# "A Novel Sparse Compositional Technique Reveals Microbial Perturbations"
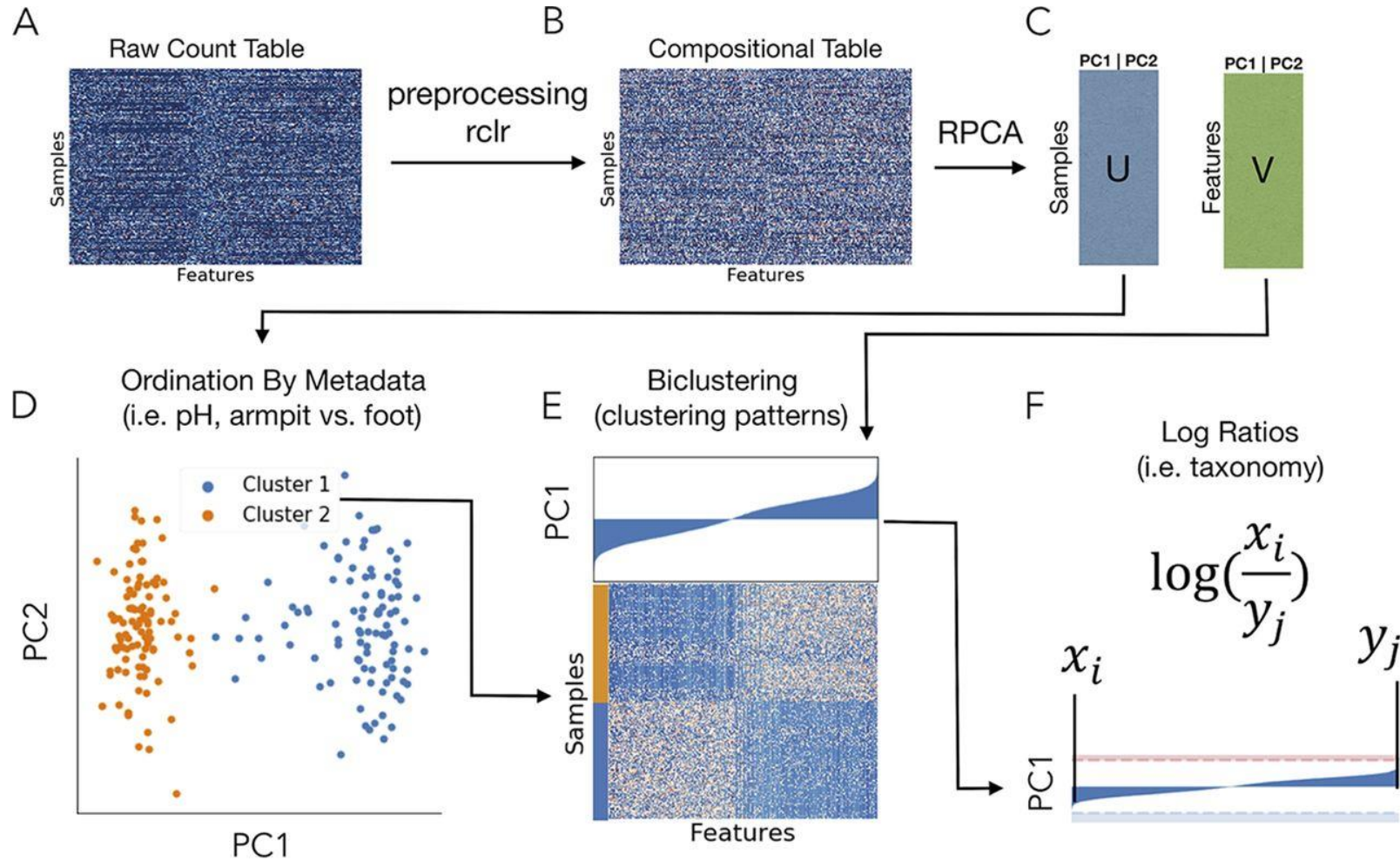
(Published in February 12th, 2019)

Authors: Cameron Martino, James T. Morton, Clarisse A. Marotz, Luke R. Thompson, Anupriya Tripathi, Rob Knight, Karsten Zengler
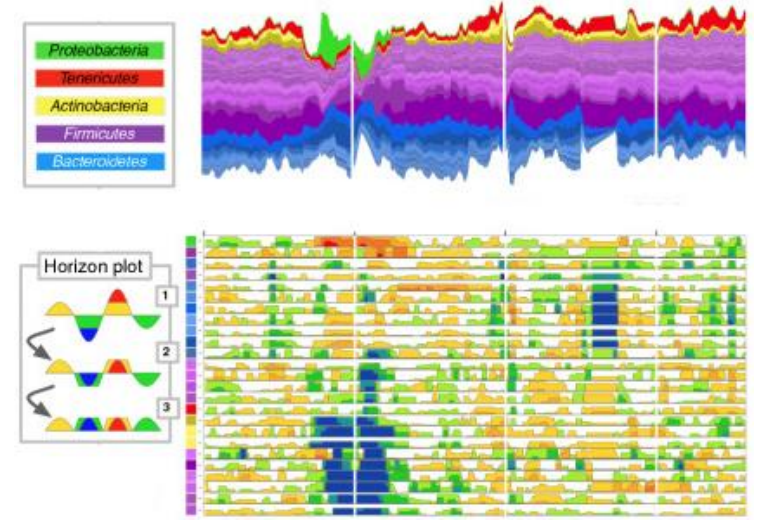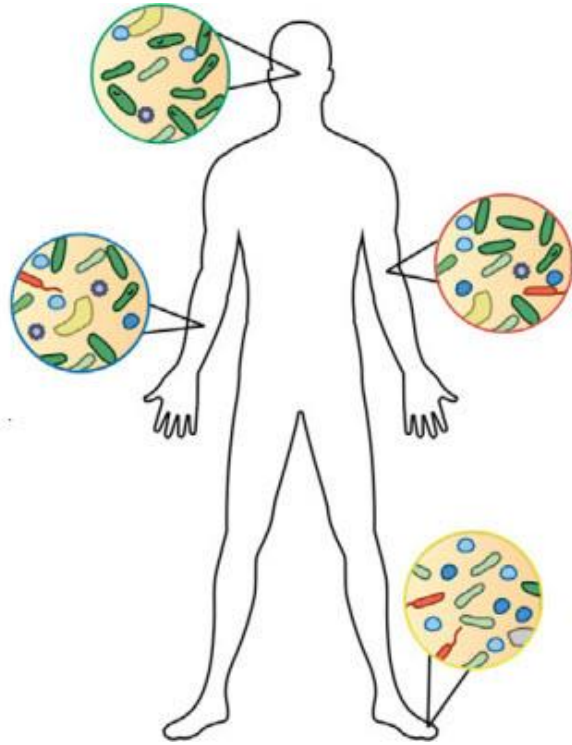
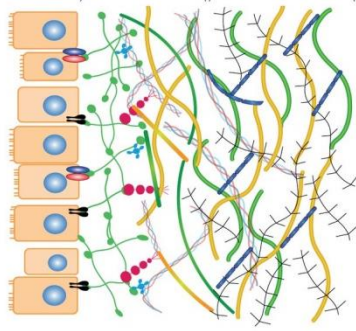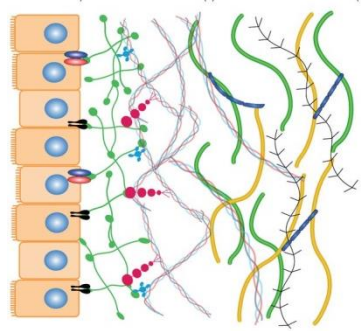Presenter: Itamar Curiel

A Novel Sparse Compositional Technique Reveals Microbial Perturbations

# ANSCTRMP?

# TL;DR



A. Raw Count Table

B. Compositional Table

C. U (Samples) | V (Features), PC1 | PC2

preprocessing rclr

RPCA

D. Ordination By Metadata (i.e. pH, armpit vs. foot)
- Cluster 1
- Cluster 2
PC2 / PC1

E. Biclustering (clustering patterns)
PC1 / Samples / Features

F. Log Ratios (i.e. taxonomy)

$$\log\left(\frac{x_i}{y_j}\right)$$

$x_i$ $y_j$

PC1

| SPECIES | FRACTION |
|---|---|
| Bacillus Thuringiensis | 0.001 |
| Campylobacter jejuni | 0.001 |
| Staphylococcus epidermidis | 0.002 |
| … | … |
| … | … |
| … | … |
| … | … |
| … | … |
| … | … |
| … | … |
| … | … |
| … | … |
| … | … |
| … | … |
| … | … |
| … | … |
| … | … |
| … | … |
| … | … |
| Methanobacterium extroquens | 0.076 |
| Bordetella bronchiseptica | 0.076 |
| Propionibacterium propionicus | 0.122 |

$x' = [0.003, 0.001, 0.002, \ldots, 0.081, 0.078, 0.122]$
(species found in this sample)

$x = [0.003, \mathbf{0}, \mathbf{0}, \mathbf{0}, 0.001, \mathbf{0}, 0.002, \ldots, 0.122, \mathbf{0}, \mathbf{0}]$
(all species)

- All non-negative!
- Values add up to 1!
- May contain zeros!

| | SPECIES | FRACTION |
|---|---|---|
| $x_1$ | Bacillus Thuringiensis | 0.003 |
| $x_2$ | Campylobacter jejuni | 0.001 |
| $x_3$ | Staphylococcus epidermidis | 0.002 |
| | … | … |
| | … | … |
| | … | … |
| | … | … |
| | … | … |
| | … | … |
| | … | … |
| | … | … |
| | … | … |
| | … | … |
| | … | … |
| | … | … |
| | … | … |
| | … | … |
| | … | … |
| | Methanobacterium extroquens | 0.081 |
| | Bordetella bronchiseptica | 0.078 |
| $x_n$ | Propionibacterium propionicus | 0.122 |

# The Problems

- High-dimensionality – our vectors are long
- Compositionality – if one species goes up, the rest go down
  - Non-normality - lack of a good "anchor" species/sample to fix compositionality
- Sparsity – lots of 0s in our data
  - Because species does not exist
  - Because species exists in small amounts and we undersampled
  - Because there's too much of the other species
  - Because of an error in the process
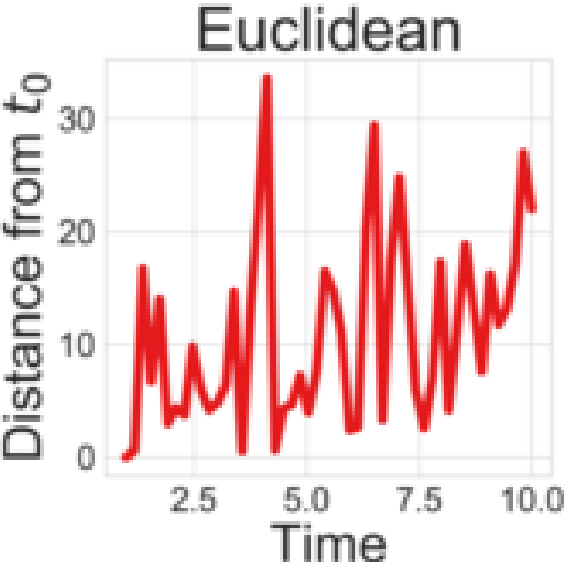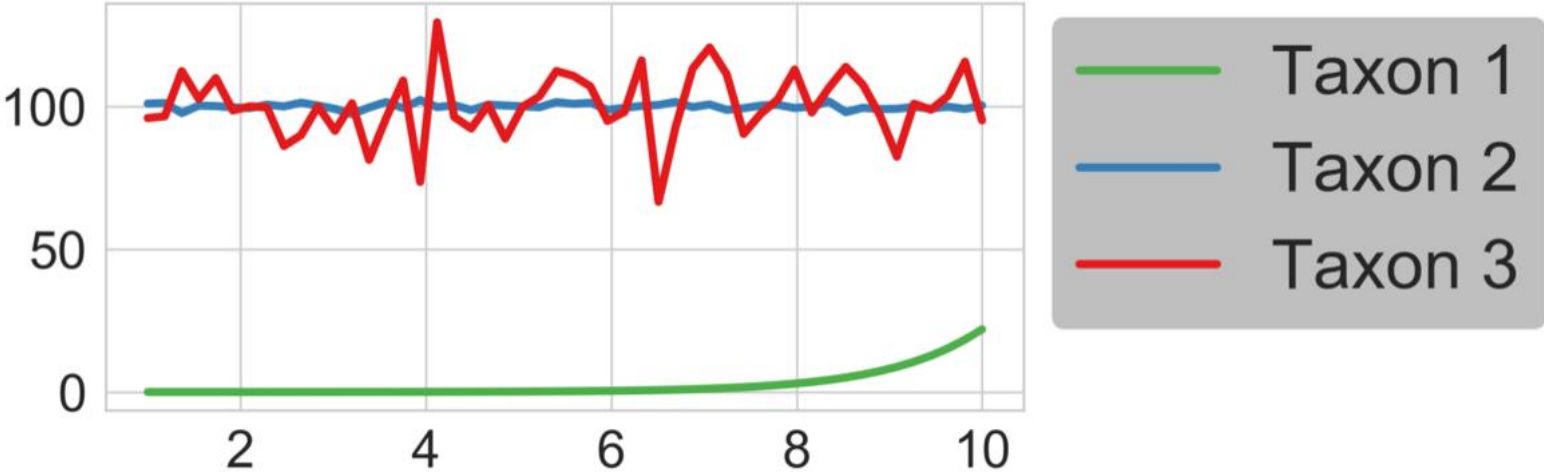- Some species are very impactful even in small amounts

# Example



Taxon 1 starts very low but increases exponentially over time

Taxon 2 starts high and remains there, stable

Taxon 3 starts high and remains there, but very unstable

We want to notice the trend of Taxon 1 increasing!

# Example

# Example

# Example

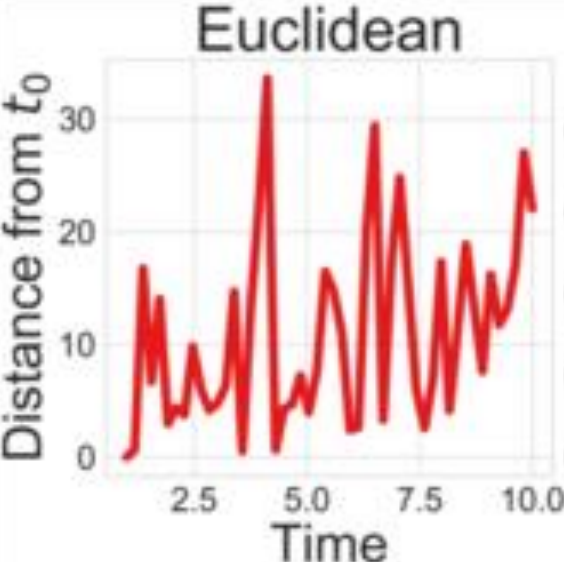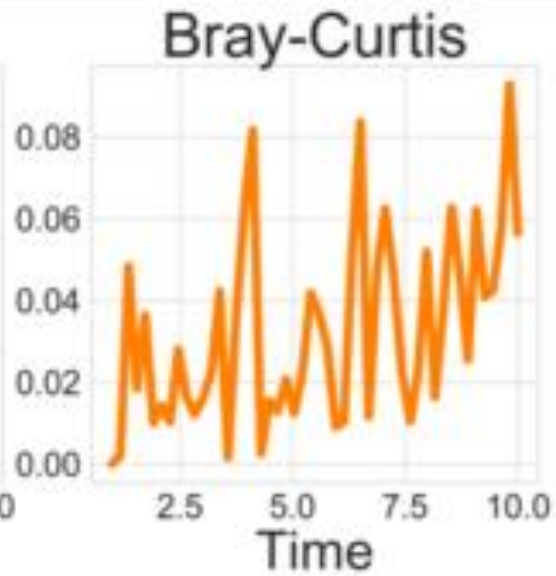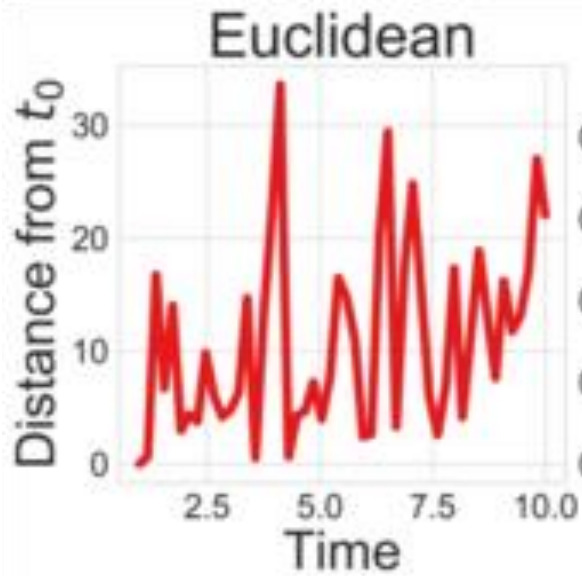**Euclidean**

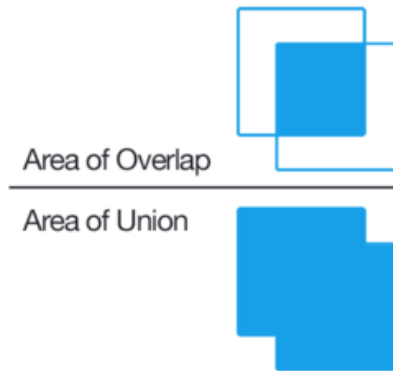$$\sqrt{\sum_{k=1}^{d}(x_k - y_k)^2}$$

(L2 norm of the delta)

**Bray-Curtis**

$$1 - \sum_{k=1}^{d}\min(x_k, y_k)$$

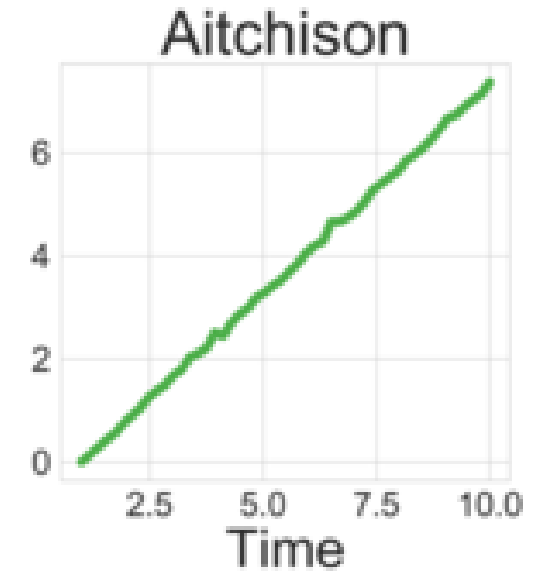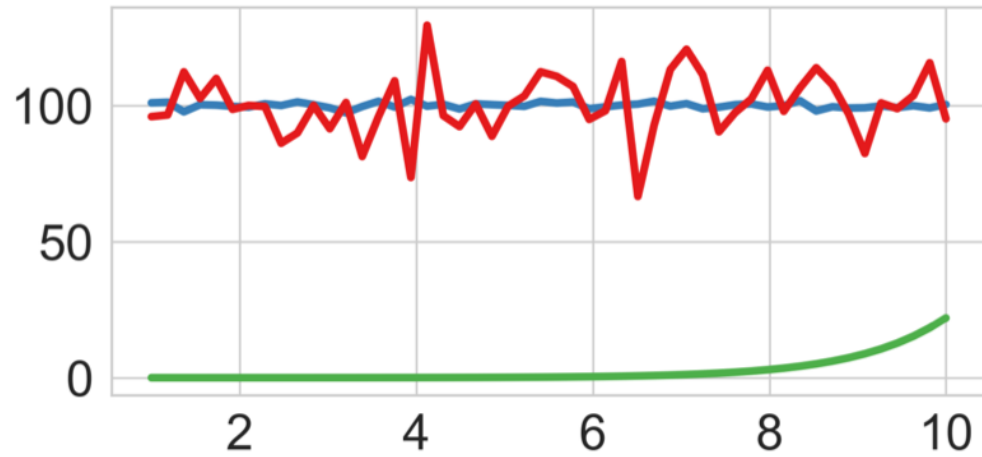(sum of lesser values for each species, assuming unit vectors. $\sum x_k = 1$)

**Jaccard**

(transforms data into binary yes/no)

$$1 - \frac{\sum_{k=1}^{d}\delta_{x_k=y_k}}{d}$$

$$\frac{\text{Area of Overlap}}{\text{Area of Union}}$$
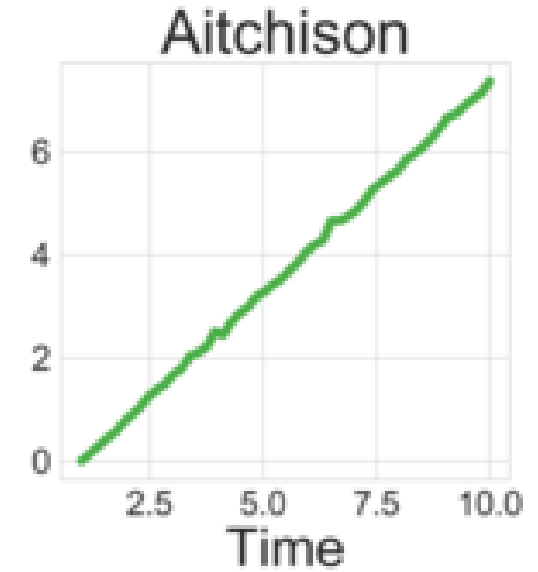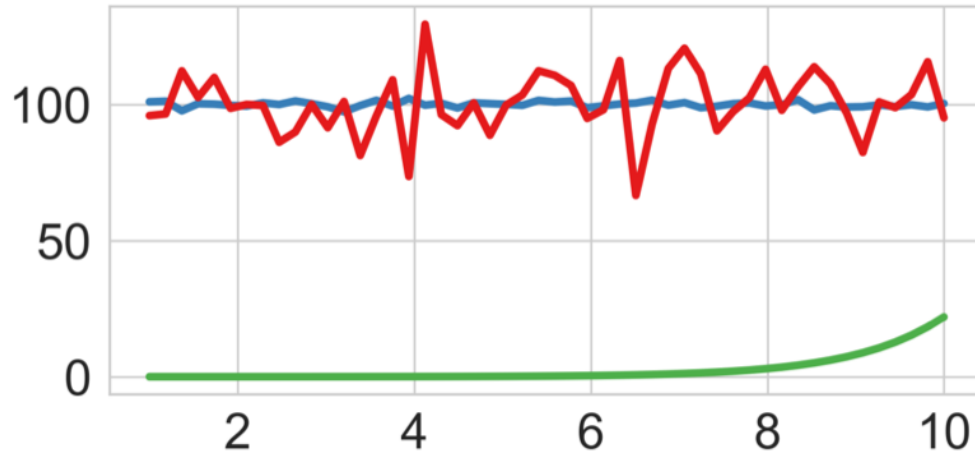
**Aitchison**

$$\sqrt{\sum_{k=1}^{d}\left(\log\frac{x_k}{g(x)} - \log\frac{y_k}{g(y)}\right)^2}$$

$$g(x) = \sqrt[d]{\prod_{k=1}^{d}x_k}$$

Aitchison

$$\sqrt{\sum_{k=1}^{d}\left(\log\frac{x_k}{g(x)}-\log\frac{y_k}{g(y)}\right)^2}$$

$$g(x) = \sqrt[d]{\prod_{k=1}^{d}x_k}$$

Aitchison



- Scale invariant!
- Perturbation invariant!
- Permutation invariant!
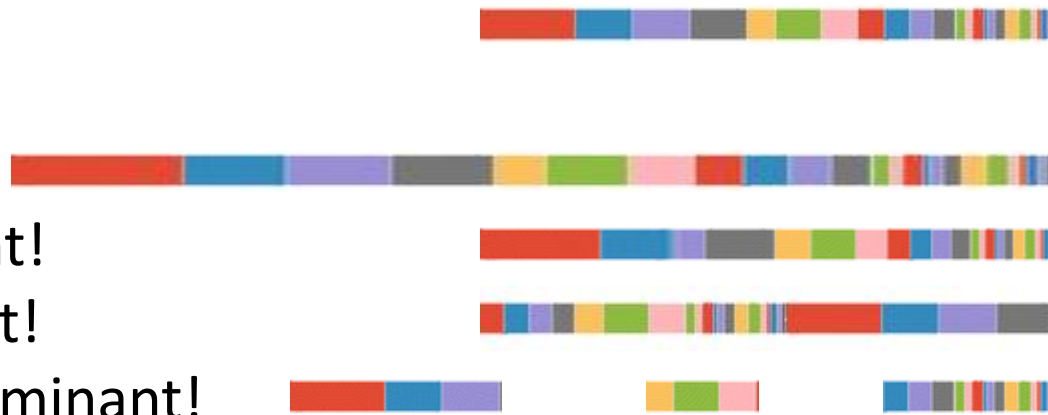- Subcompositional dominant!

$$\sqrt{\sum_{k=1}^{d}\left(\log\frac{x_k}{g(x)} - \log\frac{y_k}{g(y)}\right)^2}$$
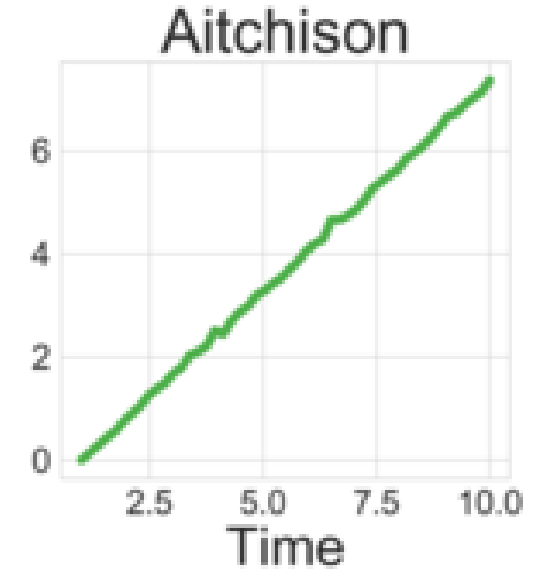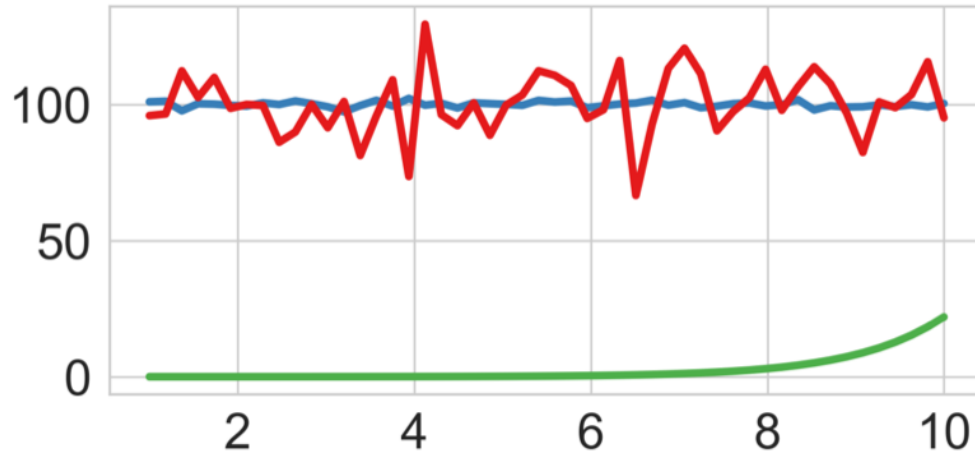
$$g(x) = \sqrt[d]{\prod_{k=1}^{d} x_k}$$

- Scale invariant!
- Perturbation invariant!
- Permutation invariant!
- Subcompositional dominant!

$$\sqrt{\sum_{k=1}^{d}\left(\log\frac{x_k}{g(x)}-\log\frac{y_k}{g(y)}\right)^2}$$

$$g(x)=\sqrt[d]{\prod_{k=1}^{d}x_k}$$

→ Fits all four requirements suggested by…. Dr. John Aitchison… ⬜

# CLR transform – Centered Log Ratio transform

$$\text{distance} = \sqrt{\sum_{k=1}^{d} \left( \log\frac{x_k}{g(x)} - \log\frac{y_k}{g(y)} \right)^2}$$

$$g(x) = \sqrt[d]{\prod_{k=1}^{d} x_k} = e^{\left(\frac{1}{d}\sum_{k=1}^{d} \log x_k\right)} = e^{\overline{\log x}}$$

# CLR transform – Centered Log Ratio transform

$$\text{distance} = \sqrt{\sum_{k=1}^{d} \left( \log \frac{x_k}{g(x)} - \log \frac{y_k}{g(y)} \right)^2} = \sqrt{\sum_{k=1}^{d} \left( \text{clr}(x_k) - \text{clr}(y_k) \right)^2}$$

$$g(x) = \sqrt[d]{\prod_{k=1}^{d} x_k} = e^{\left( \frac{1}{d} \sum_{k=1}^{d} \log x_k \right)} = e^{\overline{\log x}}$$

$$\text{clr}(x_k) = \log \frac{x_k}{g(x)} = \log x_k - \log g(x) = \log x_k - \overline{\log x}$$

$$\text{clr}(x) = \log x - \overline{\log x}$$

# CLR transform – Centered Log Ratio transform

$$\text{distance} = \sqrt{\sum_{k=1}^{d} \left( \log \frac{x_k}{g(x)} - \log \frac{y_k}{g(y)} \right)^2} = \sqrt{\sum_{k=1}^{d} \left( \text{clr}(x_k) - \text{clr}(y_k) \right)^2}$$

$$g(x) = \sqrt[d]{\prod_{k=1}^{d} x_k} = e^{\left( \frac{1}{d} \sum_{k=1}^{d} \log x_k \right)} = e^{\overline{\log x}}$$



$$\text{clr}(x_k) = \log \frac{x_k}{g(x)} = \log x_k - \log g(x) = \log x_k - \overline{\log x}$$

$$\text{clr}(x) = \log x - \overline{\log x}$$

# But wait, what about zeros...?

$$\log(0) = \textcolor{red}{\texttt{undefined}}$$

- just ignore them

- ignore them when calculating geometric mean / logarithmic average

- ignore them when calculating CLR

- call it "Robust CLR" – rCLR

$$g_r(x) = \sqrt[|\Omega_x|]{\prod_{k \in \Omega_x}^{|\Omega_x|} x_k} \qquad \text{rclr}(x_k) = \log\frac{x_k}{g_r(x)} \text{ ... or } \texttt{undefined} \text{ if } x_k = 0$$

# …Okay, but…we still have $\mathtt{undefined}$s in our data

- Pretend that they're missing values!

- Reconstruct the vectors as if these species' data is missing

- This is called "**Matrix Completion**"

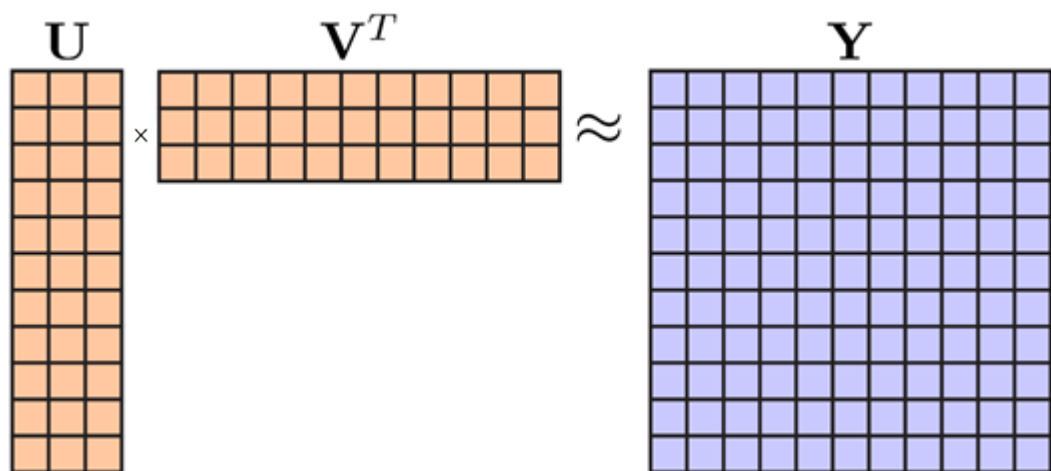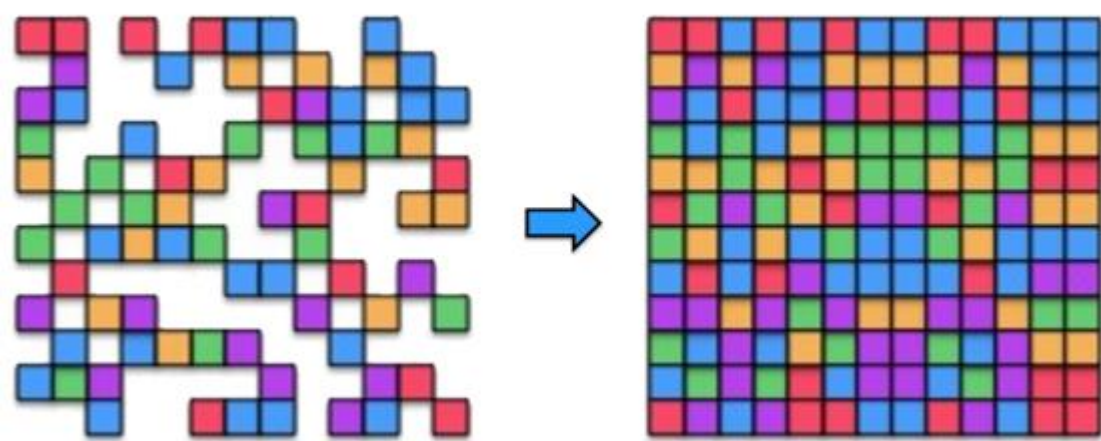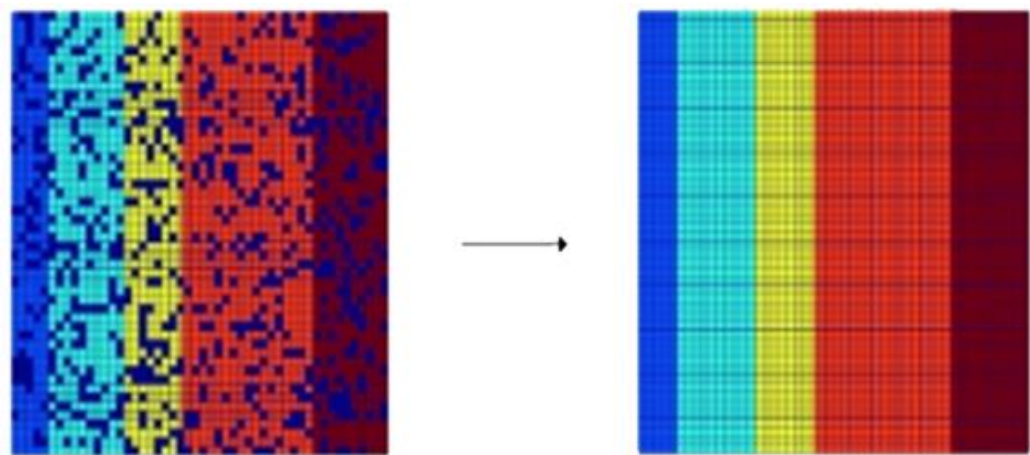(where each column in the matrix is a different sample vector)

$$\begin{pmatrix} 1 & ? & ? & 4 & ? \\ ? & 2 & 5 & ? & ? \\ ? & ? & 4 & 5 & ? \\ 5 & ? & ? & ? & 4 \end{pmatrix}$$

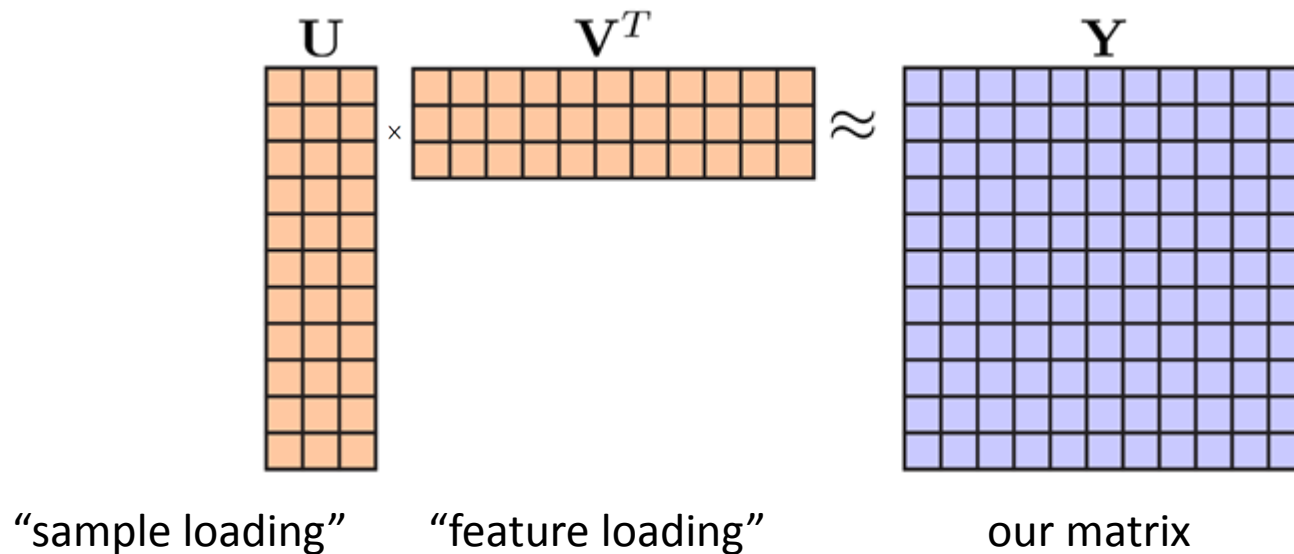|       | ![Matrix] | ![Matrix Revolutions] | ![Matrix Reloaded] | ![Metrics are coming] | ![The Tyranny of Metrics] | |
|-------|-----------|-----------------------|--------------------|-----------------------|---------------------------|---|
| Alice | 1 |   |   | 4 |   |   |
| Bob   |   | 2 | 5 |   |   |   |
| Carol |   |   | 4 | 5 |   |   |
| Dave  | 5 |   |   |   | 4 |   |
| ⋮     |   |   |   |   |   |   |

NETFLIX

$$\mathbf{U} \times \mathbf{V}^T \approx \mathbf{Y}$$

Users $\times$ Movies $\approx$ Ratings

U      $V^T$      Y

"sample loading"     "feature loading"     our matrix

In the paper, the matrix completion algorithm that was applied is called **OptSpace:**

$$\text{optimize for } \min_{U,V}(|Y - USV^T|_2^2) \quad \text{(not including new entries)}$$

# Summary

1. Transform our samples with the **Robust Centered Log Ratio transform**

2. Put all of the samples in a matrix, side by side

3. Run **Matrix Completion**

4. Result 1: a new complete matrix, representing what we "should see"

5. Result 2: two smaller matrices U and V, which represent important features and important samples

Now we can finally calculate distances between vectors in this matrix!

But instead of doing that, we'll do Ordination – Principal Component Analysis.

# PCA - Principal Component Analysis

- PCA is when you try to find order in your chaos
- You search for components (e.g. groups of species) that explain the data, name them PC1, PC2, PC3, etc. in order of how much they each explain variance in the data, and then make a graph that shows how right you are
- Usually, only PC1 and PC2 are shown, because 2D graphs are easy
- Each axis represents a component of the data that explains the difference, and has a percentage number showing how much it does
- We want percentages to be high in each axis, relative to what's left to explain (after using the other axes)

# Example

# Results

# Results

# Results



**Sequence Depth (*reads/sample*)**: 1000, 2000, 4000, 10000

**A** KL-Divergence

**B** PERMANOVA

**C** KNN Classification

**D** Positive Control Simulation

**E** Negative Control Simulation
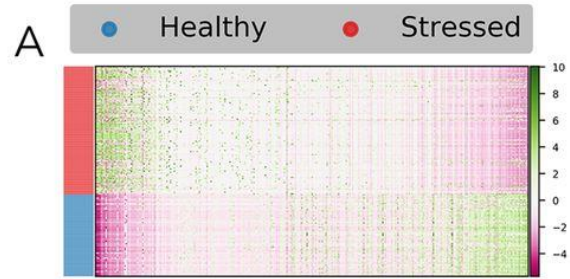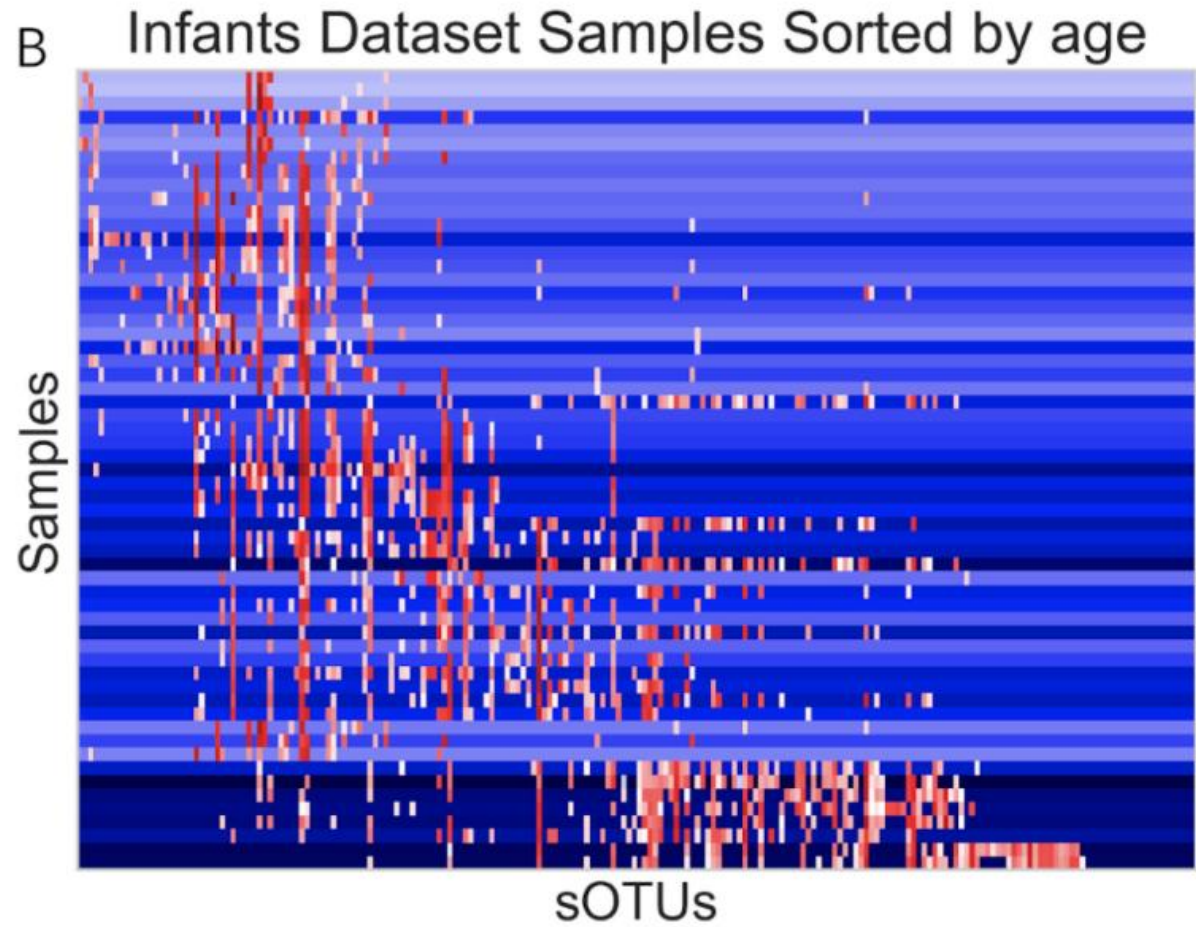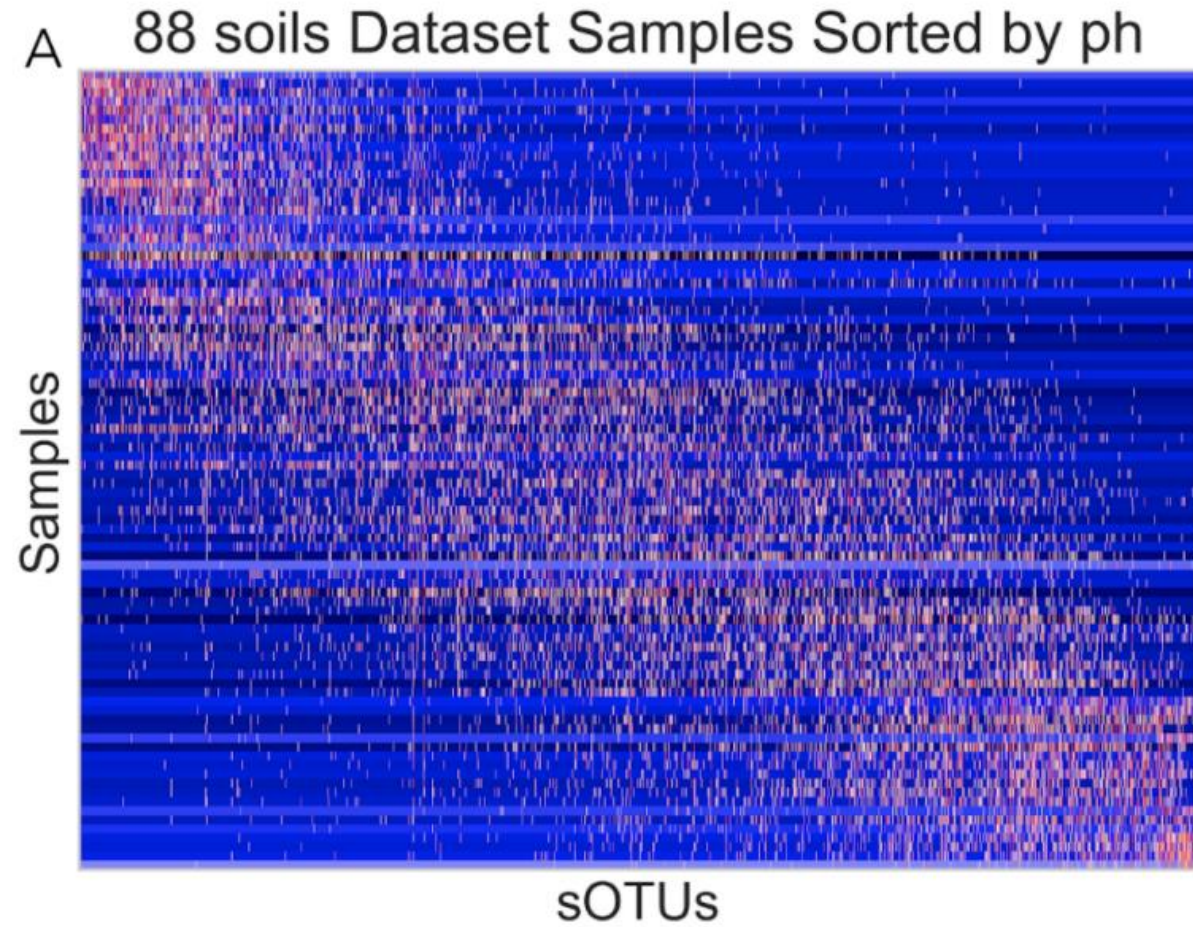
# Results

# Potential problems – high-rank matrices!



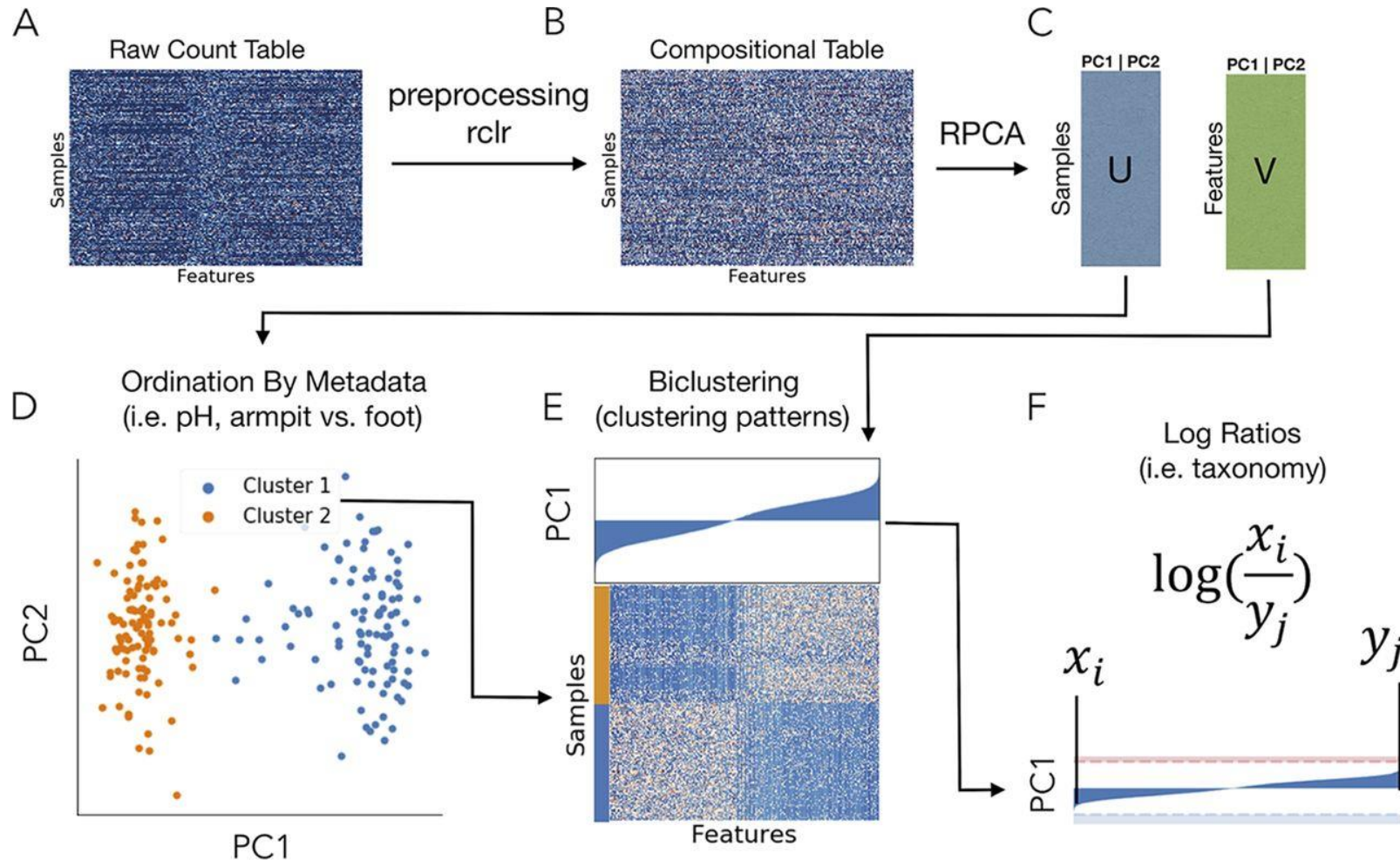Paper suggests to only use with rank 2-3, for ≤100 samples

# Conclusion

- Using Aitchison distance (with CLR transformation) is very good
- Using "robustness" and then matrix completion is pretty good on datasets where the rank is expected to be low
- Not good for "gradient-like" datasets, with a high-rank structure. It might cause misleading results
- Overfitting is possible but unlikely if we keep to low ranks (2-3)

(also the code for this is open-source, as usual)

# Discussion Points

- Is this usually useful? Do many data sets actually have a low rank?

- All of the results are on small datasets and subsets of datasets, with a small amount of samples. Does this still work with very large sample sizes? Is it computationally slow?

- Why does matrix completion work so well on zeros? It feels like dark magic.

# The End



(now you can finally understand this thing)