

ThermoRawFileParser: Modular, Scalable, and Cross-Platform RAW File Conversion

Niels Hulstaert,^{†,‡} Jim Shofstahl,[§] Timo Sachsenberg,^{||} Mathias Walzer,[⊥] Harald Barsnes,^{#,¶} Lennart Martens,^{*,†,‡} and Yasset Perez-Riverol^{*,⊥}

[†]VIB-UGent Center for Medical Biotechnology, VIB, Ghent B-9000, Belgium

[‡]Department of Biomolecular Medicine, Ghent University, Ghent B-9000, Belgium

[§]Thermo Fisher Scientific, 355 River Oaks Parkway, San Jose, California 95134, United States

^{||}Applied Bioinformatics, Department for Computer Science, University of Tuebingen, Sand 14, 72076 Tuebingen, Germany

[⊥]European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI), Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SD, United Kingdom

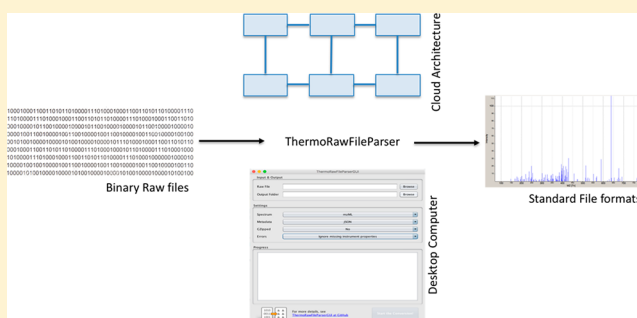
[#]Computational Biology Unit (CBU), Department of Informatics, University of Bergen, Bergen 5020, Norway

[¶]Proteomics Unit (PROBE), Department of Biomedicine, University of Bergen, Bergen 5020, Norway

Supporting Information

ABSTRACT: The field of computational proteomics is approaching the big data age, driven both by a continuous growth in the number of samples analyzed per experiment as well as by the growing amount of data obtained in each analytical run. In order to process these large amounts of data, it is increasingly necessary to use elastic compute resources such as Linux-based cluster environments and cloud infrastructures. Unfortunately, the vast majority of cross-platform proteomics tools are not able to operate directly on the proprietary formats generated by the diverse mass spectrometers. Here, we present ThermoRawFileParser, an open-source, cross-platform tool that converts Thermo RAW files into open file formats such as MGF and the HUPO-PSI standard file format mzML. To ensure the broadest possible availability and to increase integration capabilities with popular workflow systems such as Galaxy or Nextflow, we have also built Conda package and BioContainers container around ThermoRawFileParser. In addition, we implemented a user-friendly interface (ThermoRawFileParserGUI) for those users not familiar with command-line tools. Finally, we performed a benchmark of ThermoRawFileParser and msconvert to verify that the converted mzML files contain reliable quantitative results.

KEYWORDS: bioinformatics, file formats, open source, cloud, mass spectrometry, software, big data, workflows, mzML, metadata



INTRODUCTION

The field of computational proteomics is approaching the big data age,¹ driven both by a continuous growth in the number of samples analyzed per experiment as well as by the growing amount of data obtained in each analytical run. At the same time, more data are now publicly available in proteomics repositories, which in turn means that there is increasing benefit to be had from the reanalysis of millions of mass spectra^{2–5} to find new biological insights (e.g., novel variants and post-translational modifications⁵). However, in order to process these large amounts of (public) data, it is increasingly necessary to use elastic compute resources such as Linux-based cluster environments and cloud infrastructures.⁶

The development of computational proteomics tools has historically favored Microsoft Windows operating systems with tools such as ProteomeDiscover, MaxQuant,⁷ PeaksDB, and Mascot Distiller.⁸ An important driver for this bias has been

the lack of cross-platform libraries to access instrument output data files (RAW files) from major instrument providers.⁹ Several approaches have been devised to overcome this challenge, including the use of dedicated Windows machines in workflows¹⁰ for conversion to RAW data to standard file formats such as mzML,¹¹ the encapsulation of Windows tools such as ReAdW¹² and msconvert¹³ into WineHQ (http://tools.proteomecenter.org/wiki/index.php?title=Msconvert_Wine) to make these tools Linux-compatible, and even the creation of reverse-engineered RAW file readers.¹⁴

An important breakthrough was achieved in 2016 when Thermo Scientific released the first cross-platform application programming interface (API) that enables access to Thermo RAW files from all their instruments on all commonly used

Received: May 21, 2019

Published: November 22, 2019

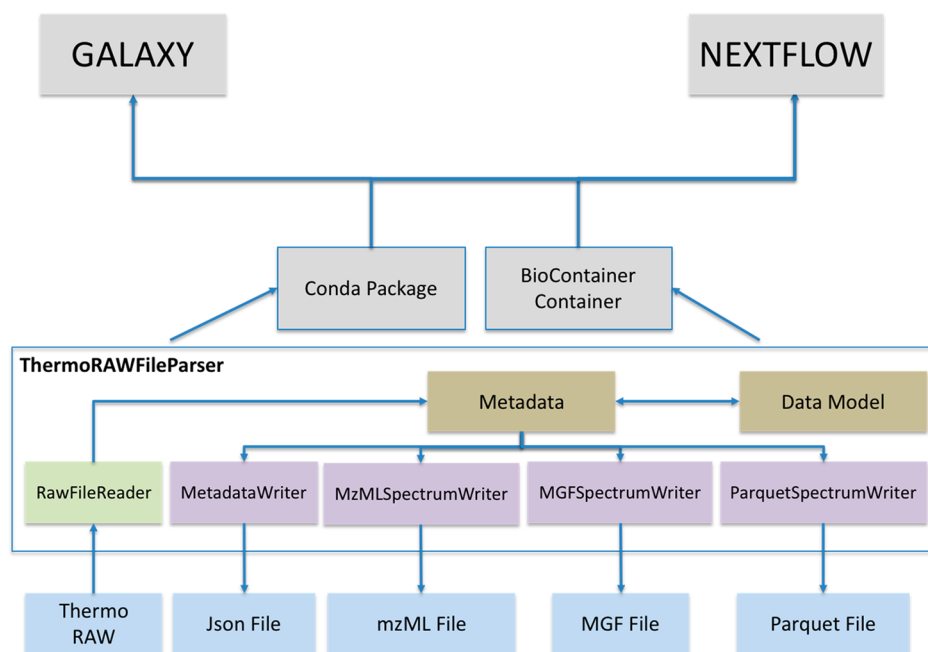


Figure 1. Modular design of ThermoRawFileParser includes exporters to MGF, mzML, Parquet, and Json Metadata. A Conda package and corresponding BioContainer is available for reuse in workflow engines such as Nextflow and Galaxy.

operating systems (e.g., Linux/Unix, Mac OS, or Microsoft Windows). Importantly, this provides the enticing possibility to move proteomics into Linux/UNIX environments, including scalable clusters and cloud environments. This library has already led to a new version of the popular MaxQuant framework that is compatible with Linux/UNIX environments,¹⁵ and it has also been incorporated into the cross-platform, cluster-oriented quantification tool moFF.¹⁶

Whereas the Thermo cross-platform library thus enables specially developed software to access Thermo Raw files on diverse operating systems, most open-source computational proteomics workflows (e.g., OpenMS,¹⁷ Galaxy-P,¹⁸ and the Trans-Proteomics pipeline (TPP)¹⁹) are based on generic, open data formats such as Mascot Generic File (MGF) or mzML. In order to allow these tools to benefit maximally from the cross-platform access to Thermo Raw files, we here present ThermoRawFileParser, an open-source, cross-platform tool that converts Thermo RAW files into open file formats such as MGF and mzML similar to other tools such as msconvert¹³ and RawTools.²⁰ To ensure the broadest possible availability and to increase integration capabilities with popular workflow systems such as Galaxy²¹ or Nextflow,²² we have also built a Conda package²³ and a BioContainers²⁴ container around ThermoRawFileParser. Finally, we performed a benchmark of ThermoRawFileParser and msconvert to verify that the converted mzML files contain reliable quantitative results. A preprint version of the manuscript has been deposited in bioRxiv.²⁵

MATERIALS

Tool Design and Integration

ThermoRawFileParser (<https://github.com/compomics/ThermoRawFileParser>) has been implemented following a modular design (Figure 1). Every file-specific exporter is implemented as an independent module, which enables easy extension to include more exporters in the future. Currently,

the tool can export to MGF (MGFSpectrumWriter), mzML (MzMLSpectrumWriter), and JSON (for the metadata only) (MetadataWriter). This modular design has already enabled the community to extend the library for other novel file formats such as Parquet (ParquetSpectrumWriter), which is designed for distributed big data processing clusters of Hadoop or Spark. The JSON export of ThermoRawFileParser can optionally be used to only extract various metadata elements (including instrument settings and scan settings; see <https://github.com/PRIDE-Archive/pride-metadata-standard>) (Figure 2). This specific feature is currently used by the PRIDE Database to reannotate thousands of RAW files with the correct instrument metadata. For peak picking, data centroiding, and noise removal, ThermoRawFileParser relies on the native methods provided by the Thermo API.

A key feature of any open-source tool is its ability to integrate with other frameworks.²⁶ We have therefore created a BioConda recipe for ThermoRawFileParser (<https://github.com/bioconda/bioconda-recipes/tree/master/recipes/thermorawfileparser>), which can be used to automatically build a Docker Container. This Docker is pushed to the BioContainer project, which in turn enables easy reuse of the tool by both the Galaxy and the Nextflow environments. As an illustration of such integration, we have developed a Nextflow workflow for the proteomics community, which converts an entire ProteomeXchange project using the ThermoRawFileParser container (<https://github.com/bigbio/nf-workflows/tree/master/thermo-convert-nf>).

In addition to the command-line tool, we have implemented a graphical user interface that makes the use of ThermoRawFileParser easier and highly intuitive, enabling the user to perform conversions of RAW files (Figure 3). The GUI includes all main options of ThermoRawFileParser and a report system to report errors during the conversion. ThermoRawFileParserGUI (Figure 3) is an open-source Java program, available in a cross-platform package that incorporates ThermoRawFileParser executables for the main operating



Figure 2. JSON representation for one file run in PXD006336 data set. The JSON file contains metadata information about the file (FileProperties), the instrument (InstrumentProperties), mass spectrometry information (MSData), sample (SampleData), and scan settings (ScanSettings).

systems. It can be downloaded from <https://github.com/componics/ThermoRawFileParserGUI>.

Benchmark Data Sets

Three different public Thermo data sets were used to compare the conversion from RAW files into mzML with the ProteoWizard msconvert tool and the ThermoRawFileParser: PXD006336 (Orbitrap Q-Exactive), PXD014346 (Orbitrap Fusion Lumos), and PXD001502 (Orbitrap Velos). We used a Nextflow workflow and the identification-free OpenMS quality control tools to benchmark different metrics such as number of spectra MS1/MS2 or number of peaks by spectrum (https://github.com/bigbio/nf-workflows/tree/master/qc-idfree_from_raw).

We used the IPRG2015 data set (<https://www.ebi.ac.uk/pride/archive/projects/PXD006336>)²⁷ to benchmark the quality of the mzML files produced by ThermoRawFileParser.

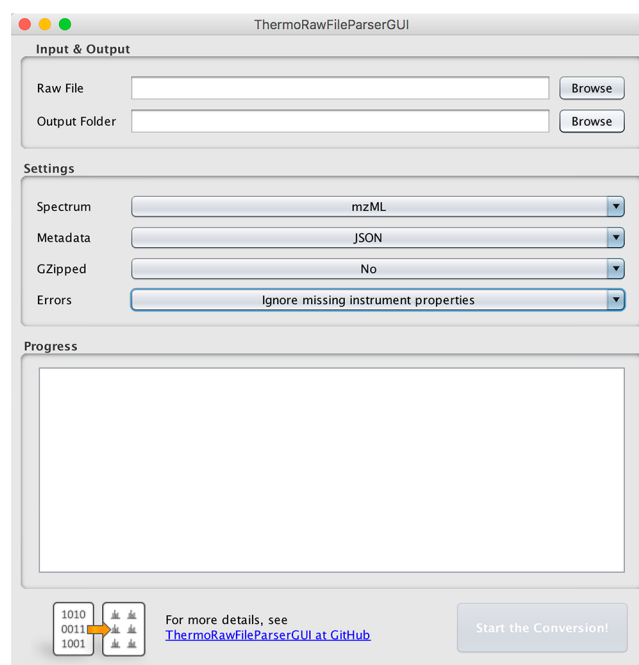


Figure 3. ThermoRawFileParserGUI provides a user-friendly user interface to convert Thermo Raw files to mzML, mgf, and metadata formats.

This data set is based on four artificially constructed samples of known composition, each containing a constant background of 200 ng of tryptic digests of *Saccharomyces cerevisiae* (ATCC strain 204508/S288c). Each sample was separately spiked in with different quantities of six individual protein digests. Samples were analyzed in three LC–MS/MS using a Thermo Scientific Q-Exactive mass spectrometer (12 runs). Both MS and MS/MS data were acquired in profile mode in the Orbitrap, with a resolution of 70 000 for MS and 17 500 for MS/MS. The MS1 scan range was 300–1650 *m/z*, the normalized collision energy was set to 27%, and singly charged ions were excluded.

Quantification Workflow

To perform the quantification benchmark using the PXD006336 data set, we built a workflow using OpenMS in which raw files were converted from Thermo Scientific RAW files to mzML using ThermoRawFileParser tool. The resulting spectra were searched using MS-GF+ (v2018.01.30),²⁸ executed via the OpenMS search engine wrapper MSGFPlusAdapter, allowing 10 ppm precursor mass tolerance and setting carbamidomethylation of cysteine as fixed and methionine oxidation as variable modification. PSMs were filtered (*q* value <5%) and used for feature detection using the semitargeted approach implemented in the OpenMS tool FeatureFinderIdentification.²⁹ Prior to identification, nonlinear retention time alignment was performed using the MapAlignerIdentification, and the identified proteins were then quantified using unique peptides only. The workflow for comparison was developed using Nextflow and BioContainers to ensure the reproducibility of the present results (<https://github.com/bigbio/nf-workflows/tree/master/benchmark-converter-nf>).

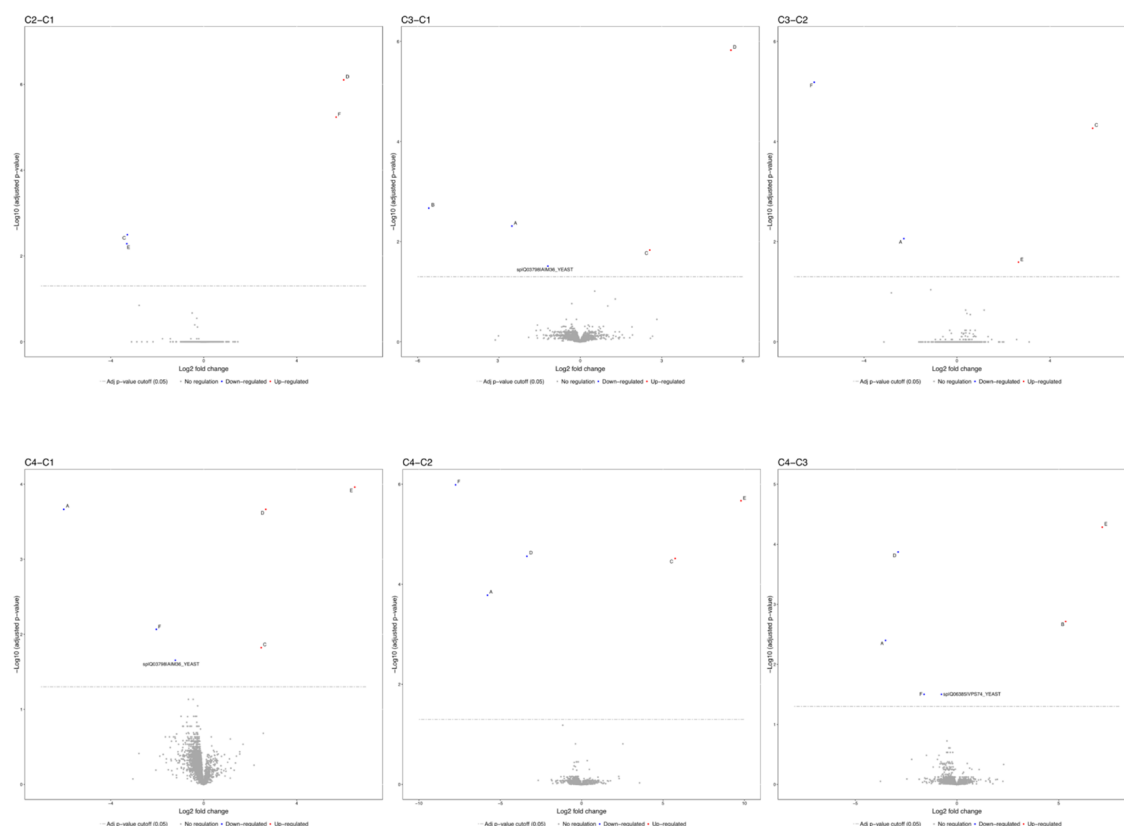


Figure 4. Volcano plot display of the results of the statistical analysis with OpenMS and MSstats of LC–MS converted using ThermoRawFileParser. y-axis: $-\log_{10} p$ value of a pairwise comparison between two samples, adjusted to control the false discovery rate in the list of differentially abundant proteins in this comparison. x-axis: log 2-fold change between two samples.

RESULTS AND DISCUSSION

We compare msconvert and ThermoRawFileParser conversion to mzML using four different metrics: number of MS1, number of MS2, MS1 peak count distribution, MS2 peak count distribution, identification map, and the precursor charge distribution. We observed no major differences between both tools (msconvert and ThermoRawFileParser) for the number of MS1/MS2 and the peak count distributions (PXD006336, [Supplementary Note 1](#); PXD014346, [Supplementary Note 2](#); PXD001502, [Supplementary Note 3](#)).

We analyzed the IPRG2015 data set (PXD006336) using OpenMS framework and MSstats.³⁰ We conducted a high-level analysis of the IPRG2015 data set to verify whether the mzML files obtained by the ThermoRawFileParser pipeline could replicate the quantification of the spike-in proteins in the sample using the approach described in the original publication.²⁷ In [Figure 4](#), high values correspond to statistically significant changes. The x-axis is the log 2-fold change between two samples, which in statistical language is sometimes called the “practical significance” of a change. Similar to best workflows reported in Choi et al.,²⁷ the estimates of log 2-fold changes among the spiked proteins (A–F) were close to the true values, whereas most background proteins did not show significant differential expression. We computed the number of false positive (FP = 3), true positive (TP = 27), and positive predictive value (PPV = 0.9) as defined by Choi et al.²⁷ If we compare this PPV (0.9) with the submissions performed in Choi et al. (Figure 3 in Choi et al.²⁷), our results are within the 10 tops out of 47 data analysis pipelines. We performed the same analysis using msconvert to

transform RAW data to mzML ([Supplementary Note 4](#)). The number of peptides and proteins identified with both workflows (msconvert and ThermoRawFileParser) were similar. In addition to msconvert, the recently published RawTools allows converting RAW files into MGF files, but it does not provide support for standard HUPO-PSI file formats such as mzML.

CONCLUSIONS

ThermoRawFileParser is an open-source software tool for the conversion of Thermo Raw files into open formats. Because of the growing need for more scalable and distributed computational proteomics approaches, ThermoRawFileParser has been designed to easily plug into large-scale workflow systems such as Galaxy or OpenMS. The current implementation also provides support for native writing into Amazon web service object stores (S3), making the tool highly portable to cloud architectures. Finally, the modular design of the library, along with its open-source nature, allows other researchers to contribute to and extend ThermoRawFileParser for new file formats in the future.

ASSOCIATED CONTENT

Supporting Information

The Supporting Information is available free of charge at <https://pubs.acs.org/doi/10.1021/acs.jproteome.9b00328>.

Quality control plots for ProteomeXchange data set PXD014346; quality control plots for ProteomeXchange data set PXD001502; benchmark results using the data set PXD006336, and OpenMS workflow ([PDF](#))

AUTHOR INFORMATION

Corresponding Authors

*E-mail: lennart.martens@UGent.be.

*E-mail: yperez@ebi.ac.uk. Phone: +44 1223 492 610. Fax: +44 1223 494 484.

ORCID

Lennart Martens: [0000-0003-4277-658X](https://orcid.org/0000-0003-4277-658X)

Yasset Perez-Riverol: [0000-0001-6579-6941](https://orcid.org/0000-0001-6579-6941)

Notes

The authors declare no competing financial interest.

ACKNOWLEDGMENTS

This work was partially supported by one ELIXIR Implementation Study. ELIXIR-EXCELERATE is funded by the European Commission within the Research Infrastructures programme of Horizon 2020, Grant Agreement No. 676559. L.M. acknowledges support from FWO research project No. G042518, and N.H. acknowledges support from FWO research project 3G0H6916 (ERA-MBT grant agreement number 604814). Y.P.-R. acknowledges the Wellcome Trust (Grant No. 208391/Z/17/Z) and Y.P.-R. and L.M. acknowledge the EPIC-XS project (Grant No. 823839), funded by the Horizon 2020 programme of the European Union. M.W. is supported by BBSRC grant [BB/P024599/1]. H.B. is supported by the Bergen Research Foundation and the Research Council of Norway. T.S. was supported by a grant from the German Federal Ministry of Education and Research (BMBF) under Grant No. 031A535A (German Network for Bioinformatics, de.NBI/CIBI). We thank Meena Choi for the review of our data.

REFERENCES

- (1) Lam, M. P. Y.; Lau, E.; Ng, D. C. M.; Wang, D.; Ping, P. Cardiovascular Proteomics in the Era of Big Data: Experimental and Computational Advances. *Clin. Proteomics* **2016**, *13*, 23.
- (2) Martens, L.; Vizcaino, J. A. A Golden Age for Working with Public Proteomics Data. *Trends Biochem. Sci.* **2017**, *42* (5), 333–341.
- (3) Perez-Riverol, Y.; Alpi, E.; Wang, R.; Hermjakob, H.; Vizcaino, J. A. Making Proteomics Data Accessible and Reusable: Current State of Proteomics Databases and Repositories. *Proteomics* **2015**, *15* (5–6), 930–949.
- (4) Vaudel, M.; Verheggen, K.; Csordas, A.; Raeder, H.; Berven, F. S.; Martens, L.; Vizcaino, J. A.; Barsnes, H. Exploring the Potential of Public Proteomics Data. *Proteomics* **2016**, *16* (2), 214–225.
- (5) Griss, J.; Perez-Riverol, Y.; Lewis, S.; Tabb, D. L.; Dianes, J. A.; Del-Toro, N.; Rurik, M.; Walzer, M. W.; Kohlbacher, O.; Hermjakob, H.; et al. Recognizing Millions of Consistently Unidentified Spectra across Hundreds of Shotgun Proteomics Datasets. *Nat. Methods* **2016**, *13* (8), 651–656.
- (6) Verheggen, K.; Barsnes, H.; Martens, L. Distributed Computing and Data Storage in Proteomics: Many Hands Make Light Work, and a Stronger Memory. *Proteomics* **2014**, *14* (4–5), 367–377.
- (7) Cox, J.; Mann, M. MaxQuant Enables High Peptide Identification Rates, Individualized p.p.b.-Range Mass Accuracies and Proteome-Wide Protein Quantification. *Nat. Biotechnol.* **2008**, *26* (12), 1367–1372.
- (8) Perkins, D. N.; Pappin, D. J.; Creasy, D. M.; Cottrell, J. S. Probability-Based Protein Identification by Searching Sequence Databases Using Mass Spectrometry Data. *Electrophoresis* **1999**, *20* (18), 3551–3567.
- (9) Martens, L.; Nesvizhskii, A. I.; Hermjakob, H.; Adamski, M.; Omenn, G. S.; Vandekerckhove, J.; Gevaert, K. Do We Want Our Data Raw? Including Binary Mass Spectrometry Data in Public Proteomics Data Repositories. *Proteomics* **2005**, *5* (13), 3501–3505.

- (10) Verheggen, K.; Maddelein, D.; Hulstaert, N.; Martens, L.; Barsnes, H.; Vaudel, M. Pladipus Enables Universal Distributed Computing in Proteomics Bioinformatics. *J. Proteome Res.* **2016**, *15* (3), 707–712.
- (11) Martens, L.; Chambers, M.; Sturm, M.; Kessner, D.; Levander, F.; Shofstahl, J.; Tang, W. H.; Römpf, A.; Neumann, S.; Pizarro, A. D.; et al. MzML—a Community Standard for Mass Spectrometry Data. *Mol. Cell. Proteomics* **2011**, *10* (1), R110.000133.
- (12) Pedrioli, P. G. A.; Eng, J. K.; Hubley, R.; Vogelzang, M.; Deutsch, E. W.; Raught, B.; Pratt, B.; Nilsson, E.; Angeletti, R. H.; Apweiler, R.; et al. A Common Open Representation of Mass Spectrometry Data and Its Application to Proteomics Research. *Nat. Biotechnol.* **2004**, *22* (11), 1459–1466.
- (13) Chambers, M. C.; Maclean, B.; Burke, R.; Amodei, D.; Ruderman, D. L.; Neumann, S.; Gatto, L.; Fischer, B.; Pratt, B.; Egerton, J.; et al. A Cross-Platform Toolkit for Mass Spectrometry and Proteomics. *Nat. Biotechnol.* **2012**, *30* (10), 918–920.
- (14) Kelchtermans, P.; Silva, A. S. C.; Argentini, A.; Staes, A.; Vandenbussche, J.; Laukens, K.; Valkenburg, D.; Martens, L. Open-Source, Platform-Independent Library and Online Scripting Environment for Accessing Thermo Scientific RAW Files. *J. Proteome Res.* **2015**, *14* (11), 4940–4943.
- (15) Sinitcyn, P.; Tiwary, S.; Rudolph, J.; Gutenbrunner, P.; Wichmann, C.; Yilmaz, S.; Hamzei, H.; Salinas, F.; Cox, J. MaxQuant Goes Linux. *Nat. Methods* **2018**, *15* (6), 401.
- (16) Argentini, A.; Staes, A.; Grüning, B.; Mehta, S.; Easterly, C.; Griffin, T. J.; Jagtap, P.; Impens, F.; Martens, L. Update on the MoFF Algorithm for Label-Free Quantitative Proteomics. *J. Proteome Res.* **2019**, *18* (2), 728–731.
- (17) Pfeuffer, J.; Sachsenberg, T.; Alka, O.; Walzer, M.; Fillbrunn, A.; Nilse, L.; Schilling, O.; Reinert, K.; Kohlbacher, O. OpenMS - A Platform for Reproducible Analysis of Mass Spectrometry Data. *J. Biotechnol.* **2017**, *261*, 142–148.
- (18) Chambers, M. C.; Jagtap, P. D.; Johnson, J. E.; McGowan, T.; Kumar, P.; Onsongo, G.; Guerrero, C. R.; Barsnes, H.; Vaudel, M.; Martens, L.; et al. An Accessible Proteogenomics Informatics Resource for Cancer Researchers. *Cancer Res.* **2017**, *77* (21), e43–e46.
- (19) Deutsch, E. W.; Mendoza, L.; Shteynberg, D.; Slagel, J.; Sun, Z.; Moritz, R. L. Trans-Proteomic Pipeline, a Standardized Data Processing Pipeline for Large-Scale Reproducible Proteomics Informatics. *Proteomics: Clin. Appl.* **2015**, *9* (7–8), 745–754.
- (20) Kovalchik, K. A.; Colborne, S.; Spencer, S. E.; Sorensen, P. H.; Chen, D. D. Y.; Morin, G. B.; Hughes, C. S. RawTools: Rapid and Dynamic Interrogation of Orbitrap Data Files for Mass Spectrometer System Management. *J. Proteome Res.* **2019**, *18* (2), 700–708.
- (21) Afgan, E.; Baker, D.; Batut, B.; van den Beek, M.; Bouvier, D.; Cech, M.; Chilton, J.; Clements, D.; Coraor, N.; Grüning, B. A.; et al. The Galaxy Platform for Accessible, Reproducible and Collaborative Biomedical Analyses: 2018 Update. *Nucleic Acids Res.* **2018**, *46* (W1), W537–W544.
- (22) Di Tommaso, P.; Chatzou, M.; Floden, E. W.; Barja, P. P.; Palumbo, E.; Notredame, C. Nextflow Enables Reproducible Computational Workflows. *Nat. Biotechnol.* **2017**, *35* (4), 316–319.
- (23) Grüning, B.; Dale, R.; Sjödin, A.; Chapman, B. A.; Rowe, J.; Tomkins-Tinch, C. H.; Valieris, R.; Köster, J. Bioconda Team. Bioconda: Sustainable and Comprehensive Software Distribution for the Life Sciences. *Nat. Methods* **2018**, *15* (7), 475–476.
- (24) da Veiga Leprevost, F.; Grüning, B. A.; Alves Aflitos, S.; Röst, H. L.; Uszkoreit, J.; Barsnes, H.; Vaudel, M.; Moreno, P.; Gatto, L.; Weber, J.; et al. BioContainers: An Open-Source and Community-Driven Framework for Software Standardization. *Bioinformatics* **2017**, *33* (16), 2580–2582.
- (25) Hulstaert, N.; Sachsenberg, T.; Walzer, M.; Barsnes, H.; Martens, L.; Perez-Riverol, Y. ThermoRawFileParser: Modular, Scalable and Cross-Platform RAW File Conversion. *bioRxiv* **2019**, DOI: [10.1101/622852](https://doi.org/10.1101/622852).
- (26) Perez-Riverol, Y.; Wang, R.; Hermjakob, H.; Müller, M.; Vesada, V.; Vizcaino, J. A. Open Source Libraries and Frameworks for

Mass Spectrometry Based Proteomics: A Developer's Perspective. *Biochim. Biophys. Acta, Proteins Proteomics* **2014**, 1844, 63–76.

(27) Choi, M.; Eren-Dogu, Z. F.; Colangelo, C.; Cottrell, J.; Hoopmann, M. R.; Kapp, E. A.; Kim, S.; Lam, H.; Neubert, T. A.; Palmblad, M.; et al. ABRF Proteome Informatics Research Group (IPRG) 2015 Study: Detection of Differentially Abundant Proteins in Label-Free Quantitative LC-MS/MS Experiments. *J. Proteome Res.* **2017**, 16 (2), 945–957.

(28) Kim, S.; Pevzner, P. A. MS-GF+ Makes Progress towards a Universal Database Search Tool for Proteomics. *Nat. Commun.* **2014**, 5, 5277.

(29) Weisser, H.; Choudhary, J. S. Targeted Feature Detection for Data-Dependent Shotgun Proteomics. *J. Proteome Res.* **2017**, 16 (8), 2964–2974.

(30) Choi, M.; Chang, C.-Y.; Clough, T.; Broudy, D.; Killeen, T.; MacLean, B.; Vitek, O. MSstats: An R Package for Statistical Analysis of Quantitative Mass Spectrometry-Based Proteomic Experiments. *Bioinformatics* **2014**, 30 (17), 2524–2526.