

SUPPORTING INFORMATION**Advancing top-down analysis of the human proteome using a benchtop quadrupole-Orbitrap mass spectrometer**

Luca Fornelli, Kenneth R. Durbin, Ryan T. Fellers, Bryan P. Early, Joseph B. Greer, Richard D. LeDuc, Philip D. Compton, and Neil L. Kelleher*

Table of contents.

Supplementary Figure S-1. Analytical SDS-PAGE gels run to visualize fractions from a GELFrEE separation of whole cell extracts of human IMR90 fibroblasts. (A), fractions resulting from two distinct lanes on a 10% GELFrEE cartridge were run adjacent to one another on the analytical gel. GELFrEE fractions (Fr) were pooled based on their elution order as indicated by black brackets and these pools were run with different data acquisition strategies according to the study design described in the main text. (B), an 8% GELFrEE cartridge was used for the fractionation of the higher molecular weight portion of the proteome (30-60 kDa). Fractions 2 to 9, used for the mass spectrometry experiments, are enclosed within the black frame.

Supplementary Figure S-2. Mass accuracy of the “medium”-resolution approach to MS¹ (i.e., use of short time-domain transients that create non-isotopically resolved data) as a function of protein molecular weight. To better illustrate how the accuracy of the average mass calculation based on non-isotopically resolved spectra varies based on the protein mass, we introduce the ‘resolution factor’, which is meant to express the FWHM resolution required to isotopically resolve spectra independently from both (i) the m/z position of each peak (as resolution decreases along the m/z axis) and (ii) the charge of a peak (given that higher charge states would require higher resolution for obtaining isotopically resolved spectra). In this way, all peaks

belonging to the charge state envelope of a certain protein should have a similar resolution factor (within experimental error). Panels A and B are based on 22824 proteoforms from the theoretical proteome of *Pseudomonas aeruginosa*. (A), a linear correlation is apparent between corresponding monoisotopic and average masses of theoretical proteoforms (ranging from ~4 to 570 kDa) (B), the difference between average and monoisotopic mass, expressed in ppm (by dividing the “delta mass” by the monoisotopic mass), remains approximately constant throughout the entire mass range, with the average value being 629.5 ppm. (C), experimental results based on short transient (8 ms) measurement of 5 proteins: ubiquitin (8.5 kDa), superoxide dismutase (16 kDa), myoglobin (17 kDa), carbonic anhydrase (29 kDa) and enolase (47 kDa). The y-axis shows the mass difference between monoisotopic (theoretical) and average mass (experimental, measured using the peak apex), expressed in ppm as in panel B; on the x-axis we report the ‘resolution factor’ for the peaks of several charge states of these five proteins, calculated as the spectral resolution expressed as full width at half maximum (FWHM) multiplied by the corresponding charge state. The plot shows that when the molecular weight of a protein approaches ~40 kDa, the peak apex corresponds closely to the position of the average mass, with very low dependence on the considered charge state. For enolase, the average difference between monoisotopic and average mass is 625.3 ppm (calculation based on 33 different charge states). Conversely, for smaller proteins the calculation of the average mass based on the peak apex results in the underestimation of the average mass itself, due to the more pronounced asymmetry of the underlying distribution of isotopomers. Note that a ‘resolution factor’ of 1 indicates that peaks are isotopically resolved (at 50% peak height), whereas larger values describe non-isotopically resolved species. These larger resolution factor values can be interpreted as the fold-increase in resolution needed to obtain isotopically resolved

charge states at FWHM. (D), zoomed-in view on the proteoform PFR20440 displayed in Figure 4. The left panel includes charge states 45-48+, while a detail of charge state 45+ is shown in the inset on the right.

Supplementary Figure S-3. An example of identifying a low-abundance proteoform via SIM marching. (A), the grey column indicates the applied isolation window (width = 3 m/z units). The precursor corresponding to the isolated, fragmented and identified proteoform is highlighted in red in this MS¹ spectrum obtained during an automated SIM march. Metrics for the protein identification and characterization of the proteoform, PFR13634, appear in the inset (upper right of the top panel). (B), graphical fragment map based on a HCD tandem mass spectrum.

Supplementary Figure S-4. Graphical output of STRING gene ontology analysis based on accession numbers of larger proteins identified by experiments run with “medium/high” data acquisition logic. The plot was obtained via the STRING graphic tool available on-line. The three main protein clusters, organized according to the KEGG pathway catalog, are indicated by red circles. The Glycolysis/Gluconeogenesis protein group (pathway ID: 00010) is represented here by 5 gene products; an additional four genes in this cluster are from the related Fructose and Mannose Metabolism pathway (ID: 00051). The gene products identified as members of the Hippo Signaling Pathway can be considered here as part of a larger group of DNA-bound proteins.

Supplementary Figure S-5. Distribution of all 1952 proteoforms identified at a 1% proteoform-level FDR. This Venn diagram includes also the 80 proteoforms that did not map to any of the 393 Accessions identified at 1% FDR cutoff at the protein level. These 80 proteoforms, corresponding to 4.1% of the total, were mostly identified

in data-dependent experiments and their average C-score was 52.6 (versus an average of 89.7 overall).

Supplementary Figure S-6. Correlation between q -values and C-scores for the proteoforms identified by AUTOPILOT high/high experiments at 1% proteoform-level FDR. (A), considering all 990 proteoforms identified in the 15 RAW files generated by AUTOPILOT, there is no clear linearity between the two scores. (B), limiting linear regression to the top 167 proteoforms with $-\text{LOG}(q\text{-value}) > 30$, the coefficient of determination R^2 does not improve significantly relative to that determined for all proteoforms. Note that the analogous graphs for datasets obtain using the two other two data acquisition modes show very similar trends (data not shown).

Supplementary Table S-1. List of identifications including entries (with assigned UniProt accession numbers) and proteoforms (with related Proteoform Record numbers) for data-dependent high/high experiments. (XLSX)

Supplementary Table S-2. List of identifications including entries (with assigned UniProt accession numbers) and proteoforms (with related Proteoform Record numbers) for AUTOPILOT high/high experiments. (XLSX)

Supplementary Table S-3. List of identifications including entries (with assigned UniProt accession numbers) and proteoforms (with related Proteoform Record numbers) for data-dependent medium/high experiments. (XLSX)

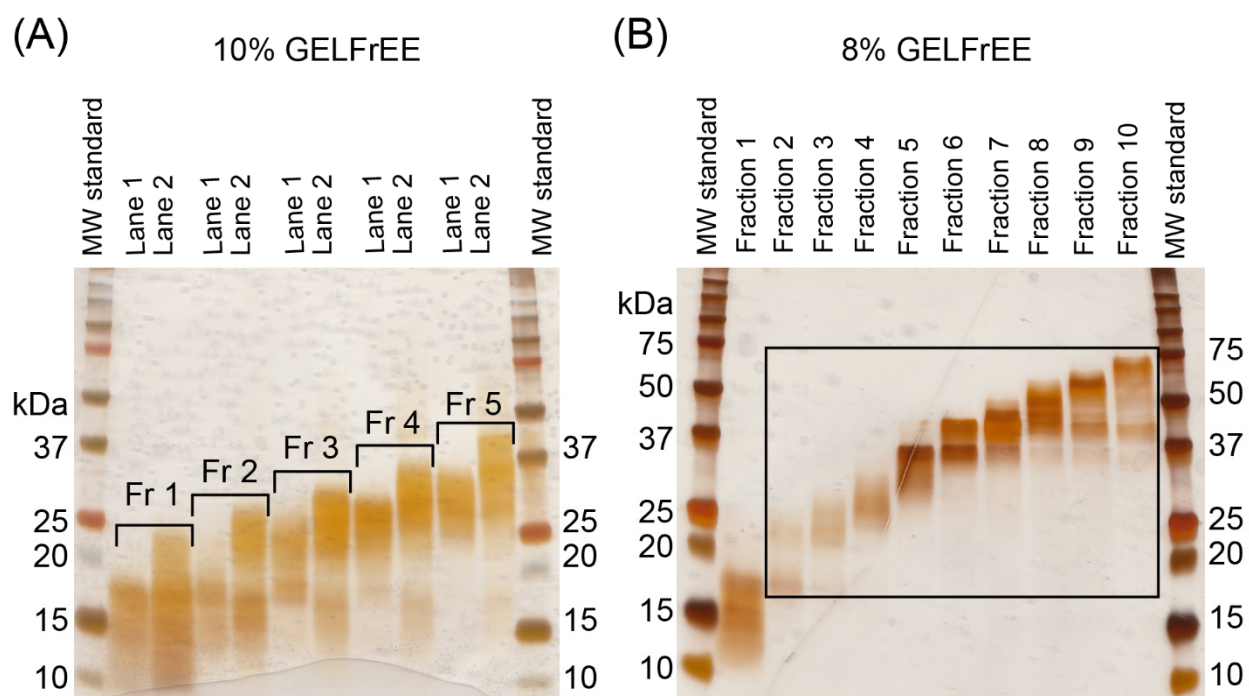
Figure S-1

Figure S-2

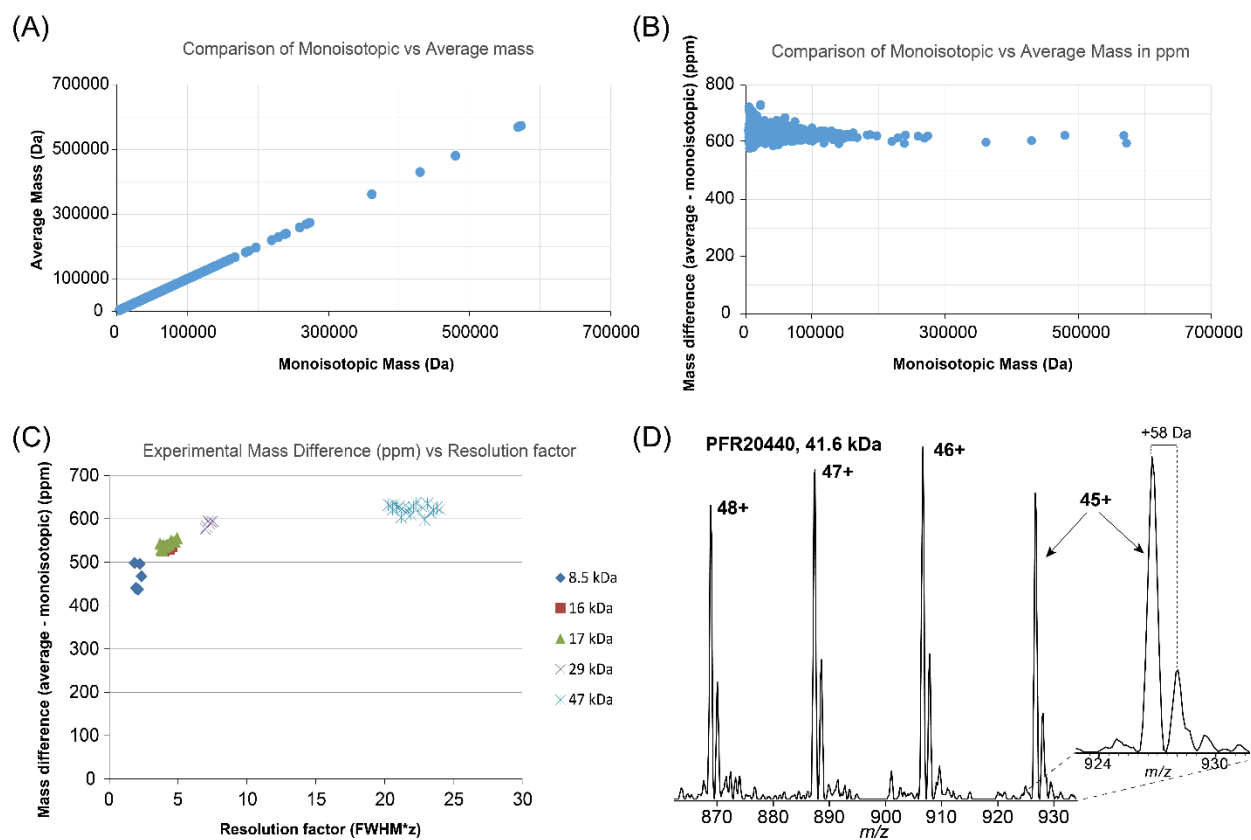


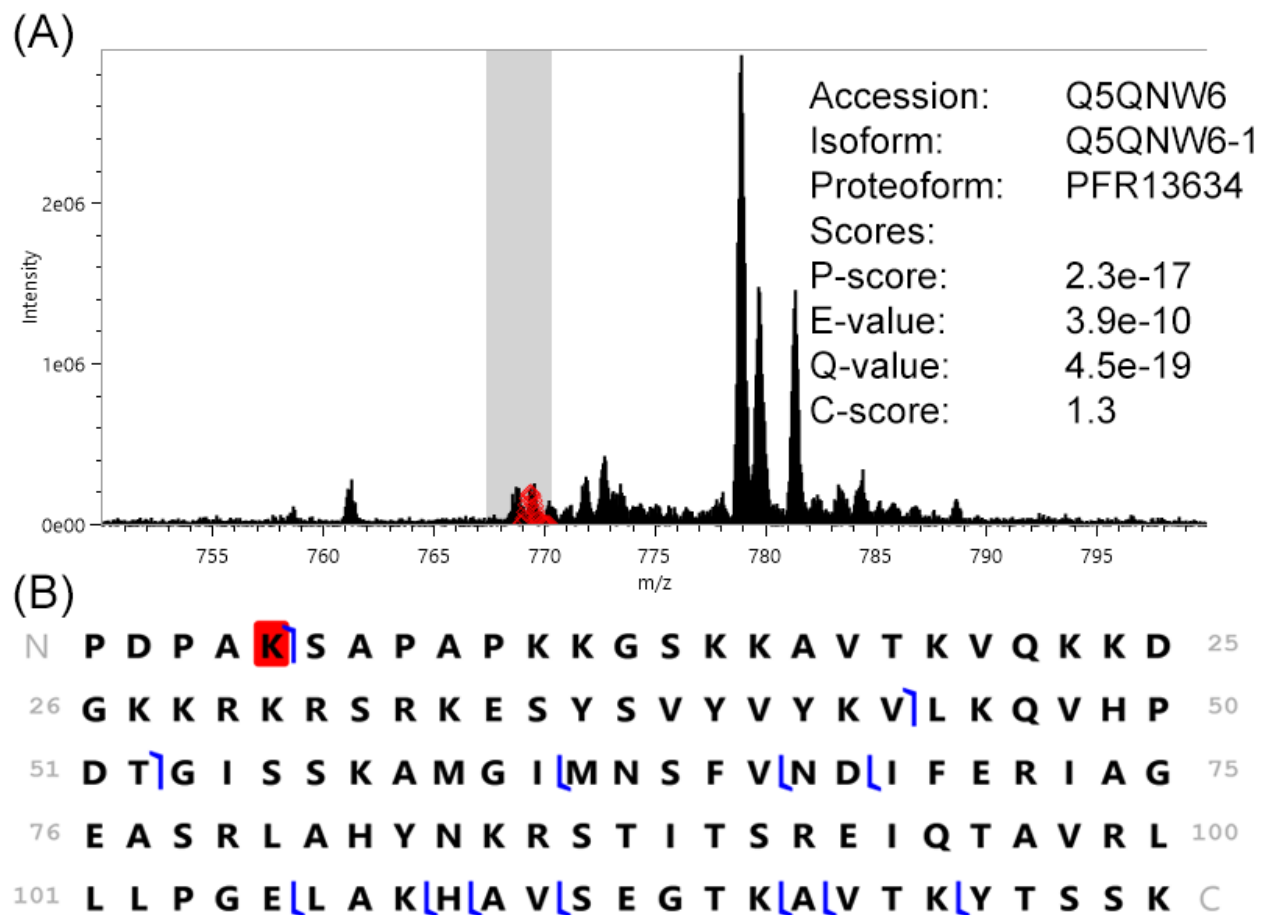
Figure S-3

Figure S-4

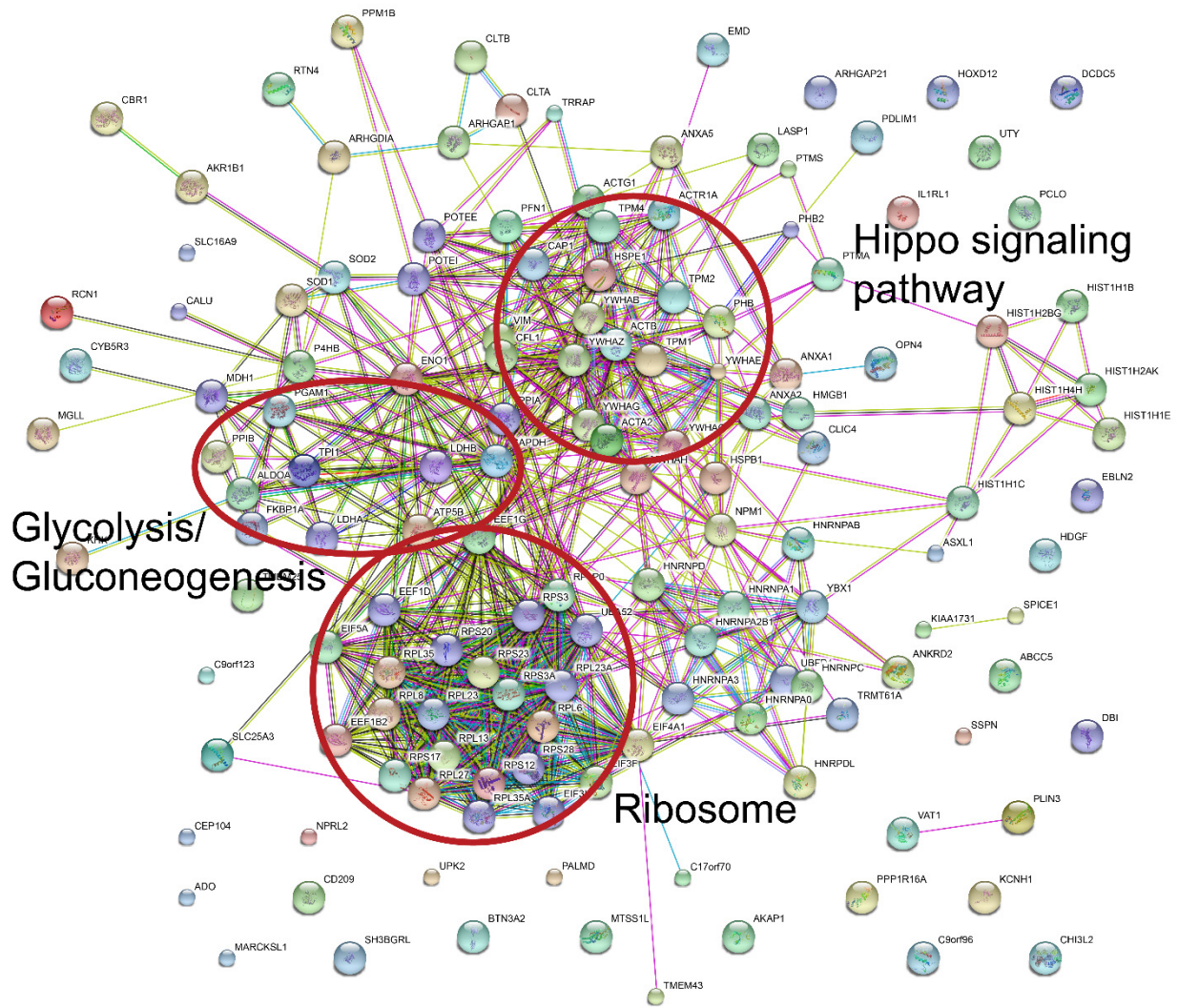


Figure S-5

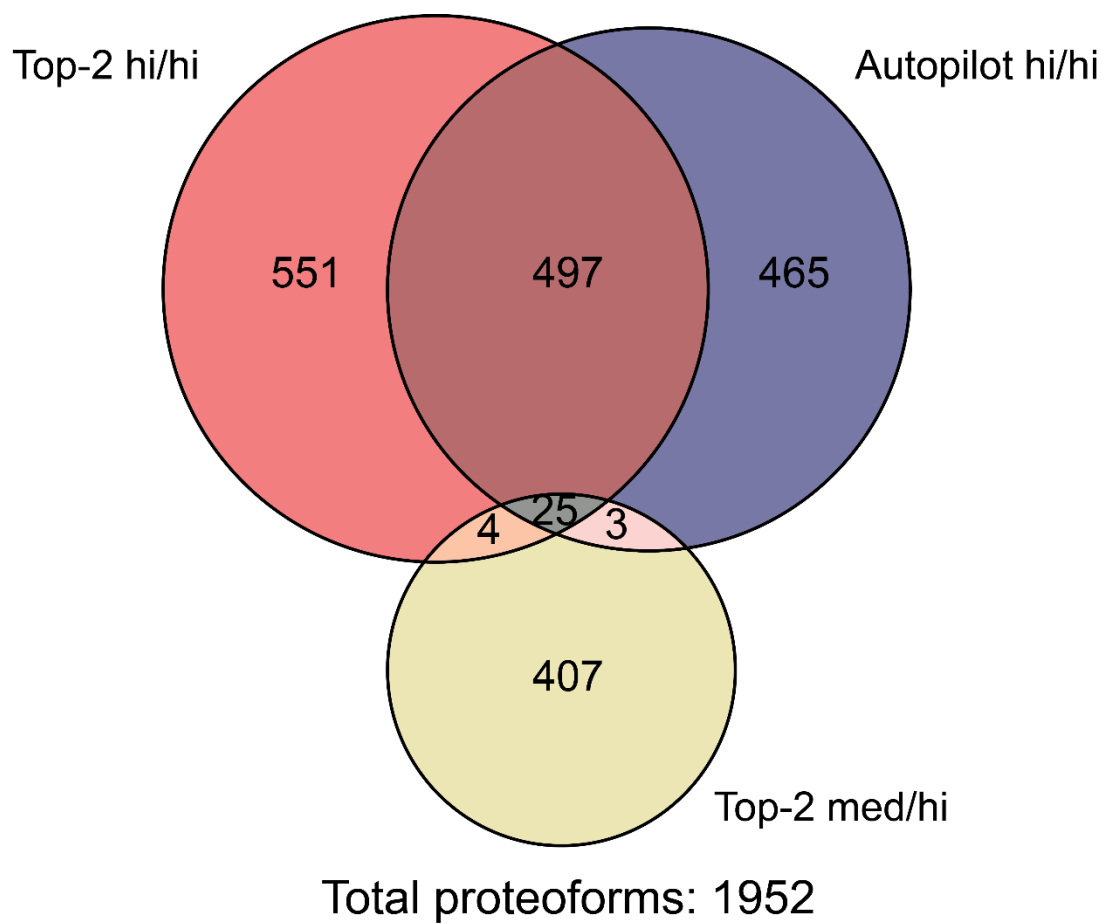


Figure S-6

