016). Missense-depleted regions in population exomes implicate ras superfamily nucleotide-binding protein alteration in patients with brain malformation. *Npj Genomic Medicine*, *1*(1), 16036. <https://doi.org/10.1038/npjgenmed.2016.36>
- Ge, X., Kwok, P.-Y., & Shieh, J. T. C. (2015). Prioritizing genes for X-linked diseases using population exome data. *Human Molecular Genetics*, *24*(3), 599–608. <https://doi.org/10.1093/hmg/ddu473>
- Gilissen, C., Hehir-Kwa, J. Y., Thung, D. T., van de Vorst, M., van Bon, B. W. M., Willemsen, M. H., ... Veltman, J. A. (2014). Genome sequencing identifies major causes of severe intellectual disability. *Nature*, *511*(7509), 344–347. <https://doi.org/10.1038/nature13394>
- Gussow, A. B., Petrovski, S., Wang, Q., Allen, A. S., & Goldstein, D. B. (2016). The intolerance to functional genetic variation of protein domains predicts the localization of pathogenic mutations within genes. *Genome Biology*, *17*(1), 9. <https://doi.org/10.1186/s13059-016-0869-4>
- Harrow, J., Denoeud, F., Frankish, A., Reymond, A., Chen, C.-K., Chrast, J., ... Guigo, R. (2006). GENCODE: producing a reference annotation for ENCODE. *Genome Biology*, *7 Suppl 1*(Suppl 1), S4.1-9. <https://doi.org/10.1186/gb-2006-7-s1-s4>
- Harrow, J., Frankish, A., Gonzalez, J. M., Tapanari, E., Diekhans, M., Kokocinski, F., ... Hubbard, T. J. (2012). GENCODE: The reference human genome annotation for The ENCODE Project. *Genome Research*, *22*(9), 1760–1774. <https://doi.org/10.1101/gr.135350.111>
- Hykes, S. (2013). Docker. San Francisco: Docker, Inc. Retrieved from <https://www.docker.com/>

- Karczewski, K. J., Weisburd, B., Thomas, B., Solomonson, M., Ruderfer, D. M., Kavanagh, D., ... MacArthur, D. G. (2017). The ExAC browser: Displaying reference data information from over 60 000 exomes. *Nucleic Acids Research*, 45(D1), D840–D845. <https://doi.org/10.1093/nar/gkw971>
- Klebe, S., Durr, A., Rentschler, A., Hahn-Barma, V., Abele, M., Bouslam, N., ... Stevanin, G. (2005). New mutations in protein kinase C γ associated with spinocerebellar ataxia type 14. *Annals of Neurology*, 58(5), 720–729. <https://doi.org/10.1002/ana.20628>
- Lal, D., May, P., Samocha, K. E., Kosmicki, J. A., Robinson, E. B., Møller, R. S., ... Daly, M. J. (2017). Gene family information facilitates variant interpretation and identification of disease-associated genes. *BioRxiv*, 159780. <https://doi.org/10.1101/159780>
- Landrum, M. J., Lee, J. M., Benson, M., Brown, G., Chao, C., Chitipiralla, S., ... Maglott, D. R. (2016). ClinVar: public archive of interpretations of clinically relevant variants. *Nucleic Acids Research*, 44(D1), D862-8. <https://doi.org/10.1093/nar/gkv1222>
- Larkin, M. A., Blackshields, G., Brown, N. P., Chenna, R., McGettigan, P. A., McWilliam, H., ... Higgins, D. G. (2007). Clustal W and Clustal X version 2.0. *Bioinformatics*, 23(21), 2947–2948. <https://doi.org/10.1093/bioinformatics/btm404>
- Lek, M., Karczewski, K. J., Minikel, E. V., Samocha, K. E., Banks, E., Fennell, T., ... MacArthur, D. G. (2016). Analysis of protein-coding genetic variation in 60,706 humans. *Nature*, 536(7616), 285–291. <https://doi.org/10.1038/nature19057>
- Lelieveld, S. H., Wiel, L., Venselaar, H., Pfundt, R., Vriend, G., Veltman, J. A., ... al., et. (2017). Spatial Clustering of de Novo Missense Mutations Identifies Candidate Neurodevelopmental Disorder-Associated Genes. *The American Journal of Human Genetics*, 8(0), 52–56. <https://doi.org/10.1016/j.ajhg.2017.08.004>
- NHLBI GO Exome Sequencing Project (ESP). (2011). Exome Variant Server. Retrieved May 14, 2015, from <http://evs.gs.washington.edu/EVS/>
- Özgül, R. K., Siemiatkowska, A. M., Yücel, D., Myers, C. A., Collin, R. W. J., Zonneveld, M. N., ... Corbo, J. C. (2011). Exome sequencing and cis-regulatory mapping identify mutations in MAK, a gene encoding a regulator of ciliary length, as a cause of retinitis pigmentosa. *American Journal of Human Genetics*, 89(2), 253–264. <https://doi.org/10.1016/j.ajhg.2011.07.005>
- Petrovski, S., Wang, Q., Heinzen, E. L., Allen, A. S., & Goldstein, D. B. (2013). Genic

Intolerance to Functional Variation and the Interpretation of Personal Genomes. *PLoS Genetics*, 9(8), e1003709. <https://doi.org/10.1371/journal.pgen.1003709>

PostgreSQL Global Development Group. (1996). PostgreSQL. PostgreSQL Global Development Group. Retrieved from <https://www.postgresql.org/>

Ronacher, A. (2010). Flask. Retrieved from <http://flask.pocoo.org/>

Rossum, G. Van, & Drake, F. L. (2010). Python Tutorial. *History*, 42(4), 1–122. https://doi.org/10.1111/j.1094-348X.2008.00203_7.x

Sifrim, A., Hitz, M.-P., Wilsdon, A., Breckpot, J., Turki, S. H. Al, Thienpont, B., ... Hurles, M. E. (2016). Distinct genetic architectures for syndromic and nonsyndromic congenital heart defects identified by exome sequencing. *Nature Genetics*, 48(9), 1060–5. <https://doi.org/10.1038/ng.3627>

Stenson, P. D., Mort, M., Ball, E. V., Evans, K., Hayden, M., Heywood, S., ... Cooper, D. N. (2017). The Human Gene Mutation Database: towards a comprehensive repository of inherited mutation data for medical research, genetic diagnosis and next-generation sequencing studies. *Human Genetics*, 136(6), 1–13. <https://doi.org/10.1007/s00439-017-1779-6>

Thomas, J., Potiekhin, O., Lauhakari, M., Shah, A., & Berning, D. (2018). *Creating Interfaces with Bulma*. (T. Mott & D. Berning, Eds.). Santa Rosa: Bleeding Edge Press. Retrieved from <https://bleedingedgepress.com/>

Traynelis, J., Silk, M., Wang, Q., Berkovic, S. F., Liu, L., Ascher, D. B., ... Petrovski, S. (2017). Optimizing genomic medicine in epilepsy through a gene-customized approach to missense variant interpretation. *Genome Research*, 27(10), 1715–1729. <https://doi.org/10.1101/gr.226589.117>

Waterhouse, A. M., Procter, J. B., Martin, D. M. A., Clamp, M., & Barton, G. J. (2009). Jalview Version 2--a multiple sequence alignment editor and analysis workbench. *Bioinformatics*, 25(9), 1189–1191. <https://doi.org/10.1093/bioinformatics/btp033>

Wiel, L., Venselaar, H., Veltman, J. A., Vriend, G., & Gilissen, C. (2017). Aggregation of population-based genetic variation over protein domain homologues and its potential use in genetic diagnostics. *Human Mutation*, (May), 1–10. <https://doi.org/10.1002/humu.23313>

Figures

Figure 1. MetaDome web server result for the gene *CDK13*. The result provided by the MetaDome web server for the analysis of gene *CDK13* with transcript ENST00000181839.4, as provided in 1.). In 2.), there is additional information that the translation of this transcript corresponds to Swiss-Prot protein Q14004. Here also various alternative visualizations can be selected. The visualization starts by default in the ‘meta-domain landscape’, a mode selectable in the graph control in 2.). The landscapes are visualized in 3.), and in the meta-domain landscape the domain regions are annotated with missense variation counts found in homologous domains as bar plots. The schematic protein representation, located at 4.), is per-position selectable, and the domains are presented as purple blocks. Selected positions are highlighted in green. The ‘Zoom-in’ section at 5.) features a selectable greyed-out copy of schematic protein representation that can zoom-in on any part of the protein. Any selected positions are in the list of selected positions in 6.). Here more information can be obtained by clicking on one of these positions. A detailed description of the functionality of each component is described in **Table 2**.

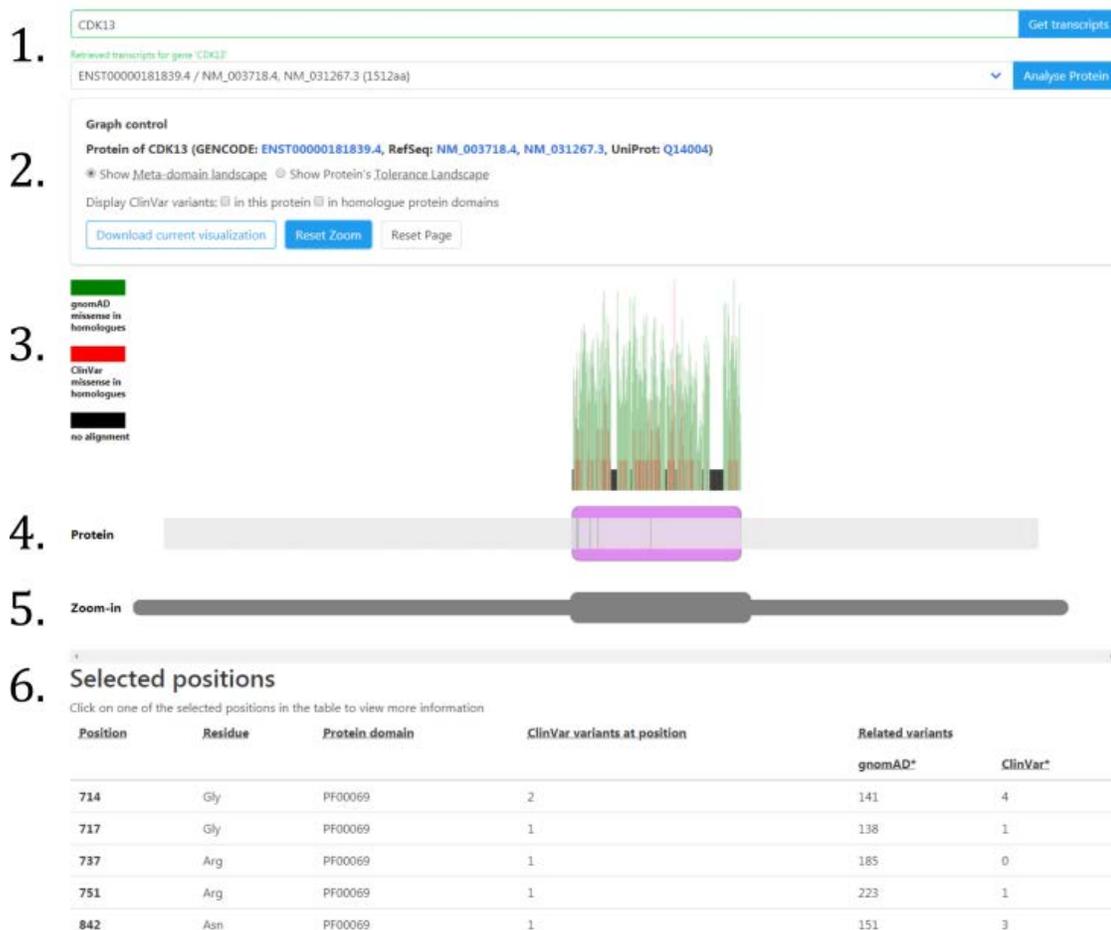
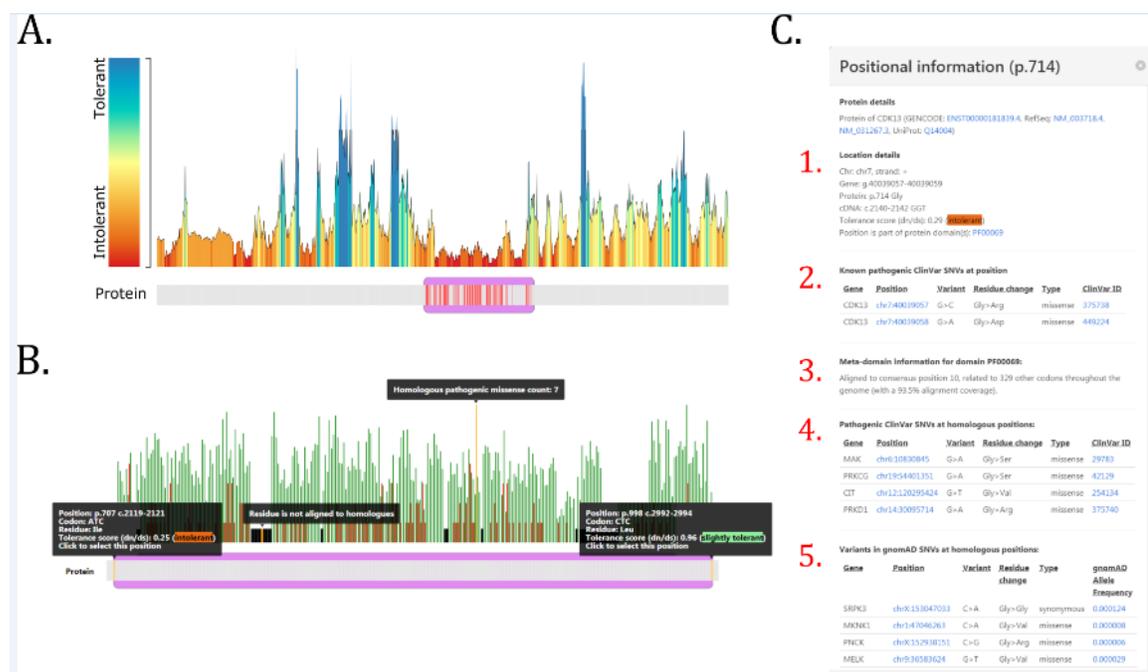


Figure 2. Examples of a MetaDome analysis for the gene *CDK13* **A.)** The tolerance landscape depicts a missense over synonymous ratio calculated as a sliding window over the entirety of the protein (**methods; Computing genetic tolerance and generating a tolerance landscape**). The missense and synonymous variation are annotated from the gnomAD dataset and the landscape provides some indication of regions that are intolerant to missense variation. In this *CDK13* tolerance landscape the Pkinase Pfam protein domain (PF00069) in purple can be clearly seen as intolerant if compared to other parts in this protein. The red bars in the schematic protein representation correspond to pathogenic ClinVar variants found in this gene and in homologous protein domains. All of these variants are contained in the intolerant region of the landscape. **B.)** A zoom-in on the meta-domain landscape for *CDK13*. The Pkinase Pfam protein domain (PF00069) is located between protein positions 707 and 998 and annotated as a purple box in the schematic protein representation. The meta-domain landscape displays a deep annotation of the protein domain: the green (gnomAD) and red (ClinVar) bars correspond to the number of missense variants found at aligned homologous positions. Unaligned positions are annotated as black bars. All of this information is displayed upon hovering over these various elements. **C.)** The positional information provides a detailed overview of a position from the ‘Selected Positions’ list, especially if that position is aligned to domain homologues. Here, for position p.Gly714 we can observe in 1.) the positional details for this specific protein position. In 2.) is any known pathogenic information for this position. We can observe here that for this position there are two known pathogenic missense variants. In 3.) meta-domain information is displayed and we can observe that p.Gly714 is aligned to consensus position 10 in the Pkinase Pfam protein domain and related to 329 other codons. This consensus position has an alignment coverage of 93.5% for the meta-domain MSA. There are also four pathogenic variants found in ClinVar on corresponding homologous positions as can be seen in 4.) and in 5.) there is an overview of all corresponding variants found in gnomAD.



Tables

Table 1. Statistics on the number of entries present in GENCODE, Swiss-Prot, and our mapping database.

Database	What	# of entries
GENCODE	Protein-coding genes	20,345
MetaDome	Protein-coding genes	19,728
GENCODE	Protein-coding transcripts	57,005
MetaDome	Protein-coding transcripts	56,319
Swiss-Prot	Canonical and isoform protein sequences	591,556
Swiss-Prot	Human canonical and isoform protein sequences	42,130
MetaDome	Gene translations identically mapped to a canonical or isoform protein sequence	42,116
MetaDome	Canonical and isoform protein sequences	33,492
MetaDome	Pfam protein domain regions	71,419
MetaDome	Unique Pfam protein domain families	5,948
MetaDome	Unique Pfam protein domain families with two or more within-human occurrences	3,334
MetaDome	Chromosome to protein position mappings	70,261,143

MetaDome	Unique chromosome positions	32,595,355
MetaDome	Unique residues (as part of a protein)	19,226,961
MetaDome	Unique protein sequences with at least one Pfam domain annotated	30,406

Table 2. Descriptions of the various functionalities on the MetaDome result page.

Component	Functionality
Gene and transcript input field (Figure 1.1)	<ul style="list-style-type: none"> • Input of gene of interest • Retrieving transcripts for gene of interest • Selecting a transcript • Starting the analysis for selected transcript
Graph control field (Figure 1.2)	<ul style="list-style-type: none"> • Toggling between different landscape representations • Reset the zoom on the landscape • Reset the web page • Toggle ClinVar variants to be displayed in the schematic protein • Download the visual representation
Landscape view (Figure 1.3)	<ul style="list-style-type: none"> • Displays the meta-domain landscape • Displays the tolerance landscape
Schematic protein (Figure 1.4)	<ul style="list-style-type: none"> • Displays a schematic representation of the gene's protein with Pfam protein domains annotated • Hovering over a position displays positional information • Clicking on a position highlights the position and adds the position to the list of 'Selected Positions' • Controls the zooming of particular parts of the

	protein (Figure 1.5)
Selected Positions (Figure 1.6)	<ul style="list-style-type: none">• Displays any positions selected in the schematic protein• Displays per selected position: if that position is part of a Pfam protein domain, any known gnomAD or ClinVar variants present at this position, and any variants that are homologously related to this position• Provides more detailed information as a pop-up when clicking on one of the positions in this list.