



HAL
open science

Merging SpecOMS and X!Tandem identification results using machine learning algorithm in i2MassChroQ

Olivier Langella, Thierry Balliau, Marlène Davanture, Filippo Rusconi, Melisande Blein-Nicolas, Dominique Tessier

► To cite this version:

Olivier Langella, Thierry Balliau, Marlène Davanture, Filippo Rusconi, Melisande Blein-Nicolas, et al.. Merging SpecOMS and X!Tandem identification results using machine learning algorithm in i2MassChroQ. ProteoAix 2023 the 3rd Joint Meeting of Spanish, French, and Portuguese Proteomics Societies, Jun 2023, Aix en Provence, France. hal-04103233

HAL Id: hal-04103233

<https://hal.science/hal-04103233>

Submitted on 19 Jun 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - ShareAlike 4.0 International License

Merging SpecOMS and X!Tandem identification results using machine learning algorithm in i2MassChroQ

Olivier Langella*¹, Balliau Thierry¹, Marlène Davanture¹, Filippo Rusconi¹, Mélisande Blein-Nicolas¹, Dominique Tessier²

olivier.langella@universite-paris-saclay.fr



1 - IDEEV - Institut diversité, écologie et évolution du vivant
GQE-Le Moulon, PAPPSO
Université Paris-Saclay, INRAE, CNRS, AgroParisTech,
F-91190 Gif-sur-Yvette, France



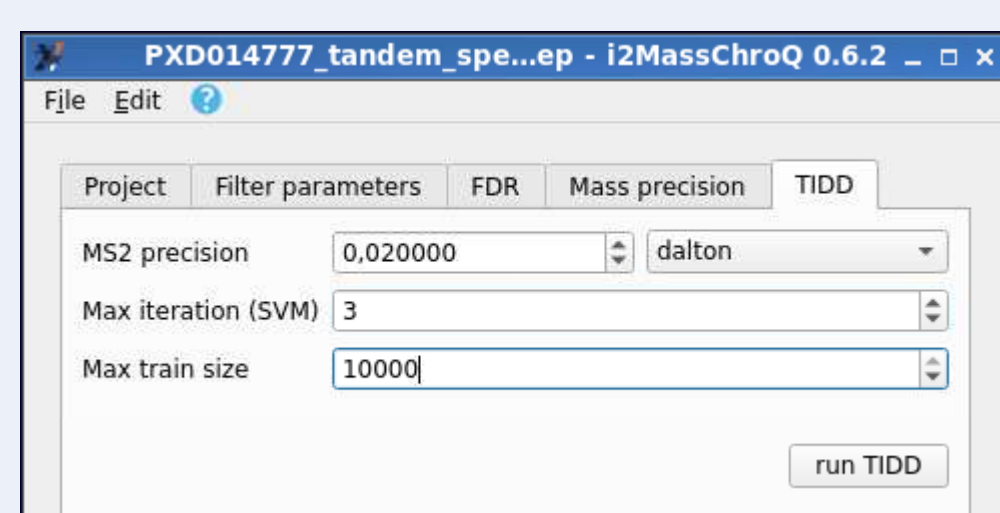
2 - INRAE, UR1268 BIA, F-44316 Nantes, France
INRAE, PROBE Research Infrastructure, BIBS Facility,
F-44316 Nantes, France



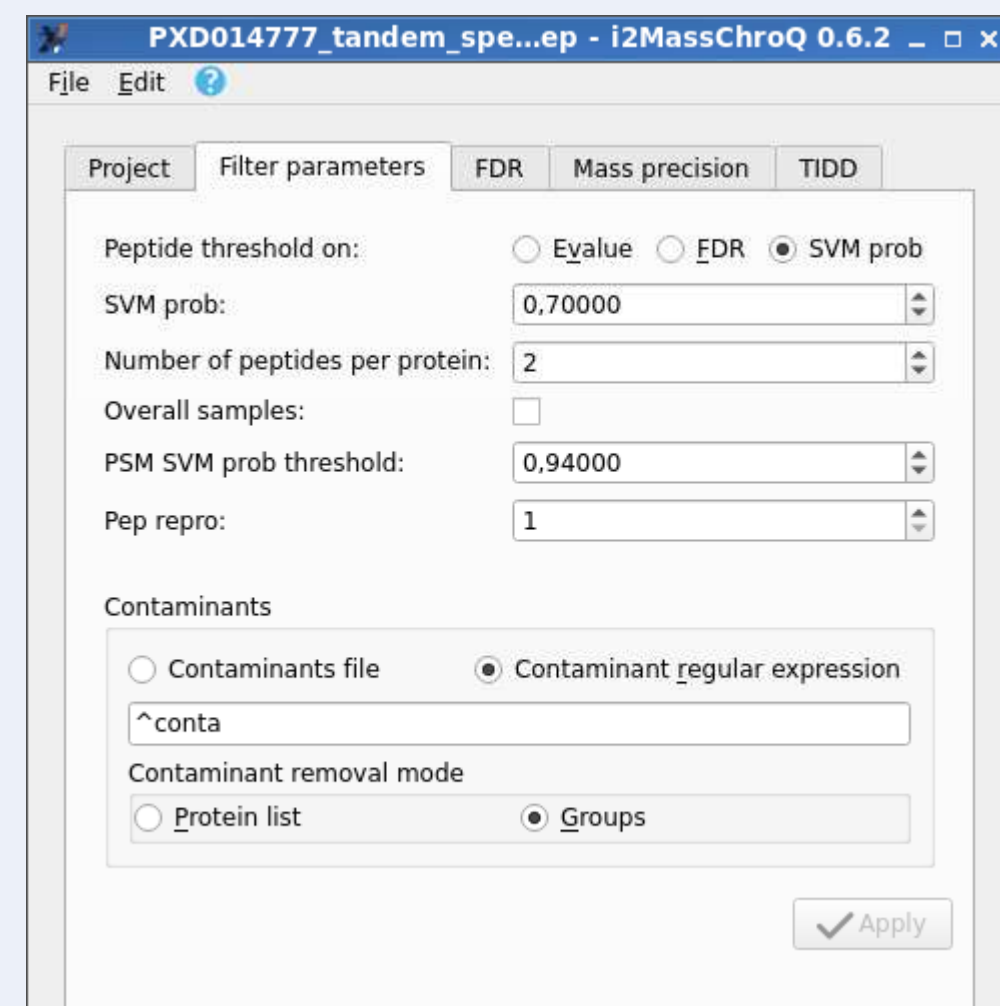
ABSTRACT

SpecOMS is an Open Modification Search method published originally in David et al. 2017. SpecOMS compares experimental spectra generated by a discovery proteomics experiment to a whole set of theoretical spectra deduced from a protein database in a few minutes on a standard workstation. The procedure yields identification results that might comprise peptides bearing unknown modifications. X!Tandem is a highly regarded classical identification engine, very efficient but unable to find peptide spectrum matches (PSMs) if the precursor ion's mass does not corresponds to the theoretical peptide's mass. Combining the results obtained by running both SpecOMS and X!Tandem reduces the proportion of undecided PSMs compared to a classical approach. However this combination approach raises a number of questions because the two engines do compute PSM scores very differently. We address this problem using a machine learning algorithm in i2MassChroQ.

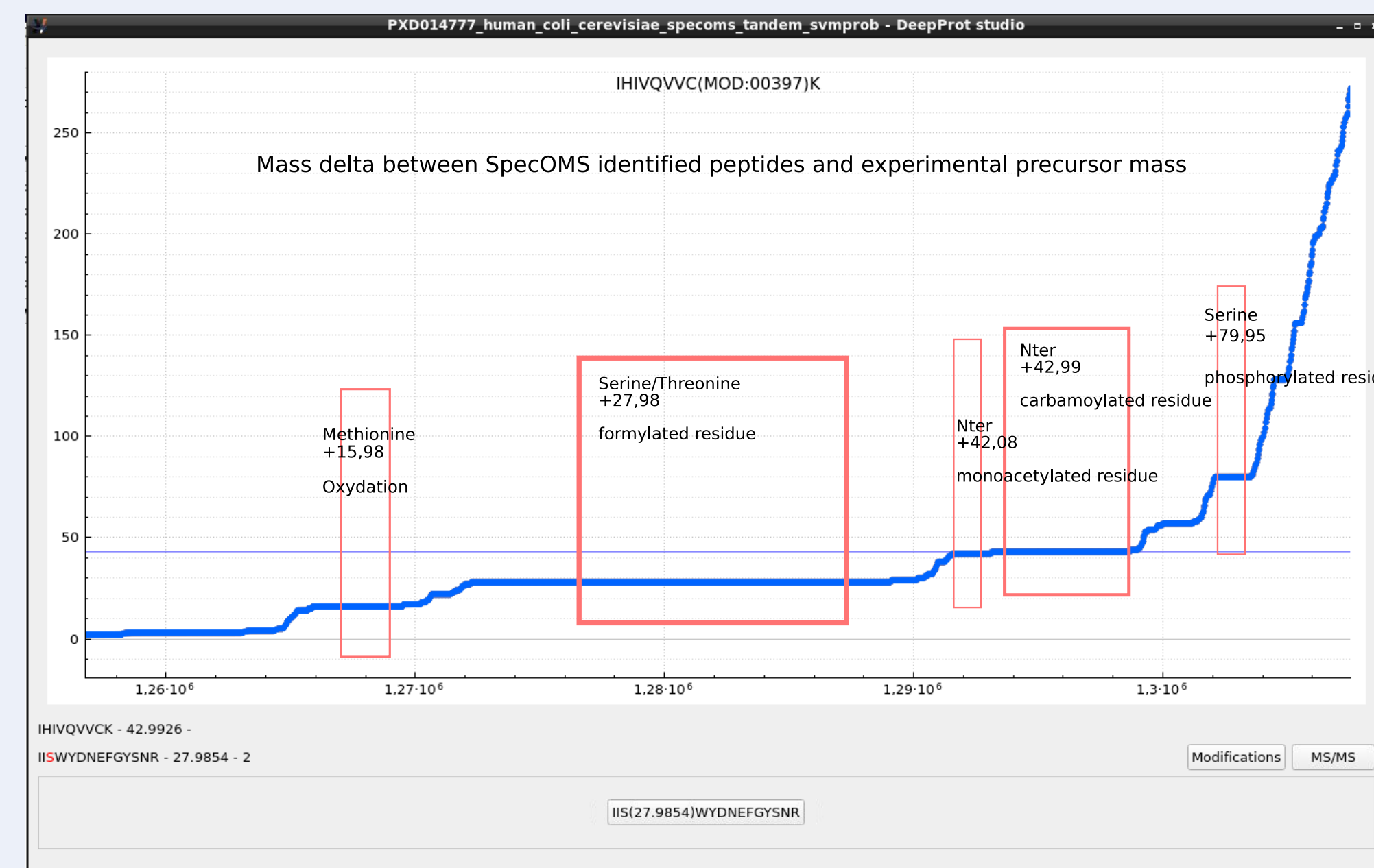
Rescoring PSMs with TIDD (Li et al. 2022)



- 1) Identification results are loaded by i2MassChroQ in a single project
- 2) i2MassChroQ rescors each PSMs with a machine learning procedure called TIDD
- 3) The new score is used to filter PSMs and proteins



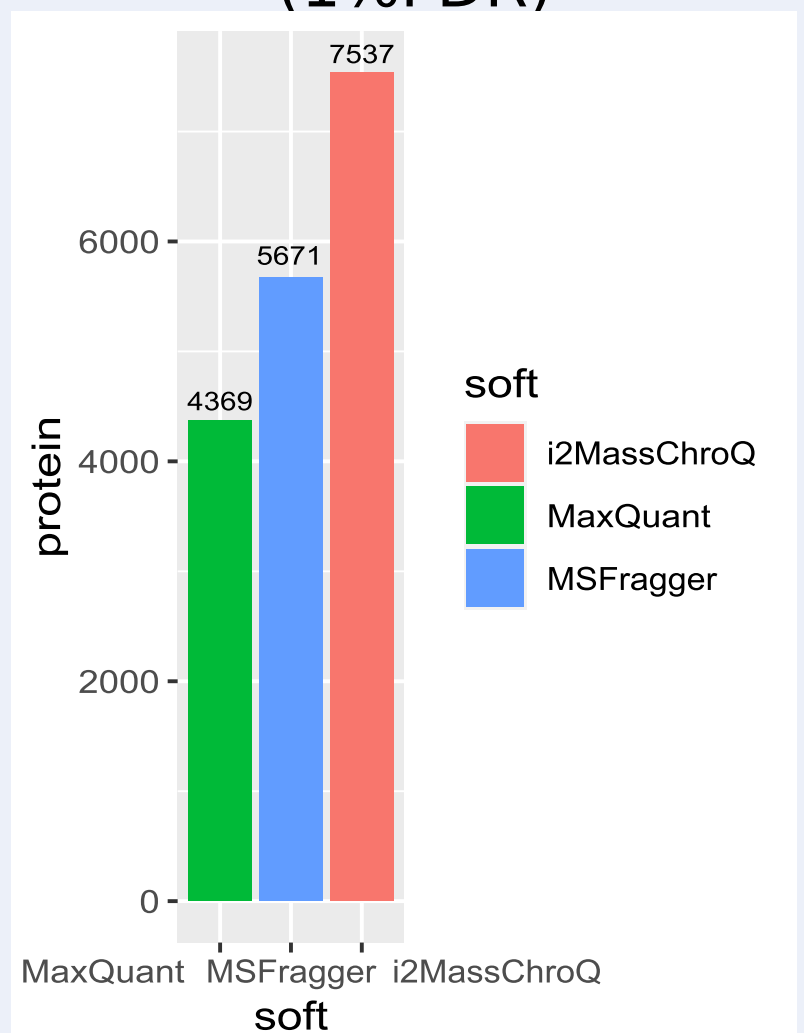
Peptide modifications browser



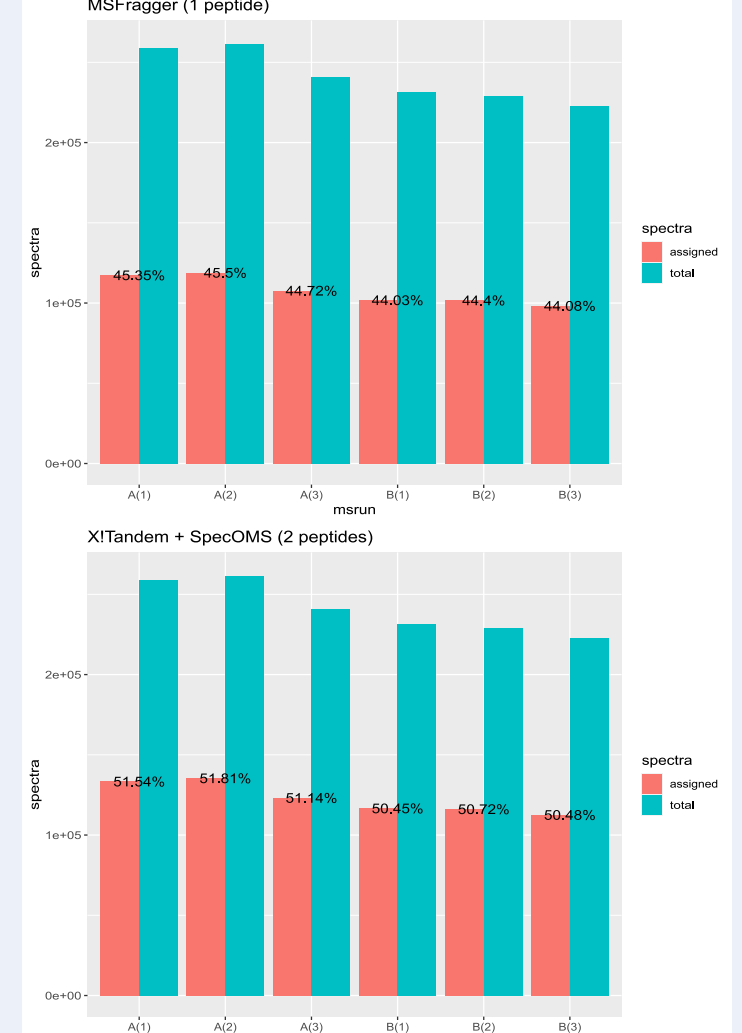
i2MassChroQ offers a dedicated user interface to visualize PSMs bearing putative unknown modifications.

Benchmarking software with timsTOF dataset PXD014777

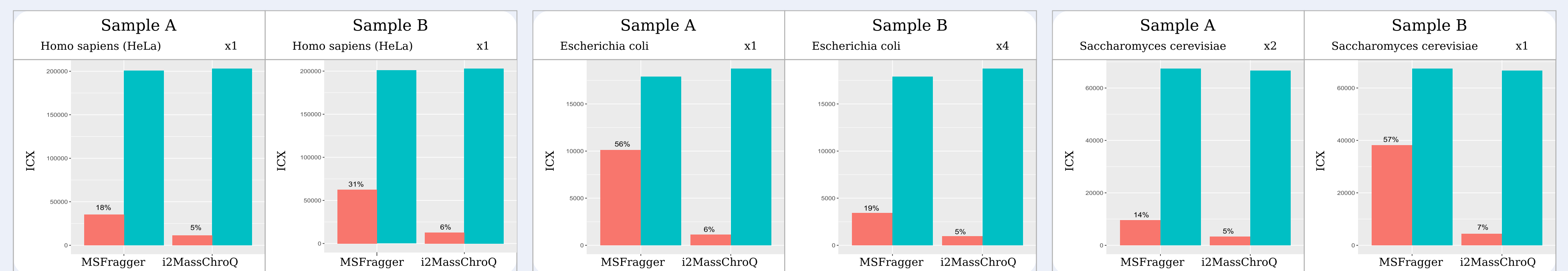
More proteins quantified in both samples (1%FDR)



More spectra assigned per MSrun (1%FDR)

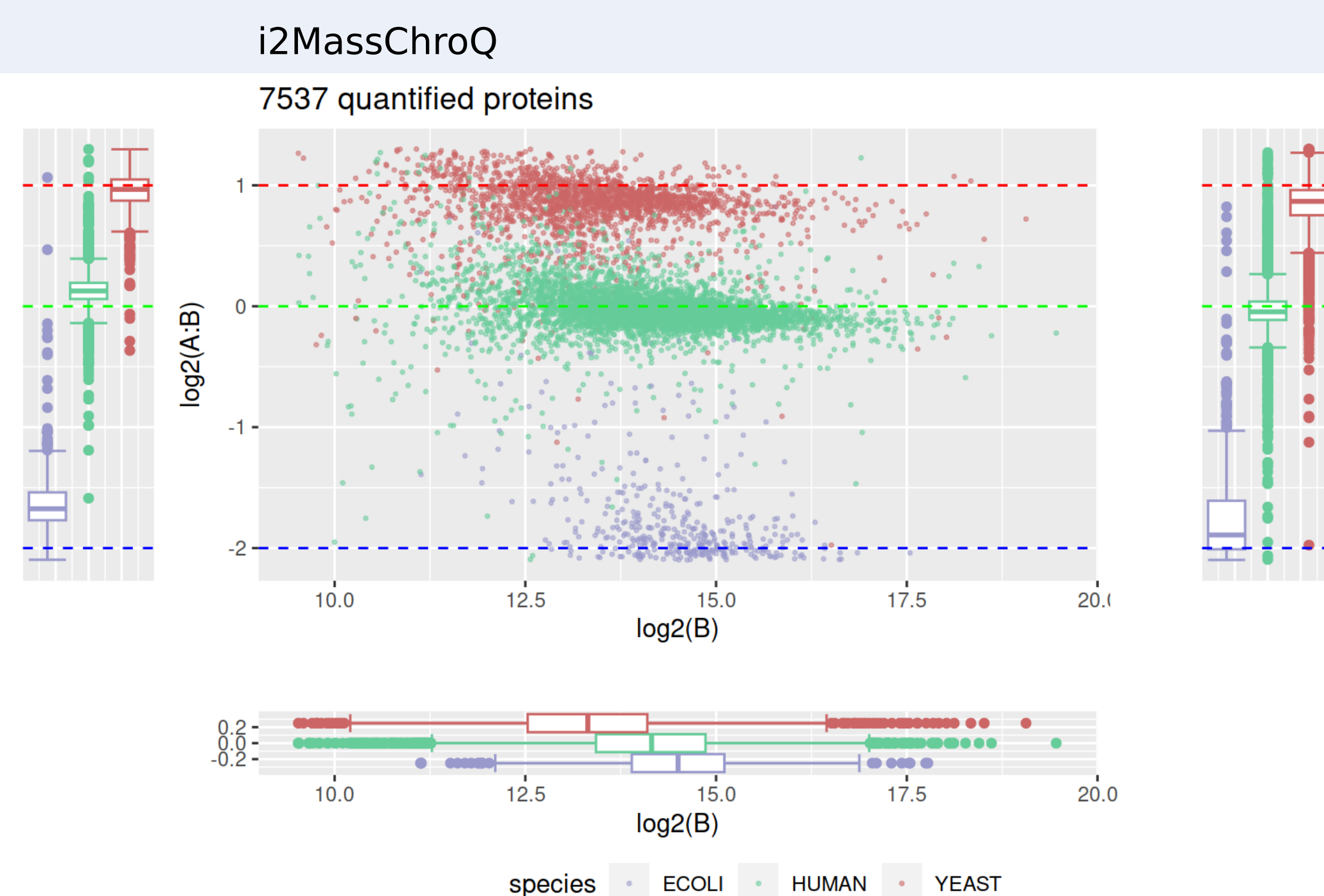
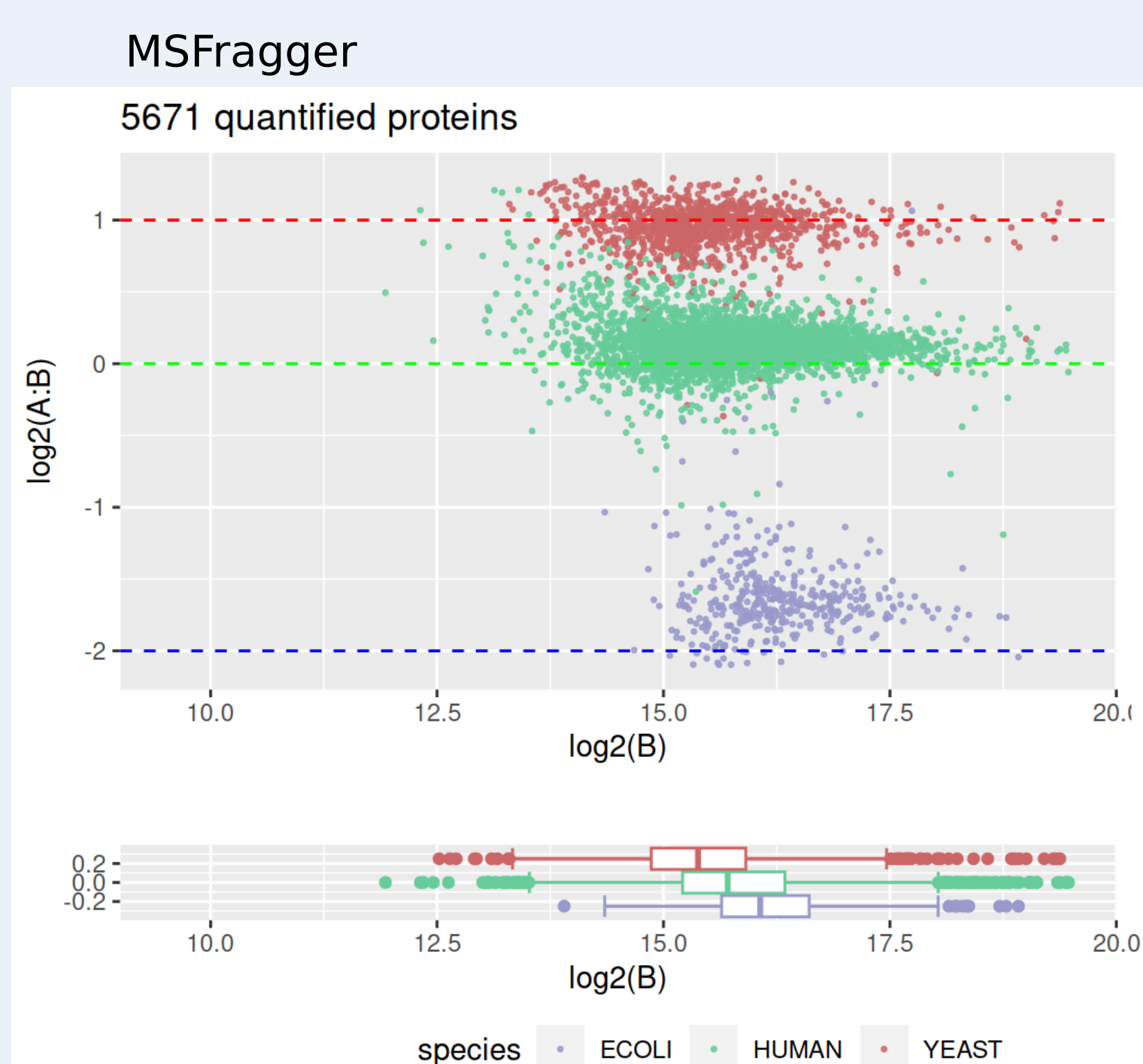


Data completeness : MSFragger 81%, i2MassChroQ 95%

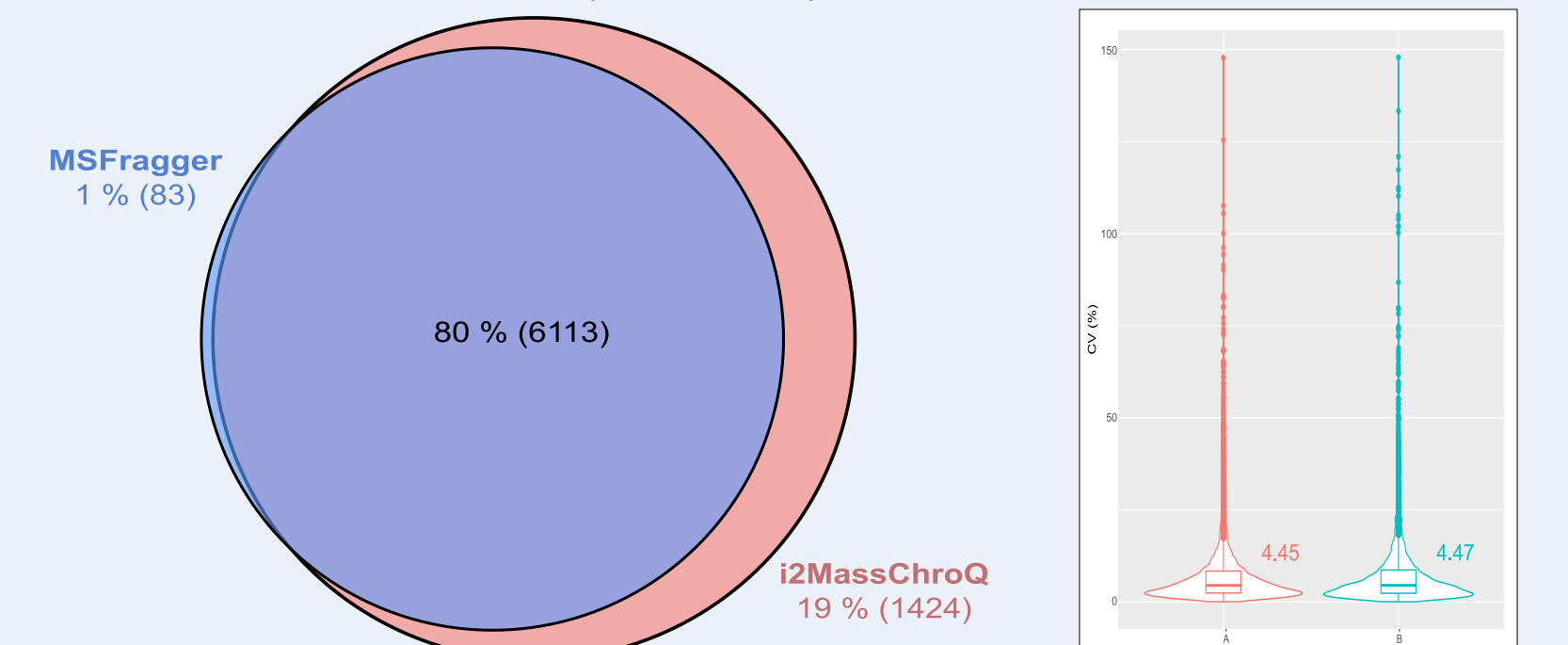


Two samples, A and B containing a mix of protein extracts from three different organisms. 3 replicates each. Samples A and B are 1:1 (*H. sapiens*), 2:1 (*S. cerevisiae*), 1:4 (*E. coli*)

Quantification accuracy



Protein quantified (FDR < 1%) i2MassChroQ CVs distributions (3 replicates)



i2MassChroQ quantifies more proteins on a wider dynamic range. SpecOMS and X!Tandem results merged with TIDD yield more peptides per protein (safer identifications, 50% or more spectra assigned per MSrun). The MassChroQ match-between-run algorithm, specifically designed for timsTOF data, lowers dramatically the proportion of missing values, especially in unfavourable contrasted conditions.

AVAILABILITY

i2MassChroQ source code and binary packages for Windows and Linux are available on the PAPPSO website (<http://pappso.inrae.fr/bioinfo/i2masschroq/>)
SpecOMS C++ identification engine source code or binary packages for Windows and Linux are available on demand to the PAPPSO facility or on the ForgeMIA project (<https://forgemia.inra.fr/pappso/specoms>).
The source code is licensed under the GNU General Public Licence v3+.

* presenting and corresponding author

this work is funded by ANR project ANR-18-CE45-0004

