

# Removing the hidden data dependency of DIA with predicted spectral libraries

Van Puyvelde, B. <sup>§a</sup>, Willems, S. <sup>§a</sup>, Gabriels, R. <sup>§b,c</sup>, Daled, S. <sup>a</sup>, De Clerck, L. <sup>a</sup>, Vande Casteele, S. <sup>a</sup>, Staes, A. <sup>b,c,d</sup>, Impens, F. <sup>b,c,d</sup>, Deforce, D. <sup>a</sup>, Martens, L. <sup>b,c</sup>, Degroeve, S. <sup>b,c</sup>, Dhaenens, M. <sup>\*a</sup>

§: Authors contributed equally

a: ProGenTomics, Laboratory of Pharmaceutical Biotechnology, Ghent University, Ghent, Belgium

b: VIB-UGent Center for Medical Biotechnology, Ghent, Belgium

c: Department of Biomolecular Medicine, Ghent University, Ghent, Belgium

d: VIB Proteomics Core, Ghent, Belgium

\*: Corresponding author; Telephone: +32 (0)9 264 83 56; E-mail: [maarten.dhaenens@ugent.be](mailto:maarten.dhaenens@ugent.be)

**Data-Independent Acquisition (DIA) generates comprehensive yet complex mass spectrometric data, which imposes the use of data-dependent acquisition (DDA) libraries for deep peptide-centric detection.**

**We here show that DIA can be redeemed from this dependency by combining predicted fragment intensities and retention times with narrow window DIA. This eliminates variation in library building and omits stochastic sampling, finally making the DIA workflow fully deterministic. Especially for clinical proteomics, this has the potential to facilitate inter-laboratory comparison.**

## Significance of the Study

Data-independent acquisition (DIA) is quickly developing into the most comprehensive strategy to analyse a sample on a mass spectrometer. Correspondingly, a wave of data analysis strategies has followed suit, improving the yield from DIA experiments with each iteration. As a result, a worldwide wave of investments in DIA is already taking place in anticipation of clinical applications. Yet, there is considerable confusion about the most useful and efficient way to handle DIA data, given the plethora of possible approaches with little regard for compatibility and complementarity. In our manuscript, we outline the currently available peptide-centric DIA data analysis strategies in a unified graphic called the DIAMond DIAGram. This leads us to an innovative and easily adoptable approach based on predicted spectral information. Most importantly, our contribution removes what is arguably the biggest bottleneck in the field: the current need for Data Dependent Acquisition (DDA) prior to DIA analysis. Fractionation, stochastic data acquisition, processing and identification all introduce bias in the library. By generating libraries through data independent, i.e. deterministic acquisition, stochastic sampling in the DIA workflow is now fully omitted. This is a crucial step towards increased standardization. Additionally, our results demonstrate that a proteome-wide predicted spectral library can surrogate an exhaustive DDA Pan-Human library that was built based on 331 prior DDA runs.

# Article

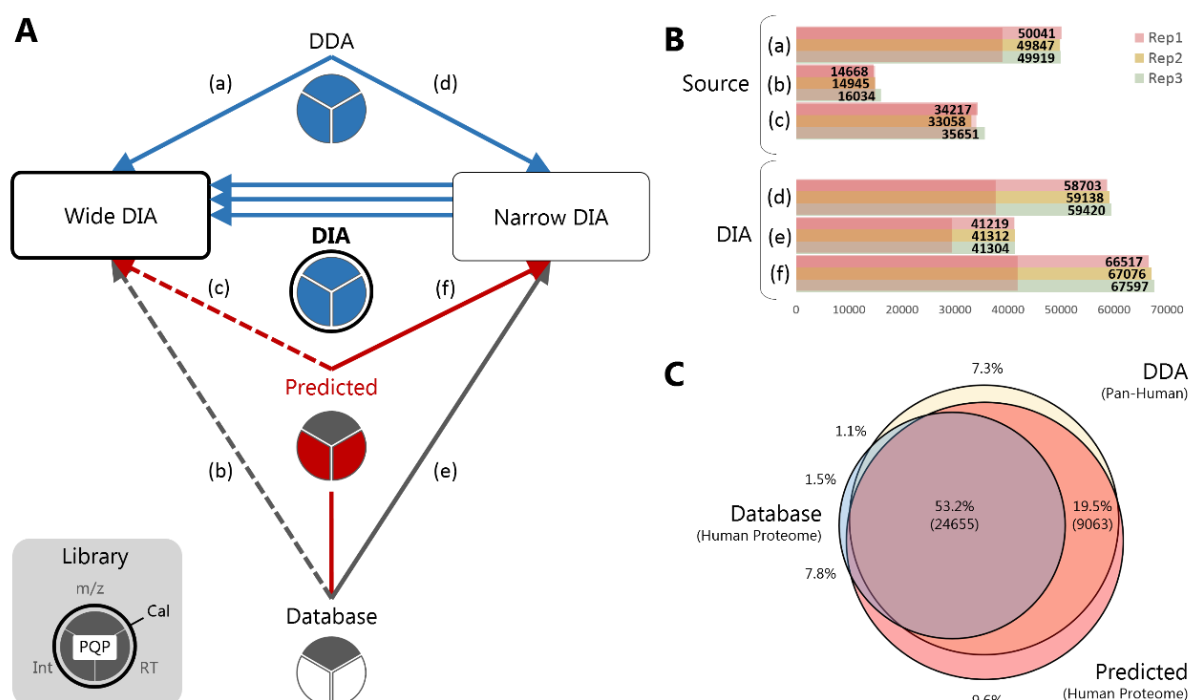
With DIA, an MS instrument regularly measures precursor ions and continuously cycles through predefined mass over charge ratio ( $m/z$ ) windows to equally regularly measure the intensity of their fragment ions throughout a liquid chromatography (LC) gradient. This is both more qualitative and quantitative than data-dependent acquisition (DDA), where precursor ions are measured intermittently while fragment ions are only measured stochastically. However, the complexity of DIA data has shown to be very challenging.

To date, the most common way to address this complexity is using previously identified peptides from DDA as targets in the DIA data. First, DDA peptide identifications are translated into a spectral library with Peptide Query Parameters (PQPs), which typically contain the sequence as well as the analytical coordinates ( $m/z$ , intensity, and retention time or RT) for the observed ions for a given peptide. These PQPs are then used to compute an evidence score for each target peptide, based on its fragment traces in DIA <sup>[1]</sup>. Ultimately, these evidence scores are supplemented with additional features, e.g. ppm and RT errors, allowing a semi-supervised machine learning algorithm to weigh and re-score the target peptides to obtain a maximum of true targets at an empirically determined False Discovery Rate (FDR) using the target-decoy approach <sup>[2][3][4]</sup>.

Unfortunately, deriving PQPs from DDA data intrinsically means transferring its limitations. In fact, fractionation, stochastic data acquisition, processing and identification introduce bias in the library and require considerable effort. This compromises inter-laboratory comparison and can even alter the biological conclusions between labs <sup>[5]</sup>. However, thanks to the availability of state-of-the-art prediction algorithms, these PQPs can now be predicted directly, setting the stage for much easier and much more reproducible peptide-centric DIA data extraction <sup>[6][7][8]</sup>.

Here, we compare the effect of using libraries from different origins on peptide-centric approaches, by assessing their qualitative and quantitative performance on a public wide window (10 - 20  $m/z$ ) DIA dataset of HeLa cells <sup>[9]</sup> (**Figure 1**). Three basic spectral libraries were used here, with PQPs derived from (a) an experimental DDA dataset, (b) a protein sequence database (FASTA), and (c) a predicted spectral dataset. Each of these three libraries can be used directly as a source library, or can be converted into a DIA library by using them first on a narrow window (2  $m/z$ ) DIA dataset of the sample. The resulting six possible libraries can all be used alike by the EncyclopeDIA software to identify and quantify wide window DIA data <sup>[9]</sup>.

In-house or public DDA source libraries are frequently built by extensive fractionation of samples. With adequate statistical control, such proteotypic libraries allow direct peptide detections in wide window DIA (**Figure 1Aa**) <sup>[10]</sup>. We illustrate this by using the publically available Pan-Human library, which contains nearly 10.000 proteins derived from 331 DDA runs on a range of human cell lines and tissues <sup>[11]</sup> (**Figure 1Ba**). To reduce the effort and variability from DDA library building, a library-free peptide-centric data analysis workflow was proposed recently <sup>[12]</sup>. Herein, the PECAN (or Walnut) scoring algorithm allows direct detection of peptides derived from a FASTA in wide window DIA data (**Figure 1Ab**). This is akin to a source library that (i) contains only peptide sequences and  $m/z$  coordinates, and (ii) lacks prior selection of proteotypic peptides. On wide window DIA data this approach thus provides a limited number of PQPs, which is not sufficient to differentiate between the high number of false targets, i.e. true negatives, and the lower number of true positives in the library <sup>[13]</sup>. This manifests as indiscernible target and decoy score distributions, resulting in a very high False Negative Rate (FNR) (**Figure 1Bb**).



**Figure 1. Peptide-centric data extraction from wide window DIA data.** (A) DIAMOND DIAgram presenting peptide-centric strategies for DIA data extraction. Peptide-centric approaches rely on libraries (central column) that contain Peptide Query Parameters (PQQs) which are derived from the peptide sequence and can additionally contain the three ion coordinates, i.e. mass to charge ratio ( $m/z$ ), Intensity (Int) and retention time (RT) (three-part pie charts). These can either be experimental (blue), theoretical (grey), or predicted (red). PQQs are used to score the evidence of peptide detections in continuous DIA data (boxes). These are supplemented with additional features of the match so that a support vector machine can weigh and re-score them to obtain a maximum of true targets at an empirically determined FDR using the target-decoy approach (arrow heads). DDA source libraries (both in-house and public) only comprise prior proteotypic peptide identifications and contain measured PQQs for all three ion coordinates. These are therefore directly applicable to quantify peptides in 10 – 20  $m/z$  wide window DIA (Wide DIA) data (a). However, when a proteome FASTA is used as a source library, sensitivity is reduced (dashed arrow), i.e. too many false negatives are produced due to the high statistical burden (b). This also holds for libraries with predicted fragment intensities (MS<sup>2</sup>PIP) and RT (Elude), albeit to a lesser extent (c). Prior 2  $m/z$  narrow window DIA (Narrow DIA) provides the specificity to remove false targets in the sample first (d)(e)(f). The DIA ion coordinates from these detections can additionally be integrated into new and calibrated PQQs (cal). These DIA libraries, called chromatogram libraries, can be derived from any source library (triple arrow). (B) Doubly and triply charged peptide detections in wide window DIA following each of the routes depicted in (A). Shading highlights the number of peptides that is detected in triplicate wide window DIA runs with at least three transitions, allowing robust quantification. (C) Comparison of the identified peptide sequences in Wide DIA for route (d), (e) and (f). The large overlap shows that all three approaches detect proteotypic peptides. Only peptides of double and triple charge that are detected in triplicate wide window DIA runs with at least three transitions are shown.

Here we propose a promising way to improve upon the FASTA source library - while still omitting prior DDA - by predicting fragment ion intensity and RT *in silico* (Figure 1Ac and Figures S1-S2). Using a spectral dataset with such predicted fragment intensities (MS<sup>2</sup>PIP) and peptide RTs (Elude) more than doubles the number of peptides detected in the wide window DIA (Figure 1Bc) [6][14]. However, considering all tryptic peptides in a Human proteome still underperforms compared to the Pan-Human DDA library, which is fully contained in the predicted spectral dataset (Figure 1Ba and 1Bc). Notably, this is not due to poor prediction because predicting only those peptides present in the Pan-Human library performs very similar to using the Pan-Human library directly (Figure S3) and the underperformance can thus only be attributed to the many false targets when using the complete database [10]. An elegant way to filter out false target peptides upfront, is by measuring a pool from every condition with staggered narrow window DIA (Figure 1Ad, 1Ae and 1Af). This reduces MS2 chimericity to DDA-like quality in a DIA setting, allowing detection with increased specificity. This accurate prior filtering makes the statistical burden of false targets in the wide window DIA surmountable again. Notably, due to instrument limitations this *Precursor*

*Acquisition Independent From Ion Count* (PACIFIC)<sup>[15]</sup> can currently only be performed by means of gas phase fractionation (GPF), i.e. sampling different *m/z* regions separately<sup>[9]</sup>. Still, the added acquisition depth and specificity allows the detection of 88k (DDA), 47k (FASTA) and 95k (predicted) peptides in six narrow window GPF DIA runs of a HeLa cell lysate (**Figure S4**). To assure that this additional filtering is accurate, we confirmed the estimated FDR by using an entrapment experiment wherein we included *Pyrococcus furiosus* proteins as false targets alongside the expected human proteins in the respective source libraries<sup>[16]</sup>. Hereby, the measured FDR for narrow window DIA filtering is 2% for the DDA, 1% for the FASTA, and 1% for the predicted source library, in accordance with the theoretically estimated FDR based on the target-decoy strategy. In the process, we can measure the identification cost of adding false targets: adding 3-6% false targets results in an average decrease of 1-2% in detections (see Entrapment section in Methods).

Additionally, the peptide detections in narrow window DIA can be translated into novel and integrated PQPs, which are calibrated to the specific LCMS system and are specific to DIA (**Figure 1A**). This approach was recently made readily applicable as chromatogram libraries: DIA libraries of narrow window DIA peptide detections comprising their calibrated PQPs<sup>[9]</sup>. Such chromatogram libraries outperform direct wide window DIA extraction for every source library. The modest gain for a DDA source library (~20%) derives mainly from PQP calibration, as only 50% of the source peptides was filtered out (**Figure 1Ba and 1Bd**). In contrast, in the FASTA source library, 98,5% of the peptides were filtered out, and RT and intensity coordinates were generated *de novo*. Taken together, this resulted in the largest gain (~170%) (**Figure 1Bb and 1Be**). Finally, the chromatogram library derived from a predicted spectral library increases the number of detections by ~100% compared to direct wide window DIA data extraction, making it the most efficient overall peptide detection strategy of the DIAMOND DIAGRAM (**Figure 1Bc and 1Bf**). The large overlap between the peptide sequences detected by all three chromatogram libraries convincingly shows that the Pan-Human library is very exhaustive and that all three chromatogram libraries mainly detect proteotypic peptides (**Figure 1C**). Peptides unique to the Pan-Human library include very high molecular masses that were not predicted, high molecular weight peptides that generate many doubly charged transitions that are not predicted by default, as well as very small peptides with inherently poor RT or fragmentation pattern predictions. Peptides that are unique to the predicted library are mainly peptides not present in the Pan-Human source library. Importantly, the PQP requirements of the source library for building chromatogram libraries on narrow window DIA are relatively liberal: the measured Pan-Human library was acquired on a TripleTOF instrument but allows wide window DIA data peptide detection on an Orbitrap instrument. The *in silico* equivalent is that 95% of the detected peptides overlap when the MS<sup>2</sup>PIP engine is trained on either Orbitrap or TripleTOF data. As a result, other fragment ion intensity predictors such as Prosit and Deep Mass<sup>[7][8]</sup> perform similarly when combined with narrow window DIA<sup>[17]</sup> (**Figure S5**).

We therefore conclude that predicted libraries are highly relevant and performant for wide window DIA identification, and that three elements of a spectral library affect its overall performance: (i) the amount of false targets included, (ii) the amount of informative PQPs, and (iii) the accuracy of PQPs on the specific instrument setup. In this study, we could show that a narrow window DIA acquisition of six GPFs combined with a predicted spectral library of the full human proteome was able to surrogate a measured DDA Pan-Human library, thus liberating the DIA workflow from any stochastic acquisition. Especially for clinical proteomics, this can facilitate inter-laboratory comparison. Importantly, the software tools MS<sup>2</sup>PIP, ELUDE and EncyclopeDIA are all instrument independent, publicly available, and mutually compatible, thus making this workflow immediately accessible to everybody interested.

## Code availability

MS<sup>2</sup>PIP, Elude and EncyclopeDIA are open source, licensed under the Apache-2.0 License, and are hosted on [https://github.com/compomics/ms2pip\\_c](https://github.com/compomics/ms2pip_c), <https://github.com/percolator/percolator> and <https://bitbucket.org/searleb/encyclopedia/wiki/Home>. All supporting material is available on <https://github.com/brvpuve/MS2PIP-for-DIA/>.

## Acknowledgments

This research was mainly funded by mandates from the Research Foundation Flanders (FWO) awarded to BVP [grant number 11B4518N], RG [grant number 1S50918N] and MD [12E9716N]. Partial funding was received through project grants from the FWO [G013916N and G042518N], from the European Union's Horizon 2020 Program under Grant Agreement 823839 [H2020-INFRAIA-2018-1], and from a PhD grant from the Flanders Agency Entrepreneurship and Innovation (VLAIO) awarded to LDC [SB-141209].

## Competing interests

The authors have declared no conflict of interest.

## Author Contributions

BVP performed all data analysis at the ProGenTomics facilities. The initial experimental design was conceived and performed by BVP, SW, MD, SDa, LDC, AS, DD and FI. RG, SDe and LM performed all machine learning predictions. MD, BVP, RG and SW wrote the draft manuscript. All authors provided critical feedback during research and writing. MD conceived the idea of using predicted libraries for DIA data extraction and supervised the project.

## References

- [1] C. Ludwig, L. Gillet, G. Rosenberger, S. Amon, B. C. Collins, R. Aebersold, *Mol. Syst. Biol.* **2018**, DOI 10.15252/msb.20178126.
- [2] L. Reiter, O. Rinner, P. Picotti, R. Huttenhain, M. Beck, M. Y. Brusniak, M. O. Hengartner, R. Aebersold, *Nat. Methods* **2011**, *8*, 430.
- [3] L. Kall, J. D. Canterbury, J. Weston, W. S. Noble, M. J. MacCoss, *Nat. Methods* **2007**, *4*, 923.
- [4] J. Telesman, H. L. Röst, G. Rosenberger, U. Schmitt, L. Malmström, J. Malmström, F. Levander, *Bioinformatics* **2015**, DOI 10.1093/bioinformatics/btu686.
- [5] E. Govaert, K. Van Steendam, S. Willems, L. Vossaert, M. Dhaenens, D. Deforce, *Proteomics* **2017**, *17*, DOI 10.1002/pmic.201700052.
- [6] R. Gabriels, L. Martens, S. Degroove, *Nucleic Acids Res.* **2019**, DOI 10.1093/nar/gkz299.
- [7] S. Gessulat, T. Schmidt, D. P. Zolg, P. Samaras, K. Schnatbaum, J. Zerweck, T. Knaute, J. Rechenberger, B. Delanghe, A. Huhmer, U. Reimer, H.-C. Ehrlich, S. Aiche, B. Kuster, M. Wilhelm, *Nat. Methods* **2019**, DOI 10.1038/s41592-019-0426-7.
- [8] S. Tiwary, R. Levy, P. Gutenbrunner, F. Salinas Soto, K. K. Palaniappan, L. Deming, M. Berndt, A. Brant, P. Cimermanic, J. Cox, *Nat. Methods* **2019**, DOI 10.1038/s41592-019-0427-6.
- [9] B. C. Searle, L. K. Pino, J. D. Egerton, Y. S. Ting, R. T. Lawrence, B. X. MacLean, J. Villén, M. J. MacCoss, *Nat. Commun.* **2018**, DOI 10.1038/s41467-018-07454-w.
- [10] G. Rosenberger, I. Bludau, U. Schmitt, M. Heusel, C. L. Hunter, Y. Liu, M. J. MacCoss, B. X.

- MacLean, A. I. Nesvizhskii, P. G. A. Pedrioli, L. Reiter, H. L. Rost, S. Tate, Y. S. Ting, B. C. Collins, R. Aebersold, *Nat. Methods* **2017**, *14*, 921.
- [11] G. Rosenberger, C. C. Koh, T. Guo, H. L. Röst, P. Kouvonen, B. C. Collins, M. Heusel, Y. Liu, E. Caron, A. Vichalkovski, M. Faini, O. T. Schubert, P. Faridi, H. A. Ebhardt, M. Matondo, H. Lam, S. L. Bader, D. S. Campbell, E. W. Deutsch, R. L. Moritz, S. Tate, R. Aebersold, *Sci. Data* **2014**, DOI 10.1038/sdata.2014.31.
- [12] Y. S. Ting, J. D. Egertson, J. G. Bollinger, B. C. Searle, S. H. Payne, W. S. Noble, M. J. MacCoss, *Nat. Methods* **2017**, *14*, 903.
- [13] N. Colaert, S. Degroeve, K. Helsens, L. Martens, *J Proteome Res* **2011**, *10*, 5555.
- [14] L. Moruz, A. Staes, J. M. Foster, M. Hatzou, E. Timmerman, L. Martens, L. Kall, *Proteomics* **2012**, *12*, 1151.
- [15] A. Panchaud, A. Scherl, S. A. Shaffer, P. D. Von Haller, H. D. Kulasekara, S. I. Miller, D. R. Goodlett, *Anal. Chem.* **2009**, DOI 10.1021/ac900888s.
- [16] M. Vaudel, J. M. Burkhardt, D. Breiter, R. P. Zahedi, A. Sickmann, L. Martens, *J Proteome Res* **2012**, *11*, 5065.
- [17] B. C. Searle, K. E. Swearingen, C. A. Barnes, T. Schmidt, S. Gessulat, B. Kuster, M. Wilhelm, *bioRxiv* **2019**, 682245.

# SUPPORTING INFORMATION

## CONTENTS

Introduction .....	1
Prediction Models: Elude and MS <sup>2</sup> PIP .....	1
Elude: Retention Time prediction .....	1
MS <sup>2</sup> PIP: intensity Prediction .....	3
Library Generation .....	4
DDA.....	4
Database (FASTA).....	4
Predicted.....	4
DIA .....	5
RAW file processing .....	5
DIA data analysis: EncyclopeDIA .....	5
FDR assessment by entrapment .....	7
References .....	9

## Introduction

DIA data has been presented as a permanent record of everything. Thus, applying our novel approach can significantly broaden the biological perspective on newly acquired as well as existing data. Using predicted spectral libraries to replace measured DDA libraries not only reduces workload and increases reproducibility; it will also facilitate the implementation of DIA into more applied fields such as clinical proteomics. Since the software tools MS<sup>2</sup>PIP, Elude and EncyclopeDIA are instrument independent, publicly available and mutually compatible, the presented workflow is accessible to everybody and directly applicable <sup>[1–3]</sup>. Therefore, we present this methods section in the form of a systematic tutorial. Briefly, both source and DIA libraries can be used in EncyclopeDIA to detect peptides in wide window DIA. However, converting source libraries into a DIA library will significantly improve the number of peptides that can be detected. This requires an additional narrow window DIA of several gas phase fractions (GPF) of a mixture of the samples. When these GPFs are acquired in the same batch as the wide window DIA, the benefit of PQP calibration is maximized.

All external resources are available on GitHub <https://github.com/brvpuyve/MS2PIP-for-DIA> for reproducibility.

## Prediction Models: Elude and MS<sup>2</sup>PIP

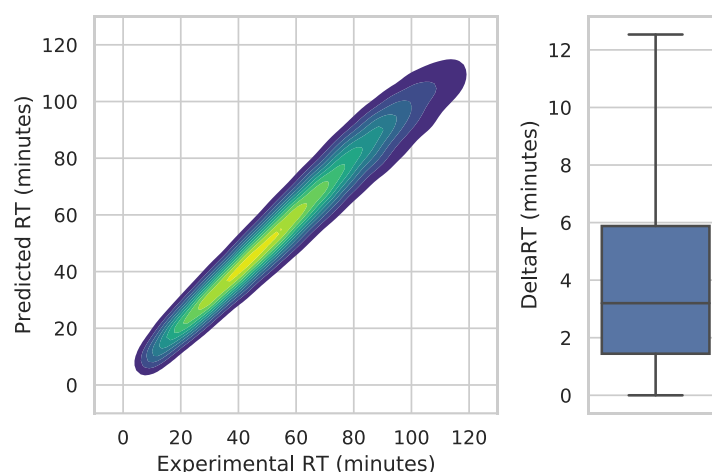
### Elude: Retention Time prediction

For RT prediction, we employed Elude (version 3.02), which is available from the Percolator GitHub repository (<https://github.com/percolator/percolator/releases>) <sup>[2]</sup>.

We trained an Elude model on the Pan-Human spectral library <sup>[4]</sup>. The spectral library was downloaded from SWATHAtlas in SpectraST SPTXT file format. The peptide sequences and their respective RTs were parsed from the SPTXT file to an MS<sup>2</sup>PIP PEPREC file using the `speclib_to_mgf.py` script, which is



available in the conversion\_tools folder of the MS<sup>2</sup>PIP GitHub repository. Out of all consensus peptide spectra built from five or more identified spectra, 10000 peptides and their mean RTs were randomly sampled for training, 10000 were randomly sampled for testing and all remaining were used for final validation of the model. The training, test and validation datasets were converted and written to the Elude input file format. Through the Elude command line interface, we trained a model with the training and test subsets. Subsequently, we used the model to predict RTs for the validation subset of the dataset. The median absolute difference in experimental and predicted RTs (DeltaRT) of the validation dataset was 3.2 minutes and 95% of the DeltaRTs were less than 12.1 minutes (Figure S1). The model predictions have a Spearman rank correlation with the validation RTs of 0.98.



**Figure S1. Evaluation of the trained Elude model.** Left: Contour plot of all predicted and experimental retention times (RTs) in minutes. Right: Boxplot of all absolute differences between experimental and predicted RT (DeltaRT) in minutes. The box displays the first (Q1), second (Q2), and third (Q3) quartiles, the whiskers display Q1 - 1.5 times the interquartile range (IQR) and Q3 + 1.5 times the IQR, respectively. Outliers are not shown.

The spectral library contains oxidation and carbamidomethylation peptide modifications. As a result, the currently trained Elude model is only able to predict RTs for unmodified peptides and peptides containing these modifications. The RTs included in the original Pan-Human SPTXT spectral library are normalized to the iRT Kit peptide sequences by SpectraST. All RT values predicted by the Elude model therefore take over this normalization. As is the case for experimental RTs, the predicted RTs are aligned to the experimental dataset by EncyclopeDIA. The Elude model file is available on our GitHub repository.

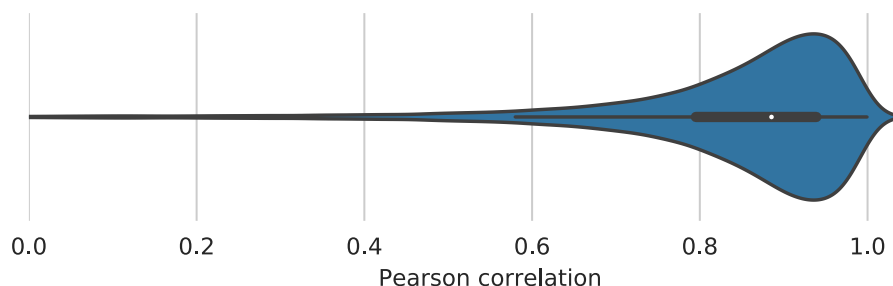


## MS<sup>2</sup>PIP: intensity Prediction

MS<sup>2</sup>PIP, the MS<sup>2</sup> Peak Intensity Predictor, first published by Degroove et al., underwent significant improvements since its initial release in 2013 <sup>[5]</sup>. Currently, a broad array of fragmentation models is available (e.g. Orbitrap-HCD, iontrap-CID, TripleTOF 5600+, ...) <sup>[1]</sup>. This gives the user the liberty to employ a model fit to the experimental setup. As both the narrow and wide window DIA datasets used in this project were obtained on a Q Exactive HF instrument (Thermo Fisher Scientific, Massachusetts, US), we employed MS<sup>2</sup>PIP's Orbitrap-HCD model, with the exception of the TT5600 model that was used for assessing PQP requirements (see main text). To further validate the application of this model, we calculated the correlations between MS<sup>2</sup>PIP predicted spectra and experimental spectra from the EncyclopeDIA DDA runs.

The HeLa DDA dataset of the EncyclopeDIA article (MassIVE MSV000082805) was imported into Progenesis Q1 for Proteomics (Nonlinear Dynamics, Newcastle upon Tyne, UK) with default parameters. The peakpicked spectra were exported as .mgf and searched with Mascot 2.6.1 against the aforementioned human FASTA. Carbamidomethylation of Cysteine and oxidation of Methionine were respectively set as fixed and variable modifications. The precursor tolerance was set to 50 ppm and the fragment tolerance was set to 0.02 Da. The search included all 2+ and 3+ precursors, allowing up to 2 tryptic missed cleavages. Afterwards, the results were reimported into Progenesis Q1 for Proteomics and converted to an .msp spectral library.

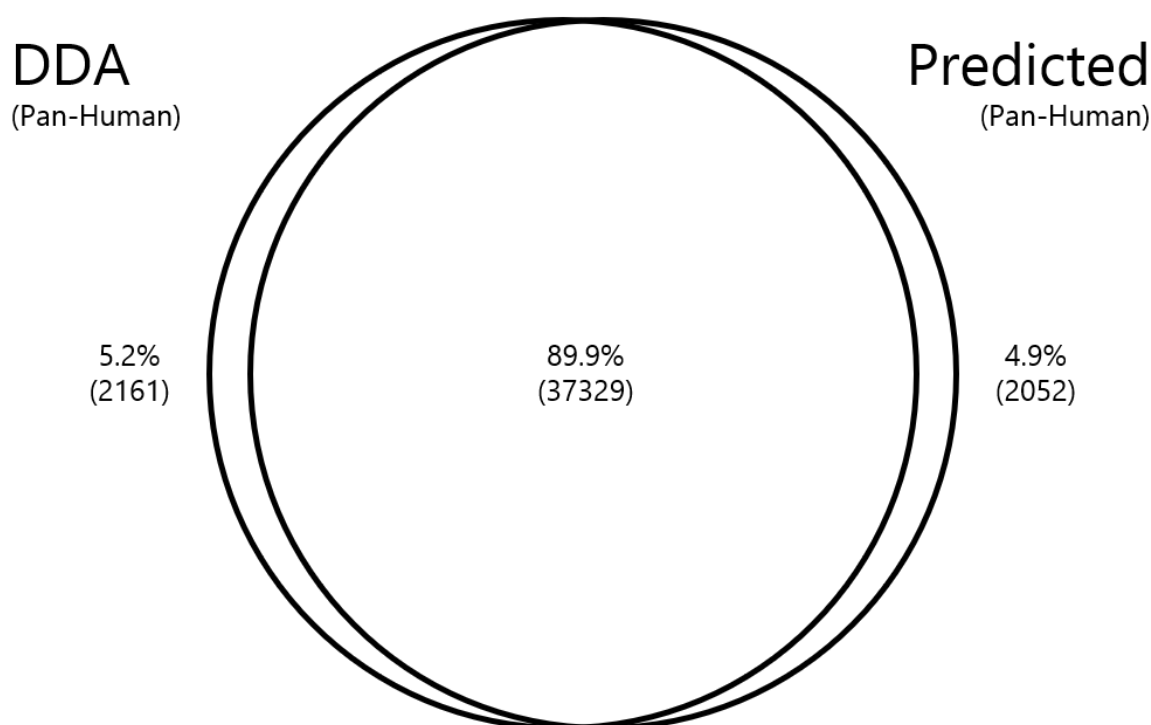
The .msp spectral library was converted back to an .mgf and an MS<sup>2</sup>PIP PEPREC input format using the `speclib_to_mgf.py` script. Both files were then run through MS<sup>2</sup>PIP with the Orbitrap-HCD model, after which Pearson correlation coefficients (PCCs) were calculated for each experimental spectrum and its prediction. This resulted in a median PCC of 0.88 with an interquartile range of {0.795297, 0.938911} (Figure S2)



**Figure S2. Pearson correlations between intensities of measured DDA and MS<sup>2</sup>PIP predicted fragments.** Violin plot showing the distribution of Pearson correlation coefficients between the MS<sup>2</sup>PIP model predictions and the experimental spectra from the EncyclopeDIA article HeLa DDA dataset.

A second experiment was performed to evaluate the performance of predicted libraries. More specifically, as was done in Gessulat et al. <sup>[6]</sup>, a clone of the Pan-Human library was produced using the HCD model and this was applied on the narrow-window DIA data, producing a chromatogram library containing 82.6k unique peptides. Afterwards, the Pan and Pan Clone chromatogram libraries were used in the peptide extraction of triplicate wide-window DIA runs. On average 63k and 62k peptides were identified at 1.0% FDR when searching the wide-window DIA data against the Pan-Human and the Pan Clone chromatogram library, respectively. The quantification reports on peptide and protein level were saved by EncyclopeDIA as .txt files and eventually imported in excel. Then, we manually filtered out all the peptide sequences with less than 3 fragment ions and those having an intensity of zero in at least one of the three replicates. The resulting reproducible peptide sequences were put in a Venn diagram to visualize the percentage overlap (Figure S3). The large overlap demonstrates i) the performance of

the HCD fragmentation model of MS<sup>2</sup>PIP and ii) the retention time prediction of ELUDE to accurately mimic the fragmentation and retention time pattern of peptide sequences acquired on a TripleTOF instrument.



**Figure S3: Overlap in peptides detected by DDA vs predicted chromatogram libraries.** All peptides in a measured Pan-Human library were cloned by predicting their fragmentation spectra using MS<sup>2</sup>PIP and their retention times using ELUDE. A DIA library from a predicted library can extract peptides equally well from wide window DIA data compared to a DDA Pan-Human source library, a logical consequence of good quality predictions.

## Library Generation

### DDA

An EncyclopeDIA .dlib version of the Pan-Human spectral library is publicly available on the EncyclopeDIA BitBucket homepage <sup>[4]</sup>. This version contains 211k unique precursors (159k unique peptide sequences). Alternatively, EncyclopeDIA accepts Skyline .BLIB, Spectronaut .csv, MaxQuant msms.txt, .TraML and .msp files.

### Database (FASTA)

Using a FASTA database does not require a separate library. More specifically, Walnut (a GUI re-implementation of the PECAN algorithm) is part of EncyclopeDIA and can directly detect peptides from DIA data using a FASTA database <sup>[7]</sup>.

Here, we used the human SwissProt proteome (UP000005640 downloaded on 12 February 2019, 20426 target sequences) downloaded as FASTA. The proteome was concatenated with the iRT FASTA obtained from the Biognosys webpage (on 12 February 2019) <sup>[8]</sup>.

### Predicted

Creating a predicted spectral library requires three steps: (i) creating an MS<sup>2</sup>PIP input PEPREC (peptide record) file from a FASTA, (ii) feeding that file to MS<sup>2</sup>PIP for predicting intensities and (iii) adding predicted retention times (RT) from Elude. For ease-of-use, we wrapped these three steps into a pipeline (fasta2speclib), that is included in the MS<sup>2</sup>PIP GitHub Repository.

MS<sup>2</sup>PIP is accessible either through the web server (<https://iomics.ugent.be/MS2PIP/>) or via a local installation ([https://github.com/compomics/MS2PIP\\_c/](https://github.com/compomics/MS2PIP_c/)). A local installation is required to use the fasta2speclib pipeline. Here, MS<sup>2</sup>PIP (version 20190130) was downloaded and installed from the MS<sup>2</sup>PIP GitHub repository, as described in the extended install instructions. For RT prediction, we employed Elude version 3.02, which is available from the Percolator GitHub repository (<https://github.com/percolator/percolator/releases>).

Briefly, the fasta2speclib pipeline makes use of Biopython to read the FASTA and uses Pyteomics for the *in silico* digestion of the protein sequences<sup>[9,10]</sup>. Next, redundant peptides and peptides not meeting the peptide length and precursor mass restrictions are removed from the peptide list. Following this step, all combinations of the requested charge states and modifications are added. Predicted spectra and RTs are then generated for all peptide-charge-modification combinations using MS<sup>2</sup>PIP and Elude. Finally, the results are written to a spectral library file (.msp, .mgf or .csv). Depending on the computational resources, a full human proteome can be predicted in just a few hours.

The fasta2speclib pipeline can be called through the command line interface as follows:

```
Python "fasta2speclib.py" [-h] [-o OUTPUT_FILENAME] [-c CONFIG_FILENAME] "fasta_filename"
```

The results presented in this manuscript were generated by predicting a spectrum for every 2+ and 3+ tryptic peptide in the aforementioned FASTA, using the pre-trained MS<sup>2</sup>PIP Orbitrap-HCD model and the Elude RT. These models are described in more detail under "Prediction models". Only tryptic peptides with a minimum length of 7 amino acid residues and a maximum precursor mass of 5000 Da were considered. Carbamidomethylation and oxidation were set as respectively fixed and variable modification, and two missed cleavages were allowed. The *in silico* spectral library was exported to an .msp file containing 3.3M precursors (between 400 – 1000 *m/z*). In the current version of MS<sup>2</sup>PIP (v20190624) the RT from Elude is automatically converted into minutes and written on a separate line in the .msp file. These predictions were performed on a Linux operated machine (Intel Xeon CPU X5670, 24 processors, 40 GB RAM) and took four hours.

## DIA

DIA libraries, called chromatogram libraries, are generated by interrogating narrow window DIA data with any of the above source libraries. Details are described under "DIA data analysis: EncyclopeDIA".

## RAW file processing

We used the publicly available dataset of the EncyclopeDIA article (MassIVE MSV000082805) of the HeLa S3 lysates to assay the different routes in the DIAMond DIAGram (**Figure 1A, boxes**). The three wide window DIA replicate runs were acquired with 25 overlapping 24 *m/z* windows and the staggered 4 *m/z* narrow window DIA data comprises six gas phase fractions (GPF) of 100 *m/z* each, together covering a 400 - 1000 *m/z* mass range. Following peak picking, these runs were demultiplexed into 12 *m/z* (wide DIA) and 2 *m/z* (narrow DIA) windows, respectively, and converted into mzML output files by MSConvertGUI with following parameters<sup>[11,12]</sup>:

```
Peak picking: Vendor specific algorithms (algorithms available for all vendors, except Waters)
Demultiplexing: overlap only with a mass error of 10 ppm
```

## DIA data analysis: EncyclopeDIA

We downloaded EncyclopeDIA from bitbucket (<https://bitbucket.org/searleb/EncyclopeDIA/downloads/?tab=downloads>) (version 0.8.2, 2019-05-21). EncyclopeDIA is a Java application developed to perform narrow- and wide window DIA data analysis. The application can be run on all three major

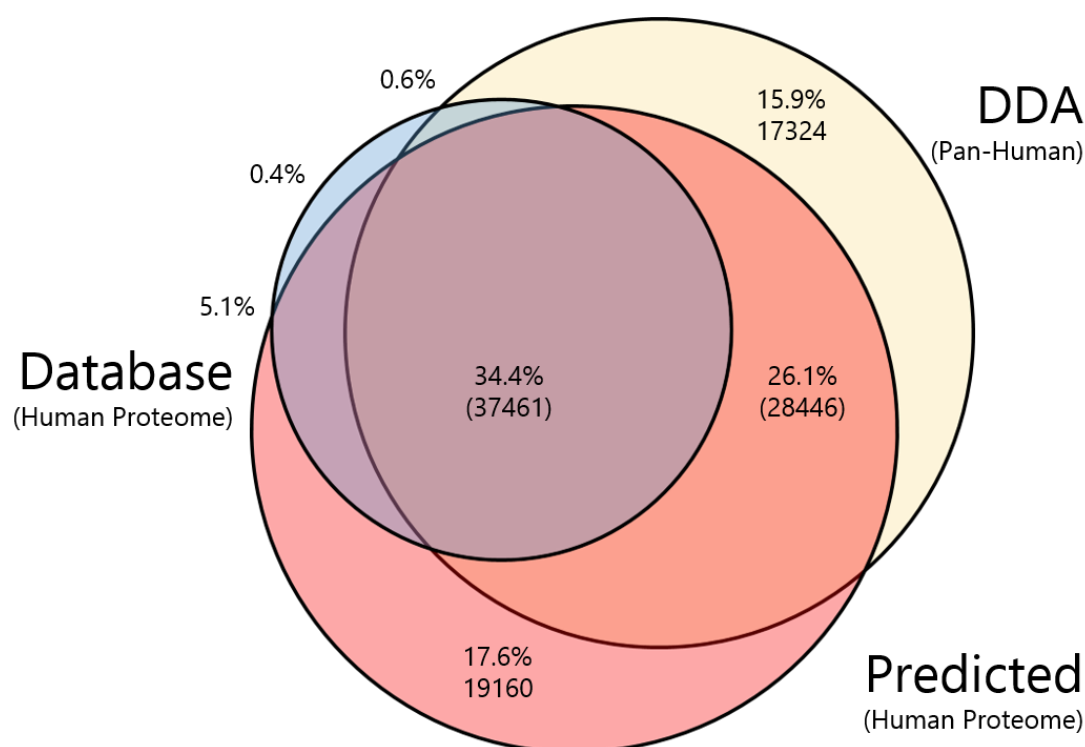
operating systems (Windows, Mac and Linux), but in this project it was used on a Windows 7 operating system (Lenovo Thinkstation, Intel Xeon E5-2620 24 processors, 128 GB ram). EncyclopeDIA was operated through the graphical user interface but also comes with a command-line interface.

For applying EncyclopeDIA on predicted spectral libraries, the .msp file is first converted into a .dlib file using the conversion tool embedded in EncyclopeDIA. EncyclopeDIA also allows the conversion of other spectrum library formats into .dlib files.

General settings in EncyclopeDIA applied to all searches in this project are as follows:

**Background:** Human\_iRT.fasta  
**Target/Decoy approach:** Normal Target/Decoy  
**Data Acquisition Type:** Non-Overlapping DIA (the narrow window DIA was already deconvoluted by MS Convert, therefore EncyclopeDIA does not need to perform an extra deconvolution. )  
**Enzyme:** Trypsin  
**Fragmentation:** CID/HCD (b- and y- fragments)  
**Precursor/Fragment/Library Mass Tolerance:** 10.0 ppm  
**Percolator Version:** v3-01  
**Number of Quantitative Ions:** 5  
**Minimum number of Quantitative Ions:** 3  
**Number of Cores:** 24 (depending on the number of CPU cores you allow/have available)

To allow direct comparison of all six routes of the DIAMond DIAGram, all libraries were trimmed upfront to retain only peptides in the 400 - 1000 *m/z* mass range. For the Pan-Human DDA library this results in 194k precursors, all charge states still included. Approximately 95% of the identified peptides on the wide window DIA were 2+ and 3+ and the other charge states were manually removed from the result file for comparison. The FASTA search was performed using Walnut, considering 2+ and 3+ precursors only. Finally, a third library was predicted by MS<sup>2</sup>PIP using the same FASTA. All three source libraries were separately used to detect peptides directly in the triplicate wide window HeLa DIA runs (**Figure 1Aa-c**). When the three source libraries were used to search the narrow window DIA data (**Figure 1Ad, 1Ae, 1Af**), this resulted in three DIA-based chromatogram libraries (.elib) of size 88k (DDA), 47k (FASTA) and the 95k (Predicted) peptides, respectively. In **Figure S4**, the overlap in peptide sequence is shown between the three chromatogram libraries. Subsequently, all three .elibs were used to search the wide window DIA data with the above parameters.



**Figure S4. Overlap in peptide detection between all chromatogram libraries.** A Venn-diagram showing the overlap in peptide sequence detections between the three DIA-based (DDA, Database and Predicted) chromatogram libraries.

**Figure 1B** depicts the number of detected peptides in each replicate as reported by EncyclopeDIA. Additionally, the peptide quantification reports were exported as .txt files and peptide sequences with at least 3 transitions and non-zero intensities in all three wide window DIA samples were selected. These are represented as the shaded portion of the bar chart. Indeed, in most settings, only confident peptides that can be quantified with robust statistics and are detected in (almost) all runs, are useful. These recurring peptides equally have more robust FDR control. For this reason, we choose to focus only on these confident peptides in Figure 1C, as depicted in the figure caption. Note that the portion of unique peptides between robust detections in Pan-Human and predicted wide window DIA is considerably lower than in the chromatogram libraries that are intrinsically representing single detection. It would be interesting to investigate what the contribution of false detections is herein. All log and result files of the searches were exported for future reference and are available on our GitHub repository.

## FDR assessment by entrapment

We validated the theoretical FDR from the target-decoy approach during chromatogram library building by performing an entrapment experiment with *Pyrococcus furiosus*. In short, this is a way to additionally validate the target-decoy FDR estimation<sup>[13]</sup>. Only peptides between 400 - 1000  $m/z$  were considered and each source library requires a different *P. furiosus* input:

- A public *P. furiosus* dataset acquired on an LTQ-Orbitrap Velos (Thermo Fisher Scientific, Massachusetts, US) was used to supplement the Pan-Human DDA library (ProteomeXchange with identifier PXD001077)<sup>[13]</sup>. Database searching was performed on the resulting .mgf file with Mascot Daemon (version 2.6.1) using following search parameters: a maximum of one missed cleavage, peptide charges 2+ to 4+, peptide mass tolerance of 10 ppm, fragment ion tolerance of 0.5 Da, carbamidomethylation of Cysteine as fixed modification and oxidation of Methionine as variable modification. The resulting .DAT file was parsed into a .BLIB using the Skyline built-

in tool BiblioSpec. The .BLIB file was parsed by EncyclopeDIA into a .dlib file. Finally, the resulting .dlib file (5.5k unique precursors) was combined with the already existing Pan-Human .dlib file of 194k peptides using EncyclopeDIA.

- For the FASTA database, we concatenated our FASTA with all 2052 *P. furiosus* UniProt entries (downloaded on June 13, 2018). Walnut parameters for library-free searching were set as described above, meaning that only 2+ and 3+ peptides without any variable modifications were considered. This translates into 168k *P. furiosus* precursors.
- For the predicted library, we converted this FASTA into a predicted *P. furiosus* spectral library using the MS<sup>2</sup>PIP Orbitrap-HCD model and our Elude RT model. Every 2+ and 3+ tryptic peptide in the proteome was predicted, with carbamidomethylation of Cysteine, and oxidation of Methionine set as respectively fixed and variable modifications. The *P. furiosus* .msp (224k precursors) was concatenated to the Human predicted .msp in EncyclopeDIA.

As decoys are generated by EncyclopeDIA, these were also appended for the *P. furiosus* proteins. All three source libraries were employed for searching the narrow window DIA data, i.e. to create a DIA-based chromatogram library. The *P. furiosus* fraction of the libraries was  $\frac{5.5k}{194k + 5.5k} \approx 3\%$ ,  $\frac{168k}{2.4M + 168k} \approx 6\%$  and  $\frac{224k}{3.3M + 224k} \approx 6\%$  respectively. To account for this differential decoy fraction, the number of *P. furiosus* detections is multiplied by the inverse of their weights, using the following formula:

$$FDR = \frac{\#PyrococcusPeptides}{\#Targets} \cdot Decoyfraction\ correction$$

This corresponds to  $\frac{56}{90k} \cdot \frac{194k + 5.5k}{5.5k} \approx 2\%$  for the DDA source library,  $\frac{19}{46k} \cdot \frac{2.4M + 168k}{168k} \approx 1\%$  for the FASTA source library and  $\frac{64}{94k} \cdot \frac{3.3M + 224k}{224k} \approx 1\%$  for the predicted source library. Note that the number of detected peptides (#targets) is slightly lower than the chromatogram libraries created without *P. furiosus* peptides (see main text). This corroborates the fact that increasing the number of false targets increases the statistical burden and thus number of false negatives, reducing the sensitivity of detection.

In the manuscript we claim the applicability of other deep learning predictors (e.g. DeepMass, Prosit) as an alternative to MS<sup>2</sup>PIP predicted libraries. To validate this claim we cloned the publicly available Pan-Human library using the Prosit webtool which is available from <https://www.proteomicsdb.org/prosit/>. Peptides containing more than 30 amino acids or with a charge state higher than 7 were manually removed from the list as this is required by Prosit. Normalized collision energy (NCE) was assumed to be 33 for all peptides. A similar clone of the Pan-Human library was made with the MS<sup>2</sup>PIP webtool using the pre-trained HCD model. After MS<sup>2</sup> peak intensity prediction, measured iRT values were parsed into both predicted libraries to remove the effect of retention time. Afterwards, the narrow window HeLa DIA data was searched against all three source libraries (Pan-Human, Prosit Clone and MS<sup>2</sup>PIP clone) separately using the settings described earlier in paragraph *DIA data analysis: EncyclopeDIA*. The results of these three searches were exported as the DDA, MS<sup>2</sup>PIP and Prosit chromatogram library, respectively. Next, three wide window HeLa DIA runs were searched with the three chromatogram libraries separately using the same settings as earlier. Again, the results were exported for further processing. The source and chromatogram libraries were converted to an OpenSWATH tsv by EncyclopeDIA, as this simplified parsing of the data. In accordance with the DIAMOND DIAGRAM we calculated PCCs for each narrow and wide window DIA experimental spectrum and its DDA, MS<sup>2</sup>PIP and Prosit source and chromatogram spectrum. Only peptides containing at least 5 transitions were considered and y1 ions were omitted.

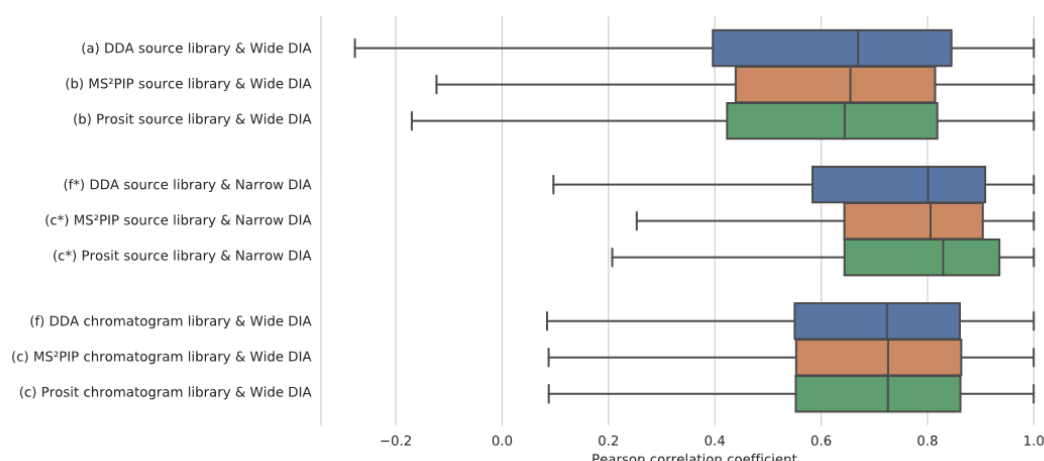


Figure S5. Boxplot showing the distribution of Pearson correlation coefficients between the experimental spectra from the Narrow and Wide-Window HeLa DIA data of the EncyclopeDIA article and the source libraries from DDA (a) or MS<sup>2</sup>PIP and ProSIT (b), as well as the chromatogram libraries derived from DDA (f) or MS<sup>2</sup>PIP and ProSIT (c). Letter annotations refer to the pathways in the DIAMOND DIAGRAM (Figure 1). The overlapping boxplots of the three chromatogram libraries in the bottom clearly illustrate that calibration through narrow window DIA eliminates prior differences in (predicted) intensities.

## References

- [1] R. Gabriels, L. Martens, S. Degroeve, *Nucleic Acids Res.* **2019**, DOI 10.1093/nar/gkz299.
- [2] L. Moruz, A. Staes, J. M. Foster, M. Hatzou, E. Timmerman, L. Martens, L. Kall, *Proteomics* **2012**, 12, 1151.
- [3] B. C. Searle, L. K. Pino, J. D. Egerton, Y. S. Ting, R. T. Lawrence, B. X. MacLean, J. Villén, M. J. MacCoss, *Nat. Commun.* **2018**, DOI 10.1038/s41467-018-07454-w.
- [4] G. Rosenberger, C. C. Koh, T. Guo, H. L. Röst, P. Kouvonen, B. C. Collins, M. Heusel, Y. Liu, E. Caron, A. Vichalkovski, M. Faini, O. T. Schubert, P. Faridi, H. A. Ebhardt, M. Matondo, H. Lam, S. L. Bader, D. S. Campbell, E. W. Deutsch, R. L. Moritz, S. Tate, R. Aebersold, *Sci. Data* **2014**, DOI 10.1038/sdata.2014.31.
- [5] S. Degroeve, L. Martens, *Bioinformatics* **2013**, 29, 3199.
- [6] S. Gessulat, T. Schmidt, D. P. Zolg, P. Samaras, K. Schnatbaum, J. Zerweck, T. Knaute, J. Rechenberger, B. Delanghe, A. Huhmer, U. Reimer, H.-C. Ehrlich, S. Aiche, B. Kuster, M. Wilhelm, *Nat. Methods* **2019**, DOI 10.1038/s41592-019-0426-7.
- [7] Y. S. Ting, J. D. Egerton, J. G. Bollinger, B. C. Searle, S. H. Payne, W. S. Noble, M. J. MacCoss, *Nat. Methods* **2017**, 14, 903.
- [8] C. Escher, L. Reiter, B. Maclean, R. Ossola, F. Herzog, J. Chilton, M. J. MacCoss, O. Rinner, *Proteomics* **2012**, 1111.
- [9] P. J. A. Cock, T. Antao, J. T. Chang, B. A. Chapman, C. J. Cox, A. Dalke, I. Friedberg, T. Hamelryck, F. Kauff, *Bioinformatics* **2009**, 25, 1422.
- [10] A. A. Goloborodko, L. I. Levitsky, M. V. Ivanov, M. V. Gorshkov, *Am. Soc. Mass Spectrom.* **2013**, 301.
- [11] P. Mallick, B. Kuster, *Nat Biotechnol* **2010**, 28, 695.



- [12] R. Adusumilli, P. Mallick, in *Proteomics Methods Protoc.* (Eds.: L. Comai, J.E. Katz, P. Mallick), Springer New York, New York, NY, **2017**, pp. 339–368.
- [13] M. Vaudel, J. M. Burkhardt, D. Breiter, R. P. Zahedi, A. Sickmann, L. Martens, *J Proteome Res* **2012**, *11*, 5065.