

MobiDB-lite 3.0: fast consensus annotation of intrinsic disorder flavours in proteins

Marco Necci¹, Damiano Piovesan¹, Damiano Clementel¹, Zsuzsanna Dosztányi² and Silvio C.E. Tosatto^{1,*}

¹Department of Biomedical Sciences, University of Padua, via U. Bassi 58/b, 35121 Padova, Italy

²MTA-ELTE Lendület Bioinformatics Research Group, Department of Biochemistry, ELTE Eötvös Loránd University, Pázmány Péter sétány 1/c, Budapest, Hungary

*To whom correspondence should be addressed.

Associate Editor: XXXXXXXX

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Abstract

Motivation: The earlier version of MobiDB-lite is currently used in large-scale proteome annotation platforms to detect intrinsic disorder. However, new theoretical models allow for the classification of intrinsically disordered regions into subtypes from sequence features associated with specific polymeric properties or compositional bias.

Results: MobiDB-lite 3.0 maintains its previous speed and performance but also provides a finer classification of disorder by identifying regions with characteristics of polyolyampholytes, positive or negative polyelectrolytes, low complexity regions or enriched in cysteine, proline or glycine or polar residues. Sub-regions are abundantly detected in IDRs of the human proteome. The new version of MobiDB-lite represents a new step for the proteome level analysis of protein disorder.

Availability: Both the MobiDB-lite 3.0 source code and a docker container are available from the GitHub repository: <https://github.com/BioComputingUP/MobiDB-lite>.

Contact: silvio.tosatto@unipd.it

1 Introduction

The identification of protein domains and sequence conservation has long been central to the annotation of proteomes (Lee *et al.*, 2005; Sonnhammer *et al.*, 1997). Many proteins, known as intrinsically disordered, have been observed to escape the typical organization of globular proteins in domains (Dunker *et al.*, 2001). A large fraction of the human proteome is devoid of domains (Mistry *et al.*, 2013) and in this ‘dark’ proteome molecular conformations are completely unknown (Perdigão *et al.*, 2015). Computational prediction of intrinsic disorder (ID) attempts to fill this gap by offering a wide array of prediction methods (Walsh *et al.*, 2012; Mészáros *et al.*, 2018) with different performances (Walsh *et al.*, 2015; Necci *et al.*, 2018). Despite many methods having been available for a long time, they had not been integrated into large-scale proteome annotation. MobiDB-lite (Necci *et al.*, 2017) was the first of such predictors to be included in InterProScan from its release 60 (Mitchell *et al.*,

2019). MobiDB-lite combines a set of complementary ID predictors in a consensus optimized on a PDB X-ray dataset (Walsh *et al.*, 2015) to limit over-prediction while balancing under-prediction.

In recent years theoretical models of ID proteins surpassed the bare distinction between disorder and structure and reached a point where classification of subtypes of disorder is possible based on sequence features (Das and Pappu, 2013; Holehouse *et al.*, 2017). Furthermore, recent evidence highlighted how ID and low sequence complexity (LC) are strictly intertwined (Mier *et al.*, 2020). For this reason, we developed a new version of MobiDB-lite, which can capture different classes of disorder and sequence features that we observed being biologically relevant in IDPs (Necci *et al.*, 2016). MobiDB-lite 3.0 is already included in the latest versions of MobiDB (Piovesan *et al.*, 2018). The new version is available as docker container and also exposes bindings to use MobiDB-lite as a python library, in compliance with the FAIR principles (Wilkinson *et al.*, 2016).

2 Implementation

MobiDB-lite disorder prediction unfolds in two steps as explained as in (Necci et al., 2017). Briefly, the first step calculates a strict majority (i.e. > 4) consensus between 8 predictors (Linding, Jensen, et al., 2003; Dosztányi et al., 2005; Walsh et al., 2012; Linding, Russell, et al., 2003; Peng et al., 2006), which in the second step is smoothed out in a process similar to dilation-erosion morphological operations and filtered to keep only regions longer than 20 residues (Necci et al., 2017). MobiDB-lite 3.0 further processes predicted disordered regions in order to achieve a finer classification in sub-regions, using a sliding window of 9 residues to assign each amino acid in the sequence of disordered regions to one of six classes based on conditions. By default each residue is only assigned to a single class by priority to simplify interpretation by non-experts and highlight the most relevant sub-regions. Overlapping sub-region assignments can be enabled by the user. Classes, in order of priority are: Polyampholyte, Positive Polyelectrolyte, Negative Polyelectrolyte, Cysteine-rich, Proline-rich, Glycine-rich, Low complexity, Polar. The first three classes reflect a classification proposed in (Das and Pappu, 2013) and were suggested to be associated with different structural and potentially functional characteristics. The latter four are assigned when the fraction of cysteines, prolines, glycines and polar residues in the sliding window is greater than 32%. Finally, Low-complexity is predicted by SEG (Wootton and Federhen, 1993). Both the sliding window size and threshold on the residue fraction in the sliding window were manually set based on a sample of biologically relevant proteins. Classification (including SEG prediction) is smoothed out in an iterative process following the same approach applied to disordered regions (Necci et al., 2017). Finally, a sub-region is reported for at least 9 residues, otherwise discarded. Both the MobiDB-lite source code and a docker container are available in the GitHub repository:

<https://github.com/BioComputingUP/MobiDB-lite>.

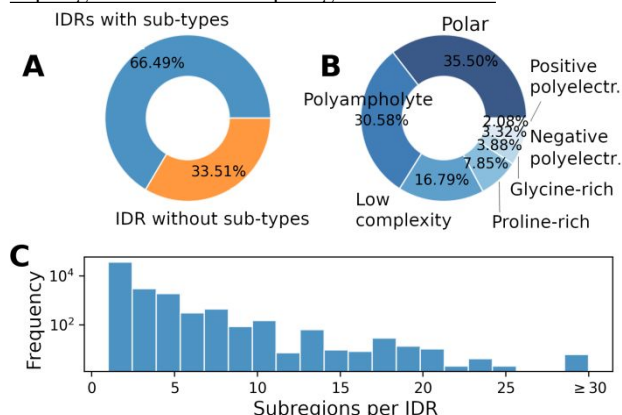


Fig. 1. Abundance of IDRs and sub-regions. MobiDB-lite results for the human proteome. (A) Fraction of IDRs with (blue) and without sub-types (orange). (B) Distribution of IDR sub-types. (C) Distribution of sub-regions detected per IDR, plotted on a logarithmic frequency scale (y-axis).

3 Use case

IDRs and sub-regions were calculated with MobiDB-lite 3.0 for the whole human proteome from UniProt (UP000005640), consisting of 74,043 amino acid sequences, of which 33,322 (44.5%) are predicted IDPs. A total of 64,484 IDRs and 71,921 sub-regions were detected. The majority of IDPs (59.9%) have just one IDR while only 6 proteins have more than 30 IDRs. Mucin-16 (14,451 residues; UniProt ID: Q8WXI7) contains 112 IDRs. Of the 64,484 IDRs detected, 21,610 (33.5%) do not have any sub-regions, while the remaining 66.5% can have 1 or more (Figure 1A). More than 66.1% of sub-regions are either

Polyampholytes or Polar (Figure 1B). The remaining sub-regions are, in order of abundance, low-complexity (16.8%), proline-rich (7.9%), glycine-rich (3.9%), negative (3.3%) and positive polyelectrolytes (2.1%). Cysteine-rich sub-regions are never detected in this dataset. Many IDRs (67.6%) have just one sub-region and the number of IDRs with more than one sub-region drops exponentially with the increase of sub-regions (Figure 1C). In only six cases an IDR has more than 30 sub-regions. Filaggrin (4,061 residues; UniProt ID: P20930) has a predicted IDR spanning from residue 255 to 3,971 hosting 105 sub-regions.

Conclusions

We have described MobiDB-lite 3.0, an improved stand-alone version achieving a finer ID classification by detecting sub-regions in predicted IDRs also available from MobiDB and InterPro. To the best of our knowledge, MobiDB-lite 3.0 is currently the only ID predictor able to sub-classify disorder and also the first ID predictor provided as a docker container.

Funding

This project has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 778247, the Italian Ministry of University and Research PRIN 2017 grant 2017483NH8 and ELIXIR, the European infrastructure for biological data.

Conflict of Interest: none declared.

References

- Das,R.K. and Pappu,R.V. (2013) Conformations of intrinsically disordered proteins are influenced by linear sequence distributions of oppositely charged residues. *Proc. Natl. Acad. Sci. U. S. A.*, **110**, 13392–13397.
- Dosztányi,Z. et al. (2005) IUPred: web server for the prediction of intrinsically unstructured regions of proteins based on estimated energy content. *Bioinformatics*, **21**, 3433–3434.
- Dunker,A.K. et al. (2001) Intrinsically disordered protein. *J. Mol. Graph. Model.*, **19**, 26–59.
- Holehouse,A.S. et al. (2017) CIDER: Resources to Analyze Sequence-Ensemble Relationships of Intrinsically Disordered Proteins. *Biophys. J.*, **112**, 16–21.
- Lee,D. et al. (2005) Identification and distribution of protein families in 120 completed genomes using Gene3D. *Proteins*, **59**, 603–615.
- Linding,R., Russell,R.B., et al. (2003) GlobPlot: Exploring protein sequences for globularity and disorder. *Nucleic Acids Res.*, **31**, 3701–3708.
- Linding,R., Jensen,L.J., et al. (2003) Protein disorder prediction: implications for structural proteomics. *Structure*, **11**, 1453–1459.
- Mészáros,B. et al. (2018) IUPred2A: context-dependent prediction of protein disorder as a function of redox state and protein binding. *Nucleic Acids Res.*, **46**, W329–W337.
- Mier,P. et al. (2020) Disentangling the complexity of low complexity proteins. *Brief. Bioinform.*, **21**, 458–472.
- Mistry,J. et al. (2013) The challenge of increasing Pfam coverage of the human proteome. *Database*, **2013**, bat023–bat023.
- Mitchell,A.L. et al. (2019) InterPro in 2019: improving coverage, classification and access to protein sequence annotations. *Nucleic Acids Res.*, **47**, D351–D360.

MobiDB-lite 3.0

- Necci, M. *et al.* (2018) A comprehensive assessment of long intrinsic protein disorder from the DisProt database. *Bioinformatics*, **34**, 445–452.
- Necci, M. *et al.* (2016) Large-scale analysis of intrinsic disorder flavors and associated functions in the protein sequence universe. *Protein Sci. Publ. Protein Soc.*, **25**, 2164–2174.
- Necci, M. *et al.* (2017) MobiDB-lite: fast and highly specific consensus prediction of intrinsic disorder in proteins. *Bioinformatics*, **33**, 1402–1404.
- Peng, K. *et al.* (2006) Length-dependent prediction of protein intrinsic disorder. *BMC Bioinformatics*, **7**, 208.
- Perdigão, N. *et al.* (2015) Unexpected features of the dark proteome. *Proc. Natl. Acad. Sci.*, 201508380.
- Piovesan, D. *et al.* (2018) MobiDB 3.0: more annotations for intrinsic disorder, conformational diversity and interactions in proteins. *Nucleic Acids Res.*, **46**, D471–D476.
- Sonnhammer, E.L. *et al.* (1997) Pfam: a comprehensive database of protein domain families based on seed alignments. *Proteins*, **28**, 405–420.
- Walsh, I. *et al.* (2015) Comprehensive large-scale assessment of intrinsic protein disorder. *Bioinformatics*, **31**, 201–208.
- Walsh, I. *et al.* (2012) ESpritz: accurate and fast prediction of protein disorder. *Bioinformatics*, **28**, 503–509.
- Wilkinson, M.D. *et al.* (2016) The FAIR Guiding Principles for scientific data management and stewardship. *Sci. Data*, **3**, 160018.
- Wootton, J.C. and Federhen, S. (1993) Statistics of local complexity in amino acid sequences and sequence databases. *Comput. Chem.*, **17**, 149–163.