

MB590-012

Microbiome Analysis

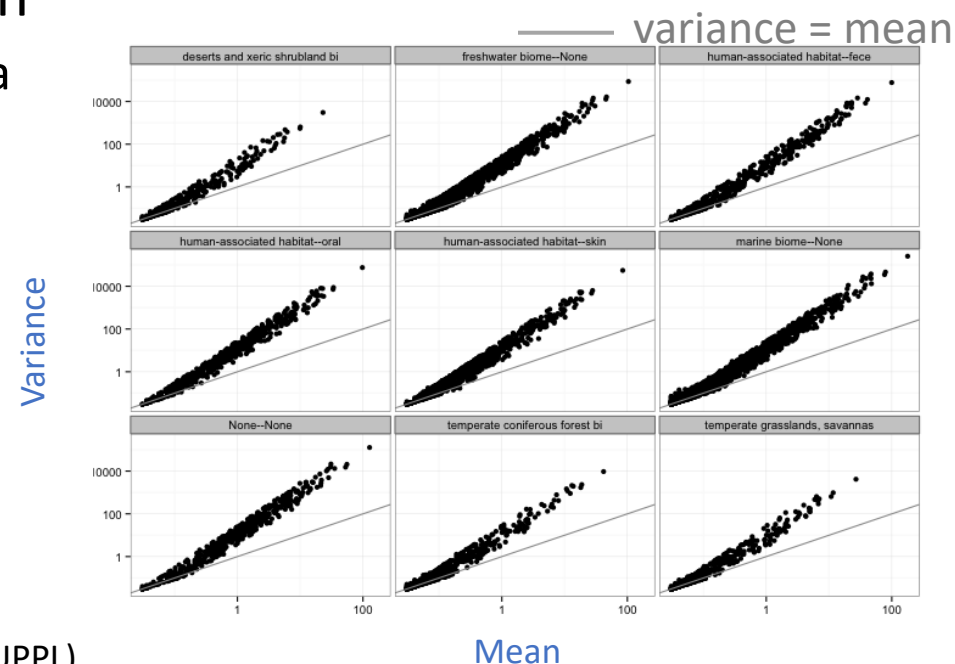
**Data considerations and
transformations**

Dr. Christine Hawkes

NC STATE UNIVERSITY

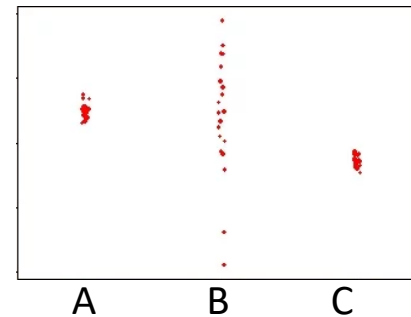
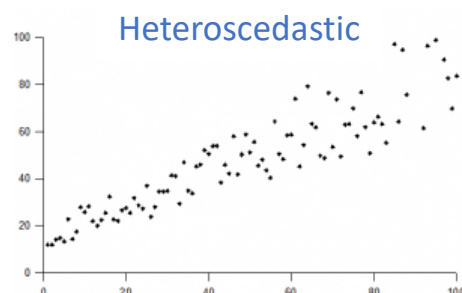
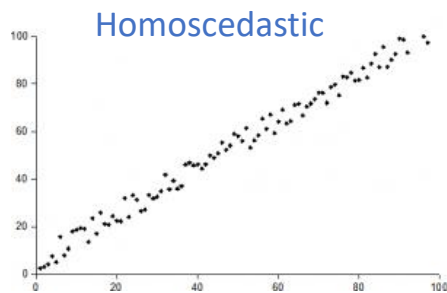
ASV matrix data

- Large number of ASVs but typically a much smaller number of samples with few replicates
- Sparse with many zeros (can be >80% zero counts!)
- Overdispersed
 - More variability than expected among biological replicates
 - Variance increases with mean
 - Example - global patterns data



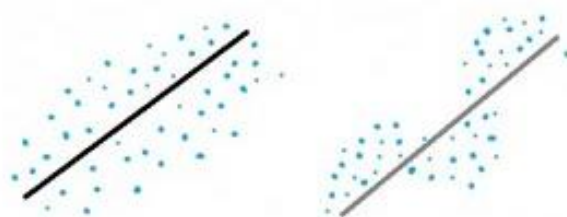
ASV matrix data

- Violate assumptions of most parametric statistics (affects parameter estimates & significance tests)
 - Linearity
 - Normality
 - Homogeneity of variance (homoscedasticity)



A vs C same var
A&C vs B diff var

- Error independence (errors are randomly distributed)



Common transformations

- Rarefaction – adjust samples to a specific count
 - *(note: I discourage you from using this approach)*
- Variance stabilizing (vst) – McMurdie & Holmes 2014
- Regularized log (rlog) – Love et al. 2014
- Centered log ratio (clr) – Gloor et al. 2017, Aitchison 1986

Rarefaction

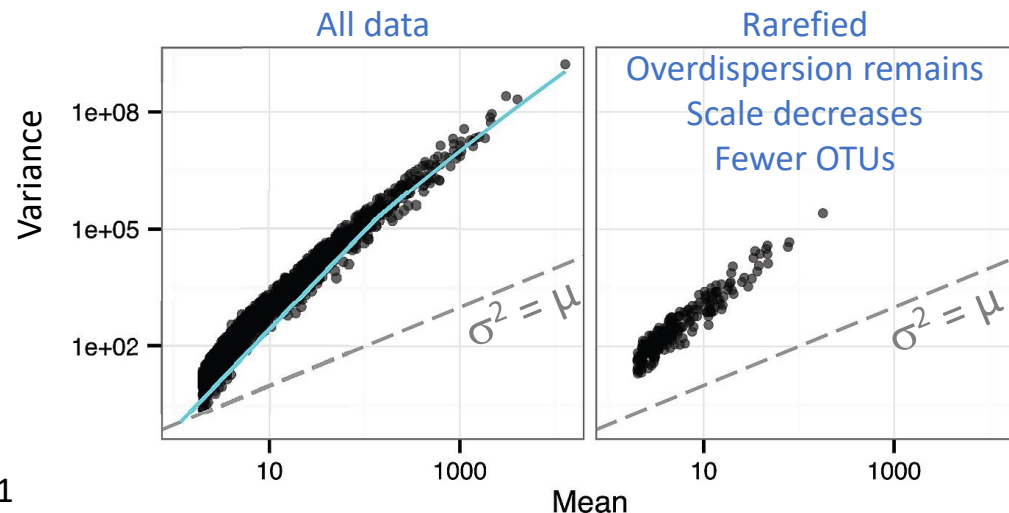
- Removes different library sizes that occur during amplification/sequencing
- Typical steps:
 - Select library size, usually the smallest in your data set
 - (Discard samples smaller than that size)
 - Randomly subsample remaining libraries without replacement to obtain the libraries of the same size in all samples

Problems with rarefaction

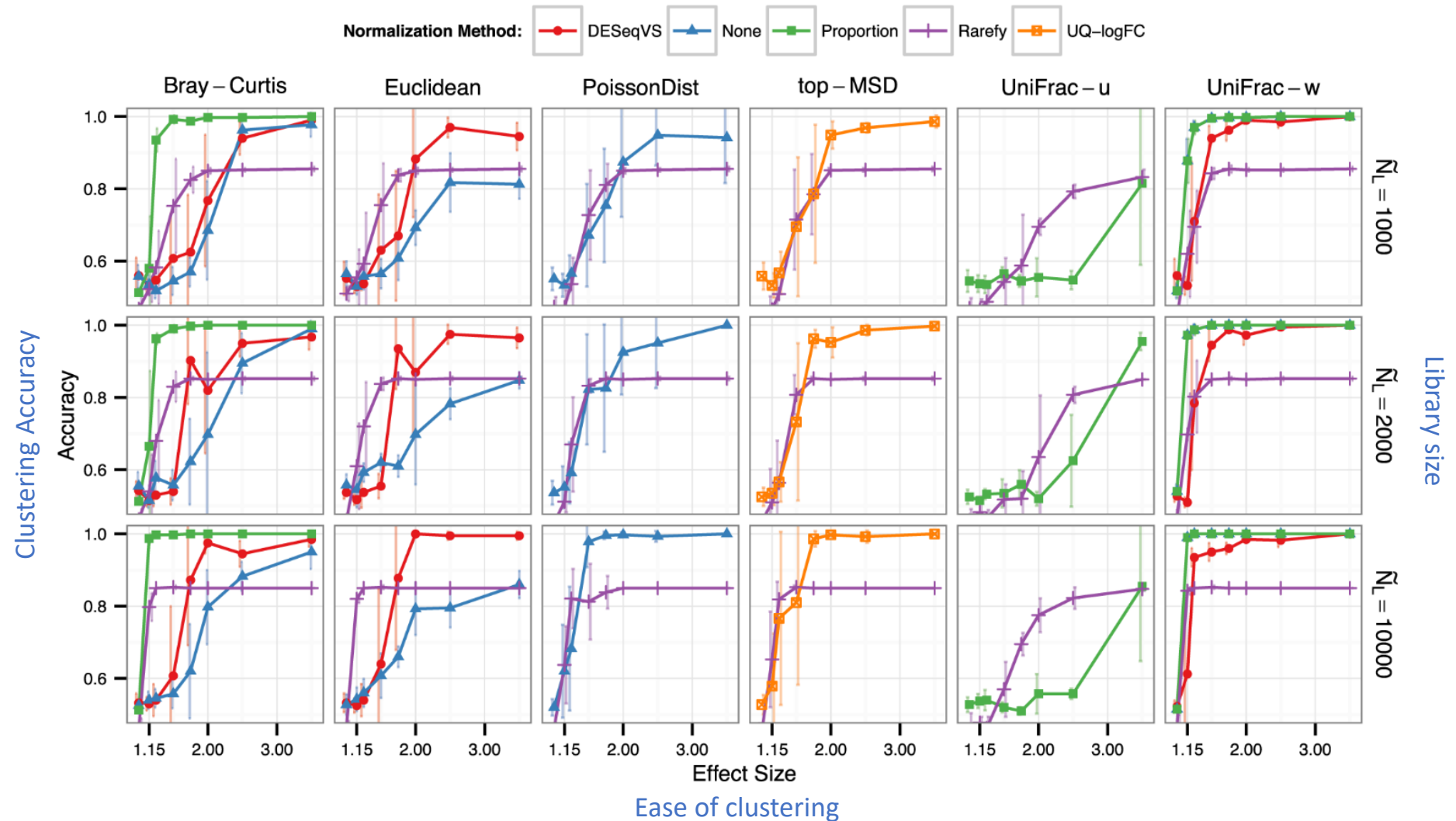
- Lack of reproducibility in subsampling
- Discards original estimates of uncertainty in abundance of each OTU by changing denominator in proportion calculations
- Inflates variance of samples to the worst value among them (smaller N in denominator) $\sigma^2 = 1/N^2 \times \sum(Y_i - \bar{Y})^2$
- Loss of information and power

Original Abundance			Rarefied Abundance		
	A	B		A	B
OTU1	62	500	OTU1	62	50
OTU2	38	500	OTU2	38	50
Total	100	1000		100	100

Standard Tests for Difference			
P-value	chi-2	Prop	Fisher
Original	0.0290	0.0290	0.0272
Rarefied	0.1171	0.1171	0.1169



Rarefaction also undermines downstream clustering (best is ~80% accuracy)



What does compositional mean?

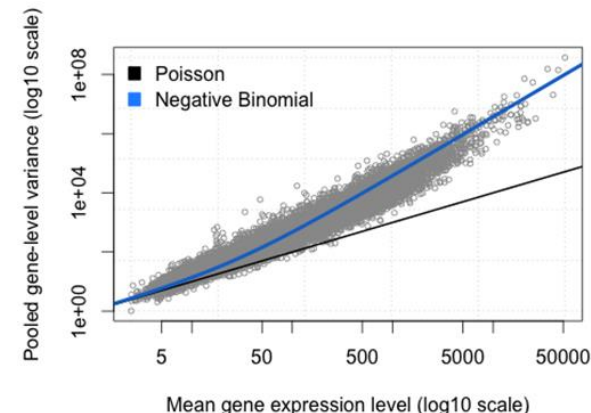
- Data provide only relative abundances
 - Proportions sum to 1
 - Not absolute abundances
- Total read number is arbitrary and constant
 - Artifact of Illumina total reads being limited by the machine
 - Can result in spurious correlations
- Reads are random subsamples of true molecules that do not necessarily reflect original sample
- Analysis of data subset is not consistent with analysis of total data (membership-dependent)
- Random sampling of variables can lead to change of conclusions (as seen with rarefaction)

How to deal with compositional data

- Initially, data were scaled to total read abundance to generate relative abundances (proportions)
 - Can skew true relationships in the data
 - Does not address problems with overdispersion
- Recent alternatives to address problems with proportions include:
 - Use of transformation methods such as vst and rlog that adjust by library size and address overdispersion
 - Log ratios transformations to preserve underlying relationships in the data

vst and rlog transformations

- Variance stabilizing (vst) – McMurdie & Holmes 2014
 - Factors out dependence of var on the mean and normalizes by library size
- Regularized log (rlog) – Love et al. 2014
 - Minimizes differences between samples for rows with small counts and normalizes by library size
- Both are run in DESeq2
 - Originally developed for gene expression data
 - Assume a negative binomial distribution (overdispersion)
 - Better fit than traditional Poisson



Variance stabilizing transform (vst)

- Addresses different library sizes, sparse data, heteroscedasticity, overdispersion
- To do this, DeSeq2 identifies the mean-variance relationship in the data, and then derives a transformation that will remove the dependence so that the variance is the same across samples (homoscedasticity)
- Approaches the \log_2 transformation for large values and square root function for small values

Regularized log transform (rlog)

- Addresses same same issues as the vst but:
- Transforms count data to log2 scale by fitting a model that includes a term for each sample plus a Bayesian prior distribution on the model coefficients (estimated from the data)

$$\log_2(q_{ij}) = \beta_{i0} + \beta_{ij}$$

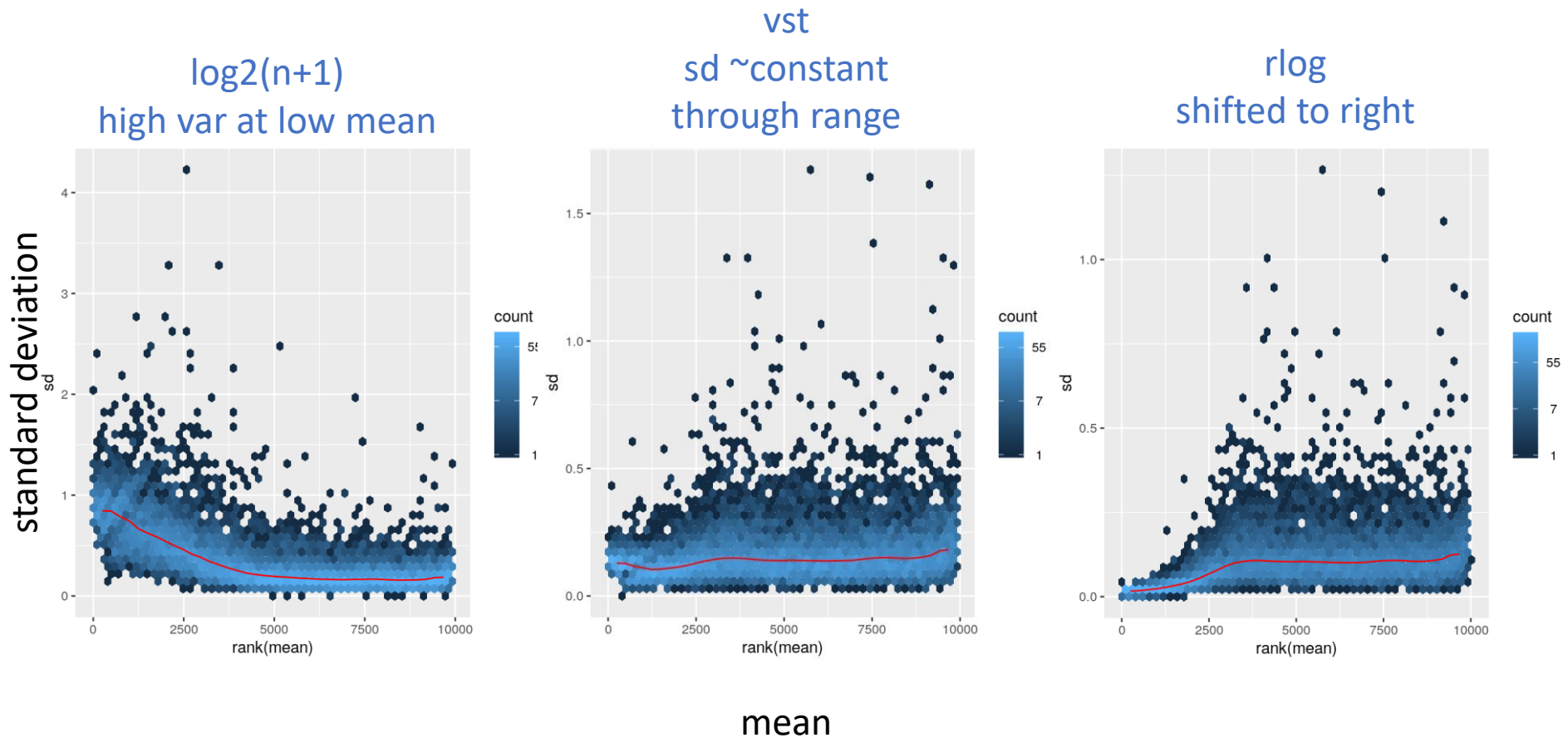
Parameter proportional to true count for ASV i in sample j

intercept

sample-specific effect, shrunk towards zero based on dispersion-mean trend in entire dataset

- Much slower than vst for large data sets (100s of samples)

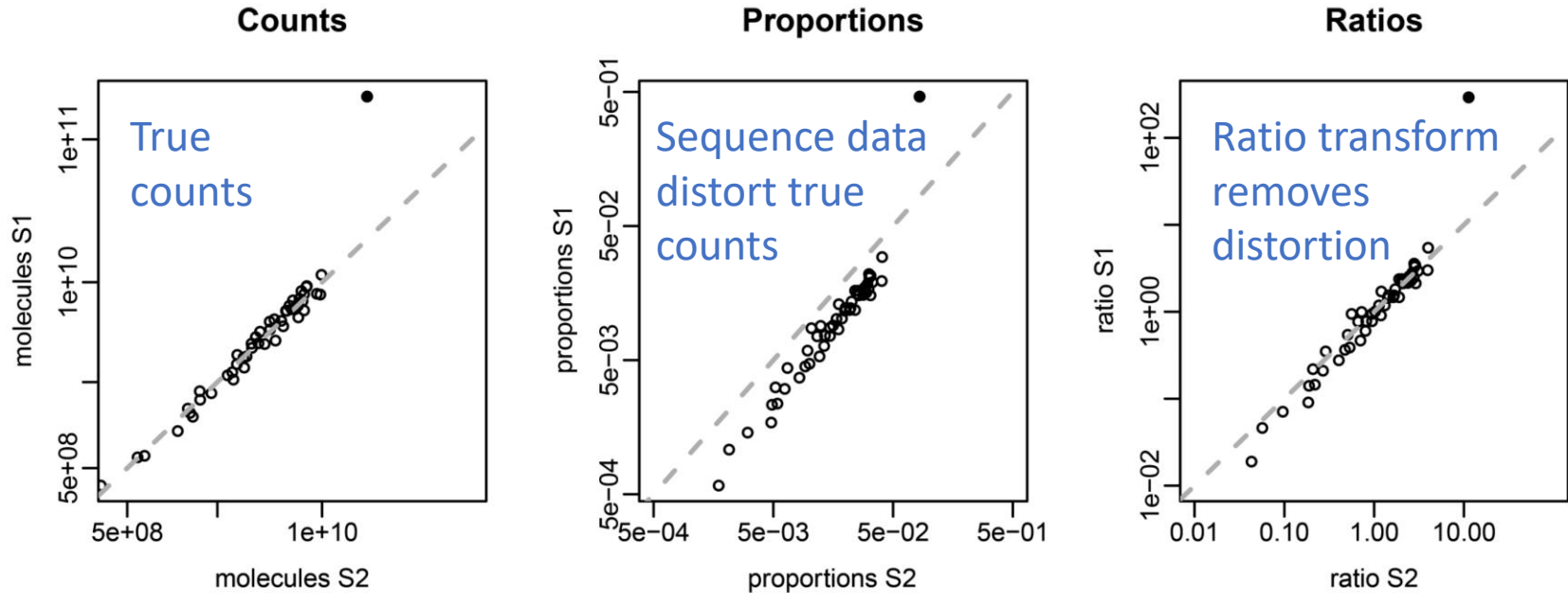
vst and rlog transformation effects on variance



vst and rlog application in DESeq

- Can transform data in context of experimental design or “blind” (intercept only)
 - Design: $\sim \text{Trt1} + \text{Trt2}$ vs. Blind: ~ 1
- Blind is inappropriate if we expect ASV counts to be explained by experimental design
 - can lead to large dispersion estimates by assuming differences caused by treatments are noise
- Dispersion estimates have 3 fit types
 - Parametric – vst on normalized counts
 - Mean – vst on negative binomial counts with fixed dispersion
 - Local – reciprocal sq root of variance of normalized counts (derived from dispersion fit), then numerically integrated, and the integral is evaluated for each count, yielding transformed values

Ratio-based transformations



Centered log ratio transform (clr)

- For D ASVs in each sample, the clr transformation is the \ln of the proportion of counts, where the denominator is the geometric mean of all counts for that sample

$$\mathbf{x}_{clr} = [\log(x_1/G(\mathbf{x})), \log(x_2/G(\mathbf{x})) \dots \log(x_D/G(\mathbf{x}))],$$
$$G(\mathbf{x}) = \sqrt[D]{x_1 \cdot x_2 \cdot \dots \cdot x_D}$$

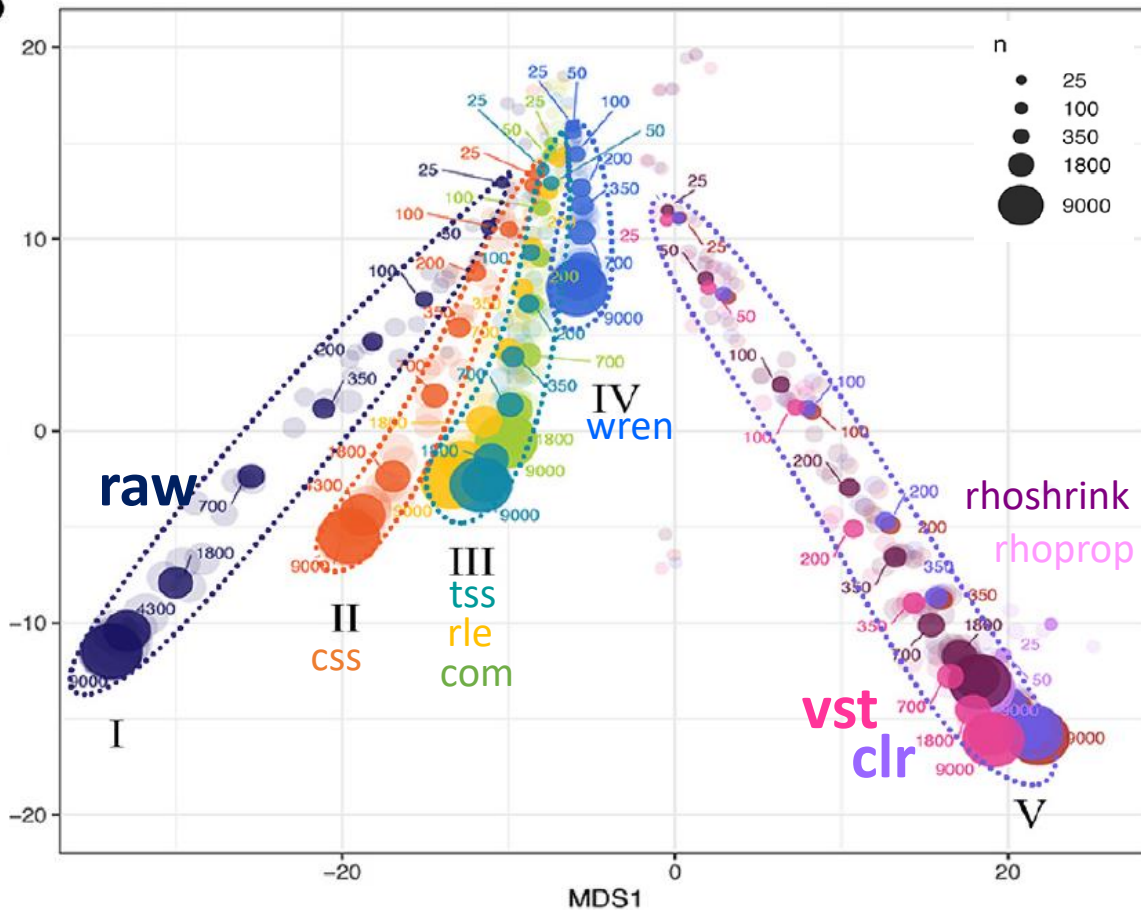
- Data remain relative (ratio-based), but are also
 - Scale invariant (unit free)
 - Sub-compositionally coherent (consistent across subsets)
- Other log ratio transforms include additive (alr), isometric (ilr), and phylogenetic isometric (philir)

What to do about zeros or negative numbers?

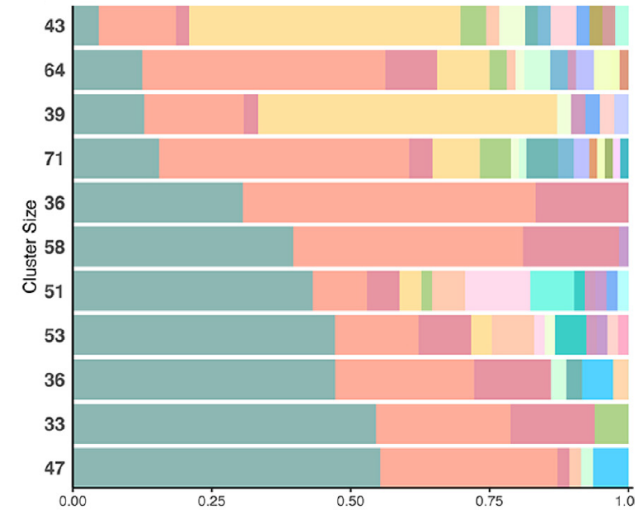
- The vst, rlog, and clr transformations may produce negative values for zeros
- This can be a problem in downstream analyses (e.g., distance matrices)
- Options
 - Set all negative values to zero
 - Add a pseudocount to all cells
 - Analyze the transformed data without further manipulation (don't calculate distances)
 - For dissimilarity analysis, use Euclidean distance on transformed data (= Aitchison distance between clr-transformed compositions)

Transformation matters to downstream analysis!

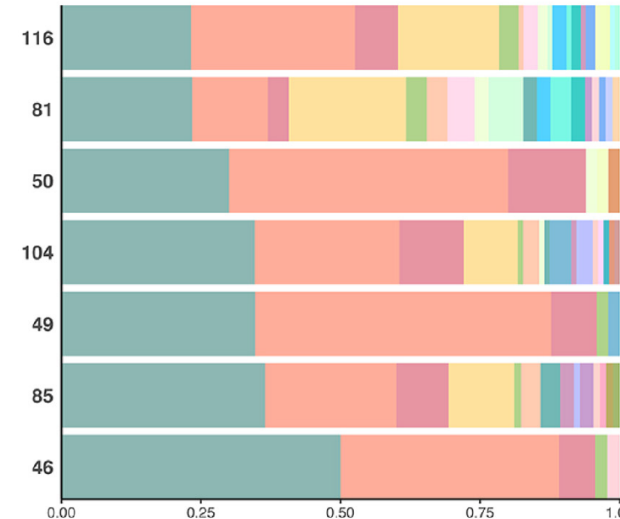
B



clr

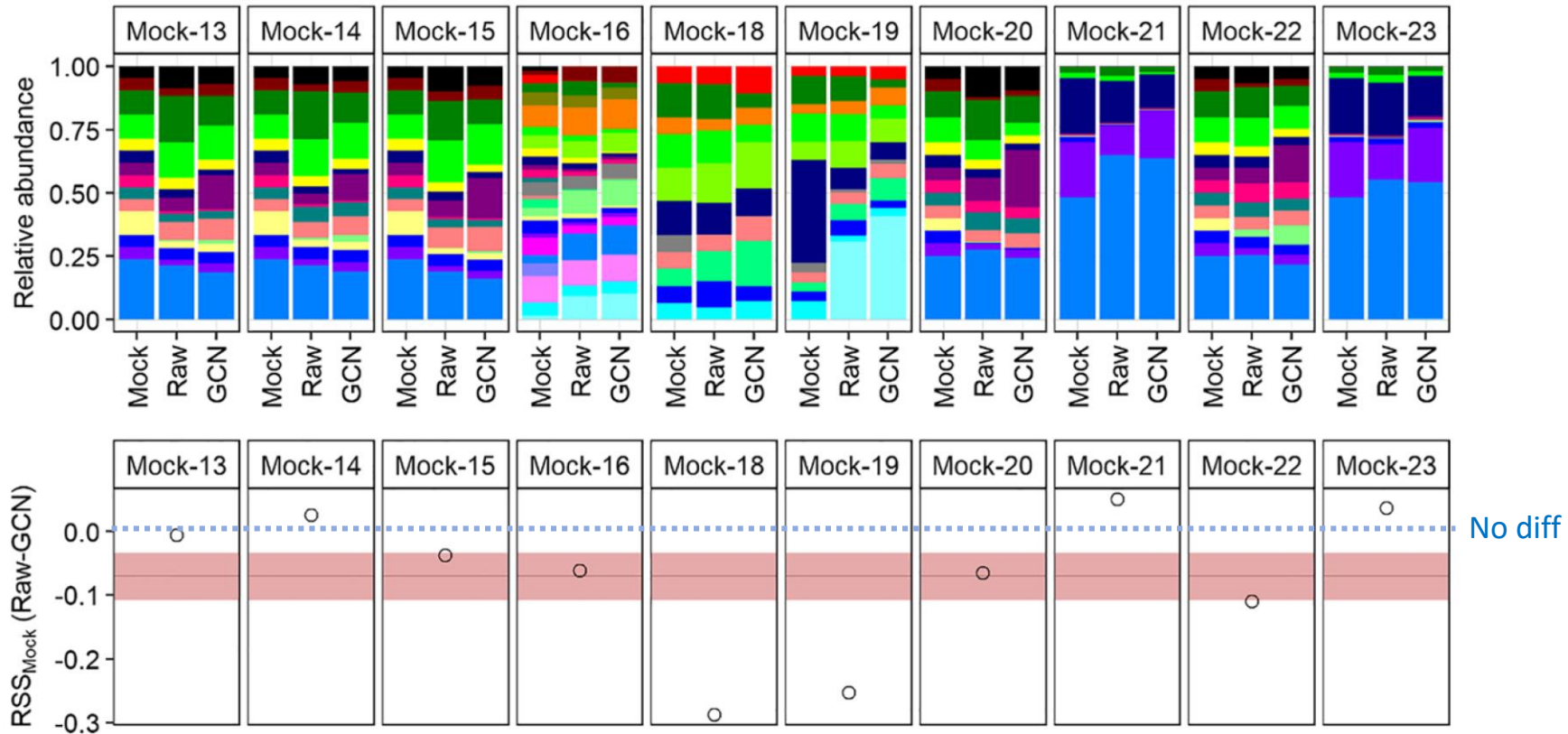


tss



Other transforms: rRNA gene copy number normalization (GCN)

Raw seq data fit actual Mock 7.1% better than GCN on avg



Bottom panel is sum of deviations between actual Mock community abundances and GCN-adjusted.

Can you get absolute abundance from microbiome data?

- Spike-in controls
- Cell flow cytometry
- qPCR
- Examples:
 - Deng, L. et al. 2019. Front. Microbiology 10: 720, DOI: 10.3389/fmicb.2019.00720
 - Hardwick et al. 2018 Nature Communications 9:3096; DOI: 10.1038/s41467-018-05555-0
 - Stammler et al. 2016 Microbiome 4:28; DOI: 10.1186/s40168-016-0175-0
 - Turlouse et al. 2016 Nucleic Acids Research 45:e23; DOI: 10.1093/nar/gkw984

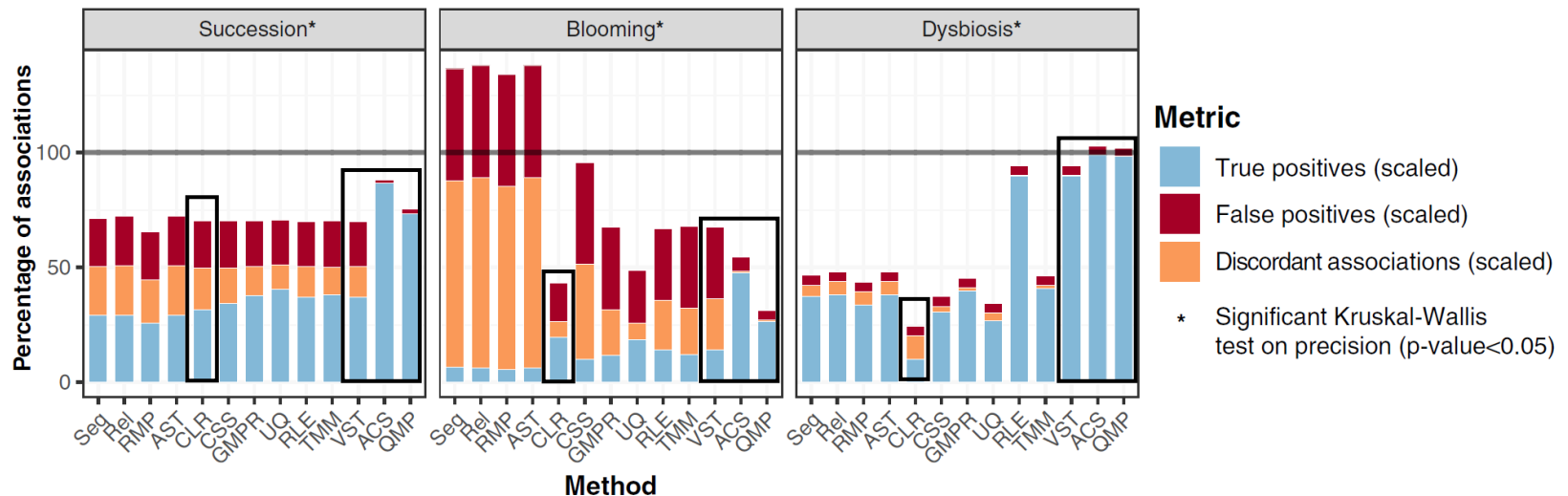
Comparing transformations to quantitative profiling

Table 2 Metagenomic data transformations benchmarked.

Method	Abbreviation	Technique	Transformation	Correction	Rarefaction	Suited for richness calculations
Raw sequencing data	Seq	–	None	–	No	Yes
Relative abundance	Rel	Computational	Relative	Sequencing depth	No	No (0:1 range)
Relative microbiome profiling	RMP	Computational	Relative	Sequencing depth	Yes	Yes
Arcsine square root	AST	Computational	Relative	Sequencing depth	No	No (0:1 range)
Centered log ratio	CLR	Computational	Compositional	Sequencing depth and compositionality	No	No (negative values)
Cumulative sum scaling	CSS	Computational	Compositional	Sequencing depth and compositionality	No	Yes (rounding data)
Relative log expression	RLE	Computational	Compositional	Sequencing depth and compositionality	No	Yes (rounding data)
Upper quantile	UQ	Computational	Compositional	Sequencing depth and compositionality	No	Yes (rounding data)
Trimmed mean of m-values	TMM	Computational	Compositional	Sequencing depth and compositionality	No	Yes (rounding data)
Geometric mean of pairwise ratios	GMPR	Computational	Compositional	Sequencing depth and compositionality	No	Yes (rounding data)
Variance-stabilizing transformation	VST	Computational	Compositional	Sequencing depth and compositionality	No	No (negative values)
Quantitative microbiome profiling	QMP	Experimental	Quantitative	Sampling depth and microbial load	Yes	Yes
Absolute count scaling	ACS	Experimental	Quantitative	Microbial load	No	Yes

Methods are categorized based on the technique applied (computational or experimental), the biases targeted (sequencing depth, sampling depth, compositionality, and/or microbial load), the inclusion of a downsizing step, and their projected suitability for richness estimations. Additionally, for study purposes, methods are broadly labeled as relative, compositional, or quantitative methods.

vst transformation was most comparable to quantitative methods

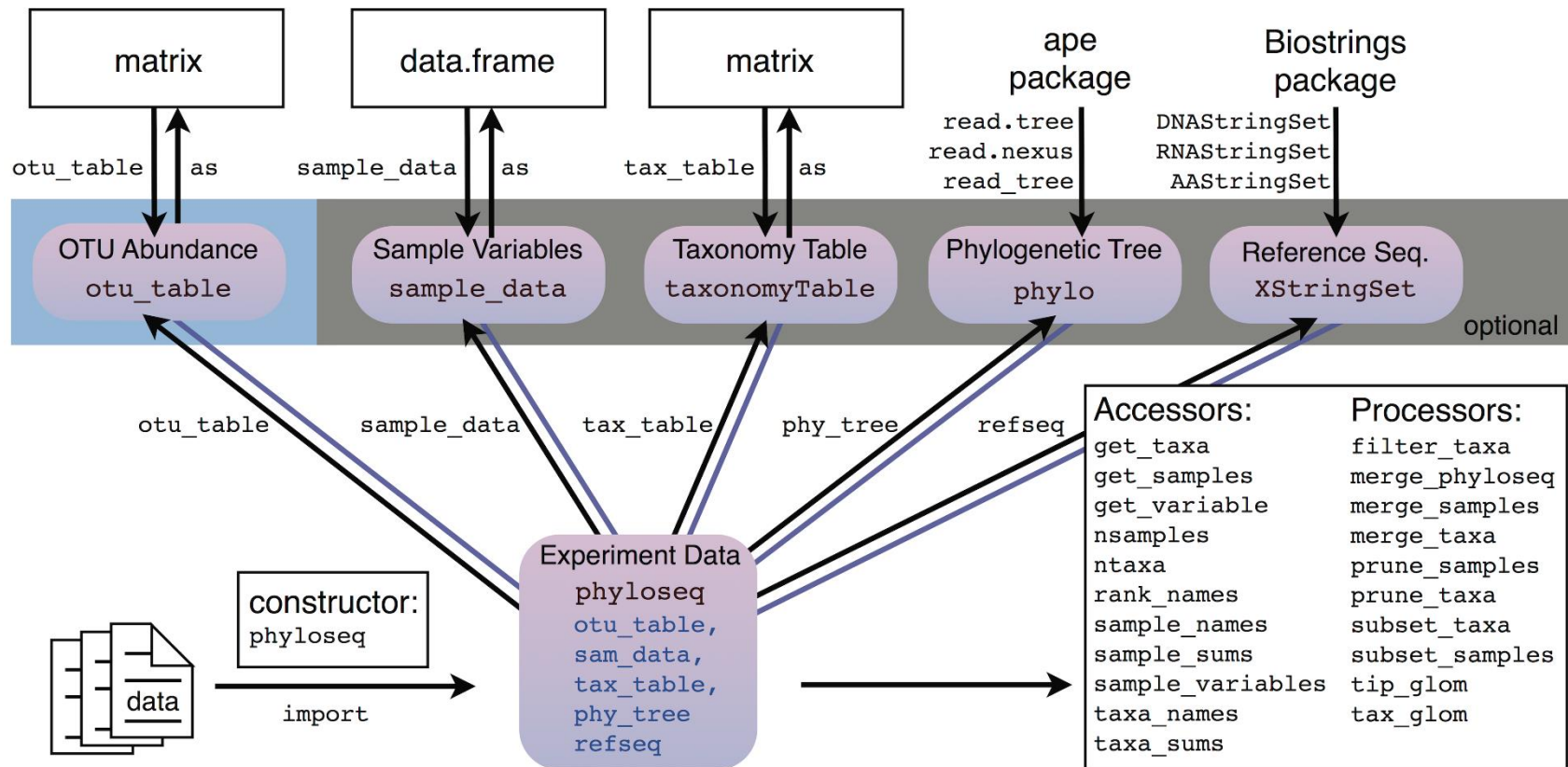


- **Quantitative methods** have the most true and fewest false positives (and no discordance)
- **vst** is most like the quantitative methods, but performance depends on dataset
- **clr** and other transforms generally perform more poorly

Different pipelines for downstream analysis for compositional data

Operation	Standard approach	Compositional approach
Normalization	Rarefaction 'DESeq'	CLR ILR ALR
Distance	Bray-Curtis UniFrac Jenson-Shannon	Aitchison
Ordination	PCoA (Abundance)	PCA (Variance)
Multivariate comparison	perManova ANOSIM	perMANOVA ANOSIM
Correlation	Pearson Spearman	SparCC SpiecEasi ϕ ρ
Differential abundance	metagenomSeq LEfSe DESeq	ALDEx2 ANCOM

Use of phyloseq for data management



Functions for building phyloseq objects

Functions for building component data objects

Function	Input Class	Output Description
<code>otu_table</code>	numeric matrix	<code>otu_table</code> object storing OTU abundance
<code>otu_table</code>	<code>data.frame</code>	<code>otu_table</code> object storing OTU abundance
<code>sample_data</code>	<code>data.frame</code>	<code>sample_data</code> object storing sample variables
<code>tax_table</code>	character matrix	<code>taxonomyTable</code> object storing taxonomic identities
<code>tax_table</code>	<code>data.frame</code>	<code>taxonomyTable</code> object storing taxonomic identities
<code>read_tree</code>	file path char	phylo-class tree, read from file
<code>read.table</code>	table file path	A matrix or <code>data.frame</code> (Std R core function)

Functions for building complex data objects

Function	Input Class	Output Description
<code>phyloseq</code>	2 or more component objects	phyloseq-class, “experiment-level” object
<code>merge_phyloseq</code>	2 or more component or phyloseq-class objects	Combined instance of phyloseq-class

Accessor functions for phyloseq objects

Function	Returns
[Standard extraction operator. works on <code>otu_table</code> , <code>sample_data</code> , and <code>taxonomyTable</code>
<code>access</code>	General slot accessor function for phyloseq-package
<code>get_taxa</code>	Abundance values of all taxa in sample 'i'
<code>get_sample</code>	Abundance values of taxa 'i' for all samples
<code>get_taxa_unique</code>	A unique vector of the observed taxa at a particular taxonomic rank
<code>get_variable</code>	An individual sample variable vector/factor
<code>nsamples</code>	Get the number of samples described by an object
<code>ntaxa</code>	Get the number of OTUs (taxa) described by an object
<code>otu_table</code>	Build or access <code>otu_table</code> objects
<code>rank_names</code>	Get the names of the available taxonomic ranks
<code>sample_data</code>	Build or access <code>sample_data</code> objects
<code>sample_names</code>	The names of all samples
<code>taxa_names</code>	The names of all taxa
<code>sample_sums</code>	The sum of the abundance values of each sample
<code>sample_variables</code>	The names of sample variables
<code>taxa_sums</code>	The sum of the abundance values of each taxa
<code>taxa_are_rows</code>	TRUE if taxa are row indices in <code>otu_table</code>
<code>tax_table</code>	A taxonomy table
<code>tre</code>	Access the tree contained in a phyloseq object