

Genome analysis

ProteoDisco: a flexible R approach to generate customized protein databases for extended search space of novel and variant proteins in proteogenomic studies

Wesley S. van de Geer ^{1,2,3,†}, Job van Riet ^{1,2,3,†} and Harmen J. G. van de Werken ^{1,3,4,*}

¹Cancer Computational Biology Center, Erasmus MC Cancer Institute, University Medical Center, 3015 GD Rotterdam, The Netherlands,

²Department of Medical Oncology, Erasmus MC Cancer Institute, University Medical Center, 3015 GD Rotterdam, The Netherlands,

³Department of Urology Erasmus MC Cancer Institute, University Medical Center, 3015 GD Rotterdam, The Netherlands and

⁴Department of Immunology, Erasmus MC Cancer Institute, University Medical Center, 3015 GD Rotterdam, The Netherlands

*To whom correspondence should be addressed.

†The authors wish it to be known that, in their opinion, the first two authors should be regarded as Joint First Authors.

Associate Editor: Can Alkan

Received on September 12, 2021; revised on November 8, 2021; editorial decision on November 23, 2021; accepted on November 26, 2021

Abstract

Summary: We present an R-based open-source software termed ProteoDisco that allows for flexible incorporation of genomic variants, fusion genes and (aberrant) transcriptomic variants from standardized formats into protein variant sequences. ProteoDisco allows for a flexible step-by-step workflow allowing for in-depth customization to suit a myriad of research approaches in the field of proteogenomics, on all organisms for which a reference genome and transcript annotations are available.

Availability and implementation: ProteoDisco (R package version $\geq 1.0.0$) is available on Bioconductor at <https://doi.org/doi:10.18129/B9.bioc.ProteoDisco> and from <https://github.com/ErasmusMC-CCBC/ProteoDisco/>.

Contact: h.vandewerken@erasmusmc.nl

Supplementary information: [Supplementary data](#) are available at *Bioinformatics* online.

1 Introduction

The rise and ease of current next-generation sequencing (NGS) techniques, coupled with reduced costs in both NGS and high-resolution mass-spectrometry, offers opportunity to incorporate sample-specific protein variants during proteomics experiments for increased accuracy and detection rates of, for instance, distinctive proteotypic peptides in bottom-up proteomics experiments. Expanding the repertoire of proteins and these proteotypic peptides can provide novel insights into disease-specific protein variants, their underlying molecular profiles and regulation, neoantigen prediction and expand our knowledge on the genetic variations encoded in proteomes (Mertins *et al.*, 2016; Nesvizhskii, 2014; Ruggles *et al.*, 2016; Vasaikar *et al.*, 2019; Wen *et al.*, 2020). This is further fueled by the standardization and publication of proteomics resources which allows for the interrogation and combination of existing datasets (Deutsch *et al.*, 2017; Zahn-Zabal *et al.*, 2020).

Rising global efforts in capturing the genetic sequences of diverse organisms, disease-related genotypes and their transcriptomes with subsequent proteome-resources warrants the implementation of a flexible yet intuitive toolset. This toolset should provide a bridge

between genomic and transcriptomic variants and their incorporation within respective protein variants (proteogenomics) using industry-standard infrastructure, such as Bioconductor (Gentleman *et al.*, 2004), and allow for flexibility in facilitating the myriad experimental settings applied in research. Therefore, we designed and developed ProteoDisco, an open-source R software-package using existing Bioconductor class-infrastructures to allow for the accurate and flexible generation of variant protein sequences and their derived proteotypic peptides from the incorporation of sample-specific genomic and transcriptomic information. In addition, we present the results of ProteoDisco and two similar open-source tools which are frequently utilized within proteogenomics [customProDB (Wang and Zhang, 2013) and QUILTS (Ruggles *et al.*, 2016)] with their performance in generating correct protein variants and respective proteotypic peptides from supplied genomic variants.

2 Approach

ProteoDisco incorporates genomic variants, splice-junctions (derived from transcriptomics) and fusion genes within provided

reference genome sequences and transcript annotations to generate their respective protein variant sequence(s). These sequences can be curated, altered and subsequently exported into a database in FASTA format for use in downstream analysis. To limit the number of generated protein variants, ProteoDisco provides filtering options based on a minimal number of distinct proteotypic (identifiable) peptides. The global workflow of ProteoDisco is summarized in six steps as depicted within Figure 1. In addition, an extended overview of how (novel) splice-junctions and gene-fusion events are incorporated is shown in Supplementary Figure S1.

To compare the accuracy of ProteoDisco against two common alternatives for proteogenomics studies [customProDB (Wang and Zhang, 2013) and QUILTS (Ruggles et al., 2016)], we utilized a manually curated dataset and two large independent proteomics studies. The manually curated dataset contained 28 genomic variants reported in COSMIC (Forbes et al., 2017) comprising multiple variant classes; synonymous and non-synonymous single-nucleotide variants (SNVs), multinucleotide variants (MNVs) and in- and out-of-frame insertions/deletions (InDels). In addition, we utilized recently published results from large-scale colon and breast cancer cohorts within the Clinical Proteomic Tumor Analysis Consortium (CPTAC) to illustrate the accuracy of ProteoDisco in generating identical proteotypic peptides as detected within these studies (Mertins et al., 2016; Wen et al., 2020). This comparison revealed that ProteoDisco correctly generated proteotypic peptides from their respective genomic variants after thorough checking and yielded the highest number of expected and reconstructed proteotypic peptides within all three datasets (Supplementary Fig. S2). This difference can be attributed to ProteoDisco's native flexibility in reference genome selection, multiple incorporation strategies, sanity-checks such as reference base verification and the correct incorporation of stop-loss variants. In total, only four enigmatic genomic variants (of three fragments) from Mertins et al. could not be reconstructed to reproduce their proteotypic peptide(s).

3 Materials and methods

3.1 Technical design of ProteoDisco

ProteoDisco was programmed within the R statistical language (v4.1.1) and built upon existing classes within the Bioconductor infrastructure (v3.13) to allow flexible inheritance and future extensions. Additional information on the usage and design of ProteoDisco can be found in the extended methodology (Supplementary File S1).

3.2 Assessment of the correct integration of genomic variants into protein variants

We generated a custom validation dataset containing established somatic variants (SNVs, MNVs and InDels; $n = 28$) and their respective protein variants as listed within COSMIC (Forbes et al., 2017) (v92; GRCh37; Supplementary Table S1). In addition, we utilized recent proteogenomics studies from the CPTAC cancer cohorts containing genomic variants and their respective *in silico* generated proteotypic peptides which had been measured and identified using high-throughput proteomics approaches (Mertins et al., 2016; Wen et al., 2020). In the Wen et al. dataset (Wen et al., 2020, CPTAC—Colon Cancer), genomic variants (and their respective proteotypic peptides) were split into sample-specific VCF-files based on the data present within their published Supplementary Data S1 (see reference Wen et al., 2020, sheet 1: 'prospective_colon_label_free_in'). The Mertins et al. dataset (Mertins et al., 2016, CPTAC—Breast Cancer) was aggregated into a single VCF-file based on the data present within their published Supplementary Table S5 (see reference Mertins et al., 2016, sheet 2: 'Variants').

Using these three datasets, we ran ProteoDisco (v0.99), customProDB (v1.30.1) and the web interface of QUILTS (v3.0; as available from <http://openslice.fenyolab.org/cgi-bin/pyquilt.cgi.pl>; accessed 13-04-2021) to generate custom protein-variant databases using uniform UCSC/RefSeq (Frankish et al., 2019, GRCh37) transcript annotations and settings. The custom protein-variant databases were generated based on two approaches within ProteoDisco. The first approach incorporated each genomic variant independently and the second allowed for the simultaneous incorporation of all genomic variants per overlapping transcript annotation, e.g. two variants on different coding exons would both be incorporated within the resulting variant protein sequence. Incorporation of all possible combinations of mutant exons yields too many combinations and is therefore not included amongst the options.

The generated variant protein sequences and respective proteotypic peptides from each customized protein-variant database were compared against the proteotypic peptides as expected from COSMIC or as detected within the respective CPTAC-studies using all three tools (Supplementary Fig. S2). For example, if ProteoDisco generated three distinct proteotypic peptides for a given genomic variant and one of those was identified within CPTAC (or COSMIC), it was counted as a concordant result.

4 Conclusion

In this article, we present ProteoDisco, a suitable, open-source and flexible suite for the generation of protein-variant databases usable in downstream proteogenomic studies and capable of correctly incorporating a diverse range of genomic variants and transcriptomic splice-junctions. We report that ProteoDisco accurately produces protein

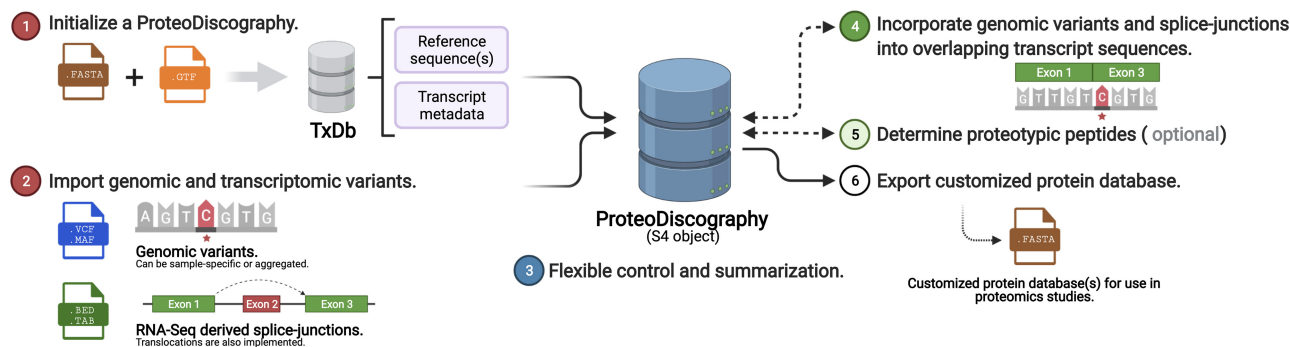


Fig. 1. Schematic overview of the ProteoDisco workflow. The global workflow of ProteoDisco can be categorized as six major steps. (1) Initialize a ProteoDiscography by utilizing custom references sequence(s) and gene-annotation(s) or using pre-existing TxDb objects. (2) Import (sample-specific) genomic variants, splice-junctions or manual sequences. Several sanity-checks are performed during importation, including the validation of matching reference nucleotide(s). (3) Dynamically view, extend, alter and customize imported records and derived sequences. (4) Incorporation of genomic variants and splice-junctions into overlapping transcript annotations, translocations between chromosomes can also be processed. The incorporation can be performed in a sample-specific manner, exon or transcript-specific manner or in an aggregated manner. (5) Cleave derived protein variant sequences and determine proteotypic peptides, per protein, which are not present within the reference protein sequences (TxDb) or additional protein databases (e.g. UniProtKB). (6) Export the derived protein variant sequences into an external protein-sequence database(s) using FASTA format

variant sequences harboring previously identified proteotypic fragments from their respective genomic variants. Further examples and use-cases can be found in the vignette of the ProteoDisco package.

4.1 Code availability

All source-code has been made available within Bioconductor (<https://doi.org/doi:10.18129/B9.bioc.ProteoDisco>) and deposited within GitHub (<https://github.com/ErasmusMC-CCBC/ProteoDisco>) under the GPL-3 license.

Acknowledgements

The authors would like to thank S. Eldeb and W.M. Hoska for their preliminary work on the integration of mutations into transcript sequences. [Figure 1](#) and [Supplementary Figure S1](#) were created using BioRender.com.

Author contributions

All authors had full access to all the data in the study and take responsibility for the integrity of the data and the accuracy of the data analysis. Study concept and design: van Riet, van de Geer, van de Werken. Acquisition of data: van Riet, van de Geer. Analysis and interpretation of data: All authors. Drafting of the manuscript: All authors. Critical revision of the manuscript for important intellectual content: All authors. Statistical analysis: All authors. Obtaining funding: None. Administrative, technical, or material support: All authors. Supervision: van de Werken. Other: None.

Funding

This research received no specific grant from any funding agency in the public, commercial, or not-for-profit sectors.

Conflict of Interest: none declared.

Data availability

The custom validation dataset (GRCh37) which has been used in the analysis as presented within this manuscript has been stored within ProteoDisco and is accessible at <https://github.com/ErasmusMC-CCBC/ProteoDisco/tree/main/inst/extdata>. COSMIC (v92; accessed on 14-04-2021) was used to derive the validation dataset (GRCh37), the external validation datasets based on CPTAC (colon and breast cancer) were generated based on the [Supplementary Data](#) published by [Wen et al. \(2020\)](#) and [Mertins et al. \(2016\)](#).

References

- Deutsch,E.W. *et al.* (2017) The ProteomeXchange consortium in 2017: supporting the cultural change in proteomics public data deposition. *Nucleic Acids Res.*, **45**, D1100–D1106.
- Forbes,S.A. *et al.* (2017) COSMIC: somatic cancer genetics at high-resolution. *Nucleic Acids Res.*, **45**, D777–D783.
- Frankish,A. *et al.* (2019) GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids Res.*, **47**, 766–773.
- Gentleman,R.C. *et al.* (2004) Bioconductor: open software development for computational biology and bioinformatics. *Genome Biol.*, **5**, R80.
- Mertins,P. *et al.* (2016) Proteogenomics connects somatic mutations to signaling in breast cancer. *Nature*, **534**, 55–62.
- Nesvizhskii,A.I. (2014) Proteogenomics: concepts, applications and computational strategies. *Nat. Methods*, **11**, 1114–1125.
- Ruggles,K.V. *et al.* (2016) An analysis of the sensitivity of proteogenomic mapping of somatic mutations and novel splicing events in cancer. *Mol. Cell. Proteomics*, **15**, 1060–1071.
- Vasaikar,S. *et al.* (2019) Proteogenomic analysis of human colon cancer reveals new therapeutic opportunities. *Cell*, **177**, 1035–1049.e19.
- Wang,X. and Zhang,B. (2013) CustomProDB: an R package to generate customized protein databases from RNA-Seq data for proteomics search. *Bioinformatics*, **29**, 3235–3237.
- Wen,B. *et al.* (2020) Cancer neoantigen prioritization through sensitive and reliable proteogenomics analysis. *Nat. Commun.*, **11**, 1759.
- Zahn-Zabal,M. *et al.* (2020) The neXtProt knowledgebase in 2020: data, tools and usability improvements. *Nucleic Acids Res.*, **48**, 328–334.