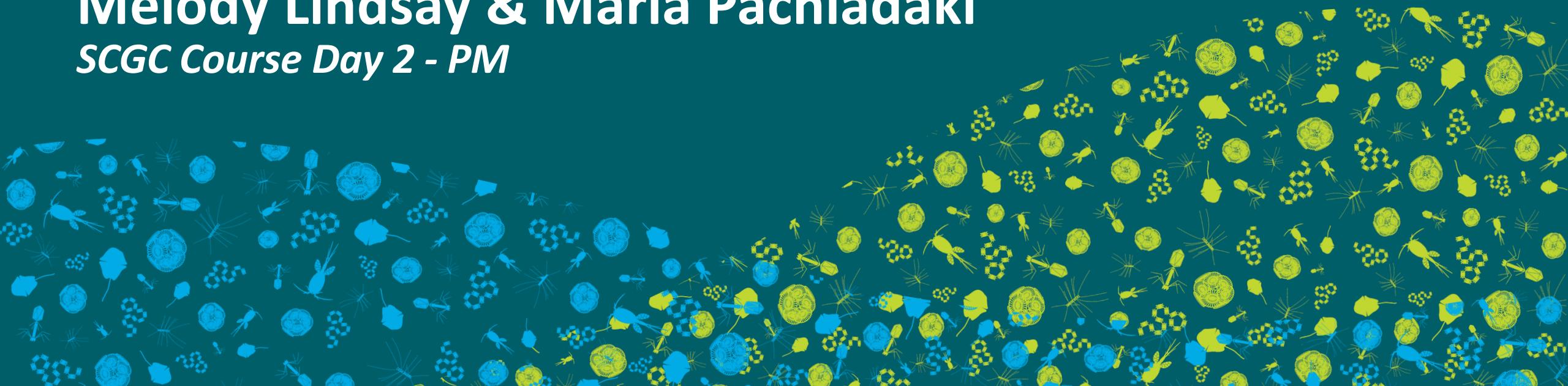


# DRAM annotation outputs *with further caveats!*

Melody Lindsay & Maria Pachiadaki

*SCGC Course Day 2 - PM*



# DRAM outputs

Now:

Annotation summaries

Let's go through the files you generated from the "annotate" and "distill" DRAM commands

What are they, how can they be used? What will be most useful to you?

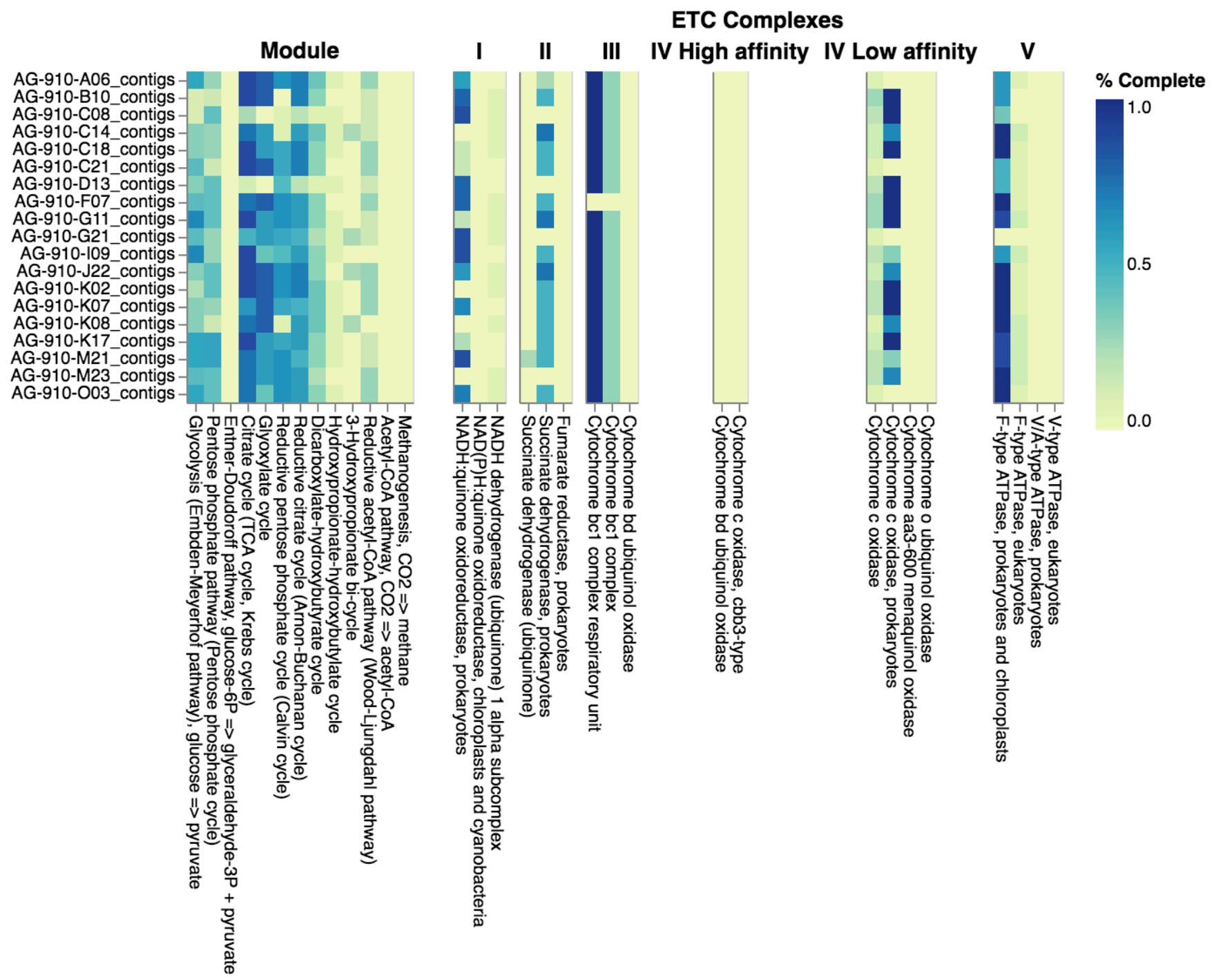
Then:

Some examples of misannotations, and ways to avoid

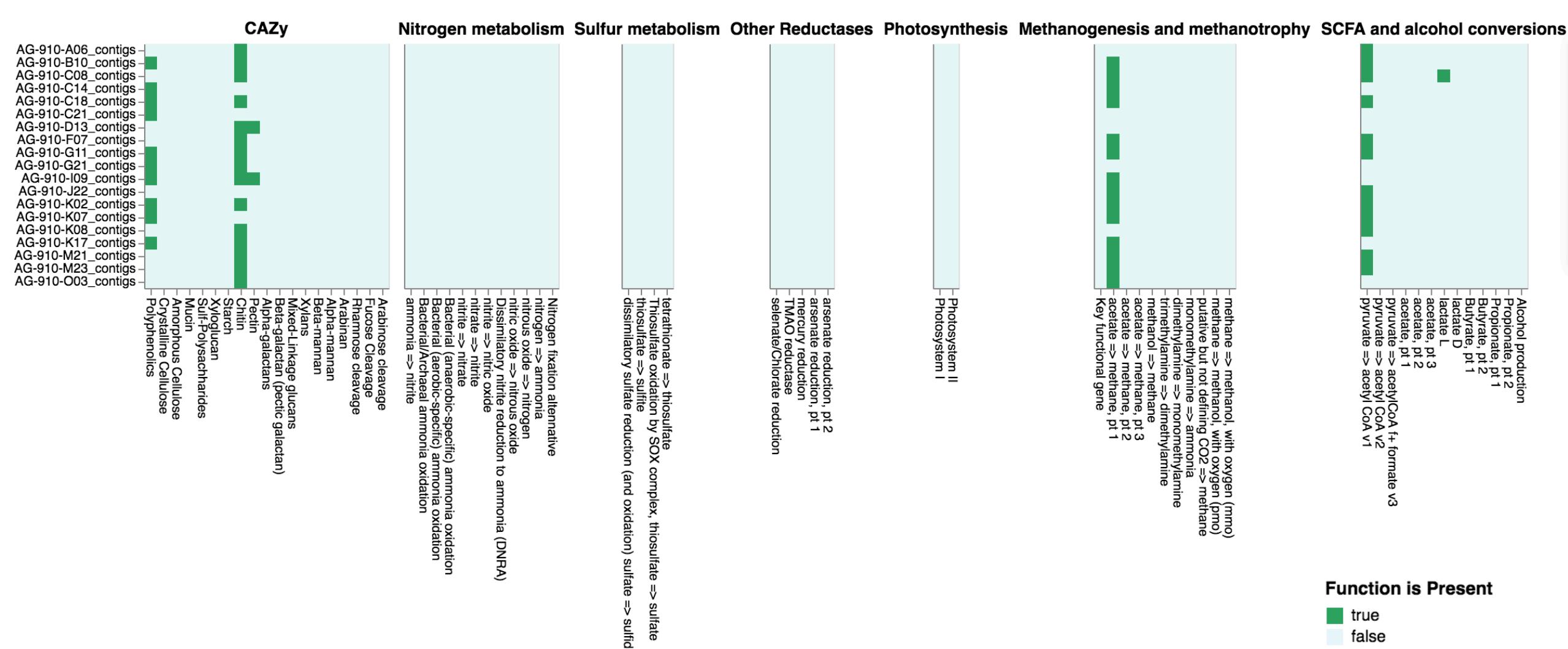
(not covering "missed" annotations, or proteins of unknown function)

Work on examples of DRAM output, manipulate and look through annotations





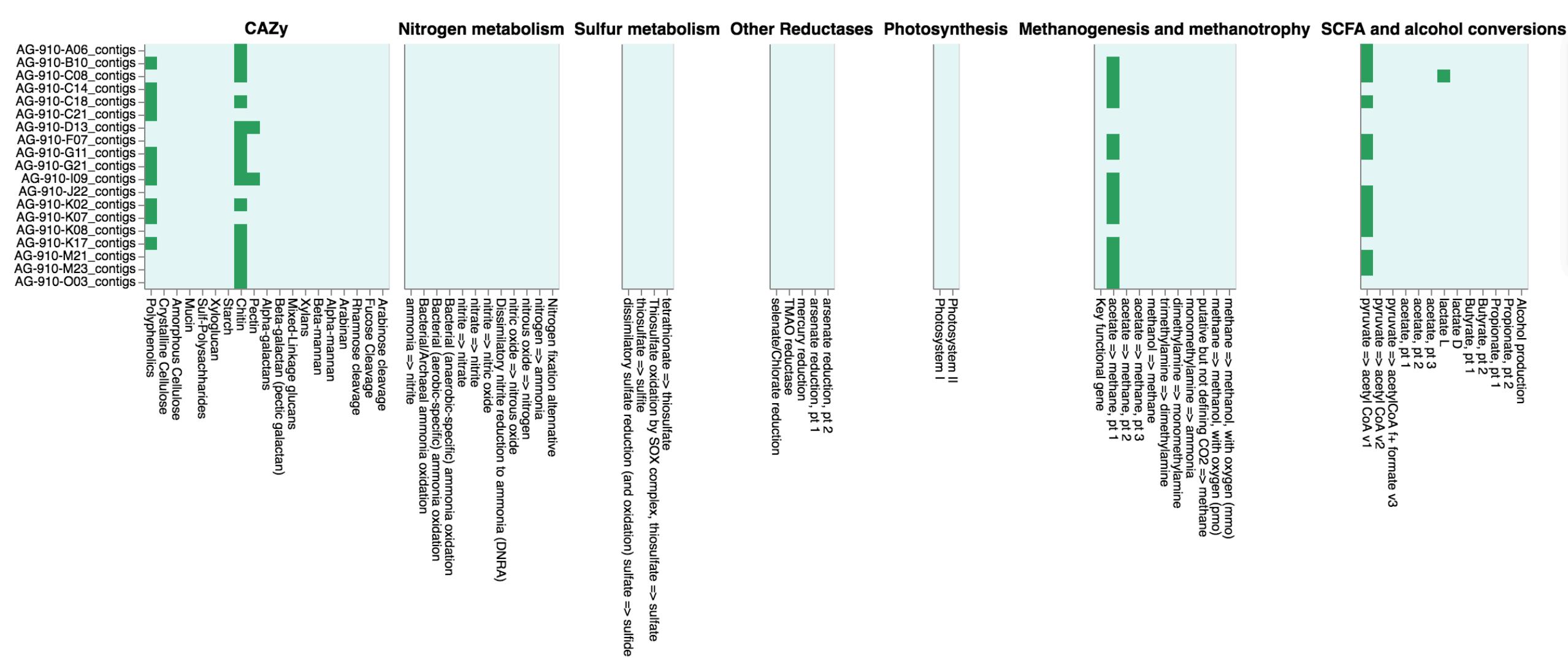
- One example of visualization (made with DRAM in the “Distill” command)
  - we won’t do a whole lot with manipulating this data, but will show you how to access what goes into a dataset like this!



# Misannotations, or genes involved in multiple pathways (ascribed to other functions)

- Rampant in databases, particularly as larger and larger sequence datasets are deposited using automated annotation programs
- How would misannotations become “the norm” for some genes?
- Do you know of any common misannotations in data from organisms you work with/might work with in the future?
- Common genes?
- Some manual curation may be needed in some datasets





Notice anything suspicious?

How would you go about verifying? (check which file??)

# Check databases (KEGG)

ORTHOLOGY: K01895	
<b>Entry</b>	K01895 KO
<b>Symbol</b>	ACSS1_2, acs
<b>Name</b>	acetyl-CoA synthetase [EC:6.2.1.1]
<b>Pathway</b>	<a href="#">map00010</a> Glycolysis / Gluconeogenesis <a href="#">map00620</a> Pyruvate metabolism <a href="#">map00630</a> Glyoxylate and dicarboxylate metabolism <a href="#">map00640</a> Propionate metabolism <a href="#">map00680</a> Methane metabolism <a href="#">map00720</a> Carbon fixation pathways in prokaryotes <a href="#">map01100</a> Metabolic pathways <a href="#">map01110</a> Biosynthesis of secondary metabolites <a href="#">map01120</a> Microbial metabolism in diverse environments <a href="#">map01200</a> Carbon metabolism
<b>Module</b>	M00357 Methanogenesis, acetate => methane

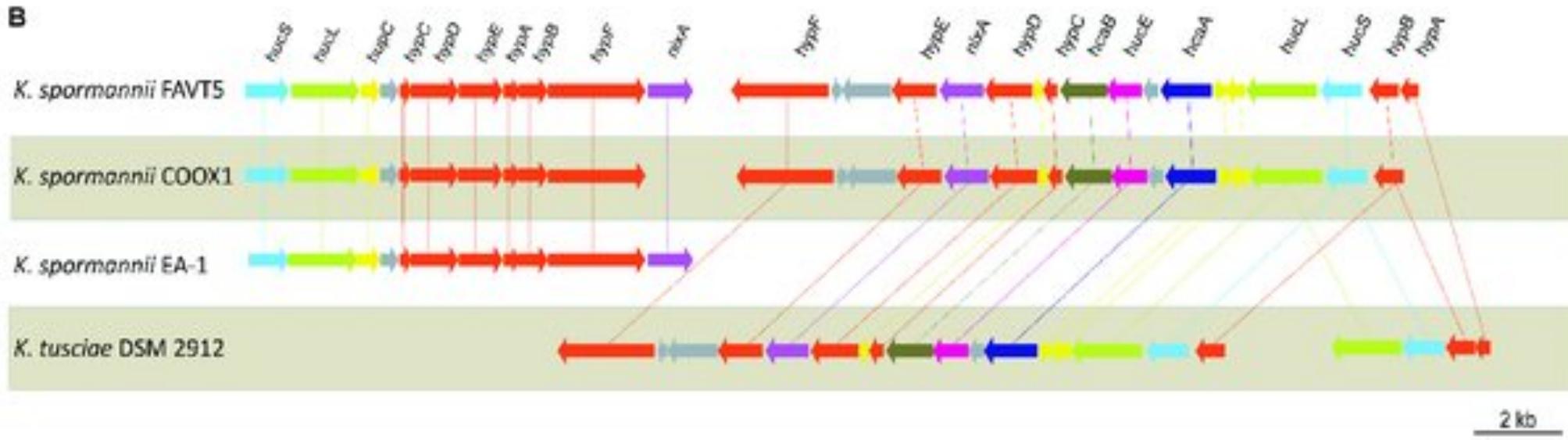
KEGG database

# Other ways to check for potential misannotations? Or divergence?

- Active sites for enzymes (do they possess key motifs)?
  - Example: cytochromes
  - Example: hydrogenases
- Specific databases of your choice
  - Examples: hydrogenases/HydDB, cytochrome P450 database, caZY,
- Do you have any commonly misannotated genes that you know of?  
Common in environments/samples/microbes that you study?

**A**

Consensus	11 motif	12 motif
<i>K. spormannii</i> FAVTS_hucl_3	X X R I C G I C G X X H S	F D O X C L V C T V
<i>K. spormannii</i> FAVTS_hucl_2	T P R I C G I C G I G H S	S Y D L C L V C T V
<i>K. spormannii</i> COOX1_hucl_3	T P R I C G I C G G S H S	S Y D S C L V C T V
<i>K. spormannii</i> COOX1_hucl_2	T P R I C G I C G G S H S	S Y D S C L V C T V
<i>K. spormannii</i> EA-1	T R R I C G I C G I G H S	S Y D L C L V C T V
<i>K. tusciae</i> BtuI_2044	T P R I C G I C G G S H S	S Y D D C L V C T V
<i>K. tusciae</i> BtuI_2524	T P R I C G I C G G S H S	S Y D S C L V C T V

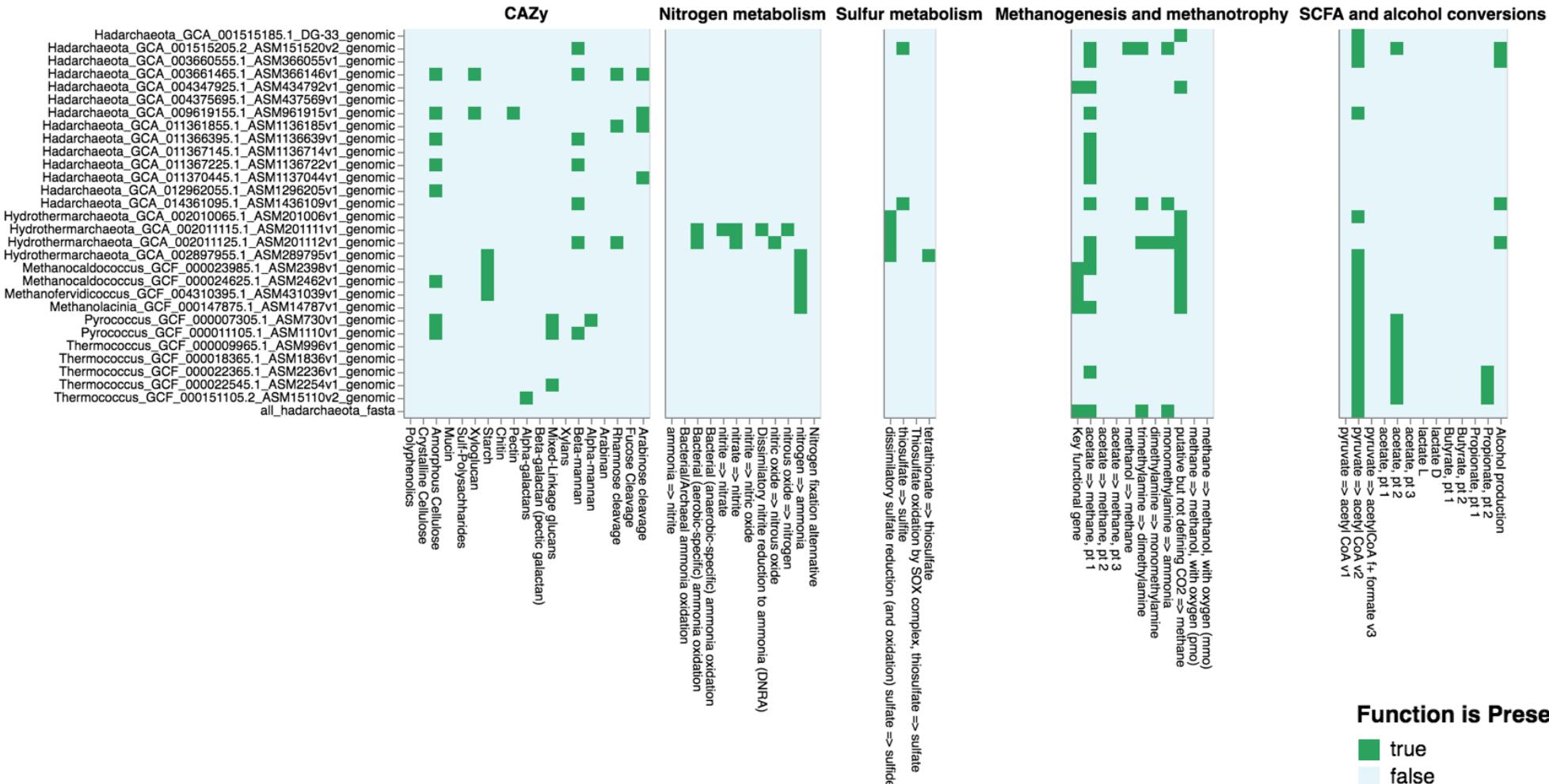
Hogendoorn  
et al., 2020**B**

- Example: Cysteine residues key for binding site of hydrogenases (for [NiFe]-hydrogenases, need paired CxxC motifs)

# Other Misannotations?

- Carbon fixation – TCA cycle, ATP citrate lyase involved in carbon fixation (appear to be TCA cycle, unless you compare with another ATP citrate lyase). **Challenge: how would you do this?**
- Arsenic reduction misannotation – one of the genes that is a regulation factor (actually in multiple pathways, need to verify that other genes are present in the arsenic reduction pathway!).  
**Challenge: how would you do this?**

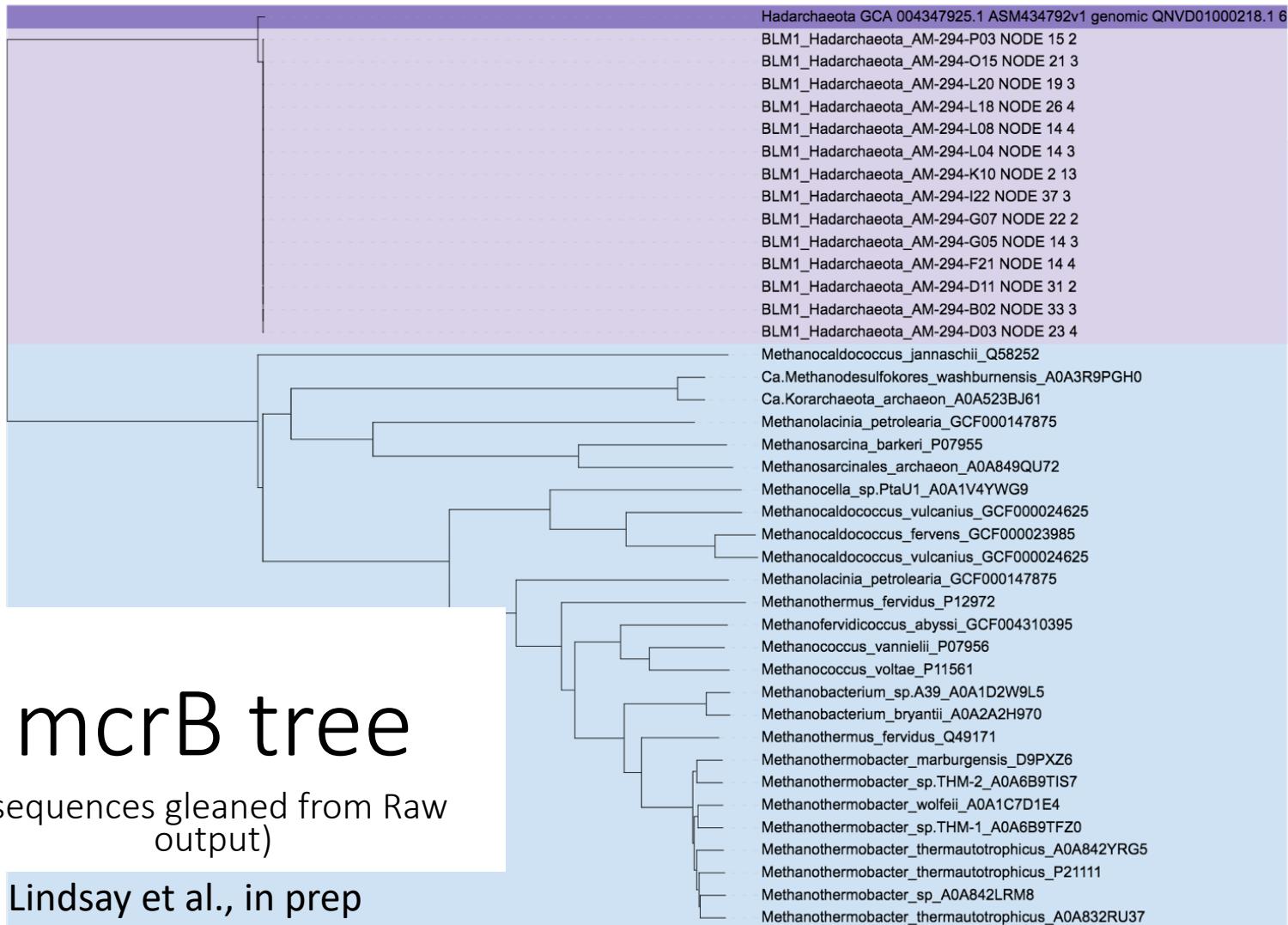
# Other ways to check for potential misannotations? Or divergence?



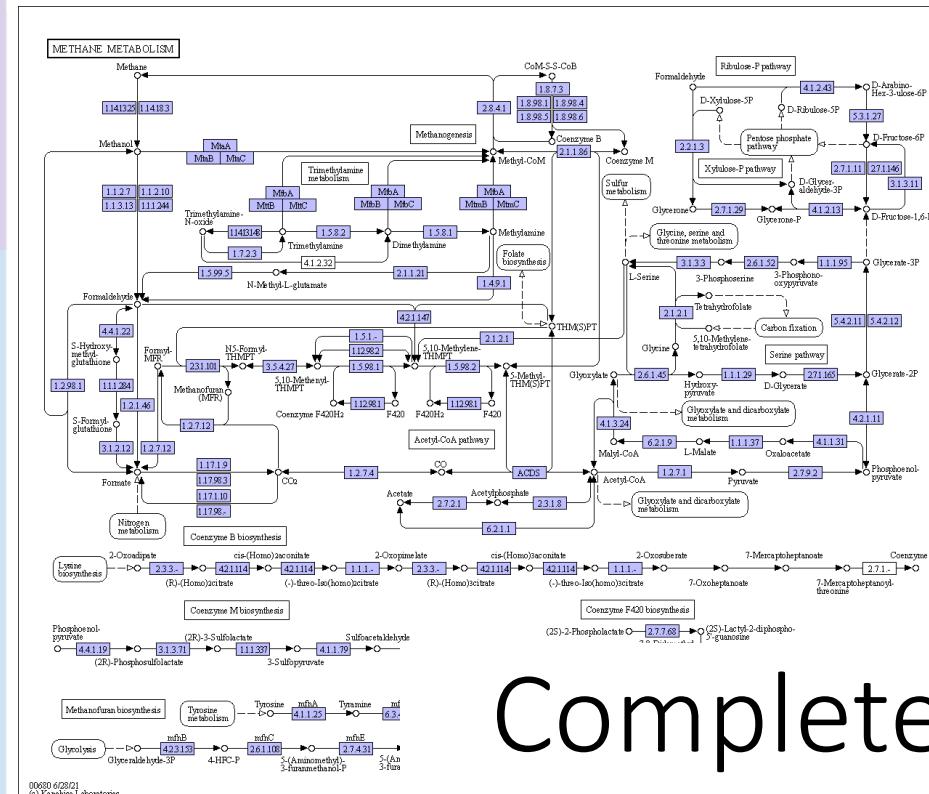
A	B	C	D	G	H	I	J	K	L	M	
1	gene_id	gene_description	module	header	294-B02_cd	294-C15_cd	294-C19_cd	294-C23_cd	294-D03_cd	294-D11_cd	294-D14_cd
174	K00320	coenzyme F420-dependent N5,N10-methenyltetrahydromethanopterin S-methyltransferase [EC:2.1.1.86] [RN:R04541]	Methanogenesis, CO <sub>2</sub> => methane	C1-methane	0	0	0	0	0	0	0
175	K00399	methyl-coenzyme M reductase [EC:2.8.4.1] [RN:R04541]	Methanogenesis, CO <sub>2</sub> => methane	C1-methane	0	0	0	0	0	0	0
176	K00401	methyl-coenzyme M reductase [EC:2.8.4.1] [RN:R04541]	Methanogenesis, CO <sub>2</sub> => methane	C1-methane	1	0	0	0	1	1	0
177	K00402	methyl-coenzyme M reductase [EC:2.8.4.1] [RN:R04541]	Methanogenesis, CO <sub>2</sub> => methane	C1-methane	0	0	0	0	0	0	0
178	K00577	tetrahydromethanopterin S-methyltransferase [EC:2.1.1.86] [RN:R04541]	Methanogenesis, CO <sub>2</sub> => methane	C1-methane	0	0	0	0	0	0	0

Distillate output

Tree scale: 1000



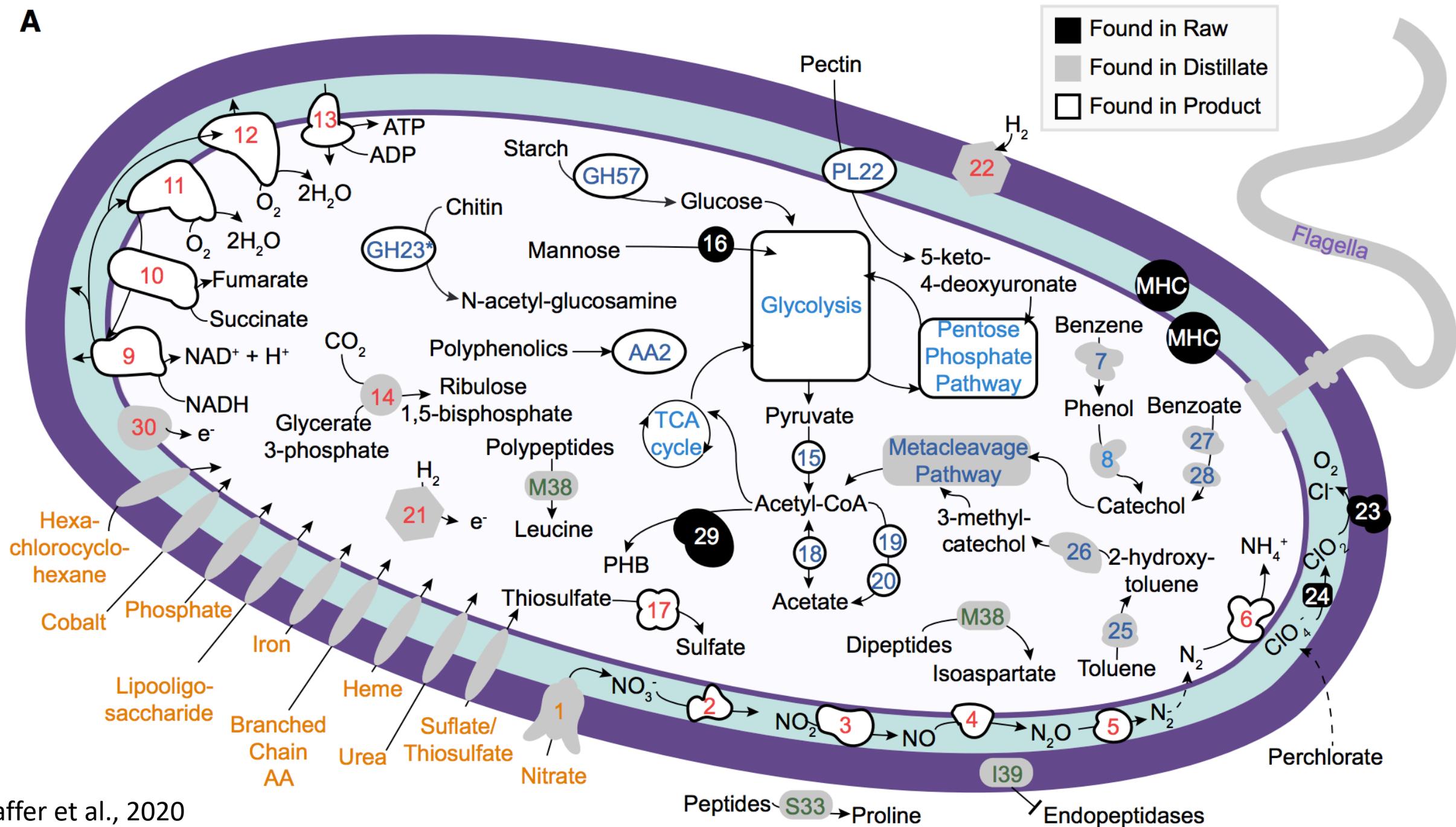
KAAS - KEGG Automatic Annotation Server  
for ortholog assignment and pathway mapping



Genome.jp

Complete pathways?

A



- What other metabolism summaries are you interested in? Try this command with your broad category of choice (as seen in product.html)
- As a reminder - you will need to parse out the excel file sheets again, in case you want to look at different broad categories (other than energy)!

Other exploratory ideas:

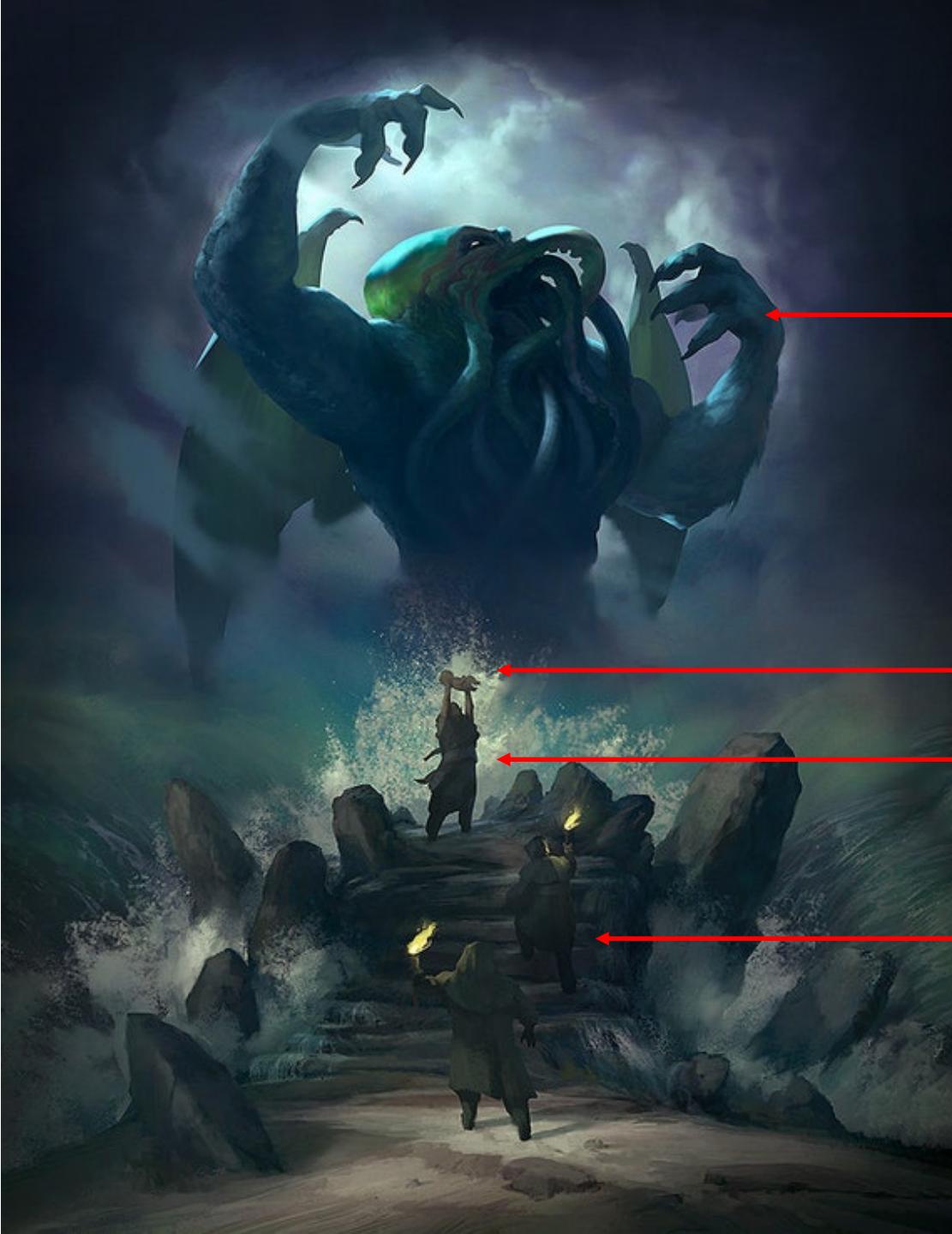
How complete are individual pathways in this dataset? (as annotated by the DRAM pipeline?)

- 1) how complete are cytochrome bc complexes from Pelagibacter SAGs? (where can you find this data?)
- 2) how complete are key carbon cycling pathways?

We also have all concatenated SAGs under /storage/lesson\_analyses/annotations/AG-910. Can you use the distill command (from DRAM github) to make summary annotations.tsv, product.html & .tsv, metabolism.xlsx, of your choice of taxa?

- 1) Flavobacterales?
- 2) Rhodobacterales?

Keep in mind that DRAM also works with MAGs and other data types.



Functional annotations/metabolic mapping of a confusing archaea?  
(which is not a methanogen!)

My sanity

Me

My postdoc advisors

<https://www.etsy.com/ie/listing/653073086/lovecraft-cthulhu-poster-cthulhu-art>