



Published in final edited form as:

Nat Methods. 2013 January ; 10(1): 54–56. doi:10.1038/nmeth.2250.

Digestion and depletion of abundant proteins improves proteomic coverage

Bryan R. Fonslow, Benjamin D. Stein, Kristofor J. Webb, Tao Xu, Jeong Choi, Sung Kyu Park, and John R. Yates III*

Department of Chemical Physiology, The Scripps Research Institute, 10550 N. Torrey Pines Rd., La Jolla, CA 92037

Abstract

Two major challenges in proteomics are the large number of proteins and their broad dynamic range within the cell. We exploited the abundance-dependent Michaelis-Menten kinetics of trypsin digestion to selectively digest and deplete abundant proteins with a method we call DigDeAPr. We validated the depletion mechanism with known yeast protein abundances and observed greater than 3-fold improvement in low abundance human protein identification and quantitation metrics. This methodology should be broadly applicable to many organisms, proteases, and proteomic pipelines.

Shotgun proteomics is a widely used approach for biological discovery.^{1,2} An integral part of the process is digestion of complex protein mixtures into peptides using proteases with high sequence specificity. As proteins in cells and tissues often exist in stable higher order structures such as protein complexes or embedded in lipid bilayers, efficient and complete digestion in solution remains a challenge and an area for continuing methodological development. A two-step digestion process for whole cell lysates employing endoprotease Lys-C digestion in 8 M urea, followed by dilution to 2 M urea and digestion with trypsin facilitated the first comprehensive analysis of the yeast proteome.³ Similarly, the use of multiple proteases either in serial or parallel analyses has improved sequence coverage of proteins.⁴⁻⁷ A chaotrope swap strategy using a molecular weight cutoff spin-filter reduces background chemical noise by removing detergent and undigested material.⁸ Aggressive strategies to digest membrane proteins for shotgun proteomics are effective for releasing peptides from the lipid bilayer for identification.^{9,10} Recently, a new protease was developed and introduced for generating larger peptides for middle-down proteomics.¹¹

Users may view, print, copy, download and text and data- mine the content in such documents, for the purposes of academic research, subject always to the full Conditions of use: http://www.nature.com/authors/editorial_policies/license.html#terms

*jyates@scripps.edu, Ph: 858-784-8862, Fax: 858-784-8883.

AUTHOR CONTRIBUTIONS

B.R.F. designed experiments, performed experiments, analyzed data, and wrote the paper. B.D.S. prepared HEK cell lysates and provided conceptual advice. K.J.W. prepared yeast lysates. T.X., J.C., and S.K.P developed software for data analysis. J.R.Y. wrote the manuscript and provided conceptual guidance.

COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

The digestion of complex protein mixtures, however, is often biased by the presence of high abundance proteins. High abundance proteins produce a corresponding excess of tryptic peptides, which can also be further digested by trypsin's endoprotease activity,¹² creating proteolytic background. An excess of high abundance peptides necessitates more chromatographic fractionation, limits dynamic range in the mass spectrometer, and, in turn, biases identification to high abundance proteins in shotgun proteomics.¹³ Common strategies to address the abundance challenge include affinity depletion and enrichment of proteins with antibody arrays or ligand libraries and prefractionation of proteins and peptides.¹⁴ Even with these strategies, the broad dynamic range and the large, varied number of high- and mid-abundance proteins between different sample and cell types present a challenge for the analysis of low abundance proteins. Protease digestion of proteins to peptides can be described by Michaelis-Menten kinetics. The rate of the digestion of a protein (v) is a function of the substrate concentration ($[S]$), the maximum rate of reaction under substrate saturated conditions (V_{max}), and the substrate concentration (K_M) at half V_{max} :

$$v = \frac{V_{max} [S]}{K_M + [S]} \quad (1)$$

Thus, the rate of trypsin digestion of a protein lysate is defined primarily by the protein concentration in relation to the K_M of trypsin. For digestion of a single protein these factors affect the digestion time. For a complex protein mixture, these factors also affect the relative rates at which proteins will be digested based on their relative abundances. We derived an equation (Supplementary Note 1) to describe this phenomenon where the digestion rate of an individual low abundance protein (P_i) is defined by that protein's concentration ($[P_i]$) and the total protein concentration ($[P_T]$):

$$v = \frac{V_{max} [P_i]}{K_M + [P_T]} \quad (2)$$

Briefly, this equation illustrates that the rate of digestion of proteins is dependent on the concentration of each individual protein within the protein lysate and the relationship between the total protein concentration and K_M . In fact, this phenomenon is similar to competitive inhibition of an enzyme. The primary difference is that high abundance "inhibitory" proteins form a peptide product whereas a competitive inhibitor simply dissociates from the enzyme. However, the preference for other proximal tryptic sites on the same high abundance protein likely contributes most to the non-linear inhibitor-like effect.

We exploited these digestion phenomena (Supplementary Note 2) to address both the abundance-dependent digestion of proteins and abundance-dependent sampling of peptides by mass spectrometers, with a method we call DigDeAPr. Briefly, 1 mg of proteome (approximately ten times that typically analyzed by Multidimensional Protein Identification Technology (MudPIT) liquid chromatography – tandem mass spectrometry (LC-MS/MS)) is digested to $85 \pm 10\%$ completion under trypsin- and diffusion-limited conditions in the presence of 2 M urea (Fig. 1a and Online Methods). High abundance proteins are selectively

of more new peptides per run (Supplementary Fig. 4e–f). Notably, there were only minor changes to the quality scores between theoretical and experimental peptide spectra for all peptide abundances (Supplementary Fig. 5a–b). These results indicate that improving identification comprehensiveness with DigDeAPr did not adversely affect the quality or confidence of peptides also easily identified in Control runs. Similarly, changes to spectral counts through DigDeAPr did not adversely affect the reproducibility of spectral count quantitation for proteins with fewer than 100 spectral counts in comparison to Control runs (Fig. 2f). Improvements to this and other protein quantitation metrics such as precursor and fragment ion intensities, precursor ion signal-to-noise ratio (S/N), and chromatographic peak area were found (Fig. 3a–d, Supplementary Fig. 6, and Supplementary Data) and are described further (Supplementary Note 4).

DigDeAPr directly addresses the main challenges of analyzing whole proteomes by selective digestion based on protein abundance to improve the dynamic range of analysis in an unbiased manner (Supplementary Notes 5–6 and Supplementary Figs. 7–8). Because it relies solely on the K_M of a protease and the natural abundance of proteomes, it should be broadly applicable to other organisms, proteases, and proteomic pipelines to improve proteomic sequence coverage. Our method currently uses ten-fold more protein mass than typical comprehensive proteomic analyses, but further optimizations of conditions and the use of higher sensitivity mass spectrometers should make it applicable to mass-limited samples as well. Although we purposely changed the absolute abundance of proteins within a sample using DigDeAPr, the spectral count reproducibility was similar to Control runs, indicating that relative ratios of isotopically labeled protein pairs should remain unchanged as with current protease digestions methodologies. Thus DigDeAPr should also be applicable to quantitative proteomic pipelines using metabolic or chemical labeling strategies.

METHODS

Reagents and Chemicals

Unless otherwise noted all chemicals were purchased from Thermo Fisher Scientific. Deionized water (18.2 M Ω , Barnstead) was used for all preparations.

Growth, isolation, and lysis of log phase yeast

S288C *S. cerevisiae* strain was obtained from ATCC. 250 mL of log phase cells were grown at 30 °C in YPD media (1% bacto-yeast extract, 2% bacto-peptone, 2% dextrose) to an optical density of 0.6 at 600 nm. The culture was harvested by centrifugation at 3,000 \times g for 5 min at 4 °C and washed twice with 10 mL of sterile water. The resulting pellet was snap frozen in liquid nitrogen and placed in –80°C until lysis. The YeastBuster protein extraction reagent (Novagen) was used to lyse cell pellets. The procedure was identical to the manufacturer’s protocol with the addition of 0.5 g of 0.5 mm zirconia beads (RPI Research) per 1 gram of cell pellets. During the 15 min incubation time the lysates were vortexed three times for 30 seconds with one minute rest on ice between cycles. Protein concentration was determined using a non-interfering protein assay kit (Calbiochem).

Cell growth and lysis

Human embryonic kidney cells, HEK 293T, were grown in Dulbecco's Modified Eagle Medium (Mediatech) supplemented with 10% Fetal Bovine Serum Certified (Invitrogen) to 90% confluency in a 5% CO₂ incubator at 37 °C. For collection, plates were washed twice with 20 mL Dulbecco's Phosphate Buffered Saline (-Mg⁺, -Ca⁺) (Invitrogen). Following washing, 1 mL of DPBS containing 1X complete protease inhibitors - EDTA free (Roche) was added to each plate. Cells were lifted from dish surface using Cell Lifter (Corning) and collected into 1.7 mL microcentrifuge tube. Cells were lysed using a probe sonicator at 4 °C, where three cycles of 10 pulses were utilized per sample with 30 seconds on ice between each pulse cycle to offset heating. Lysates were centrifuged at 145,000 × g for 45 minutes. The supernatant was collected as the soluble fraction and used for all subsequent experiments.

Digestion and depletion of abundant proteins

Proteins (~1 mg) were digestion depleted by first denaturing and reducing in 250 μL 8 M urea, 100 mM Tris(hydroxyethylamine) pH 8.5, and 5 mM tris(2-carboxyethyl)phosphine for 30 min. Cysteine residues were acetylated with 10 mM iodoacetamide for 15 min in the dark. The sample was diluted to 1 mL (2 M urea) with 100 mM Tris(hydroxyethylamine) pH 8.5. A 20 μL aliquot was taken for protein quantitation. Trypsin (25 ng, Promega) was added at a 25,000:1 protein:protease mass ratio along with CaCl₂ to 1 mM for a 12 hr diffusion-limited digestion at 37 °C. Digests were transferred to regenerated cellulose 10,000 molecular weight cutoff centrifugal filters (Amicon Ultra-4, ULTRACEL 10K, Millipore) and spun at 2.5K × g for 30 min at 4 °C until 100 – 200 μL remained in the filter. A 20 μL aliquot was taken from the flow through for protein quantitation. The cellulose filter was rinsed with 250 μL 8 M urea, 100 mM Tris(hydroxyethylamine) pH 8.5, then diluted to 2 M urea with 750 μL with 100 mM Tris(hydroxyethylamine) pH 8.5. The digest was spun again to 100 – 200 μL. A 20 μL aliquot was taken from the digestion depleted sample for protein quantitation. Protein quantitation was performed in duplicate using BCA analysis (Micro BCA Protein Assay Kit, Pierce) on aliquots taken during digestion and depletion. The protein masses were calculated to ensure mass balance and quantify the extent of digestion depletion using the following equation:

$$m_{\text{protein, total}} = m_{\text{protein, depleted}} + m_{\text{peptide, depletion}} + m_{\text{protein, filter}} \quad (1)$$

where $m_{\text{peptide, total}}$ is protein mass before digestion depletion, $m_{\text{protein, depleted}}$ is protein mass after digestion depletion, $m_{\text{peptide, depletion}}$ is the peptide mass from the spin-filter flow through, and $m_{\text{peptide, filter}}$ is the peptide mass retained on the spin-filter membrane. Complete protein digestion of digestion depleted samples were continued by transferring the remaining protein solution (100 – 200 μL) to a centrifuge tube, washing the spin-filter membrane twice with 50 μL 8 M urea - 100 mM Tris(hydroxyethylamine) pH 8.5, diluting the protein solution to 2 M urea - 100 mM Tris(hydroxyethylamine) pH 8.5, adding 2 μg trypsin and CaCl₂ to 1 mM for an overnight digestion at 37 °C. Peptides were stored at -80 °C until the day of analysis. On the day of analysis peptide samples were acidified to 5% formic acid and spun at 18,000 × g.

Control protein digestion

Proteins (~100 µg) were digested by first denaturing and reducing in 60 µL 8 M urea, 100 mM Tris(hydroxyethylamine) pH 8.5, and 5 mM tris(2- carboxyethyl)phosphine for 30 min. Cysteine residues were acetylated with 10 mM iodoacetamide for 15 min in the dark. The sample was diluted to 2 M urea with 100 mM Tris(hydroxyethylamine) pH 8.5. Trypsin (2 µg as 0.5 µg/µL) was added at a 1:100 protease:protein ratio along with CaCl₂ to 1 mM for an overnight digestion at 37 °C. Peptides were stored at –80 °C until the day of analysis. On the day of analysis peptide samples were acidified to 5% formic acid and spun at 18,000 × g.

Multidimensional Protein Identification Technology (MudPIT) analysis

Capillary columns were prepared in-house for LC-MS/MS analysis from particle slurries in methanol. An analytical RPLC column was generated by pulling a 100 µm ID/360 µm OD capillary (Polymicro Technologies, Inc) to 5 µm ID tip. Reverse phase particles (Jupiter C18, 4 µm dia., 90 Å pores, Phenomenex) were packed directly into the pulled column at 800 psi until 15 cm long. The column was further packed, washed, and equilibrated at 100 bar with buffer B followed by buffer A. A MudPIT trapping column was prepared by creating a Kasil frit at one end of an undeactivated 250 µm ID/360 µm OD capillary (Agilent Technologies, Inc.), then successively packed with 2.5 cm strong cation exchange particles (Luna SCX, 5 µm dia., 100 Å pores, Phenomenex) and 2.5 cm reverse phase particles (Aqua C18, 5 µm dia., 125 Å pores, Phenomenex). The Kasil frit was prepared by briefly dipping a 20 cm capillary in well-mixed 300 µL Kasil 1624 (PQ Corporation) and 100 µL formamide, curing at 100 °C for 4 hrs, and cutting the frit to ~2 mm in length. The MudPIT trapping column was equilibrated using buffer A for 15 min at 400 bar. Peptide samples (~100 µg) were loaded onto columns at 400 bar. MudPIT and analytical columns were assembled using a zero-dead volume union (Upchurch Scientific).

LC-MS/MS analysis was performed using an Agilent 1200 HPLC pump and Thermo LTQ-Orbitrap XL using an in-house built electrospray stage. Electrospray was performed directly from the analytical column by applying the ESI voltage at a tee (150 µm ID, Upchurch Scientific) directly downstream of a 1:1000 split flow used to reduce the flow rate to 250 nL/min through the columns. Ten-step MudPIT experiments were performed with consecutive application of 0, 10, 15, 20, 25, 30, 40, 50, 60, 70, 85, and 100% buffer C for 5 min at the beginning of each 2 hr gradient. The repetitive 2 hr gradients were from 100 % buffer A to 60% buffer B over 70 min, up to 100% B over 20 min, held at 100% B for 10 min, then back to 100% A for a 10 min column re-equilibration. HPLC buffers (Honeywell) were 5% acetonitrile 0.1% formic acid (A), 80% acetonitrile 0.1% formic acid (B), and 500 mM ammonium acetate 0.1% formic acid pH 6.0 (C). Precursor scanning in the Orbitrap XL was performed from 300 – 2000 m/z with the following settings, respectively: 5×10^5 target ions, 50 ms maximum ion injection time, and 1 microscan. Data-dependent acquisition of MS/MS spectra with the LTQ on the Orbitrap XL were performed with the following settings: collision-induced dissociation on the 8 most intense ions per precursor scan, 30K automatic gain control target ions, 100 ms maximum injection time, 35% normalized collision energy, and 1 microscan. Dynamic exclusion settings used were as follows: repeat count, 1; repeat duration, 30 second; exclusion list size, 500; and exclusion duration, 60

second. All raw data is available as Thermo. RAW files at <http://fields.scripps.edu/published/DigDeAPr2012/>

Data analysis

Protein and peptide identification and comparison were done with Integrated Proteomics Pipeline (IP2, <http://www.integratedproteomics.com/>). Tandem mass spectra were extracted to MS1 and MS2 files from raw files using RawExtract 1.9.9.¹⁶ MS/MS spectra were searched against a combined UniProtKB/Swiss-Prot and UniProtKB/VarSplic human database with reversed sequences using ProLuCID.¹⁷ Human protein entries were extracted and combined from the complete UniProtKB Swiss-Prot and VarSplic databases downloaded at ftp://ftp.uniprot.org/pub/databases/uniprot/current_release/knowledgebase/complete/ on 11/8/2010. The spectral search space included all fully-, half-, and non-tryptic peptide candidates within a 50 ppm window surrounding the peptide candidate precursor mass. Carbamidomethylation (+57.02146) of cysteine was considered as a static modification. Peptide candidates were filtered to 0.1% FDR and proteins candidates to 1% FDR using DTASelect^{18, 19} with a 10 ppm peptide precursor mass window and statistical consideration of peptide tryptic status and mass accuracy. Spectral count, XCorr, CN and summed fragment ion intensities were extracted from DTASelect results. Precursor intensities and S/N for identified peptides were extracted from MS1 files using in-house software.²⁰ Chromatographic peak areas were extracted with Census.²¹ Protein physicochemical properties were calculated using an in-house script.²² Calculations and log₂ comparisons of protein and peptide spectral counts and peptide XCorr, CN, precursor intensity, S/N, peak area, and fragment ion intensity values were performed using Microsoft Excel (Supplementary Data).

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgments

This project was supported by the National Center for Research Resources (5P41RR011823-17), National Institute of General Medical Sciences (8P41GM103533-17), National Institute of Digestive and Diabetes and Kidney Disease (R01DK074798), National Heart, Lung, and Blood Institute (RFP-NHLBI-HV-10-5), and the National Institute of Mental Health (R01MH067880). We thank Jeffrey N. Savas, Claire M. Delahunty, and Jolene K. Diedrich for comments on the manuscript.

References

1. Cravatt BF, Simon GM, Yates JR 3rd. *Nature*. 2007; 450:991–1000. [PubMed: 18075578]
2. Nilsson T, et al. *Nature methods*. 2010; 7:681–685. [PubMed: 20805795]
3. Washburn MP, Wolters D, Yates JR 3rd. *Nat Biotechnol*. 2001; 19:242–247. [PubMed: 11231557]
4. MacCoss MJ, et al. *Proceedings of the National Academy of Sciences of the United States of America*. 2002; 99:7900–7905. [PubMed: 12060738]
5. Choudhary G, Wu SL, Shieh P, Hancock WS. *Journal of proteome research*. 2003; 2:59–67. [PubMed: 12643544]
6. Swaney DL, Wenger CD, Coon JJ. *Journal of proteome research*. 2010; 9:1323–1329. [PubMed: 20113005]
7. Tran BQ, et al. *Journal of proteome research*. 2011; 10:800–811. [PubMed: 21166477]

8. Manza LL, Stamer SL, Ham AJ, Codreanu SG, Liebler DC. *Proteomics*. 2005; 5:1742–1745. [PubMed: 15761957]
9. Wu CC, MacCoss MJ, Howell KE, Yates JR 3rd. *Nat Biotechnol*. 2003; 21:532–538. [PubMed: 12692561]
10. Blonder J, Chan KC, Issaq HJ, Veenstra TD. *Nature protocols*. 2006; 1:2784–2790. [PubMed: 17406535]
11. Wu C, et al. *Nature methods*. 2012; 9:822–824. [PubMed: 22706673]
12. Picotti P, Aebersold R, Domon B. *Mol Cell Proteomics*. 2007; 6:1589–1598. [PubMed: 17533221]
13. Liu H, Sadygov RG, Yates JR 3rd. *Anal Chem*. 2004; 76:4193–4201. [PubMed: 15253663]
14. Jmeian Y, El Rassi Z. *Electrophoresis*. 2009; 30:249–261. [PubMed: 19101934]
15. Liebler DC, Ham AJ. *Nature methods*. 2009; 6:785. author reply 785–786. [PubMed: 19876013]
16. McDonald WH, et al. *Rapid Commun Mass Spectrom*. 2004; 18:2162–2168. [PubMed: 15317041]
17. Xu T, et al. *Mol Cell Proteomics*. 2006; 5:S174.
18. Tabb DL, McDonald WH, Yates JR 3rd. *Journal of proteome research*. 2002; 1:21–26. [PubMed: 12643522]
19. Cociorva D, DLT, Yates JR. *Curr Protoc Bioinformatics*. Chapter 13(Unit 13–14):2007.
20. Wong CC, Cociorva D, Venable JD, Xu T, Yates JR 3rd. *J Am Soc Mass Spectrom*. 2009; 20:1405–1414. [PubMed: 19467883]
21. Park SK, Venable JD, Xu T, Yates JR 3rd. *Nature methods*. 2008; 5:319–322. [PubMed: 18345006]
22. Fonslow BR, et al. *Journal of proteome research*. 2011; 10:3690–3700. [PubMed: 21702434]

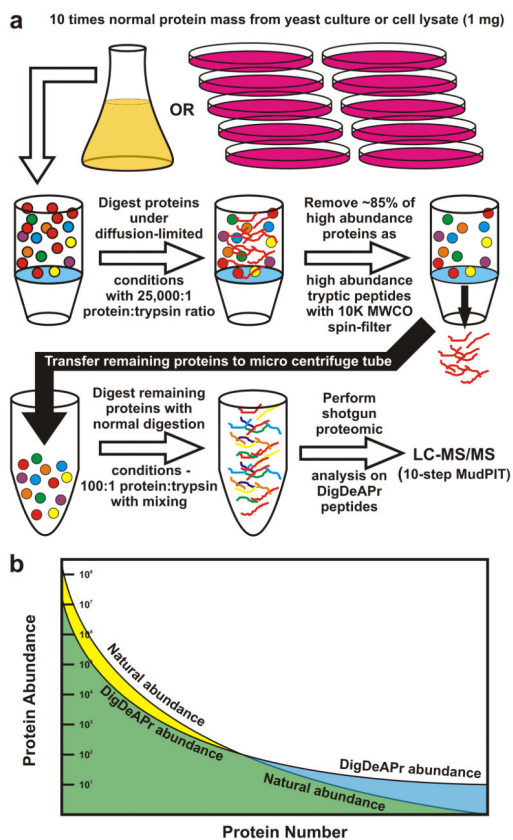


Figure 1. Schematic and description of the DigDeAPr method. **(a)** Ten times the normal mass of proteins are digested under trypsin- and diffusion-limited conditions for removal of abundant proteins as abundant peptides with a MWCO spin-filter. The remaining proteins are digested as normal for LC-MS/MS MudPIT analysis. **(b)** DigDeAPr changes the abundance profile of the proteome by starting with ten times more protein and digesting away ~85% of the most abundant proteins. The higher abundant proteins are preferentially digested by trypsin and depleted as peptides (yellow region), reducing their natural abundance to their DigDeAPr abundance (green region). By using ten times the protein mass to start, the DigDeAPr abundance of low abundance proteins should be ten times higher than their natural abundance (blue region).

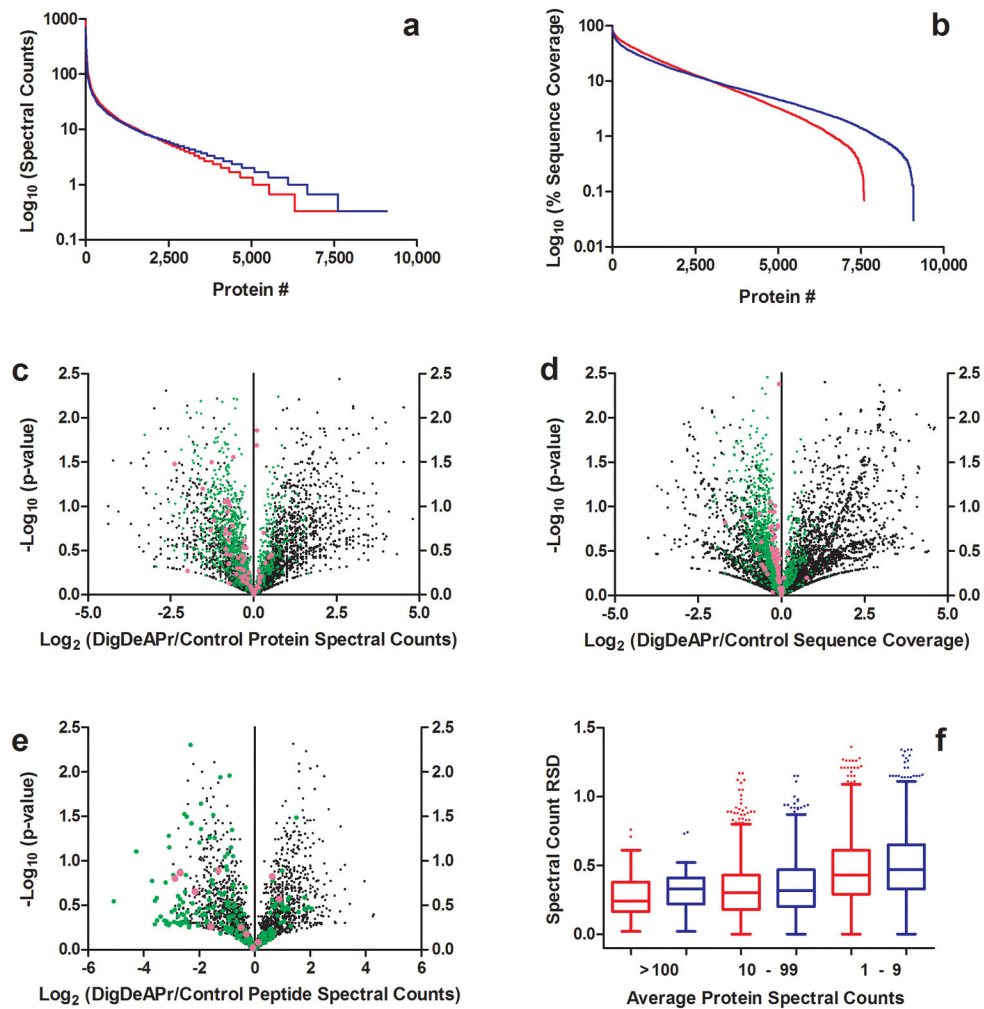
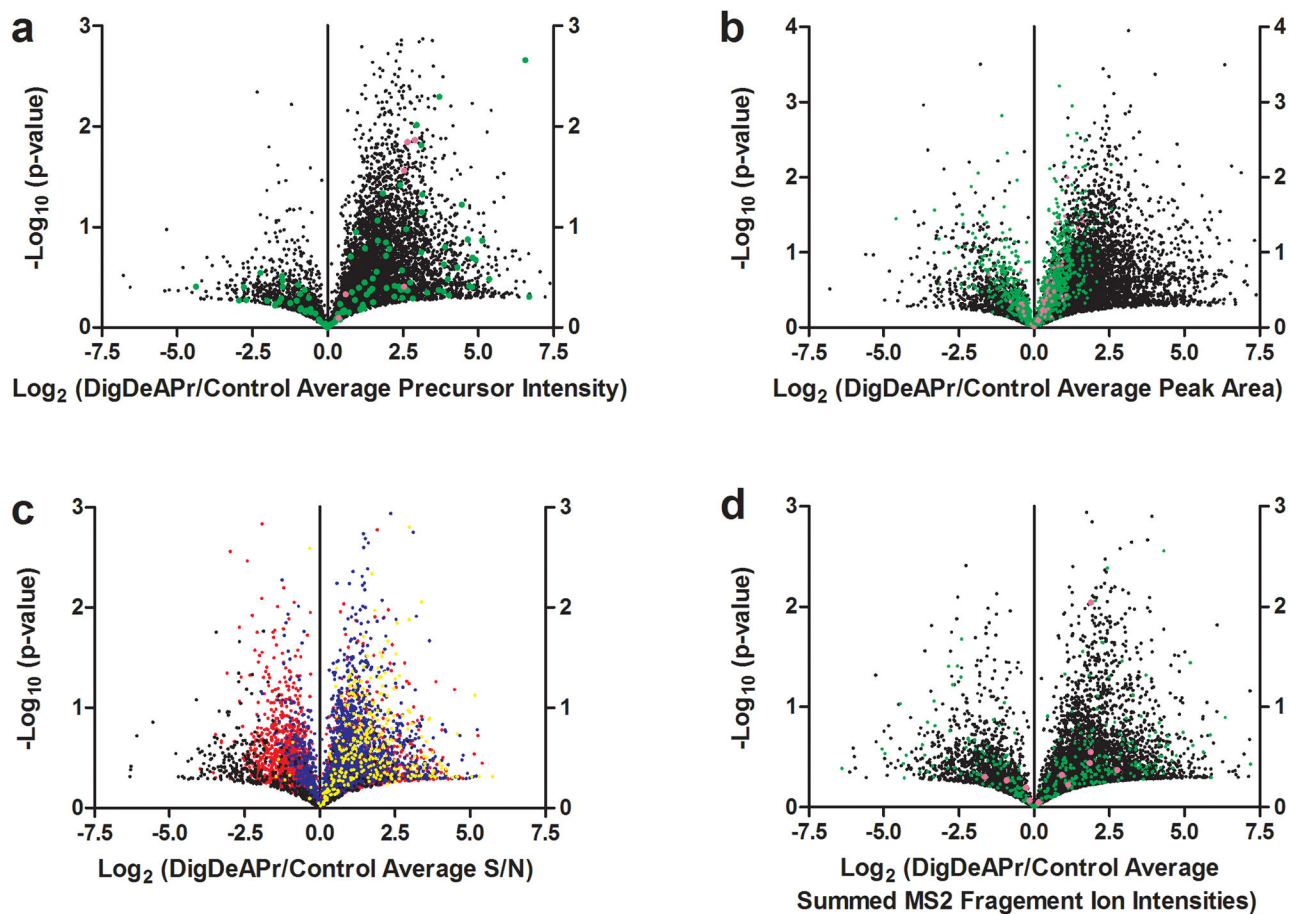


Figure 2. Analysis of proteomic metric improvements from DigDeAPr on HEK cell lysates. Rank abundance plots of proteins based on (a) spectral count and (b) sequence coverage from triplicate DigDeAPr (blue) and triplicate Control (red) runs. Volcano plots of the log₂ ratio of the average from triplicate DigDeAPr (c) protein spectral counts, (d) protein sequence coverage, and (e) peptide spectral counts and the corresponding average from triplicate Control runs plotted against the p-value. Data points are plotted based on average spectral counts from triplicate Control runs: 1–9 spectral counts (black), 10–99 (green), and more than 100 (magenta). (f) Spectral count reproducibility comparison between Control (red) and DigDeAPr (blue) runs based on average protein spectral counts and their relative standard deviation (RSD).

**Figure 3.**

Analysis of proteomic metrics relevant to MS- and MS/MS-based quantitation. Volcano plots of the \log_2 ratio of the average from triplicate DigDeAPr peptide (a) precursor intensity, (b) chromatographic peak area averaged by protein, (c) precursor S/N, and (d) summed MS/MS fragment ion intensities and the corresponding average from triplicate Control runs plotted against the p-value. Data points in (a), (b), and (d) are plotted based on average spectral counts from triplicate Control runs: 1–9 spectral counts (black), 10–99 (green), and more than 100 (magenta). Data points in (c) are plotted based on average S/N from triplicate Control runs: 1–9 (yellow), 10–19 (blue), 20–99 (red), and greater than 100 (black).