



Quantitative Insights Into Microbial Ecology



#### PICRUSt

- Predict metagenome functional content from marker gene

# Workshop 11: Metagenomics Analysis

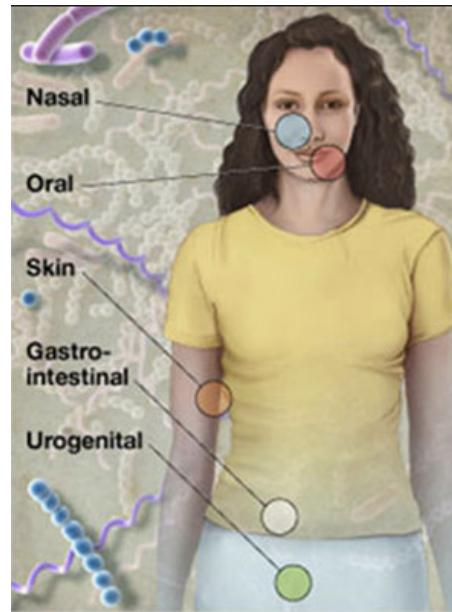
*Shi, Baochen  
Department of Pharmacology, UCLA*

# The Microbiome/Microbiota

---



# The Human Microbiome



<http://hmpdacc.org/>

HMP 2008-2013, funded by NIH

- characterized the microbial communities at several different sites on the human body
  - Initial 16S & mWGS metagenomic studies to generate an estimate of the complexity of the microbial community at each body site, providing initial answers to the questions of whether there is a "core" microbiome at each site
  - Demonstration projects to determine the relationship between disease and changes in the human microbiome
- 
- Development of a reference set of 3,000 isolate microbial genome sequences
  - Development of new tools and technologies for computational analysis
  - establishment of a data analysis and coordinating center (DACC), and resource repositories

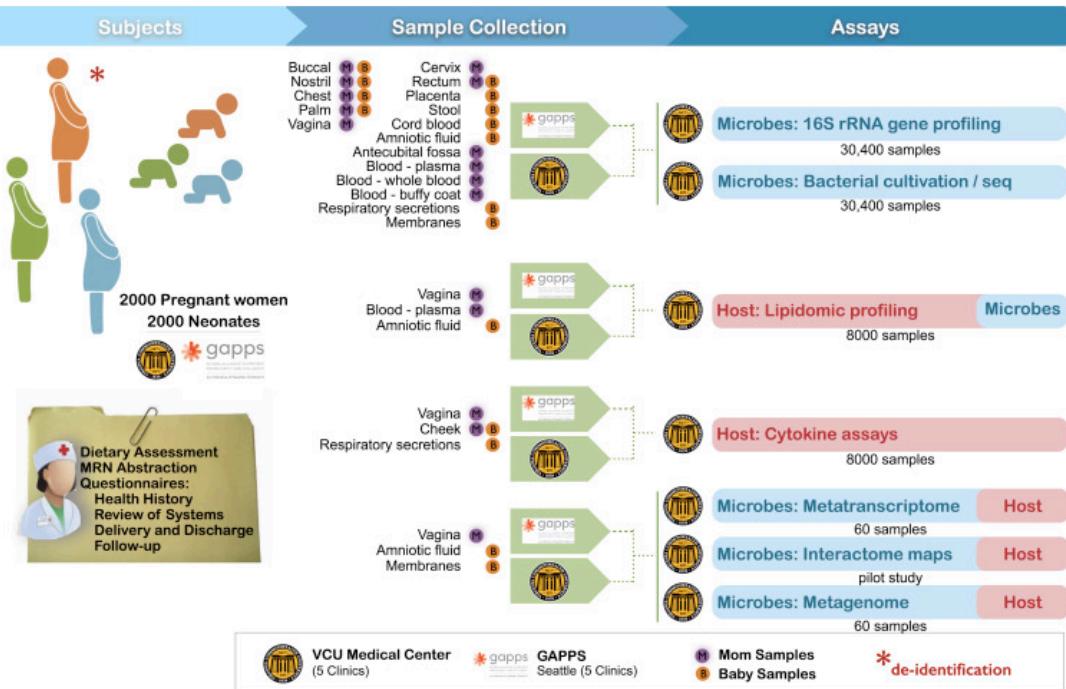
# The Human Microbiome



*Cell Host Microbe. 2014 Sep 10;16(3):276-89*

iHMP creates integrated longitudinal datasets of biological properties from 3 different cohort studies of microbiome-associated conditions using multiple "omics" technologies.

## Project 1: Pregnancy & Preterm Birth



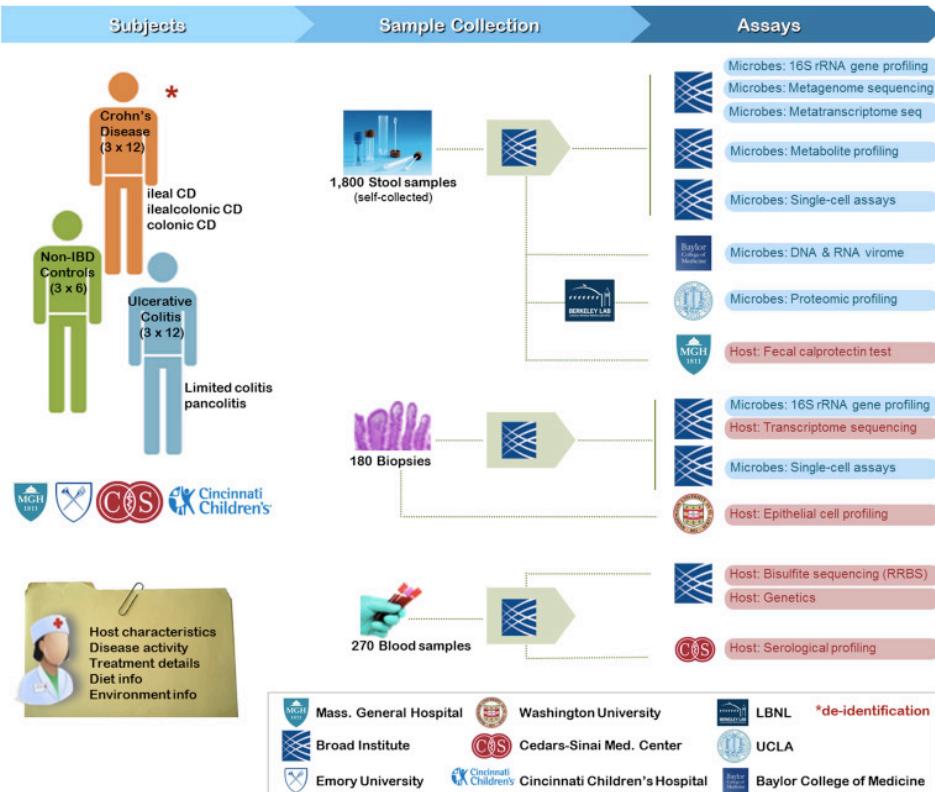
# The Human Microbiome



*Cell Host Microbe. 2014 Sep 10;16(3):276-89*

iHMP creates integrated longitudinal datasets of biological properties from 3 different cohort studies of microbiome-associated conditions using multiple "omics" technologies.

## Project 2: Onset of Inflammatory Bowel Disease (IBD)



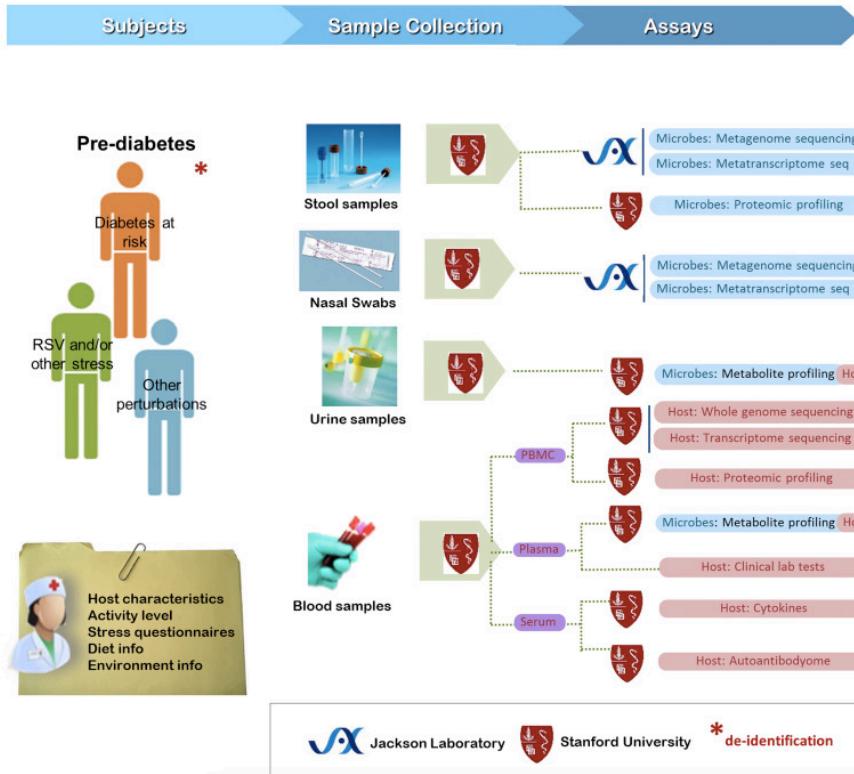
# The Human Microbiome



*Cell Host Microbe. 2014 Sep 10;16(3):276-89*

iHMP creates integrated longitudinal datasets of biological properties from 3 different cohort studies of microbiome-associated conditions using multiple "omics" technologies.

## Project 3: Onset of Type 2 Diabetes



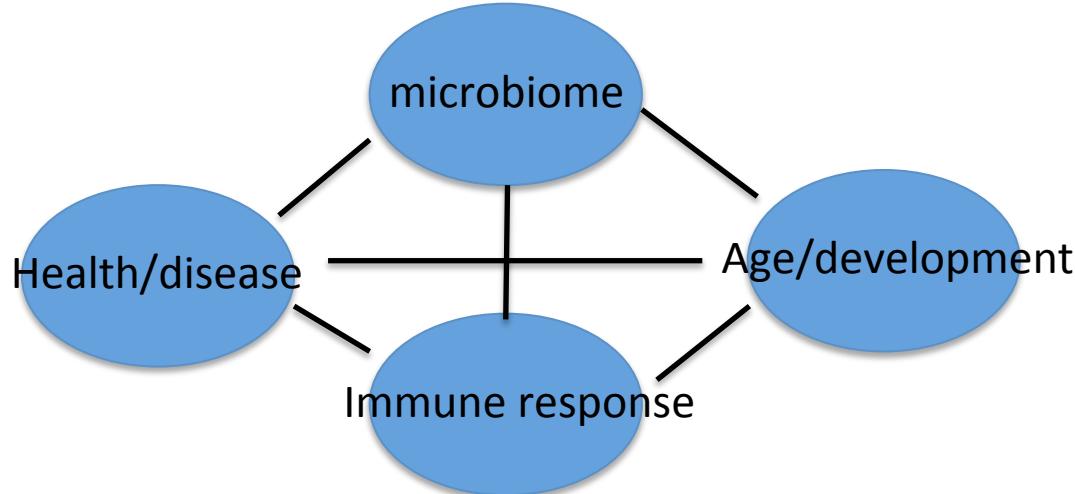
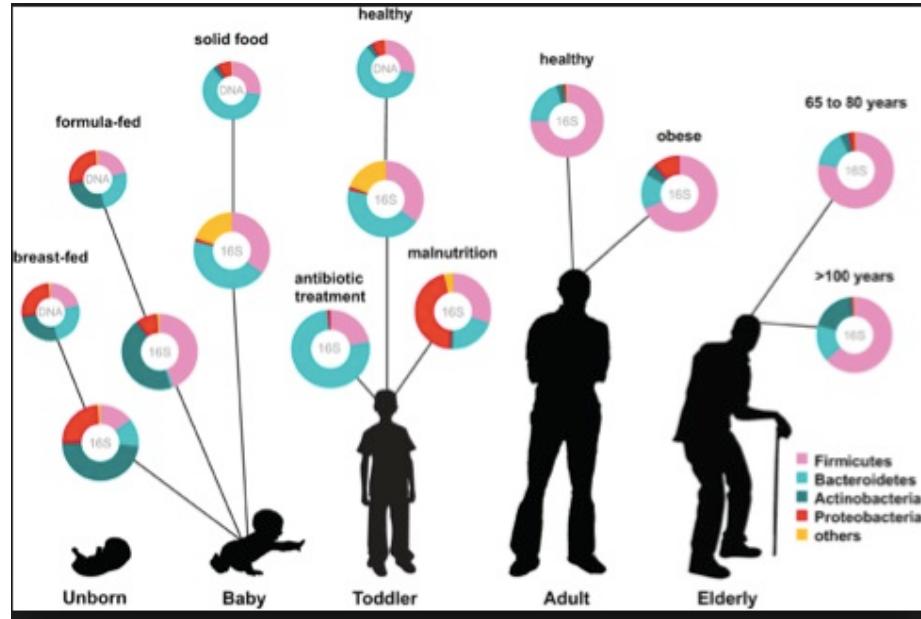
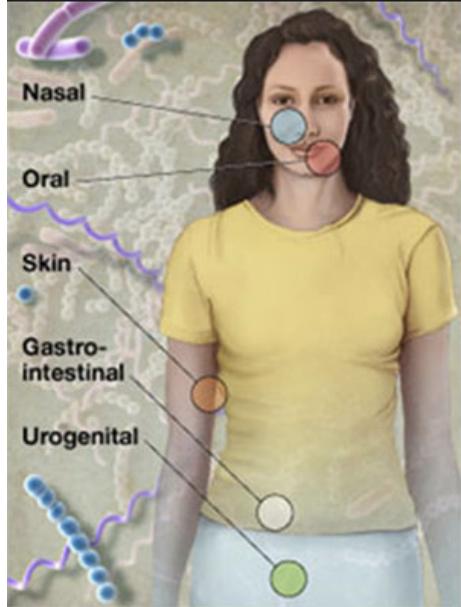
# The Human Microbiome

---

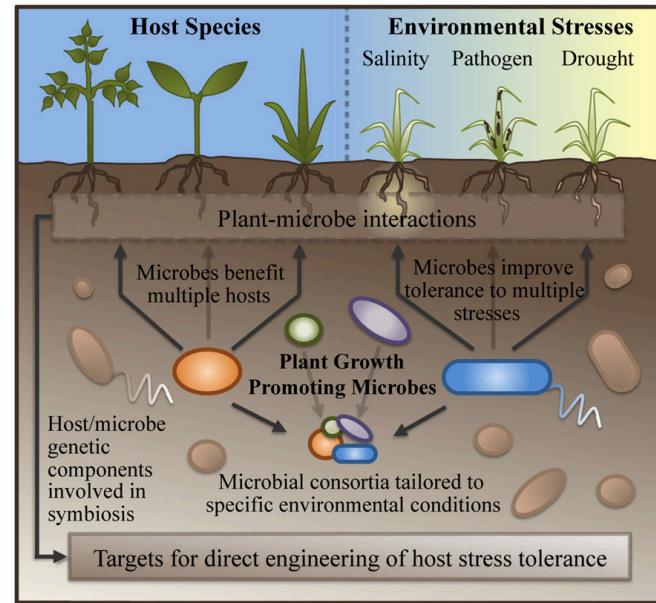


- The human gut microbiome in health and disease
- mWGS

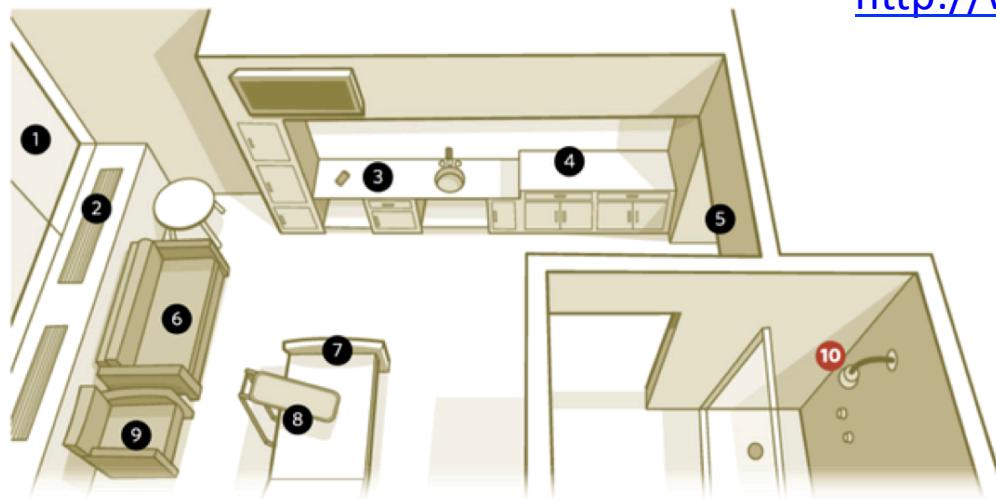
# The Human Microbiome



# The Earth/Environmental Microbiome



<http://www.earthmicrobiome.org/>





Quantitative Insights Into Microbial Ecology



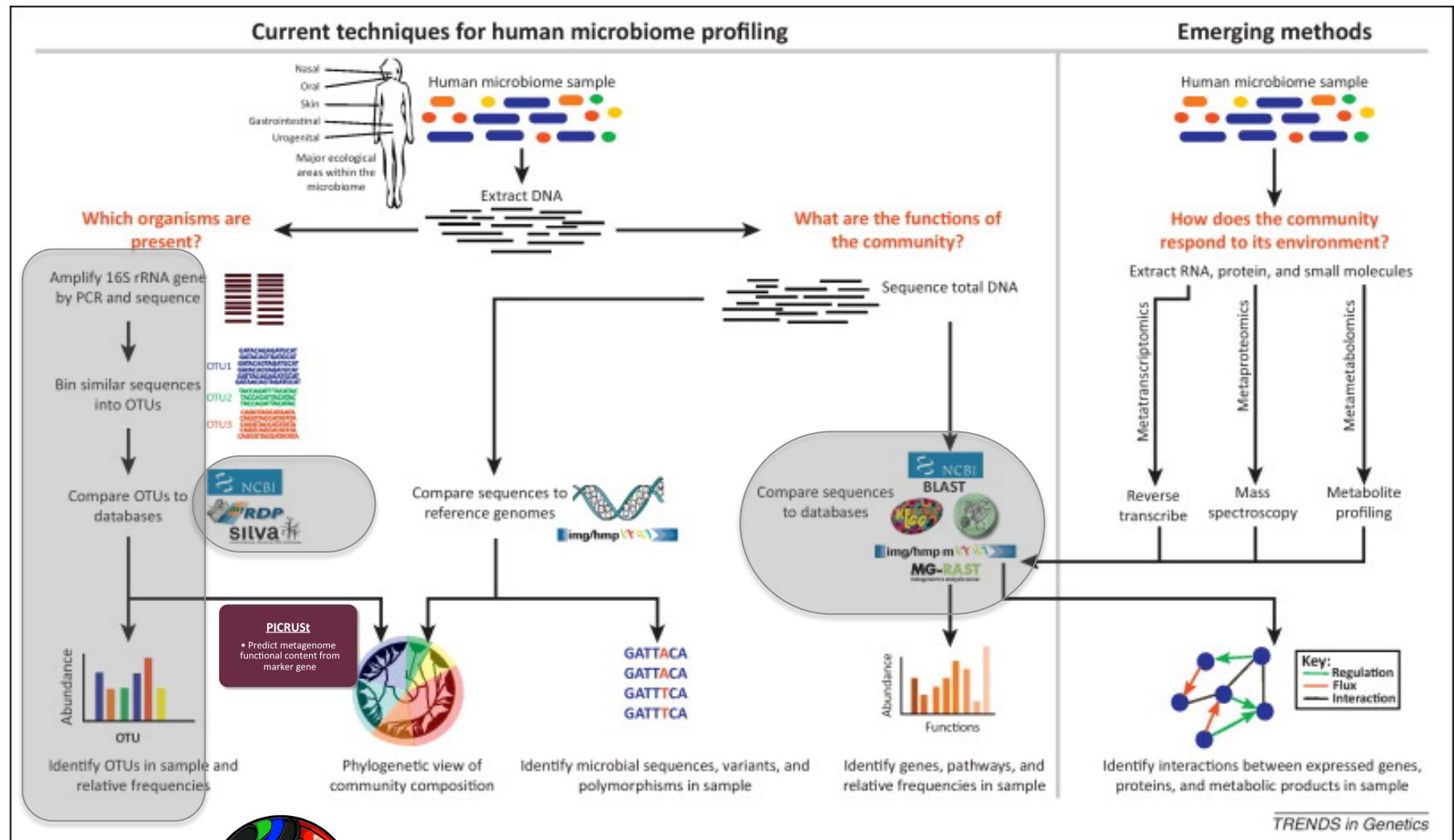
#### PICRUSt

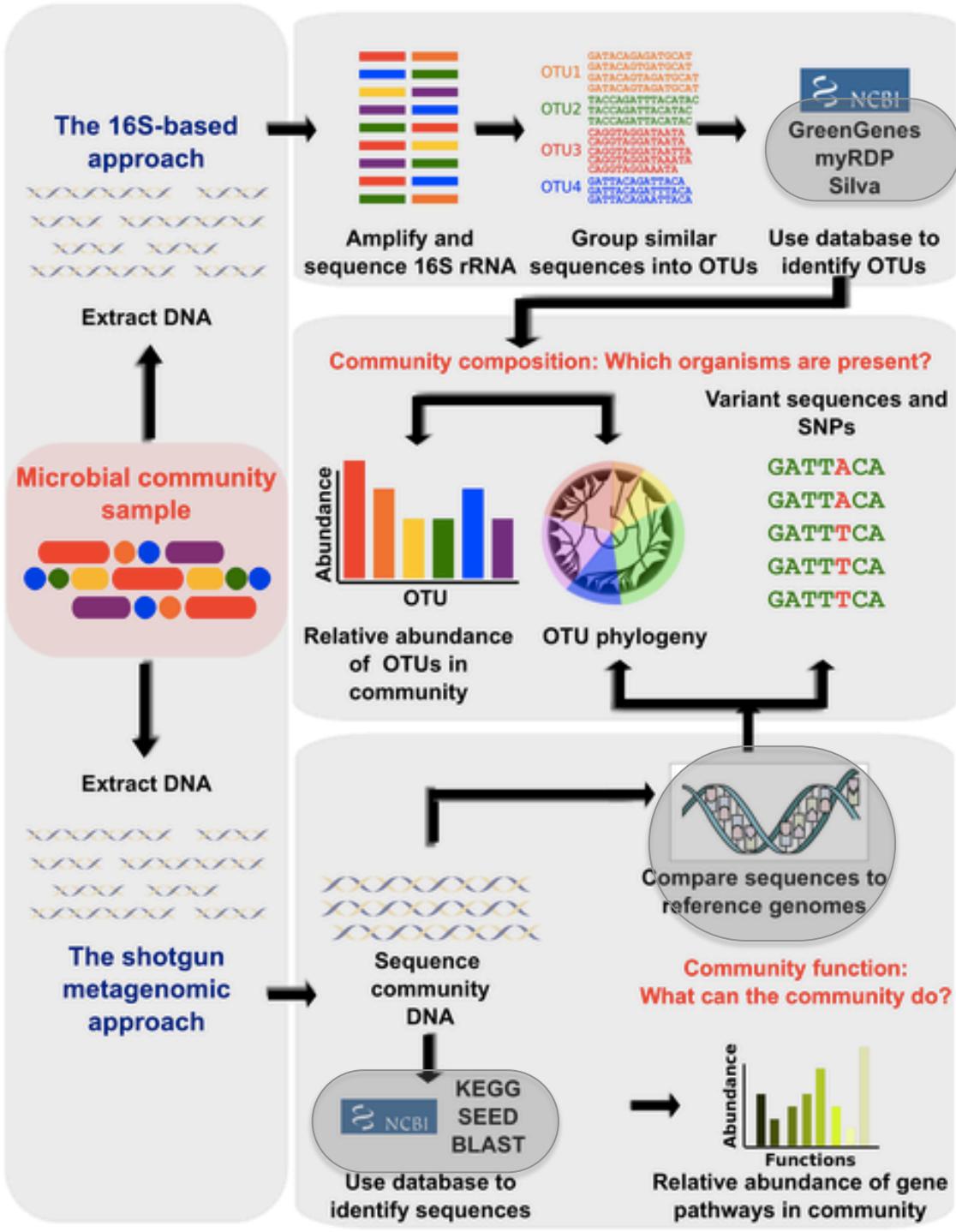
- Predict metagenome functional content from marker gene

# Workshop 11: Metagenomics Analysis

*Shi, Baochen  
Department of Pharmacology, UCLA*

# The Microbiome study





**qIIME**  
Quantitative Insights Into Microbial Ecology

VS.

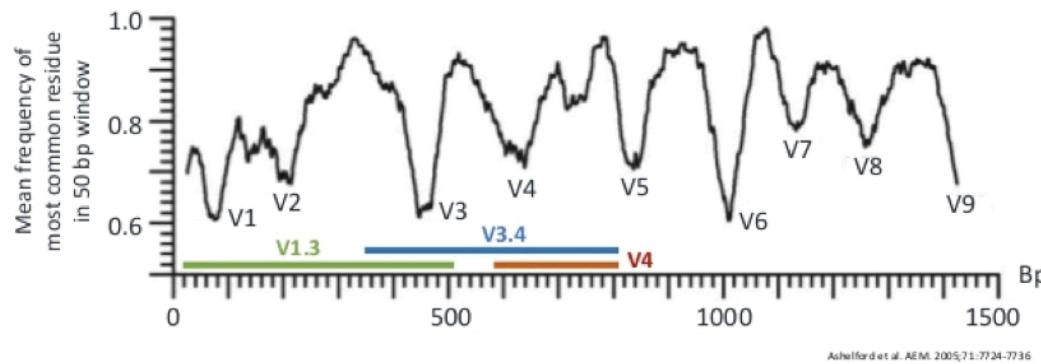
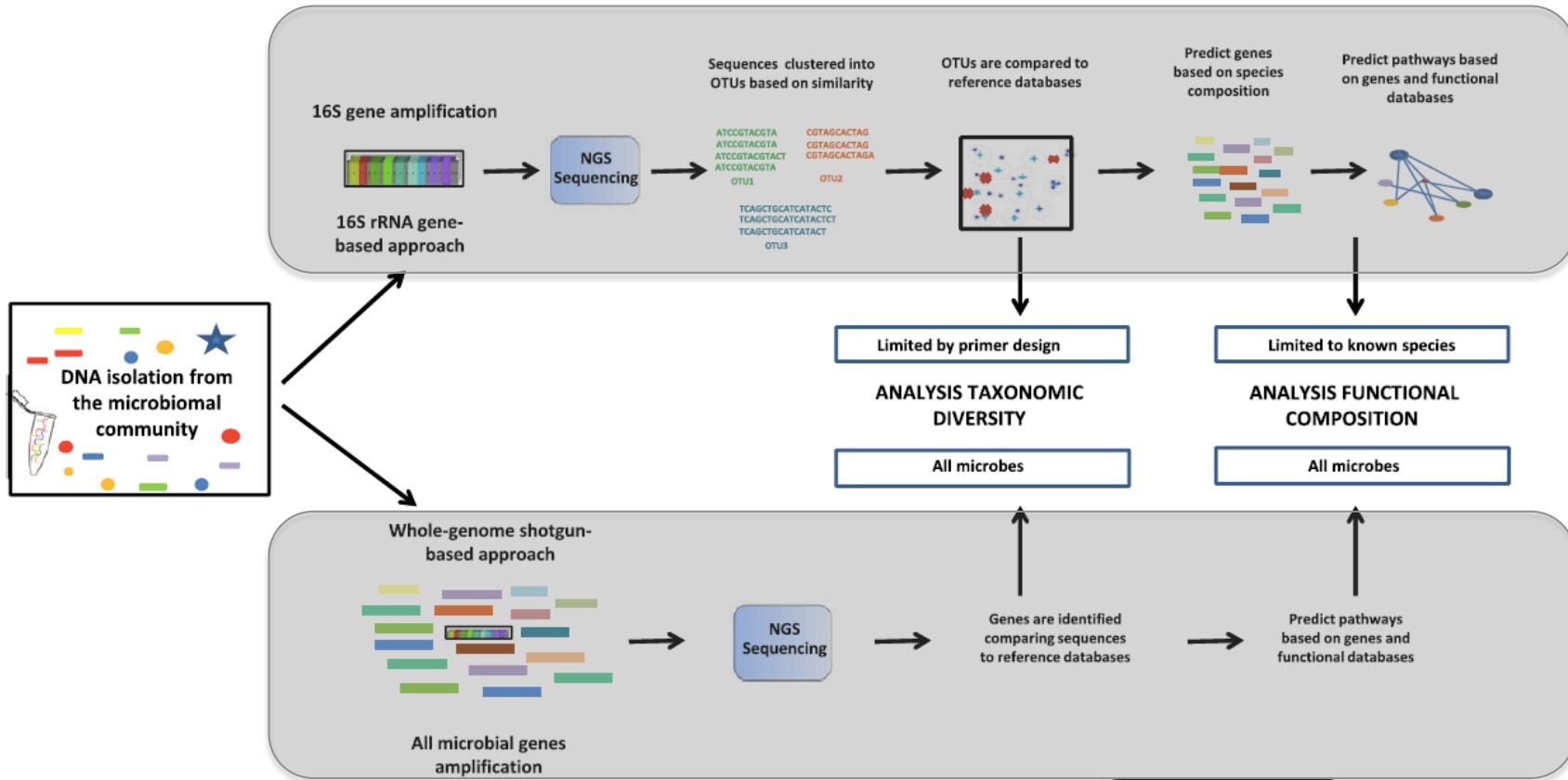


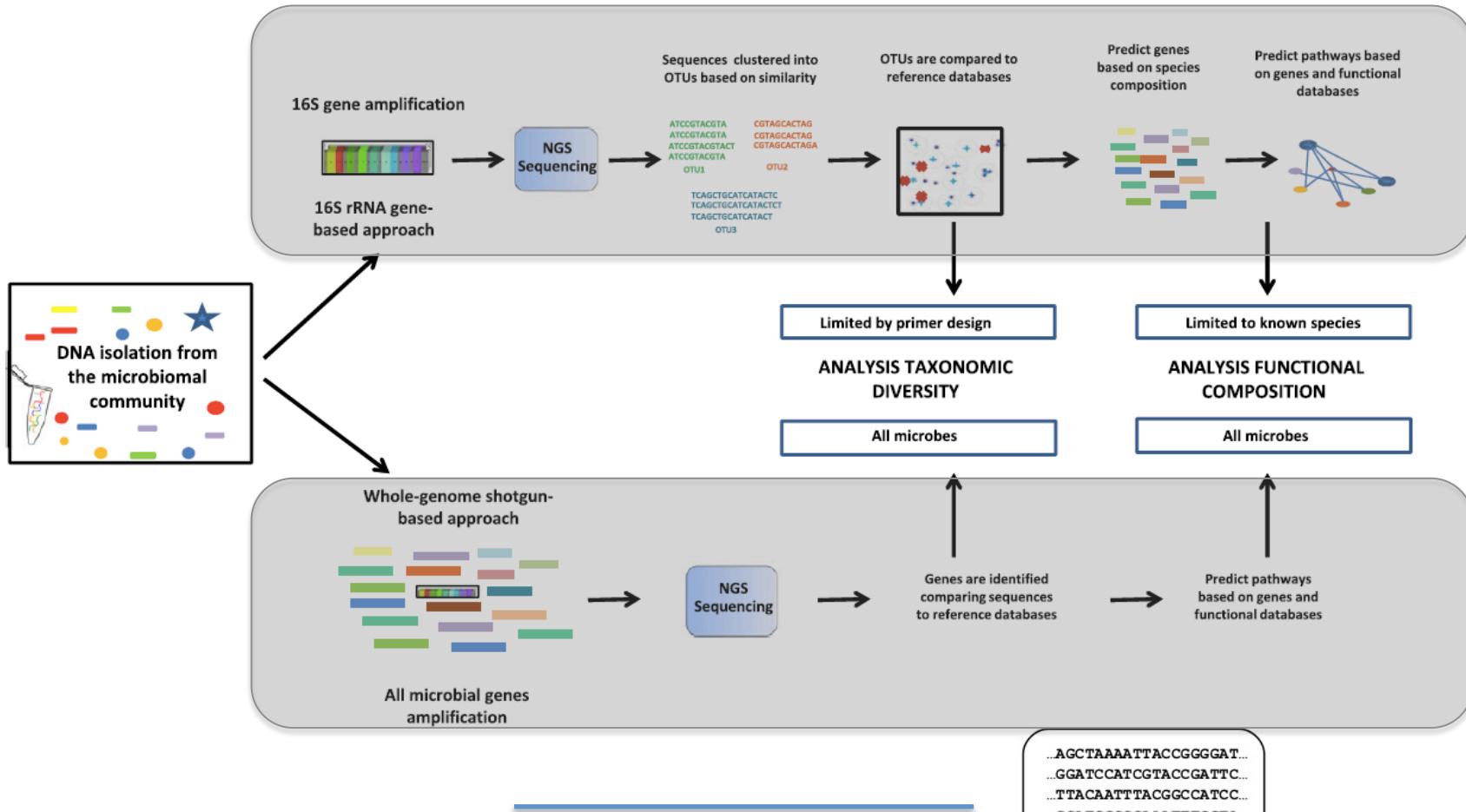
Mothur

**PICRUSt**

- Predict metagenome functional content from marker gene

**MG-RAST**  
metagenomics analysis server





Reference-based method

**MG-RAST**  
metagenomics analysis server

de novo assembly

**Ray**      **SOAP**

..AGCTAAAATTACCGGGGAT...  
..GGATCCATCGTACCGATT...  
..TTACAATTTCAGGCCATCC...  
..CCATGGCCGAATTTCGTA...  
..CCATGCGATCGATCGGAAT...  
...

Shotgun reads (~ 800 bases)

↓ Assembly

Contig      Contig      Singleton

Non-redundant sequence of microbial DNA

↓ Intensive analysis of the sequences by bioinformatics

# Outlines

---

The workshop: 2 hours per day over 3 days.

Day 1.

- i) how to perform the 16S rRNA-based analysis using bioinformatics pipeline QIIME

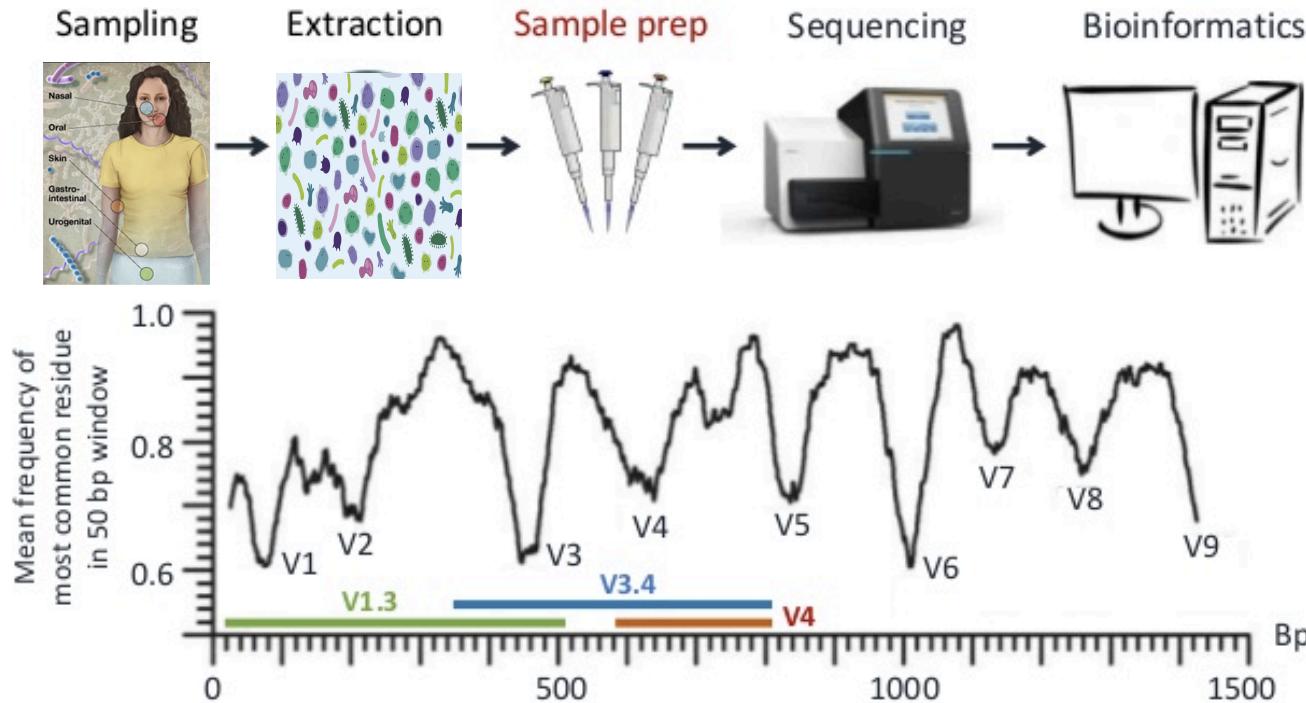
Day 2.

- i) Introduce statistical analyses in QIIME

Day 3.

- i) Functional analyses of the microbiome
- ii) Open Q&A

## Typical workflow



Ashelford et al. AEM. 2005;71:7724-7736

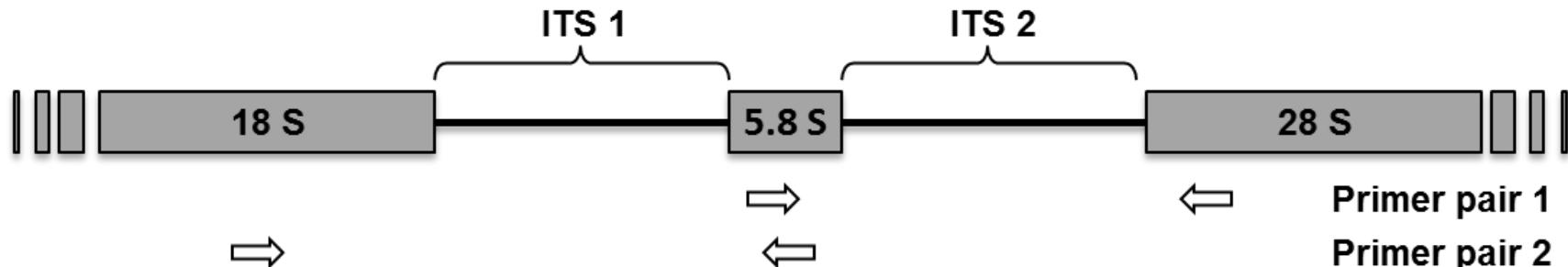
Universal, unbiased, distinguishable → 16S rRNA gene

V1-V3: human skin/vaginal microbiome

V3-V5: human gut microbiome

V4: human gut microbiome or soil microbiome

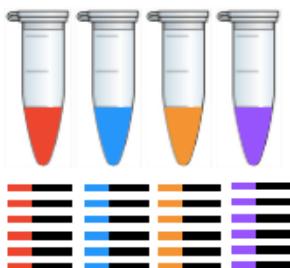
# Outlines



fungal ITS sequences



**Extract DNA and PCR amplify with barcoded primer**



Pool amplicons



**Pyrosequence**  
amplicons using 454's  
GS FLX instrument

Screen, assign  
sequences to samples  
using barcode

```
>AGTGAGAGAAGCAGGGTCGTAATGTT . . .
>AGTGCAGTCGTAGGGTCGTAATGCG . . .
>AGTGCAGTCGTAGGGTCGTAATGTA . . .
>AGTGGATGCTCTAGGGTCGTAATGCA . . .
>AGTGTACGGTGAGGGTCGTAATGGG . . .
>AGTGGATGCTCTAGGGTCGTAATGTT . . .
>AGTGTACGGTGAGGGTCGTAATGCC . . .
>AGTGAGAGAAGCAGGGTCGTAATCAC . . .
. . .
```

# Outlines



We will demonstrate 16S amplicon analysis using QIIME

- a) Run QIIME on Hoffman2 or local installation
- b) Sequence data preparation
- c) Operational Taxonomic Units (OTU) picking, Taxonomic assignment & inferring phylogeny
- d) microbiome diversity analyses

# Useful UNIX Commands: Covered in Workshop 1 & 2

- Where am I?                                      pwd
- Current directory                                 ./
- Home directory                                     ~/
- Change directory                                  cd ~/data
- Move up one level                                 cd ..
- List files in folder                               ls
- Look at a file                                     less fileName
- Copy a file   cp ~/data/file ~/otherdir/
- Delete a file                                       rm fileName
- Delete a directory                                 rmdir ~/dirName/
- Move a file   mv ~/data/file ~/otherdir/file
- Secure copy   scp user@host1:dir/file user@host2:dir/file
- Compress a file                                     gzip -c file > file.gz
- Uncompress a file                                 gunzip file.gz
- Make a new folder                                 mkdir data2
- Count lines in a file                               wc -l fileName

# **Questions/In Doubt/Lost?**

- Unix manual: man functionX
- Google is your best friend!
- Any of your friendly QCB fellows

## a) Run on Hoffman2

---



A shared account on mac for students:

login: workshop

password: NGS\_Analysis

Login in hoffman2:

```
ssh biosbc@hoffman2.idre.ucla.edu  
qrsh
```

Qiime is in :

```
cd /u/local/apps
```

Load Qiime:

```
module load qiime
```

# Logging into Hoffman2

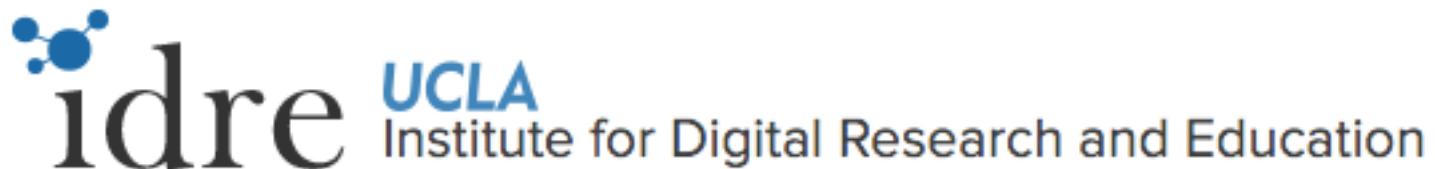
ssh biosbc@hoffman2.idre.ucla.edu

```
$ ssh biosbc@hoffman2.idre.ucla.edu  
biosbc@hoffman2.idre.ucla.edu's password:  
Welcome to the Hoffman2 Cluster!
```

change to your user ID

## a) Run on Hoffman2

---



A shared account on mac for students:

login: workshop

password: NGS\_Analysis

Login in hoffman2:

```
ssh biosbc@hoffman2.idre.ucla.edu
```

```
qrsh
```

Qiime is in :

```
cd /u/local/apps
```

Load Qiime:

```
module load qiime
```

# Useful Tools on Hoffman2

ls /u/local/apps/

abaqus	blcr	emacs	grads	libgtextutils	mumps	picard-tools	scilab	tvmet
abaqusdocs	blitz	espresso	graphicsmagick	libtool	muscle	plink	scons	uclust
accelrys	bmapuclatools	fastphase	graphviz	libxc	mygroup	pop-c++	sentaurus	udunits
Accelrys	boost	fasttree	gromacs	lmdi	mysql	povray	shapeit	usearch
activeperl	boost-jam	fastx_toolkit	gsl	loni_pipeline	namd	pplicer	shapelib	usr_lib_Gloverride
ActiveTcl	bowtie	fbat	hadoop	lumerical	ncl	preseq	skampi	valgrind
adina	bowtie2	fe-safe	handbrake	lynx	nco	proj.4	slots	vasp
Adobe	bwa	ffmpeg	harminv	mach	netbeans	protobuf	snappy	vcftools
affymetrix	caffe	fft2w	haskell	mafft	netcdf	pypy	solar	vegas
AFNI	casava	fft3w	hdf	manorm	newbler	pypy3	sortmerna	ViennaRNA
aida	cdbtools	flex	hdf5	maple	ngsplot	python	soxr	vim
amber	cd-hit	fltk	homer	maq	nlopt	qcachegrind	splicetrap	visit
ampliconnoise	cernlib	freebayes	hyperworks	maqview	numpy	qchem	spm	vmd
anaconda	cern_root	freesurfer	iaida	mathematica	nwchem	qgfe	stata	votca
annovar	clearcut	freetds	idl	matio	ocaml	qhull	stressappstest	vtk
ansys_inc	clhep	fribidi	idr	matlab	octave	qiime	structure	wannier90
ant	cln	fsl	igraph	mats	omssa	qiime data	subversion	weblogo
antlr	cmake	gamess	ilog	mecab	oommf	qrupdate	suitesparse	wolfram
armadillo	common	gatk	IM-IMA	meep	openbabel	qsub	sundials	x264
arpack	comsol	gatk-queue	impute	meld	opencv	qt	superlu_dist	x86_open64
atlas	conan	gaussian	imsl	merlin	openfoam	quantiSNP	swarm	xerces-c
atomeye	condor	gaussview	infernal	metis	openmotif	queue.logs	sysstat	xfDTD
atompaw	consed	gcta	inspect	mfold	opensees	R	szip	xfig_OLD
autoconf	cp2k	gdal	installation	microbiomeutil	osiris	raxml	t500	xilinx-vivado
automake	cpmd	gdc	intel	migrate	osu-micro-benchmarks	rdp_classifier	tau	xmd
bamtools	ctffind	geant4	isight	minirosetta	papi	relion	tcl	xmedcon
bcftools	cufflinks	genetorrent	jags	mira	paraview	remcom	tcsh	xpdf_OLD
bcl2qfastq	cytoscape	geos	jam	molcas	parmetis	repeatmasker	tecplot360	yaml
beagle	ddd	gflags	java	molden	parsec	RepeatModeler	texlive	yasm
beaglecall	ddplot	gftp	jaxodraw	molpro	parsiinsert	resmap	tmux	zlib
beagle_utilities	dejagnu	ginac	jmol	mono	pcre	rosetta	tophat	
bedtools	dirac	git_OLD	julia	mopac	pdt	rstudio	toscastructure	
bfast	diskusage	globalarrays	kepler	mothur	pennncnv	rsync	totalview	
bioscope	dl_poly	glog	lammps	mpc	perl_modules	rtax	trans-ABYSS	
blas	dx	gmp	lapack	mpfr	petsc	ruby	treemix	
blast	eclipse	gnuplot	lapack++	mpiblast	pgsql	sage	trilinos	
blast+	eigen	gpac	ldope	mplayer	phantompeakqualtools	samtools	trinity	
blat	eigensoft	grace	leveldb	mrt	phase	scalapack	tt	

## a) Run on Hoffman2

---



A shared account on mac for students:

login: workshop

password: NGS\_Analysis

Login in hoffman2:

```
ssh biosbc@hoffman2.idre.ucla.edu  
qrsh
```

Qiime is in :

```
cd /u/local/apps
```

Load Qiime:

```
module load qiime
```

# On Hoffman2 Interactively

#grab a Hoffman2 session:

**qrsh**

or

**qrsh -now n -l i,mem=5G,time=02:00:00,exclusive=TRUE -pe shared 6**

Request exact configuration

Reserve the whole node

Request 6 cores

time limit of the session  
(default is 2 hours, 24hrs max)

Memory requested per core (default is 1G)

## a) MacQIIME local installation

---

QIIME have a lot of dependencies .

The quickest way to get started using QIIME VirtualBox or MacQIIME

<http://qiime.org/1.9.0/install/install.html>

MacQIIME:

a1) Download MacQIIME (<http://www.wernerlab.org/software/macqiime/macqiime-installation>)

```
tar -xvf MacQIIME_*.tgz
```

a2) install MacQIIME

Copy macqiime folder to root, and then copy "macqiime" script to /usr/bin/

```
./install.s
```

source the environment variables

```
source /macqiime/configs/bash_profile.txt
```

Tutorial, test data & test script:

/macqiime/QIIME/

## a) QIIME VirtualBox installation

---



QIIME VirtualBox (Linux, windows, Mac):

- a1) Download and install oracle VirtualBox ([http://qiime.org/1.9.0/install/virtual\\_box.html](http://qiime.org/1.9.0/install/virtual_box.html))
- a2) Download the QIIME Virtual Box
- a3) Create a new virtual machine

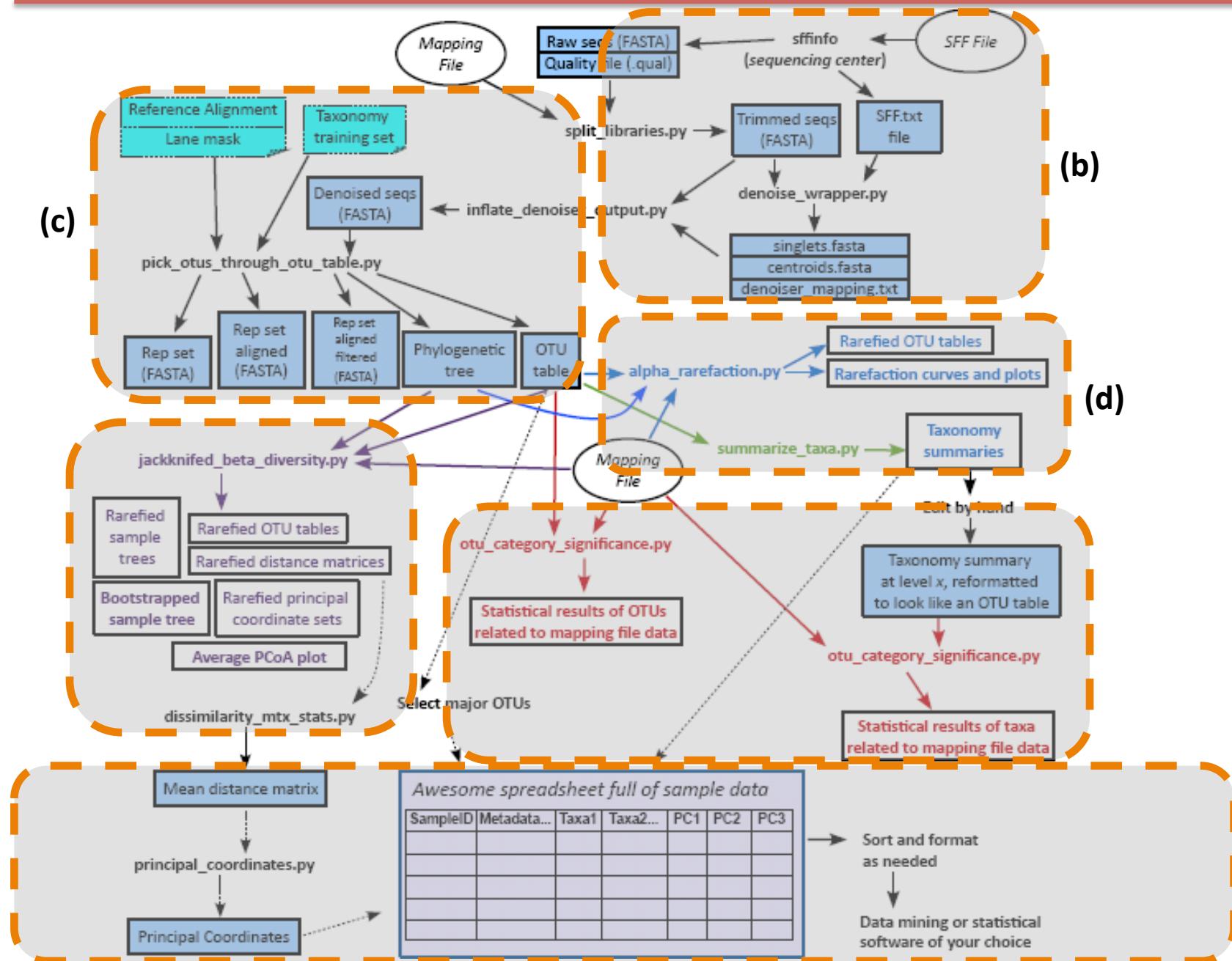
# Outlines

---

We will demonstrate steps for

- a) Run QIIME on Hoffman2 or local installation
- b) Sequence data preparation
- c) Operational Taxonomic Units (OTU) picking, Taxonomic assignment & inferring phylogeny
- d) microbiome diversity analyses

# Flowchart



# Flowchart

1. SFF (raw 454 data, optional)

2. fasta/qual files

3. demultiplexing/quality filtering

**(b) Sequence data preparation**

4. OTU picking

5. representative sequences

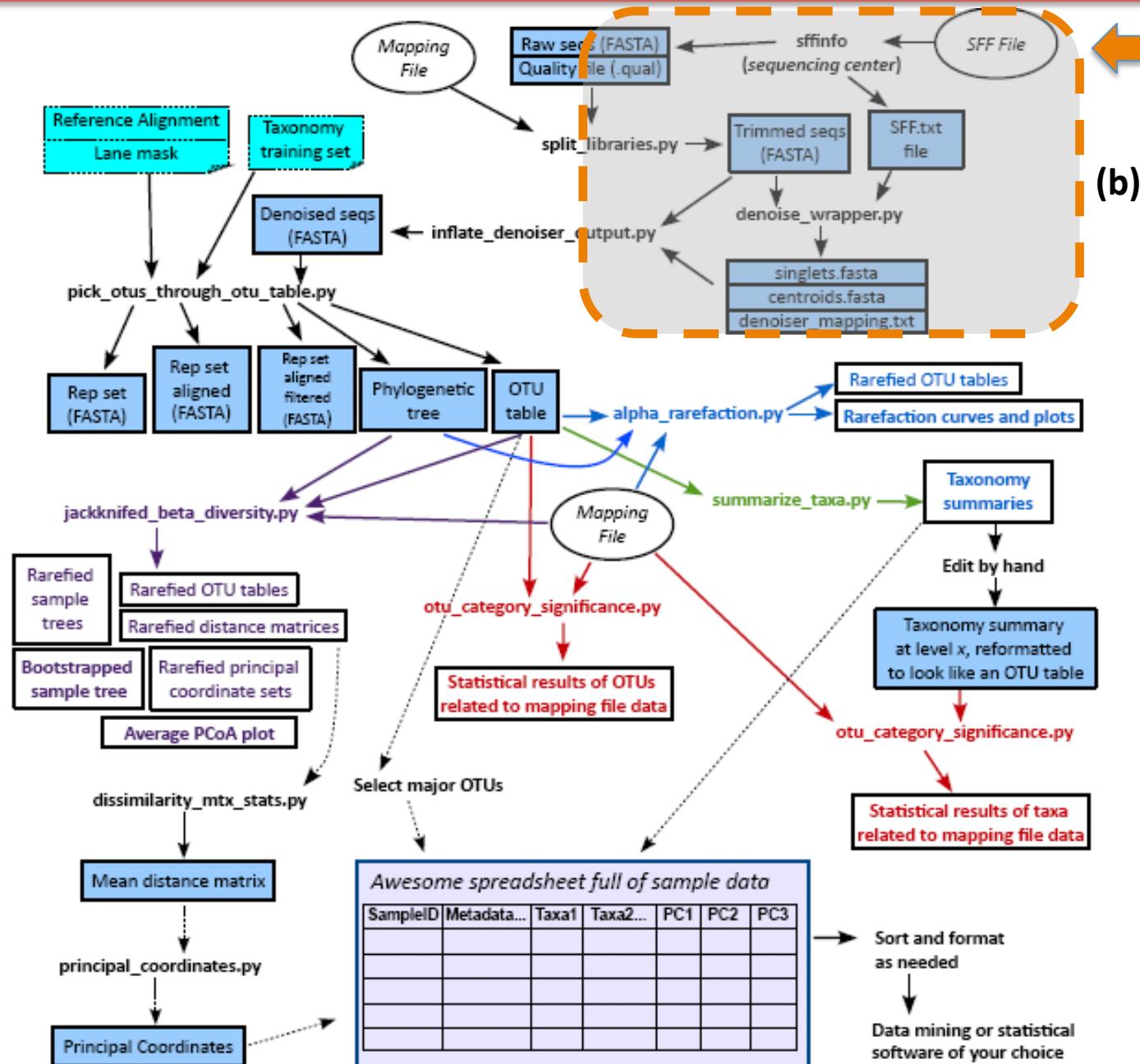
6. taxonomic assignments/tree building

**(c) Operational Taxonomic Units (OTU) picking,  
Taxonomic assignment & inferring phylogeny**

7. OTU table and downstream processing

**(d) microbiome diversity analyses**

# Flowchart



## b) Data preparation

---

We will start out with raw sequencing data generated on **454**

- Sequences (.fna): 454 generated FASTA file.
- Quality Scores (.qual): 454 generated quality score file
- experimental data about the samples (Mapping File) generated by user.

#SampleID	BarcodeSequence	LinkerPrimerSequence	Treatment	DOB	Description
#Example mapping file for the QIIME analysis package.					
PC.354	AGCACGAGCCTA	YATGCTGCCTCCCGTAGGAGT	Control	20061218	Control_mouse_I.D._354
PC.355	AACTCGTCGATG	YATGCTGCCTCCCGTAGGAGT	Control	20061218	Control_mouse_I.D._355
PC.356	ACAGACCACTCA	YATGCTGCCTCCCGTAGGAGT	Control	20061126	Control_mouse_I.D._356
PC.481	ACCAGCGACTAG	YATGCTGCCTCCCGTAGGAGT	Control	20070314	Control_mouse_I.D._481
PC.593	AGCAGCAGTTGT	YATGCTGCCTCCCGTAGGAGT	Control	20071210	Control_mouse_I.D._593
PC.607	AACTGTGCGTAC	YATGCTGCCTCCCGTAGGAGT	Fast	20071112	Fasting_mouse_I.D._607
PC.634	ACAGAGTCGGCT	YATGCTGCCTCCCGTAGGAGT	Fast	20080116	Fasting_mouse_I.D._634
PC.635	ACCGCAGAGTCA	YATGCTGCCTCCCGTAGGAGT	Fast	20080116	Fasting_mouse_I.D._635
PC.636	ACGGTGAGTGTC	YATGCTGCCTCCCGTAGGAGT	Fast	20080116	Fasting_mouse_I.D._636

On Hoffman2, copy all the files you need to scratch folder

```
cd $Scratch
```

```
pwd
```

```
cd /u/home/b/biosbc/ change to your ID
```

```
cd /u/local/apps/qiime/1.8.0/examples/qiime_tutorial/
```

```
cp /u/local/apps/qiime/1.8.0/examples/qiime_tutorial/* /u/home/b/biosbc/
```



Quantitative Insights Into Microbial Ecology

## b) Data preparation

---

Standard flowgram format (SFF) is a binary file format used to encode results of 454 pyrosequencing

```
process_sff.py -i Fasting_Example.sff -f -o output_dir
```

-f --make\_flowgram

- Sequences (.fna): 454 generated FASTA file.
- Quality Scores (.qual): 454 generated quality score file

## b) Data preparation



### Mapping File

#SampleID	BarcodeSequence	LinkerPrimerSequence	Treatment	DOB	Description
#Example mapping file for the QIIME analysis package.					
PC.354	AGCACGAGCCTA	YATGCTGCCTCCCGTAGGAGT	Control	20061218	Control_mouse_I.D._354
PC.355	AACTCGTCGATG	YATGCTGCCTCCCGTAGGAGT	Control	20061218	Control_mouse_I.D._355
PC.356	ACAGACCACACTCA	YATGCTGCCTCCCGTAGGAGT	Control	20061126	Control_mouse_I.D._356
PC.481	ACCGAGCGACTAG	YATGCTGCCTCCCGTAGGAGT	Control	20070314	Control_mouse_I.D._481
PC.593	AGCAGCACITGT	YATGCTGCCTCCCGTAGGAGT	Control	20071210	Control_mouse_I.D._593
PC.607	AACTGTGCGTAC	YATGCTGCCTCCCGTAGGAGT	Fast	20071112	Fasting_mouse_I.D._607
PC.634	ACAGAGTCGGCT	YATGCTGCCTCCCGTAGGAGT	Fast	20080116	Fasting_mouse_I.D._634
PC.635	ACCGCAGAGTCA	YATGCTGCCTCCCGTAGGAGT	Fast	20080116	Fasting_mouse_I.D._635
PC.636	ACGGTGAGTGTC	YATGCTGCCTCCCGTAGGAGT	Fast	20080116	Fasting_mouse_I.D._636

assign the multiplexed reads to samples based on their nucleotide barcode (*demultiplexing*)

```
split_libraries.py -m Fasting_Map.txt -f Fasting_Example.fna -q Fasting_Example.qual -o split_library_output
```

## b) Data preparation

---

- **split\_library\_log.txt** : summary of demultiplexing and quality filtering, including the number of reads detected for each sample and a brief summary of any reads that were removed due to quality considerations.
- **histograms.txt** : tab-delimited file shows the number of reads at regular size intervals before and after splitting.
- **seqs.fna** : fasta formatted

This step also performs quality filtering based on the characteristics of each sequence, removing any low quality or ambiguous reads.

```
split_libraries.py -m Fasting_Map.txt -f Fasting_Example.fna -q Fasting_Example.qual -o split_library_output
```

-l, --min\_seq\_length

Minimum sequence length, in nucleotides [default: 200]

-t, --trim\_seq\_length

Calculate sequence lengths after trimming primers and barcodes [default: False]

-s, --min\_qual\_score

Min average qual score allowed in read [default: 25]

## b) Data preparation



Illumina runs:

demultiplexing of Fastq sequence data where barcodes and sequences are contained in two separate **fastq** files.

Demultiplex and quality filter (at Phred  $\geq$  Q20) one lane of Illumina fastq data

```
split_libraries_fastq.py -i fastq.gz -b barcode.fastq.gz --rev_comp_mapping_barcodes -o outdir/ -m Fasting_Map.txt -q 19
```

# Flowchart

