1    Title: ProteoClade: a taxonomic toolkit for multi-species and metaproteomic analysis

2    Authors: Arshag D. Mooradian†, Sjoerd van der Post†, Kristen M. Naegle&, Jason M. Held†‡§*

3

4    †Department of Medicine, Washington University School of Medicine in St. Louis, St. Louis, MO 63110,

5    USA

6    &Department of Biomedical Engineering, University of Virginia, Charlottesville, VA, 22904, USA

7    ‡Department of Anesthesiology, Washington University School of Medicine in St. Louis, St. Louis, MO

8    63110, USA

9    §Siteman Cancer Center, Washington University School of Medicine in St. Louis, St. Louis, MO 63110,

10   USA

11

12   *Corresponding Author: Jason M. Held, Washington University School of Medicine,

13   Campus Box 8076, 660 South Euclid Avenue, Saint Louis, MO 63110. Phone: 314-747-9738. Email:

14   jheld@wustl.edu.

15 **Abstract:**

16 We present ProteoClade, a Python toolkit that performs taxa-specific peptide assignment, protein infer-

17 ence, and quantitation for multi-species proteomics experiments. ProteoClade scales to hundreds of

18 millions of protein sequences, requires minimal computational resources, and is open source, multi-

19 platform, and accessible to non-programmers. We demonstrate its utility for processing quantitative

20 proteomic data derived from patient-derived xenografts and its speed and scalability enable a novel *de*

21 *novo* proteomic workflow for complex microbiota samples.

22

23 **Main Text:**

24 The goal of metaproteomic and multispecies proteomic studies is to characterize the proteomes of sam-

25 ples containing multiple, comingled species, which can provide insight into the complex interactions at

26 the interface between organisms. Proteomic analysis of these samples can quantify thousands of pro-

27 teins from hundreds of species in a single mass spectrometry experiment[1], characterize education of

28 stromal tissue by patient-derived xenografts (PDXs)[2], and extensively characterize the human oral mi-

29 crobiome[3].

30 Metaproteomic and multispecies data analyses depend on the ability to integrate reference protein se-

31 quence databases, taxonomic lineages, *in silico* proteolytic digestion, peptide identification, and quanti-

32 tation. These studies universally perform bottom-up analysis, where proteins are digested into peptides

33 with a protease, and therefore require assignment of peptides to proteins based on their taxonomic

34 specificity. Several software tools provide one or more of these features, but have practical and tech-

35 nical limitations that render them unable to facilitate complete analysis pipelines of quantitative prote-

36 omics data and scale to the rapidly increasing number of available reference protein sequences[4,5]. With

37 regard to annotating peptides to taxa, Unipept is a commonly used taxonomic annotation tool that can

38 provide access to the entire UniProt sequence repository, provides web-based visualizations and a

39    command line interface, and was demonstrated to annotate peptides orders of magnitude faster than a

40    prior UniProt-based application, Peptide Match[6,7]. However, the Unipept database is unavailable for the

41    end-user to customize, is restricted to a fixed set of assumed experimental parameters such as protease,

42    cannot be used with custom protein databases such as those generated by sequencing, and lacks capa-

43    bility for protein quantitation which limits its utility for analyzing many experimental data sets. Addi-

44    tionally, the Unipept database was generated using high performance computing resources which poses

45    a technical challenge as the number of sequences in UniProt grows exponentially[5,8]. Other tools offering

46    more complete metaproteomic pipelines such as MetaProteomeAnalyzer (MPA) support database gen-

47    eration, peptide spectral matching, and taxonomic parsing. However, these tools restrict the user to

48    bundled open-source targeted database search engines which scale poorly when using large reference

49    database sizes such as the entirety of UniProt[9]. MPA also lacks support for post-translational modifica-

50    tions and common MS2-based quantitation approaches, limiting its applicability.

51    To overcome these limitations, we developed ProteoClade, an open-source Python library that enables

52    flexible, rapid, and easy taxon-specific quantification for proteomic experiments. ProteoClade utilizes

53    standard taxonomic and protein sequence repositories and is optimized for large databases. Proteo-

54    Clade is the first tool to 1) enable users to generate and search customizable, *in silico* digested peptide-

55    to-taxa mapped databases that can scale to the entire UniProt database with the optional inclusion of

56    user-specified reference protein sequences; and 2) provide a novel *de novo* workflow to efficiently an-

57    notate peptides sequenced without defining the taxonomic composition *a priori*. Additionally, Proteo-

58    Clade allows the user to choose their preferred commercial or open-source search engine, as well as

59    preserves MS1-, MS2-, and spectral count-based quantitation from the experiment to calculate gene-

60    level, taxon-specific quantitative results. Together, ProteoClade uniquely enables fast, taxon-specific

61    quantitation of database-targeted multispecies experiments as well *de novo* searches enhanced by large

62    databases for complex metaproteomic experiments. The ProteoClade software is freely available at

63    **http://github.com/HeldLab/ProteoClade** .

64    ProteoClade retrieves complete taxonomic lineages from the NCBI, and interfaces directly with the Uni-

65    Prot API to download and concatenate reference protein sequence databases based on the organism IDs

66    and database parameters supplied by the user. To assign peptide-spectral matches from mass spectra

67    search engines to both taxon and gene identifiers, ProteoClade creates a SQLite database (ProteoClade

68    Database: PCDB) by digesting reference proteomes *in silico* with proteolytic parameters that can be cus-

69    tomized according to the user's experimental conditions (**Fig. S1**). The PCDB addresses scalability issues

70    by storing peptide sequences as hashed integer values which compresses the average peptide storage

71    requirement by 62.8%. ProteoClade leverages multiple CPUs for parallel processing to speed both the

72    PCDB creation and search functions, and indexes the database to quickly assign peptides to taxon-

73    specific genes. ProteoClade's implementation is further detailed in the STAR methods.

74    Metaproteomic experiments involving communities of many organisms present substantial computa-

75    tional challenges since more sequence information available to aid in the identification and quantitation

76    of the experimental data requires more computational resources. Thus, it is imperative to efficiently

77    store and parse these data for taxonomic assignment and protein quantitation in light of the dramatic

78    increase in the number of organisms with proteomic annotation.  We evaluated ProteoClade's scalability

79    by creating a tryptic PCDB from the November 2018 release of UniProt containing 140.2 million protein

80    entries (UniProt PCDB). The resultant database contained 10.7 billion peptides, 5.02 billion of which

81    were unique**,** from 1,040,460 organisms. PCDB creation, including indexing, took only 11 hours on mod-

82    est consumer-grade hardware and resulted in a 515GB file (**Fig. S2a**), demonstrating that ProteoClade

83    enables the use of large, customizable peptide databases.

84    We compared the database indexing and taxonomic annotation features of Unipept, MetaProteomeAn-

85    alyzer (MPA), and ProteoClade to highlight ProteoClade's optimizations and improvements over prior

86      taxonomic tools. For peptide database creation and indexing, ProteoClade uses 32-fold less RAM and 6-

87      fold less time than Unipept, and > 60-fold less RAM and > 120 fold less time than MPA (**Fig. S2a, S2b**).

88      We found that ProteoClade annotates experiments at 8.8x the speed of Unipept and preserves quantita-

89      tive information with the ability to sum the peptides' quantitation to the gene level as is common in

90      multispecies experiments, while MPA lacked the ability to annotate peptides outside of database-

91      targeted searches (**Fig. S2c**). ProteoClade's technical optimizations enable taxonomic annotation and

92      quantitation of peptides at a speed and scale that exceed previously used tools.

93      To demonstrate ProteoClade's ability to perform integrated, species-specific quantitation of multi-

94      species samples, we analyzed publicly available TMT-labeled global proteomics data from Patient-

95      derived xenograft (PDX) lines in which six triple-negative breast cancer tumors were each grafted into

96      three mice, but were originally analyzed without considering mouse-specific peptides[10]. PDXs are a bur-

97      geoning, mixed-species model of tumor biology[11] in which the stromal microenvironment of tumors,

98      comprised of fibroblasts, immune-related cells, and vasculature, is originally human but is replaced after

99      several passages by murine cells. Thus, species-specific proteomic analysis of PDX data can simultane-

100     ously and independently characterize the tumor and the invasive murine stroma to examine how tu-

101     mors remodel stromal proteomes in a process known as stromal education[2].

102     ProteoClade was used to create a customized, concatenated human and mouse UniProtKB/Swiss-Prot

103     database for both peptide-spectral matching with MaxQuant and creation of a PCDB for species-specific

104     taxonomic annotation and filtering (**Fig. 1a**). In addition, ProteoClade's quantitation module flexibly al-

105     lows for species-specific summation of TMT reporter ions to each peptide's assigned gene symbol to

106     generate quantitative proteomic maps of the tumor and stromal proteomes.

107     Taxon-specific peptide assignment with a tool such as ProteoClade is important in multi-species samples

108     since peptides with shared amino acid sequences between human and mouse may bias proteomic quan-

109     titation of PDXs. By digesting peptides *in silico*, we found that 71.1% of genes in the human proteome

110     produce tryptic peptides with sequences identical to their murine homolog, and thus most human and

111     mouse proteins are potentially susceptible to quantitative interference by the presence of the other or-

112     ganism. We examined the consequences of normalizing and quantifying the human and mouse compo-

113     nents of the data by comparing a naive informatic assumption that only one organism was present to

114     utilizing ProteoClade for species-specific peptide filtering and protein quantification (**Fig. 1b**). The rela-

115     tive protein abundance varied by more than 1.5-fold for 262 genes in the human-specific data and 891

116     genes in the murine-specific data, highlighting the large bias that the inclusion or exclusion of species-

117     shared peptides can have on the quantitation. Overall, genes share a higher proportion of tryptic pep-

118     tides have larger differences in gene-level quantitation indicating that taxonomic interference has a sub-

119     stantial effect on the resulting data (**Fig. 1b**).

120     ProteoClade greatly facilitated species-specific quantitation of this PDX dataset[10], identifying 2,326 mu-

121     rine proteins in the microenvironment that are significantly altered by patient-derived tumors, the most

122     expansive set of tumor-educated stromal proteins to date. The murine proteins clustered by biological

123     replicate, confirming that tumor-intrinsic factors drive tumors to persistently educate the stromal pro-

124     teome as previously observed (**Fig. 1c**)[2]. ANOVA revealed that stromal education by tumors is wide-

125     spread, with 77.1**%** of stromal proteins differentially regulated by the embedded tumors. Differentially-

126     regulated stromal proteins combined from this analysis and a prior study (n = 3,015 proteins) were en-

127     riched for proteins in the extracellular matrix, cytoskeletal processes, and myeloid-derived immune

128     components which are important contributors to tumor growth and metastasis (**Fig. 1d**)[2,10].

129     Beyond two-species systems, we examined ProteoClade's applicability to large-scale metaproteomics

130     workflows using the entire UniProt repository, which enables taxonomic annotation to all potential or-

131     ganisms in the sample (**Fig. 2a**). One important consideration for metaproteomic analysis is that all se-

132    quence databases are biased by the inclusion and exclusion criteria selected by their curators.

133    Overrepresentation of peptidomes from certain genera and peptide sequence similarity between relat-

134    ed organisms can result in misattributing taxon-specific peptides due to errors in mass spectrometry-

135    based detection and sequencing. We establish that UniProt has a substantial sequence bias from the

136    perspective of bottom-up proteomics experiments by quantifying the number of peptides contained in

137    the UniProt PCDB for each genus as well as the proportion of peptide sequences that were unique (**Fig.**

138    **2b**). By mapping all 5.02 billion unique peptides back to their respective genera, we found that taxonom-

139    ic redundancy from the inclusion of thousands of bacterial strains results in vast overrepresentation of

140    tryptic peptides from several genera, including *Streptomyces, Pseudomonas,* and *Bacillus*. These

141    overrepresented genera are likely to be identified in nearly every genus-specific proteomic analysis re-

142    gardless of their presence or absence in the sample, highlighting the importance of user-customization

143    to database generation. We provide a table of these data to help inform when selecting a narrower da-

144    tabase scope for certain genera in UniProt may be appropriate (**Table S1**).

145    *De novo* proteomic analysis is advantageous for metaproteomics since it does not require pre-specifying

146    which organisms are present, which is a requirement for the database-targeted MS2 search approaches

147    universally utilized in metaproteomics. ProteoClade integrates with *de novo* searches since it can taxo-

148    nomically annotate and quantify *de novo* search results in a unique workflow that can identify organisms

149    and proteins without *a priori* specification, obviating the need for 16s rRNA sequencing or metagenomic

150    assembly for proteomics studies. ProteoClade enables this unique informatic workflow by assigning mil-

151    lions of candidate peptide sequences from *de novo* searched spectra unbiasedly to all possible organ-

152    isms, which ProteoClade makes possible without the use of high performance computing resources (**Fig.**

153    **S2**).

154     We evaluated using ProteoClade for *de novo* peptide assignment of a large oral microbiome proteomic

155     data set[3] previously analyzed with a standard database-targeted approaches by coupling the UniProt

156     PCDB with *de novo* MS2 assignment (**Fig. 2a**)[3]. ProteoClade parsed 8.13 million peptide-spectral match

157     candidates to annotate 1.83 million MS/MS scans in 1.94 hours (261 spectra per second), making this

158     time-efficient even when challenged with large data sets and reference databases. We identified genus-

159     unique peptides for 39 genera from the oral microbiome in this dataset without prior specification of

160     the bacterial taxa present in the sample, despite the fact that our criteria for identification included ge-

161     nus-specificity in the context of more than 72,000 genera compared to the 108 genera present in the

162     Human Oral Microbiome Database used for the original publication (**Fig. 2c**)[12].

163     We next developed an approach to serially check every candidate sequence for each MS2 spectrum until

164     a match in the PCDB was found, and compared these results to an approach which only considers the

165     top scoring candidate for each spectrum. As *de novo* search engines lack a reference database and do

166     not provide false discovery correction for peptide sequencing, we additionally used ProteoClade to gen-

167     erate a reverse 'decoy' sequence PCDB, and annotated the peptide sequence candidates using the pa-

168     rameters we chose for our forward search. We then compared the proportion of annotations to the

169     forward database to the reverse database for each genus, which is the same strategy that is used for

170     most targeted database searches, but is novel to *de novo* searches[13]. The combination of considering

171     multiple candidate amino acid sequences for every spectrum and offsetting false discovery with a decoy

172     database increased the number of identified genus-unique human peptides by 32.6% in the forward

173     PCDB annotations, while we observed no increase in the reverse, decoy PCDB annotations (**Fig. 2c**). This

174     demonstrates that ProteoClade can identify more biologically-valid peptide sequences than would oth-

175     erwise be offered by upstream *de novo* search software. For several bacterial genera, including *Prevotel-*

176     *la* and *Lactobacillus*, we found an increase in both the forward and reverse PCDB annotations, indicating

177     that while the number of peptide candidate sequences increased, those candidates are likely false posi-

178    tives. By comparing all forward and reverse annotations, we observed a threshold (>7 Δforward-reverse

179    genus specific peptides) above which only genera known to have been identified in the human oral cavi-

180    ty are present. We identified human and 14 bacterial genera above this threshold without prior specifi-

181    cation, including *Streptococcus*, which was previously reported as the most abundant genus in the oral

182    microbiome (**Fig. 2d**)[14].

183    ProteoClade enables taxon-specific analysis of targeted database and *de novo* proteomic experiments. It

184    functions on all major operating systems and operates at a speed and scale that enable fast and novel

185    forms of proteomic data processing using consumer hardware. We expect ProteoClade's ease of use,

186    applicability to a broad set of biological model systems such as PDXs and metaproteomics, and its novel

187    integration with *de novo* spectral searches[4] provides a unique and powerful tool to researchers perform-

188    ing quantitative proteomic analysis of multiple mixed species.

189

190    1.    Zhang, X. *et al.* Deep Metaproteomics Approach for the Study of Human Microbiomes. *Anal.*

191          *Chem.* **89**, 9407–9415 (2017).

192    2.    Wang, X. *et al.* Breast tumors educate the proteome of stromal tissue in an individualized but

193          coordinated manner. *Sci. Signal.* **10**, (2017).

194    3.    Grassl, N. *et al.* Ultra-deep and quantitative saliva proteome reveals dynamics of the oral

195          microbiome. *Genome Med.* **8**, 44 (2016).

196    4.    Starr, A. E. *et al.* Proteomic and Metaproteomic Approaches to Understand Host-Microbe

197          Interactions. *Anal. Chem.* **90**, 86–109 (2018).

198    5.    UniProt Consortium. UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Res.* **47**,

199          D506–D515 (2019).

200    6.    Gurdeep Singh, R. *et al.* Unipept 4.0: Functional Analysis of Metaproteome Data. *J. Proteome Res.*

201          **18**, 606–615 (2019).

202    7.    Mesuere, B. *et al.* Unipept web services for metaproteomics analysis. *Bioinformatics* **32**, 1746–

203          1748 (2016).

204    8.    Mesuere, B. *et al.* The Unipept metaproteomics analysis pipeline. *Proteomics* **15**, 1437–1442

205          (2015).

206    9.    Tanca, A. *et al.* Evaluating the Impact of Different Sequence Databases on Metaproteome

207          Analysis: Insights from a Lab-Assembled Microbial Mixture. *PLoS One* **8**, 1–14 (2013).

208    10.   Mundt, F. *et al.* Mass Spectrometry-Based Proteomics Reveals Potential Roles of NEK9 and

209          MAP2K4 in  Resistance to PI3K Inhibition in Triple-Negative Breast Cancers. *Cancer Res.* **78**, 2732–

210          2746 (2018).

211    11.   Hidalgo, M. *et al.* Patient-derived xenograft models: an emerging platform for translational

212          cancer research. *Cancer Discov.* **4**, 998–1013 (2014).

213    12.   Chen, T. *et al.* The Human Oral Microbiome Database: a web accessible resource for investigating

214          oral microbe taxonomic and genomic information. *Database (Oxford).* **2010**, baq013 (2010).

215    13.   Elias, J. E. & Gygi, S. P. Target-decoy search strategy for increased confidence in large-scale

216          protein identifications by mass spectrometry. *Nat. Methods* **4**, 207–214 (2007).

217    14.   Bik, E. M. *et al.* Bacterial diversity in the oral cavity of 10 healthy individuals. *Isme J.* **4**, 962

218          (2010).

1    **Materials and Methods**

2

3    ProteoClade implementation and testing

4    ProteoClade was developed using Python 3.6 (64-bit) and has been tested on Windows 7, MacOS High

5    Sierra (10.13), and Ubuntu Linux (18.04).  All functions have been run and timed using a computer with

6    an Intel i7-2600k processor, 16 gigabytes of RAM, and a Samsung 860 EVO solid state drive. Taxonomic

7    rankings are retrieved directly from the NCBI FTP servers and are assembled as a pickled Python

8    dictionary object which is stored in RAM and allows the software to rapidly assign higher level

9    taxonomies from organism taxonomy cross-reference (OX) IDs with $O(1)$ time complexity. Protein

10   sequences are retrieved using UniProt's REST API to enable the selection of specific combinations of taxa

11   or database sources (SwissProt, TrEMBL, and/or Reference) using the OX IDs.

12

13   The ProteoClade Database (PCDB) and all analyses we performed used the default digest parameters of

14   peptides that range from 7-55 amino acids in length, trypsin with C-terminal proline allowed (trypsin/p),

15   protein N-terminal methionine excision, leucine-isoleucine interconversion, and all combinations up to

16   two missed cleavages; however, we enable the user to select from any variation of these rules, including

17   alternative built-in or custom proteases and alternative ranges of amino acid lengths. Peptides are

18   stored as hashed integers in the PCDB in order to reduce the on-disk storage size.  The PCDB creation

19   process was optimized with the use of multiprocessing for computing the hashed peptide integers and

20   integrating them into the database. Similarly, peptide indexing makes use of multiple threads which

21   reduces the total time required for database creation, and is saved as a separate step at the end of

22   database creation to allow the database creation time to scale linearly with the number of protein

23   sequences inserted into it.

24

1

25    Peptide annotation was implemented as a multithreaded process and pushes the bottleneck of

26    ProteoClade's performance to the input-output operations per second (IOPS) of the harddrive.

27    Annotation can be performed at any combination of taxa the user requests. Quantitative information

28    from search engine outputs is preserved during annotation, and is assigned to gene symbols on the basis

29    of filters supplied by the user. These filters include specifying a taxonomic level to determine the

30    taxonomic uniqueness of a sequence, inclusion and exclusion lists for taxa that the use may want to

31    include or exclude, and a default taxon that can be used to assign quantitative information in the event

32    that multiple taxa are present but only one is of interest to the user, all of which allow the user to create

33    taxon-specific data sets with flexible criteria. Detailed documentation for user-facing functions is

34    included at https://proteoclade.readthedocs.io.

35

36    ProteoClade was benchmarked against MetaProteomeAnalyzer (MPA Portable v. 1.9) and Unipept (v.

37    1.4.1) to assess database creation and taxon annotation speeds. Benchmarking hardware was the same

38    as described above for ProteoClade's testing. For MPA, time and RAM requirements for peptide indexing

39    were monitored over a 24 hour period of attempting to search a simple mass spectra file using the

40    November 2018 UniProt repository, and figures listed in Fig. S2 are based on the projected estimates, as

41    the software had consumed 14 GB out of 16 GB available RAM in the first 24 hours while only processing

42    the first 2 million (out of 140 million) protein sequences. Unipept times and RAM requirements are

43    taken from the latest Unipept publication[1] as the Unipept team reported using a high performance

44    computer, and additionally database creation is not a feature of the software that is available for end

45    users to perform. For taxa annotation speed, three files of 5,000 peptides were randomly drawn from

46    the database, and were taxonomically annotated using Unipept's "pept2taxa" and ProteoClade's

47    "annotate_peptides" functions.

48

2

49    PDX proteomic search

50    Raw spectra files for the Mundt et al. study[2] were downloaded from the Clinical Proteomics Tumor

51    Analysis Consortium (CPTAC) data portal (https://cptac-data-portal.georgetown.edu/cptac/public).

52    ProteoClade was used to retrieve and concatenate the January 2018 human and mouse SwissProt

53    proteome references into a single FASTA file. A PCDB was generated using the same concatenated

54    database with default settings to enable quick annotation.  Spectra were searched using MaxQuant

55    1.6.0.16 with the following parameters: the instrument acquisition settings were set to the default

56    Orbitrap parameters and the protease selected was trypsin/p. Cysteine carbamidomethylation was set

57    as a fixed modification, with protein N-terminal acetylation and methionine oxidation set as variable

58    modifications. TMT 6-plex MS2 ions were used for quantification. The PSM and peptide FDRs were set to

59    0.01.

60

61    ProteoClade PDX annotation and quantitation

62    ProteoClade was used to annotate the resultant peptide ("peptides.txt") files by organism and gene

63    symbol. For each of the mouse and human tissue perspectives, data were analyzed two ways: 1) one in

64    which peptides were included assuming only the organism of interest was present, and 2) one in which

65    both organisms were assumed present for data normalization but separated into species-specific data

66    sets. For method 1, ProteoClade filtered peptides into a data set in which the organism of interest only

67    had to be a plausible assignment for each peptide and the other organism was excluded, simulating a

68    targeted proteomics search in which only a single organism's proteome was used as a reference.

69    ProteoClade assigned peptide sequences and summed the MS2 intensities to gene symbols and then

70    data were normalized by taking the relative intensities of each TMT channel compared to the internal

71    reference pool for each TMT-plex used in the experiment, log2-transforming the data, centering the

72    data by each channel's median, and dividing the relative intensities of each channel by the channel's

3

73    standard deviation. For method 2, ProteoClade assigned peptides to organisms only if the peptide was

74    unique to that organism in the context of the combined human and mouse proteomes. Downstream

75    processing was similar to method 1, but the human and mouse genes quantified using species-unique

76    peptides were separated to yield two distinct data sets after gene assignment and data normalization.

77    Mouse and human data processing approaches were compared using the absolute difference between

78    methods 1 and 2, and plotted using matplotlib 3.0.2.

79

80    To obtain a theoretical perspective of tryptic peptide similarity, the human and mouse Swiss-Prot

81    (release 2018_01) reference databases were digested using a modified version of ProteoClade's digest

82    function. This yielded raw (i.e., non-compressed) peptide sequences which could be compared across

83    homologs using their respective gene symbols.

84

85    Differential gene expression of PDX stroma

86    The species-specific murine data set from Mundt et al.[2] was selected for differential gene expression

87    analysis. For each PDX WHIM sample, TMT channels corresponding to 2 hour vehicle, 50 hour vehicle,

88    and "washout" were used as replicates. Genes were compared across WHIM samples using one-way

89    analysis of variance (ANOVA), and the FDR was calculated and set to a limit of 0.05 using the Benjamini-

90    Hochberg method (RStudio 1.0.153). Data were graphed using ggplot.

91

92    Pathway analysis of PDX stroma

93    Differentially-regulated stromal gene lists from both PDX studies[2,3] as determined by FDR-corrected

94    ANOVA were combined for pathway analysis. Mouse genes were further filtered by removing plasma[4]

95    and abundant erythrocyte proteins identified by proteomics[5] prior to downstream analysis. The web-

96    based gene set analysis toolkit WebGestalt[6] was used for overrepresentation analysis of the stromal

4

97     proteins differentially regulated by PDX tumors in either dataset (2,2293 genes) versus those that were

98     not (671 genes) using the default parameters, *Homo sapiens* as the organism, redundant datasets, and

99     the affinity propagation option for redundancy reduction.

100

101     <u>Genus-level peptidome diversity</u>

102     Taxonomic representation and sequence diversity at the tryptic peptide level was calculated using a

103     modified version of ProteoClade's PCDB module. A tryptic digest of the complete Swiss-Prot and TrEMBL

104     sequences (release 2018_11) using ProteoClade's default parameters was performed, but rather than

105     assign peptides to organisms, non-compressed peptides were assigned directly to genera. Each peptide

106     in the resultant database was checked for genus specificity, and the number of unique and shared

107     peptides for each genus was tallied using custom scripts. Data for the number of unique peptides and

108     fraction shared between genera were plotted using matplotlib.

109

110     <u>Oral microbiome *de novo* search</u>

111     Raw spectra from the Grassl et al. study[7] were retrieved from ProteomeXchange under ID PXD003028. A

112     *de novo* search was performed using PEAKS Studio X (10) with the following parameters: the parent

113     mass error tolerance was set to 10 ppm, the fragment mass error tolerance was set to 0.05 Da, the

114     enzyme was set to trypsin, and MS2 fragmentation was set to higher energy collision-induced

115     dissociation (HCD). Cysteine carbamidomethylation was set as a fixed modification and methionine

116     oxidation was set as a variable modification. ProteoClade was used to make a PCDB of all 140.2 million

117     Swiss-Prot and TrEMBL (release 2018_11) sequences for annotation using its default parameters as

118     described in the implementation section.

119

5

120 ProteoClade *de novo* annotation and genus identification

121 The searched results file ('all de novo candidates.csv') was filtered to only include candidate PSMs with a

122 minimum average local confidence (ALC) score of 50. ProteoClade was used to annotate these PSM

123 candidates using the "Database-Constrained" method, and assigned peptides to the species, genus, and

124 superkingdom ranks. For the Database-Constrained method, candidates for each MS/MS spectra were

125 ordered by descending confidence score, and the sequences were serially checked against the UniProt

126 PCDB until a match, if any, was found. Peptides that did not belong to either to the bacteria

127 superkingdom or human species were removed prior to further processing. Peptides with post-

128 translational modifications were combined with their unmodified sequences for spectral counting, and

129 peptides were only kept in the data set if they were unique to a single genus and had a minimum of two

130 spectral counts across all samples. For FDR control, a PCDB with reversed protein sequences was

131 created with ProteoClade, and the annotation steps were repeated against this database.

132

133 Genus-unique assignments from the *de novo* search that were compared to Grassl et al.'s original

134 targeted database search if the genera were found in both data sets. Additionally, a comparison of

135 ProteoClade's top-scoring PSM annotation option to its Database-Constrained option was made. Genus-

136 specific assignments were ranked by taking the difference between their forward and reversed PCDB

137 annotations. Plots were made using Plotly.

138

139 References:

140 1.    Gurdeep Singh, R. *et al.* Unipept 4.0: Functional Analysis of Metaproteome Data. *J. Proteome Res.*

141       **18**, 606–615 (2019).

142 2.    Mundt, F. *et al.* Mass Spectrometry-Based Proteomics Reveals Potential Roles of NEK9 and

143       MAP2K4 in  Resistance to PI3K Inhibition in Triple-Negative Breast Cancers. *Cancer Res.* **78**, 2732–

6

144        2746 (2018).

145    3.    Wang, X. *et al.* Breast tumors educate the proteome of stromal tissue in an individualized but

146        coordinated manner. *Sci. Signal.* **10**, (2017).

147    4.    Tu, C. *et al.* Depletion of abundant plasma proteins and limitations of plasma proteomics. *J.*

148        *Proteome Res.* **9**, 4982–4991 (2010).

149    5.    Kakhniashvili, D. G., Bulla, L. A. J. & Goodman, S. R. The human erythrocyte proteome: analysis by

150        ion trap mass spectrometry. *Mol. Cell. Proteomics* **3**, 501–509 (2004).

151    6.    Wang, J., Vasaikar, S., Shi, Z., Greer, M. & Zhang, B. WebGestalt 2017: a more comprehensive,

152        powerful, flexible and interactive gene set enrichment analysis toolkit. *Nucleic Acids Res.* **45**,

153        W130–W137 (2017).

154    7.    Grassl, N. *et al.* Ultra-deep and quantitative saliva proteome reveals dynamics of the oral

155        microbiome. *Genome Med.* **8**, 44 (2016).
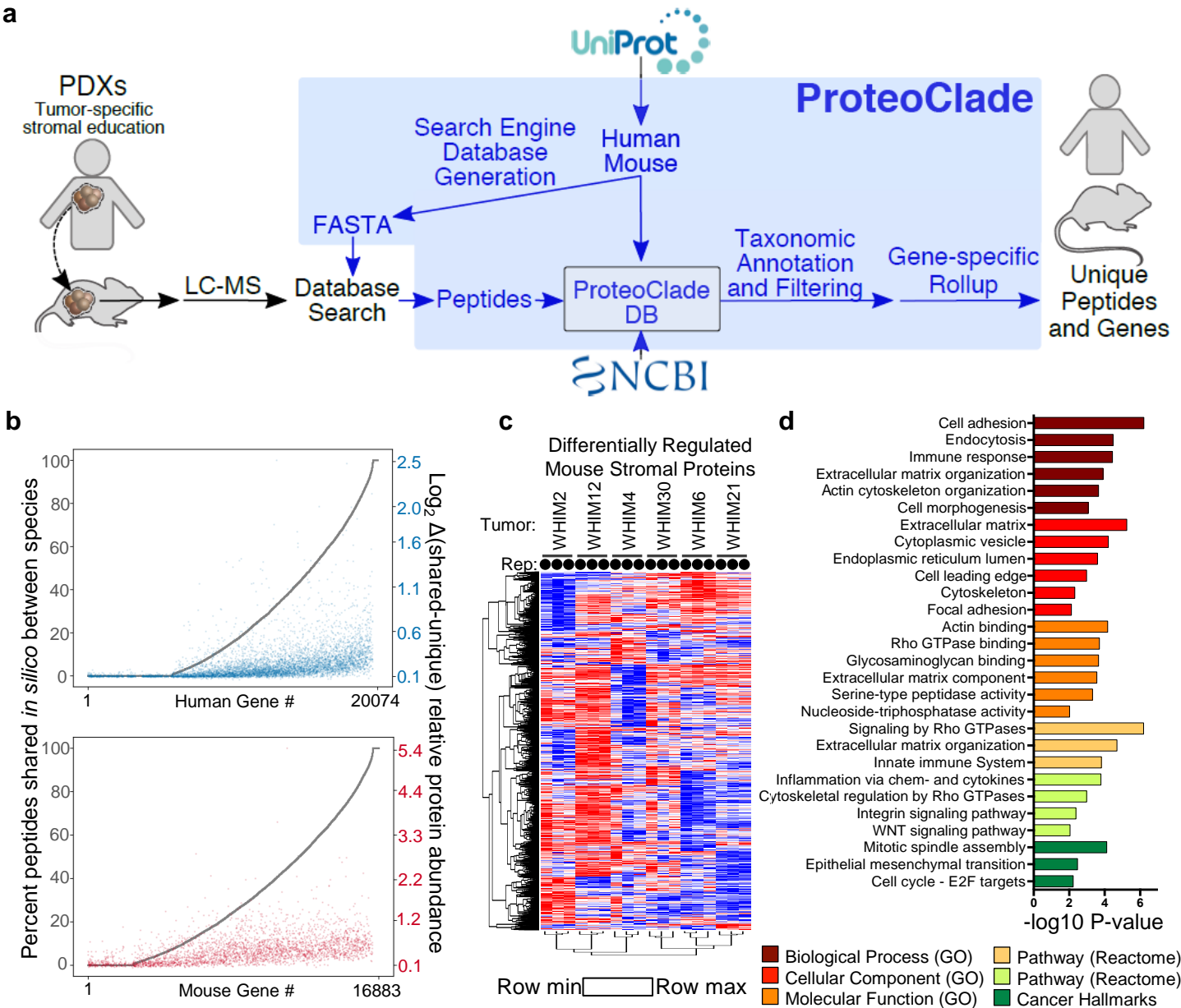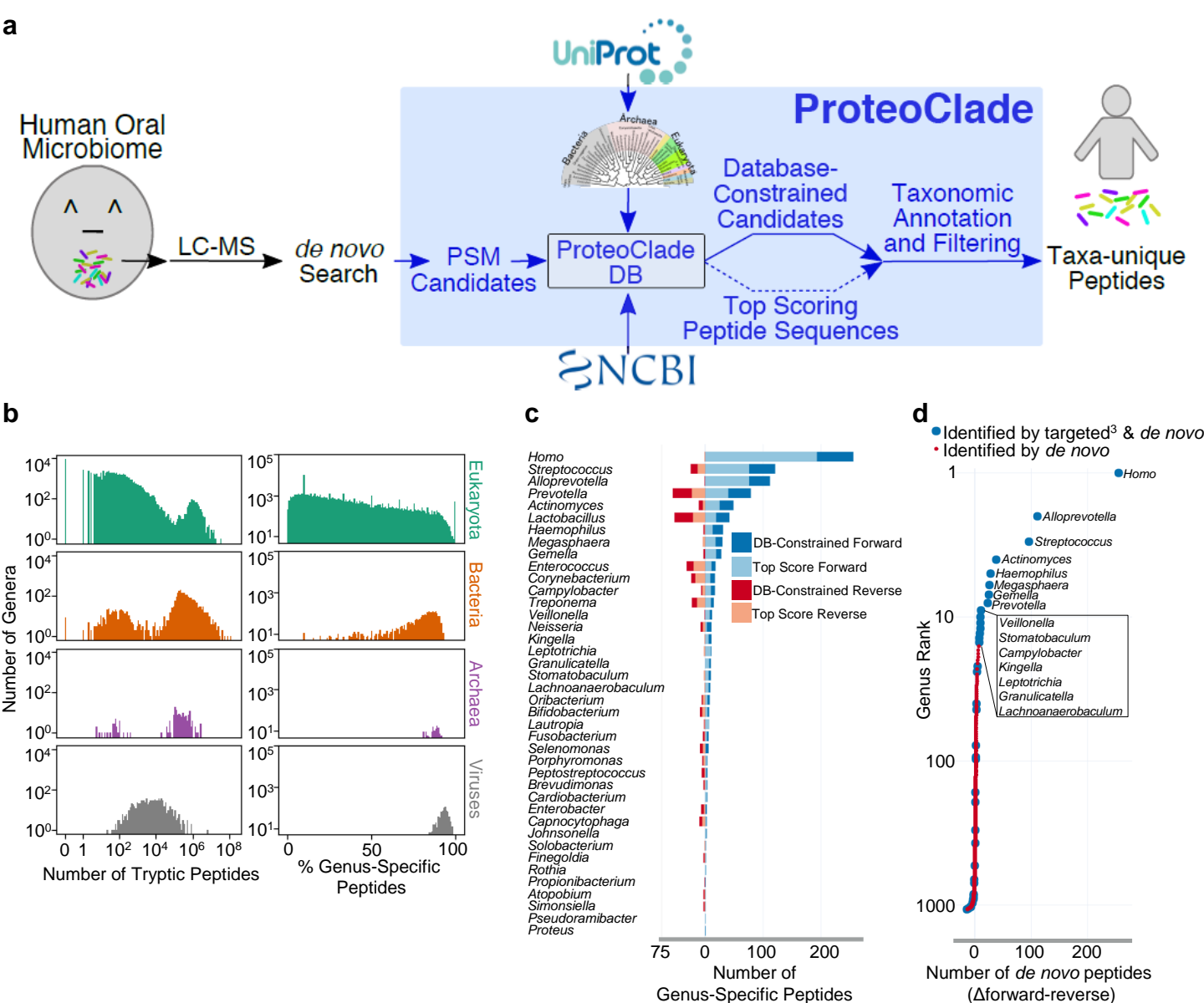
156

# Figure 1

Figure 1: **ProteoClade easily interrogates the interaction between comingled species, including PDX tumors' education of their microenvironment. a)** The ProteoClade workflow for targeted database searches. ProteoClade downloads and concatenates reference FASTA sequence databases and builds a ProteoClade database for fast taxon-specific protein quantitation. **b)** Comparison between taxon-specific and taxon-shared data processing approaches for a PDX data set. The theoretical fraction of shared peptides for each proteome (grey) correlates with the observed quantitative difference between the approaches. A baseline difference between these two analytic approaches was present due to how the assumptions made when assigning peptides to genes affected data normalization. **c)** A significant portion of the murine stromal proteome is differentially regulated by tumors. 3 biological replicates 'Rep' per tumor, indicated with circles. **d)** Pathway enrichment analysis of differentially regulated stromal proteins indicates an enrichment in proteins involved in cellular adhesion and the immune response.

# Figure 2

**Figure 2: ProteoClade enables identification of taxa-specific peptides in metaproteomic samples by *de novo* sequencing. a)** The schematic for ProteoClade's *de novo* pipeline supports the entire UniProt database and the conversion of peptide-spectral candidates to taxon-specific peptides. **b)** A characterization of the entire UniProt peptidome reveals the overrepresentation in the number of peptides contributed by bacterial genera (left), and the significant peptide diversity of microorganisms present in the repository (right). **c)** 39 bacterial genera from the Human Oral Microbiome Database were identified by the *de novo* search in addition to human. DB-Constraned: ProteoClade enhances *de novo* searches by providing the ability to constrain sequence candidates to a database for biological plausibility. "Reverse": ProteoClade controls the false discovery rate by generating a reversed sequence database. **d)** Comparing the forward and reverse UniProt databases resulted in human and 14 oral bacterial genera being identified by the *de novo* approach. Additional organisms were annotated but these annotations lacked statistical confidence.

## MainData

| SQL Field | Type | Description |
|---|---|---|
| RowID | Integer | Row Identifier |
| Protein | Text | Full protein sequence |
| Organism | Integer | Organism Taxon Identifier |
| Gene | Text | Gene Symbol |

## Reference

| SQL Field | Type | Description |
|---|---|---|
| RowID | Integer | Row Identifier |
| HashPeptide | Integer | Hashed Peptide Sequence |
| ProtRowID | Integer | RowID of Protein |

## DBParams

| SQL Field | Type | Description |
|---|---|---|
| RowID | Integer | Row Identifier |
| min_length | Integer | Minimum Peptide Length |
| max_length | Integer | Maximum Peptide Length |
| m_cleave | Boolean | N-terminal Methionine Excision |
| li_swap | Boolean | Leucine to Isoleucine Conversion |
| rule | Text | Proteolysis Rule (Regular Expression) |
| date | Text | Database Creation Date |

# Figure S1

Figure S1: **Database schema for the ProteoClade Database (PCDB)**. The database stores information across three tables: "MainData" stores complete protein sequences, organisms, and genes; "Reference" contains all peptide information and a key back to the protein table (black arrow); "DBParams" stores all database parameters at the time of database creation. The indexed column, "HashPeptide," is indicated in red.

**a** *In silico* Peptide Digestion and Indexing

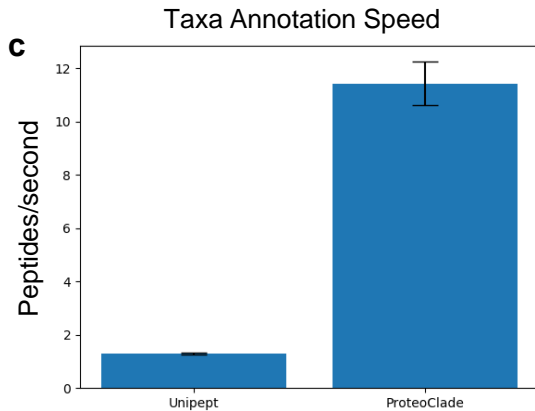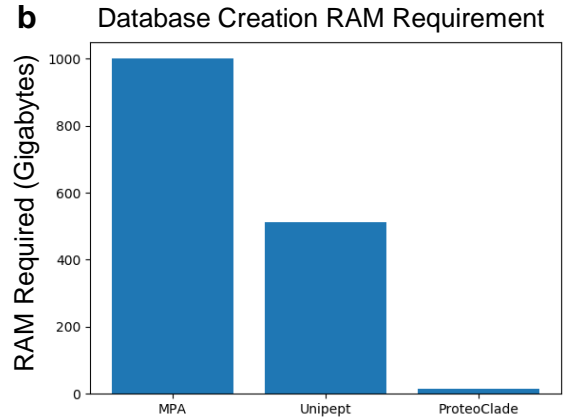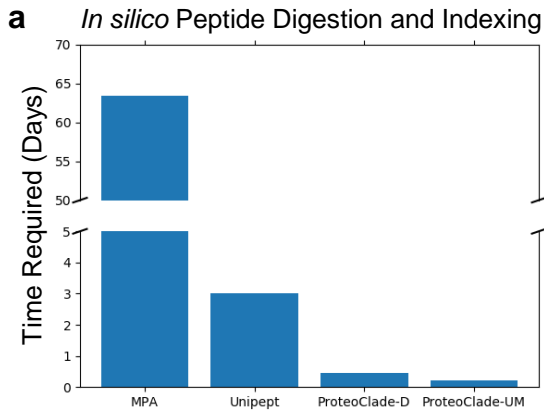**b** Database Creation RAM Requirement

**c** Taxa Annotation Speed

Figure S2

Figure S2: **Technical comparisons between metaproteomic analysis tools. a)** *In silico* digestion of peptides and indexing for ProteoClade's default settings (ProteoClade-D) is faster than previous tools when using the entire UniProt repository. A database using the same parameters as Unipept (ProteoClade-UM) was faster still, due to the absence of missed cleaved peptides. **b)** Database RAM requirements for ProteoClade enable users to generate large databases without using high performance computers. **c)** Annotating all taxa for experimental results is 8.8x faster for ProteoClade than Unipept.