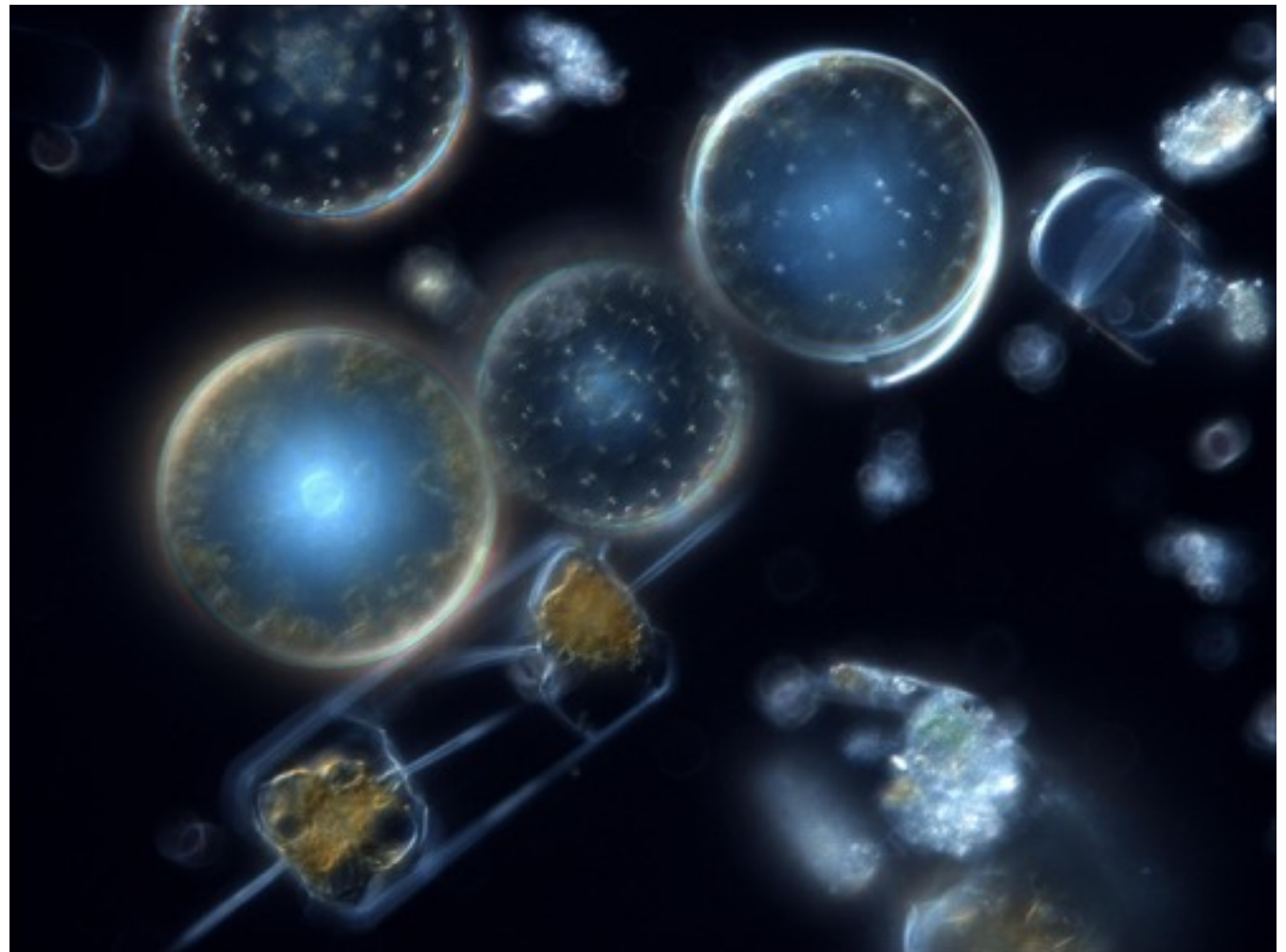


Analysis of large-scale patterns in phytoplankton diversity

Sophie Clayton

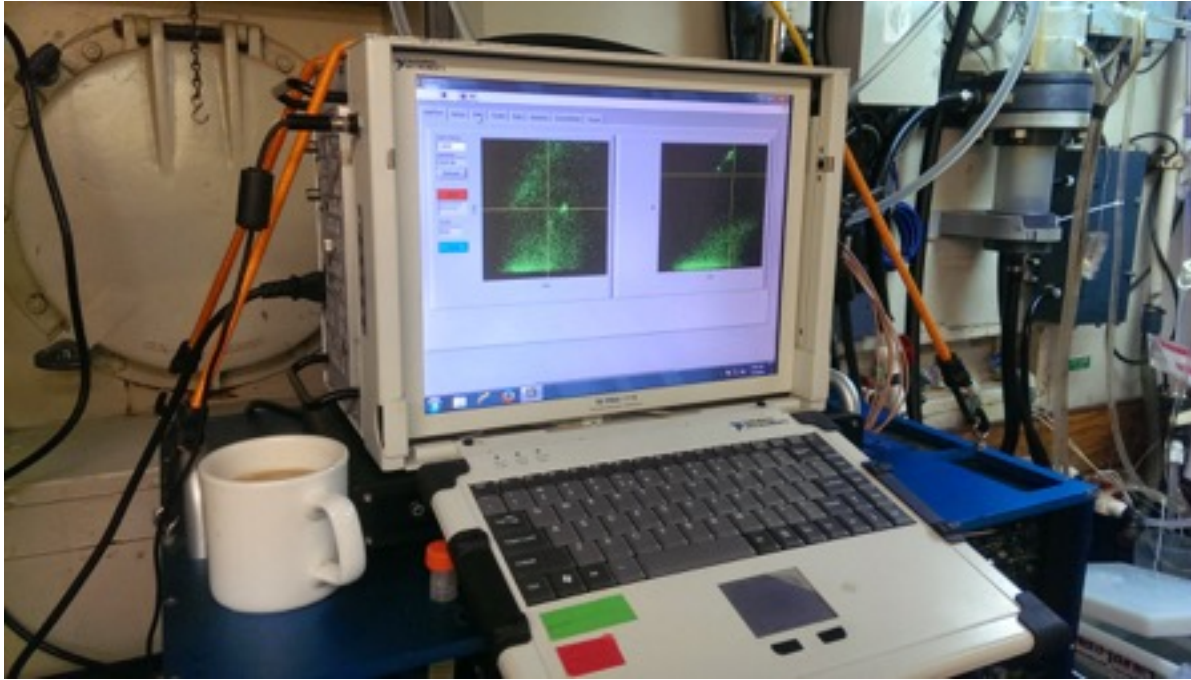
sclayton@uw.edu

December, 2014



School of Oceanography
University of Washington

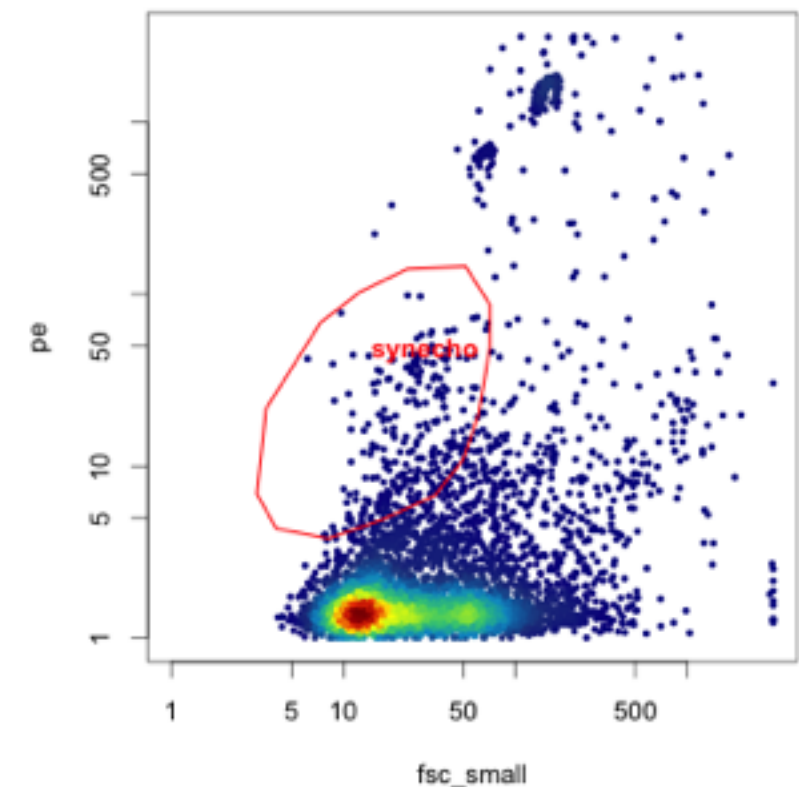
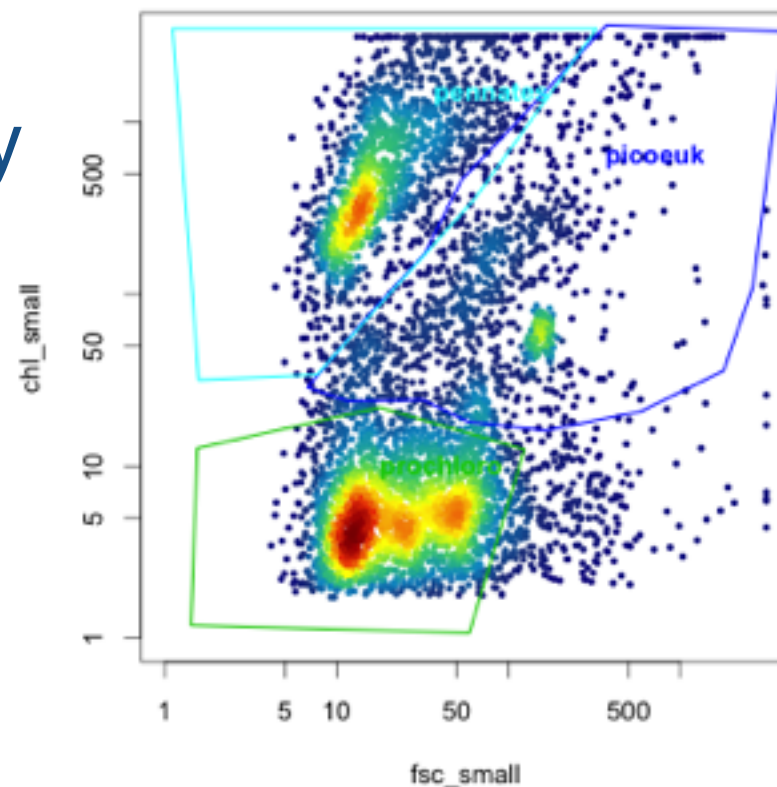
SeaFlow flow cytometer



SeaFlow in action aboard the R/V Melville

- Phytoplankton identified by their size and pigments.
- Beads used as internal standard.
- Estimate diversity from cytogram.

- Measures properties of particles in seawater.
- Data collected continuously and stored every 3 minutes.

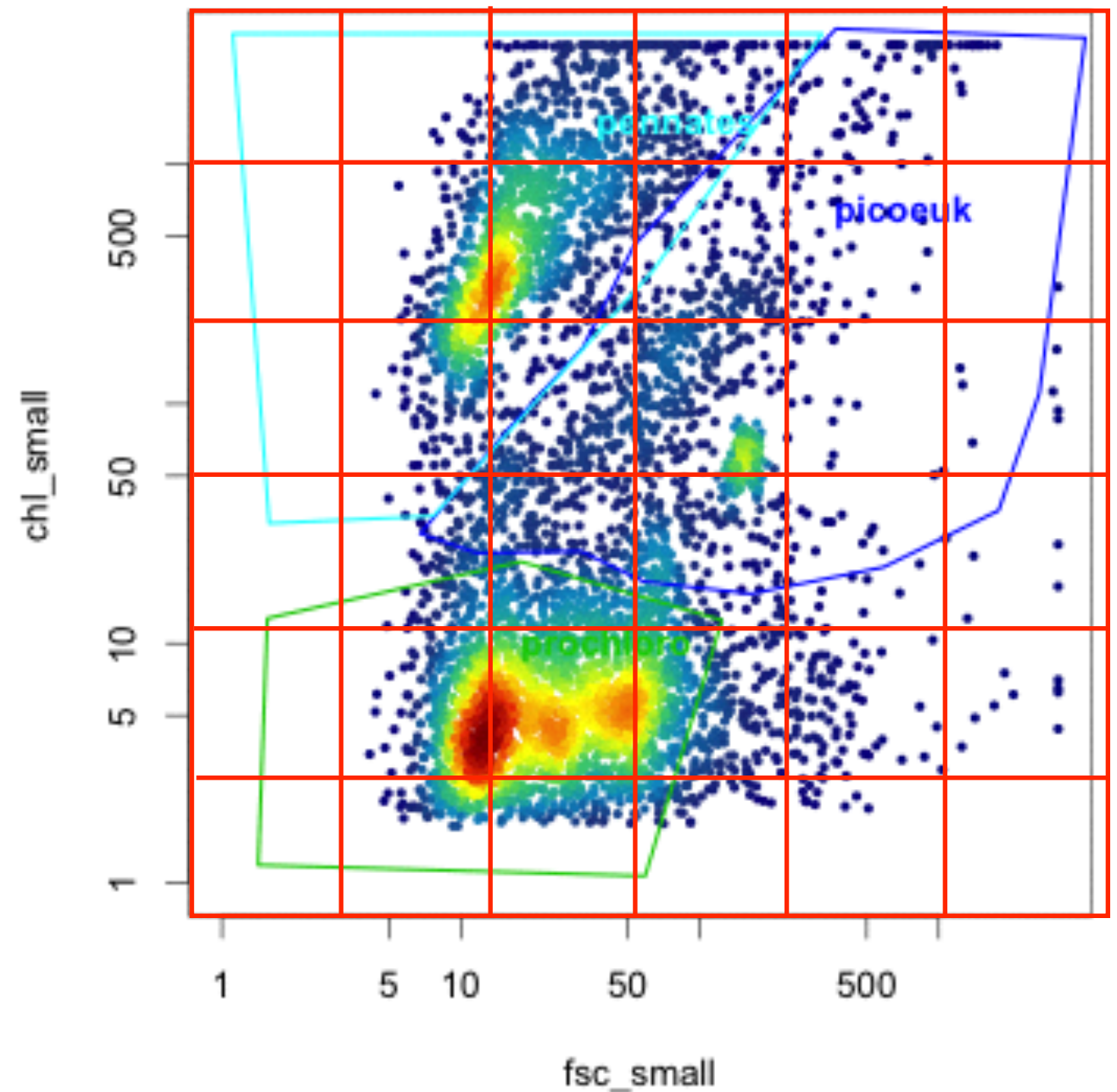


Example cytogram with population assignment

Diversity index from SeaFlow data

Taxonomic determinations of diversity are very labour-intensive.

Diversity indices are related to the area of the cytogram “occupied” by particles - doesn’t require formal identification.



Particles in each cytogram are binned to estimate diversity of observed population

The SeaFlow dataset: overview

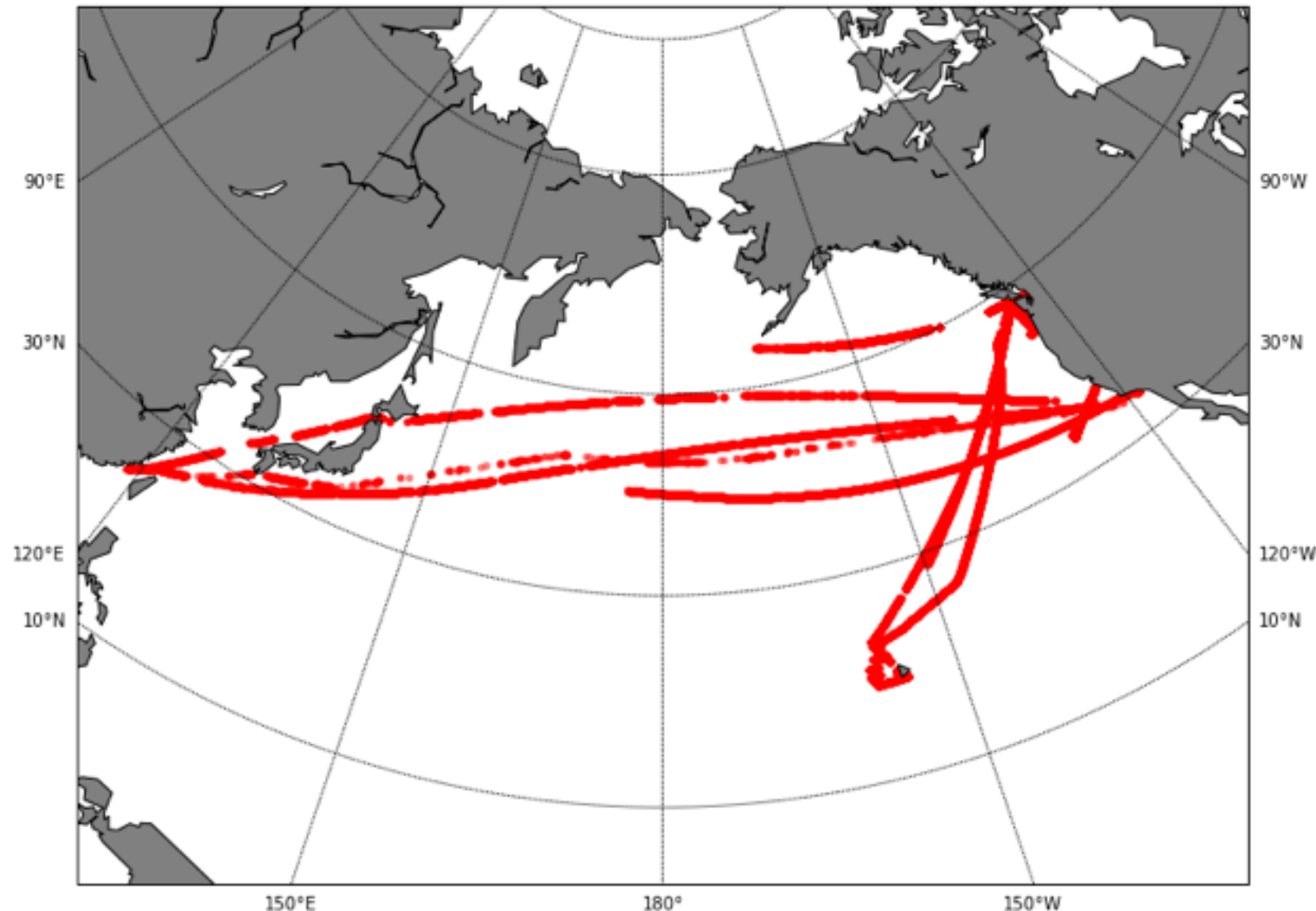
Goal: to explore patterns in phytoplankton diversity and their relationship to the environment in the North Pacific.

Data from 18 cruises in the North Pacific

- > 500GB data
- 1.7B particles
- > 50,000 data points

Data types:

- OPP - particle properties
- VCT - population assignment
- SDS - environmental data



Map of all SeaFlow data points

The SeaFlow dataset: Myria and SQLshare

- OPP and VCT data already uploaded to Myria
- SQLShare for interpolating SDS data onto OPP timestamp



- OPP, VCT and SDS data can be joined using Cruise, Day and File_ID.
- Myria used for filtering data and computationally-intensive queries (e.g. re-scaling, binning)

The screenshot shows the Myria web interface. At the top is a navigation bar with "Myria", "Editor", "Queries", and "Datasets". A notification banner at the top right says "There is an upcoming reservation for research experiments. The reservation will begin about 16 hours from now. For more information, please check the calendar." The main area is divided into two panels. The left panel is a code editor with the text "Write your code here, perhaps starting from one of the examples at the right." and a MyriaL query. The right panel shows the "Get the schema of a dataset" section with a search form and a table of dataset attributes.

```
1 div = scan(diversity_untrans_max_scaled_16bins);
2 chl = scan(armbrustlab:seaflo:tot_chl_byfile);
3 TS = scan(armbrustlab:seaflo:all_sds_v2);
4
5 div_TS = select b.N0, b.N, b.J, c.tot_chl, t.T, t.S,
6               t.LON, t.LAT, t.day, t.Cruise, t.file
7             from div as b, chl as c, TS as t
8             where b.Cruise = t.Cruise
9                  and t.Cruise = c.Cruise
10                  and b.Day = t.day
11                  and t.day = c.Day
12                  and int(b.File_Id) = t.file
13                  and t.file = int(c.File_Id);
14
15 store(div_TS, armbrustlab:seaflo:div_chl_sds);
```

Execute the Query Parse Myria JSON

Query Language MyriaL

Developer Options

Get the schema of a dataset

Use the search form to retrieve the schema of a dataset which includes its columns and column types.

armbrustlab:seaflo:good_opp_vct_v4

More details: armbrustlab:seaflo:good_opp_vct_v4

Name	Type
time	INT_TYPE
pulse_width	INT_TYPE
d1	INT_TYPE
d2	INT_TYPE
fsc_small	INT_TYPE
fsc_perp	INT_TYPE
fsc_big	INT_TYPE
pe	INT_TYPE
chl_small	INT_TYPE
chl_big	INT_TYPE
Cell_Id	LONG_TYPE
Cruise	STRING_TYPE
Day	STRING_TYPE
File_Id	STRING_TYPE
pop	STRING_TYPE

MyriaDB vs. R: calculating diversity

[illegible]

```
-- Load the existing dataset
AllData = scan(armbrustlab:seaflo:good_opp_vct_v4);

-- Assign a linear value into one of 16 bins 0..15
-- N.B.: // is integer division
def makebins(x): x//(pow(2, 16)/16);

-- For each cruise & sample (day + file_id)
-- break the 3-D cytogram given by forward scatter,
-- chlorophyll, and phycoerythrin into a 16x16x16
-- bin space and count the number of cells in each bin.
AllDataBinned = select Cruise, Day, File_Id,
                      makebins(fsc_small) as fsc_bin,
                      makebins(chl_small) as chl_bin,
                      makebins(pe) as pe_bin,
                      count(*) as num_particles
                      from AllData;

-- Compute the Richness N0 as the number of full bins
Richness = select Cruise, Day, File_Id,
                  count(*) as richness
                  from AllDataBinned;

store(Richness, richness_untrans);
```

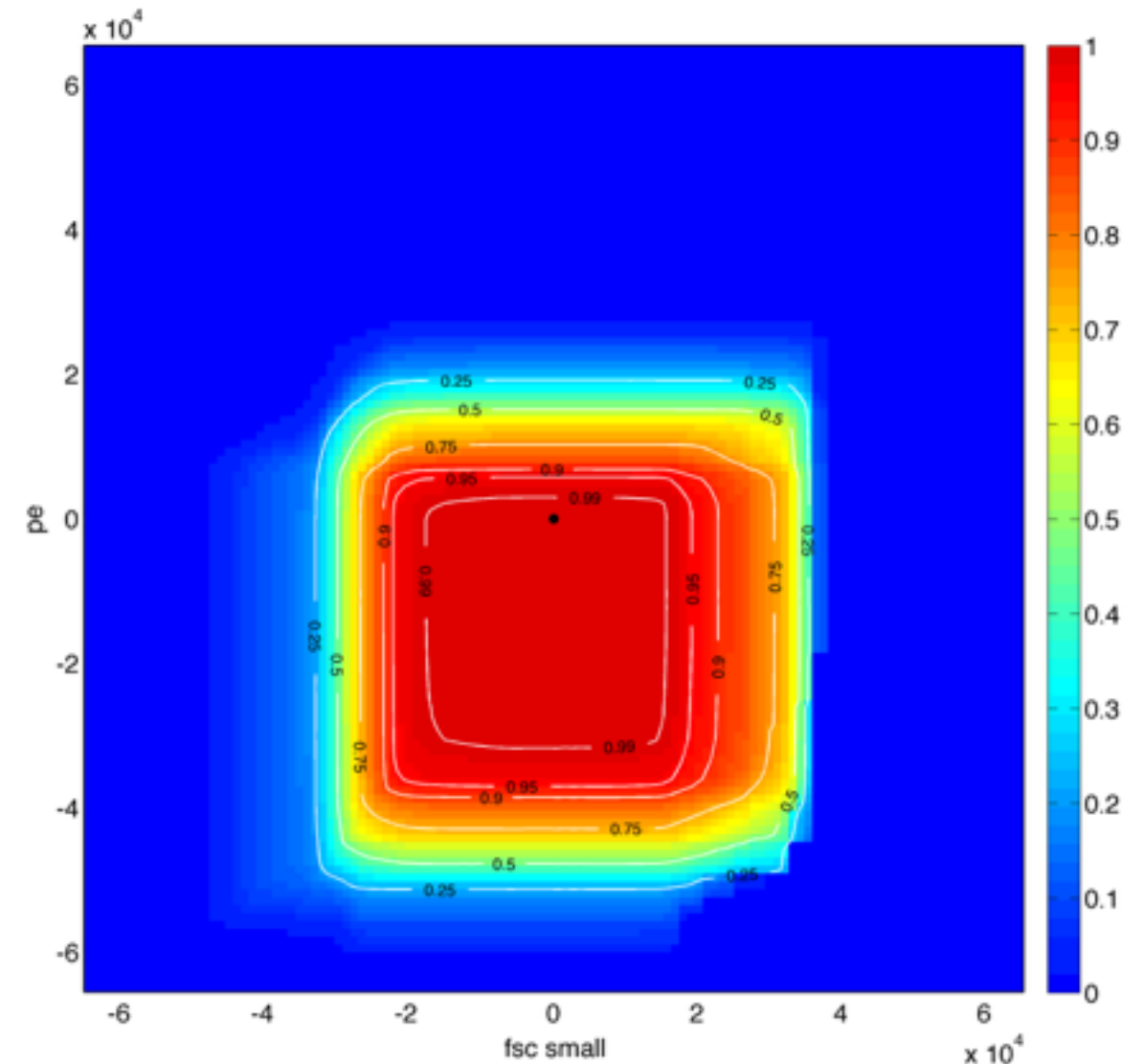
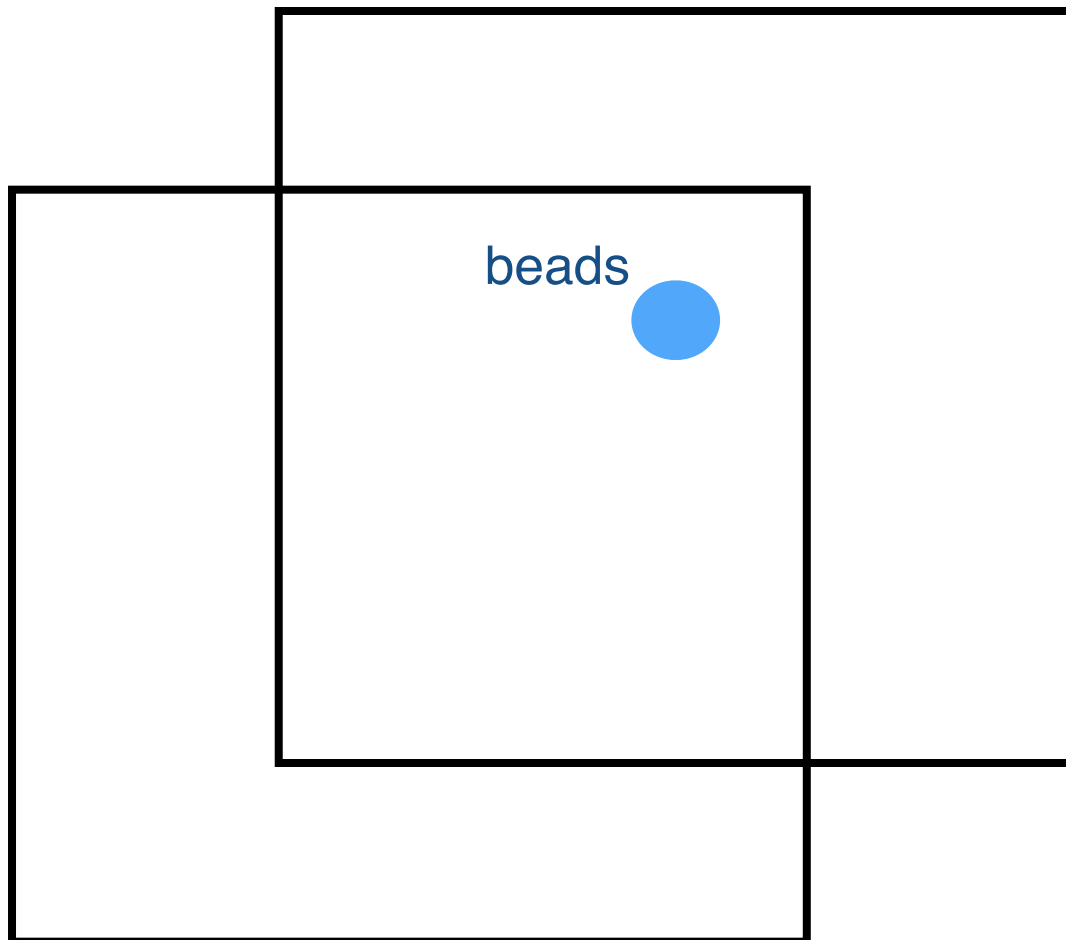
Run time: several hours

Run time: ~10 minutes

Normalizing SeaFlow data

Need to account for variable SeaFlow settings, and different instruments.

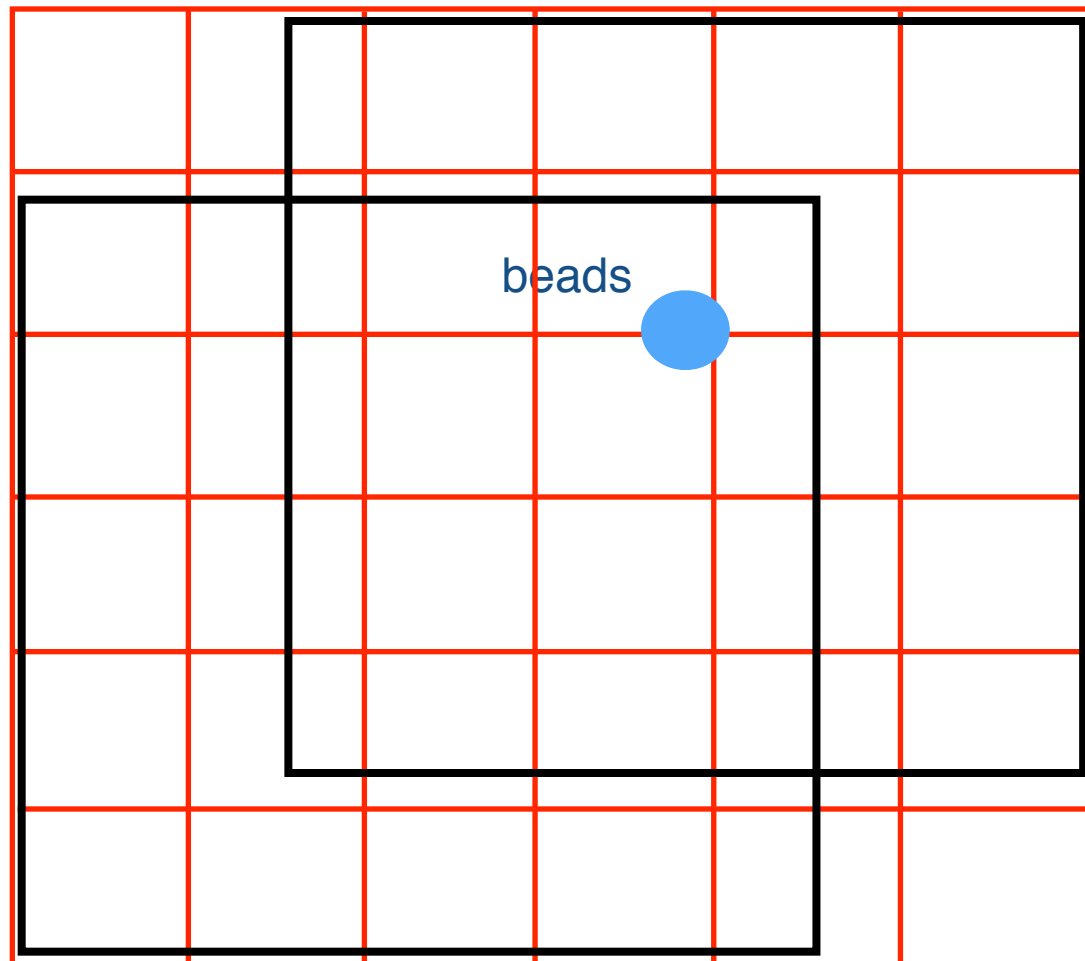
Use beads to standardize.



Proportion of file coverage around bead position

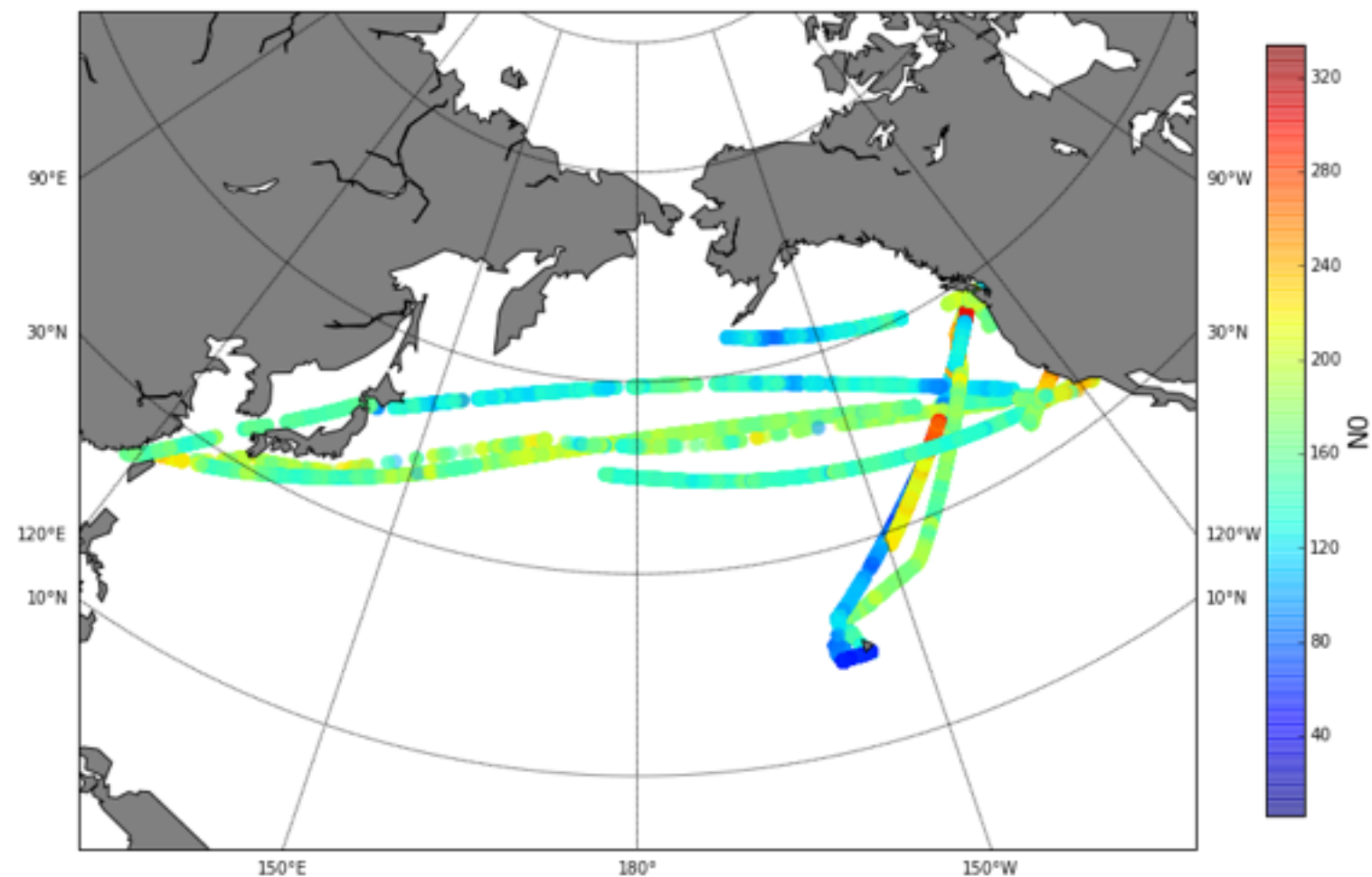
Particle properties are re-scaled to be relative to the value of the bead properties in each file.

Normalized estimates of diversity



Example of how the normalized data is binned

Diversity indices calculated over re-gridded normalized OPP data.

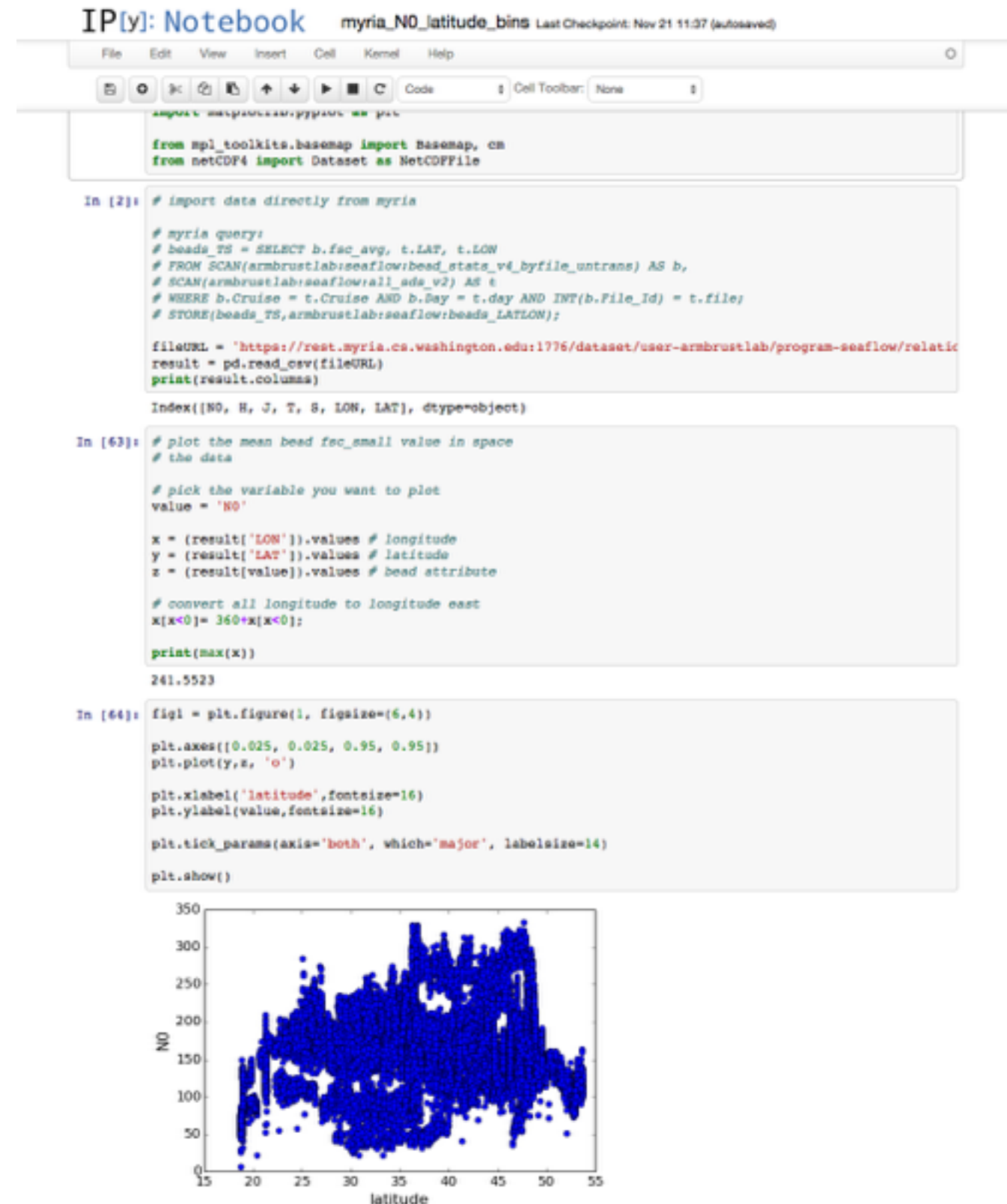


Map of bead normalized cytometric diversity

Further analysis and visualization

Download data directly from Myria REST server into Python.

Document analysis and visualization in iPython notebooks stored in GitHub repo.

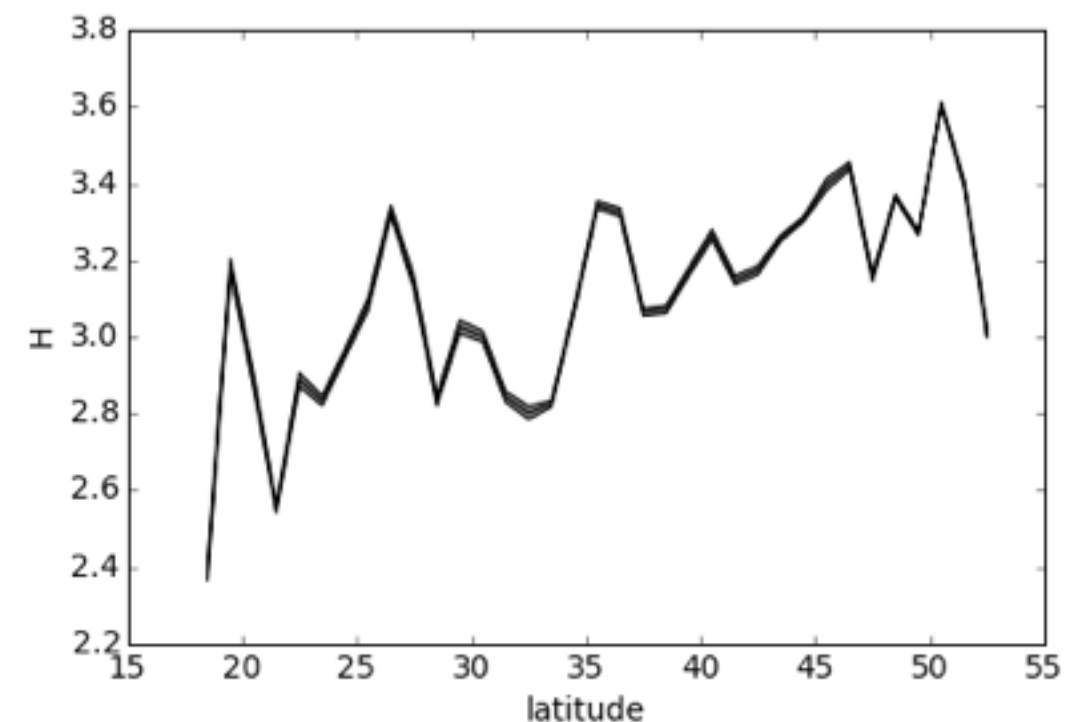
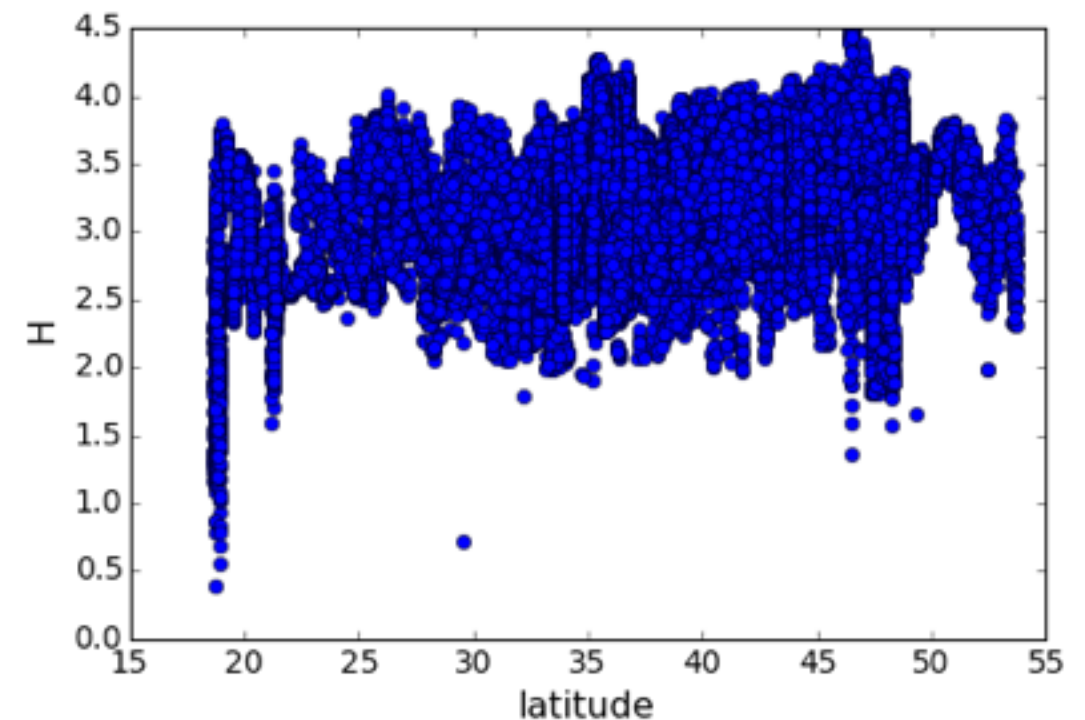


Conclusions & Outlook

We have developed tools to work with the SeaFlow dataset in Myria and Python.

Next steps:

- explore seasonal patterns in data
- split data into coastal vs. open ocean sets
- sensitivity analysis of the binning parameters for the diversity indices



Zonally averaged cytometric diversity

Thanks!

- Dan Halperin
- Bill Howe
- eScience Fall 2014
Incubator staff and
participants



sophieclayton commented 10 days ago

just tried this query:

```
good_opp_vct = scan(armbrustlab:seafLOW:good_opp_vct_v4);  
  
day = select Cruise, File_Id, substr(Day, 5, len(Day)) as yearday  
      from good_opp_vct;  
  
store(day, armbrustlab:seafLOW:yearday_opp_vct);
```

and got this error:

Error 400 (Bad Request): Error 400 (Bad Request): The second argument of substr has to be INT.

Since when is 5 not an integer?!...



dhalperi commented 10 days ago

Owner

wtf? I will look into this.