

ProteoFlux: a modality-aware framework for downstream analysis of quantitative proteomics data

Dariusz Mollet¹ and Alexander Schmidt¹

¹ Biozentrum, University of Basel, Basel, Switzerland

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#)
- [Repository](#)
- [Archive](#)

Editor: [↗](#)

Submitted: 17 February 2026

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](https://creativecommons.org/licenses/by/4.0/)).

Summary

ProteoFlux is an open-source Python framework for downstream analysis of quantitative proteomics data, operating on outputs generated by upstream peptide and protein identification and quantification tools. It provides a unified, modality-aware analysis workflow supporting protein-centric proteomics, peptide-centric analyses, and phosphoproteomics with optional protein-level covariate adjustment. ProteoFlux emphasizes explicit data semantics, transparent preprocessing, and reproducible statistical modeling, making it well suited for routine analysis and re-analysis in core facility and high-throughput settings. Results are stored in standardized, fully annotated AnnData objects that preserve raw data, intermediate representations, and statistical outputs. By decoupling analysis from visualization, ProteoFlux enables flexible downstream exploration and reporting, with ProteoViewer provided as an example consumer of the analysis-ready outputs.

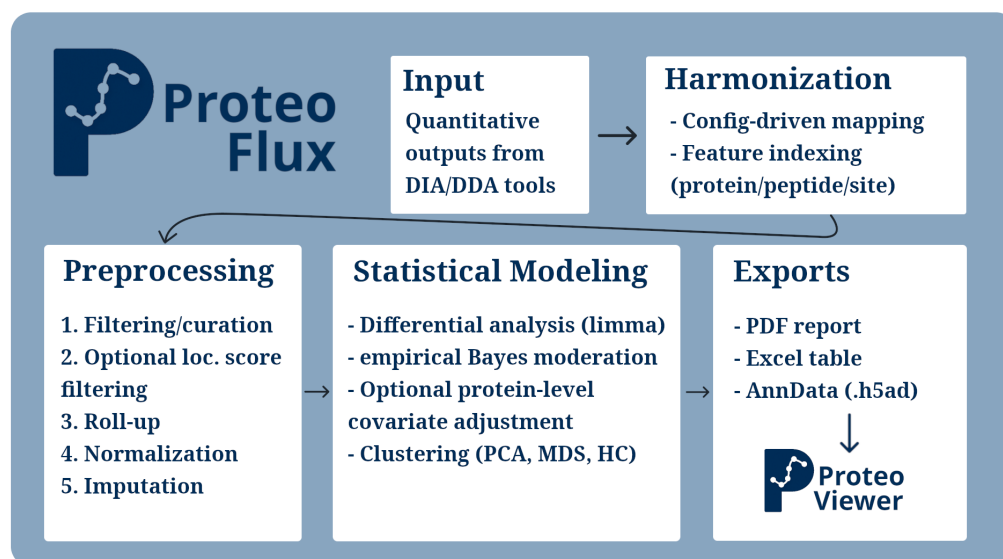


Figure 1: Figure 1: Overview of the ProteoFlux workflow. Quantitative outputs from upstream identification tools are harmonized via configuration into a canonical representation, followed by preprocessing, statistical modeling, and standardized exports for downstream visualization and reporting.

Statement of need

Modern quantitative proteomics experiments routinely generate measurements across multiple biological modalities, including proteins, peptides, and post-translational modifications. While

20 upstream identification and quantification tools have matured, downstream analysis remains
21 fragmented and difficult to standardize across experiment types and data sources.

22 In practice, downstream proteomics analysis is often modality-specific and tightly coupled to
23 individual tools or vendor ecosystems. Preprocessing, statistical modeling, and visualization
24 are frequently intertwined, making analyses hard to inspect, modify, or reproduce consistently.
25 Re-analysis of existing datasets, as well as analysis of new experiments following similar designs,
26 often requires re-implementing custom scripts or adapting workflows originally developed for a
27 different context. This variability complicates comparison across experiments and increases the
28 risk of unintentional analytical bias.

29 These challenges are particularly acute for proteomics core facilities and advanced users, who
30 routinely process large numbers of related experiments under comparable designs. Maintaining
31 consistent preprocessing and statistical assumptions across analyses is critical to limit analytical
32 variance and avoid selective parameter tuning. This requires software that makes preprocessing
33 and modeling decisions explicit, produces stable and inspectable outputs, and clearly separates
34 analysis from visualization and reporting.

35 ProteoFlux was developed to address these needs by providing a unified downstream framework
36 with explicit configuration, deterministic input handling, early error detection, and transparent
37 reporting at each stage of the pipeline.

38 State of the field

39 Several established tools exist for downstream proteomics analysis, including MSstats and
40 related R-based workflows (Kohler et al. 2023; Michalik et al. 2025; Zhu et al. 2020). These
41 tools provide robust statistical modeling frameworks and are widely adopted.

42 However, adapting such workflows across modalities or integrating advanced preprocessing and
43 covariate-aware analyses often requires custom scripting. In many environments, preprocessing,
44 statistical modeling, and visualization are implemented within partially overlapping toolchains,
45 complicating reuse and consistent application across heterogeneous experiments.

46 ProteoFlux was designed to provide a single modality-aware downstream framework that
47 standardizes harmonization, preprocessing, statistical modeling, and output generation across
48 experiment types. Rather than replacing established statistical methodology, it integrates
49 established approaches (e.g., limma-based modeling and directLFQ) within an explicitly
50 configured Python pipeline and exports standardized analysis objects intended for downstream
51 consumers, including ProteoViewer.

52 Software Design

53 Overall Design Philosophy

54 ProteoFlux is a downstream-only framework. It operates on quantitative outputs produced
55 by upstream identification and quantification tools and does not perform peptide or
56 protein identification or signal extraction. This scope avoids overlap with vendor- or
57 search-engine-specific software and focuses on reproducible preprocessing and statistical
58 modeling.

59 Experimental semantics such as conditions, replicates, feature identities, and quantitative signals
60 are defined explicitly through configuration and interpreted deterministically. Inconsistent
61 inputs trigger informative errors or explicit downgrade modes (e.g. pilot-study handling), rather
62 than silent fallbacks. Harmonization, preprocessing, statistical modeling, and outputs are
63 strictly separated.

64 Input Handling and Harmonization

65 ProteoFlux supports long- and wide-format quantitative tables from tools such as Spectronaut,
66 FragPipe, and DIA-NN (Bruderer et al. 2015; Demichev et al. 2020; Yu et al. 2020). Input
67 handling is driven by user-defined configuration files specifying column semantics, experimental
68 layout, and analysis type.

69 Wide-format inputs are reshaped deterministically. When available, external annotation files
70 define authoritative sample metadata; otherwise, identifiers are derived from configuration-
71 defined patterns.

72 For peptide-centric and phosphoproteomics workflows, feature identities are constructed
73 during harmonization. Modified peptide sequences are normalized by removing non-defining
74 variable modifications. Phosphosite identifiers combine parent protein identifiers with absolute
75 modification positions. Localization probabilities are retained for filtering.

76 Phosphoproteomics analysis is performed primarily at the phosphosite level. When the same site
77 appears in multiple peptide contexts, measurements are collapsed using an explicit aggregation
78 rule. For multi-site peptides, an “explode multisite” policy assigns quantitative values to each
79 constituent site, acknowledging non-independence (Müller et al. 2025). Alternatively, a retain
80 mode preserves multi-site peptides as single features.

81 Filtering, quantification, and normalization

82 Filtering is applied at the precursor level prior to aggregation. Quality criteria (e.g., PEP,
83 q-values, precursor evidence counts) are applied conservatively when available. Contaminant
84 filtering and optional low-intensity censoring are recorded explicitly. ProteoFlux does not filter
85 features based on missingness across samples; all detected features are retained to maintain
86 rectangular matrices and enable contrast-consistent tracking.

87 Precursor intensities are aggregated to peptides and proteins using explicit roll-up strategies.
88 Protein quantification can be performed via summation or via directLFQ (Ammar et al. 2023),
89 with user-defined inclusion criteria. All quantification parameters are recorded in metadata.

90 Normalization is applied after log transformation by default, using conservative median-based
91 scaling (Callister et al. 2006; Cox et al. 2014; Karpievitch et al. 2012).

92 Imputation

93 Missing values are retained throughout preprocessing and handled explicitly at the imputation
94 stage. ProteoFlux provides multiple imputation strategies (Shah et al. 2017; Wei et al. 2018).
95 The default method, LC-ConMed, implements a conservative, condition-aware left-censoring
96 strategy for log-scale proteomics data.

97 LC-ConMed performs imputation conditional on experimental grouping. If at least a user-
98 defined minimum number of observations is present (default: one), missing values are sampled
99 around the within-condition median using variance derived from pooled within-group variance
100 across all conditions. Sampling width is constrained and imputed values are clipped to the
101 empirical inter-quantile range of observed values for that condition, preserving condition
102 structure while stabilizing variance for moderated modeling.

103 If insufficient observations are available within a condition, missing values are imputed near
104 a global limit of detection estimated from the lower tail of the intensity distribution. These
105 values are shifted below the detection limit, jittered with small variance, and constrained to
106 remain below the global threshold, ensuring consistent treatment of left-censored features.

107 All intermediate matrices (non-imputed, log-transformed, normalized, imputed) are preserved
108 as separate AnnData layers.

109 Statistical modeling

110 Differential analysis uses a limma-based linear modeling framework via the inmoose Python
111 port (Ritchie et al. 2015; Colange et al. 2025), fitting expression \sim condition models without
112 intercept. All pairwise contrasts are supported. Raw and empirical Bayes–moderated statistics
113 are retained.

114 In phosphoproteomics, optional protein-level covariate adjustment is implemented via a two-
115 stage residualization approach equivalent to ANCOVA without interaction (Wu et al. 2011).
116 Phosphosite intensities are first regressed against matched protein-level abundances with a
117 fixed slope of one, reflecting the assumption of a proportional (1:1) relationship between
118 protein abundance and phosphosite signal; residuals are then analyzed using the same limma
119 framework. Adjustment is suppressed when protein-level measurements are insufficient to
120 support regression. For transparency, the decomposition raw log2FC \sim adjusted log2FC +
121 covariate contribution is stored.

122 Pilot-study scenarios (e.g., single condition or replicate) are detected and formal testing
123 disabled.

124 Fully imputed contrasts are detected on a per-feature, per-contrast basis using the pre-
125 imputation intensity matrices. For such contrasts, statistical outputs are neutralized explicitly
126 (zero fold change, unit p-value) prior to multiple-testing correction, ensuring consistent and
127 interpretable downstream results.

128 Clustering and dimensionality reduction

129 ProteoFlux computes principal component analysis (PCA), multidimensional scaling (MDS), and
130 hierarchical clustering on normalized expression matrices during analysis and stores embeddings
131 and linkage structures directly in the AnnData object. Missingness-based clustering is computed
132 separately to capture structure driven by data completeness.

133 For unsupervised analysis, matrices are fully imputed using a conservative, context-independent
134 left-censoring strategy to provide a contrast-independent representation suitable for quality
135 control and exploratory visualization. This imputation layer is used exclusively for clustering
136 and embedding and does not affect differential testing.

137 Outputs and provenance

138 The primary output is a compressed AnnData (.h5ad) object containing:

- 139 ▪ Experimental design and metadata
- 140
- 141 ▪ Raw, log-transformed, normalized, and imputed matrices
- 142
- 143 ▪ Differential expression results (raw and empirical Bayes–moderated statistics)
- 144
- 145 ▪ Residualized matrices and covariate components where applicable
- 146
- 147 ▪ PCA and MDS embeddings
- 148
- 149 ▪ Hierarchical clustering results
- 150
- 151 ▪ Metadata describing preprocessing, modeling parameters, software versions, and analysis
152 mode

153 Intermediate representations are preserved as separate layers to enable inspection without
154 re-running the pipeline.

155 Additional exports include:

- A structured Excel results file integrating annotations, fold changes, statistical measures, missingness summaries, and intensities, with additional tabs for identification statistics and peptide tables.
- A static PDF report summarizing preprocessing decisions, filtering outcomes, normalization settings, and quality metrics, including hierarchical clustering and volcano plots. The report layout is configurable.

Visualization Decoupling

ProteoFlux writes all preprocessing results, statistical models, embeddings, and derived quantities into a single AnnData object, separating analysis from visualization.

ProteoViewer is a companion application that consumes ProteoFlux outputs without performing additional preprocessing or statistical computation. It can be deployed as a local desktop application or as a server-based web service.

ProteoViewer enables interactive exploration of:

- Quality control and preprocessing diagnostics
- Principal component analysis and multidimensional scaling embeddings
- Hierarchical clustering
- Differential expression results (including volcano plots and peptide trends)
- Feature-level trends across samples, conditions, and contrasts

Because clustering, embeddings, and statistical modeling are precomputed within ProteoFlux, visualization reflects exactly the analyzed data and introduces no additional analytical steps. Results can also be accessed directly in custom Python workflows without ProteoViewer.

Research impact statement

ProteoFlux is used for routine downstream analysis within a proteomics core facility at the University of Basel. The framework supports high-throughput LC-MS experiments spanning protein-centric, peptide-centric, and phosphoproteomics workflows.

Analyses are configured via reusable templates enabling semi-automated generation of experiment-specific configurations while maintaining consistent preprocessing and statistical assumptions.

ProteoFlux-generated outputs are distributed to internal and external collaborators in standardized formats, enabling reproducible re-analysis and downstream visualization without re-running the computational pipeline.

AI Usage Disclosure

Portions of this software and accompanying documentation were developed with assistance from generative AI tools (OpenAI ChatGPT, GPT-5 family models) for tasks including code refactoring suggestions, documentation drafting, CI scaffolding, and language editing of the manuscript.

All AI-assisted outputs were reviewed, validated, and edited by the human authors. The conceptual design, statistical methodology, software architecture, and final technical decisions

199 were made solely by the authors. The authors assume full responsibility for the correctness,
200 originality, licensing, and integrity of all submitted materials.

201 **Acknowledgements**

202 We thank members of the Proteomics Core Facility for feedback and routine testing of the
203 software.

DRAFT