



Ultrafast metaproteomics for quantitative assessment of strain isolates and microbiomes

Elizaveta Kazakova^a, Mark Ivanov^a, Tomiris Kusainova^a, Julia Bubis^a, Valentina Polivtseva^b, Kirill Petrikov^b, Vladimir Gorshkov^c, Frank Kjeldsen^c, Mikhail Gorshkov^a, Yanina Delegan^b, Inna Solyanikova^{b,d}, Irina Tarasova^{a,*}

^a V. L. Talrose Institute for Energy Problems of Chemical Physics, N. N. Semenov Federal Research Center of Chemical Physics, Russian Academy of Sciences, 119334 Moscow, Russia

^b Skryabin Institute of Biochemistry and Physiology of Microorganisms, Pushchino Scientific Center for Biological Research, Russian Academy of Sciences, 142290 Pushchino, Russia

^c Department of Biochemistry and Molecular Biology, University of Southern Denmark, DK-5230 Odense M, Denmark

^d Belgorod National Research University, 308015 Belgorod, Russia

ARTICLE INFO

Keywords:
Microbiological methods
Bacteria
Microbiome
Metaproteomics
Proteomics
Metabolic pathways
Biodegradation
LC-MS/MS

ABSTRACT

Background: Microbial communities are essential in human health and environmental regulation, but present a challenge for the analytical science due to their diversity and dynamic range. Tandem mass spectrometry provides functional insights on microbial life cycle, but is time-consuming. MALDI TOF excels in rapid species identification, but not functional assessment. To address critical challenges in human health and environmental sustainability, microbiology needs advanced mass spectrometry methods and bioinformatic tools enabling both rapid identification and accurate assessment of functional activity of microbial communities.

Results: We show for the first time that both identity and functional activity of microorganisms and their communities can be accurately determined in experiments as short as 7 min per sample, using the basic Orbitrap MS configuration without peptide fragmentation. The approach was validated using strain isolates, mock microbiomes composed of bacteria spiked at known concentrations and human fecal microbiomes. Our new bioinformatic algorithm identifies the bacterial species with an accuracy of 95 %, when no prior information on the sample is available. Microbiome composition was resolved at the genus level with the mean difference between the actual and identified components of 12 %. For mock microbiomes, Pearson coefficient of up to 0.97 was achieved in estimates of strain biomass change. By the example of *Rhodococcus* biodegradation of *n*-alkanes, phenols and its derivatives, we showed the accurate assessment of functional activity of strain isolates, compared with the standard label-free and label-based approaches.

Significance: Our approach makes microbial proteomics fast, functional and insightful using the Orbitrap instruments even without employing peptide fragmentation technology. The approach can be applied to any microorganisms and can take a niche in routine functional assessment of microbial pathogens and consortiums in clinical diagnostics together with MALDI TOF MS and 16S rRNA gene sequencing.

1. Introduction

Metaproteomics is a growing area of research [1] devoted to monitoring the physiology of biological communities, for example, microbial ones, at the level of their proteomes [2]. The metaproteomic studies are focused on determining the species diversity in a community, studying various stages of the life cycle of microorganisms, their response to

stress, the functional characterization of microbiomes, etc. [3]. Being a young field, metaproteomics possess challenges that must be overcome before using for routine research [2,4]. Currently, the research efforts are focused on accuracy, sensitivity, and reproducibility of the analysis, development of bioinformatic solutions, and standardization of workflows [5].

Proteomic analysis of microbiomes relies on liquid chromatography

* Corresponding author at: 38 Leninsky pr., bld 2, 119334 Moscow, Russia.
E-mail address: iatarasova@yandex.ru (I. Tarasova).

<https://doi.org/10.1016/j.microc.2024.111823>

Received 27 June 2024; Received in revised form 24 September 2024; Accepted 30 September 2024

Available online 5 October 2024

0026-265X/© 2024 Elsevier B.V. All rights are reserved, including those for text and data mining, AI training, and similar technologies.

tandem mass spectrometry (LC-MS/MS) that allows the identification of thousands of protein molecules in a single experiment. A common approach is data-dependent acquisition (DDA) mass spectrometry, in which the success of peptide identification strongly depends on the choice of the most abundant ions. Reportedly, the performance of DDA decreases with increasing the sample complexity [6]. The time-consuming DDA analysis with insufficient peptide fragmentation and protein sequence coverage, and heavy bioinformatic processing of metaproteomic data still present a gap for more efficient solutions. Another acquisition method recently applied in metaproteomics is data independent acquisition (DIA) [7–9]. DIA allows collecting fragment spectra for all peptide ions without pre-selection of a fixed number of top abundant ones, and, thus, promises better sequence coverage and protein quantitation [10]. The third approach can be the use of short chromatographic gradients with acquisition of only MS1 spectra without fragmentation. It has been shown that interpretation of LC-MS1 data based on alignments of peptide feature intensities without peptide and protein identification allows large-scale screening of microbiomes and characterization of changes in MS1 profiles, depending on drug treatment [11]. Widely implemented in clinics, MALDI-TOF MS allows identification of the genus and species by matching mass spectra against spectral libraries of known microorganisms [12], but the strain resolution with this technique is problematic [13]. Identifying bacteria by matching the experimental LC-MS1 data against *in silico* generated taxon-specific tryptic peptide masses compiled based on UniProt, SwissProt, and TrEMBL databases, can resolve species and, sometimes, strain identity [14]. However, the functional characterization of microbiomes still cannot be assessed without the identification of proteins.

Identification of proteins relies on protein databases which represents one of the hardest issues in MS-based metaproteomics due to its huge size. The choice of the protein database affects the sensitivity of peptide identification because the larger the search space, the higher the chance of random and incorrect matches. The optimal database should only contain proteins that are present in the sample [15]. From this viewpoint, the protein database constructed using metagenome data or restricted using 16S rRNA gene analysis is the best. However, they are costly and require large sample amounts, as it is important to obtain high-quality sequences of the genome. Attempts to fix the FDR problem typically include selecting top organisms or peptide sequences within a certain window of *m/z* values followed by repeating the search using the restricted database [16–21]. The MetaProteomeAnalyzer [22] implements a search with four different search engines and the integration of the results. Following the search, the identified proteins are grouped into so-called meta-proteins and annotated with protein-level information.

This study presents an adaptation of the novel ultrafast MS1-based method called DirectMS1 [23,24] to the characterization of bacterial cultures and microbiomes. Using DirectMS1 we demonstrate the blind identification (with no prior knowledge of the bacteria in the sample) of strain isolates and microbial communities composition using two-step database search against the microbial part of SwissProt and TrEMBL (62 Gb database size), the evaluation of fold changes in biomass of microorganisms between the metaproteomic samples, and the assessment of metabolic activity of *Rhodococcus* species responding to a change in growth conditions and cold shock.

2. Methods

2.1. Public datasets

Public dataset (<https://zenodo.org/record/3573994>) collected in DDA mode for 19 well-characterized bacterial strains [14] was used to validate the two-step search algorithm for blind identification of individual microorganisms against the microbial part of SwissProt and TrEMBL database (accessed on Jan 2020).

Public dataset (PRJEB42466) for human fecal microbiome, collected in DDA mode, was used for the testing of the blind search algorithm to identify composition of human fecal microbiome [5].

2.2. Strains, genomes and proteogenomic databases

Strains used in the study were isolated from different sources and characterized as previously described (*Rhodococcus erythropolis* X5 [25], *Rhodococcus opacus* S8 [26], *Rhodococcus qingshengii* 7B [27], *Rhodococcus qingshengii* VT6 [28], *Priestia aryabhatai* 25 [29], *Gordonia amicalis* 6-1 [30], *Gordonia alkanivorans* 135 [31,32], *Gordonia polyisoprenivorans* 135 [33], *Rhodococcus opacus* 1CP [34], *Rhodococcus opacus* 3D [35]).

To construct protein databases using genome data, NCBI PGAP 6.6 [36] was used. Annotations of proteins of interest were manually checked using UniProt [37].

2.3. Cell cultivation to characterize metabolic activity changes in response to stress

R. opacus 1CP was cultivated in the presence of four different carbon sources (glucose, benzoate, phenol, 4-chlorophenol) with six biological replicates per condition. The cultivation was carried out in 750-mL Erlenmeyer flasks with 200 mL of the mineral medium of the following composition (g/L): Na₂HPO₄ – 0.7, K₂HPO₄ – 0.5, NH₄NO₃ – 0.75, MgSO₄·7H₂O – 0.2, MnSO₄ – 0.001, FeSO₄ – 0.02 with the addition of one of the following sources of carbon: phenol 500 mg/L, benzoate 500 mg/L, 4-chlorophenol 100 mg/L, glucose 10 g/L. The cultivation was performed for 48–72 h in a shaking incubator (HZQ-111) at 24 °C and 220 rpm. During cultivation time, the carbon sources (phenol and 4-chlorophenol) were periodically injected into the medium after consumption of previously added. The disappearance of substrates was determined by characteristic absorption spectra in the 230–290 nm region. Culture growth was monitored with an optical absorption coefficient at 560 nm.

R. erythropolis X5 cells were grown in an orbital shaking incubator at 180 rpm at two different temperatures (6 °C and 28 °C) for 3 and 11 days, respectively, with three biological replicates. The cultivation was carried out in 750 mL Erlenmeyer flasks with 100 mL of the modified Evans mineral salts medium supplemented with 2 mL of *n*-hexadecane. The composition of the cultivation medium was as follows (per liter): K₂HPO₄, 8.71 g; 5 M solution of NH₄Cl, 1 mL; 0.1 M solution of Na₂SO₄, 1 mL; 62 mM solution of MgCl₂, 1 mL; 1 mM solution of CaCl₂, 1 mL; 0.005 mM solution of (NH₄)₆Mo₇O₂₄, 1 mL; trace element solution, 1 mL. The value of pH was adjusted to 7.5 by concentrated HCl. Composition of the trace element solution in 1 % water solution of HCl (g/L): ZnO, 0.41 g; FeCl₃, 2.9 g; MnCl₂, 1.28 g; CuCl₂, 0.13 g; CoCl₂, 0.26 g; H₃BO₃, 0.06 g.

2.4. Cell cultivation for model microbiomes

Cells were cultivated for 72 h in a shaking incubator (HZQ-111) at 27 °C and 250 rpm in Erlenmeyer flasks with 100 mL of Evans medium. The composition of the medium was as follows (g/L or mL/L): K₂HPO₄ 8.71 g, 5 M NH₄Cl 1 mL, 0.1 M Na₂SO₄ 1 mL, 62 mM MgCl₂ 1 mL, 1 mM CaCl₂ 1 mL, 0.005 mM (NH₄)₆Mo₇O₂₄·4H₂O, micronutrients solution 1 mL (as follows in g/L): ZnO 0.41 g, FeCl₃·6H₂O 5.4 g, MnCl₂·4H₂O 2 g, CuCl₂·2H₂O 0.17 g, CoCl₂·6H₂O 0.48 g, H₃BO₃ 0.06 g, (pH 7.0) – with the addition of 200 µL 10⁸ CFU inoculum and glucose (1 % v/v for *R. opacus* 1CP and 2 % v/v for *R. erythropolis* X5, *G. alkanivorans* 135, *G. amicalis*, *P. aryabhatai* 25) as the sole source of carbon and energy.

2.5. Model microbiome compositions

Model microbiomes were prepared by mixing tryptic peptide solutions to test the performance of DirectMS1 method. The composition of

each sample is described in Tables 1–3.

2.6. Sample preparation for mass spectrometry-based proteomics

For mix models I and III (Tables 1 and 3), cell pellets were resuspended in 150 μ L lysis buffer (50 mM ammonium bicarbonate (ABC), 10 % ACN, 0.1 % ProteaseMax (Promega, USA)). Cells incubated for 45 min at room temperature, after that they were boiled for 10 min at 95 °C, and then cooled down on ice. To lyse cells, the samples were sonicated for 10 min (1 s on 1 s off) on 60 % amplitude (QSonica Q125, Newtown, Connecticut, USA). Samples were centrifuged for 7 min on 10,000 g, then supernatant was taken, and protein concentration was measured by BCA kit (Thermo Scientific, Germany). 10 mM dithiothreitol (DTT) was added for 20 min at 56 °C, followed by incubation with iodoacetamide (IAA) for 30 min in the dark at room temperature. Lys-c was added in 1:100 m/m ratio and samples incubated for 2 h at 37 °C then trypsin was added in 1:50 m/m ratio and samples incubated overnight at 37 °C. Digest was stopped by adding 1 % trifluoroacetic acid (TFA), then samples were desalted with OASIS HLB cartridges (Waters, USA) and dried in a vacuum concentrator.

For mix model II (Table 2) and characterization of metabolic activity changes, cell pellets were resuspended in a 150 μ L lysis buffer (100 mM ABC, 4 % SDS, 10 mM DTT). Cells were boiled for 10 min at 95 °C and then cooled down on ice. To lyse cells, the samples were sonicated for 10 min (1 s on 1 s off) on 60 % amplitude (QSonica Q125, Newtown, Connecticut, USA). Samples were centrifuged for 10 min on 10,000 g, then the supernatant was taken. 100 μ L of supernatant was purified and concentrated with chloroform–methanol precipitation. Dry pellets were resuspended in 50 μ L of 4 M urea in 50 mM ABC, urea concentration was further reduced to 1 M with 50 mM ABC. Protein concentration was measured by a BCA kit (Thermo Scientific, Germany). 10 mM DTT was added for 20 min at 56 °C, followed by incubation with IAA for 30 min in the dark at room temperature. Trypsin was added in a 1:50 m/m ratio and samples were incubated overnight at 37 °C. Digest was stopped by adding 1 % TFA, then samples were desalted with OASIS HLB cartridges (Waters, USA) and dried in a vacuum concentrator.

2.7. LC-MS1 data acquisition

The LC-MS experiments were performed using Orbitrap Q Exactive HF-X mass spectrometer (Thermo Scientific, San Jose, CA, USA) coupled to UltiMate 3000 LC system (Thermo Fisher Scientific, Germering, Germany) and Orbitrap Fusion Lumos mass spectrometer (Thermo Scientific, San Jose, CA, USA) coupled to UltiMate 3000 LC system (Thermo Fisher Scientific, Germering, Germany) and equipped with FAIMS Pro interface. Trap column μ -Precolumn C18 PepMap100 (5 μ m, 300 μ m, i.d. 5 mm, 100 Å) (Thermo Fisher Scientific, USA) and self-packed analytical column (Reprosil-Pur 3 μ m, 75 μ m i.d., 5 cm length) were employed for separation. Mobile phases were as follows: (A) 0.1 % formic acid (FA) in water; (B) 80 % ACN, 0.1 % FA in water. The gradient was from 5 % to 35 % phase B in 4.8 min at 1500 nL/min. Total method time including column washing and equilibration was 7.5 min. Field asymmetric ion mobility spectrometry (FAIMS) separations were performed with the following compensation voltages (CV) –50 V, –65 V, and –80 V in a stepwise mode during LC-MS analysis. Data acquisition was performed

Table 1

Model I: five low fold-change mixtures of peptides (ng) from 3 to 5 strains, total of 500 ng of the sample per injection.

Strain	M1.1	M1.2	M1.3	M1.4	M1.5
<i>Rhodococcus opacus</i> 1CP	167	250	167	100	71
<i>Rhodococcus erythropolis</i> X5	–	–	–	100	143
<i>Priestia aryabhatai</i> 25	167	125	250	100	71
<i>Gordonia amicalis</i> 6-1	–	–	–	100	143
<i>Gordonia alkanivorans</i> 135	167	125	83	100	71

Table 2

Model II: three mixtures of peptides (ng) from 8 strains, total of 1 μ g of the sample per injection.

Strain	M2.1	M2.2	M2.3
<i>Rhodococcus opacus</i> 1CP	16	125	16
<i>Rhodococcus opacus</i> 3D	63	8	31
<i>Rhodococcus opacus</i> S8	251	63	4
<i>Rhodococcus qingshengii</i> 7B	31	502	251
<i>Rhodococcus qingshengii</i> VT6	125	251	502
<i>Rhodococcus erythropolis</i> X5	4	16	125
<i>Gordonia alkanivorans</i> 135	8	4	63
<i>Gordonia polyisoprenivorans</i> 135	502	31	8

Table 3

Model III: four ABRF-like mixtures of peptides (ng) from five strains, total of 1 μ g of the sample per injection. Benchmark design was suggested earlier [38].

Strain	ABRF1	ABRF2	ABRF3	ABRF4
<i>Rhodococcus opacus</i> 1CP	417	96	146	3
<i>Rhodococcus erythropolis</i> X5	353	26	19	104
<i>Priestia aryabhatai</i> 25	96	3	631	88
<i>Gordonia amicalis</i> 6-1	71	1	97	803
<i>Gordonia alkanivorans</i> 135	64	873	107	1

in MS1-only mode. Full MS scans were acquired in a range from m/z 375 to 1500 at a resolution of 120,000 at m/z 200 with AGC target of 4e5, 1 microscan, and 50 ms maximum injection time. Samples were resuspended in phase A and quantities of 1 μ g were loaded per injection.

2.8. LC-MS2, data dependent acquisition

MS/MS-based analysis of samples was performed using Orbitrap Lumos Fusion mass spectrometer (Thermo Fisher Scientific, San Jose, CA, USA) and Orbitrap Q Exactive HF mass spectrometer (Thermo Fisher Scientific, San Jose, CA, USA) coupled to UltiMate 3000 LC system. Trap column μ -Precolumn C18 PepMap100 (5 μ m, 300 μ m, i.d. 5 mm, 100 Å) (Thermo Fisher Scientific, USA) and self-packed analytical column (Inertsil 2 μ m, 75 μ m i.d., 25 cm length) were employed for separation. Mobile phases were as follows: (A) 0.1 % FA in water; (B) 95 % ACN, 0.1 % FA in water.

For TMT-based quantitation of *R. erythropolis* X5, the gradient from 5 % to 30 % phase B for 114 min at a flow rate of 300 nL/min was used. Data was acquired in top20 mode. Full MS scans were acquired in a range from m/z 300 to 1400 at a resolution of 60,000 at m/z 200 with AGC (Automatic Gain Control) target of 3e6, 1 microscan, and 50 ms maximum injection time. Precursor ions were isolated in a 1.4 m/z window and accumulated for a maximum of 100 ms or until the 1e5 AGC target was reached.

For label free quantitation of *R. erythropolis* X5, the gradient from 2 % to 40 % phase B for 60 min at a flow rate of 300 nL/min was used. Full MS scans were acquired in a range from m/z 300 to 1400 at a resolution of 60,000 at m/z 200 with AGC target of 3e6, 1 microscan, and 50 ms maximum injection time. Precursor ions were isolated in a 1.4 m/z window and accumulated for a maximum of 50 ms or until the AGC target of 1e5 charges was reached. Precursors of charge states from 2+ to 6+ (inclusive) were scheduled for fragmentation. To save instrument time for label free DDA analysis, biological replicates were pooled and each pooled sample was measured in triplicate.

2.9. Data processing

Peptide feature detection was performed using the Biosaur2 software [39]. The proteomic search engine ms1searchpy [24] was used for protein identification. Searches were performed against the SwissProt + TrEMBL database (accessed on 06 Dec 2020) of bacterial proteins and strain-specific databases based on annotations of bacterial genomes.

Table 4

Summary on databases and target-decoy FDRs used to characterize different samples. *group-specific FDR, target-decoy approach.

Sample type	NCBI: txid	Database	# proteins in fasta	FDR
19 strain isolates	All	SwissProt + TrEMBL		0.05
Fecal human microbiome	All	SwissProt + TrEMBL		0.05*
Model I (<i>R. opacus</i> 1CP)	37919	SwissProt + TrEMBL	7769	0.05
<i>R. erythropolis</i> X5	–	Proteogenomic	6407	
<i>G. alkanivorans</i> 135	–	Proteogenomic	5855	
<i>G. amicalis</i> 6-1	1220574	SwissProt + TrEMBL	4528	
<i>P. aryabhatai</i> 25)	1358420	SwissProt + TrEMBL	6393	
Model II (<i>R. opacus</i> 1CP)	37919	TrEMBL	7769	0.05
<i>R. opacus</i> 3D	–	Proteogenomic	8870	
<i>R. opacus</i> S8	–		8102	
<i>R. erythropolis</i> X5	–		6407	
<i>R. qingshengii</i> 7B	–		6277	
<i>R. qingshengii</i> VT6	–		6789	
<i>G. alkanivorans</i> 135	–		5855	
<i>G. polyisoprenivorans</i> 135)	–		6248	
Model III (<i>R. opacus</i> 1CP)	37919	SwissProt + TrEMBL	7769	0.05
<i>R. erythropolis</i> X5	–	Proteogenomic	6407	
<i>P. aryabhatai</i> 25	–	Proteogenomic	5855	
<i>G. amicalis</i> 6-1	1220574	SwissProt + TrEMBL	4528	
<i>G. alkanivorans</i> 135)	1358420	SwissProt + TrEMBL	6393	
Strain isolate response 1	37919	SwissProt + TrEMBL	7769	0.01
<i>R. opacus</i> 1CP				
Strain isolate response 2	–	Proteogenomic	6407	0.01
<i>R. erythropolis</i> X5				

Table 4 summarizes the origin of databases used for every type of sample. Mass tolerance for precursors in all data was ± 10 ppm, fragment mass tolerance for DDA data was 0.01 Da. Carbamidomethylation of cysteine residues was fixed modification, no variable modifications and no missed cleavages were allowed in MS1 search. For DDA data, the variable methionine oxidation and tmt10plex tags (where applicable) were allowed. The other search parameters were used by default. The detailed information on search algorithm, target-decoy method used to control false positive protein identifications, on matching the same protein/peptide feature to multiple species and its quantitation is provided in the [supplementary materials \(Appendix A–C\)](#).

To estimate fold changes of microorganism biomasses between samples of model microbiomes I and II, protein identification was made against the pooled proteogenomic databases (**Table 4**). Quantitation was performed with DirectMS1Quant [40].

To characterize changes in metabolic activity of *R. opacus* 1CP, the DirectMS1Quant was used for quantitation using the following parameters: differentially regulated proteins satisfy Benjamini Hochberg FDR < 0.05, fold change (FC) threshold was two standard deviations of \log_2 FC distribution; intensity normalization by 1000 quantified peptides with maximal intensities was applied. Functional annotation and gene ontology analysis was performed using STRING db [41]. The annotation is publicly available at <https://version-12-0.string-db.org/organism/STRG0A76PPK>.

R. erythropolis X5 was quantified using Diffacto [42] and QRePS [43]. Differentially regulated proteins were selected using the following

criteria: $|\log_2FC| > 1.2$, $|\log_{10}FDR| > Q3 + 1.5 * IQR$, where Q3 and IQR are 3rd quartile and interquartile range [43]. Functional annotation and gene ontology analysis were performed using STRING db. The annotation is freely available at <https://version-12-0.string-db.org/organism/STRG0A44VNI>. **Tables S1 and S2** containing quantitation results for *R. opacus* 1CP and *R. erythropolis* X5 are provided in [supporting information](#).

3. Results

3.1. Algorithm for two-stage blind database search identifies genus and species of a bacteria from LC-MS1 data

3.1.1. Algorithm for two-stage search against SwissProt + TrEMBL database

The standard DirectMS1 analysis using ms1searchpy search engine cannot handle protein databases with more than few hundred thousands of proteins. To solve that issue, we first used a preliminary search to identify the most probable strains presented in the sample using the whole bacterial database SwissProt + TrEMBL as input candidates. The algorithm for this search is based on a simple matching of theoretical peptide m/z values with experimental MS1 spectra. Based on the identifications from the preliminary search, the shortened database is composed of the most probable species. This shortened database is further used for accurate protein identification using the standard ms1searchpy search. Details on the algorithm are provided in the [Supporting information \(Appendix A, Fig. S1\)](#). As an example, the runtime for full analysis on a computer with Intel(R) Core(TM) i7-7700HQ CPU @ 2.80 GHz (4 cores) takes 4 h:3 h for SwissProt + TrEMBL parsing and calculation of theoretical m/z (performed once), 6 min for preliminary search and 1 h for standard search against the shortened database (a sum for 5 sample replicates).

3.1.2. Two-stage blind search against SwissProt + TrEMBL demonstrates high accuracy in the identification of isolated bacterial species

To estimate the efficiency of the proposed two-stage identification approach, public LC-MS/MS data for 19 strain isolates was used [14]. The results of blind identification are presented in **Table 5** and **Table S3**. **Table S3** shows the summary on the refined database sizes, number of identified proteins and the percentages of proteins corresponding to the correct species and to the correct genus. In **Table 5**, the top identified microorganism matched the analyzed strain at the level of genus + species in 95 %, and 32 % of those cases matched the correct strain ("Level of match" column of **Table 5**). Analysis of top-3 identified microorganisms reveals that 11 of 19 cases showed the undoubted leader supported by at least 10 times higher numbers of identified proteins than the top-second and top-third taxons. The remaining eight cases demonstrated less than 2-fold difference in numbers of protein identifications distributed between the top-first and top-second organisms. This observation corresponds to the identification of taxons of the same species or species group and means a high similarity of proteomes. *E. coli* was the only strain identified at the level of the genus due to the classification within the ete3 toolkit [44]. The *E. coli* strain K-12 (NCBI: txid83333, 4518 proteins in SwissProt) is not reported with "NCBITaxa().get_descendant_taxa()" function for *E. coli* species (NCBI:txid562) which leads to it being excluded from the search space. At the next step, *E. coli* species database is also excluded in blind search due to its excessive size (1501705 proteins). Reportedly, differentiation between *Escherichia coli* and *Shigella* [45] or *Bacillus cereus* and *Bacillus anthracis* [46] is also problematic for MALDI TOF that requires development of specific solutions. The proposed blind search algorithm can be successfully applied for identification of individual species and strains using fast proteome profiling data. In addition, **Table S3** shows the database sizes. On average, 1.0 % of the proteins from the selected database were identified. Among these protein identifications, 97 % and 66 % belong to the correct genus and species, respectively. These results show that our

Table 5

Summary of blind identification of bacteria against SwissProt + TrEMBL database shows correct match between the analyzed and top identified strains at the level of species groups, species and strains. Level of match was determined using phylogenetic trees available at <https://phylot.biobyte.de/>. Percentage is defined as (#proteins identified per species/total identified proteins) * 100 %. Strains taken in multiple replicates were analyzed calculating the mean number of identifications across replicates.

Strain analyzed	NCBI: txid	1st taxon by #proteins, %	2nd taxon by #proteins, %	3rd taxon by #proteins, %	Level of match
<i>Acinetobacter baumannii</i> DSM 30007	470 575584	<i>Acinetobacter baumannii</i> (45 %)	<i>Acinetobacter</i> sp. FDAARGOS_559 (41 %)	<i>Acinetobacter</i> sp. 25977_8 (6 %)	strain
<i>Bacillus cereus</i> DSM 31	1396 226900	<i>Bacillus anthracis</i> (51 %)	<i>Bacillus</i> sp. S66 (28 %)	<i>Bacillus</i> sp. B13 (5 %)	Species group strain
<i>Bacillus velezensis</i> DSM 23117/FZB42	492670 326423	<i>Bacillus velezensis</i> FZB42 (32 %)	<i>Bacillus velezensis</i> (19 %)	<i>Bacillus</i> sp. VMFN-A1 (18 %)	strain
<i>Burkholderia cepacia</i> ATCC 25416	292 983594	<i>Burkholderia reimsii</i> (22 %)	<i>Burkholderia</i> sp. LS-044 (16 %)	<i>Burkholderia lata</i> (14 %)	Species group strain
<i>Burkholderia thailandensis</i> DSM 13276	57975 271848	<i>Burkholderia thailandensis</i> (65–68 %)	<i>Burkholderia oklahomensis</i> (6–8 %)	<i>Burkholderia</i> sp. MSMB1552 (6–8 %)	strain
<i>Burkholderia thailandensis</i> E125	57975	<i>Burkholderia thailandensis</i> (67–70 %)	<i>Burkholderia</i> sp. MSMB1552 (8–12 %)	<i>Burkholderia oklahomensis</i> (6–8 %)	species
<i>Burkholderia thailandensis</i> E131	57975	<i>Burkholderia thailandensis</i> (67–68 %)	<i>Burkholderia</i> sp. MSMB1552 (7 %)	<i>Burkholderia oklahomensis</i> (1–9 %)	species
<i>Burkholderia thailandensis</i> E153	57975	<i>Burkholderia thailandensis</i> (65–68 %)	<i>Burkholderia</i> sp. MSMB1552 (7–13 %)	<i>Burkholderia oklahomensis</i> (8 %)	species
<i>Burkholderia thailandensis</i> LMG 20219	57975	<i>Burkholderia thailandensis</i> (67–71 %)	<i>Burkholderia</i> sp. MSMB1552 (8–11 %)	<i>Burkholderia oklahomensis</i> (1–8 %)	species
<i>Burkholderia oklahomensis</i> DSM 21774	342113	<i>Burkholderia oklahomensis</i> (78–83 %)	<i>Burkholderia lata</i> (4 %)	<i>Burkholderia</i> sp. ABCPW 14 (3 %)	strain
<i>Citrobacter freundii</i> DSM 30039	546 1006003	uncultured <i>Citrobacter</i> sp. (57 %)	<i>Citrobacter freundii</i> complex sp. CFNIH9 (14 %)	<i>Citrobacter</i> sp. LUTT5 (8 %)	Species group strain
<i>Enterococcus faecalis</i> DSM 20371	1351	<i>Enterococcus faecalis</i> (92 %)	<i>Enterococcus faecium</i> (2 %)	<i>Enterococcus haemoperoxidus</i> ATCC BAA-382 (1 %)	strain
<i>Mycobacteroides abscessus</i> DSM 44196	36809 561007	<i>Mycobacteroides abscessus</i> ATCC 19977 (57–58 %)	<i>Mycobacteroides abscessus</i> subsp. Abscessus (18–19 %)	<i>Mycobacteroides franklinii</i> (11 %)	strain
<i>Escherichia coli</i> DSM 3871 (K12 W3110 derivative)	562 83333	<i>Escherichia fergusonii</i> ATCC 35469 (20 %)	<i>Escherichia</i> sp. KTE114 (18 %)	<i>Shigella boydii</i> (12 %)	genus
<i>Pseudomonas aeruginosa</i> ATCC 27853	287	<i>Pseudomonas aeruginosa</i> (93–94 %)	<i>Pseudomonas</i> sp. ATCC 13867 (2 %)	<i>Pseudomonas nitroreducens</i> (1 %)	species
<i>Staphylococcus epidermidis</i> DSM 1798	1282	<i>Staphylococcus epidermidis</i> (85 %)	<i>Staphylococcus</i> sp. HMSC34G04 (4 %)	<i>Staphylococcus</i> sp. RIT622 (2 %)	species
<i>Staphylococcus aureus</i> DSM 20231/NCTC 8532	1280 1241616	<i>Staphylococcus aureus</i> (95 %)	<i>Staphylococcus pasteurii</i> (1 %)	<i>Staphylococcus argenteus</i> (1 %)	species
<i>Vibrio cholerae</i> NIH 41	666	<i>Vibrio cholerae</i> (77 %)	<i>Vibrio mimicus</i> (8 %)	<i>Vibrio metoecus</i> (6 %)	species
<i>Yersinia pseudotuberculosis</i> DSM 8992	633	<i>Yersinia wautersia</i> (41 %)	<i>Yersinia pseudotuberculosis</i> (28 %)	<i>Yersinia similis</i> (5 %)	species

method cannot identify the accurate list of homologous species present in the sample, except the top-ranked one which is definitely presented in the sample. However, these results do not mean that method cannot detect the difference in the quantitation between homologous species, which will be discussed below in the manuscript in quantitation section.

3.2. Two-stage blind database search identifies relative composition of microbial community at level of genus

3.2.1. Model mixtures of soil bacteria

To study the performance of a two-stage blind search to identify microbiome composition, the soil bacteria mixed at known concentrations (Tables 1–3) were used. The main motivation was to investigate how blind searches resolve the metaproteomic samples containing different species of the same genus. The results of blind identification of sample compositions are summarized in Table 6. The actual and estimated compositions of each model microbiome were compared at two levels: genus and species. In blind identification, the mean difference between the actual and identified components of microbiome composition was 12 % at the level of genus. The difference between the actual and estimated compositions depended on the microorganisms under study. In our example, the content of *Rhodococcus* was easily identified

and often overestimated, while *Gordonia* was identified in all samples (even though its content was below 15 %), but always underestimated. *Priestia* was identified with high accuracy within a few percent, but only if its actual fraction in the sample was higher than 30 %. Below 30 %, *Priestia* was not identified at all. At the level of species, the mean difference between the estimated and actual content was 15 %. However, the number of missing identifications of microorganisms increased from 20 % at the level of genus to 40 % at the level of species. Thus, the current realization of blind search algorithm for identification of microbiome composition provides the most complete information at the level of genus.

3.2.2. Fecal microbiome

To test the algorithm in identifying the composition of the real microbiome, the public dataset for human fecal microbiome was used [14]. At the level of family-specific FDR < 0.05 (Appendix B), the blind search algorithm identified four families in the fecal sample (sample “F04” [14]: *Oscillospiraceae* (heterotypic synonym for *Ruminococcaceae*), *Lachnospiraceae*, *Eubacteriaceae* and *Clostridiaceae* (Fig. 1a). Fig. S2 shows the top 15 organisms identified in the sample at the first stage of the blind search. These top 15 organisms are the representatives of the mentioned above families. In the original publication [14], the families

Table 6

Summary of blind identification of microbial compositions against SwissProt + TrEMBL database. Sample labels and actual compositions correspond to Tables 1-3. Percentage in actual compositions is defined as (peptide mass per genus (species)/total peptide mass) * 100 %. Percentage in estimated compositions is defined as (#proteins per genus (species)/total #proteins) * 100 %. Empty cells mean missing identification of microorganisms from a given genus or species. The mean difference between actual and estimated composition was calculated as $\sum |\text{Actual \%} - \text{Estimated \%}|/N$, where $N = 33$ is the number of rows in the column "actual composition, % by mass". Missings were imputed by zero.

Sample	GENUS		SPECIES	
	Actual composition, % by mass	Estimated composition, % by #proteins	Actual composition, % by mass	Estimated composition, % by #proteins
M1.1	<i>Rhodococcus</i> (33 %) <i>Priestia</i> (33 %) <i>Gordonia</i> (33 %)	<i>Rhodococcus</i> (38.4 %) <i>Priestia</i> (39.7 %) <i>Gordonia</i> (15.4 %)	<i>R. opacus</i> (33 %) <i>P. aryabhatai</i> (33 %) <i>G. alkanivorans</i> (33 %)	<i>R. opacus</i> (5.3 %) <i>P. aryabhatai</i> (39.7 %)
M1.2	<i>Rhodococcus</i> (50 %) <i>Priestia</i> (25 %) <i>Gordonia</i> (25 %)	<i>Rhodococcus</i> (80.5 %) <i>Gordonia</i> (9.9 %)	<i>R. opacus</i> (50 %) <i>P. aryabhatai</i> (25 %) <i>G. alkanivorans</i> (25 %)	<i>R. opacus</i> (6.1 %)
M1.3	<i>Rhodococcus</i> (33 %) <i>Priestia</i> (50 %) <i>Gordonia</i> (17 %)	<i>Rhodococcus</i> (37.2 %) <i>Priestia</i> (50.3 %) <i>Gordonia</i> (6.6 %)	<i>R. opacus</i> (33 %) <i>P. aryabhatai</i> (50 %) <i>G. alkanivorans</i> (17 %)	<i>R. opacus</i> (5.9 %) <i>P. aryabhatai</i> (48.4 %)
M1.4	<i>Rhodococcus</i> (40 %) <i>Priestia</i> (20 %) <i>Gordonia</i> (40 %)	<i>Rhodococcus</i> (64.3 %) <i>Gordonia</i> (26.2 %)	<i>R. opacus</i> (20 %) <i>R. erythropolis</i> (20 %) <i>P. aryabhatai</i> (20 %) <i>G. amicalis</i> (20 %) <i>G. alkanivorans</i> (20 %)	<i>R. opacus</i> B4 (3.5 %) <i>R. erythropolis</i> PR4 (12.7 %) <i>G. amicalis</i> NBRC 100051 = JCM 11271 (3.8 %) <i>G. alkanivorans</i> NBRC 16433 (13.7 %)
M1.5	<i>Rhodococcus</i> (43 %) <i>Priestia</i> (14 %) <i>Gordonia</i> (43 %)	<i>Rhodococcus</i> (64.3 %) <i>Gordonia</i> (29.7 %)	<i>R. opacus</i> (14 %) <i>R. erythropolis</i> (29 %) <i>P. aryabhatai</i> (14 %) <i>G. amicalis</i> (29 %) <i>G. alkanivorans</i> (14 %)	<i>R. opacus</i> B4 (3.5 %) <i>R. erythropolis</i> (25 %) <i>R. erythropolis</i> PR4 (5.9 %) <i>G. alkanivorans</i> NBRC 16433 (22.0 %)
M2.1	<i>Rhodococcus</i> (49 %) <i>Gordonia</i> (51 %)	<i>Rhodococcus</i> (64.8 %) <i>Gordonia</i> (30.0 %)	<i>R. opacus</i> (33 %) <i>R. qingshengii</i> (16 %) <i>R. erythropolis</i> (0.4 %) <i>G. alkanivorans</i> (0.8 %) <i>G. polyisoprenivorans</i> (50 %)	<i>R. opacus</i> (31.8 %) <i>R. opacus</i> B4 (2.8 %) <i>R. erythropolis</i> PR4 (10.8 %) <i>G. polyisoprenivorans</i> NBRC 16320 = JCM 10675 (23.8 %)
M2.2	<i>Rhodococcus</i> (96.5 %) <i>Gordonia</i> (3.5 %)	<i>Rhodococcus</i> (95.6 %) <i>Gordonia</i> (0.9 %)	<i>R. opacus</i> (20 %) <i>R. qingshengii</i> (75 %) <i>R. erythropolis</i> (1.6 %) <i>G. alkanivorans</i> (0.4 %) <i>G. polyisoprenivorans</i> (3 %)	<i>R. opacus</i> (19.5 %) <i>R. opacus</i> B4 (2.4 %) <i>R. qingshengii</i> (21.4 %) <i>R. erythropolis</i> (31.1 %) <i>R. erythropolis</i> PR4 (7.5 %)
M2.3	<i>Rhodococcus</i> (92.9 %) <i>Gordonia</i> (7.1 %)	<i>Rhodococcus</i> (95.2 %) <i>Gordonia</i> (1.1 %)	<i>R. opacus</i> (5.1 %) <i>R. qingshengii</i> (75.1 %) <i>R. erythropolis</i> (13 %) <i>G. alkanivorans</i> (6 %) <i>G. polyisoprenivorans</i> (0.8 %)	<i>R. opacus</i> B4 (1.1 %) <i>R. qingshengii</i> (28.5 %) <i>R. erythropolis</i> (43.7 %) <i>R.erythropolis</i> PR4 (11.0 %)
ABRF1	<i>Rhodococcus</i> (77 %) <i>Priestia</i> (10 %) <i>Gordonia</i> (13 %)	<i>Rhodococcus</i> (89.2 %) <i>Gordonia</i> (3.1 %)	<i>R. opacus</i> (42 %) <i>R. erythropolis</i> (35 %) <i>P. aryabhatai</i> (10 %) <i>G. amicalis</i> (7 %) <i>G. alkanivorans</i> (6 %)	<i>R. opacus</i> B4 (4.8 %) <i>R. erythropolis</i> (32.3 %) <i>R. erythropolis</i> PR4 (8.9 %)
ABRF2	<i>Rhodococcus</i> (12.6 %) <i>Priestia</i> (0.3 %) <i>Gordonia</i> (87.1 %)	<i>Rhodococcus</i> (18.7 %) <i>Gordonia</i> (73.9 %)	<i>R. opacus</i> (9.6 %) <i>R. erythropolis</i> (3 %) <i>P. aryabhatai</i> (0.3 %) <i>G. amicalis</i> (0.1 %) <i>G. alkanivorans</i> (87 %)	<i>R. opacus</i> B4 (3.1 %) <i>G. amicalis</i> NBRC 100051 = JCM 11271 (4.8 %) <i>G. alkanivorans</i> NBRC 16433 (46.9 %)
ABRF3	<i>Rhodococcus</i> (17 %)	<i>Rhodococcus</i> (21.8 %)	<i>R. opacus</i> (15 %) <i>R. erythropolis</i> (2 %)	<i>R. opacus</i> B4 (5.6 %) <i>R. erythropolis</i> PR4 (1.9 %)

(continued on next page)

Table 6 (continued)

Sample	GENUS		SPECIES	
	Actual composition, % by mass	Estimated composition, % by #proteins	Actual composition, % by mass	Estimated composition, % by #proteins
	<i>Priestia</i> (63 %) <i>Gordonia</i> (20 %)	<i>Priestia</i> (65.0 %) <i>Gordonia</i> (6.0 %)	<i>P. aryabhatai</i> (63 %) <i>G. amicalis</i> (10 %) <i>G. alkanivorans</i> (10 %)	<i>P. aryabhatai</i> (65.0 %)
ABRF4	<i>Rhodococcus</i> (10.7 %) <i>Priestia</i> (8.8 %) <i>Gordonia</i> (80.5 %)	<i>Rhodococcus</i> (32.4 %) <i>Gordonia</i> (63.7 %)	<i>R. opacus</i> (0.3 %) <i>R. erythropolis</i> (10.4 %) <i>P. aryabhatai</i> (8.8 %) <i>G. amicalis</i> (80.4 %) <i>G. alkanivorans</i> (0.1 %)	<i>R. opacus</i> (0.7 %) <i>R. erythropolis</i> PR4 (11.1 %) <i>G. amicalis</i> NBRC 100051 = JCM 11271 (12.6 %)

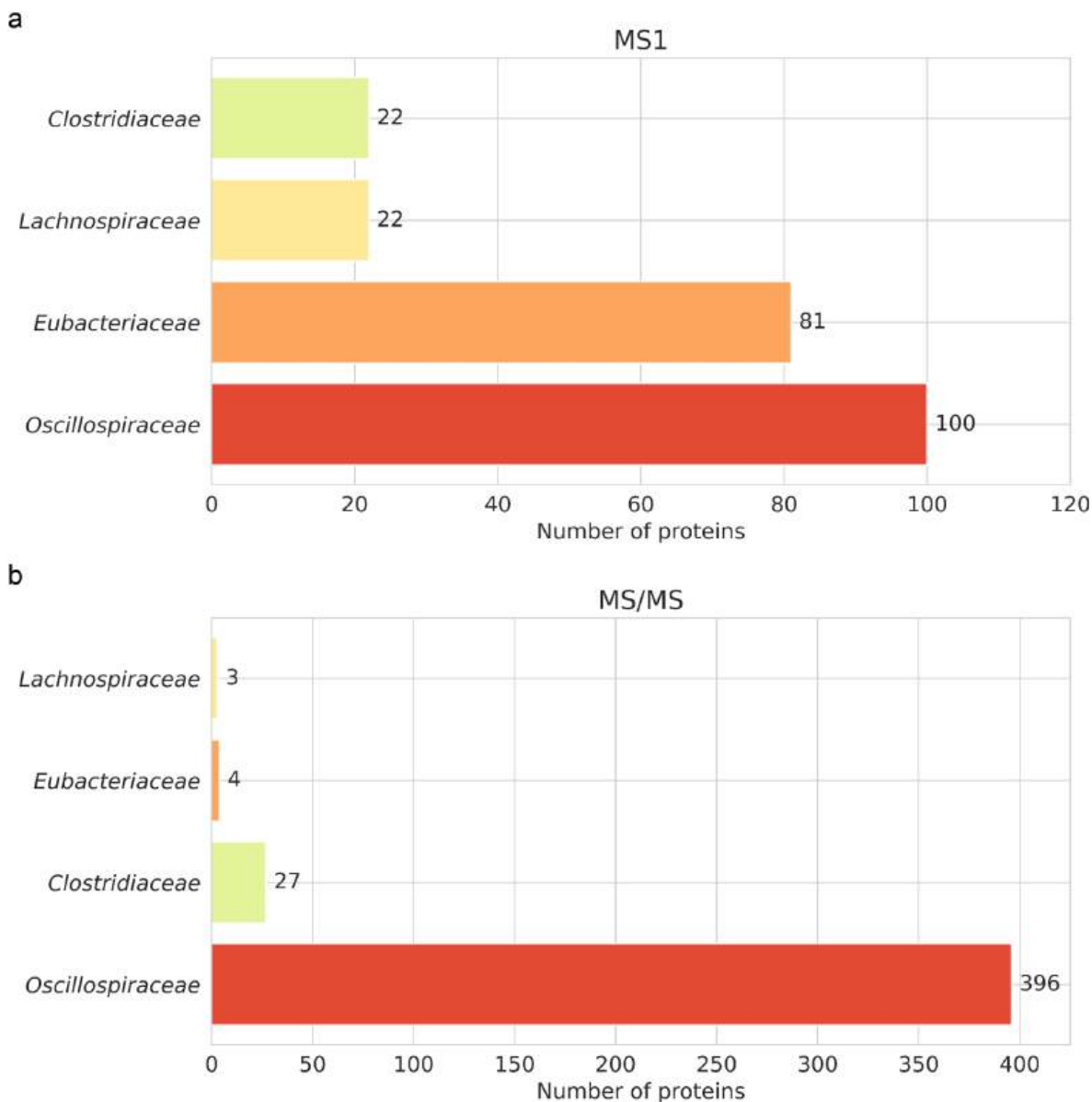


Fig. 1. Comparison of the performance of DirectMS1 and standard LC-MS/MS DDA in establishing fecal microbiome composition. Sample: labeled in the original study as F04 [14]. Family-specific FDR is 0.05 (estimated using the target-decoy approach, Appendix B).

Oscillospiraceae (heterotypic synonym for *Ruminococcaceae*), *Lachnospiraceae*, *Eubacteriaceae* and *Clostridiaceae* have been identified as the most abundant (approximately 90 % of the sample).

For comparison with MS/MS-based identification, we re-analyzed

MS/MS spectra from single shot LFQ DDA run of fecal sample against the shortened protein database composed after the first stage of blind search (Fig. S1-VI). Taxonomic distribution (Fig. 1b) revealed the top family *Oscillospiraceae* that coincided with the results derived from the

MS1-spectra (Fig. 1a). *Lachnospiraceae*, *Eubacteriaceae* and *Clostridiaceae* bacterial families corresponded together to less than 8 % of the proteins identified from MS/MS (Fig. 1b). MS1 analysis revealed that these families compose approximately 56 % of the sample (10.0 %, 10.0 %, and 36.0 %, respectively). This result agrees well with the results presented across all datasets in the original publication [14]. We suggest that our MS1-based bioinformatic workflow demonstrates higher sensitivity than MS/MS-based identification for data collected in single shot LFQ DDA analysis of not fractionated samples. Thus, the proposed algorithm can be applicable for the identification of both strain isolates and complex samples such as the real human gut microbiome.

3.3. DirectMS1 allows correct determination of the fold change of species biomass between microbiome samples

Model microbiomes I, II and III (Tables 1–3) were used to analyze strain abundance variation between different samples with the following procedure. After quantitation at peptide level, the fold changes of all quantified peptides without decoys were used to plot the baseline density distribution of matched peptides (Fig. 2a). Here, the density is the value of the probability density function of \log_2FC at the bin, normalized so that the integral over the range is 1. Then, fold change distribution was built for all peptides from each strain of the model microbiome (“strain” distribution). The baseline was subtracted from the “strain” distribution and the corrected density distribution was analyzed (Fig. 2b–d). See Appendix C and Fig. S3 for details on the procedure. A weighted average of fold changes for the corrected distribution provides an approximation of the biomass fold change for the selected strain, where weights are non-negative corrected densities.

Fig. 3 summarizes estimates of fold changes in strain biomass between samples and compares it to actual mass ratios. For all samples from the Model I, the experimentally measured fold changes coincide well with the actual values (Pearson correlation coefficient of 0.97, Fig. 3a). It means that species mixed at ratios 1:1, 1:2, 1:3 and composing at least from 15 to 20 % of a microbiome are typically well resolved with our method. For Model II and III, a decrease in fold change predictions was observed (Fig. 3b, c). Model II represented the extreme case of uniform composition of the same species (*R. opacus* 1CP, *R. opacus* 3D, *R. opacus* S8, *R. erythropolis* X5, *R. qingshengii* 7B, *R. qingshengii* VT6, *G. alkanivorans* 135, *G. polyisoprenivorans* 135) mixed in a range of concentrations differing by up to two orders (Table 2 in M&M). This complex case resulted in a decrease of Pearson correlation to 0.81 due to *Gordonia* species composing just 3.5 % (M2.2) and 7.1 % (M2.3) of the model microbiome (Fig. 3b).

Model III had compositions differing at the level of species (*R. opacus* 1CP, *R. erythropolis* X5, *P. aryabhatai* 25, *G. amicalis* 6-1, *G. alkanivorans* 135) (Table 3 in M&M) and covering the similar range of concentrations as samples from Model II. Model III is less complex compared with Model II that results in more accurate fold change estimations ($R = 0.93$, Fig. 3c).

For the low amounts (species or genera content within a few percent) in one of the samples, the missing values can occur and subsequent imputation of peptide abundance results in a decreased accuracy. Fig. S4 illustrates such a case with the bacterium *P. aryabhatai* 25 from Model III ABRF2/ABRF3, for which the peptide content was 3 ng and 631 ng, respectively. If the total biomass of microorganisms is low in both samples, the method could not quantify enough peptides to construct the distribution. This case is illustrated with the bacterium *R. erythropolis* X5

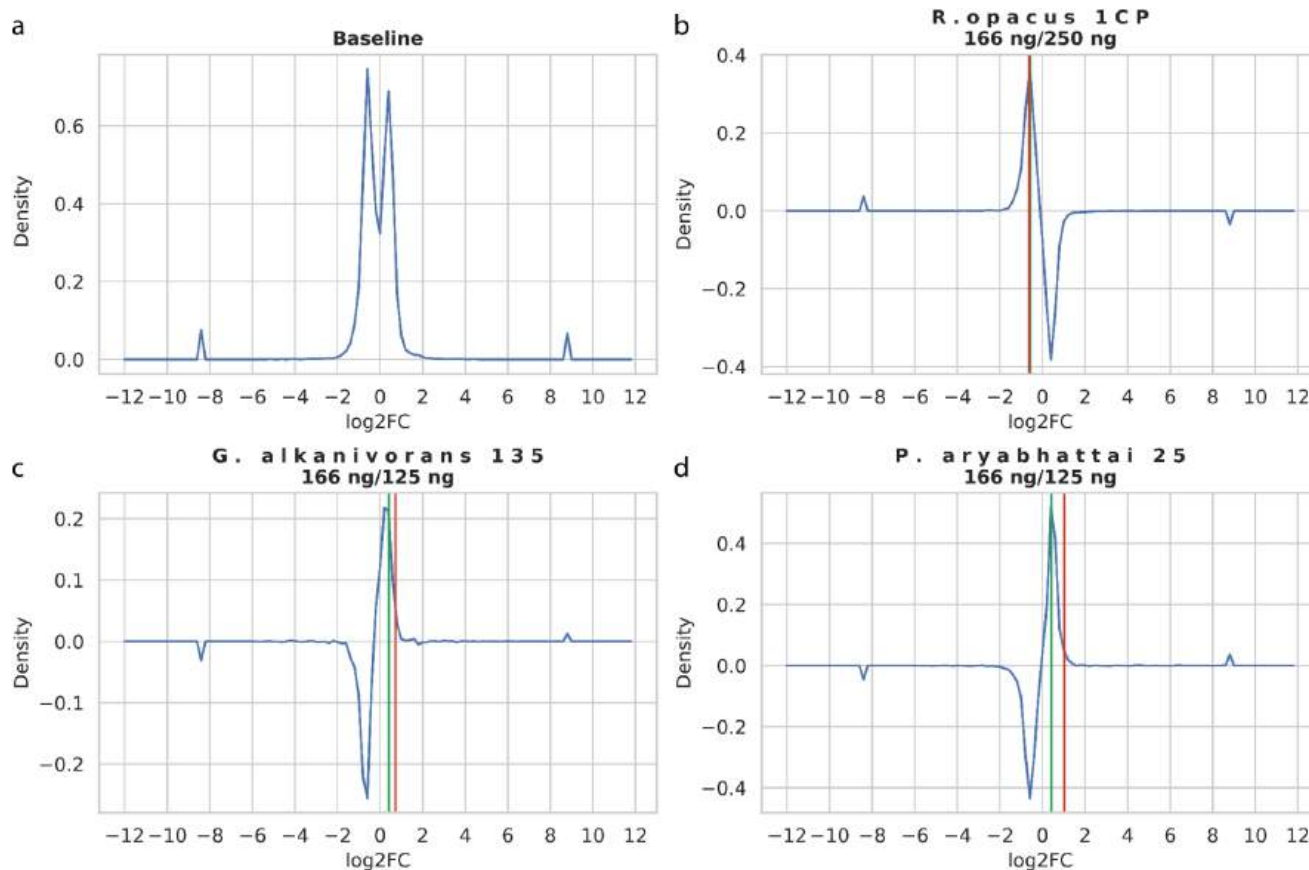


Fig. 2. The densities of \log_2FC distribution for M1.1/M1.2 comparison. The small outlier bins correspond to the cases when fold change was not estimated due to missing intensity values in one of the samples (they were imputed by the maximum (or minimum) fold change). The green line stands for the actual ratio of peptide masses (given in subtitles) and the red line stands for the calculated ratio, \log_2 -scaled. FC is the fold change in strain biomass between model microbiomes. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

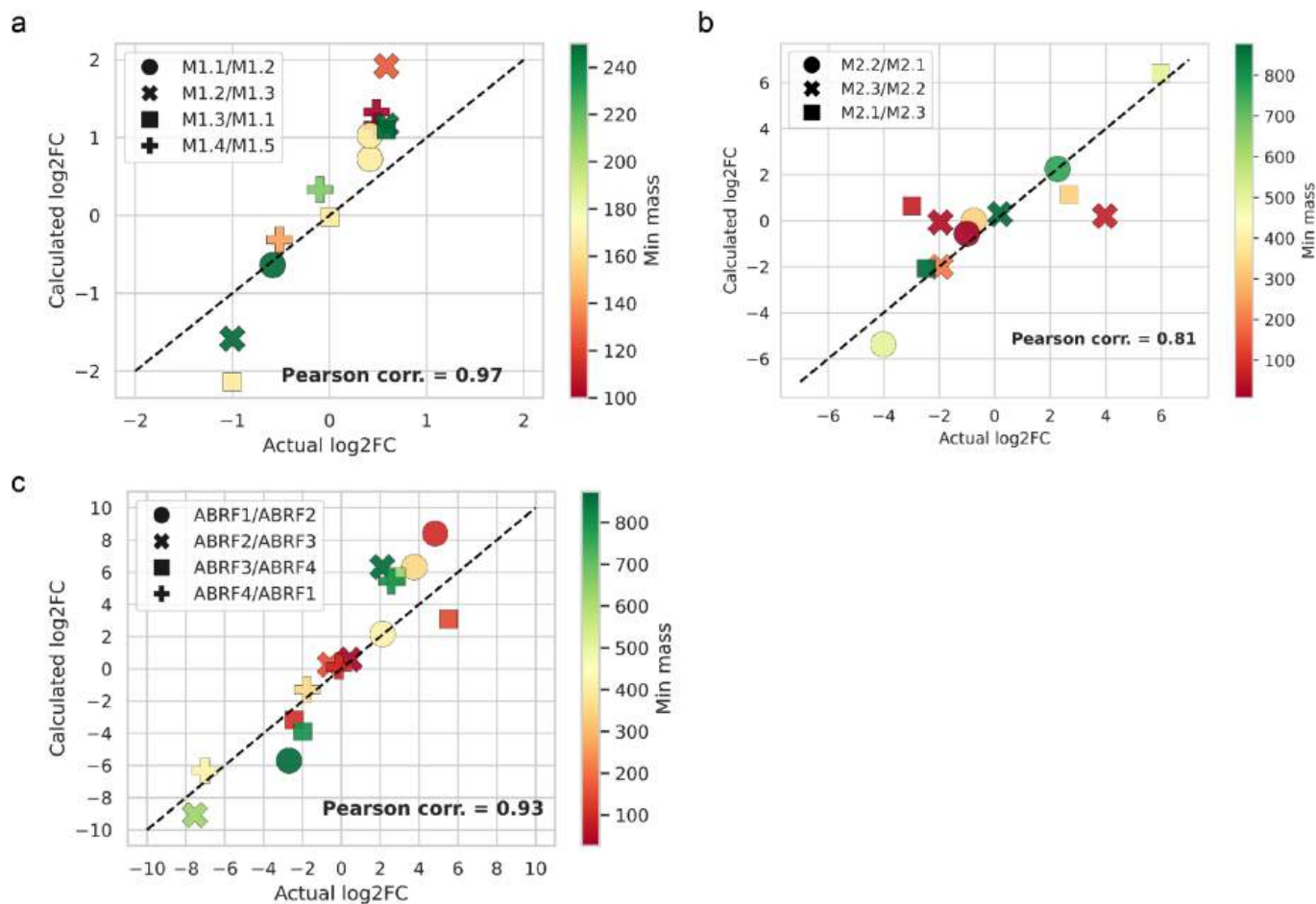


Fig. 3. Correlation of actual and experimentally measured fold changes in strain biomass for Model I (a), Model II (b), Model III (c). Y: fold change estimate from experimental data, X: actual ratios, in log₂-scale. Legends correspond to labels in Tables 3–5. The color corresponds to the maximum of two compared species biomasses, in nanograms.

in Fig. S4, for which there was no fold change estimation provided.

Another case we explored is the fold change predictions for mixtures with high overlap of the identified peptides coming from different species. For example, the fold change distributions for *G. alkanivorans* 135 and *G. amicalis* 6-1 from Model III ABRF2/ABRF3 coincide (Fig. S4) and the fold change estimate matches the sum of the peptide masses of both strains. This trend was observed for all mixtures analyzed. Next example is the Model II M2.2/M2.3 comparison demonstrating an absence (or undetectable concentrations) of unique peptides required for distinguishing the *R. opacus* strains, as well as *R. erythropolis* and *R. qingshengii* species (Fig. S5). From a practical viewpoint, it means that the measured fold change between the components of complex microbiomes matches the sum of the biomasses of all strains identified by the shared peptides.

3.4. DirectMS1 allows the characterization of changes in metabolic pathways and biodegradation ability of strain isolates

The performance of DirectMS1 to characterize metabolic response of bacteria to the stimulus was tested using two pollutant-degrading actinobacteria: (1) *Rhodococcus opacus* 1CP grown in presence of glucose, benzoate, phenol, and 4-chlorophenol; and (2) *Rhodococcus erythropolis* X5 grown at 28 °C and 6 °C in presence of *n*-hexadecane. Figs. 4 and 5 summarize the results of quantitation analysis.

Fig. 4 summarizes the key differences in the proteomes of *Rhodococcus opacus* 1CP grown in the presence of phenol, 4-chlorophenol and benzoate, compared with glucose. *R. opacus* 1CP is the actinobacterium

able to degrade phenol and their derivatives [34]. In our study, the presence of phenol resulted in upregulation of enzymes involved in phenol degradation: three different phenol hydroxylases ((**pheA1(1)**), **pheA1(2)**, and **pheA1(3)**), and 2-oxopent-4-enoate hydratase (**R1CP_00930**). These enzymes are also induced in the presence of benzoate and 4-chlorophenol, but with smaller magnitude (Table S1). The data obtained are in full agreement with previously obtained real-time PCR results [47]. Using specific primers for the small subunit of all three phenol hydroxylases of strain 1 CP growing in the presence of phenol, the gene activation has been shown; an increase for phe A1(3) was approximately 2000 times.

Enzyme cluster consisting of upregulated 3-methyl-2-oxobutanoate hydroxymethyltransferase (**panB**), catechol 1,2-dioxygenases (**catA1/catA2**), 4-nitrophenol 2-monooxygenase (**nphA1**), metapyrocatechase (**pheB**), and Fe-ADH domain-containing protein were specific to phenol. For degradation of 4-chlorophenol, we observed no specific features, the upregulated enzymes were the same as for phenol conditions. For benzoate conditions, we found specific pattern of enzymes degrading benzoate: benzene 1,2-dioxygenase (**A8787_2660**), benzoate 1,2-dioxygenase subunit alpha (**benA**), benzoate 1,2-dioxygenase small subunit (**benB**), 1,6-dihydroxycyclohexa-2, 4-diene-1-carboxylatedehydrogenase (**xylL**), and ferrienterobactin-binding periplasmic protein (**fepB**).

In total, the results of ultrafast proteomics profiling were consistent with the expected changes in metabolic activity and can serve as a valuable support in characterization of pollutant degradation activity of different microorganisms.

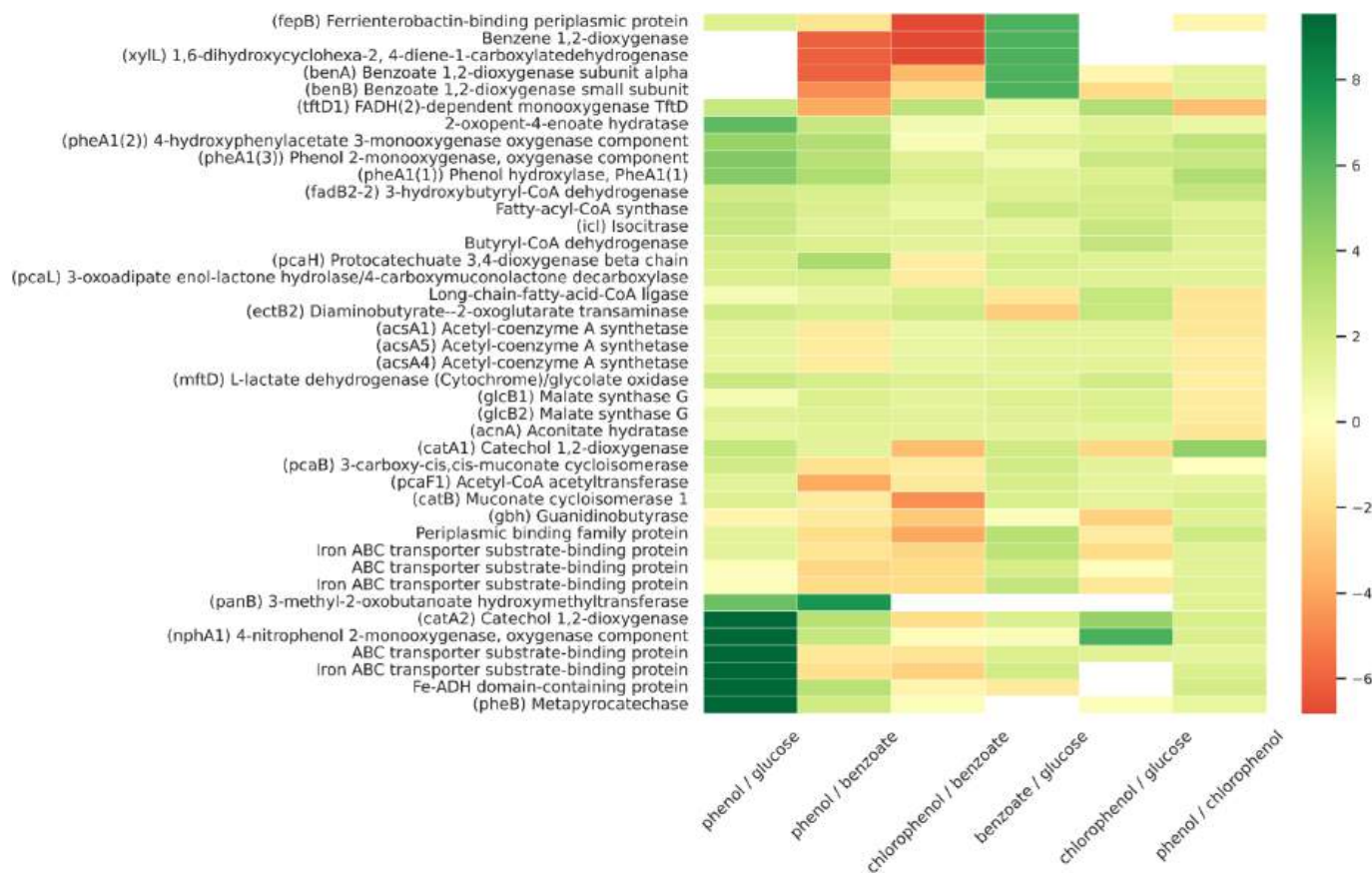


Fig. 4. Proteomic response of actinobacteria *Rhodococcus opacus* 1CP to the presence of aromatic compounds: benzoate, phenol, and 4-chlorophenol, compared with glucose. Quantitation: DirectMS1Quant [40]. Protein selection corresponds to the most enriched biological processes (GO score ≥ 6). Colorbar corresponds to \log_2FC .

Fig. 5 shows response of the bacterium *Rhodococcus erythropolis* X5 grown on *n*-hexadecane to temperature change from 6 °C to 28 °C. *R. erythropolis* X5 is the psychrotrophic bacterium that possesses a wide range of catabolic activities within 4 °C–28 °C temperature range. Specifically, this strain can degrade *n*-alkanes which constitute the main fraction of oil pollution [25]. In our study, the proteome of X5 strain grown on *n*-hexadecane at 6 °C was compared with its proteome at 28 °C. At 6 °C, we observed proteins related to all steps of *n*-alkane degradation including oxidation to alcohols (QEX08599, QEX09356, etc.), oxidation to aldehydes (QEX08400, QEX09399, etc.), oxidation to fatty acids and fatty acid metabolism (QEX08594, QEX12285, QEX08597, etc.), exopolysaccharide production (QEX09055, QEX10519, QEX12651, QEX11399, etc.), and iron transport (QEX13740, QEX13741, QEX09379, etc.) (Table S2, Appendix D). Besides enzymes involved in *n*-alkane biodegradation, the chaperons/chaperonins (QEX09859, QEX13151, QEX09986, QEX09823) and stress/protection proteins (QEX09757, QEX11651, QEX11169) were upregulated at 6 °C. Also, a significant loss of 30S/50S ribosomal proteins was observed at cold conditions (Fig. S6). Numerous studies report that loss of ribosomal proteins in bacteria is the reaction to stress and is necessary for survival under nutrient- or energy-limited conditions [48–51]. GO analysis revealed that translation, gene expression, and protein folding were significantly affected at cold conditions (Fig. 5a). Proteins involved in DNA and RNA processing, regulation of transcription and translation are shown in Fig. 5b and Appendix D.

At 28 °C, we observed the abundant protein group corresponding to degradation of *n*-alkanes, but expressed from other parts of genome: *n*-alkanes to alcohols (QEX08601, QEX08447), alcohols to aldehydes (QEX10266, QEX10225, QEX10036), aldehydes to fatty acids/fatty acid metabolism (QEX12404, QEX10265, QEX10256, and many others), exopolysaccharide production (QEX12160, QEX13966, QEX14217,

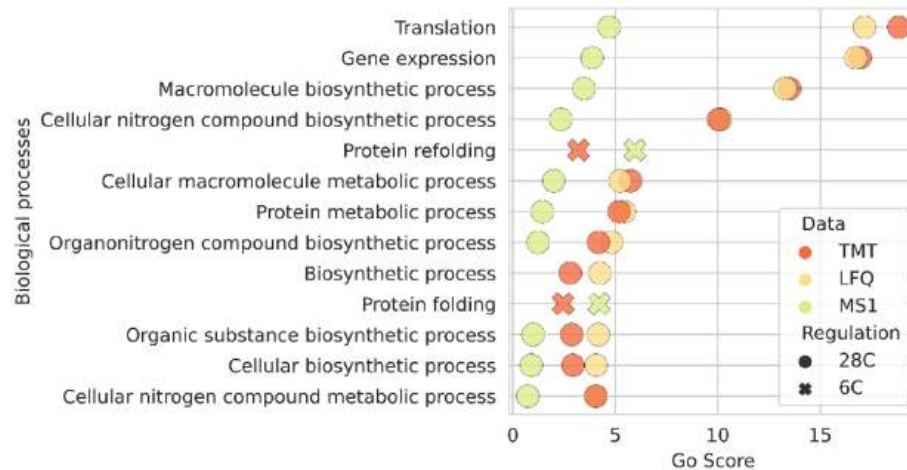
etc.), and iron transport (QEX11646, QEX12370, QEX11669, etc.) (Table S2, Appendix D). This observation suggests that different parts of the *R. erythropolis* X5 genome were transcribed to translate the enzymes needed for *n*-alkane degradation under optimal and cold conditions. This transcriptional behavior can be a conserved mechanism of coping with stress.

To confirm the technical reproducibility of the results, a comparison was made between label-free DirectMS1, DDA LFQ and TMT-based quantitation. Fig. 5a shows the enriched biological processes obtained using the three quantitation methods. The results of the analysis show that the set of enriched biological processes for all methods has a high overlap. The protein content behind the top GO enrichments, measured using different quantitation workflows, is shown in Fig. 5b. Fig. 5c compares the duration of the experimental methods for the used quantitation methods. DirectMS1 requires less experimental time providing comparable performance.

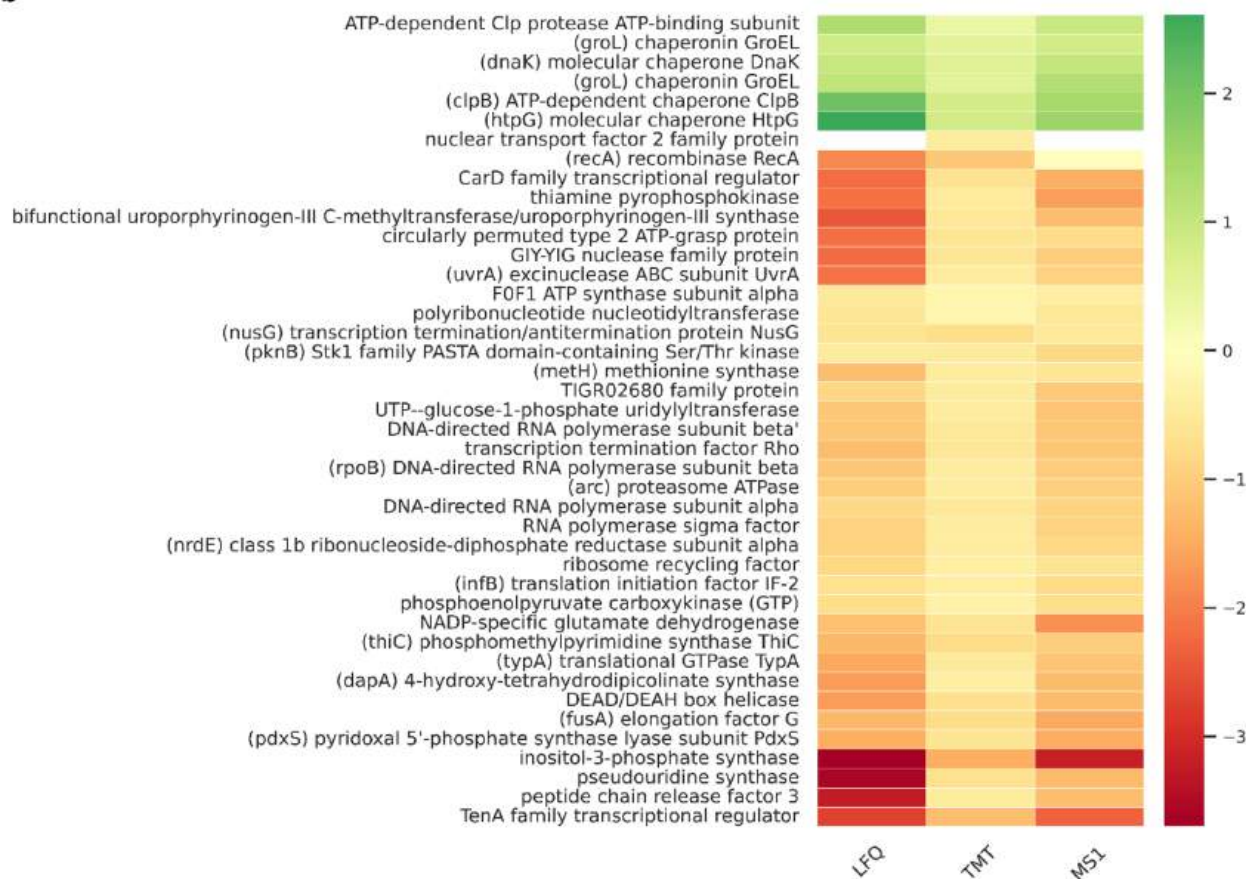
4. Discussion

We demonstrate that the DirectMS1 method for fast proteomic profiling can be effectively used for functional characterization of strain isolates and quantitative assessment of microbial communities. Our algorithm for two-stage blind database search discriminated bacteria to species level from LC-MS1 data with high accuracy of 95 %, and 32 % of those cases were matched to the correct strains. The algorithm identifies the composition of microbial communities at the level of genus and provides quantitative estimates in strain biomass. DirectMS1 for metaproteomics takes a position between MALDI TOF MS-based identification of microorganisms and LC-MS/MS-based proteomics, in terms of fast identification and quantitative assessment of strain isolates and microbiomes at reasonable costs.

a



b



c

Method	Time per run, min	Total time, min
DirectMS1, LC-MS1 (MS1)	7.5	45
Label Free Quantitation, DDA (LFQ)	87	522
Tandem mass tag labeling, DDA (TMT)	137	137

Fig. 5. Response of actinobacteria *Rhodococcus erythropolis* X5 grown on *n*-hexadecane to temperature change from 6 °C to 28 °C: (a) enriched biological processes; (b) \log_2FC estimated using different quantitation methods; (c) instrument time required for the compared workflows. Quantitation: Diffacto [42] (LFQ, MS1-only acquisition), MSFragger + Scavager + NSAF (LFQ, DDA), and MSFragger + Scavager + Diffacto (DDA, TMT labeling). GO Score = $E * \log_{10}FDR$, where E is the enrichment of biological processes, FDR is the statistical significance of the GO enrichment corrected for multiple comparisons using Benjamini-Hochberg method [52].

MALDI-TOF MS, 16S rRNA gene sequencing, and whole-genome sequencing are widely used for microorganism identification [53–56]. Accuracy of any identification method depends on the number of strain-specific features measured: the larger the number of features, the higher the level of discrimination. However, closely related species (*Yersinia pestis*, *Yersinia pseudotuberculosis*, *Yersinia wautersii* [57–59], *Citrobacter* strains [54], *Bacillus cereus* and *Bacillus anthracis* [46], *Escherichia coli*, *Escherichia fergusonii* and *Shigella* [45,60,61] with nearly identical genomes and 16S rRNA genes, are difficult to differentiate using either method. This suggests that combining techniques, including whole genome sequencing and LC-MS-based ultrafast proteomics, should be beneficial for accurate identification of closely related strains.

The 5-minute LC-MS1-based proteome profiling of strain isolates responding to environmental conditions demonstrated high performance in quantitation (Figs. 4 and 5). We detected the key enzymes degrading aromatic compounds, such as phenol, benzoate and 4-chlorophenol, that are in full agreement with previous studies [62,63]. We measured the key enzymes involved in the oxidation of *n*-alkanes to alcohols (i.e. alkane monooxygenases and cytochrome P450 family proteins); alcohols oxidation to aldehydes (alcohol dehydrogenases); aldehydes oxidation to fatty acids (aldehyde dehydrogenases); fatty acid metabolism (acyl-CoA dehydrogenases, fatty acid-CoA synthases, fatty acid synthesis proteins); exopolysaccharide production (glycosyl-transferases, glucose dehydrogenases, fructose synthases, mannose dehydrogenases, etc.); and iron transport (ferredoxins, ferredoxin reductases, ABC transporter proteins) [64]. We conclude that DirectMS1 was proved as a valuable approach to functional annotation of strain isolates responding to the environment.

CRedit authorship contribution statement

Elizaveta Kazakova: Visualization, Validation, Software, Investigation, Formal analysis, Data curation. **Mark Ivanov:** Writing – review & editing, Software, Methodology, Investigation. **Tomiris Kusainova:** Visualization, Resources, Investigation, Data curation. **Julia Bubis:** Validation, Methodology, Investigation. **Valentina Polivtseva:** Writing – review & editing, Resources, Investigation. **Kirill Petrikov:** Writing – review & editing, Resources, Investigation. **Vladimir Gorshkov:** Writing – review & editing, Investigation. **Frank Kjeldsen:** Writing – review & editing, Investigation. **Mikhail Gorshkov:** Writing – review & editing, Project administration, Funding acquisition. **Yanina Delean:** Writing – review & editing, Methodology, Investigation. **Inna Solyanikova:** Writing – review & editing, Resources, Investigation. **Irina Tarasova:** Writing – review & editing, Writing – original draft, Supervision, Project administration, Methodology, Conceptualization.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

LC-MS1 and LC-MS/MS data are available at proteomexchange.org (PXD050587, PXD050761, PXD050807, PXD050887). Code for blind search, identification of microbial sample composition and assessment of species biomass change between samples is available at GitHub (<https://github.com/kazakova/Metaproteomics-DirectMS1>).

Acknowledgements

I.T. and M.G. thank Prof. Victor Zgoda from “Human Proteome” Core Facility at the Institute of Biomedical Chemistry for help with implementation of DirectMS1 method on the Orbitrap FTMS system.

Funding

The study (developing DirectMS1 and bioinformatic solutions for metaproteome profiling, study design, sample preparation, data processing, etc.) was supported by the Russian Science Foundation grant no. 20-14-00229.

F.K. and V.G. acknowledge generous grants to the VILLUM Center for Bioanalytical Sciences (VILLUM Foundation grant no. 7292), PRO-MS: Danish National Mass Spectrometry Platform for Functional Proteomics (grant no. 5072-00007B), and the Novo Nordisk Foundation (INTEGRA, NNF20OC0061575) for support of instrumentation infrastructure at the University of Southern Denmark.

Appendix A. Supplementary material

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.microc.2024.111823>.

References

- [1] A. Ascandari, S. Aminu, N.E.H. Safdi, A. El Allali, R. Daoud, A bibliometric analysis of the global impact of metaproteomics research, *Front. Microbiol.* 14 (2023).
- [2] M.M. Kleiner, Much more than measuring gene expression in microbial communities, *mSystems* 4 (2019) e00115-19.
- [3] T. Van Den Bossche, et al., The Metaproteomics Initiative: a coordinated approach for propelling the functional characterization of microbiomes, *Microbiome* 9 (2021) 243.
- [4] N. Miura, S. Okuda, Current progress and critical challenges to overcome in the bioinformatics of mass spectrometry-based metaproteomics, *Comput. Struct. Biotechnol. J.* 21 (2023) 1140–1150.
- [5] T. Van Den Bossche, et al., Critical Assessment of MetaProteome Investigation (CAMP): a multi-laboratory comparison of established workflows, *Nat. Commun.* 12 (2021) 7305.
- [6] F. Baquero, C. Nombela, The microbiome as a human organ, *Clin. Microbiol. Infect. off. Publ. Eur. Soc. Clin. Microbiol. Infect. Dis.* 18 (Suppl 4) (2012) 2–4.
- [7] S. Long, et al., Metaproteomics characterizes human gut microbiome function in colorectal cancer, *npj Biofilms Microbiomes* 6 (2020) 14.
- [8] J. Zhao, et al., Data-independent acquisition boosts quantitative metaproteomics for deep characterization of gut microbiota, *npj Biofilms Microbiomes* 9 (2023) 4.
- [9] S. Pietilä, T. Suomi, L.L. Elo, Introducing untargeted data-independent acquisition for metaproteomics of complex microbial samples, *ISME Commun.* 2 (2022) 51.
- [10] R. Lou, W. Shui, Acquisition and analysis of DIA-based proteomic data: a comprehensive survey in 2023, *Mol. Cell. Proteomics MCP* 23 (2024) 100712.
- [11] C.M.A. Simopoulos, et al., MetaProClust-MS1: an MS1 profiling approach for large-scale microbiome screening, *mSystems* 7 (2022) e0038122.
- [12] S. Schubert, M. Kostrzewa, MALDI-TOF MS in the microbiology laboratory: current trends, *Curr. Issues Mol. Biol.* 23 (2017) 17–20.
- [13] T.R. Sandrin, J.E. Goldstein, S. Schumaker, MALDI TOF MS profiling of bacteria at the strain level: a review, *Mass Spectrom. Rev.* 32 (2013) 188–217.
- [14] P. Lasch, A. Schneider, C. Blumenschein, J. Doellinger, Identification of microorganisms by liquid chromatography-mass spectrometry (LC-MS1) and in silico peptide mass libraries, *Mol. Cell. Proteomics MCP* 19 (2020) 2125–2139.
- [15] R. Heyer, et al., Challenges and perspectives of metaproteomic data analysis, *J. Biotechnol.* 261 (2017) 24–36.
- [16] P. Jagtap, et al., A two-step database search method improves sensitivity in peptide sequence matches for metaproteomics and proteogenomics studies, *Proteomics* 13 (2013) 1352–1357.
- [17] A. Bassignani, et al., Benefits of iterative searches of large databases to interpret large human gut metaproteomic data sets, *J. Proteome Res.* 20 (2021) 1522–1534.
- [18] P. Kumar, et al., A sectioning and database enrichment approach for improved peptide spectrum matching in large, genome-guided protein sequence databases, *J. Proteome Res.* 19 (2020) 2772–2785.
- [19] X. Zhang, et al., MetaPro-IQ: a universal metaproteomic approach to studying human and mouse gut microbiota, *Microbiome* 4 (2016) 31.
- [20] J. Xiao, et al., Metagenomic taxonomy-guided database-searching strategy for improving metaproteomic analysis, *J. Proteome Res.* 17 (2018) 1596–1605.
- [21] D. Beyter, M.S. Lin, Y. Yu, R. Pieper, V. Bafna, ProteoStorm: an ultrafast metaproteomics database search framework, *Cell Syst.* 7 (2018) 463–467.e6.
- [22] T. Muth, et al., The MetaProteomeAnalyzer: a powerful open-source software suite for metaproteomics data analysis and interpretation, *J. Proteome Res.* 14 (2015) 1557–1565.
- [23] M.V. Ivanov, et al., DirectMS1: MS/MS-free identification of 1000 proteins of cellular proteomes in 5 minutes, *Anal. Chem.* 92 (2020) 4326–4333.
- [24] M.V. Ivanov, et al., Boosting MS1-only proteomics with machine learning allows 2000 protein identifications in single-shot human proteome analysis using 5 min HPLC gradient, *J. Proteome Res.* 20 (2021) 1864–1873.
- [25] Y. Delean, et al., Complete genome sequence of *Rhodococcus erythropolis* X5, a psychrotrophic hydrocarbon-degrading biosurfactant-producing bacterium, *Microbiol. Resour. Announc.* 8 (2019) e01234-19.

- [26] Y. Delegan, et al., Complete genome analysis of *Rhodococcus opacus* S8 capable of degrading alkanes and producing biosurfactant reveals its genetic adaptation for crude oil decomposition, *Microorganisms* 10 (2022) 1172.
- [27] L. Iminova, et al., Physiological and biochemical characterization and genome analysis of *Rhodococcus qingshengii* strain 7B capable of crude oil degradation and plant stimulation, *Biotechnol. Rep. Amst. Neth.* 35 (2022) e00741.
- [28] Y. Delegan, et al., Complete genome sequence of *Rhodococcus qingshengii* VT6, a promising degrader of persistent pollutants and putative biosurfactant-producing strain, *Microbiol. Resour. Announc.* 11 (2022) e0117921.
- [29] T.Z. Esikova, T.O. Anokhina, T.N. Abashina, N.E. Suzina, I.P. Solyanikova, Characterization of soil bacteria with potential to degrade benzoate and antagonistic to fungal and bacterial phytopathogens, *Microorganisms* 9 (2021) 755.
- [30] E. Frantsuzova, Y. Delegan, A. Bogun, D. Sokolova, T. Nazina, Comparative genomic analysis of the hydrocarbon-oxidizing dibenzothiophene-desulfurizing *Gordonia* strains, *Microorganisms* 11 (2022) 4.
- [31] Y. Delegan, L. Valentovich, A. Vetrova, E. Frantsuzova, Y. Kocharovskaya, Complete genome sequence of *Gordonia* sp. 135, a promising dibenzothiophene- and hydrocarbon-degrading strain, *Microbiol. Resour. Announc.* 9 (2020) e0145019.
- [32] Y. Delegan, Y. Kocharovskaya, E. Frantsuzova, R. Streletskii, A. Vetrova, Characterization and genomic analysis of *Gordonia alkanivorans* 135, a promising dibenzothiophene-degrading strain, *Biotechnol. Rep. Amst. Neth.* 29 (2021) e00591.
- [33] E. Frantsuzova, et al., Complete genome sequence of *Gordonia polyisoprenivorans* 135, a promising degrader of aromatic compounds, *Microbiol. Resour. Announc.* 12 (2023) e0005823.
- [34] E.V. Emelyanova, I.P. Solyanikova, Evaluation of phenol-degradation activity of *Rhodococcus opacus* 1CP using immobilized and intact cells, *Int. J. Environ. Sci. Technol.* 17 (2020) 2279–2294.
- [35] T.O. Anokhina, et al., Alternative naphthalene metabolic pathway includes formation of ortho-phthalic acid and cinnamic acid derivatives in the *Rhodococcus opacus* Strain 3D, *Biochem. Biokhimiia* 85 (2020) 355–368.
- [36] T. Tatusova, et al., NCBI prokaryotic genome annotation pipeline, *Nucleic Acids Res.* 44 (2016) 6614–6624.
- [37] E. Coudert, et al., Annotation of biologically relevant ligands in UniProtKB using ChEBI, *Bioinform. Oxf. Engl.* 39 (2023) btac793.
- [38] M. Choi, et al., ABRF Proteome Informatics Research Group (iPRG) 2015 study: detection of differentially abundant proteins in label-free quantitative LC-MS/MS experiments, *J. Proteome Res.* 16 (2017) 945–957.
- [39] D.A. Abdrakhimov, et al., Biosaur: an open-source Python software for liquid chromatography-mass spectrometry peptide feature detection with ion mobility support, *Rapid Commun. Mass Spectrom.* RCM 20 (2021) e9045.
- [40] M.V. Ivanov, et al., DirectMS1Quant: ultrafast quantitative proteomics with MS/MS-free mass spectrometry, *Anal. Chem.* 94 (2022) 13068–13075.
- [41] D. Szklarczyk, et al., STRING v10: protein-protein interaction networks, integrated over the tree of life, *Nucleic Acids Res.* 43 (2015) D447–D452.
- [42] B. Zhang, M. Pirmoradian, R. Zubarev, L. Käll, Covariation of peptide abundances accurately reflects protein concentration differences, *Mol. Cell. Proteomics MCP* 16 (2017) 936–948.
- [43] E.M. Kazakova, et al., Proteomics-based scoring of cellular response to stimuli for improved characterization of signaling pathway activity, *Proteomics* 23 (2023) e2200275.
- [44] J. Huerta-Cepas, F. Serra, P. Bork, ETE 3: reconstruction, analysis, and visualization of phylogenomic data, *Mol. Biol. Evol.* 33 (2016) 1635–1638.
- [45] A. Paauw, et al., Rapid and reliable discrimination between *Shigella* species and *Escherichia coli* using MALDI-TOF mass spectrometry, *Int. J. Med. Microbiol. IJMM* 305 (2015) 446–452.
- [46] P. Lasch, et al., Identification of highly pathogenic microorganisms by matrix-assisted laser desorption ionization-time of flight mass spectrometry: results of an interlaboratory ring trial, *J. Clin. Microbiol.* 53 (2015) 2632–2640.
- [47] N.S. Egozarian, et al., Removal of phenol by *Rhodococcus opacus* 1CP after dormancy: insight into enzymes' induction, specificity, and cells viability, *Microorganisms* 12 (2024) 597.
- [48] H. Cheng-Guang, C.O. Gualerzi, The ribosome as a switchboard for bacterial stress response, *Front. Microbiol.* 11 (2020) 619038.
- [49] R. Njenga, J. Boele, Y. Öztürk, H.-G. Koch, Coping with stress: how bacteria fine-tune protein synthesis and protein transport, *J. Biol. Chem.* 299 (2023) 105163.
- [50] S. Zhang, J.M. Scott, W.G. Haldenwang, Loss of ribosomal protein L11 blocks stress activation of the *Bacillus subtilis* transcription factor sigma(B), *J. Bacteriol.* 183 (2001) 2316–2321.
- [51] W.M. El-Sharoud, Ribosome inactivation for preservation: concepts and reservations, *Sci. Prog.* 87 (2004) 137–152.
- [52] Y. Benjamini, Y. Hochberg, Controlling the false discovery rate: a practical and powerful approach to multiple testing, *J. R. Stat. Soc. B. Methodol.* 57 (1995) 289–300.
- [53] R. Franco-Duarte, et al., Advances in chemical and biological methods to identify microorganisms-from past to present, *Microorganisms* 7 (2019) 130.
- [54] H.-L. Kwak, et al., Development of a rapid and accurate identification method for citrobacter species isolated from pork products using a matrix-assisted laser-desorption ionization time-of-flight mass spectrometry (MALDI-TOF MS), *J. Microbiol. Biotechnol.* 25 (2015) 1537–1541.
- [55] M.T. Caudill, K.A. Brayton, The use and limitations of the 16S rRNA sequence for species classification of anaplasma samples, *Microorganisms* 10 (2022) 605.
- [56] M. Vargha, Z. Takáts, A. Konopka, C.H. Nakatsu, Optimization of MALDI-TOF MS for strain level differentiation of *Arthrobacter* isolates, *J. Microbiol. Methods* 66 (2006) 399–409.
- [57] B. Feng, et al., Effective discrimination of *Yersinia pestis* and *Yersinia pseudotuberculosis* by MALDI-TOF MS using multivariate analysis, *Talanta* 234 (2021) 122640.
- [58] K. Trebesius, D. Harmsen, A. Rakin, J. Schmelz, J. Heesemann, Development of rRNA-targeted PCR and in situ hybridization with fluorescently labelled oligonucleotides for detection of *Yersinia* species, *J. Clin. Microbiol.* 36 (1998) 2557–2564.
- [59] C. Savin, et al., Genus-wide *Yersinia* core-genome multilocus sequence typing for species identification and strain characterization, *Microb. Genomics* 5 (2019) e000301.
- [60] R.H. Dahal, Y.-J. Choi, S. Kim, J. Kim, Differentiation of *Escherichia fergusonii* and *Escherichia coli* isolated from patients with inflammatory bowel disease/ischemic colitis and their antimicrobial susceptibility patterns, *Antibiot. Basel Switz.* 12 (2023) 154.
- [61] R. Liu, et al., Genomic characterization of two *Escherichia fergusonii* isolates harboring *mcr-1* gene from farm environment, *Front. Cell. Infect. Microbiol.* 12 (2022) 774494.
- [62] J. Nešvera, L. Rucká, M. Pátek, Catabolism of phenol and its derivatives in bacteria: genes, their regulation, and use in the biodegradation of toxic pollutants, *Adv. Appl. Microbiol.* 93 (2015) 107–160.
- [63] L. Rucká, J. Nešvera, M. Pátek, Biodegradation of phenol and its derivatives by engineered bacteria: current knowledge and perspectives, *World J. Microbiol. Biotechnol.* 33 (2017) 174.
- [64] K. Lacz, et al., Metabolic responses of *Rhodococcus erythropolis* PR4 grown on diesel oil and various hydrocarbons, *Appl. Microbiol. Biotechnol.* 99 (2015) 9745–9759.