

Metapathways Workshop: Functional Space Reduction

Analyzing metagenomic data from
different oceanic provinces

Metapathways: Steps

1. Quality Control and ORF prediction
2. ORF annotation
 - A. Functional Annotation
 - B. Taxonomic Annotation
3. Analyses
4. ePGDB construction

Metapathways Input

➤ .fasta file

File Path ▾ : ~/MetaPathways_1_0/input/jgi_illumina/4096453.combine

4096453.combined_unique.fa

```
1 >Sequence0000000001
2 GGTGCTGAGGTAATTCATCTTGGTCATAATCGCTCAGTGCATGAAATAC
3 AAGGTGGTCACATGGAGTTTTTCAAATATATGTATGATATGCTTGAAG/
4 TCCAGAGGAAGTAAAAGAACTAATGAATTATGGTATTACCAGAATTTA
5 ATGGAGAAATCTGATTTTCCACAGGAGAAAGTCTAAAAATTGATTTA
6 CAGAAAATTATCCTGAAATAGCAGCACCCATGATAAAAGAGCTGATAA
7 TGCTGGAAAATCTTCACTTGTGGATGAAATTGTAAGACGCTTTCTTGA
8 AAATCTGGAGGTGCACTTCTTGGCGATCGTATAAGGATGAATGCAATT
9 GTATTTCAAATTATGTAAAGGATGCTGAAATAATTGCTCAGGCCGCAG
10 AATAATTGATCATAGTGATGTTTCATTCTATGTTATGACCCCTGAATA
```

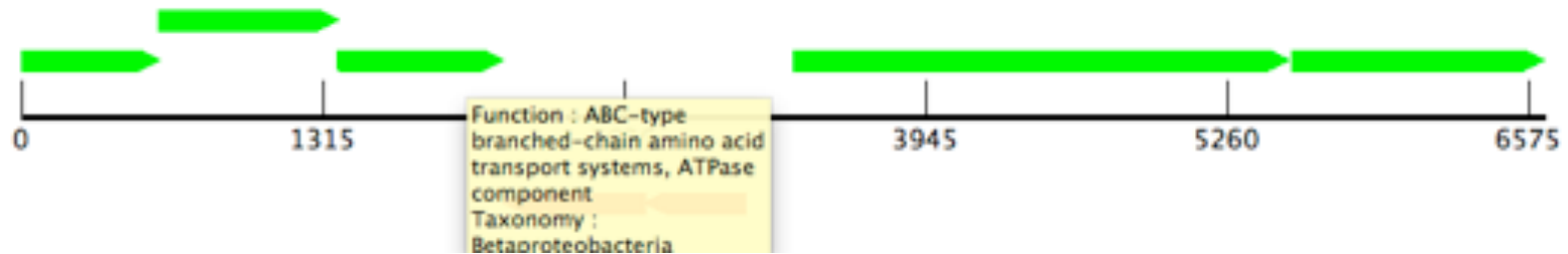
➤ Template Param File

File Path ▾ : ~/MetaPathways_1_0/input/jgi_illumina/Sep_2013/illumina_param.txt

illumina_param.txt

```
1 ##V.1 do not remove this line
2 # MetaPathways v1.0
3 # Kishori M. Konwar, Niels W. Hanson
4 # Parameter File
5
6 INPUT:format fasta
7 # e.g. fasta gbk-annotated gbk-unannotated gff-annotated gff-unannotated
8
9 # Quality Control parameters
10 quality_control:min_length 180
11 quality_control:delete_replicates yes
12
13 # ORF prediction parameters
14 orf_prediction:algorithm prodigal
15 orf_prediction:min_length 60
16
17 # ORF annotation parameters
18 annotation:algorithm last
19 # e.g. blast or last
20 annotation:dbs metacyc-v5-2011-10-21,kegg-pep-2011-06-18,cog-2007-10-30,refseq_protein_2012-11-13,MDM_SAG_proteins
21 # e.g. annotation:dbs cog,kegg,refseq,metacyc
22 annotation:min_bsr 0.4
23 annotation:max_evalue 0.000001
24 annotation:min_score 20
25 annotation:min_length 60
26 annotation:max_hits 5
27
28 # rRNA annotation parameters
29 # e.g. rRNA:refdbs
30 rRNA:refdbs GREENGENES_gg16S-2012-11-06,SSURef_111_NR_tax_silva-2012-11-06,LSURef_111_tax_silva_2012
31 rRNA:max_evalue 0.000001
32 rRNA:min_identity 20
33 rRNA:min_bitscore 50
34
35 # pathway tools parameters
36 ptools_settings:taxonomic_pruning no
37
38 # grid settings
39 grid_engine:batch_size 200
40 grid_engine:max_concurrent_batches 400
41 grid_engine:walltime 10:00:00
42 grid_engine:RAM 10gb
43 grid_engine:user myusername
44 grid_engine:server mygrid.domain.com
45
46 # pipeline execution flags
47 # e.g. yes, skip, redo
48 metapaths_steps:PREPROCESS_FASTA yes
49 metapaths_steps:ORF_PREDICTION yes
50 metapaths_steps:GFF_TO_AMINO yes
51 metapaths_steps:FILTERED_FASTA yes
52 metapaths_steps:COMPUTE_REFSCORE yes
53 metapaths_steps:BLAST_REFDB yes
54 metapaths_steps:PARSE_BLAST yes
55 metapaths_steps:SCAN_rRNA yes
56 metapaths_steps:STATS_rRNA yes
57 metapaths_steps:SCAN_tRNA yes
58 metapaths_steps:ANNOTATE yes
59 metapaths_steps:PATHOLOGIC_INPUT yes
60 metapaths_steps:GENBANK_FILE yes
61 metapaths_steps:CREATE_SEQUIN_FILE yes
62 metapaths_steps:CREATE_REPORT_FILES yes
63 metapaths_steps:MLTREEMAP_CALCULATION skip
64 metapaths_steps:MLTREEMAP_IMAGEMAKER skip
65 metapaths_steps:PATHOLOGIC redo
```

1. Quality Control and ORF prediction



- Finds regions in nucleotide sequence which code for an open reading frame (ORF)
- ORFs predicted using PRODIGAL
- Conversion from nucleotide seq to AA sequence
e.g. **ATG** → **MET**

Default lengths: 180 nucleotides / 60 AAs

*/Output: /pre-processed
/orf_prediction*

2. ORF annotation

A. Functional Annotation

B/BLAST to compare AA sequences in your query to proteins in user-defined reference databases

e.g. Databases used in NESAP Illumina Analysis:

KEGG, COG, RefSeq, MetaCyc, MDM-SAG-proteins

Also could use CAZy, EggNog, or any other protein database

2. ORF annotation

B. Taxonomic Annotation

Nucleotide sequences are queried against reference nucleotide databases (e.g. SILVA and Greengenes) to identify ribosomal genes in sample metagenomes

- Functional and taxonomic info are combined to generate input files for ePGDB creation

/Output: /genbank (.annotated.gff) → ePGDB creation
/blast_results (.blast.parsed.txt)
/results/rRNA (rRNA.stats.txt)
/results/annotation_tables (.fxn_and_taxa_table.txt)

3. Analyses

A. tRNA Scan to identify relevant tRNAs

/Output: /results/tRNA

B. Least Common Ancestor for taxonomic binning

/Output: /results/LCA

C. ML TreeMap

/Output: /results/mltreemap

D. ePGDB creation

Input file: .annotated.gff

/Output: /ptools/

Current Metapathways Outputs

① Individual Sample Level

① Auxiliary Sample Statistics

- Nucleotide length distribution, AA length distributions

② Functional and Taxonomic Table

③ Taxonomic distribution of sample in NCBI tree

④ Functional characterization at KEGG & COG levels

★ ⑤ Pathway Table

② Multi-Sample Comparison Level

★ ① Master pathway table

Sample Pathway Table

- Created after pipeline run using:

extract_pathway_table_from_pgdb.pl

/Output: /results/pgdb/XXX.basepathways.txt

File Path ▾ : ~/MetaPathways_1_0/output/fos_ends/Sep13_2013/a4_10/results/pgdb/a4_10.basepathways.txt



a4_10.basepathways.txt

	PWY_NAME	PWY_COMMON_NAME	NUM_REACTIONS	NUM_COVERED_REACTIONS	ORF_COUNT						
1	PWY0-1312	acetate formation from acetyl-CoA I 2	1	3	a4_10_46_0	a4_10_1224_0	a4_10_5932_0				
2	PWY-6961	L-ascorbate degradation II (bacterial, aerobic)	8	1	3	a4_10_1904_1	a4_10_11058_0	a4_10_7394_0			
3	PWY-5794	malonate degradation I (biotin-independent)	3	1	3	a4_10_5943_0	a4_10_11134_0	a4_10_2927_0			
4	PWY-5162	2-oxopentenoate degradation 3	2	6	a4_10_342_1	a4_10_11189_0	a4_10_8802_0	a4_10_9514_0	a4_10_802_0	a4_10_37	

ROWS

COLUMNS

pathways

- 
- 
- (1) Pathway short name
 - (2) Pathway long name
 - (3) # of reactions needed to complete a pathway (from MetaCyc)
 - (4) # of distinct reactions covered in a sample
 - (5) # of ORFs found in a pathway in a sample
 - (6) Names of ORFs found in each pathway in each sample

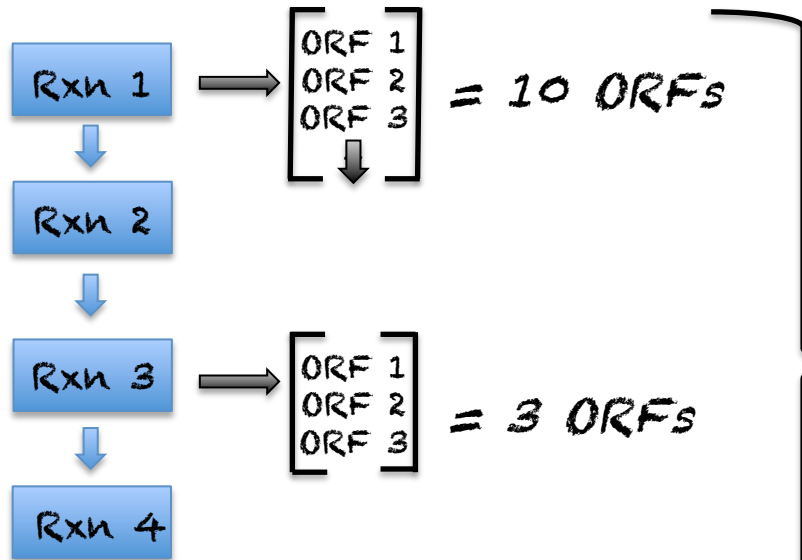
Example: Sample Pathway Table

ROWS
pathways

COLUMNS

- (1) Pathway short name
- (2) Pathway long name
- (3) # of reactions needed to complete a pathway (from MetaCyc)
- (4) # of distinct reactions covered in a sample
- (5) # of ORFs found in a pathway in a sample
- (6) Names of ORFs found in each pathway in each sample

Example: sample A, pathway X



For Sample A, pathway X

Column (3) = 4

Column (4) = 2

Column (5) = 13

Column (6) = ORF 1, ORF 2,

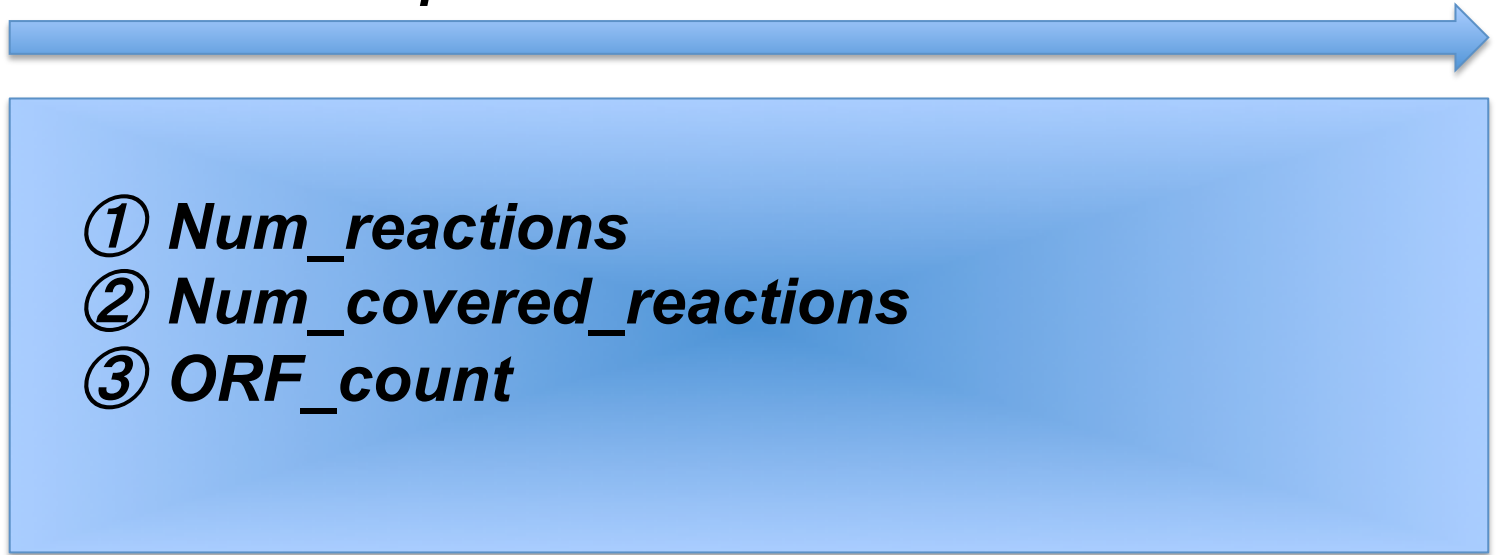
Total 13 ORFs covering 2 distinct reaction steps

Master Pathway Table

- Created after pipeline run using:
./make_table.pl

ROWS
pathways

COLUMNS - samples



Text files that you can import to program of choice for analysis (e.g. R, Matlab, Excel)

A pathways level analyses of the NESAP metagenomes

Run through Metapathways:

- ☑ 49 Saanich and Line P fosmid ends
- ☑ 91 Saanich and Line P Illumina samples

Template Param File Statistics:

LAST algorithm

Libraries used: *KEGG-pep-2011-06-18*

COG 2007-10-30

Metacyc v5 2012-10-21

RefSeq protein 2012-11-13

MDM-SAG proteins

Master Pathway Table

- Created after pipeline run using:
/make_table.pl

ROWS
pathways

COLUMNS - *samples*



① <i>Num_reactions</i>		
② <i>Num_covered_reactions</i>		
③ <i>ORF_count</i>		

Matlab Workflow I

1. Import master tables
2. Plot number of nucleotide sequences & ORFs before, after, and difference for each sample
 - *Calculate variability in ORFs after across the samples to decide effect of relativizing by library size*
3. Calculate variability in pathway length (number of reactions) for each pathway in dataset
 - *Done to decide if need to relativize by pathway length*

Note: #2 & 3 only matter if using quantitative data

Matlab Workflow II

To get a quick sense of the data:

4. Calculate variability in:

- (i) Total ORFs in the samples down the pathways
- (ii) Total ORFs in the pathways across the samples

5. Venn diagrams of shared and exclusive pathways

Matlab Workflow III

Statistical Analyses:

5. Make distance/dissimilarity matrix
6. NMS (aka NMDS)
7. Monte Carlo to check dimensions on NMS
8. Cluster Analysis
9. Define groups using either NMS or cluster analysis
10. Run ISA to determine pathways which are causing the groups to differentiate

Hallam Lab Timeseries: NESAP Metagenomes

- **49 Fosmid End Libraries:**

Saanich = 19 samples from 2004-2007

Line P = 30 samples from 2009-2010

Unassembled

➤ Can use quantitative ORF counts

- **91 Illumina Samples:**

Saanich = 48 samples from 2009-2011

Line P = 43 samples from 2008-2010

Assembled by JGI using Velvet

➤ Cannot use quantitative ORF counts at present

Hallam Lab Timeseries: NESAP Metagenomes

- **49 Fosmid End Libraries:**

Saanich = 19 samples from 2004-2007

Line P = 30 samples from 2009-2010

Unassembled

➤ Can use quantitative ORF counts

- **91 Illumina Samples:**

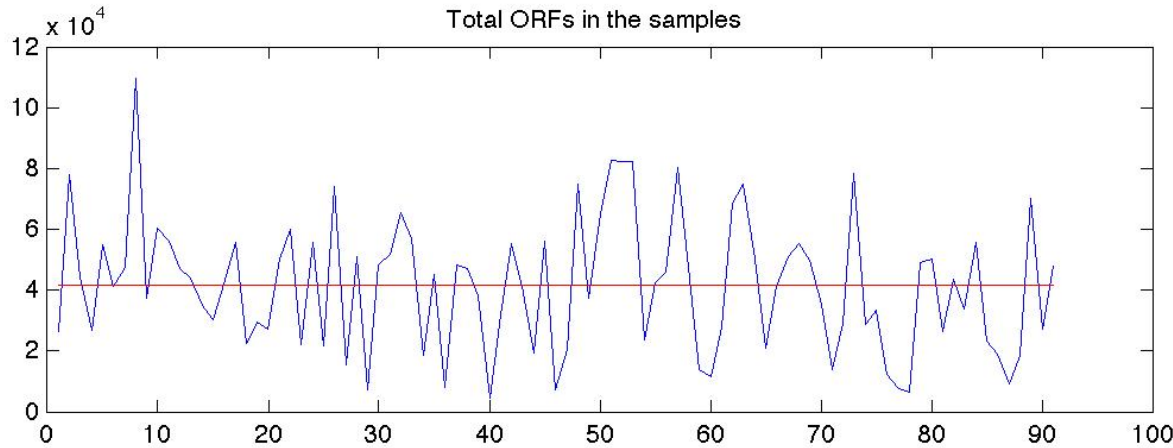
Saanich = 48 samples from 2009-2011

Line P = 43 samples from 2008-2010

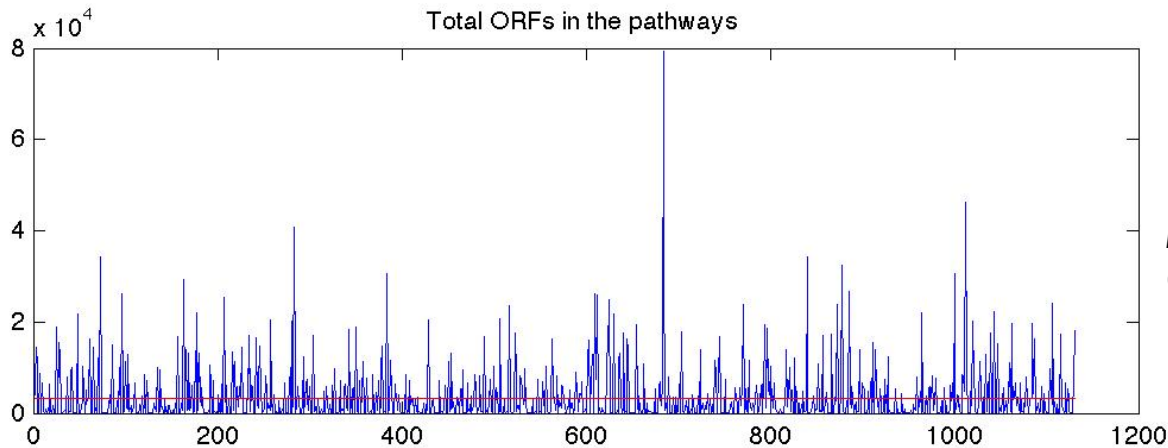
Assembled by JGI using Velvet

➤ Cannot use quantitative ORF counts at present

NESAP Illumina Dataset: Quick Look



Coefficient of Variation (CV)
 $100 * (\text{Std} / \text{Mean})$: 52%
Variability among the samples based on total ORFs is moderate



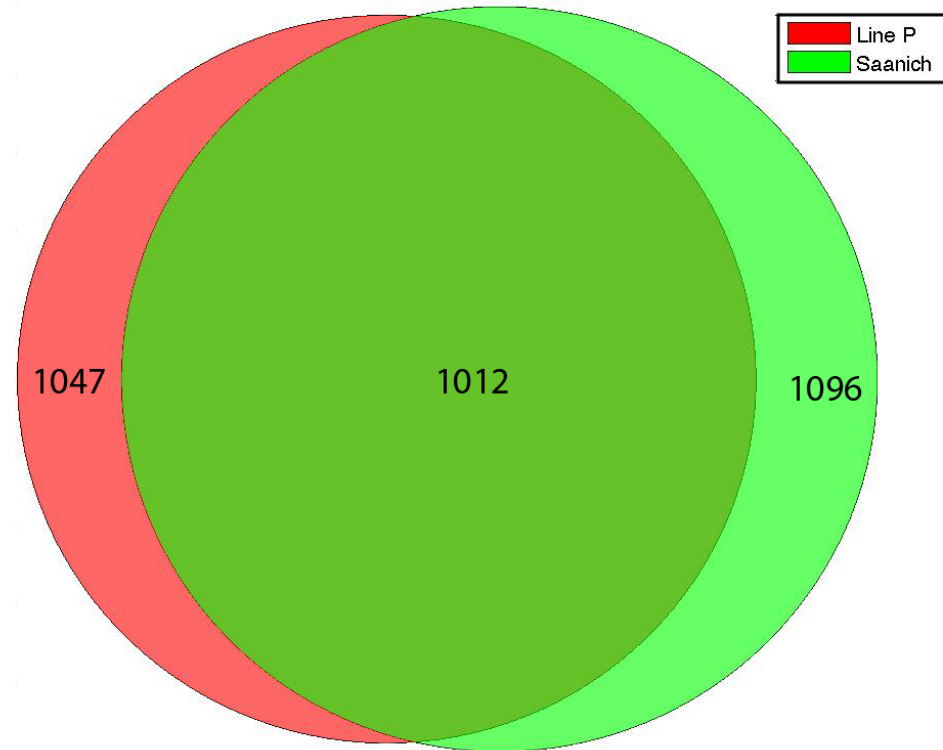
CV = 182%
Variability among the pathways based on total ORFs is large

- 1131 pathways resolved
- ~2000 pathways in Metacyc
- Maximum resolved using MetaPathways is 1239
- Illumina dataset is resolving 91.3% of max using MetaPathways

NESAP Illumina Dataset: Shared & Exclusive Pathways

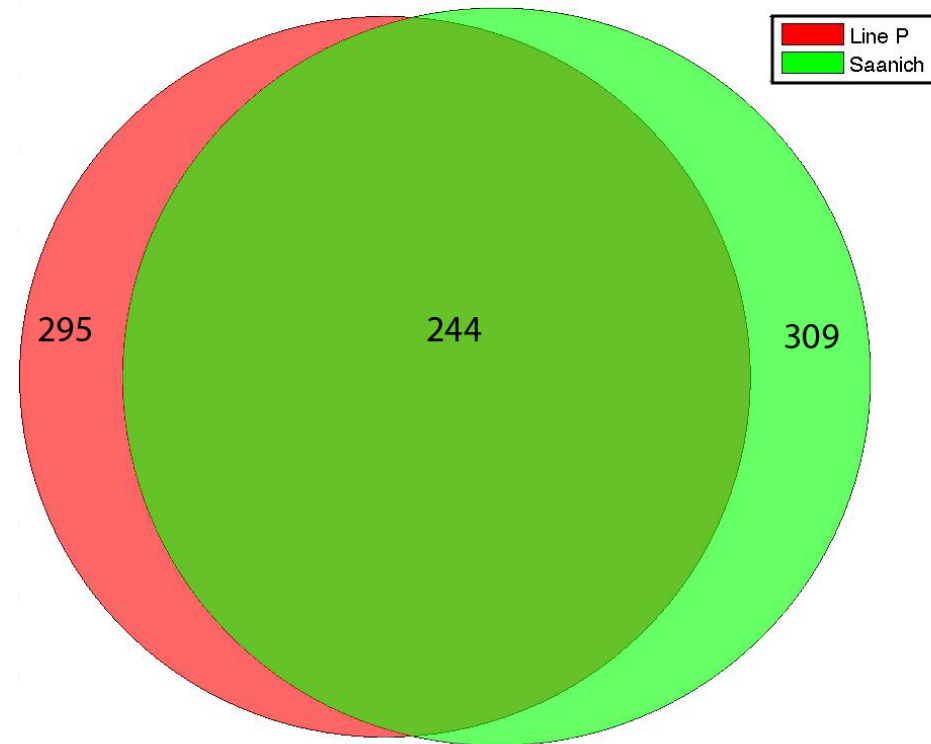
Pathway Distribution:

1012 pathways shared at some point



Core Pathways:

244 pathways always present

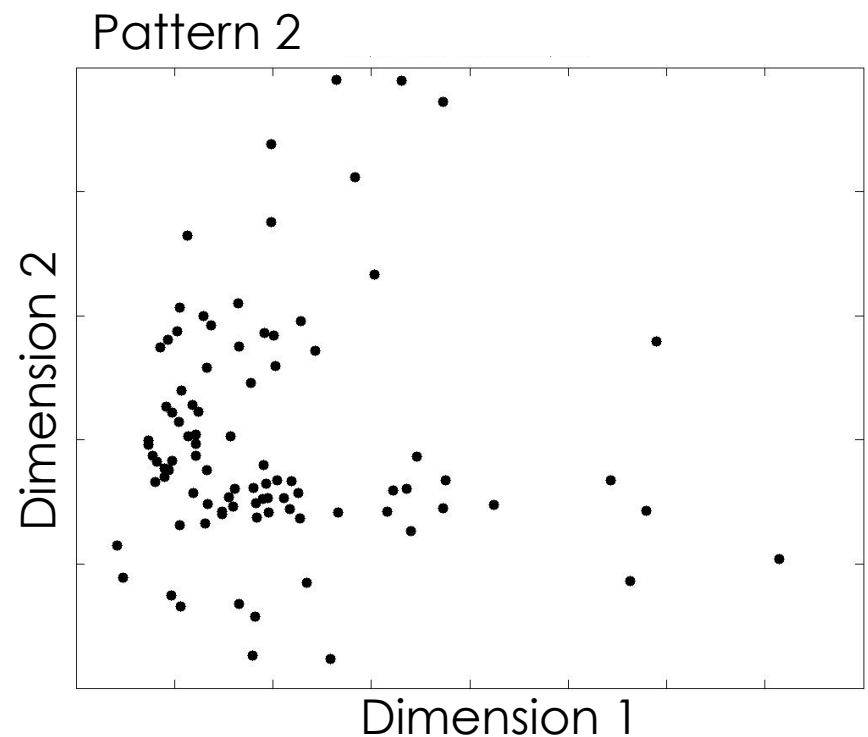
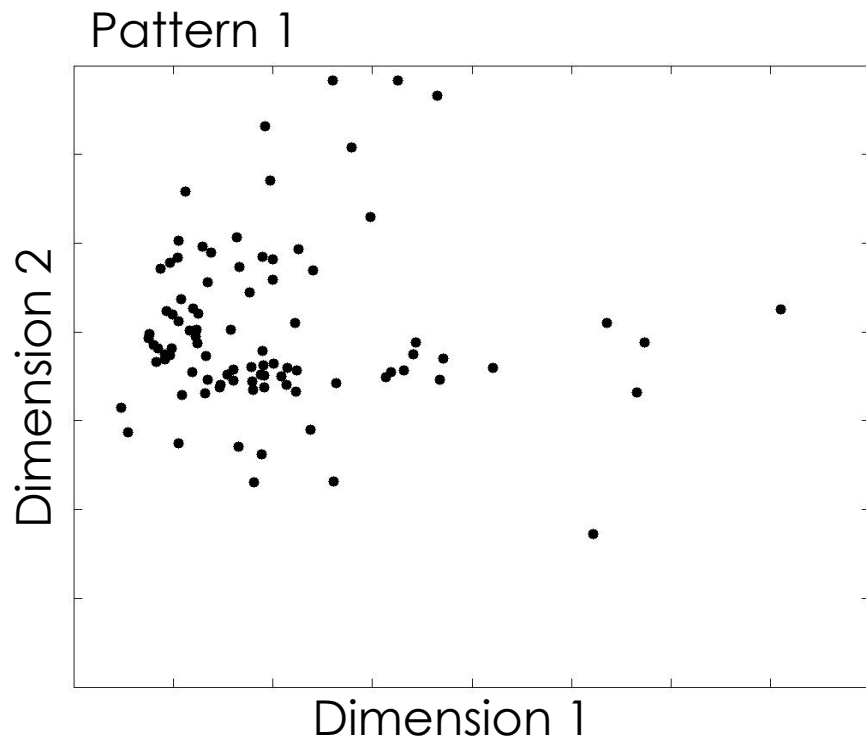


- Follows **119 pathways are NOT shared**
- Of which **35 are exclusive to Line P**
- And **84 are exclusive to Saanich**

NESAP Illumina Dataset: NMS

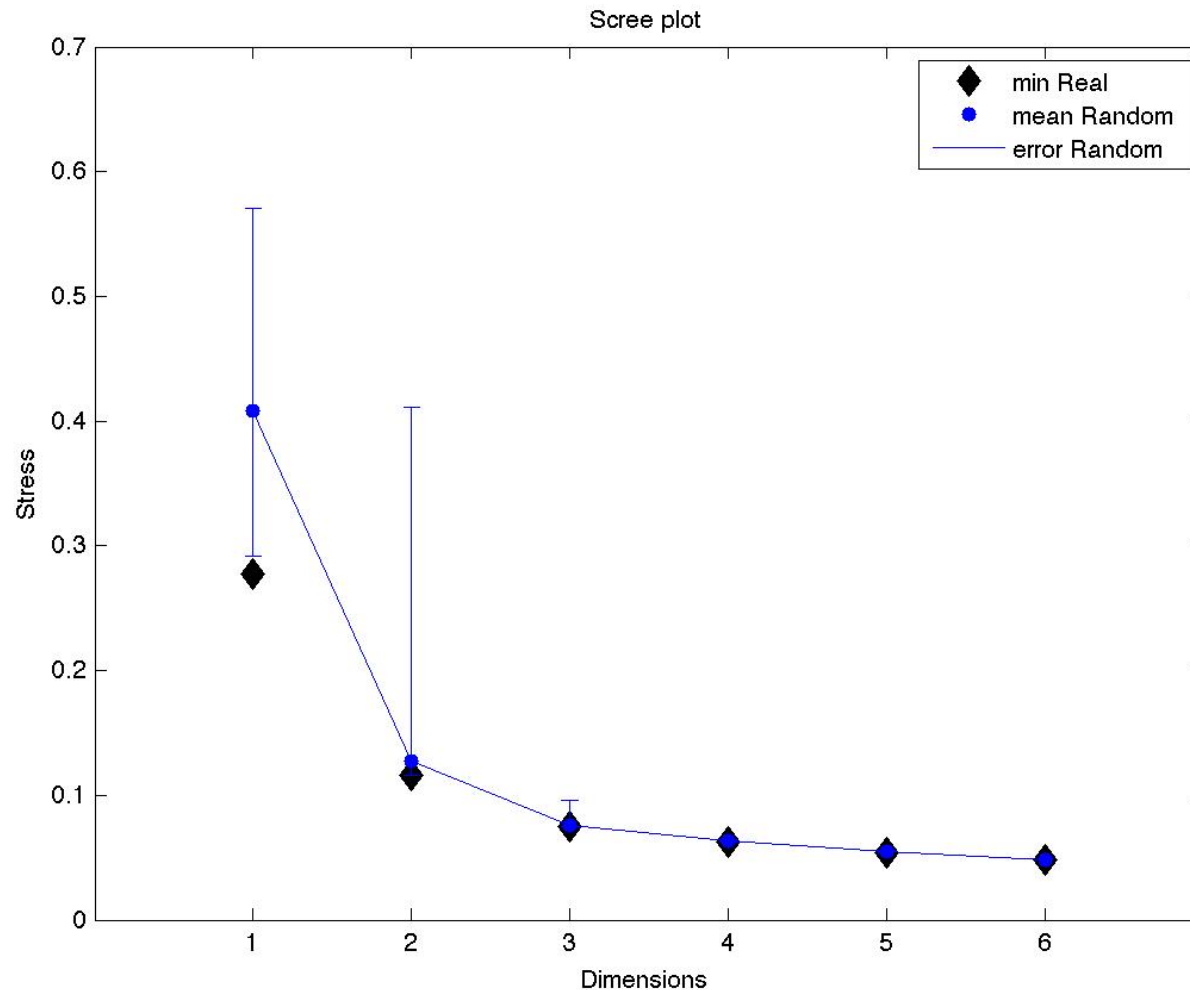
Many NMS runs consistently reveal a stable configuration, resulting in 1 of 2 similar patterns

➤ Average stress $1 \times 100 \sim 12$ (good/fair)



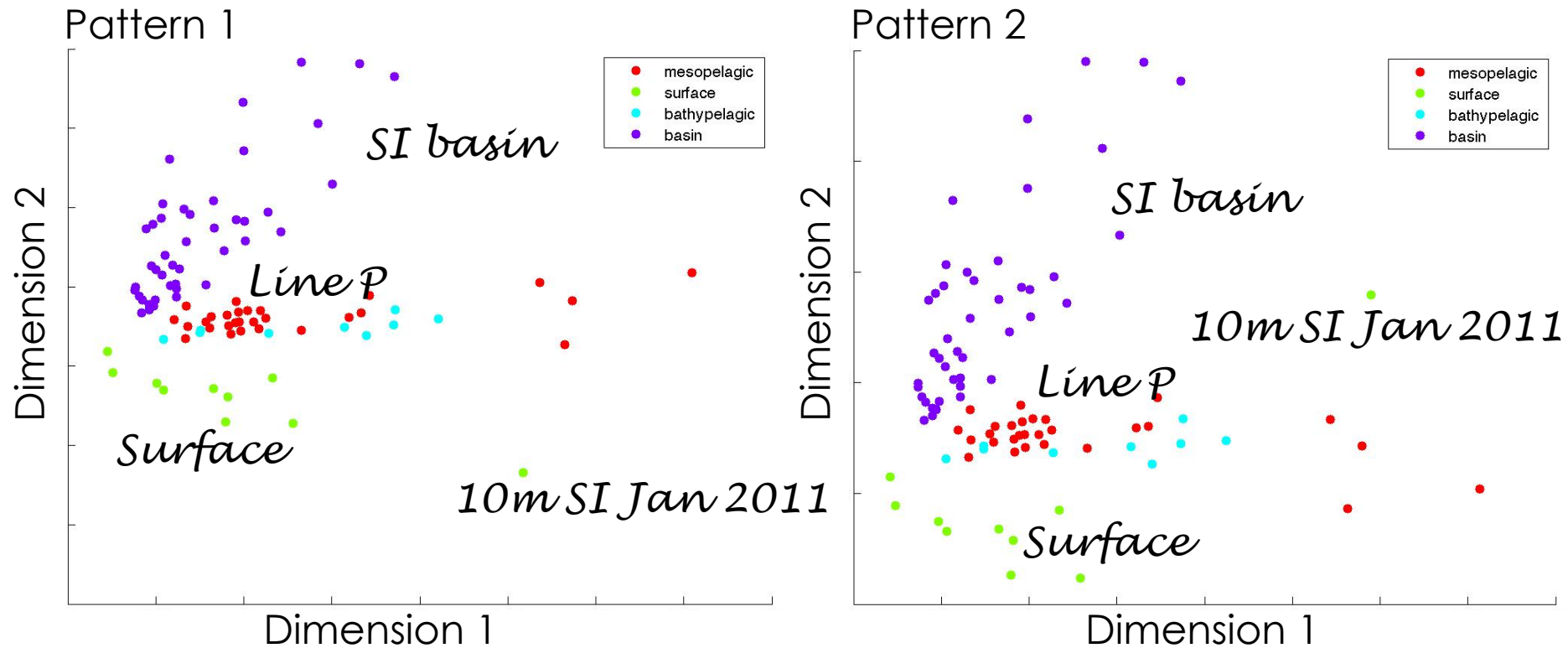
NESAP Illumina Dataset: NMS

Monte Carlo reveals 2 dimensions for jgi illumina produces the greatest reduction in stress:



NESAP Illumina Dataset: NMS

Depth

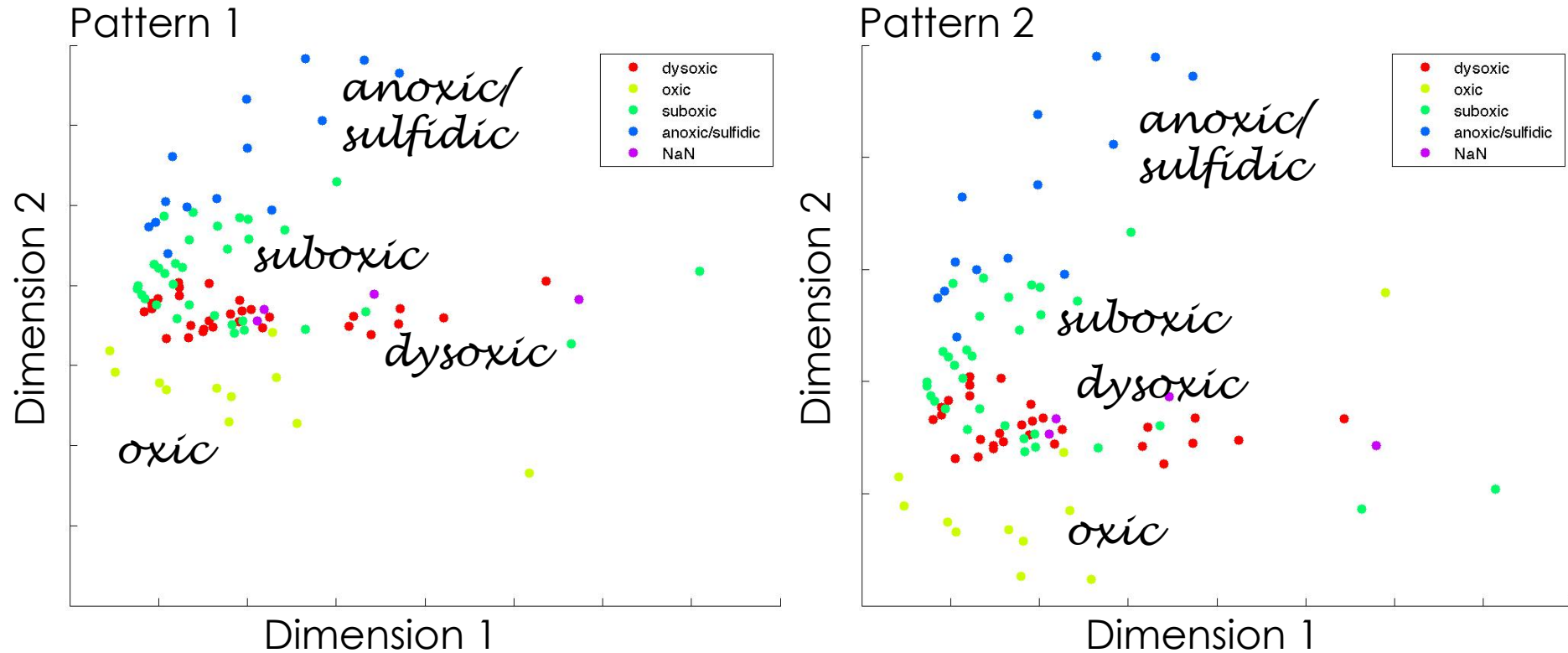


Patterns 1 & 2 give the same results

✓ *Pattern with depth*

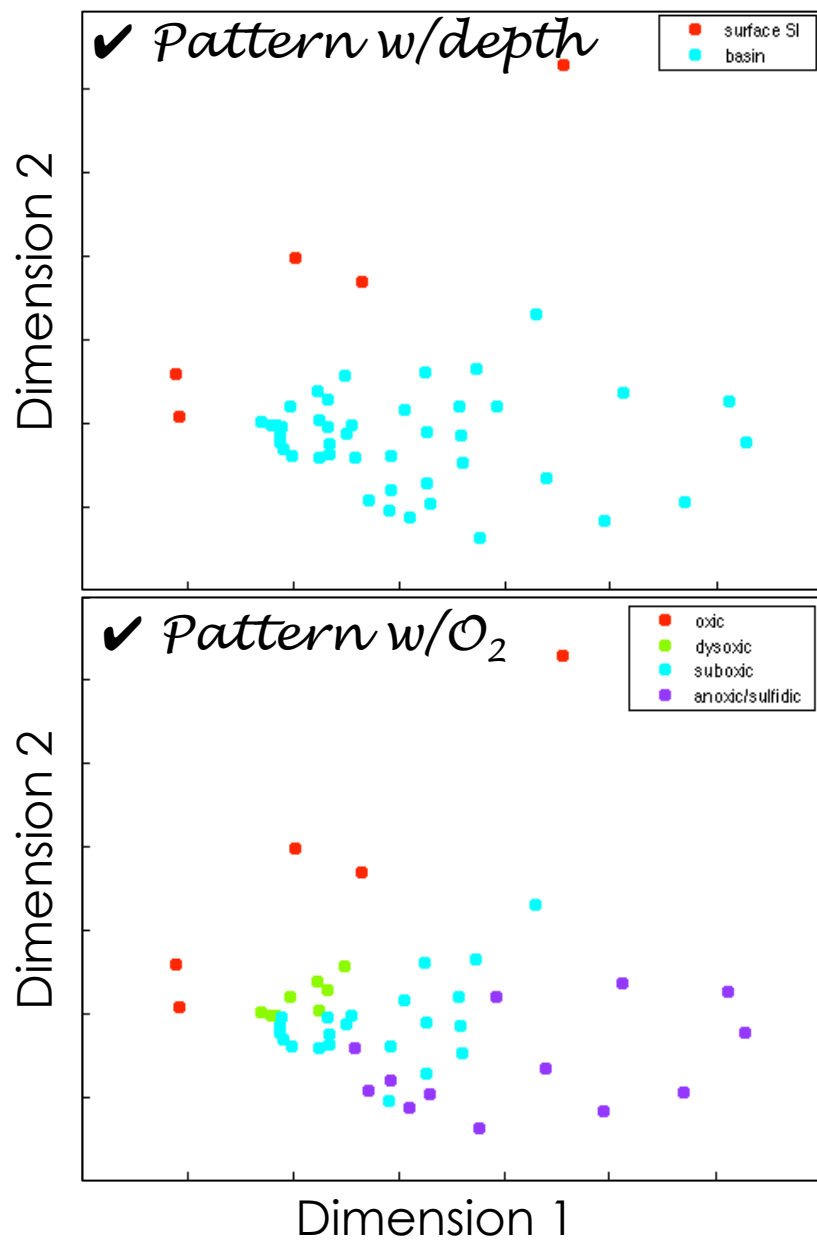
Pathways are different between samples at different depth classes

NESAP Illumina Dataset: NMS Oxygen

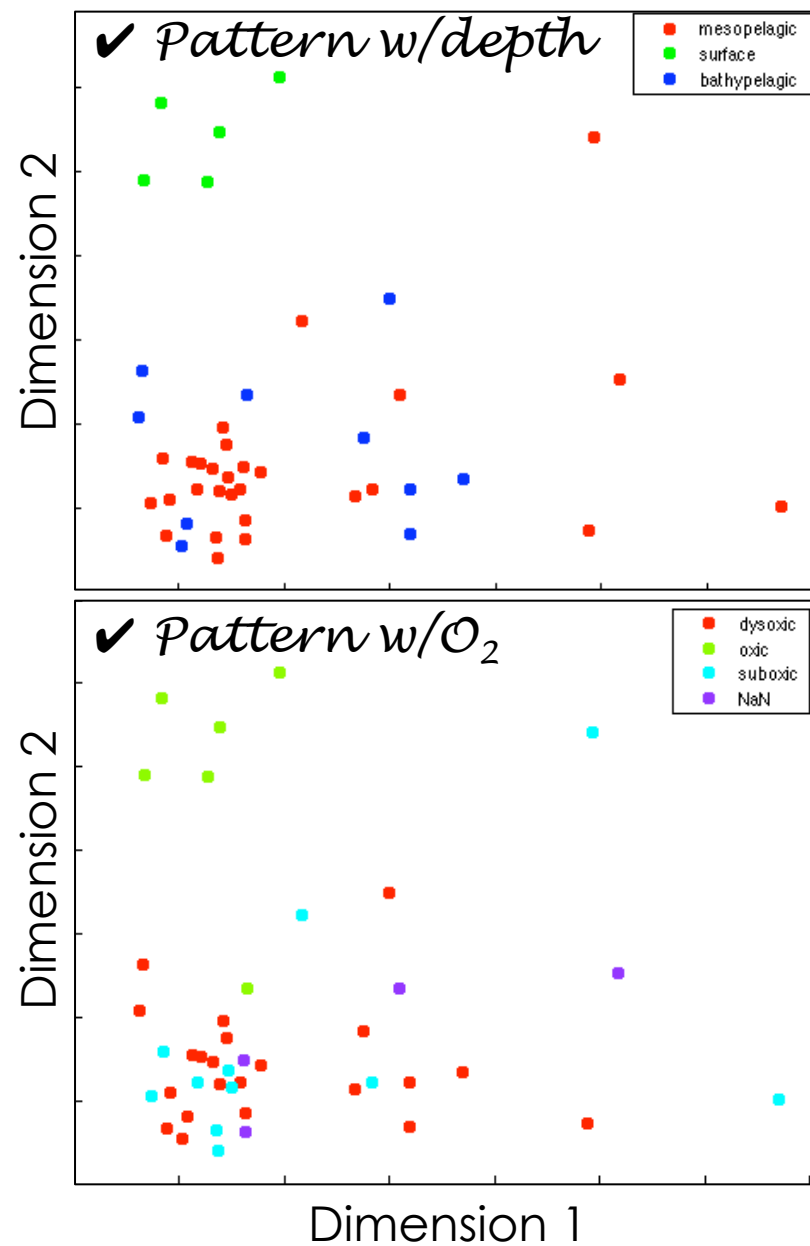


Patterns 1 & 2 give the same results
✓ *Pattern with oxygen*

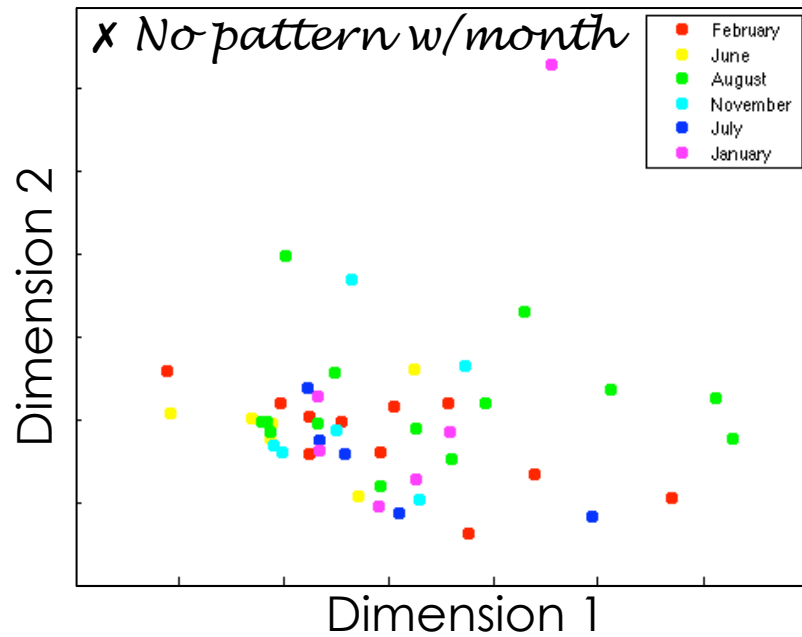
Saanich



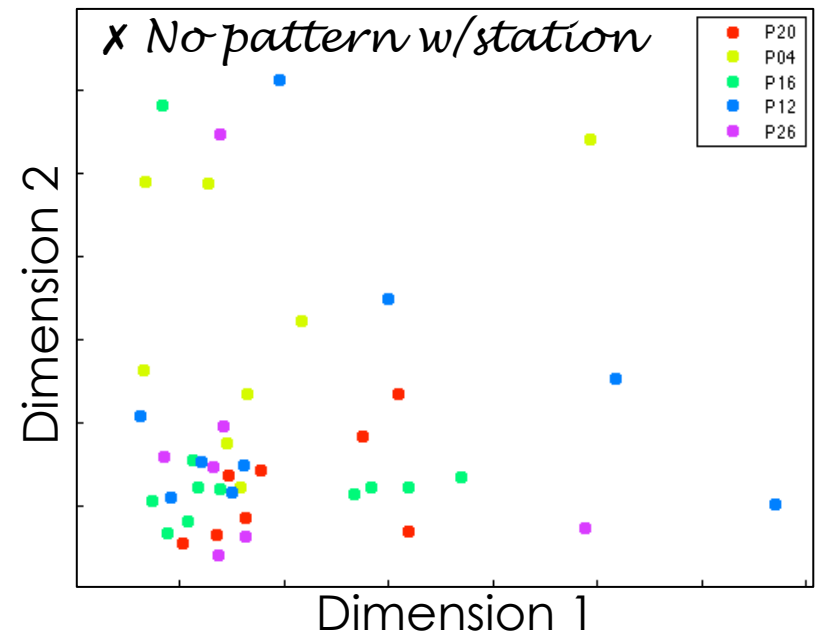
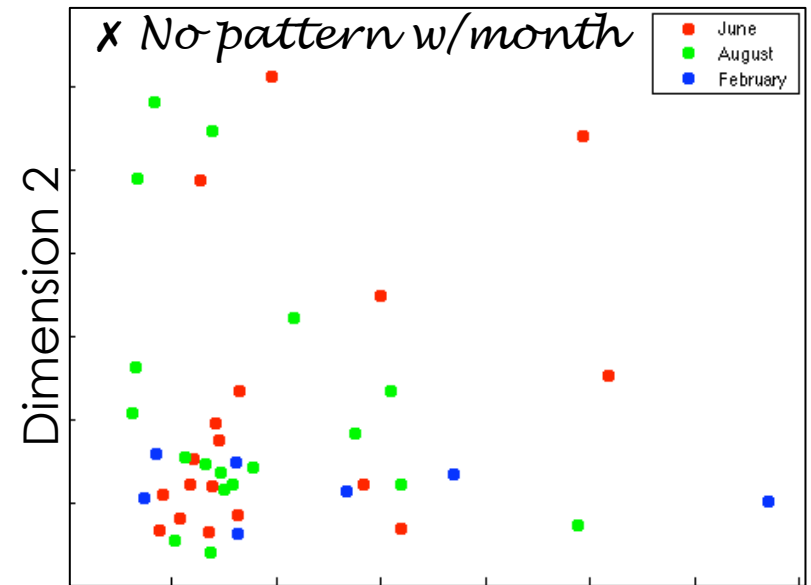
Line P



Saanich



Line P



Summary of NMS Trends for NESAP Pathways

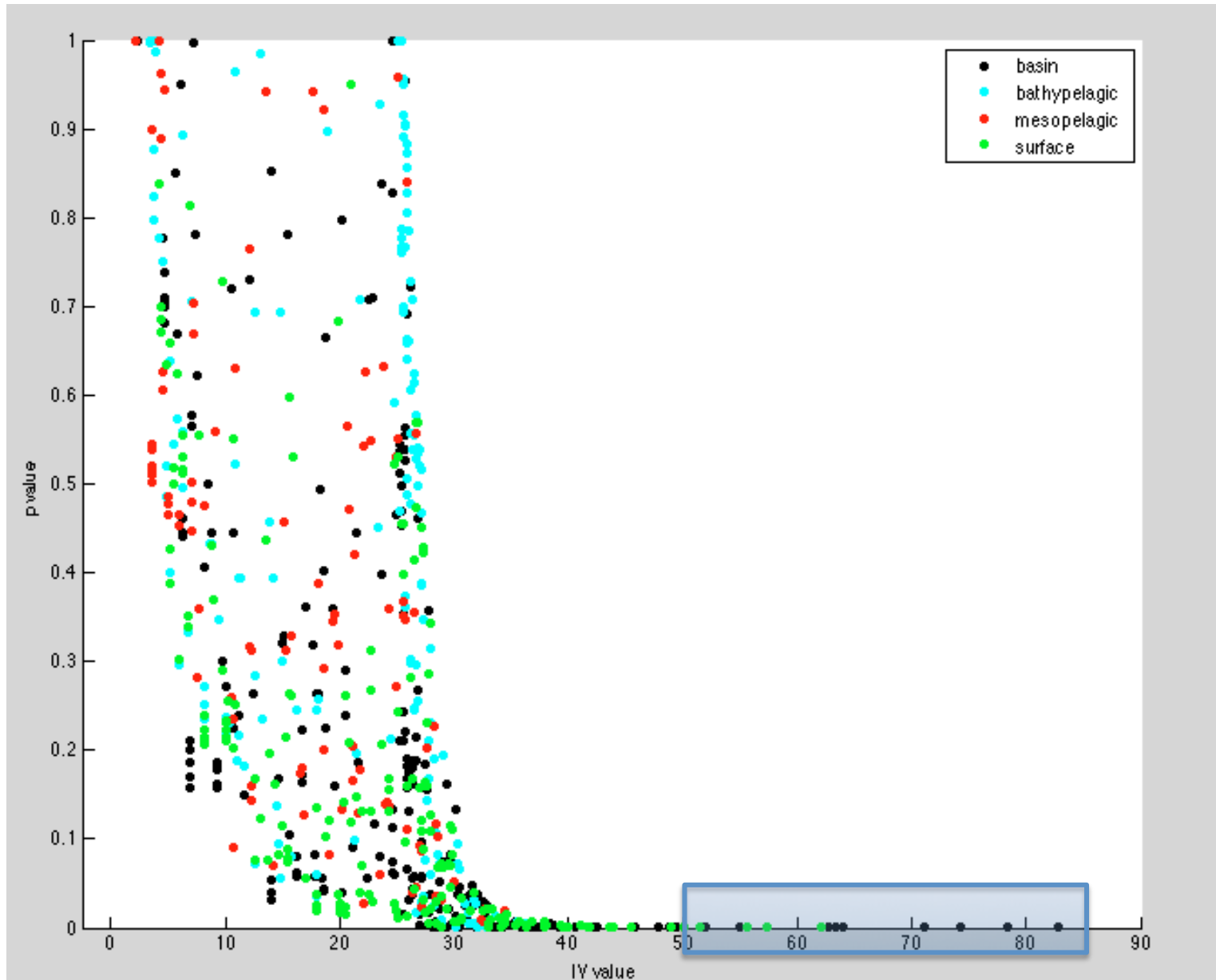
Groups	Depth	O ₂	Year	Month	Station	NO ₃ ⁻	H ₂ S
LP	✓	✓	x	x	x	possibly	N/A
SI	✓	✓	x	x	N/A	possibly	✓
LP + SI	✓	✓	x	x	N/A	possibly	N/A

- Pathways are different between samples at different depths & between samples with different oxygen concentrations
- Statistical basis for using groups for ISA based on 4 depth and oxygen classes

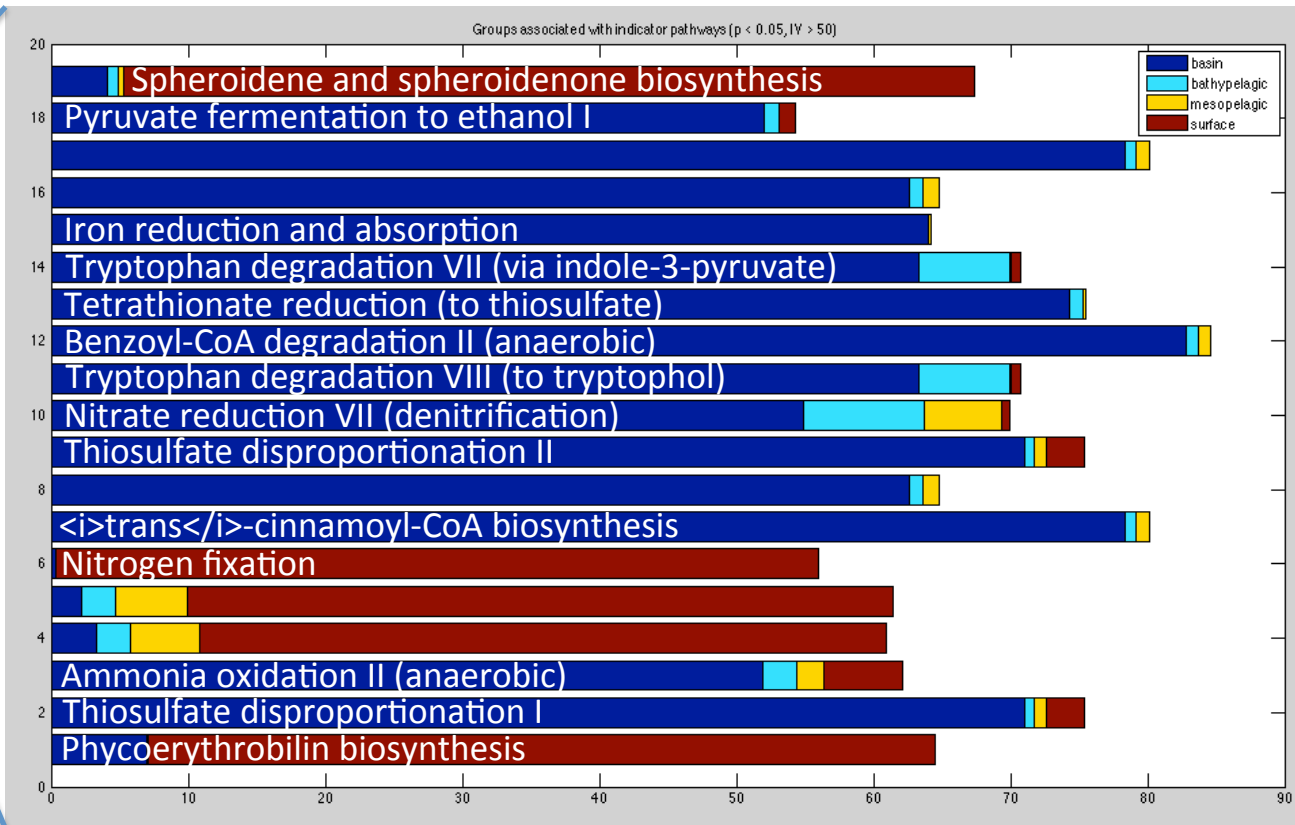
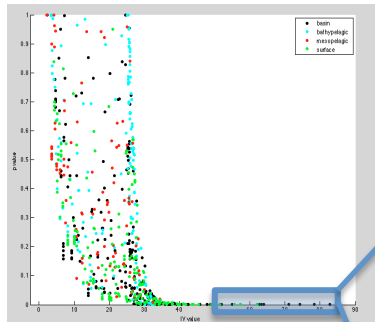
NESAP Indicator Species Analyses

- Will hopefully tell us which of the 1131 pathways are causing the groups to differentiate
 - To be a perfect “indicator” pathway for a group, a pathway needs to be both *exclusive* (always present in that group) and *faithful* (only present in that group)
-
- *We'll see how successful this analysis is considering the amount of shared pathways!*
 - *As a caveat, this analysis will be more successful with quantitative data*

Indicator Species Analysis (ISA)



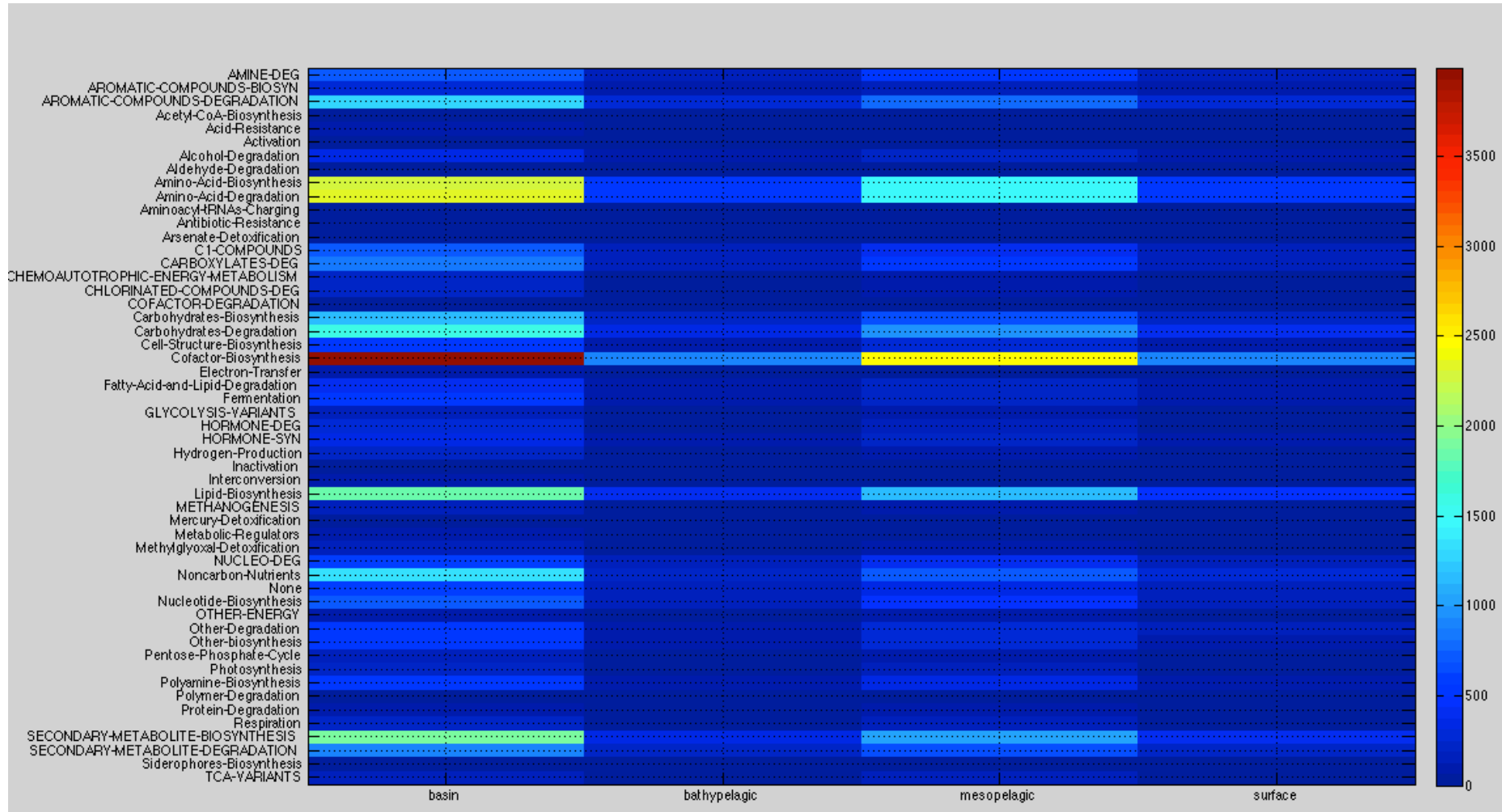
Indicator Species Analysis (ISA)



19 'indicator' pathways: 5 for the surface, 14 for the basin

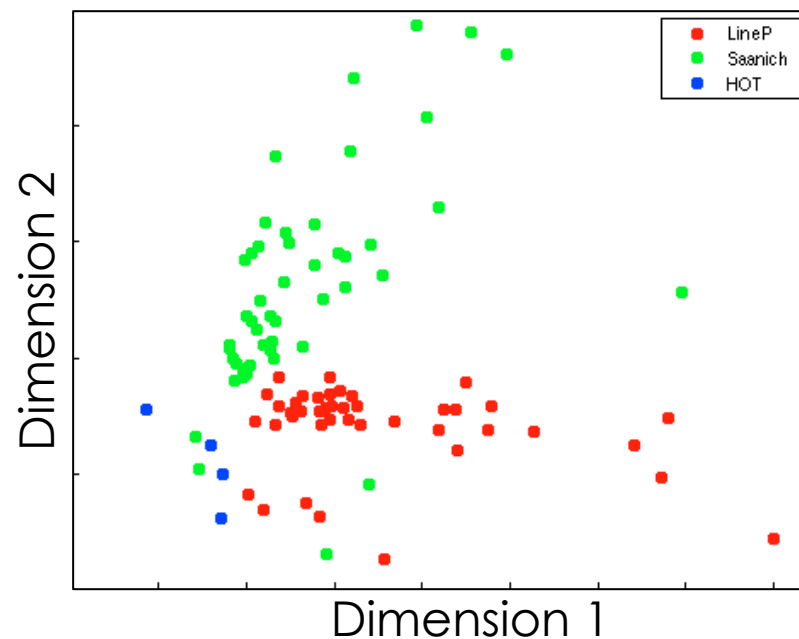
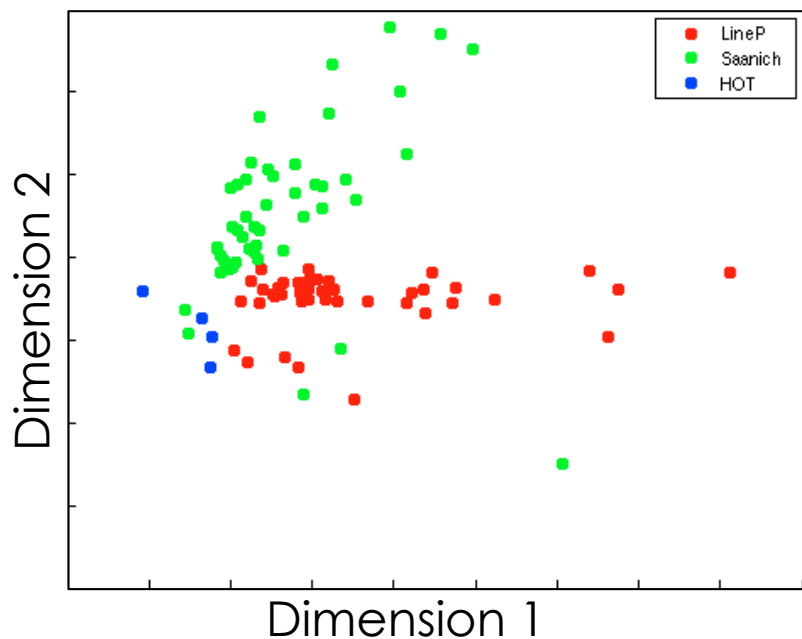
We can consider pathways major metabolic traits associated with two end-member environments (aerobic sunlit, anaerobic dark)

Quick Heatmaps



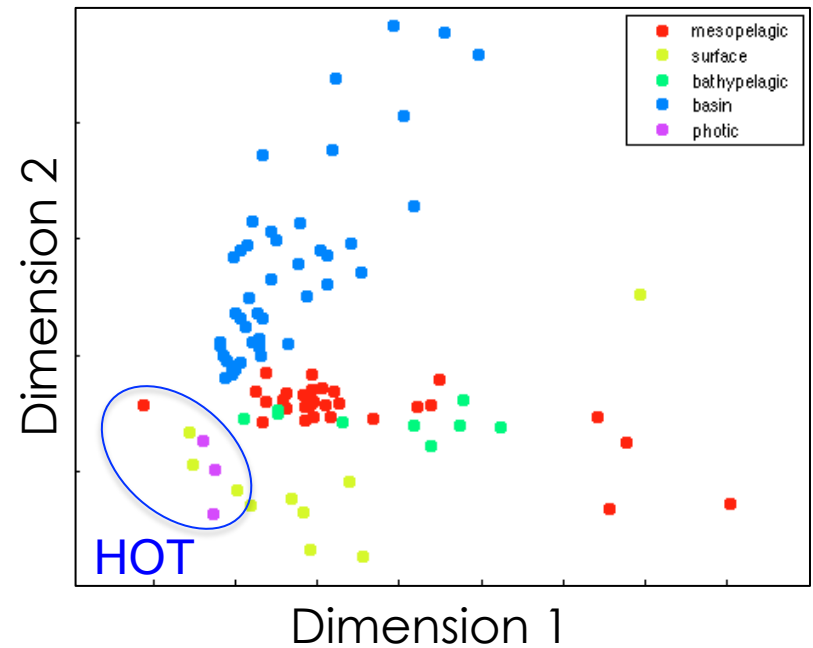
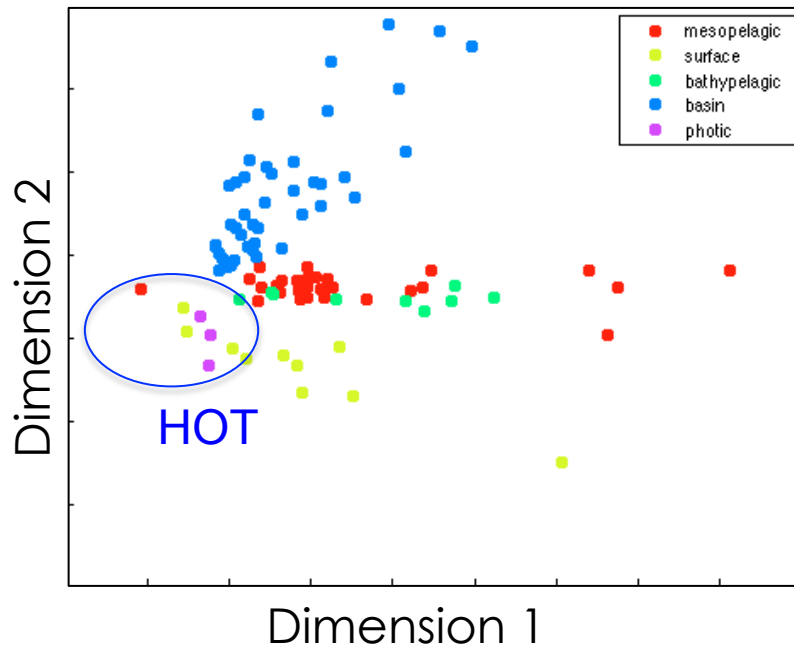
Adding in HOT...

- HOT 454 data from Cruise 186 October 2006
- 4 samples: 25m, 75m, 110m, 500m
- Combined master table (SI + LP + HOT):
 - 1144 pathways (+13 unique to HOT)



- Again, NMS reveals same patterns 1 & 2
- Stress is also similar, ~12 (good/fair)

Adding in HOT...



- HOT samples grouping with surface NESAP
- *Pathways in HOT most similar to those also found in NESAP surface*

Questions?