

CHIMERYYS: An AI-Driven Leap Forward in Peptide Identification

Martin Frejno¹; Daniel P Zolg¹; Tobias Schmidt¹; Siegfried Gessulat¹; Michael Graber¹; Florian Seefried¹; Magnus Rathke-Kuhnert¹; Samia Ben Fredj¹; Shyamnath Premnadh¹; Kai Fritzsche²; Frank Berg²; Waqas Nasir²; David Horn³; Bernard Delanghe²; Christoph Henrich²; Bernhard Kuster⁴; Mathias Wilhelm⁴ ¹MSAID GmbH, Garching b.München, Germany; ²Thermo Fisher Scientific (Bremen) GmbH, Bremen, Germany; ³Thermo Fisher Scientific, San Jose, CA; ⁴Technical University Munich, Freising, Germany

ABSTRACT

Purpose: Chimeric spectra represent a substantial challenge for bottom-up proteomics data analysis. Here, we describe CHIMERYYS™, a novel, highly scalable, cloud-native, microservice-based and artificial intelligence-powered search algorithm that rethinks the analysis of tandem mass spectra from the ground up and deconvolutes chimeric spectra based on predicted fragment ion intensities.

Methods: We performed comparative analyses of standard HeLa tryptic digests that were acquired on various mass spectrometry platforms using different gradient lengths and isolation widths, as well as *in-silico* generated and publicly available datasets from various organisms using Sequest HT™, the Precursor Detector Node, INFERYYS™ Rescoring [1] and CHIMERYYS™ as implemented in a pre-release version of Thermo Fisher™ Proteome Discoverer™ 3.0 software.

Results: CHIMERYYS doubles peptide identifications in classical data-dependent acquisition (DDA) datasets compared to Sequest HT and increases the number of identified peptides per protein by 2.5-fold on average, which translates to ~2 PSMs per spectrum and an identification rate of >80%. Entrapment analyses suggest that the CHIMERYYS score set is well-calibrated and dilution experiments confirm that peptides unique to CHIMERYYS follow the expected ratio distribution. Experiments based on simulated chimeric spectra establish that CHIMERYYS has a sensitivity of >90%. Using CHIMERYYS enables more efficient data acquisition strategies, as both wider isolation windows and shorter gradients can be used to generate more PSMs in a shorter timeframe.

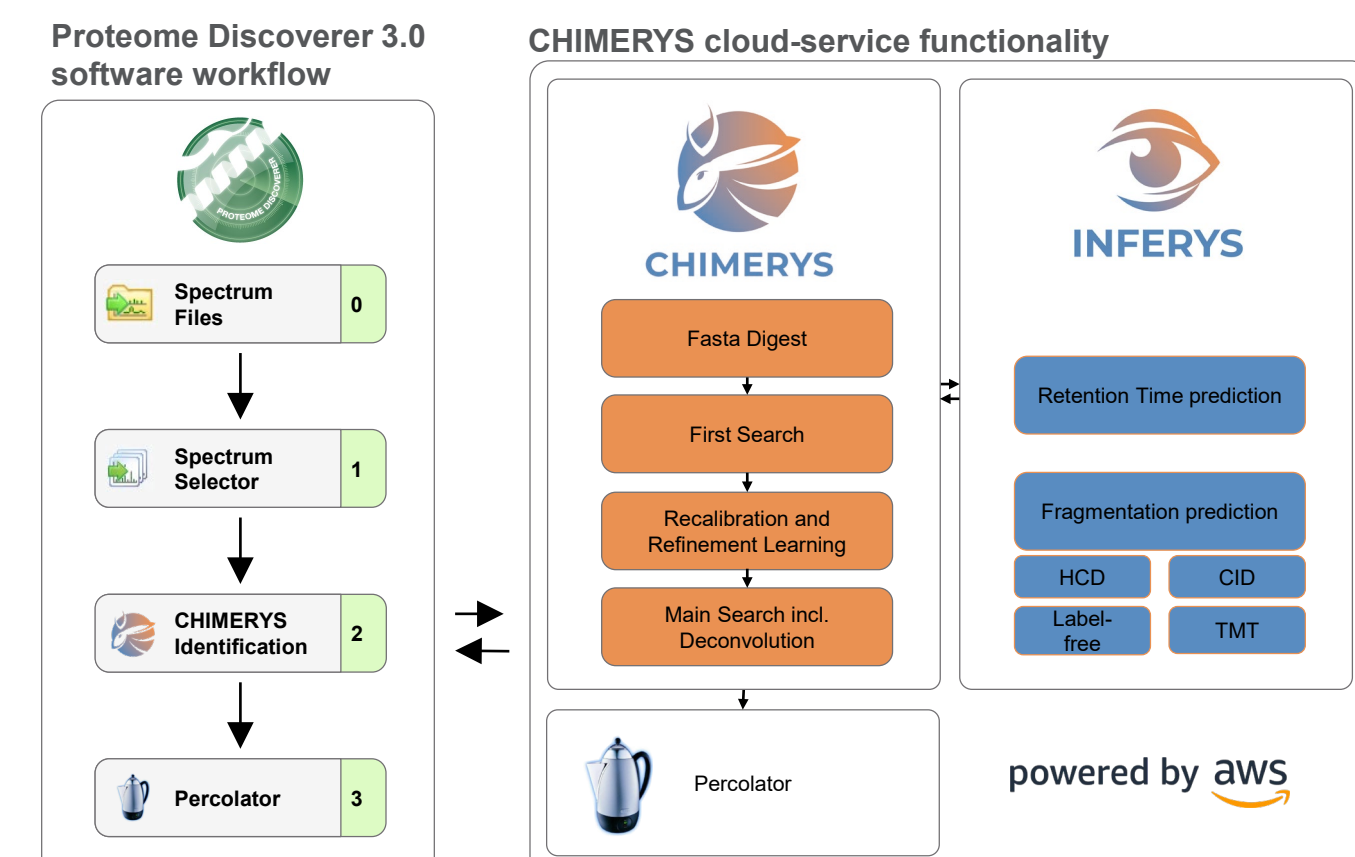
INTRODUCTION

Matching peptide sequences to tandem mass spectra is integral to bottom-up proteomics. Chimeric spectra are estimated to constitute >40% of DDA data [2], violating the assumption that one spectrum represents one peptide. Some search engines allow multi-pass searches or duplicate chimeric spectra for several possible precursors, but few account for the fact that the measured intensities of (isobaric) fragment ions may be the sum of multiple peptides. This introduces errors and leaves valuable information unused, resulting in far fewer peptide identifications than contained in the data. Here, we describe CHIMERYYS, a new AI-based search algorithm that rethinks the analysis of tandem mass spectra from the ground up. It routinely doubles the number of peptide identifications in comparison to classical search algorithms and reaches identification rates of >80%.

MATERIALS AND METHODS

Data Analysis

CHIMERYYS is a cloud-native search algorithm that uses accurate predictions of peptide fragment ion intensities and retention times provided by the deep learning framework INFERYYS 2.0. Based on an initial coarse search, INFERYYS performs data-driven model refinement to maximize prediction accuracy. Tandem mass spectra are analyzed without pre-processing or candidate selection using features detected in precursor mass spectra. Instead, all candidates in the isolation window of a given tandem mass spectrum are considered simultaneously and compete for measured fragment ion intensity in one concerted step. CHIMERYYS aims to explain as much measured intensity with as few candidate peptides as possible, resulting in the deconvolution of chimeric spectra. Peptide-spectrum match (PSM)-level false discovery rate (FDR)-control is performed using Percolator [3]. CHIMERYYS profits from cloud-based parallelization and is available through a node in Thermo Scientific™ Proteome Discoverer™ 3.0 software.



RESULTS

CHIMERYYS doubles peptides identification in classical DDA datasets

CHIMERYYS' deconvolution algorithm identifies peptide precursors hidden in chimeric spectra of DDA data files. Here, a digest of a HeLa cell lysate was analyzed using a 1-hour gradient on a Thermo Scientific™ Orbitrap Exploris™ 480 mass spectrometer and processed in Proteome Discoverer software using Sequest HT and CHIMERYYS. The results demonstrate a more comprehensive data analysis when using CHIMERYYS: over 80% of all MS2 spectra were matched to one or more peptide precursors and the average number of PSMs per spectrum substantially increases.

Figure 1. Number of PSMs, peptide and protein groups for a HeLa cell lysate

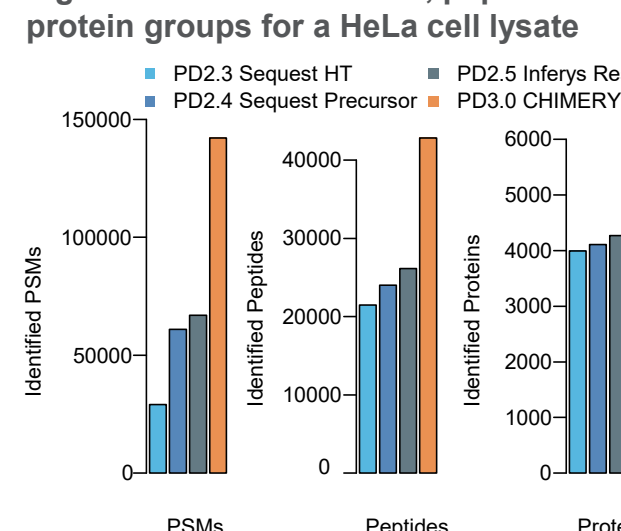
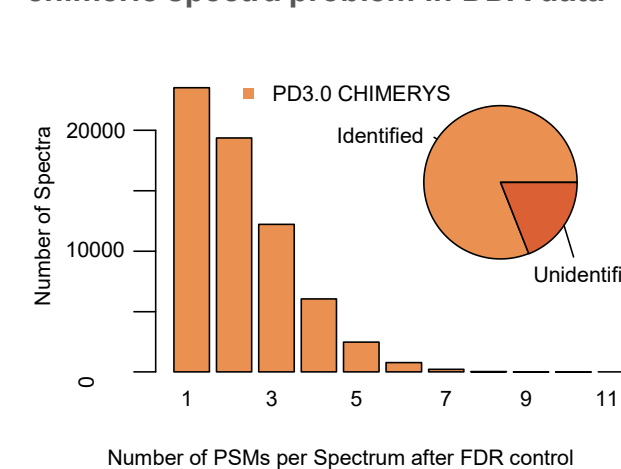


Figure 3. Number of PSMs per spectrum and identification rate achieved by CHIMERYYS demonstrates the extend of the chimeric spectra problem in DDA data

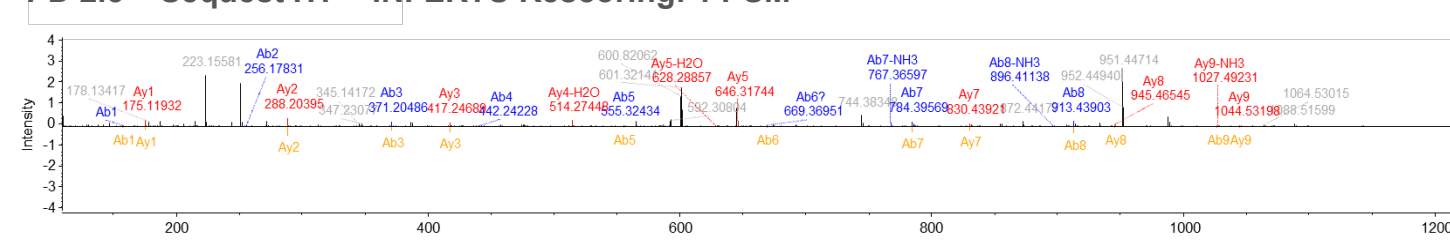


Accurate deconvolution by CHIMERYYS unlocks information hidden in chimeric spectra

CHIMERYYS deconvolutes MS2 spectra by considering all relevant peptide precursors for a given spectrum simultaneously, which then compete for the available experimental intensity in a single step. This results in the identification of several PSMs from chimeric spectra. Using the Proteome Discoverer Spectrum Viewer functionality with direct connection to INFERYYS 2.0, users can visualize the proportional contributions of the individual peptides for every single MS2 spectrum in a mirror plot.

Figure 5. Mirror plot of an experimental spectrum and PSMs identified by Sequest HT + INFERYYS Rescoring (top panel) or CHIMERYYS (bottom panel) at 1% FDR. While INFERYYS Rescoring identifies only one peptide, CHIMERYYS identifies three additional peptides, resulting in a drastically increased explained intensity of the experimental spectrum.

PD 2.5 – Sequest HT + INFERYYS Rescoring: 1 PSM



PD 3.0 – CHIMERYYS: 4 PSMs

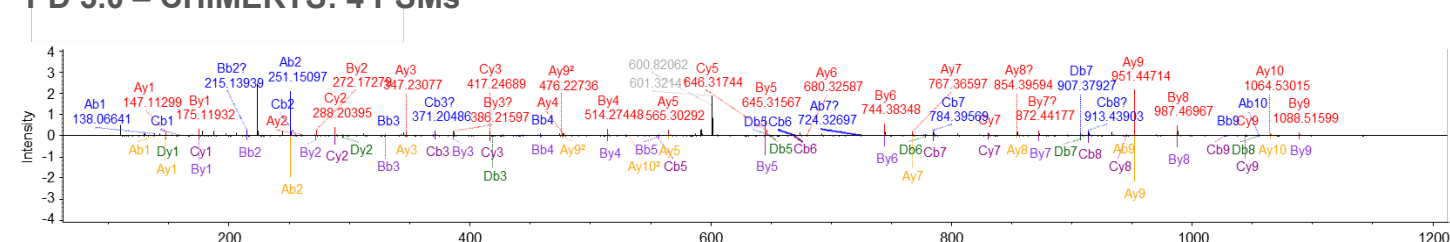


Figure 2. Overlap of peptide identifications by the different search engines and processing workflows

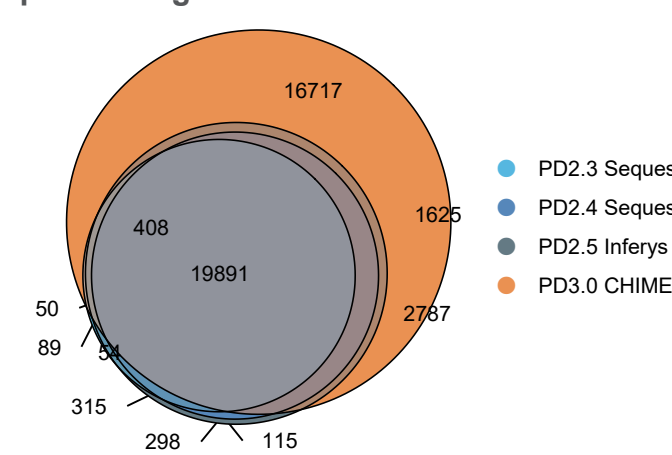
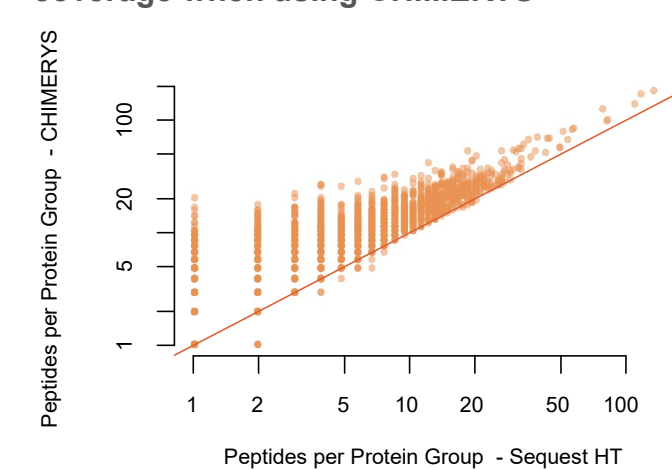


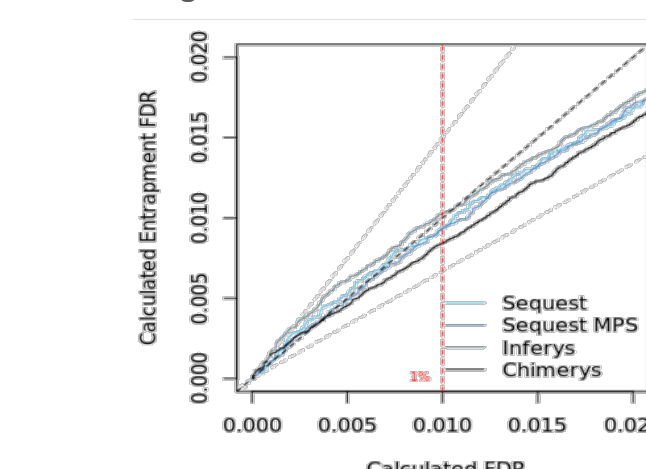
Figure 4. Number of peptides per protein group identified by CHIMERYYS or Sequest HT demonstrates the increase in sequence coverage when using CHIMERYYS



Validation of CHIMERYYS results using entrapment searches

Double-decoy approaches enable the calculation of an entrapment FDR and are common benchmarking methods to determine the correctness of FDR estimations. Here, we utilized a 3x shuffled human database as an entrapment database to demonstrate the accuracy of the PSM-level FDR calculation performed by Percolator on CHIMERYYS search results.

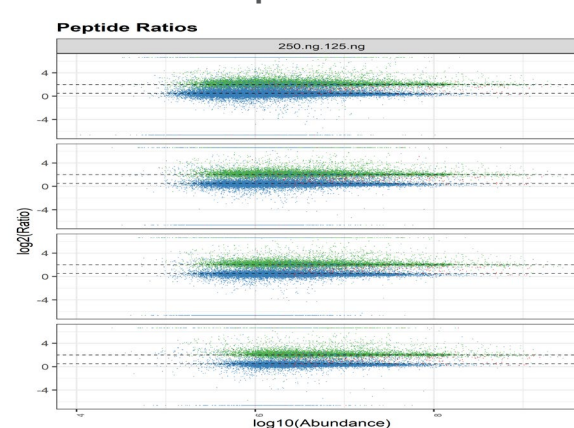
Figure 6. Entrapment FDR vs. calculated FDR analysis using a 3x shuffled human entrapment database across different search engines and workflows.



CHIMERYYS increases the number of accurately quantifiable peptides

Due to the increased analysis depth and comprehensive identification of PSMs and peptides, CHIMERYYS aids in the accurate quantification of label free data sets. We demonstrate this using a two organism dilution series and compare the quantification results using the Minor feature detector node. This demonstrates that CHIMERYYS produces more, especially lower abundant quantified peptides and proteins. In this case, 75% more correctly quantified proteins compared to Sequest HT.

Figure 8. Quantification of Peptides/Proteins from a HeLa/Yeast dilution row experiment.



CHIMERYYS demonstrates an exquisite sensitivity in simulation experiments

To validate CHIMERYYS, we developed an *in-silico* chimeric spectra system (ICS) that spikes *in-silico* generated chimeric spectra into raw files, which can then be used as a ground truth dataset to evaluate search algorithms. Briefly, the system selects seed MS2 spectra with high-confidence identifications from a prior database search from a raw file and convolutes them with several predicted MS2 spectra. To create realistic chimeric data, predicted spectra are derived from peptides with a precursor *m/z* value within the isolation window of the seed MS2 spectrum and a similar predicted retention time. The created raw file is then submitted to both CHIMERYYS and Sequest HT. Using this system, we demonstrate the sensitivity of CHIMERYYS, which recovers >91% of the *in-silico* chimeric spectra in the convoluted data.

Figure 10. Schema of the ICS system for generating a ground-truth dataset containing *in-silico* chimeric spectra

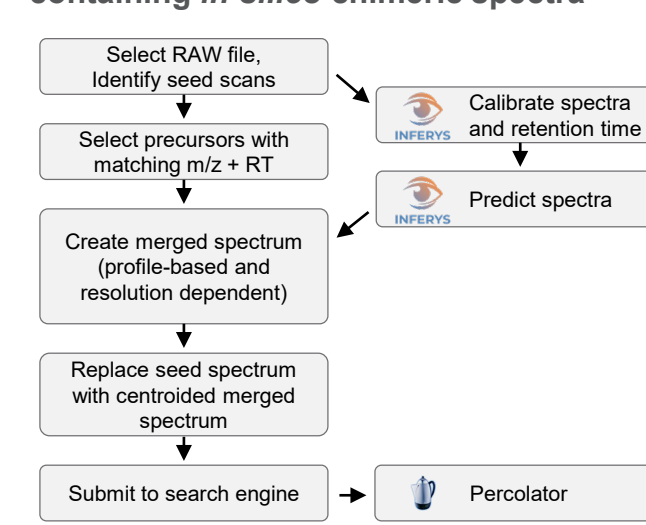


Figure 7. Double-log plot of the data shown in Figure 6, visualizing the low FDR region.

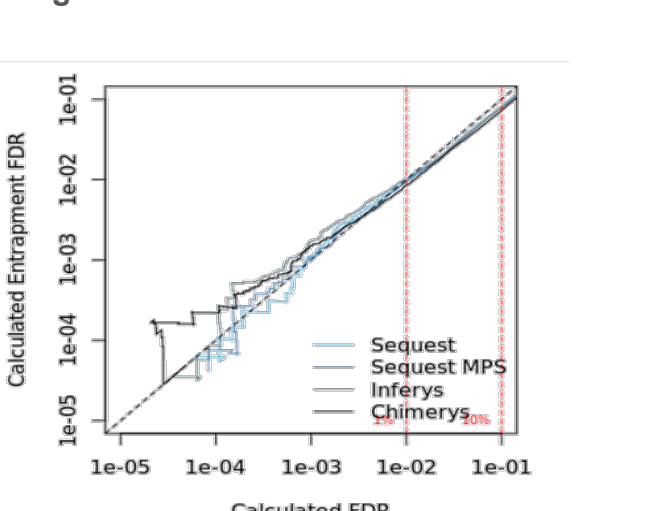


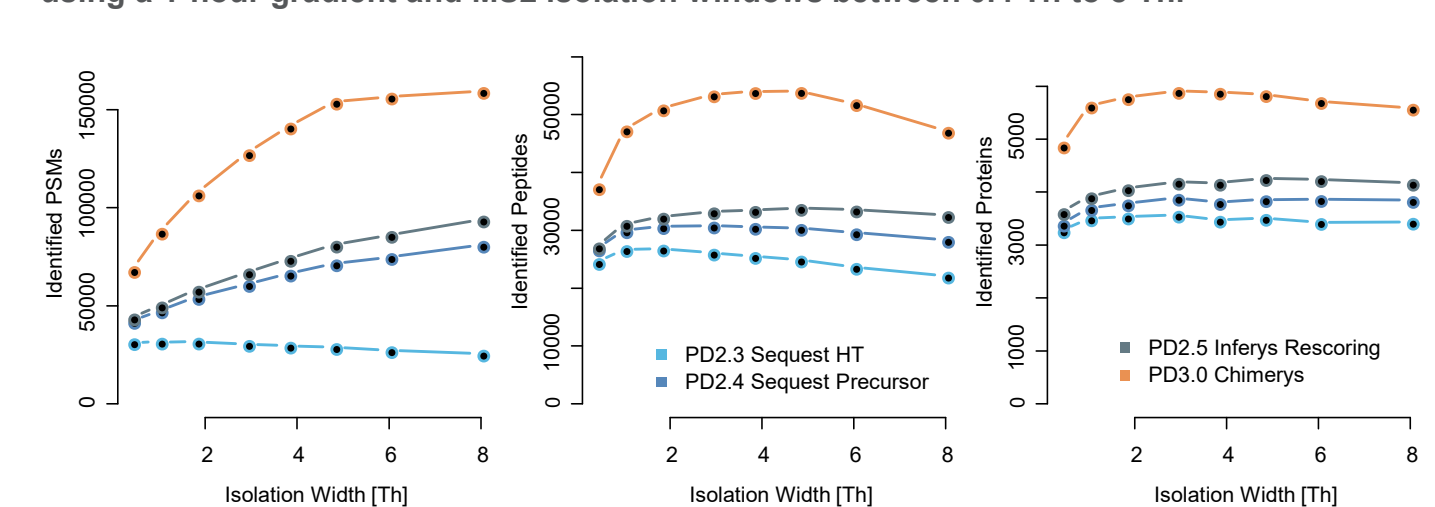
Figure 9. Distribution of quantitative Yeast protein ratios from dilution experiment, (correct: 0.5 * expected < r < 2 * expected)



CHIMERYYS enables optimized acquisition settings and profits of increased MS2 complexity

CHIMERYYS' deconvolution algorithm is optimized for highly complex samples resulting in convoluted MS2 spectra. Hence, it allows for optimizing data acquisition settings to increase measurement efficiency by identifying more proteins per unit time. Here, we demonstrate that CHIMERYYS enables wider DDA isolation windows that result in more chimeric MS2 spectra, providing more identifications while keeping the gradient length constant.

Figure 12. Number of PSMs, peptide and protein groups identified of a DDA HeLa cell lysate digest acquired on a Thermo Scientific™ Orbitrap Eclipse™ Tribrid™ mass spectrometer using a 1-hour gradient and MS2 isolation windows between 0.4 Th to 8 Th.



CHIMERYYS increases the throughput of measurements by allowing shorter gradients

CHIMERYYS uniquely deciphers complex samples and MS2 spectra, enabling shorter acquisition times and gradients for LC-MS/MS measurements without losing peptide or protein information in comparison to Sequest HT. Here, we demonstrate how CHIMERYYS identifies the same number of peptides and protein groups in 1/3 of the measurement time. Shorter gradients increase the gap between Sequest HT and CHIMERYYS using separation times from 8 to 60 min on a classical HeLa cell lysate using a Thermo Scientific™ Vanquish NEO™ liquid chromatography system.

Figure 13. Number of PSMs, peptide and protein groups identified by CHIMERYYS or Sequest HT from digests of a HeLa cell lysate acquired on an Orbitrap Exploris 480 MS with gradient lengths ranging from 8 to 60 minutes.

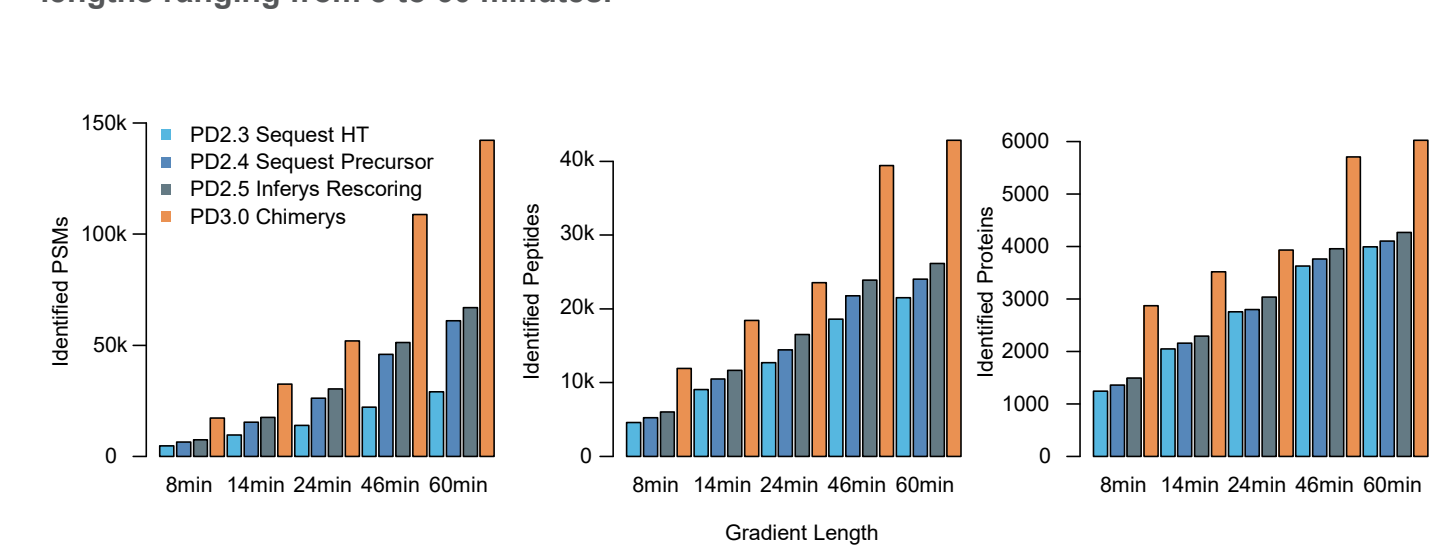
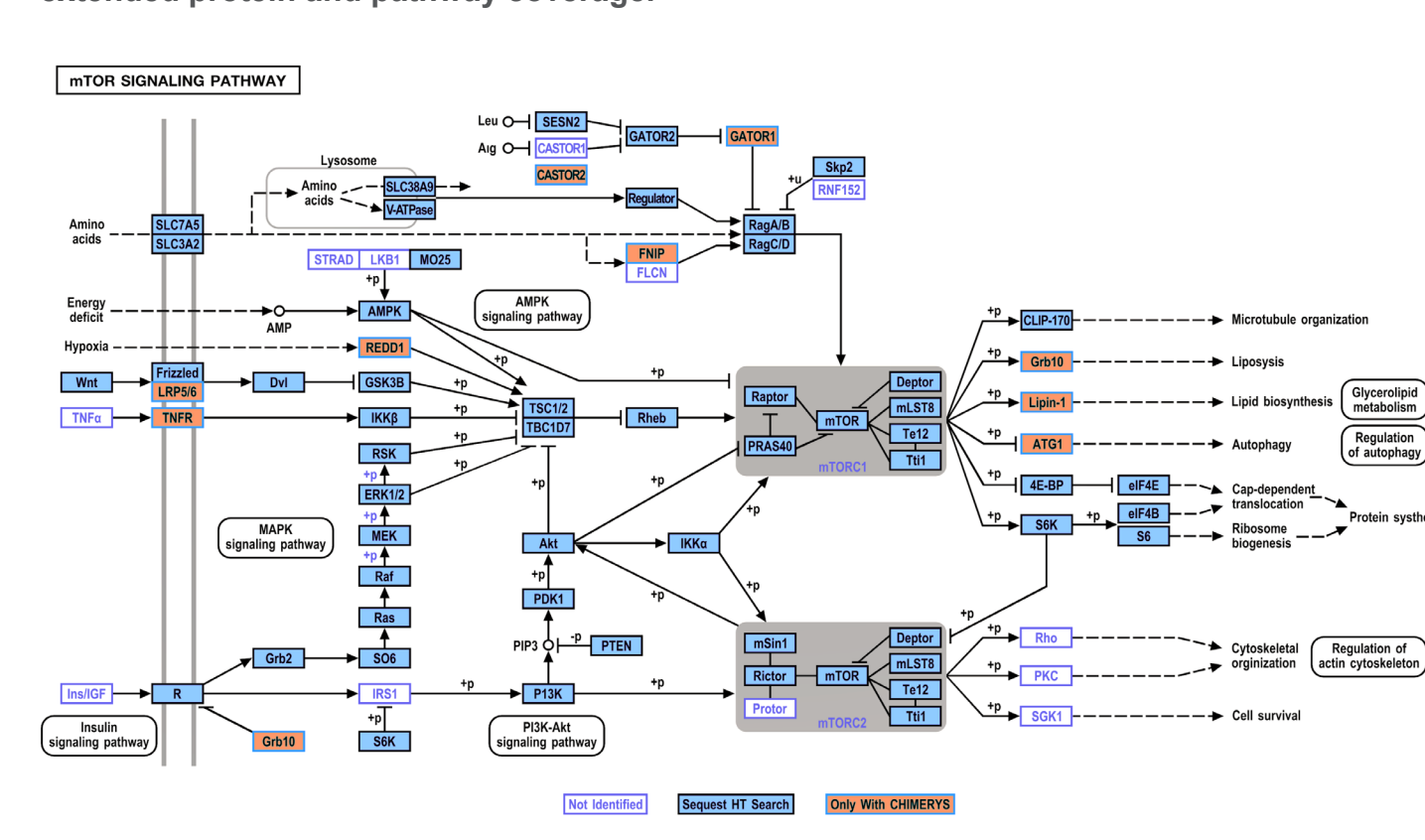


Figure 14. Proteins of the mTOR signaling pathway identified by CHIMERYYS or Sequest HT show the potential for new biological insight that can be generated using CHIMERYYS and extended protein and pathway coverage.



CHIMERYYS outperforms Sequest HT on data from various organisms and complexity

CHIMERYYS is fueled by predictions from INFERYYS 2.0 that are independent of the sample source under investigation. Paired with its resilience with respect to highly complex data, CHIMERYYS is well-equipped to handle fractionated or non-fractionated measurements from organisms from all kingdoms of life [4] and less complex samples like body fluids [5]. Here, we demonstrate its capabilities on a selection of publicly available data.

Figure 15. Protein groups identified by CHIMERYYS and Sequest HT for a fractionated *Arabidopsis thaliana* proteome; raw data from PRIDE Project PXD019483 [4]

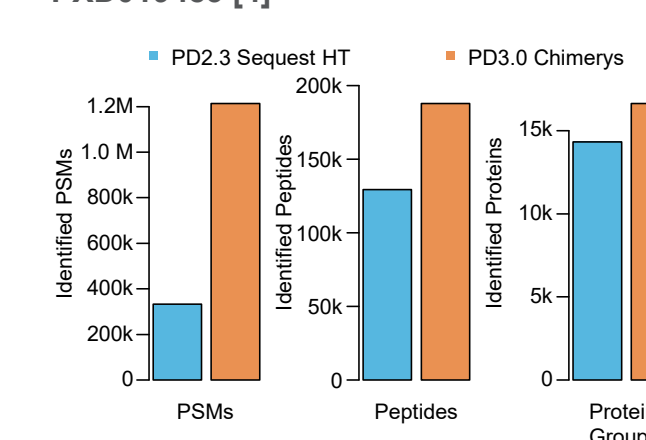
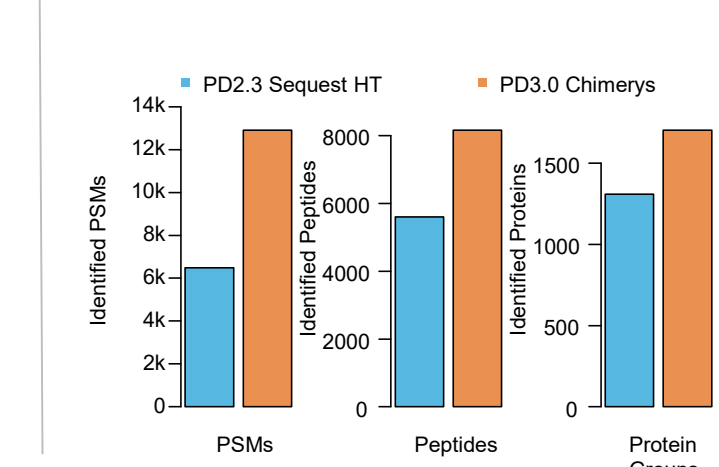


Figure 16. Protein groups identified by CHIMERYYS and Sequest HT for a single 30 min Urine proteome file; raw data from PRIDE Project PXD015087 [5]



CONCLUSIONS

- CHIMERYYS is an innovative, cloud-native search algorithm that uses AI-based predictions to deconvolute chimeric spectra and is fully integrated into Proteome Discoverer 3.0 software
- Using CHIMERYYS results in drastically increased numbers of PSM, peptide and protein group identifications, higher sequence coverage and more confident quantification
- CHIMERYYS excels at analyzing complex samples, enabling more efficient measurements, advanced acquisition settings and shorter gradients to enhance proteomic throughput, productivity and efficiency

REFERENCES

- Zolg, DP; Gessulat, S; Paschke, C, Frejno M, et al. INFERYYS rescoring: Boosting peptide identifications and scoring confidence of database search results. *Rapid Commun Mass Spectrom*. 2021;e9128. <https://doi.org/10.1002/rcm.9128>
- Dorfer V; Maltsev S; Winkler S; Mettler K. CharmerT: Boosting Peptide Identifications by Chimeric Spectra Identification and Retention Time Prediction. *Journal of Proteome Research* 2018 17 (8), 2581-2589. <https://doi.org/10.1021/acs.jproteome.7b00836>
- The M; MacCoss MJ; Noble WS; Käll L. Fast and Accurate Protein False Discovery Rates on Large-Scale Proteomics Data Sets with Percolator 3.0. *J Am Soc Mass Spectrom*. 2016;27(11):1719-1727. <https://doi.org/10.1007/s13361-016-1460-7>
- Müller JB; Geyer, PE; Colaço, A.R, Mann M; et al. The proteome landscape of the kingdoms of life. *Nature* 582, 592–596 (2020). <https://doi.org/10.1038/s41586-020-2402-x>
- Bian, Y; Zheng, R; Bayer, FP; Kuster B; et al. Robust, reproducible and quantitative analysis of thousands of proteomes by micro-flow LC–MS/MS. *Nat Commun* 11, 157 (2020). <https://doi.org/10.1038/s41467-019-13973-x>

ACKNOWLEDGEMENTS

The authors would like to thank the alpha and beta testers and especially Prof. Dr. Bernhard Kuster and Karl Mechtler for the evaluation and validation of CHIMERYYS. The authors wish to thank numerous scientists and colleagues at MSAID and Thermo Fisher Scientific for fruitful discussions and technical assistance.

TRADEMARKS/LICENSES

© 2021 Thermo Fisher Scientific Inc. All rights reserved. CHIMERYYS™, INFERYYS™ and MSAID® and are trademarks of MSAID GmbH. SEQUEST is a trademark of the University of Washington. All other trademarks are the property of Thermo Fisher Scientific and its subsidiaries. This information is not intended to encourage use of these products in any manner that might infringe the intellectual property rights of others.