



# De Novo Library Construction for Metaproteomic Analysis of Human Datasets

Rajczewski, A.T.\*, Wagner, R.†, Mehta, S.\*, Jagtap, P.D.\*, Griffin, T.J.\*

\*Department of Biochemistry, Molecular Biology, and Biophysics, University of Minnesota †Minnesota Supercomputing Institute, University of Minnesota

## 1 Motivation

- **Metaproteomics** uses bottom-up mass spectrometry to gain taxonomic and functional information on microbial communities.
- As with all bottom-up proteomics, metaproteomics requires a FASTA file of anticipated species present
- When interrogating mass spectral data for metaproteomic insights, there are often no metagenomic data available to describe the environment

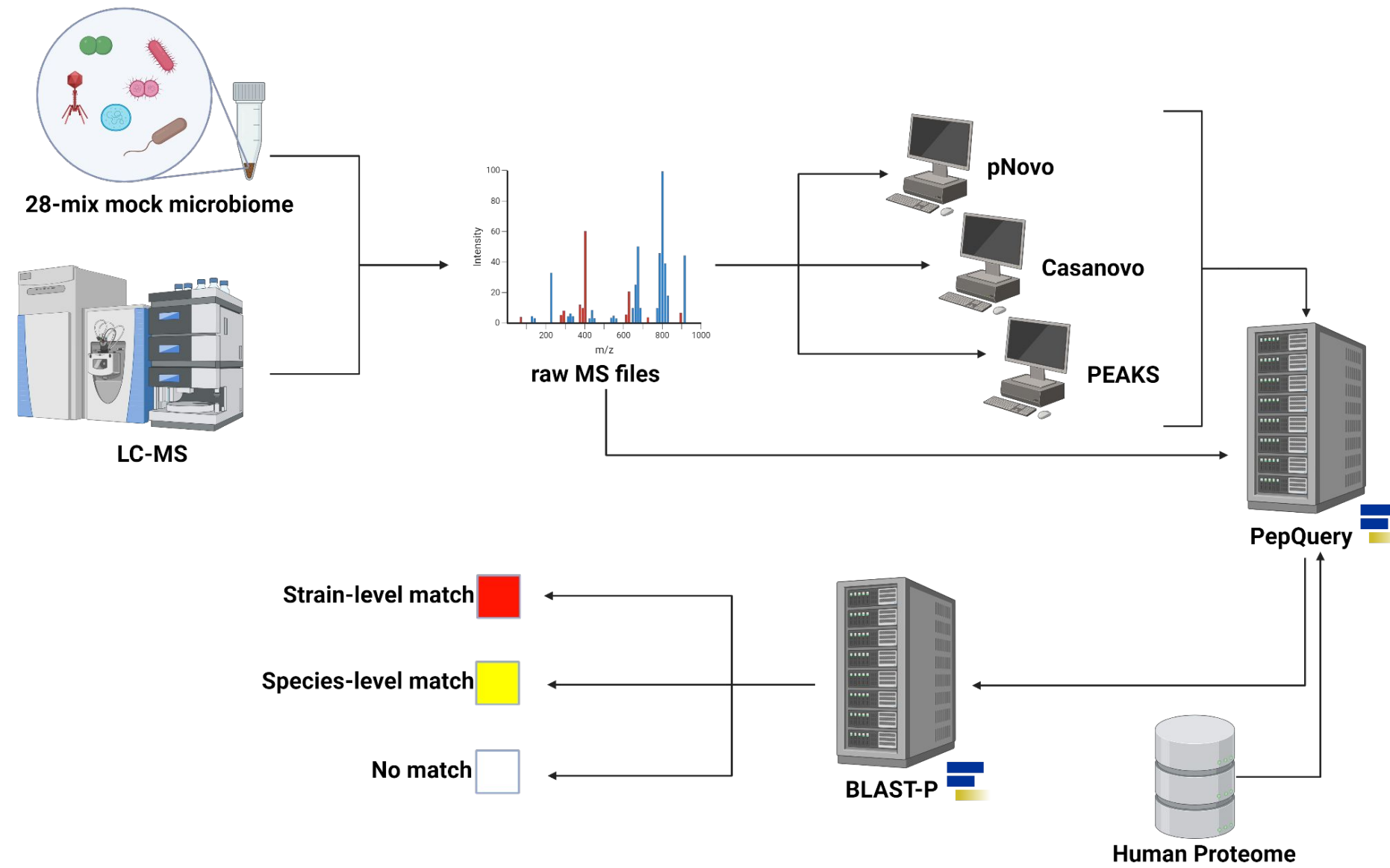
## 2 Hypothesis

- *De novo* sequencing of mass spectra to yield peptide sequencing has made great strides in recent years
  - Numerous open-source and for-purchase options exist
- ***De novo sequencing, when paired with technologies to remove contaminating sequences and identify contributing species, can be used to build databases for metaproteomic analyses***

## 3 Samples & Methods

- **Samples:** 1µg replicates of constructed microbial communities
- **LC:** Ultimate 3000 UHPLC in C18 nano mode with a 90-minute gradient
- **MS:** QExactive Orbitrap Quadrupole Hybrid Mass Spectrometer
- **Urinary proteome mass spectrometry data** from Yu *et al.* DOI: 10.7150/thno.16086
- **Casanovo, pNovo, PEAKS:** utilized for *de novo* sequencing
- **PepQuery, BLAST-P, and Uniprot** were utilized to refine and construct a database
- **MaxQuant** software suites utilized for quantitative proteomics analysis

## 4 De Novo Sequencing Software Selection

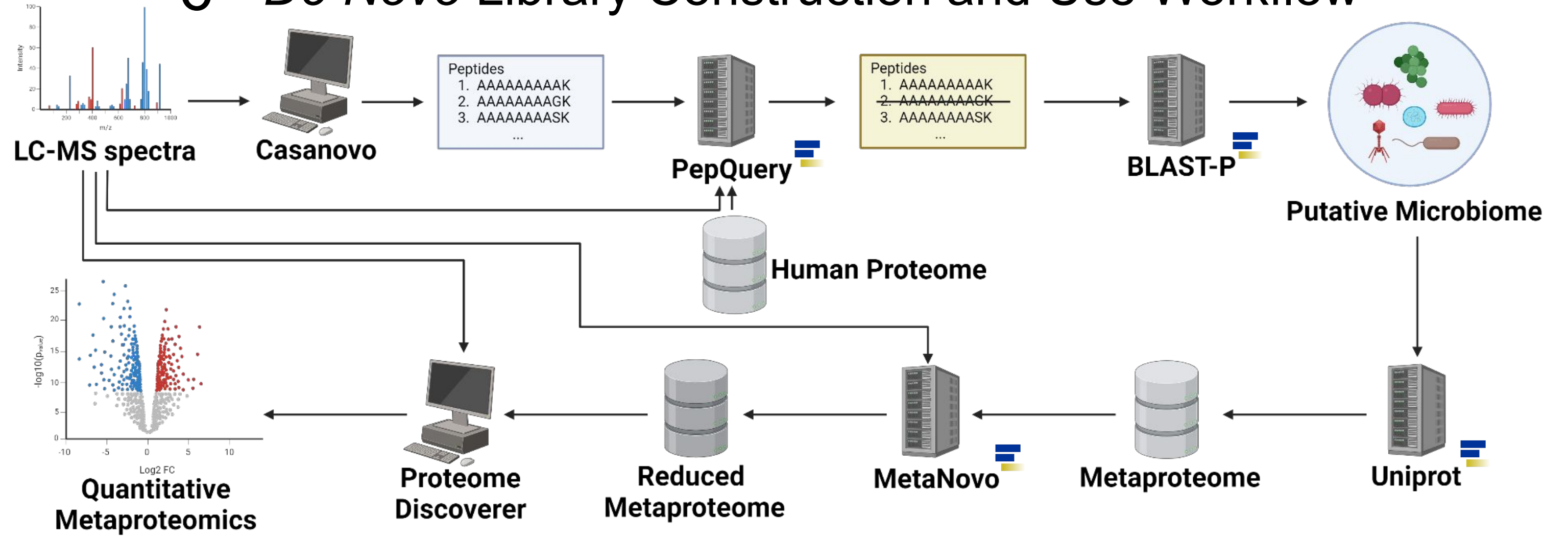


## 5 Casanovo for De Novo Sequencing

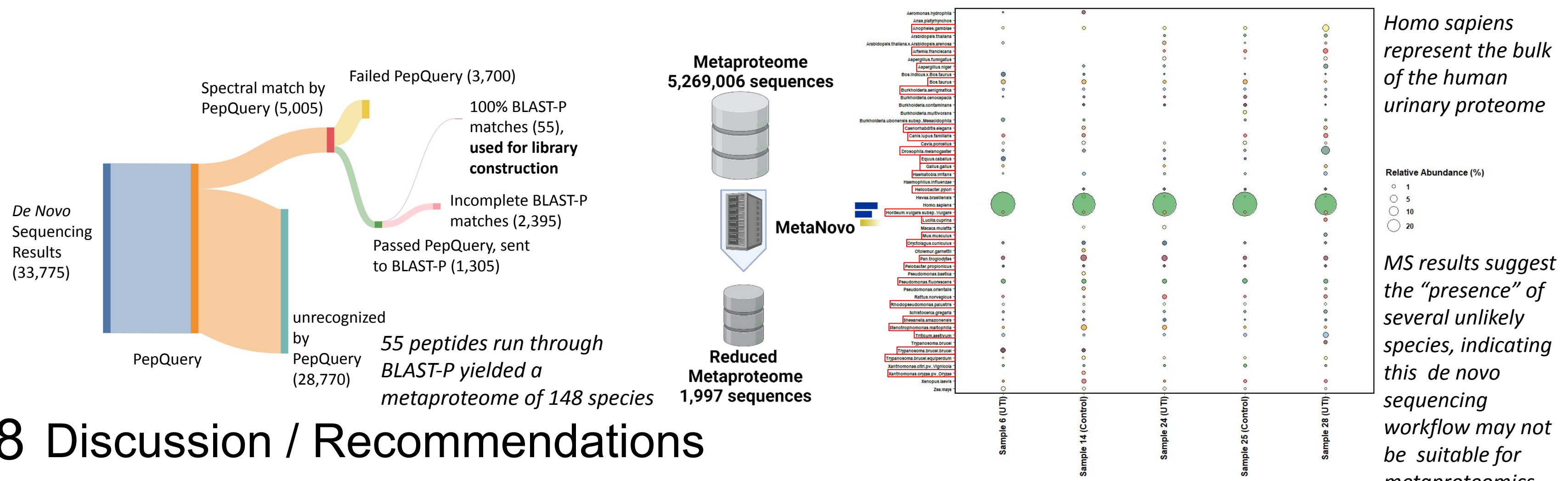
	PEAKS	Casanovo	pNovo	Protein abundance %
<i>Salmonella enterica typhimurium</i> LT2	Red	Red	Red	33.774
<i>Cupriavidus metallireducens</i> CH34	Red	Red	Red	15.519
<i>Stenotrophomonas maltophilia</i> SeITE02	Yellow	Yellow	Yellow	8.021
<i>Pseudomonas fluorescens</i> ATCC 13525, Type strain	Yellow	Yellow	Yellow	6.696
<i>Escherichia coli</i> K12 with Flac+ Plasmid	Yellow	Yellow	Yellow	5.788
<i>Agrobacterium tumefaciens</i> NTL4	Yellow	Yellow	Yellow	5.647
<i>Chlamydomonas reinhardtii</i>	Red	Red	Red	3.996
<i>Rhizobium leguminosarum</i> bv. <i>Viciae</i>	Yellow	Yellow	Yellow	3.172
<i>Pseudomonas denitrificans</i> ATCC 13867	Red	Red	Red	2.871
<i>Staphylococcus aureus</i> subsp. <i>aureus</i> <i>Rosenbach</i>	Yellow	Yellow	Yellow	2.606
<i>Thermus Thermophilus</i> HB27	Red	Red	Red	1.68
<i>Roseobacter</i> sp. AK199	Red	Red	Red	1.596
<i>Chromobacterium violaceum</i> CV026	Yellow	Yellow	Yellow	1.259
<i>Pseudomonas pseudoalcaligenes</i> KF707	Red	Red	Red	1.165
<i>Alteromonas macleodii</i> ATCC 27126	Red	Red	Red	0.954
<i>Desulfovibrio vulgaris</i> Hildenborough	Red	Red	Red	0.946
<i>Paracoccus denitrificans</i> ATCC 17741	Red	Red	Red	0.922
<i>Nitrososphaera viennensis</i>	Red	Red	Red	0.819
<i>Bacillus subtilis</i> 168	Red	Red	Red	0.788
<i>Nitrosomonas ureae</i> Nm10	Red	Red	Red	0.543
<i>Burkholderia xenovorans</i> LB400	Red	Red	Red	0.433
<i>Nitrospira multiformis</i> ATCC 25196	Red	Red	Red	0.209
Phage M13	Red	Red	Red	0.147
Phage P22 (HT105)	Red	Red	Red	0.106
Phage F0	Red	Red	Red	0.088
Phage ES18 (H1)	Red	Red	Red	0.088
Phage F2	Red	Red	Red	0.084
<i>Nitrosomonas europaea</i> ATCC 19718	Red	Red	Red	0.082

Strain-level (red) and species-level (yellow) matches of the three de novo sequencing softwares are presented here. Casanovo was selected for use with human datasets due to its high number of matches and potential for integration into Galaxy

## 6 De Novo Library Construction and Use Workflow



## 7 De Novo Metaproteomic Analysis



## 8 Discussion / Recommendations

- Casanovo was selected for *de novo* library construction
- Of the 33,775 peptides generated by *de novo* sequencing, 55 peptides showed 100% sequence alignment in BLAST-P
- *De novo* sequencing posits unlikely species in the urinary metaproteome
- Other *de novo* sequencing platforms will be tested in future studies
- Future experiments will automate the workflow in the Galaxy suite.

## 9 Acknowledgements and Funding

This research was supported in part by the National Institutes of Health grant P01 CA138338. Andrew Rajczewski was supported by an NIH biotechnology training grant T32GM008347 from the NIH National Institute of General Medical Sciences. Workflows generated via biorender.com