

Metagenomics Workshop

Overview/Discussion

Lab Meeting

Aug 10, 2016

What was EDAMAME 2016?

Explorations in **D**ata **A**nalyses for
Metagenomic **A**dvances in **M**icrobial **E**cology

Michigan State University,
Kellogg Biological Station
July 10-20, 2016

Learning Goals

Ashley Shade edited this page 20 days ago · 7 revisions

Tutorials Organized by Learning Goals

- Computing literacy
 - [Shell](#)
 - [tmux for remote sessions](#)
 - [Getting started with GitHub](#)
 - [Computing workflows for biologists - paper](#)
- Cloud computing
 - [Amazon EC2 start-up](#)
 - [File transfer to the EC2](#)
 - [Getting started with the EC2 - from Angus](#)
- Microbial amplicon analysis
 - [Assessing sequencing quality with FastQC \(bonus: intro to automation, installing software on the EC2\)](#)
 - [Subsampling a large amplicon dataset for developing an analysis workflow](#)
 - [Firing up the QIIME AMI](#)
 - [QIIME workflow overview](#)
 - [QIIME tutorial](#)
 - [mothur workflow](#)
- Microbial shotgun metagenome analysis
 - [Examples of installing mg tools on the EC2](#)
 - [Demo: using seqtk for subsampling a large metagenome dataset for developing an analysis workflow](#)
 - [metaG sequencing preprocessing, quality control, and trimming](#)
 - [Digital normalization](#)
 - [Assembly with MEGAHit](#)
 - [Evaluating assembly](#)
 - [Estimating abundance from metagenomes](#)
 - [Binning assemblies](#)
 - [Annotation of assembled reads](#)
 - [Xander for targeted gene assembly](#)
- Ecological Statistics with R
 - [R basics](#)
 - [Visualizations Demo](#)
 - [Quick intro to R for comparative \(beta\) diversity](#)
- Using Databases
 - [Local BLAST](#)
 - [Using APIs to access NCBI and MG-RAST data](#)
 - [Getting data from NCBI](#)

Major Workshop Themes

Microbial ecology concepts

Metagenomic concepts & tools

Data management & sharing tools

Major Workshop Themes

Microbial ecology concepts

Metagenomics concepts & tools

Data management & sharing tools

Microbial Ecology Concepts

OTUs = **O**perational **T**axonomic **U**nits

Common measures:

richness – number of OTUs present

evenness – abundance of OTUs

composition – taxonomic assignment of OTUs

Microbial Ecology Concepts (cont.)

Diversity measures:

within sample/location - includes richness & evenness, aka alpha diversity

between samples/location – aka beta diversity

regional - includes α & β , aka gamma diversity

See A. Shade's preprint: *Diversity is the question, not the answer*
<https://peerj.com/preprints/2287/>

Vocabulary

OTU

reads (paired vs single, raw vs assembled)

sequencing coverage

reference database

PCR/amplification

Major Workshop Themes

Microbial ecology concepts

Metagenomic concepts & tools

Data management & sharing tools

Activity Intro:

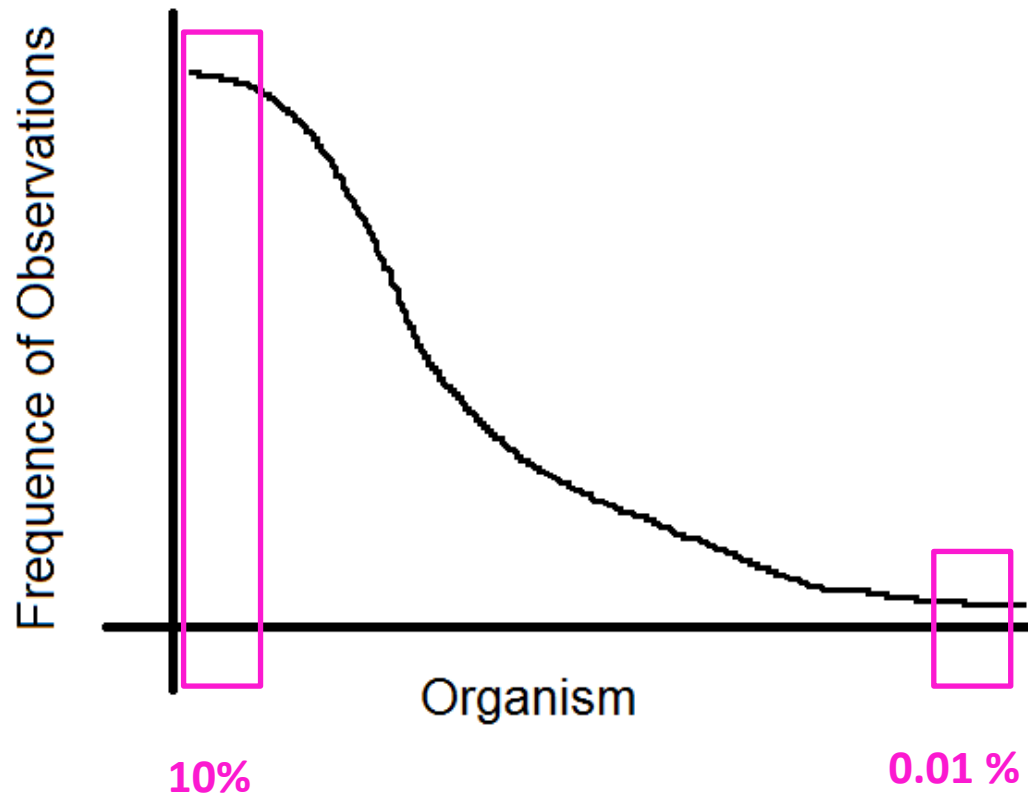
Before sequencing, you can do some back of the envelope calculations.

These can help you:

- Get the sequencing info you need
- Save money
- Tailor your experimental design

Activity:

Cornell Genomics Specs: <http://www.biotech.cornell.edu/brc/genomics-facility/services/next-generation-sequencing>



Activity:

Cornell Genomics Specs: <http://www.biotech.cornell.edu/brc/genomics-facility/services/next-generation-sequencing>

You decide to do some shotgun sequencing at the Cornell Genomics Facility using their HiSeq 2500 instrument (on “Rapid Run” mode), which has an output of 35 Gbp per sequencing lane. Assuming the average genome size of an organism you’re looking for is 5 Mbp and you want to make sure you see it at least 50 times, what percentage cutoff of the community would you be able to survey?

Try these steps:

1. Multiply 5 Mbp x 50 – this is the number of base pairs you’ll be looking for
2. Convert #1 and the HiSeq output to bp to make things easier
3. Divide #1 (in bp) by the output per lane (in bp) - this multiplied by 100 is the percentage cutoff

Now try these with your partner:

1. If you wanted to see the organism 100 times, how would this change the percentage cutoff with one lane?
2. What would you do to survey percentage cutoff equal to that of the practice problem? that you did in the practice question?
3. How would you adjust it to see a rare organism (0.01% of the community) 50 times?

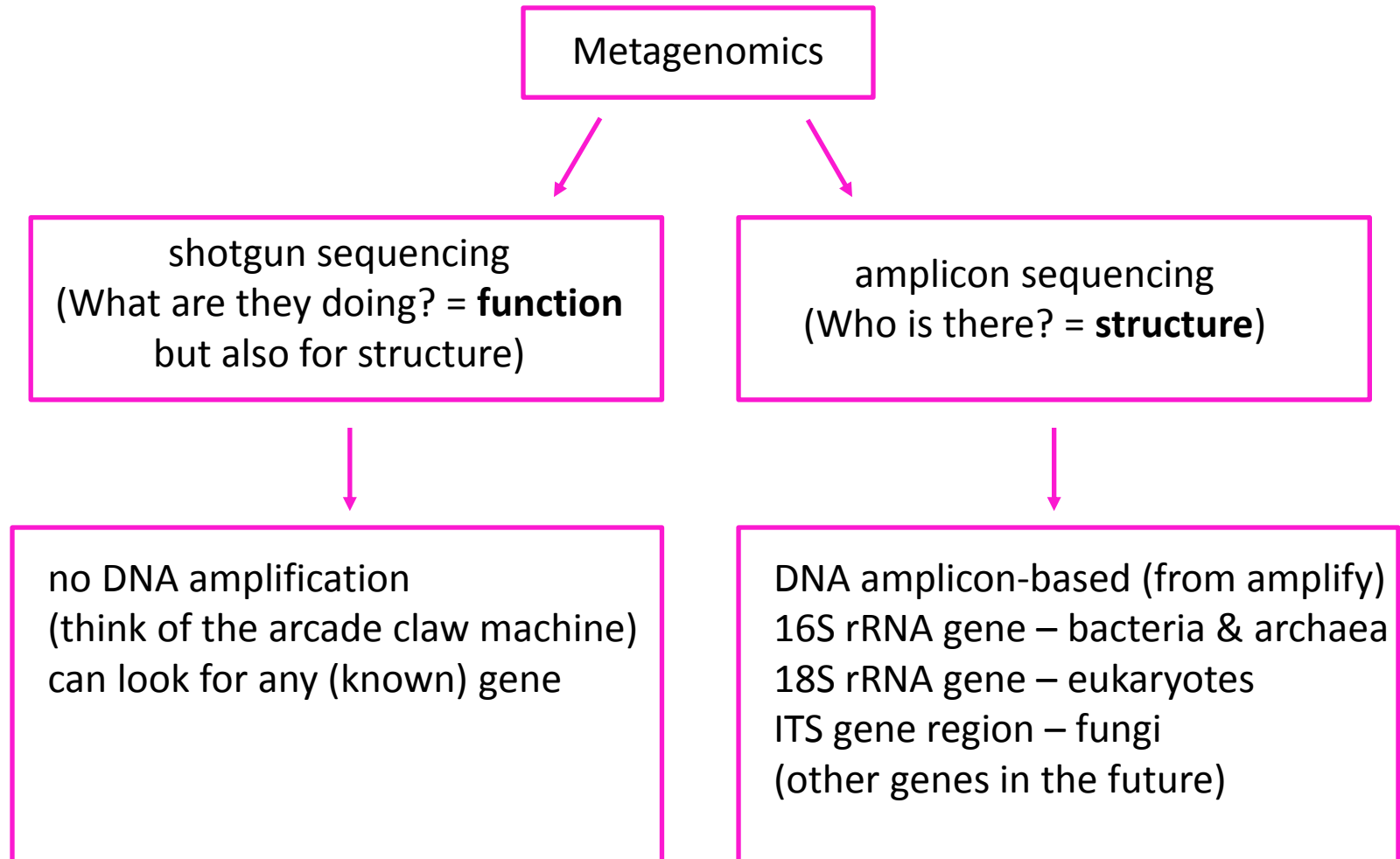
We might ask...

How does microbial community **structure** and/or **function** affect some sort of **phenotype** or change along a gradient (space, time, perturbation, environment)?

Example:

Is P availability correlated with patterns in ppk gene abundance and diversity across a soil moisture gradient?

Intro to Metagenomics



phenotype – How would we measure this?

Intro to Metagenomics

Metagenomics

```
graph TD; A[Metagenomics] --> B[shotgun sequencing<br/>(What are they doing? = function<br/>but also for structure)]; A --> C[amplicon sequencing<br/>(Who is there? = structure)]; B --> D[outputs:<br/>gene table – gene counts/“species”<br/>annotation table – link gene to taxa<br/>metadata – e.g. soil moisture, etc.]; C --> E[outputs we’ll need:<br/>OTU table – OTU counts/“species”<br/>annotation table – link OTU to taxa<br/>metadata – e.g. soil moisture, etc.];
```

2

shotgun sequencing
(What are they doing? = **function**
but also for structure)

outputs:
gene table – gene counts/“species”
annotation table – link gene to taxa
metadata – e.g. soil moisture, etc.

1

amplicon sequencing
(Who is there? = **structure**)

outputs we’ll need:
OTU table – OTU counts/“species”
annotation table – link OTU to taxa
metadata – e.g. soil moisture, etc.

Amplicon Sequencing

Who's there?/Structure

https://github.com/SchlossLab/MiSeq_WetLab_SOP/blob/master/MiSeq_WetLab_SOP_v4.md

Qiime tutorial: https://github.com/edamame-course/Amplicon_Analysis/blob/master/final/2016-07-13-QIIME1.md

mothur tutorial: http://www.mothur.org/wiki/MiSeq_SOP

Two popular programs:
Qiime & mothur

Main points:

- Replication (>5 biological reps, 9 is best)
- Include mock communities (to report errors)
- Use paired-end reads (to reduce errors)
- MiSeq runs are inexpensive (compared to HiSeq)

Shotgun Sequencing

What are they doing?/Function

https://github.com/edamame-course/Metagenome/blob/master/edamame_metagenomics_overview.pdf?raw=true

Many programs

Same main points (from slide 19) apply

Computational developments make it possible

Shotgun Sequencing (cont.)

What are they doing?/Function

Xander tutorial: <https://github.com/edamame-course/Xander/blob/master/Xander.md>

Xander

- assembly based on a functional gene of interest
- download files from FunGene
(<http://fungene.cme.msu.edu/>) or work with RDP
to make your own gene repository

Shotgun Sequencing (cont.)

What are they
Xander tutori:

Biogeochemical cycles

gene—contributor
amoA_AOA—Feifei Liu
amoA_AOB—RDP
buk—RDP
but—RDP
cbh1—Cheryl Kuske
chb—Fan Yang
cooS—Fan Yang
cydA—Rachel Morris
dsrA—Alexander Loy/Michael Wagner
dsrB—Alexander Loy/Michael Wagner
exc1—Fan Yang
fixN—Rachel Morris
glx—Qichao Tu
hydA—Fan Yang
lcc_ascomycetes—Chris Wright
lcc_basidiomycetes—Chris Wright
ligE—Ryan Penton
lip—Qichao Tu
mcrA—Blaz Stres
mmoX—Qichao Tu
mnp—Qichao Tu
nag3—Fan Yang
napA—Laurent Philippot
narG—Laurent Philippot
nifD—RDP
nifH—RDP
nirA—RDP
nirB—RDP
nirK—Tracy Teal
nirS—Veronica Gruntzig
norB—Gesche Braker
nosZ—Blaz Stres
nosZ_atypical_1—Robert Sanford
nosZ_atypical_2—Robert Sanford

ileS—Scott Santos/Howard Ochman
lepA—Scott Santos/Howard Ochman
leuS—Scott Santos/Howard Ochman
pyrG—Scott Santos/Howard Ochman
recA—Scott Santos/Howard Ochman
recG—Scott Santos/Howard Ochman
rplB—Scott Santos/Howard Ochman
rpoB—Scott Santos/Howard Ochman

Biodegradation

gene—contributor
alkB—Gerben Zylstra/Elyse Rodgers-Vieira
benA—Stephan Gantner
bph—Gerben Zylstra
bphA1—Stephan Gantner
bphA2—Stephan Gantner
BSH—Robert Stedtfeld
carA—Shoko Iwai
cntA—Robert Stedtfeld
cutC—Robert Stedtfeld
dbfA1—Shoko Iwai
dxnA—Shoko Iwai
dxnA-dbfa1—Tim Johnson
HSDH—Robert Stedtfeld
npah—Gerben Zylstra
p450—Gerben Zylstra/Elyse Rodgers-Vieira
ppah—Gerben Zylstra
PSA—Robert Stedtfeld

Metal Cycling

gene—contributor
arsA—PFAM
arsB—PFAM
arsC—PFAM
arsD—PFAM

Sequencing Data

Comes to you as a .fasta or .fastq file

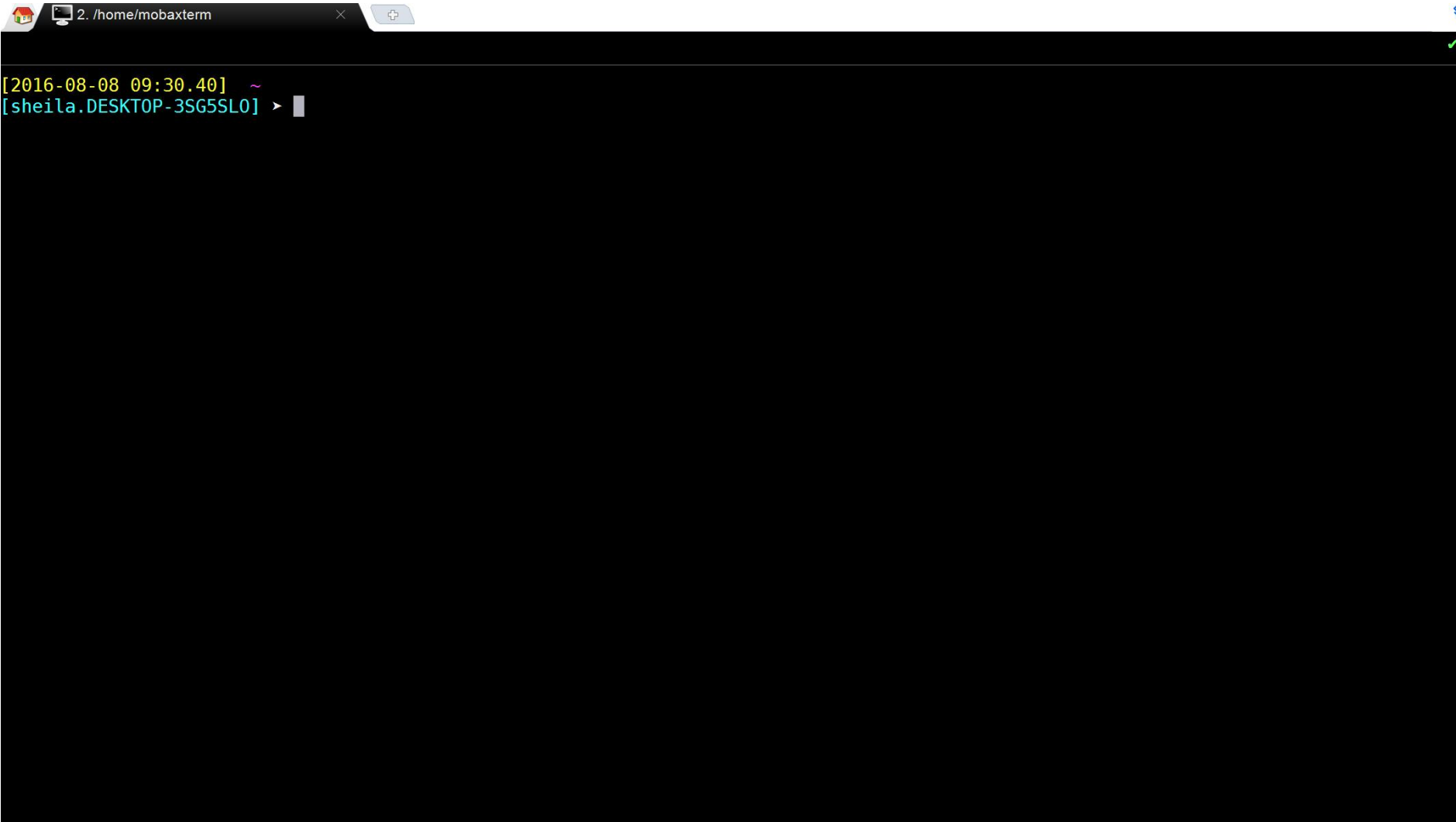
(4 lines per read, millions+ of reads per file)

```
@SMS01_R1
```

```
CCCTTCTTGTCTTCAGCGTTTCTCC
```

```
+
```

```
::3,,,,,,,,,7,,,,,88
```

A screenshot of a MobaXterm terminal window. The title bar at the top shows a small house icon, a monitor icon, and the text "2. /home/mobaxterm". The terminal area has a black background. The first line of text is "[2016-08-08 09:30.40] ~" in yellow. The second line is "[sheila.DESKTOP-3SG5SL0] > " in cyan, followed by a white cursor bar.

```
[2016-08-08 09:30.40] ~  
[sheila.DESKTOP-3SG5SL0] > 
```

Amazon Cloud Computing

More info: <https://aws.amazon.com>

The screenshot displays the AWS Management Console interface for the 'Launch Instance Wizard'. The browser address bar shows the URL: <https://console.aws.amazon.com/ec2/v2/home?region=us-east-1#LaunchInstanceWizard>. The console header includes the AWS logo, navigation menus for 'Services' and 'Edit', and user information for 'Sheila Saia' in the 'N. Virginia' region. The wizard progress bar indicates the current step is '1. Choose AMI'. The main content area lists several AMIs, with 'Ubuntu Server 14.04 LTS (HVM), SSD Volume Type' highlighted by a pink rectangular box. Other visible AMIs include Amazon Linux, Red Hat Enterprise Linux, and SUSE Linux. Each AMI entry shows its icon, name, description, root device type, virtualization type, and a 'Select' button. A 'Cancel and Exit' link is located at the top right of the wizard steps.

AMI Name	AMI ID	Root Device Type	Virtualization Type	Architecture
Amazon Linux AMI 2016.03.3 (HVM), SSD Volume Type	ami-6869aa05	ebs	hvm	64-bit
Red Hat Enterprise Linux 7.2 (HVM), SSD Volume Type	ami-2051294a	ebs	hvm	64-bit
SUSE Linux Enterprise Server 12 SP1 (HVM), SSD Volume Type	ami-b7b4fedd	ebs	hvm	64-bit
Ubuntu Server 14.04 LTS (HVM), SSD Volume Type	ami-2d39803a	ebs	hvm	64-bit
Microsoft Windows Server 2012 R2 Base	ami-74a73263	ebs	hvm	64-bit

Major Workshop Themes

Microbial ecology concepts

Metagenomic concepts & tools

Data management & sharing tools

Paper Discussion

Link: <http://journals.plos.org/plosbiology/article?id=10.1371/journal.pbio.1002303>

1. Take-away messages for hydrologists?
2. Are you using an reproducible workflows?
What works and what doesn't?
3. Anything else you agree/disagree with?

Specific Data Mgmt. Tools

Reproducible Research

- GitHub
- R/RStudio (scripting)
- (markdown, Jupyter notebooks)

Version Control

- Git

Info Sharing/Open Source

- GitHub Repositories
- Mendeley paper groups
- Etherpad

GitHub & Git

help guide: <https://guides.github.com/activities/hello-world/>

The screenshot shows the GitHub repository page for 'JoFlow' by 'SoilWaterLab'. The repository name 'JoFlow' is highlighted with a pink box. The repository description is 'JoFlow: Simple, regionally parameterized lumped watershed model — Edit'. The repository statistics show 3 commits, 2 branches (highlighted with a pink box), 0 releases, and 1 contributor. The 'Clone or download' button is highlighted with a pink box. The repository files list includes 'Lumped_VSA_model' and 'README.md'. The README content is visible below the files list.

GitHub, Inc. (US) | <https://github.com/SoilWaterLab/JoFlow> | Search

This repository | Search | Pull requests | Issues | Gist

SoilWaterLab / **JoFlow** | Unwatch 6 | Star 0 | **Fork 0**

Code | Issues 0 | Pull requests 0 | Wiki | Pulse | Graphs | Settings

JoFlow: Simple, regionally parameterized lumped watershed model — Edit

3 commits | **2 branches** | 0 releases | 1 contributor

Branch: master | New pull request | Create new file | Upload files | Find file | **Clone or download**

chelseakmorris Update README.md | Latest commit c825a92 on May 10

File	Commit	Time
Lumped_VSA_model	Create Lumped_VSA_model	3 months ago
README.md	Update README.md	3 months ago

README.md

JoFlow

Simple, regionally parameterized lumped watershed model originally written by Josephine Archibald. For reference see J.A. Archibald, B.P. Buchanan, D.R. Fuka, C.B. Georgakakos, S.W. Lyon, M.T. Walter, A simple, regionally parameterized model for predicting nonpoint source areas in the northeastern US, Journal of Hydrology: Regional Studies, Volume 1, July 2014, Pages 74-91, ISSN 2214-5818, <http://dx.doi.org/10.1016/j.ejrh.2014.06.003>. Model code is located on master branch.

© 2016 GitHub, Inc. | Terms | Privacy | Security | Status | Help | Contact GitHub | API | Training | Shop | Blog | About

GitHub & Git

help guide: <https://guides.github.com/activities/hello-world/>

The screenshot shows the GitHub web interface for the repository `SoilWaterLab / JoFlow`. The browser address bar displays the URL `https://github.com/SoilWaterLab/JoFlow/branches`. The repository page includes a header with the repository name, a search bar, and navigation links for Pull requests, Issues, and Gist. Below the header, the repository's default branch is `master`, updated 3 months ago by `chelseakmorris`. The active branches section lists the `Isotope` branch, also updated 3 months ago by `chelseakmorris`. A pink box highlights the `Isotope` branch name, and another pink box highlights the `New pull request` button next to it. The footer of the page shows the GitHub logo and links for Contact GitHub, API, Training, Shop, Blog, and About.

Command Line GitHub & Git

workshop guide: <https://github.com/edamame-course/Github/blob/master/Tutorial.md>

```
posh~git ~ 2016-tutorials [master]
C:\Users\sheila\Documents\GitHub> ls

Directory: C:\Users\sheila\Documents\GitHub

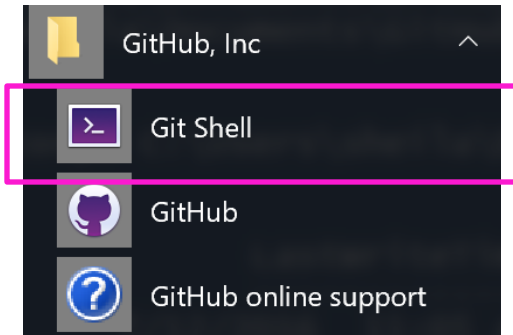
Mode                LastWriteTime         Length Name
----                -
d-----          7/12/2016 11:40 AM             2016-tutorials
d-----          7/14/2016  4:21 PM             Jones_R_Files
d-----          7/18/2016  1:28 PM             PAPER_LeeSorensen_inprep

C:\Users\sheila\Documents\GitHub> cd .\2016-tutorials
C:\Users\sheila\Documents\GitHub\2016-tutorials [master ≡]> ls

Directory: C:\Users\sheila\Documents\GitHub\2016-tutorials

Mode                LastWriteTime         Length Name
----                -
d-----          7/12/2016 11:05 AM             HandoutsResources
d-----          7/12/2016 11:05 AM             images
d-----          7/12/2016 11:05 AM             lectures
d-----          7/12/2016 11:48 AM             test
d-----          7/12/2016 11:05 AM             Xander
-a----          7/12/2016 11:05 AM           5193 databasesExercise.md
-a----          7/12/2016 11:19 AM           199 README.md

C:\Users\sheila\Documents\GitHub\2016-tutorials [master ≡]>
```



> git status

> git pull

> git add

> git commit

> git push

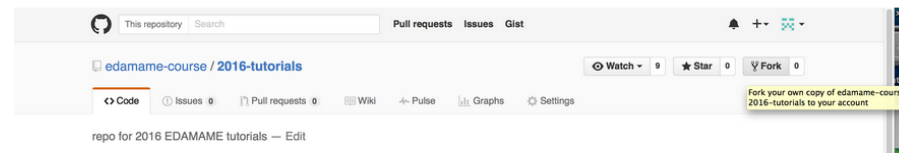
Tutorials Hosted on GitHub

Using markdown (.md)

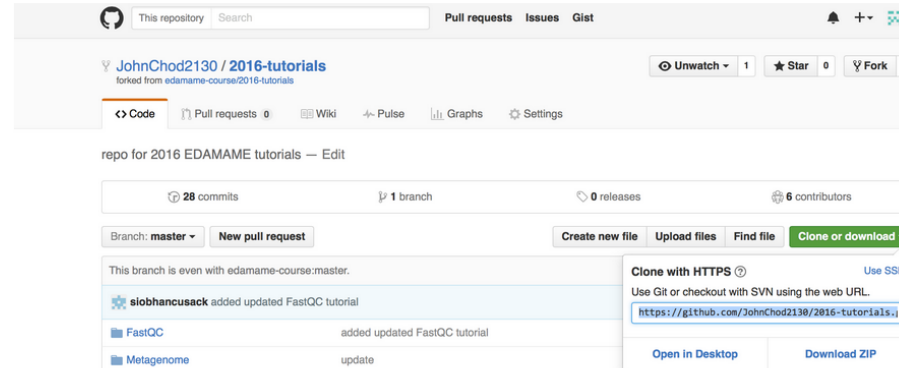
3. Fork and Clone a repository

The main benefit of forking a repository is that it allows you to copy an existing repository and then have the ability to make your own edits to the repository without changing the original repository.

- Navigate to the EDAMAME [2016-tutorials](#) repository.
- On the upper right-hand side there is a box labeled "Fork". Click on that.



- You'll be re-directed to your github account. The repository is now on your github account but we still need to clone the repository so we have local access to the files.
- In the new window, on the right hand side there is a box labeled "Clone or download".



- Click on that and then copy the link.
- Choose a local directory that you want this repository to be added to. Change into that directory and use git clone with the URL just copied:
- **Note that the below command will not work for you because you need to appropriately edit the URL**

```
git clone https://github.com/**YourGitHubName**/2016-tutorials.git
```

- Directions for 'git clone' can also be found at [GitHub](#).
- This protocol can be used to clone any public repository. For EDAMAME repos, you can `pull` to get the most up-to-date materials from GitHub, but you cannot `push` to edit those resources and have your edits tracked to the main repository, because you are not part of the EDAMAME team. (More details on these commands below).

R Markdown

help guide: <http://rmarkdown.rstudio.com/lesson-1.html>

RStudio > File > New File > R Markdown...

The screenshot displays the RStudio interface with a new R Markdown file, 'edamameTest_README.Rmd', open. The left pane shows the source code, and the right pane shows the rendered HTML output. The console at the bottom shows the output of the R code chunks.

Source Code (Left Pane):

```
1 ---
2 title: "Test README"
3 author: "Sheila S"
4 date: "August 6, 2016"
5 output: html_document
6 ---
7
8 {r setup, include=FALSE}
9 knitr::opts_chunk$set(echo = TRUE)
10
11 ## R Markdown
12
13 This is an R Markdown document. Markdown is a simple formatting
14 details on using R Markdown see <http://rmarkdown.rstudio.com>
15
16 When you click the **Knit** button a document will be generated
17 R code chunks within the document. You can embed an R code chunk
18 like this:
19
20 {r cars}
21 summary(cars)
22
23 ## Including Plots
24
25 You can also embed plots, for example:
26
27 {r pressure, echo=FALSE}
28 plot(pressure)
29
30 Note that the `echo = FALSE` parameter was added to the code chunk
31 to prevent printing the R output to the console.
```

Rendered HTML (Right Pane):

Test README

Sheila S
August 6, 2016

R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.


When you click the Knit button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

```
summary(cars)
```

##	speed	dist
## Min.	: 4.0	Min. : 2.00
## 1st Qu.	:12.0	1st Qu.: 26.00
## Median	:15.0	Median : 36.00
## Mean	:15.4	Mean : 42.98
## 3rd Qu.	:19.0	3rd Qu.: 56.00
## Max.	:25.0	Max. :120.00

Including Plots

You can also embed plots, for example:



Content type 'application/zip' length 2221261 bytes (2.1 MB)
downloaded 2.1 MB

package 'base64enc' successfully unpacked and MD5 sums checked
package 'rmarkdown' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
C:\Users\sheila\AppData\Local\Temp\RtmpaQWStc\downloaded_packages

Etherpad

Create a new Etherpad: <https://public.etherpad-mozilla.org/>
Etherpad info: <https://github.com/ether/etherpad-lite#installation>

The screenshot shows a web browser window with the address bar displaying <https://public.etherpad-mozilla.org/p/EDAMAME-2016>. The Etherpad editor interface includes a top toolbar with text formatting options (Bold, Italic, Underline, Strikethrough), list creation, undo/redo, and a search bar. The main editing area contains the following text:

```
1 Etherpad for EDAMAME 2016
2
3 http://www.edamamecourse.org
4
5 Schedule
6 https://github.com/edamame-course/2016-tutorials/wiki/Schedule-EDAMAME-2016
7
8
9 Day 1 7/11/16
10
11 FOR PC USERS:
12   Download Mobaxterm here: http://mobaxterm.mobatek.net/MobaXterm\_Setup\_9.1.msi
13
14 ATM location: https://goo.gl/maps/qVqz1z3qTyG2
15
16 hi EVERYONE EDAMAME ROCKS! - Everyone, feel free to join the note taking
17
18 why we are here: microbes running the show. we HAVE to learn and UNDERSTAND them
19 *response from participant*
20 -antagonistic interaction: plant vs pathogen,
21 -symbiotic: gut microbiome, environmental effect on ecosystem fxn
22
23 traits of microbial communities:
24 -communities consists of fraction of live microbes vs dna from dead? so, how we can distinguish that?
25
26 *deep coverage: amt of data per samples
27
28 Make a list of questions
29
30 What are your questions for understanding a system?
31
32 THE AWESOME TABLE:
33 Déjio: How diet affects microbiome and disease propensity in tadpoles
34 Danielle: How social interactions shape individual preen gland microbiome in dark-eyed juncos (birds); especially in mated pairs
35 Betsy: Plant-microbe interaction – plants that fix their own nitrogen have a higher requirement for phosphate in the soil. How do we identify community members which solubilize
36 Phos. le: assigning identity to function without culturing individual organism.
37 Liz: cheeses – can you use any aspect of the metadata/sample type to predict microbial community. Interactions between microbial taxa
```

On the right side, there is a chat window titled "Chat" with a "Enter your name" input field. The chat history shows the following messages:

- disappeared)
- Taylor: do you see the classroom? 10:53
- BE@ND-Mauna: No, Pat's face 10:54
- Siobhan: how about now? 10:55
- BE@ND-Chissals: There we go! 10:55
- BE@ND-Chissals: thanks! 10:55
- BE@ND-Mauna: Thanks guys! 10:55
- BE@ND-Chissals: Question for Patrick: Sorry if this is a silly question, but what are some ways you can estimate the resources your job/program will use? (Besides reading the documentation - like if it's something you wrote yourself?) 11:25
- BE@ND-Chissals: thx for asking! 11:27

Jupyter Notebooks

Installation & help: <http://jupyter.readthedocs.io/en/latest/install.html>

PAOsoil_intialcalcs

This script removes unnecessary columns from the original data (i.e. "Notes") and also removes NA values. It also runs initial calculations, joins dataframes, and exports these new dataframes for other PAOsoil applications.

Contact: Sheila Sala (sms493@cornell.edu)

Last Updated: June 8, 2016

Loading libraries

```
In [39]: library(dplyr) #database mgmt
```

Setting directory and loading in files

```
In [40]: # Clear workspace
rm(list = ls())

# Need to set directory to location of the USDA_SoilDB text files
setwd('C:/Users/Sheila/Documents/Cornell/Research/PhD/soil_sampling/USDA_SoilDB/USDA_SoilDB_TextFiles')

# Import files
fileList=read.table('fileList.txt',header=TRUE,sep="\t",fill=TRUE)
for (i in 1:nrow(fileList))
  assign(as.character(fileList[i,1]),read.table(as.character(fileList[i,2]),header=TRUE,sep="\t",fill=TRUE),envir=.GlobalEnv)

# Check to see if the text files loaded
ls()

"cacI2Results" "fe2Results" "fileList" "I" "loiResults" "oxResults" "pHResults" "sampleEvent" "seqDNAInfo" "siteInfo" "tempResults"
"totalexResults" "vmcResults"
```

Removed "Notes" columns & NA's

```
In [41]: siteInfoSel=select(siteInfo,SiteID:State) %>% na.omit
sampleEventSel=select(sampleEvent,SampleID:Month) %>% na.omit
cacI2ResultsSel=select(cacI2Results,SiteID,SampleID:SRppmCalibRaw) %>% na.omit
fe2ResultsSel=select(fe2Results,SiteID,SampleID:Fe2ppmRawCalib) %>% na.omit
oxResultsSel=select(oxResults,SiteID,SampleID:OxppmCalibRawI8x) %>% na.omit
totalexResultsSel=select(totalexResults,SiteID,SampleID:TotppmCalibRawI20x) %>% na.omit
tempResultsSel=select(tempResults,SiteID,SampleID:TempC) %>% na.omit

# Take only first reps
# for vmc second rep was to check probe for one month
# for pH and Loi was done by accident for one month
vmcResultsSel=vmcResults %>% filter(RepVmc=="1") %>% select(SiteID,SampleID,VMCcalibperc) %>% na.omit
pHResultsSel=pHResults %>% filter(RepPH=="1") %>% select(SiteID,SampleID:pH) %>% na.omit
```


Other Key Themes

Research goals/questions shape your path

Try to focus on hypothesis driven studies

Planning ahead

No need to reinvent the wheel

(wrt. ecological tools & computer programs)

Use statistical tests to relate who/what with process

Other Resources

EDAMAME Tutorials by Subject

<https://github.com/edamame-course/2016-tutorials/wiki/Learning-Goals>

EDAMAME Schedule (with tutorial links)

<https://github.com/edamame-course/2016-tutorials/wiki/Schedule-EDAMAME-2016>

MiSeq (Amplicon Sequencing) SOP

https://github.com/SchlossLab/MiSeq_WetLab_SOP/blob/master/MiSeq_WetLab_SOP_v4.md

Illumina Sequencing Video (paired reads)

<https://www.youtube.com/watch?v=womKfikWlXM>

de Bruijn graph explanation (used in Xander):

http://www.cs.jhu.edu/~langmea/resources/lecture_notes/assembly_dbg.pdf

Project Templates

http://projecttemplate.net/getting_started.html

R Studio Cheatsheets

<https://www.rstudio.com/resources/cheatsheets/>

Not So Standard Deviation podcast

<https://soundcloud.com/nssd-podcast>

Extras

HiSeq 2500

[Overview](#)
[System](#)
[Applications](#)
[Publications](#)
[Featured Researchers](#)
[Kits](#)
[Specifications](#)
[Workflow](#)
[Literature](#)
[Options & Accessories](#)
[Scientific Data](#)
[Technology](#)
[Software](#)
[Support](#)

Interested in receiving newsletters, case studies, and information on new applications? Enter your email address below.

* First Name:

* Last Name:

* Email:

* Area of Interest:

Select...

* Job Function:

Select...

* Country:

Select...

Sign Up

HiSeq 2500 Specifications

HiSeq System Performance Parameters

High Output Run Mode*

HISEQ SBS V4 SPECIFICATIONS				TRUSEQ SBS V3		
Read length	Dual Flow Cell	Single Flow Cell	Dual Flow Cell Run Time	Dual Flow Cell	Single Flow Cell	Dual Flow Cell Run Time
1×36	128-144 Gb	64-72 Gb	29 hrs	95-105 Gb	47-52 Gb	2 days
2×50	360-400 Gb	180-200 Gb	2.5 days	270-300 Gb	135-150 Gb	5.5 days
2×100	720-800 Gb	360-400 Gb	5 days	540-600 Gb	270-300 Gb	11 days
2×125	900-1 Tb	450-500 Gb	6 days	N/A	N/A	N/A
Reads Passing Filter (8 lanes per flow cell)	Up to 4 billion single read or 8 billion paired-end reads	Up to 2 billion single read or 4 billion paired-end reads		Up to 3 billion single read or 6 billion paired-end reads	Up to 1.5 billion single read or 3 billion paired-end reads	
Quality	Greater than 85% of bases above Q30 at 2×50 bp Greater than 80% of bases above Q30 at 2×100 bp Greater than 80% of bases above Q30 at 2×125 bp			Greater than 85% of bases above Q30 at 2×50 bp Greater than 80% of bases above Q30 at 2×100 bp		

*Install specifications based on Illumina PhiX control library at supported cluster densities (between 610-678 K clusters/mm² passing filter using TruSeq v3 Kits or 870-930 K clusters/mm² passing filter using HiSeq v4). Run times for high output mode correspond to sequencing only. Performance may vary based on sample quality, cluster density, and other experimental factors.

Rapid Run Mode*

HISEQ RAPID SBS KIT V2 SPECIFICATIONS			
Read length	Dual Flow Cell	Single Flow Cell	Dual Flow Cell Run Time
1×36	18-22 Gb	9-11 Gb	7 hr
2×50	50-60 Gb	25-30 Gb	16 hr
2×100	100-120 Gb	50-60 Gb	27 hr
2×150	150-180 Gb	75-90 Gb	40 hr
2×250	250-300 Gb	125-150 Gb	60 hr
Reads Passing Filter (2 lanes per flow cell)	Up to 600 million single read or 1.2 billion paired-end reads	Up to 300 million single read or 600 million paired-end reads	
Quality	Greater than 85% of bases above Q30 at 2×50 bp Greater than 80% of bases above Q30 at 2×100 bp Greater than 75% of bases above Q30 at 2×250 bp		

*Install specifications based on Illumina PhiX control library at supported cluster densities (between 700-820 K clusters/mm² passing filter using

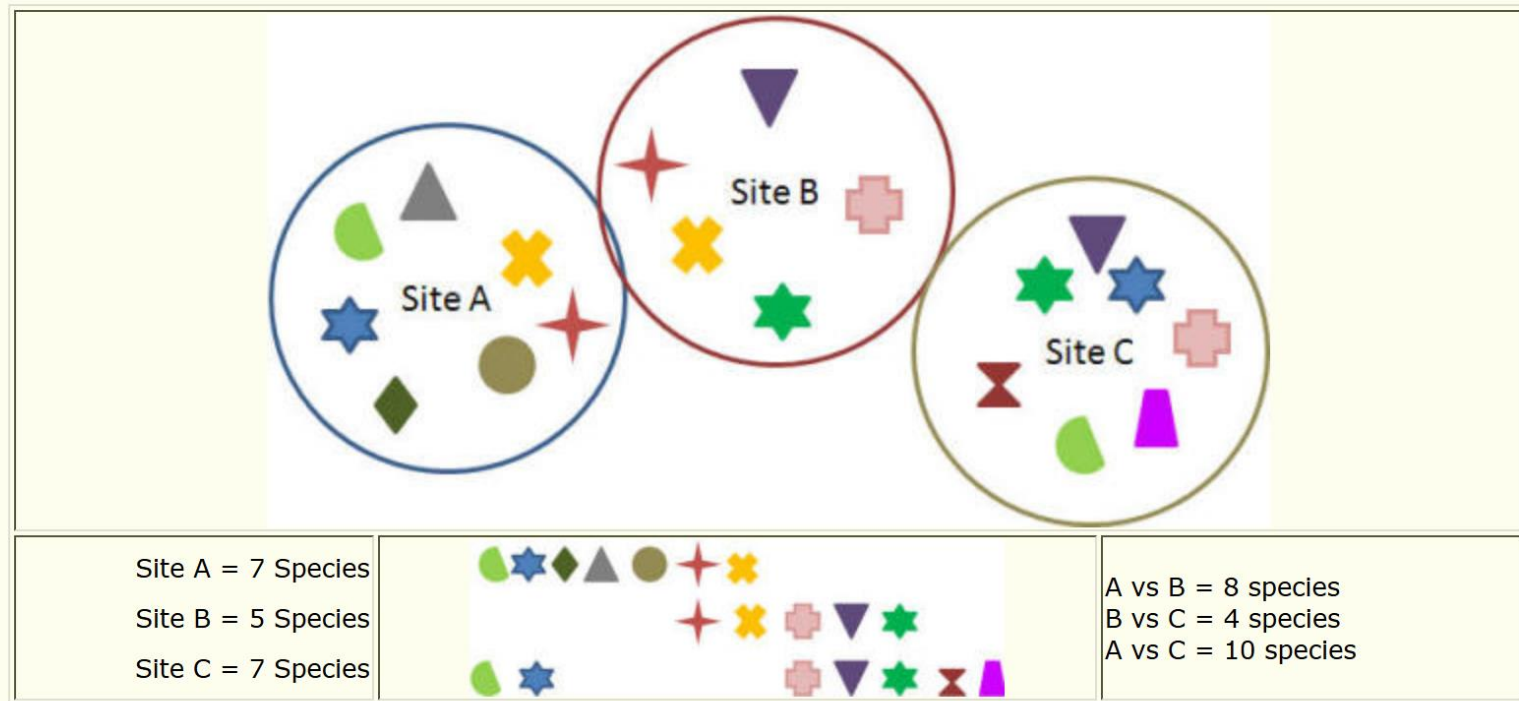
Biodiversity Can be Expressed at Several Scales

Biodiversity can be measured and monitored at several spatial scales.

Alpha Diversity = richness and evenness of individuals within a habitat unit. For example in the figure below, **Alpha Diversity** of Site A = 7 species, Site B = 5 species, Site C = 7 species.

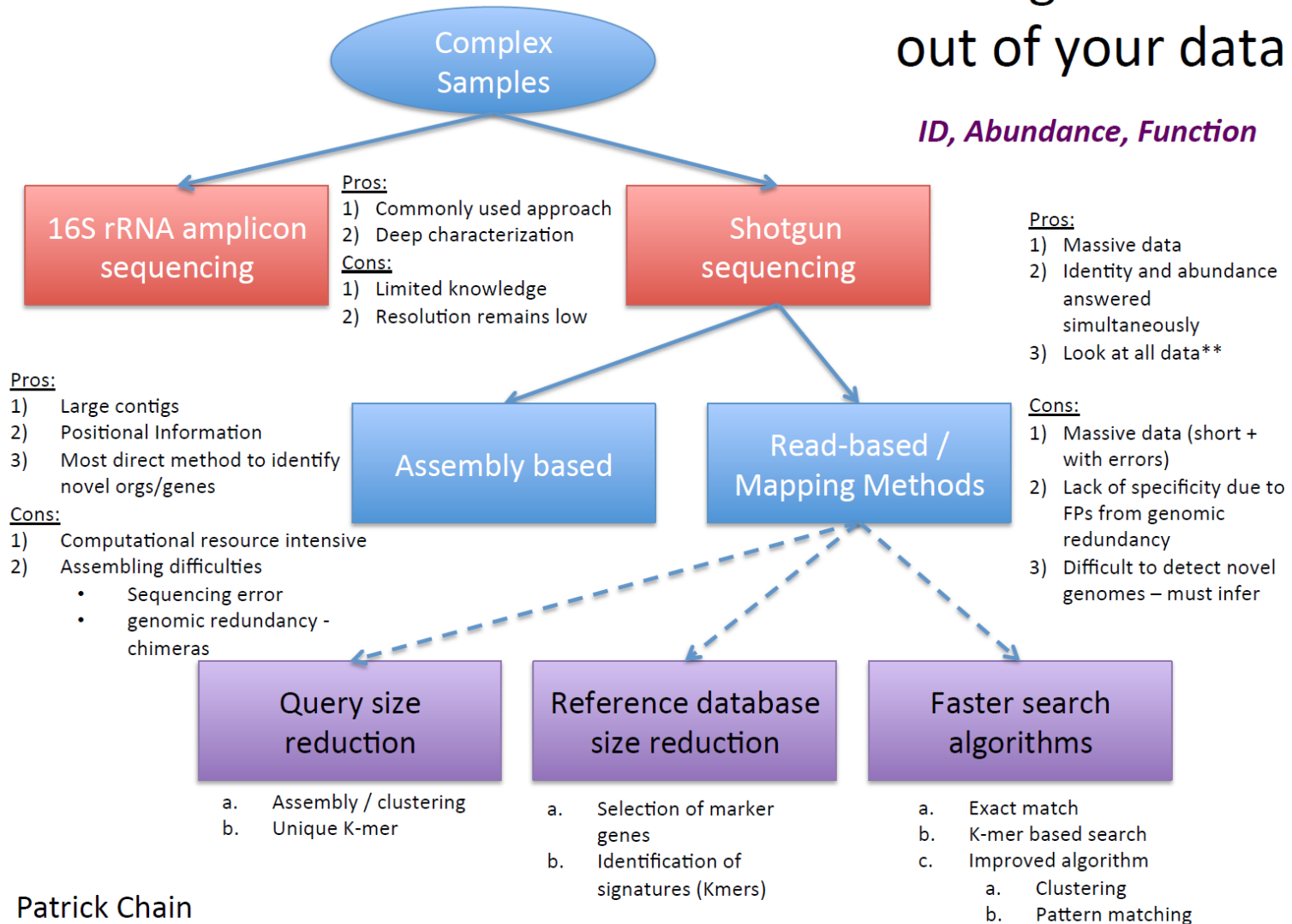
Beta Diversity = expression of diversity between habitats. In the example below, the greatest **Beta Diversity** is observed between Site A and C with 10 species that differ between them and only 2 species in common.

Gamma Diversity = landscape diversity or diversity of habitats within a landscape or region. In this example, the gamma diversity is 3 habitats with 12 species total diversity.



Getting the most out of your data

ID, Abundance, Function



Patrick Chain