

Investigation of Protein Sequence Variants in a Natural Microbial Community by Top-Down Proteomics

Brian Erickson¹; Mark Lefsrud¹; Nathan VerBerkmoes¹; Steven Singer²; Michael Thelen²; Jillian Banfield³; Robert Hettich¹

¹Oak Ridge National Laboratory, Oak Ridge, TN; ²Lawrence Livermore Natl. Lab, Livermore, CA; ³University of California, Berkeley, CA

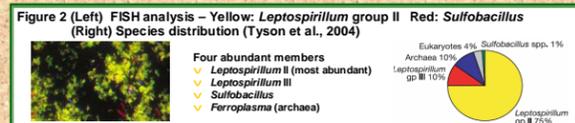


OVERVIEW

- The complexity of a natural microbial community presents unique and technical challenges for proteomics based discoveries.
- "Bottom-up" and "top-down" proteomics techniques can be integrated to identify the protein members and their isoforms within the extracellular portion of a natural microbial community.
 - These protein isoforms can include sequence variants and various post-translational modifications (PTMs).
- Within the community, a series of cytochrome related proteins were identified.
 - FTICRMS measurement along with IRMPD fragmentation revealed variant versions of the genome-predicted cytochrome-579.
 - Characteristic sequence tags from each cytochrome were used to identify the predicted genome cytochrome versus the variant cytochrome
- "Top-down" LC-MS identified 333 non-redundant proteins.
 - The identifications included 262 proteins with predicted methionine truncations and 36 with signal peptide cleavages.

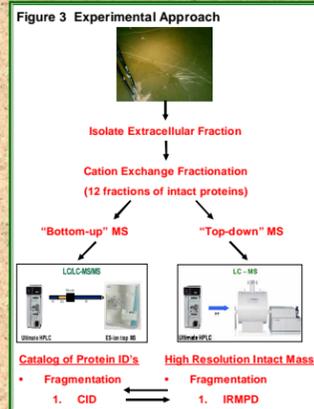
INTRODUCTION

- Proteome characterization of microbial species from natural environmental communities is challenged by the fact that these consortia are not clonal but rather encompass clades (i.e. bacterial cousins).
 - This leads to a substantial amount of strain variation in related proteins, precluding the easy identification of peptides by database searching algorithms in many cases.
 - Combining the advantages of high-throughput "bottom-up" MS with the high resolution intact mass measurements of "top-down" analyses will provide a comprehensive picture of the proteins present.
 - Furthermore, the identification of post-translation modifications can provide significant insight into protein function.
 - The identification of signal peptides can provide support for the localization and functional role of the proteins.
- Acid Mine Drainage (AMD)**
- Iron Mountain Mine, CA (near Redding, CA)
 - Biofilm exists in pH < 1, -42°C, Molar/sub-molar concentrations of Fe, Zn, Cu, As
 - Genome sequence (Tyson et al., 2004, *Nature*)
 - Proteome Characterization (Ram et al., 2005, *Science*)



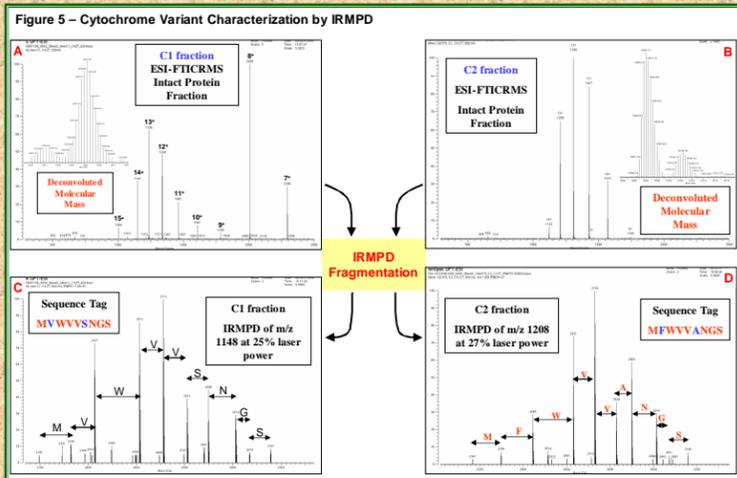
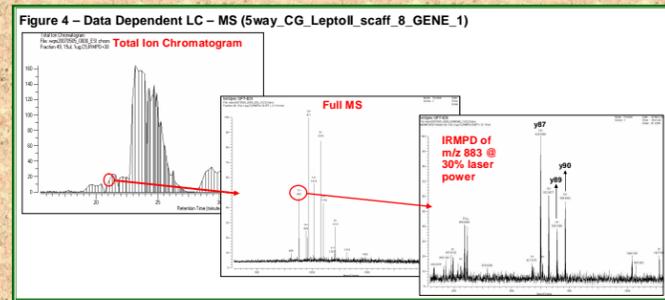
EXPERIMENTAL

- Sample Preparation**
 - The extracellular portion of whole biofilms were isolated by centrifugation (without cell lysis).
 - The resulting sample was then further separated by cation exchange, resulting in 12 fractions.
 - Each fraction was divided in half by volume for analysis by "bottom-up" and "top-down" MS.
- "Bottom-up"**
 - Samples were denatured, reduced, and proteolytically digested into peptides with trypsin.
 - Peptides were measured via a 1D (RP)-LC-MS/MS pipeline utilizing a linear ion trap (Thermo Finnigan LTQ).
- "Top-down"**
 - Intact proteins were analyzed by a 1D-LC-MS pipeline with all measurements on a 9.4T IonSpec FTICRMS.
- Protein ID**
 - "Bottom-up" spectra were searched with DBDigger and filtered with DTASelect. (Tabb et al., *Analytical Chemistry*, 2005)
 - "Top-down" spectra were identified with PTMSearchPlus.



RESULTS

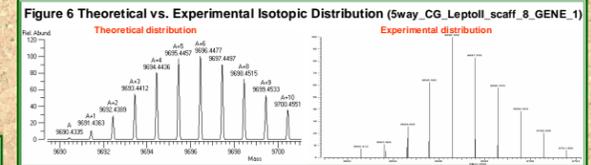
- Given the unique and hostile environment of the AMD biofilm, it is expected that a diverse and functionally important complement of proteins will be present in the extracellular portion. Identifying not only the proteins present, but also the degree of protein modification can provide valuable insight into their individual specific functions as well as how the microbial community exists as a whole.
- LC-FTMS**
 - On-line HPLC allowed for spatial protein separation of the complex biofilm samples. Parent ion selection was achieved in a data dependent manner utilizing dynamic exclusion as well as ion intensity for selection. (Figure 4)
- Cytochrome identification**
 - Two separate off-line HPLC fractions (C1 and C2) consisting of cytochrome 579 were isolated and examined by ESI-FTICRMS for interrogation of the intact proteins.
 - The molecular masses for these two variants of cytochrome 579 were distinct (Figure 5A/B)
 - Dissociation of these intact proteins was achieved with infrared multiphoton dissociation (IRMPD), and generated fragment ions that could infer sequence tag information (Figure 5C/D).
 - Through DNA sequencing it has been determined that the genome predicted cytochrome and the newly identified sequence variant share 88% sequence identity.



Protein Identification

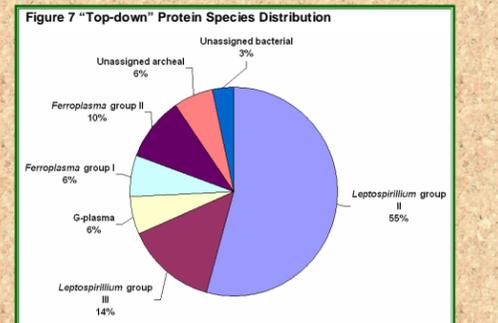
- "Top-down" analysis identified 333 proteins while "bottom-up" analysis identified 215 proteins. Within the identifications, a large number of hypothetical proteins were found. The top-down measurements revealed 262 proteins that contained an N-terminal methionine truncation, and 36 proteins that contained a signal peptide cleavage. (Table 1)
- Protein identifications were verified at high resolution by comparing theoretical vs. experimental isotopic distributions. (Figure 6)
- "Top-down" identifications include proteins from *Leptospirillum* III, *Ferroplasma*, *G-plasma*, as well as unassigned archaeal and bacterial species. The high mass resolution and accuracy (low ppm) increases the confidence of the protein assignments. (Table 2)
- Spectra were searched with an N-terminal methionine truncation, resulting in hundreds of putative identifications. Examination of the second amino acid residue provided biological support for these identifications. (Table 3)
- The mass spectra were also searched with an appended database containing protein sequences that include signal peptide cleavage. (Table 4)
- The composite protein identifications were analyzed for species of origin and resulted in a distribution correlating well with previous studies of the biofilm. As before, *Leptospirillum* group II was the most dominant. (Figure 7)

"Top-down" Identifications	333
"Bottom-up" Identifications	215
Methionine Truncation	262
Signal Peptide Cleavage	36



Protein Name	Protein Description	Calculated Mass (Da)	Observed Mass (Da)	PPM
URA_Leptoll_scaff8_8135_GENE_48	Uncharacterized conserved protein	10084.01	10084.01	0.0
Leu2_scaff_25_GENE_2	Hypothetical protein	12146.017	12146.017	0.0
URA_Leptoll_scaff8_8241_GENE_855	Hypothetical protein	13116.951	13116.951	0.0
Leu2_scaff_15_GENE_30	Cytoplasmic ATPase synthase	15409.448	15409.448	0.0
URA_Leptoll_scaff_27	S-Glycine cleavage system H protein (glyoxylase)	14942.508	14942.527	0.1
gpl_scaff_305_GENE_2	Dihydrogenase (cobalt-thiamine aldehyde dehydrogenase)	26433.295	26433.289	0.1
Sway_CG_Leptoll_scaff_8_GENE_1	Hypothetical protein	9996.871	9996.873	0.2
URA_Leptoll_scaff8_8241_GENE_189	Hypothetical protein	14841.528	14841.525	0.2
URA_Leptoll_scaff8_8241_GENE_279	Hypothetical protein	12052.827	12052.825	0.2
gpl_scaff_126_GENE_5	P-5-nucleoside (including N-terminal domain of PolI)	14841.528	14841.525	0.2
Sway_CG_Leptoll_scaff_74_GENE_21	Transposase	12844.489	12844.501	0.3
URA_Leptoll_scaff8_8002_GENE_84	Uncharacterized product	13174.130	13174.134	0.3
Sway_CG_Leptoll_scaff_302_GENE_2	Hypothetical protein	12844.489	12844.501	0.3
URA_Leptoll_scaff8_8135_GENE_10	Hypothetical protein	8973.303	8973.307	0.5
Leu2_scaff_248_GENE_1	Uncharacterized product	21187.384	21187.371	0.6
Leu2_scaff_217_GENE_2	Hydroxylase	11849.268	11849.264	0.3
Leu2_scaff_118_GENE_2	Hypothetical protein	12053.139	12053.139	0.0
URA_Leptoll_scaff_25	Hypothetical protein	14841.414	14841.427	0.9
gpl_scaff_289_GENE_2	RNA-binding protein Rplp	24928.712	24928.736	1.0
URA_Leptoll_scaff_19	Transposase	28423.808	28423.809	0.0
URA_Leptoll_scaff8_8002_GENE_7	Hypothetical protein	20099.213	20099.209	1.3
URA_Leptoll_scaff8_8135_GENE_1	Hypothetical protein	19172.482	19172.488	4.4
URA_Leptoll_scaff8_8135_GENE_107	Proteinase	19144.862	19144.869	1.4
Sway_CG_Leptoll_scaff_59_GENE_15	Cytoplasmic ATPase synthase	17242.411	17242.405	1.7
Sway_CG_Leptoll_scaff_59_GENE_15	Cytoplasmic ATPase synthase	18062.600	18062.628	1.6
Leu2_scaff_276_GENE_5	Ribosomal protein	17340.985	17340.985	1.7
URA_Leptoll_scaff8_8135_GENE_9	Transposase	16454.881	16454.909	1.7
Sway_CG_Leptoll_scaff_76_GENE_9	Transposase	6913.264	6913.408	2.1
URA_Leptoll_scaff8_8002_GENE_180	Mg-dependent chitinase	28793.881	28793.885	2.2
Leu2_scaff_485_GENE_1	Transposase	18482.059	18482.098	2.3
URA_Leptoll_scaff8_8048_GENE_295	Transposase	16876.303	16876.237	3.9
URA_Leptoll_scaff8_8241_GENE_177	Transposase	18849.772	18849.848	4.0
URA_Leptoll_scaff_25	Transposase	14315.071	14315.149	4.9
URA_Leptoll_scaff8_8002_GENE_48	Proteinase	16888.035	16888.234	11.9
gpl_scaff_406_GENE_2	Proteinase	26617.464	26617.201	15.3
Sway_CG_Leptoll_scaff_3_GENE_8	Multiple antibiotic resistance	21400.764	21401.148	18.4

Protein Name	Protein Description	2nd AA after Met	Calculated Mass (Da)	Observed Mass (Da)	PPM
URA_Leptoll_scaff8_8241_GENE_854	Hypothetical protein	E	1426.8326	1426.8326	0.0
Leu2_scaff_899	Conserved hypothetical protein	E	7146.2639	7146.262764	0.0
URA_Leptoll_scaff_10	Hypothetical protein	T	9379.8495	9379.849481	0.0
Sway_CG_Leptoll_scaff_88_GENE_19	Hypothetical protein	S	20311.1217	20311.12139	0.0
gpl_scaff_196_GENE_16	Hypothetical protein	N	13714.1218	13714.12151	0.0
Leu2_scaff_118_GENE_2	Inorganic pyrophosphatase	K	19045.8488	19045.84796	0.1
Sway_CG_Leptoll_scaff_180_GENE_15a	Ribosomal protein S4	S	5169.8811	5169.88095	0.1
URA_Leptoll_scaff8_8241_GENE_488	Hypothetical protein	N	5100.0811	5100.08087	0.1
URA_Leptoll_scaff8_8048_GENE_207	Uncharacterized product	G	15410.3741	15410.37197	0.1
Sway_CG_Leptoll_scaff_27_GENE_15	Hypothetical protein	A	12862.0015	12862.00084	0.2
URA_Leptoll_scaff8_8002_GENE_162	Leucylaminoacylase	K	18910.482	18910.48586	0.2
Leu2_scaff_89_GENE_13	Hypothetical protein	A	5409.7953	5409.79468	0.2
URA_Leptoll_scaff8_8048_GENE_252a	Hypothetical protein	G	8695.7711	8695.77381	0.2
URA_Leptoll_scaff8_8241_GENE_230	Ribosomal protein S27	K	50075.3453	50075.34461	0.2
gpl_scaff_131_GENE_15	Hypothetical protein	D	13028.9278	13028.92287	0.2
URA_Leptoll_scaff8_8135_GENE_11	Hypothetical protein	Y	4489.8804	4489.87725	0.2
URA_Leptoll_scaff8_8892_GENE_102	Ribosomal protein S17	K	10583.4391	10583.44312	0.2
Leu2_scaff_119_GENE_15	Hypothetical protein	K	18910.482	18910.48586	0.2
URA_Leptoll_scaff8_8241_GENE_197	Uncharacterized product	A	18871.817	18871.81065	0.2
Sway_CG_Leptoll_scaff_824_GENE_3	Cold shock protein	A	22140.9453	22140.94454	0.2
URA_Leptoll_scaff_2	Hypothetical protein	K	8612.8411	8612.83536	0.2
Leu2_scaff_48_GENE_27	Superfamily II DNA and RNA helicases	D	15316.1597	15316.14987	0.7
URA_Leptoll_scaff_4	Hypothetical protein	A	8612.8411	8612.83536	0.7
Leu2_scaff_48_GENE_27	Hypothetical protein	I	6207.1567	6207.15151	0.8
URA_Leptoll_scaff_37	Hypothetical protein	A	5885.655	5885.64843	0.8
URA_Leptoll_scaff8_8135_GENE_2	Hypothetical protein	S	15206.8824	15206.88168	0.8
URA_Leptoll_scaff8_8135_GENE_128	Uncharacterized product	S	15206.8824	15206.88174	0.8
URA_Leptoll_scaff_34	Hypothetical protein	Y	10822.921	10822.9162	1.0
Leu2_scaff_303_34	Hypothetical protein	P	13027.7655	13027.76622	1.0
Leu2_scaff_179a	Hypothetical protein	A	1642.7816	1642.7816	1.0
URA_Leptoll_scaff_25a_GENE_6	Ribosomal protein L1	Y	10822.921	10822.9162	1.0
Leu2_scaff_43_GENE_9	Hypothetical protein	Y	14996.8348	14996.81408	1.5
URA_Leptoll_scaff_1	Hypothetical protein	S	7656.18772	7656.18356	1.5
Sway_CG_Leptoll_scaff_27_GENE_50	Hypothetical protein	R	10877.1177	10877.13415	1.6
URA_Leptoll_scaff8_703_GENE_426	Hypothetical protein	R	10877.1177	10877.13415	1.6
Leu2_scaff_126_GENE_3	Hypothetical protein	D	8693.1206	8693.120504	1.6



DISCUSSION

- Methionine truncation
 - N-terminal methionine truncation is an essential modification thought to occur in ~55-70% of proteins, especially extracellular species. Expected residues following the methionine include A, C, G, P, S, T, and V. (Gigliore et al., 2004, *Cell. Mol. Life Sci.*)
 - Within this study, "top-down" MS identified 262 proteins with possible methionine truncation.
- Signal Peptide Identification
 - Cellular location is an important functional diagnostic tool. The presence of a cleaved signal peptide supports the extracellular location of the protein and can provide insight into the functional role of the protein.
 - Utilizing an "in-house" developed script, all predicted proteins in the biofilm database were searched for possible signal peptides using SignalP (Bendtsen et al., 2004, *J. Mol. Biol.*)
 - Proteins predicted to contain a putative signal peptide were appended to the database with the signal peptide sequence removed along with "SigP" appended to the gene name. "Top-down" spectra were searched with the appended signal peptide database using PTMSearchPlus (*in-house software*, V. Kertesz).
 - Initial results indicate 36 proteins are predicted to have signal peptide cleavage.
 - For all 36 predicted signal peptide cleaved proteins, the genome predicted full length protein was not identified, increasing confidence in the signal peptide prediction.

CONCLUSIONS

- Variations in amino acid sequence between proteins can be distinguished by high resolution FTICRMS along with IRMPD fragmentation.
- "Top-down" LC-MS identified 333 proteins. Additionally, 262 proteins were identified to contain methionine truncation.
- Combining *in-silico* prediction methods along with "top-down" LC-MS resulted in 36 proteins predicted to contain signal peptide cleavage.
- Possible post-translational modifications can be identified in a high-throughput manner with the use of LC-FTMS.
- The abundance of extracellular proteins appear to correlate well with the abundant microbial membership of the biofilm.

ACKNOWLEDGMENTS

- This research sponsored by U.S. Department of Energy, Genomics:Genomes-to-Life Program under contract DE-AC05-00OR-22725 with Oak Ridge National Laboratory, managed and operated by UT-Battelle, LLC.
- B. Erickson acknowledges financial support from the UTK-ORNL GST program.
- Vilmos Kertesz (ORNL) is acknowledged for development and implementation of the PTMSearchPlus software algorithm for integrating top-down and bottom-up MS data.