

Technology Developments for Characterizing the Complex Protein Datasets of Microbial Communities

Christine Shook^{1,3}; Melissa Thompson^{2,3}; Nathan VerBerkmoes³; Abijeet Borole³; Manesh Shah³; Frank Larimer³; Brian Davison³; Robert Hettich³

¹University of Tennessee, Knoxville, TN; ²ORNL-UTK GST Graduate School, Oak Ridge, TN; ³Oak Ridge National Lab, Oak Ridge, TN

OVERVIEW

- Determining the physiological states of microbes within complex communities and environments is an important analytical and biological problem for MS-based proteomics.
- Recent studies have shown that the genomic sequence of major microbial constituents of communities can be mapped out in both simple microbial communities (Tyson, Nature 2004) and more complex environmental samples (Venter, Science 2004).
- The next step is characterizing protein complexes and mapping the proteome of both the individual organism and the organism as a member of a community combining community genomics and proteomics for a proteogenomic approach (Ram, Science 2005).
- “Shotgun” proteomics, directly assessing the protein composition of complex mixtures, is currently the method of choice.
- The goal of this study is to assess and improve MS-based “shotgun” proteomics technology for analyzing simple known microbial mixtures. The purpose is to improve current techniques and develop analytical procedures for analyzing highly complex communities and environmental samples.

INTRODUCTION

Recently, there have been strong efforts to develop techniques for genomic sequencing and annotation of microbial communities (metagenomics). With the potential of partial or near complete microbial genomes obtained from environmental samples, along with the rapid proliferation of isolated microbial genomes, systems biology in microbial communities by combining genomic, transcriptomic, proteomic and metabolic studies may be possible in the near future. Our current studies seek to develop and demonstrate metaproteome analysis techniques for mixed cultures, establishing the basis for “whole community proteomics”. Currently our major focus is to use controlled, simple microbial mixtures of four species (*E. coli*, *S. cerevisiae*, *R. palustris*, and *S. oneidensis*) to develop MS-based proteomics methods, as well as the proteome bioinformatics tools for detailed analysis of sequenced microbial communities. Our goal is to provide “deep” and “wide” proteome measurements of complex microbial mixtures. Some of the considerations for any microbial community proteome project are highlighted below.

Considerations for Proteome Analysis of Microbial Communities

- Level of DNA sequence availability and quality of annotation for species in that community
- Diversity and dynamic range of species in the community
- Quantity and quality of community available (i.e. how much total protein can be obtained)
- Inter- and intra- species relationship at the amino acid level

EXPERIMENTAL

Cell Growth and Production of Protein Fractions:

- The four microbes (Figure 1) used in this study (*E. coli*, *R. palustris*, *S. cerevisiae*, and *S. oneidensis*) were all grown individually and mixed after cell harvesting at appropriate concentrations based on wet cell weight. Figure 2 illustrates the different mixtures analyzed in this study.
- R. palustris* was used as the target organism and grown under three different metabolic states: photoheterotrophic, chemoheterotrophic, and nitrogen fixation and then mixed with the other three species at 25% (Mix 1). Comparisons were then made between nitrogen fixation vs. photoheterotrophic and chemoheterotrophic vs. photoheterotrophic.
- Mixed cell pellets were washed twice with Tris buffer and disrupted with sonication. A crude and membrane fraction were separated by ultracentrifugation at 100,000g for 1 hour. Protein fractions were quantified and digested with sequencing grade trypsin using standard protocol. The digested mixed proteomes were desalted via solid phase extraction and concentrated to ~10µg/µL and frozen at -80°C until analysis.

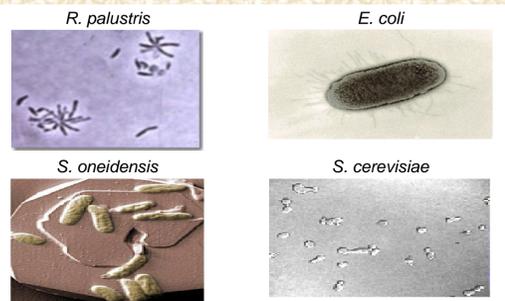


Figure 1 Test species for Artificial Mixture

Organism	Mix 1	Mix 2	Mix 3	Mix 4	Mix 5
<i>E. coli</i>	25%	32%	33%	33%	33%
<i>S. cerevisiae</i>	25%	32%	33%	33%	33%
<i>R. palustris</i>	25%	5%	1%	0.2%	0.0%

Figure 2 Concentrations in Artificial Mixtures

LC-MS/MS Analysis and Database Searching:

- The tryptic digestions of the crude and membrane fractions from all mixtures were analyzed via 24-hour two-dimensional LC-MS/MS with a split-phase MudPIT column as described in McDonald, W.H., JMS, 2002. The LC-MS system was composed of an Ultimate HPLC (LC Packings) and an LCQ-DECA XP ion trap mass spectrometer (Thermo Finnigan) or a Surveyor pump (Thermo Finnigan) and LTQ linear trapping quadrupole (Thermo Finnigan). Each analysis required ~500µg of starting protein material from either the crude or membrane fraction.
- The resultant MS/MS spectra from the mixed proteomes were searched with SEQUEST (Thermo Finnigan) against the four databases highlighted below. The resultant files were filtered and sorted with DTASelect (Tabb, Journal of Proteome Research 2004) and unique and non-unique peptide identifications were extracted with in-house perl scripts. Reverse database searches were employed to test false positive levels.

Databases Used to Search MS/MS Spectra

- DB4 - the four species in the mixtures:
 - E. coli*
 - R. palustris*
 - S. cerevisiae*
 - S. oneidensis*
- DB13 - the four species plus 1 plant and 8 other microbes
- DB large - 261 species with 1,011,612 total proteins
- DB4 rev - DB4 plus DB4 in reversed sequence

Initial Mixture Studies:

- The initial experiments for this study have focused on testing current 2D-LC-MS/MS methodologies to analyze an artificial 4 microbe mixture of *E. coli*, *R. palustris*, *S. cerevisiae* and *S. oneidensis*.
- Two aspects of community proteomics were analyzed in this study:
 - To determine at what level functionally meaningful results could be obtained from a target species (*R. palustris*) whose concentration was decreased over a range of 25% to 0.2%. All experiments were conducted with current 2D-LC-MS/MS technologies.
 - To test the impact of database size on the identification of unique peptides from the four microbial species in the sample. Unique peptides are defined as those peptides unique to a given protein and specific to a species in a given database.

- Figure 3 highlights the results from the 4 species database with decreasing concentrations of *R. palustris*.
- Figure 4 highlights the results from the 13 species database with decreasing concentrations of *R. palustris*.
- Figure 5 highlights results from the 261 species database of the 25% mixture.

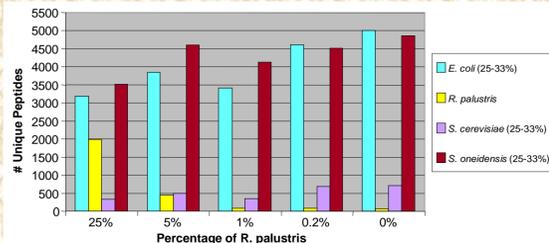


Figure 3 Unique Peptides Identified with DB4

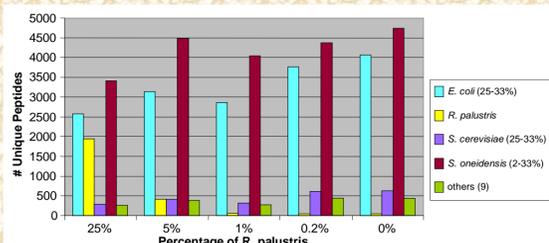


Figure 4 Unique Peptides Identified with DB13

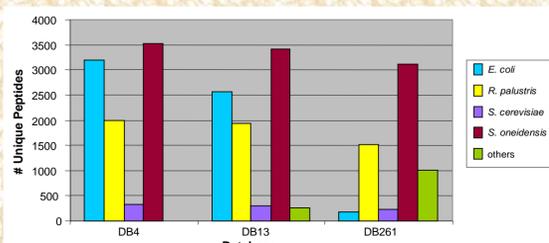


Figure 5 Database comparison of the 25% mixture

RESULTS AND DISCUSSION

R. Palustris Proteome from Initial Mixture Studies

- The functional categories of the identified proteome from *R. palustris* analyzed from the different mixtures were compared with our *R. palustris* baseline proteome (VerBerkmoes et al. submitted J. Bacteriology). Figure 6 illustrates this comparison: top-left pie chart is *R. palustris* baseline proteome and the other three pie charts are the 25%, 5%, & 1% mixtures.
- Figure 7 highlights the results from a single protein identified in the mixtures. This protein puhA is a critical subunit of the photoreaction center and is indicative of photosynthesis.
- Figures 8 and 9 illustrate diagnostic tandem mass spectra of peptides from puhA at the different concentrations of *R. palustris* in the sample, ranging down to 0.2%.

Figure 6 *R. palustris* Functional Categories

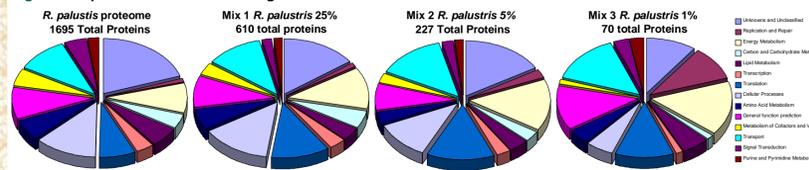


Figure 7 PuhA, a Component of the Photosynthetic Reaction Center from *R. palustris*

Mix	Charge State	Xcorr	DelCN	Peptide
Mix 1 44.70%	3	4.1787	0.5624	K.TVPSTNSDRPNVALTPAAPWPGAPFVPTGNPFADGVGPGSYAQR.A
•	3	5.8592	0.4355	R.ADVPELGLDNLPIVPLR.A
•	2	5.6884	0.5109	R.ADVPELGLDNLPIVPLR.A
•	2	3.788	0.4484	R.ADVPELGLDNLPIVPLRAAK.G
•	1	2.7334	0.2724	R.YLEVEAKS
•	4	4.0343	0.5105	R.VLLPVFPALINDPFGK.V
•	3	5.8444	0.5249	R.VLLPVFPALINDPFGK.V
•	2	3.2296	0.6166	R.VLLPVFPALINDPFGKVSVDAIR.G
•	2	4.5225	0.5718	K.VSVDAIRGDDQFAGVPTTSKGDQVSK.L
•	3	3.8473	0.3707	K.VSVDAIRGDDQFAGVPTTSKGDQVSK.L
Mix 2 42.70%	2	2.5409	0.3636	K.IGVPPADPK.T
•	3	3.5846	0.4472	K.TVPSTNSDRPNVALTPAAPWPGAPFVPTGNPFADGVGPGSYAQR.A
•	3	5.8143	0.4663	R.ADVPELGLDNLPIVPLR.A
•	2	5.6008	0.4809	R.ADVPELGLDNLPIVPLR.A
•	2	4.1629	0.369	R.ADVPELGLDNLPIVPLRAAK.G
•	2	2.9884	0.485	R.VLLPVFPALINDPFGK.V
•	2	3.7545	0.4465	R.GDQFAGVPTTSKGDQVSK.L
Mix 3 17.30%	3	5.1517	0.583	K.TVPSTNSDRPNVALTPAAPWPGAPFVPTGNPFADGVGPGSYAQR.A
Mix 4 7.10%	2	4.4881	0.4467	R.ADVPELGLDNLPIVPLR.A

Figure 8 Diagnostic Peptide at 1% *R. palustris*

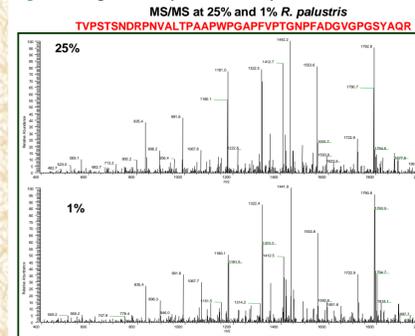
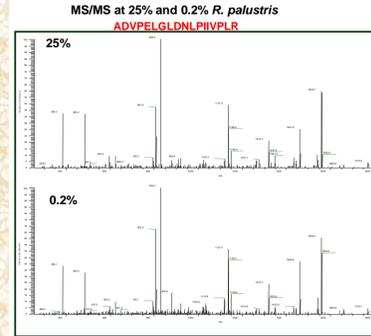


Figure 9 Diagnostic Peptide at 0.2% *R. palustris*



Comparison of *R. palustris* Metabolic States in a Simple Microbial Mixture

- Experiments were performed to test the ability to detect differential metabolic states of *R. palustris* in the presence of the three other microbes (all at 25%):
 - Aerobic (Chemoheterotrophic)
 - Anaerobic (Photoheterotrophic)
 - Nitrogen Fixation
- Table 1a: number of proteins identified at the 1- and 2-peptide level in all three growth states
- Table 1b: proteins identified for each species in the LTQ Nitrogen Fixation experiment
- Table 2 number of unique peptides found in the nitrogen fixation state with the associated percentage of false positives identified using the reverse database (DB4 rev). Results for the other states were similar.
- Table 3: Differential expression of some nitrogen fixation proteins from the nitrogen fixation vs. photoheterotrophic metabolic states.
 - strongest differences highlighted in red
 - blue highlighting represents known nitrogen fixation proteins
- Table 4 highlights unique peptides detected from the NifK beta-chain protein.
- Requirements for differentially expressed proteins: at least 4 unique peptides or more than 30% sequence coverage difference

Table 1a Protein Totals with Conservative and Liberal Filters

Metabolic State	Experiment	1 Peptide	2 Peptide
Photoheterotrophic	LCQ Run 1	2679	1448
	LCQ Run 2	2593	1357
	LTQ	5391	3339
Chemoheterotrophic	LCQ Run 1	2615	1335
	LCQ Run 2	2306	1180
	LTQ	5517	3354
Nitrogen Fixation	LCQ Run 1	2588	1420
	LCQ Run 2	2164	1177
	LTQ	5762	3650

Table 1b Species Protein Counts for Nitrogen Fixation LTQ

Nitrogen Fixation	<i>E. coli</i>	1182
	<i>R. pal</i>	1274
	<i>S. oneidensis</i>	763
	<i>S. cerevisiae</i>	431
	Total	3650

Table 3 Nitrogen Fixation vs. Photoheterotrophic Activated Genes

Gene	LCQ1	LCQ2	LTQ	LCQ1	LCQ2	LTQ	Function
RP4274	0	0	0	38.4	38.4	31.4	glnH2 GlnK, nitrogen regulatory protein P-II
RP41927	0	0	0	32	32	57.7	hypothetical protein
RP41928	0	0	0	20.5	20.5	61.6	Herediton-like protein (25-25)
RP4225	0	26.4	31.8	31.8	68.2	69.1	conserved unknown protein
RP4293	0	0	9.1	0	0	45.5	lexA SOS response transcriptional repressor lexA
RP42011	0	0	9.9	0	0	58.5	unknown protein
RP4369	0	0	0	28.6	28.6	45.2	amino acid uptake ABC transporter periplasmic solute-binding protein
RP4409	0	0	0	12.5	12.5	27.5	glnH II glutamine synthetase II
RP4402	0	0	0	19.4	19.4	19.4	FluX heterodimer like protein, fluX
RP4403	0	0	0	23.4	23.4	52	FluX nitrogen fixation protein, fluX
RP4404	0	0	0	13.6	13.6	25.5	FluX electron transfer flavoprotein alpha chain protein fluX
RP4405	0	0	0	13.9	13.9	59.0	FluX electron transfer flavoprotein beta chain fluX
RP4408	0	0	0	0	11.2	11.2	nifS2 nitrogenase cofactor synthesis protein nifS
RP4414	0	0	0	14.9	14.9	48.4	DUF295
RP4415	0	0	0	5.3	0	33.3	nifK nitrogenase molybdenum-iron protein nifK
RP4416	0	0	0	45.9	45.9	61.1	nifK nitrogenase molybdenum-iron protein beta chain, nifK
RP4419	0	0	0	37	37	63.2	nifD nitrogenase molybdenum-iron protein alpha chain, nifD
RP4420	0	0	0	49.3	49.3	60.7	nifH nitrogenase iron protein, nifH
RP4423	0	0	0	68.2	68.2	69.7	fluX, nif conserved hypothetical protein

Table 4 Unique Peptides Identified in NifK

RP4416	RP4416	RP4416	RP4416	RP4416	RP4416
1	1	1	1	1	1
2	2	2	2	2	2
3	3	3	3	3	3
4	4	4	4	4	4
5	5	5	5	5	5
6	6	6	6	6	6
7	7	7	7	7	7
8	8	8	8	8	8
9	9	9	9	9	9
10	10	10	10	10	10
11	11	11	11	11	11
12	12	12	12	12	12
13	13	13	13	13	13
14	14	14	14	14	14
15	15	15	15	15	15
16	16	16	16	16	16
17	17	17	17	17	17
18	18	18	18	18	18
19	19	19	19	19	19
20	20	20	20	20	20
21	21	21	21	21	21
22	22	22	22	22	22
23	23	23	23	23	23
24	24	24	24	24	24
25	25	25	25	25	25
26	26	26	26	26	26
27	27	27	27	27	27
28	28	28	28	28	28
29	29	29	29	29	29
30	30	30	30	30	30
31	31	31	31	31	31
32	32	32	32	32	32
33	33	33	33	33	33
34	34	34	34	34	34
35	35	35	35	35	35
36	36	36	36	36	36
37	37	37	37	37	37
38	38	38	38	38	38
39	39	39	39	39	39
40	40	40	40	40	40
41	41	41	41	41	41
42	42	42	42	42	42
43	43	43	43	43	43
44	44	44	44	44	44
45	45	45	45	45	45
46	46	46	46	46	46
47	47	47	47	47	47
48	48	48	48	48	48
49	49	49	49	49	49
50	50	50	50	50	50
51	51	51	51	51	51
52	52	52	52	52	52
53	53	53	53	53	53
54	54	54	54	54	54
55	55	55	55	55	55
56	56	56	56	56	56
57	57	57	57	57	57
58	58	58	58	58	58
59	59	59	59	59	59
60	60	60	60	60	60
61	61	61	61	61	61
62	62	62	62	62	62
63	63	63	63	63	63
64	64	64	64	64	64
65	65	65	65	65	65
66	66	66	66	66	66
67	67	67	67	67	67
68	68	68	68	68	68