# Different routes you could take to analyze "18S" datasets

Questions are:
What do you lose in the analysis?  What form of payment? How much flexibility / conformity?

Engage with a consultant:
James White Resphera biosciences: paid consultant.  Does not require publication/ requires $.

Engage with a publicly financed project -- like the one we are working on.

Upload your data to a website:
https://www.arb-silva.de/ngs/

Use a AWS instance of qiime or mothur (or other software).
https://forum.qiime2.org/t/basic-ec2-aws-qiime2-install-help/5655
OR
Use a 'server' in some other context.
What the above approaches involve is not working with the data directly on your own computer.

However we can use qiime / mothur / usearch if you are able to install them
(and can deal with command line).

Other options
Commercial software
CLC genomics has a java implimentation of _____ and if you have $ you can use it.  Sometimes you c

MacVector: Concatenate reference database and map reads to it.

Within data

Observations
or 'draws'

Binning
Clustering
Trimming ends
Organizing sample names
Assembling
Sorting
Removing chimeras
...
Removing other artifacts
Removing background

Comparing to reference
Putting a name on it

Calculate frequency based on OTU
Resolving OTU from each other?
List of OTU
Think of this as the diversity
in the sample.

Animating that result

Usearch is great and self-referential (largely)

What is the cost of a 'wrong' sequence?
Chimera
Other artifact

What we did last week:

I handed you a set of __ sequences (drawn from Tara Oceans project)
We took one and used blastn and saw it had a 'good' --
          very high identity match across the entire sequence to a Diplonemid.
We then digressed to command line
          Look at sequence in 'text editor' WYSIWYG
Then attempt to run mothur
Essentially the next steps were like boiling down gravy -- reducing redundant sequences,
calculating dominant or most common sequences
•••Three window view -- mothur window, list of files in a window of the folder, text editor of results•••
We then could take the list of sequences and run blast on the command line.
If asked, we could summarize the results in different ways, with different levels of detail and specificity.
Started with 65, removed to 35 based on x, y and z.
Determined that Diplonemid was dominant sequence insample (and thus likely in sample of water).
We skipped trimming the sequence based on a reference alignment, chimera removal,
and display of results.

Start of yeast sequence

End of yeast sequence

5'

3'

*Saccharomyces cerevisiae*
(U53879)
1.Eukaryota 2.eukaryote crown group
3.Fungi/Metazoa group 4.Fungi
5.Ascomycota 6.Saccharomycetales
7.Saccharomycetaceae 8.Saccharomyces
May 2004

SSU ITS1 5.8S ITS2 LSU

100%

_O. magnificus_

50%

_Histioneis_ sp.

_Dinophysis_ sp.

0Kb     1Kb     2Kb     3Kb   3.55Kb

*Ornithocercus magnificus* FTL83
*Ornithocercus magnificus* FTL 1
*Ornithocercus magnificus* FTL 77
*/* *Ornithocercus magnificus* FTL 123
99/94 *Ornithocercus magnificus* CBC4 7
*/* *Ornithocercus quadratus* GS1a 41
91/92 *Ornithocercus steinii* FTL 203
88/90 *Histioneis* sp. FTL 70
*Histioneis* sp. FTL 62
97/95
*Dinophysis* sp. FPIP
*Dinophysis caudata* FTL 69
98/* *Dinophysis caudata* CBC4 8
*/* *Dinophysis caudata* FTL 93
*/* *Dinophysis acuminata* CBC4 L3
93/80 *Dinophysis acuta* CBC4 L10
*/* *Dinophysis* sp. FTL 124
*Dinophysis* sp. FTL ?
99/99 *Phalacroma rapa* CBC4 201
*/* *Phalacroma rapa* CBC4 L5
*/* *Phalacroma rapa* FTL 67
*/* *Phalacroma* sp. FTL 61
*/* *Phalacroma* sp. FTL 110
96/92 *Phalacroma* sp. CBC4 L128
*/* *Phalacroma rotundatum* FTL 71
*/* *Phalacroma rotundatum* FTL L10
*Phalacroma rotundatum* FTL 121
*Prorocentrum micans* CCMP 1589
*Prorocentrum minimum* SERC

0.03

*/97/99 Ornithocercus magnificus CBC4L7
Ornithocercus magnificus FTL1
99/88/66 Ornithocercus magnificus FTL83
Ornithocercus magnificus FTL123
Ornithocercus magnificus FTL77
gs1a 23
gs1a 65
cbc4 68
cbc4 18
gs1a 46
gs1a 5
gs1a 70
ffl 46
gs1a 42
gs1a 43
gs1a 77
cbc3 47
gs1a 21
gs1a 16
cbc4 14
*/96/* gs1a 37
gs1a 61
gs1a 92
gs1a 59
gs1a 26
gs1a 11
Ornithocercus steinii FTL 203
gs1a 38
gs1a 49
gs1a 30
cbc4 46
gs1a 39
gs1a 35
gs1a 51
cbc4 52
gs1a 52
*/94/* Ornithocercus quadratus GS1a 41
gs1a 10
gs1a 36
cbc4 61
gs1a 79
gs1a 93
gs1a 20

AU= 0.001

AU= 0.201

*/71/51

*Ornithocercus* group

gs1a 27
74/65/92 Histioneis sp. FTL 62
gs1a 76
92/73/56 gs1a 54
*/84/92 Histioneis sp. FTL 70
gs1a 64
gs1a 73
gs1a 63
gs1a 19
gs1a 71
81/90/* gs1a 55
gs1a 24
gs1a 33

95/56/-

*/*/97

*Histioneis* group

gs1a 28
Dinophysis caudata FTL93
Dinophysis caudata CBC4L8
Dinophysis caudata FTL69
ffl 18
Dinophysis sp. FPIP
*/94/97 Dinophysis caudata AY040584
Dinophysis caudata AF318241
Dinophysis caudata AF318240
*/93/95 Dinophysis tripos AY040585
Dinophysis odiosa AY259241
Dinophysis tripos AY259242
*/81/82 Dinophysis fortii AB355151
*/93/* Dinophysis acuta AY277645
66/79/99 Dinophysis acuta AY277648
*/*/- Dinophysis acuminata CBC4 L3
Dinophysis acuminata CBC4 L10
96/72/80 Dinophysis dens AY040571
Dinophysis dens AY040572
Dinophyceae AF318257
Dinophysis acuta AY040570
Dinophysis acuminata AY040573
Dinophysis acuminata AY040579
Dinophysis acuminata AF414683
Dinophysis sacculus AY040580
Dinophysis sacculus AY040581
Dinophysis norvegica AY259239
Dinophysis acuminata AY277638
cbc4 17

AU= 0.055

92/65/63

*Dinophysis* group

gs1a 13
gs1a 74
cbc4 19
*/97/* Dinophysis sp. FTL 124
gs1a 89
cbc4 59

*/99/99

Dinophysis sp. FTL ?

99/73/83 gs1a 45
ffl 66
*/99/98 ffl 69
FTL35
Phalacroma sp. FTL110
ffl 7
ffl 41
ffl 71
*/92/* cbc4 13
*/95/99 gs1a 78
*/63/73 */96/95 Phalacroma sp. FTL61
gs1a 66
*/*/95 Phalacroma rapa CBC4 L5l
/56/66 Phalacroma rapa CBC4 L201
Phalacroma rapa FTL 67
*/68/89 gs1a 17
cbc4 66
cbc4 44
cbc4 48
*/57/90 cbc4 25
*/96/96 gs1a 69
*/93/94 gs1a 14
cbc4 35
*/68/91 */99/99 cbc4 45
gs1a 18
*/99/* gs1a 31
gs1a 4
cbc4 62
92/76/90

AU= 0.00001

*/85/*

Dinophysis rotundata AJ506979
Phalacroma rotundatum FTL71
Dinophysis rotundata AF318235
*/95/93 Phalacroma rotundatum FTLL10
Phalacroma rotundatum FTL121
gs1a 25
gs1a 7
*/90/96 gs1a 53
*/96/99 cbc4 40
89/63/- */80/94 gs1a 6
*/87/78 cbc4 11
cbc4 58
Phalacroma sp. CBC4 L128

* /89/98

*Phalacroma* group

ffl 68
Prorocentrum minimum SERC
Prorocentrum micans CCMP1589

0.09