

You must choose, but choose wisely:

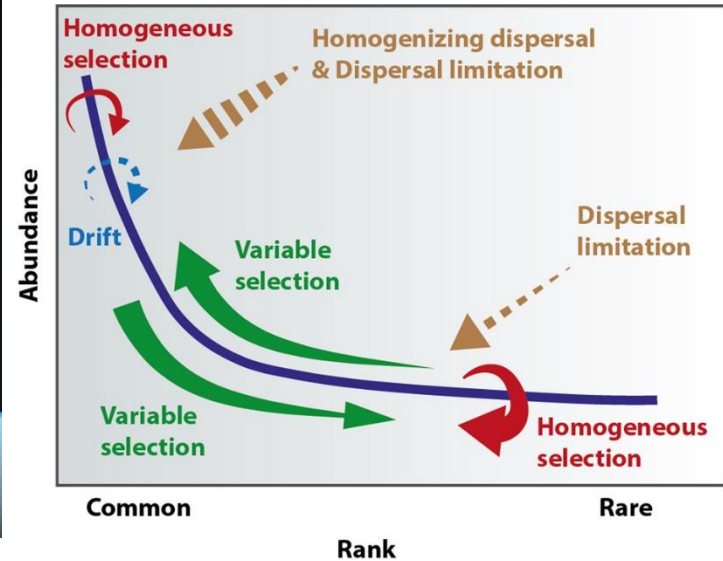
Model-based approaches for microbial community analysis and data integration

Marcio F. A. Leite

M.Leite@nioo.knaw.nl



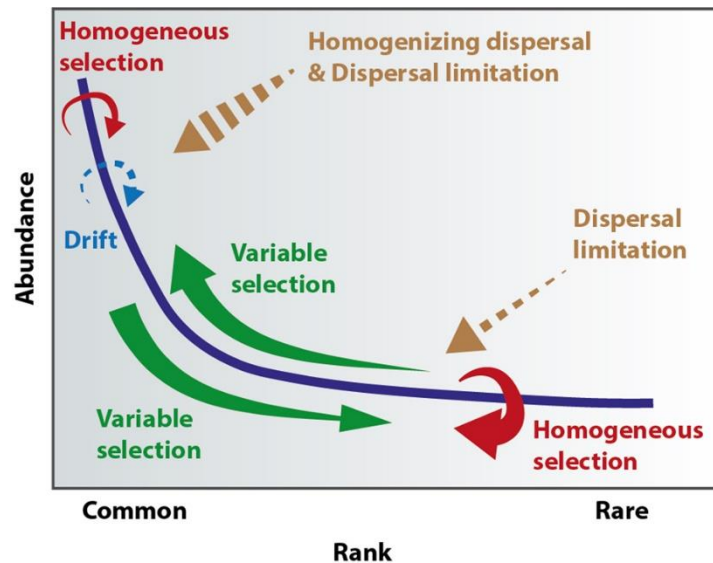
What we want/hope to know from metagenomics sequencing (MGS) data?



Jia, Dini-Andreote, Salles. *ISME COMMUN.* 2, 96 (2022).

What we want to know from sequencing data?

- Community composition
- Diversity
- Differential abundance
- Co-occurrence



Jia, Dini-Andreote, Salles. *ISME COMMUN.* 2, 96 (2022).

A still from the movie 'The Matrix' showing Keanu Reeves as Neo standing in a gun store. He is holding a handgun in his right hand. The background is filled with shelves of various firearms, including handguns and rifles, arranged in rows. The lighting is somewhat dim, and the overall tone is serious.

**We need
dissimilarity
metrics**

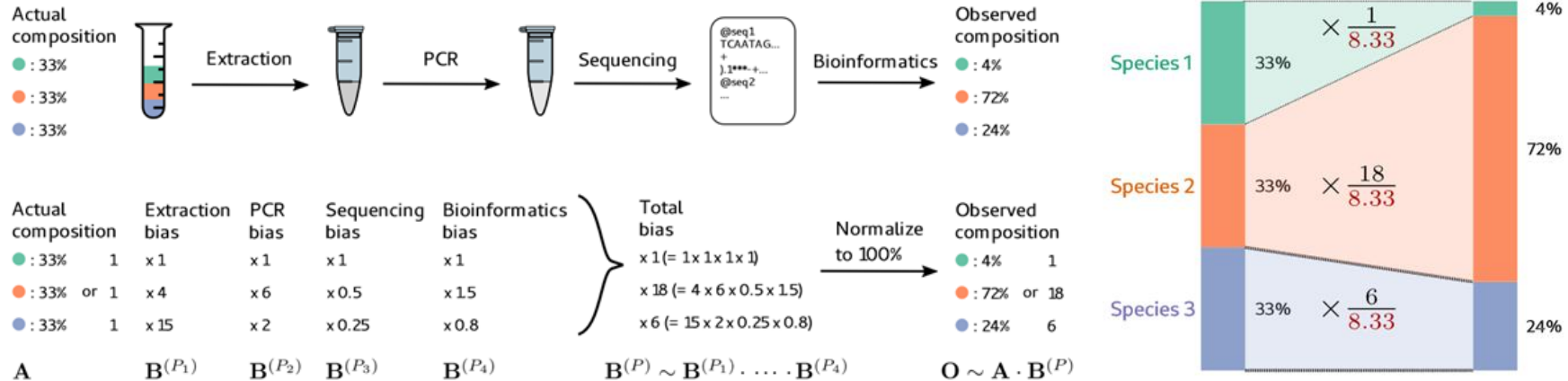
Lots of them!



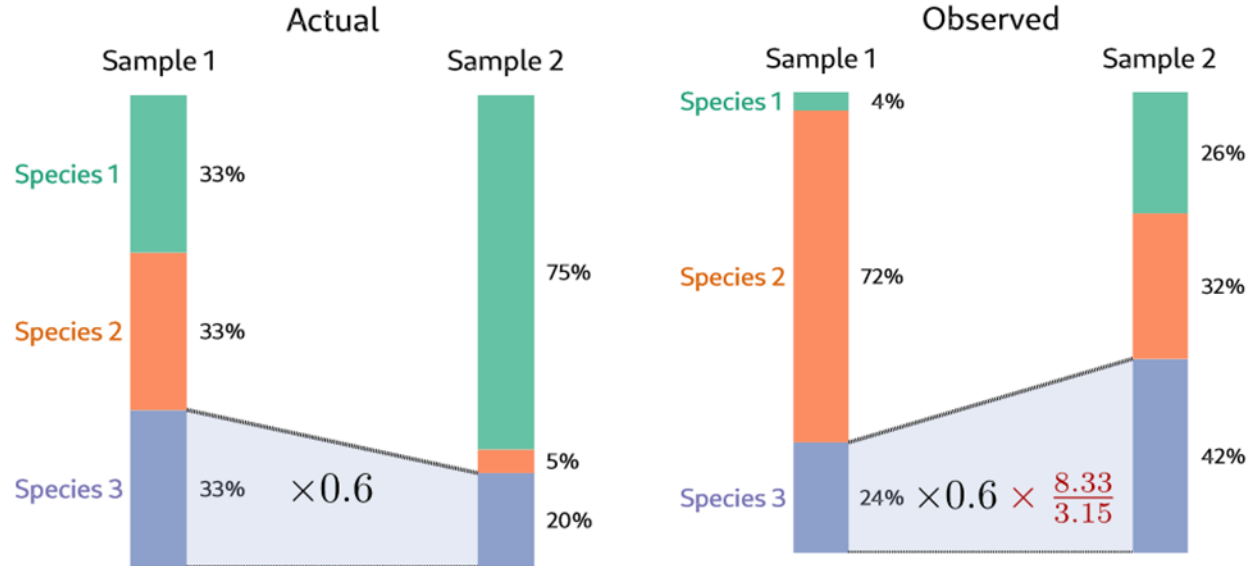
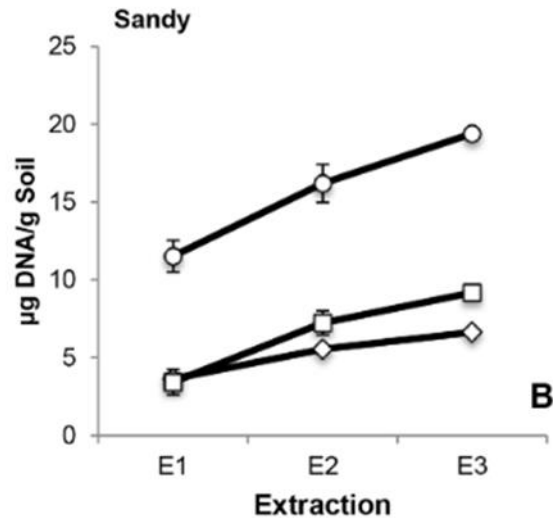
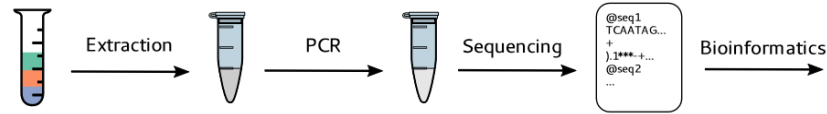
TSS TMM DESeq2 ZIP
CSS ANCOM-BC ANCOM
CLR ALDEx2 ZINB corncob

Data essence precedes data existence: unique challenges of microbiome data

MGS bias + Compositionality + sparsity + high variability (overdispersion)



MGS bias

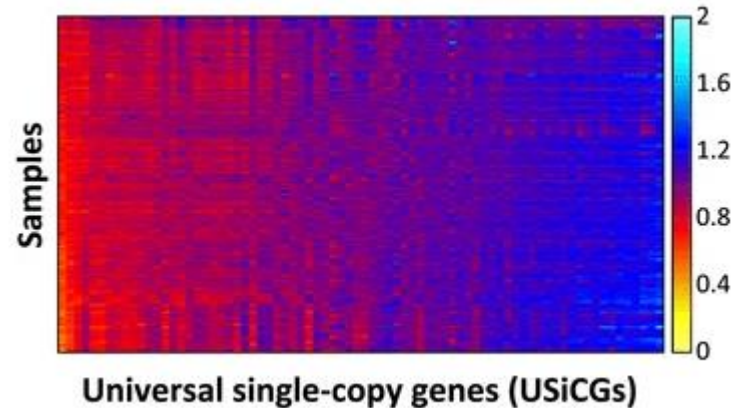
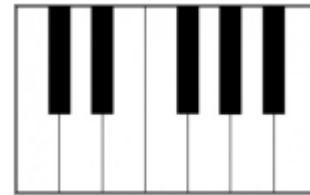


Diversity



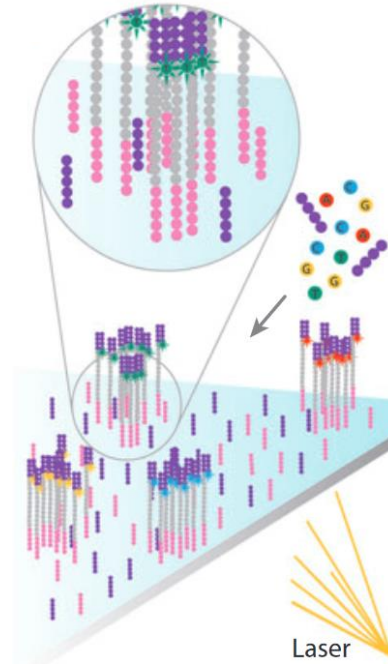
α -Diversity
+
 β -Diversity

MUSiCC



Compositionality

Sequencing is not counting.

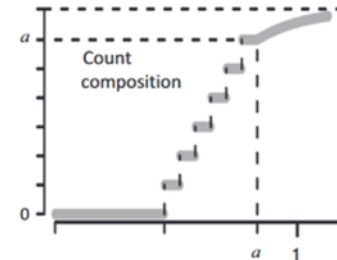


System Specifications

↔
540 Mb–15 Gb
OUTPUT RANGE

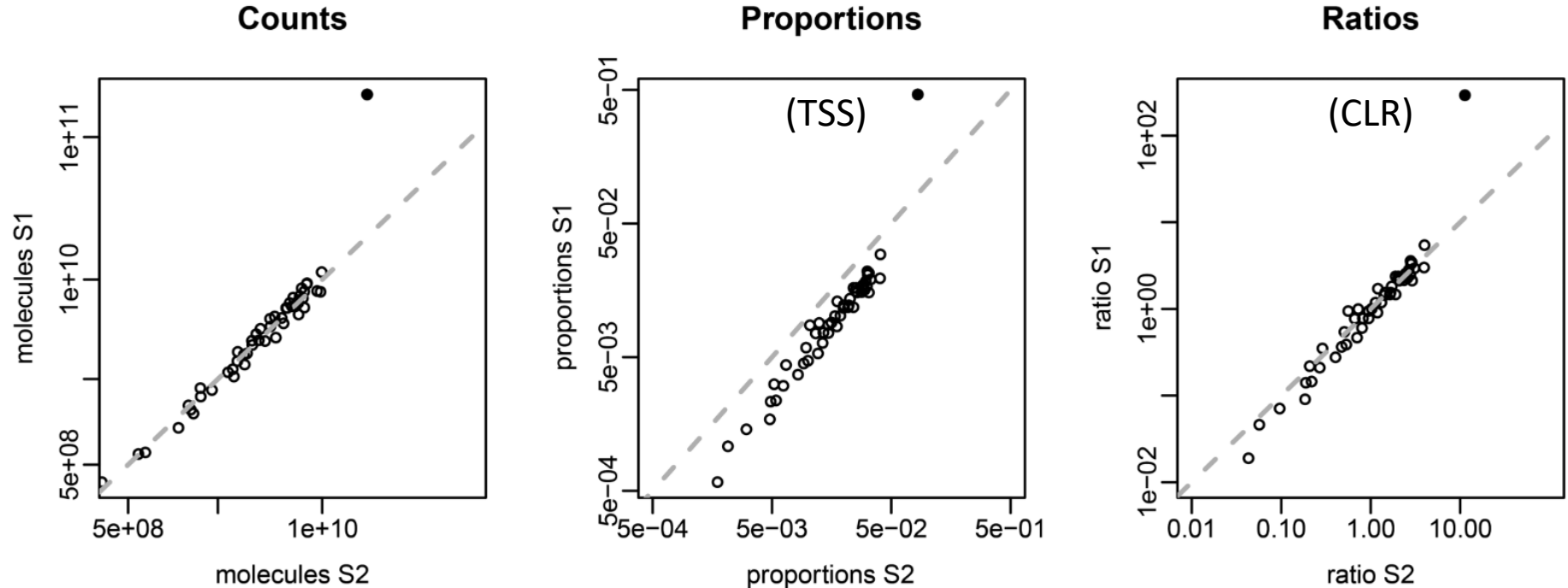
1-25 million
READS PER RUN

↗
2 × 300 bp
MAX READ LENGTH

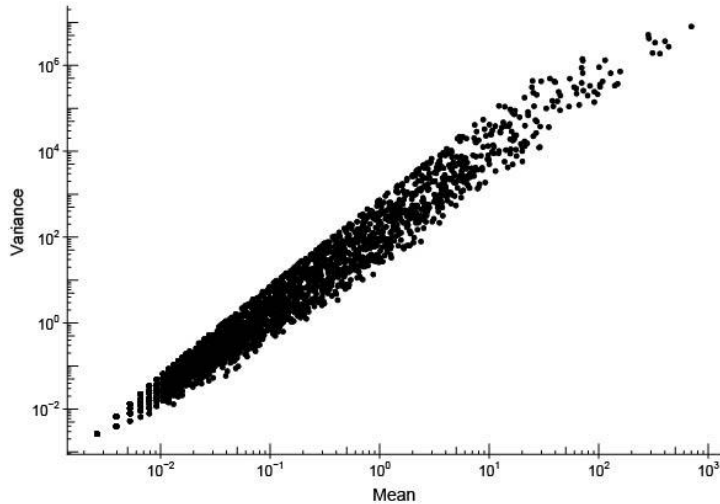


Compositionality

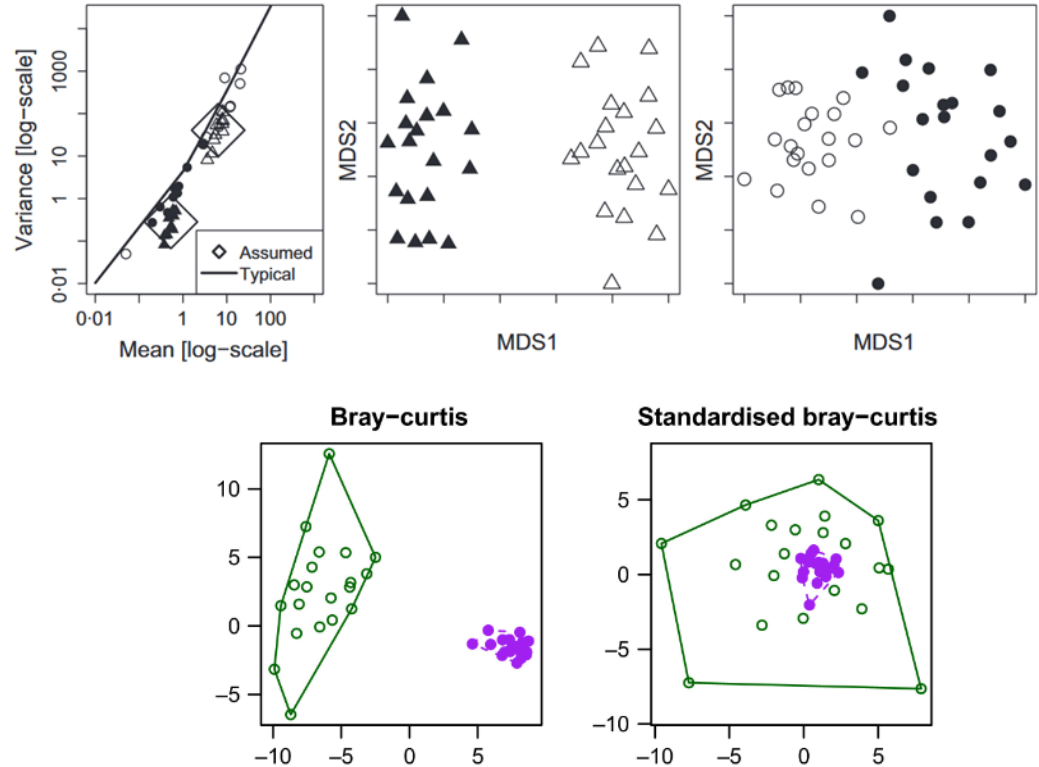
Relative abundance (proportions) is prone to errors.



Sparsity + Overdispersion



Leite & Kuramae. *SBB*. 2020.



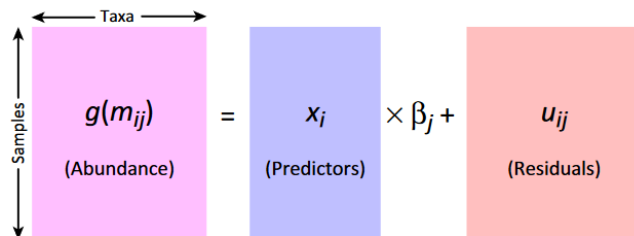
Warton et al. *Methods Ecol. Evol.* 2012.

Warton & Hui. *Methods Ecol. Evol.* 2017.


What we need... (the Plato's method)

- i. Investigate microbial community under the **compositional constraint**;
- ii. Determine the influence of **sequencing bias**;
- iii. Control data **sparsity and overdispersion**;
- iv. Account for the influence of **rare taxa**;
- v. **Integrate additional information** (e.g., traits, environmental covariates, experimental design);
- vi. Disentangle biotic interaction from environmental response.

So far... Model-based approaches

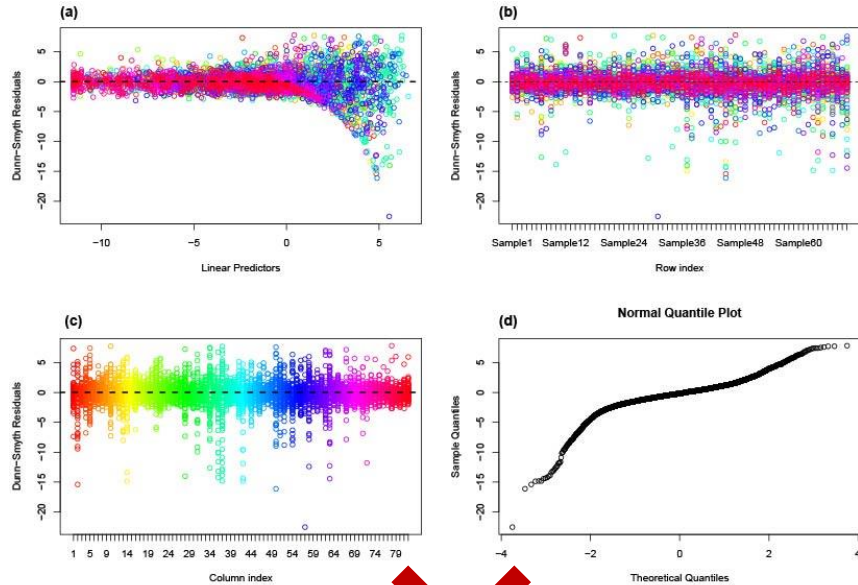


The diagram illustrates a model-based approach equation. It features three colored boxes: a pink box on the left labeled $g(m_{ij})$ (Abundance), a blue box in the middle labeled x_i (Predictors), and a red box on the right labeled u_{ij} (Residuals). The equation is $g(m_{ij}) = x_i \times \beta_j + u_{ij}$. Above the pink box, a double-headed arrow labeled 'Taxa' indicates the horizontal dimension. To the left of the pink box, a double-headed arrow labeled 'Samples' indicates the vertical dimension.

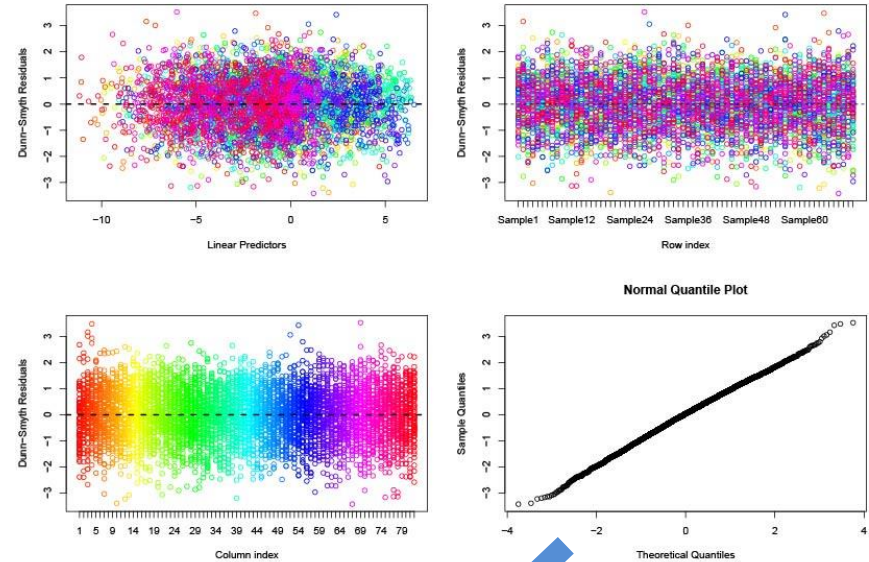
- 
- ✓ Compositional constraint;
 - ✓ MGS bias;
 - ✓ Overdispersion via mean-variance relationship;
 - ✓ Rare taxa;
 - ✓ Integrate additional information (e.g, traits, environmental covariates, experimental design);
 - ✓ Disentangling interaction.
 - ✓ Bonus: Unobserved/Latent variables

Assumptions first

Poisson Distribution

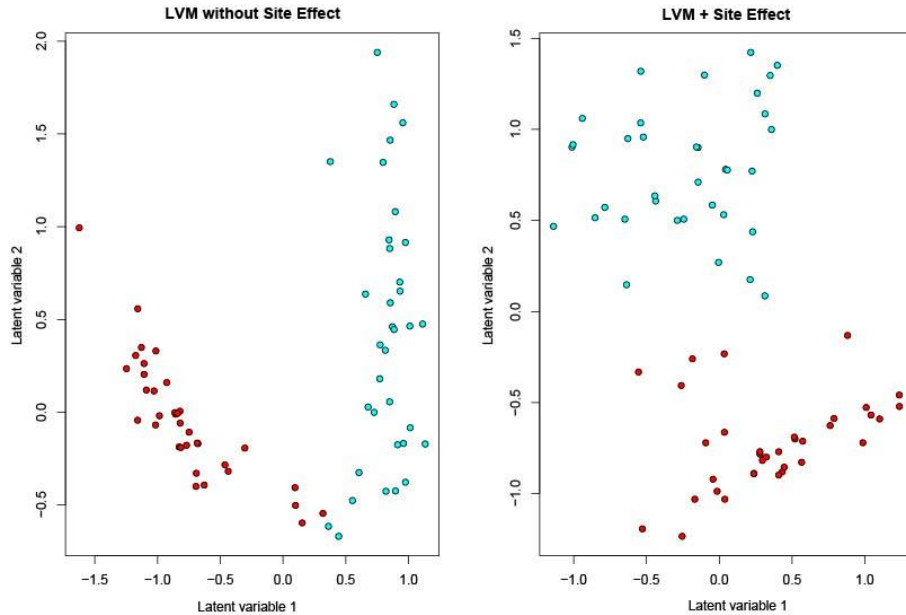


Negative Binomial Distribution



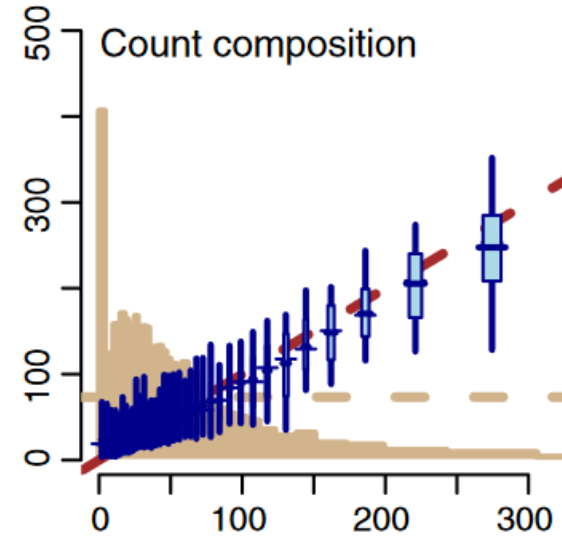
Assumptions first

Compositionality via site effects



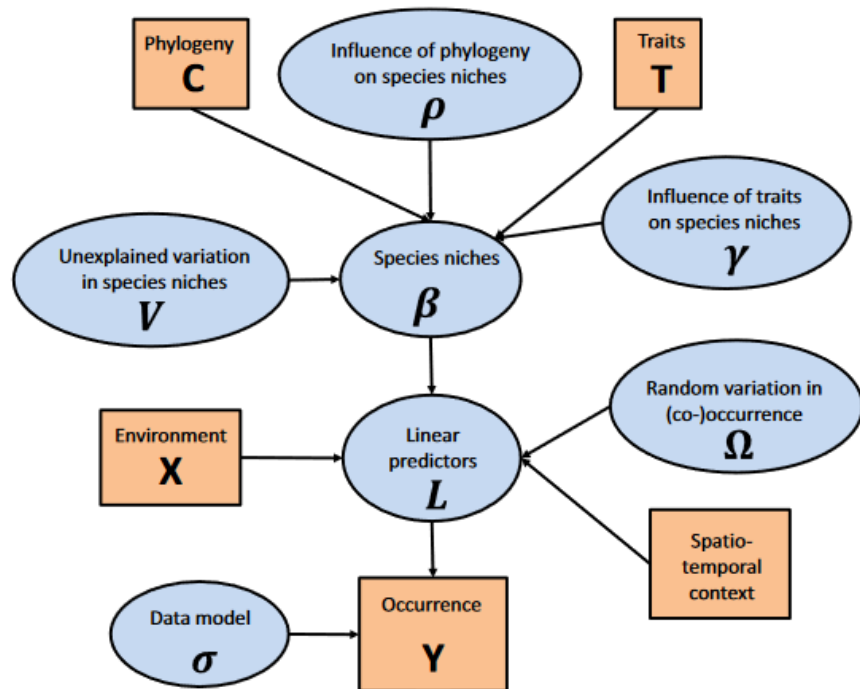
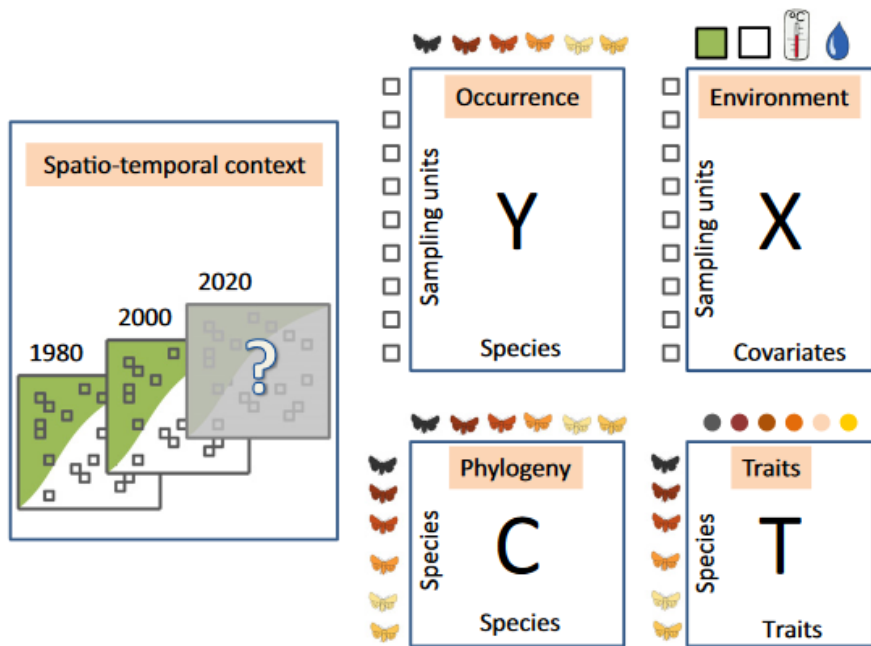
Leite & Kuramae. *SBB*. 2020.

Composition scale + Sampling effort

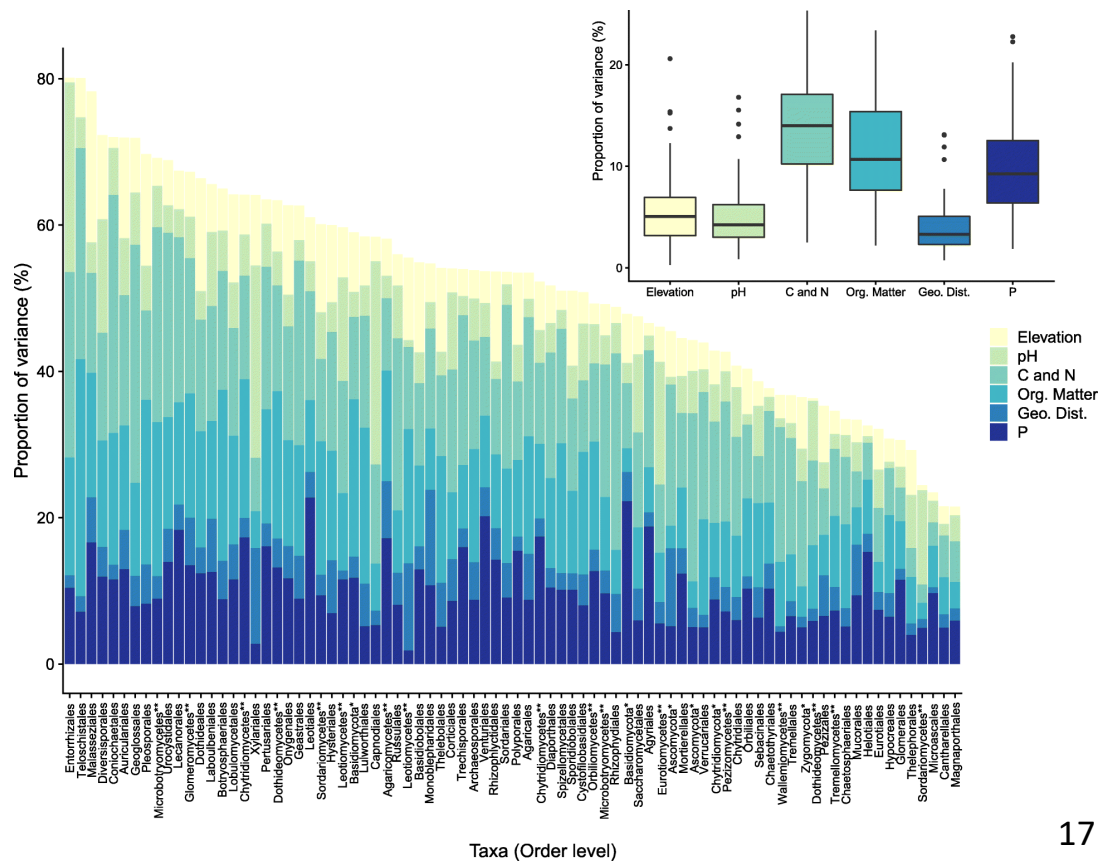
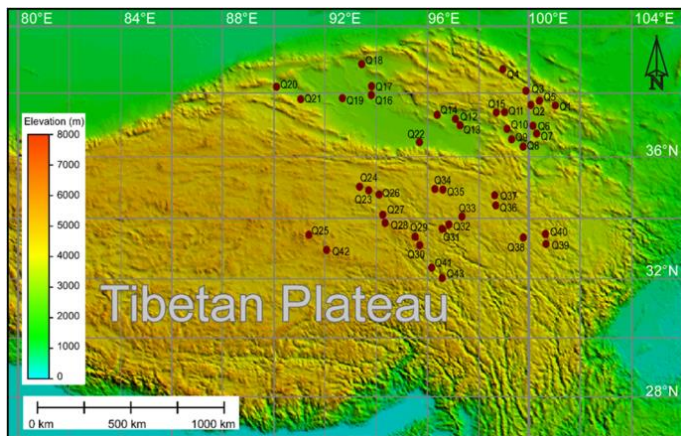


Clark et al. *Ecol. Monogr.* 2017.

Joint and hierarchical models



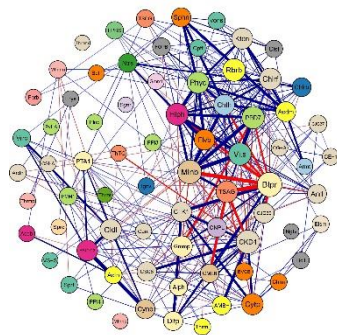
Joint and hierarchical models



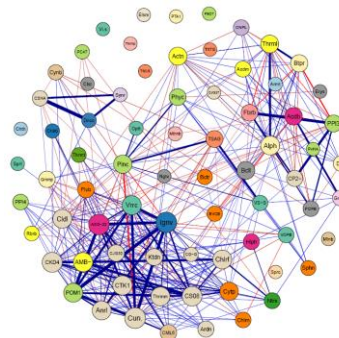
Disentangling correlations



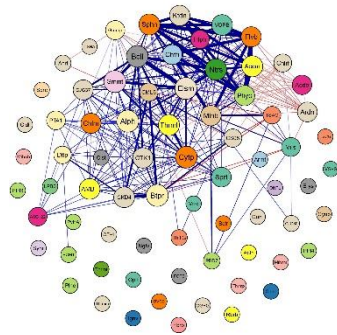
Plant-Soil driven interactions



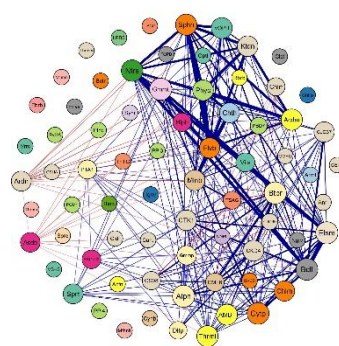
Plant-Soil driven interactions



Bacteria-Bacteria interactions



Bacteria-Bacteria interactions

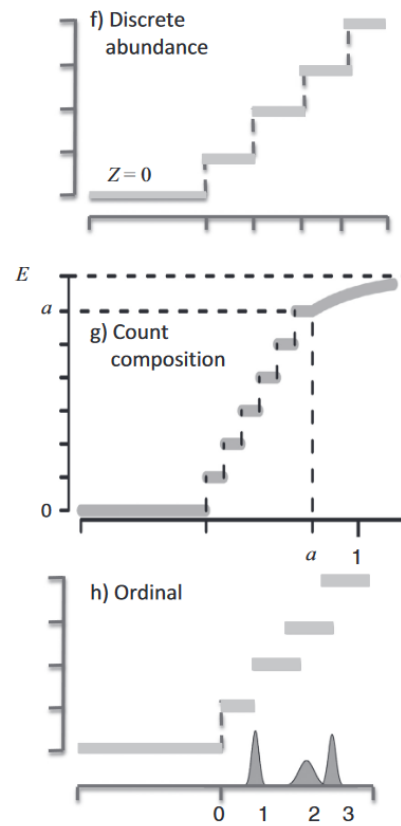
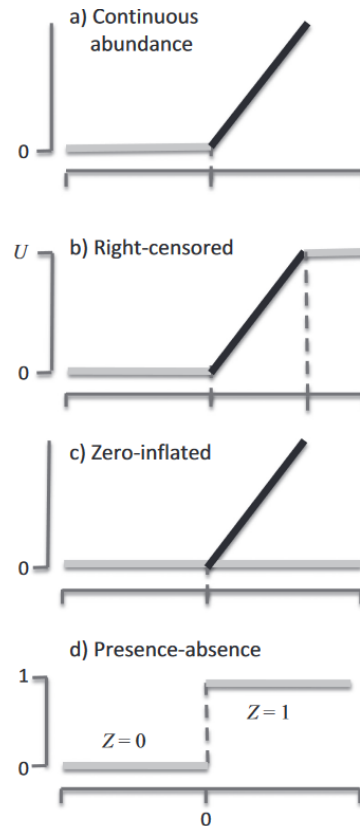


Data Integration

- How to integrate data that has different types?

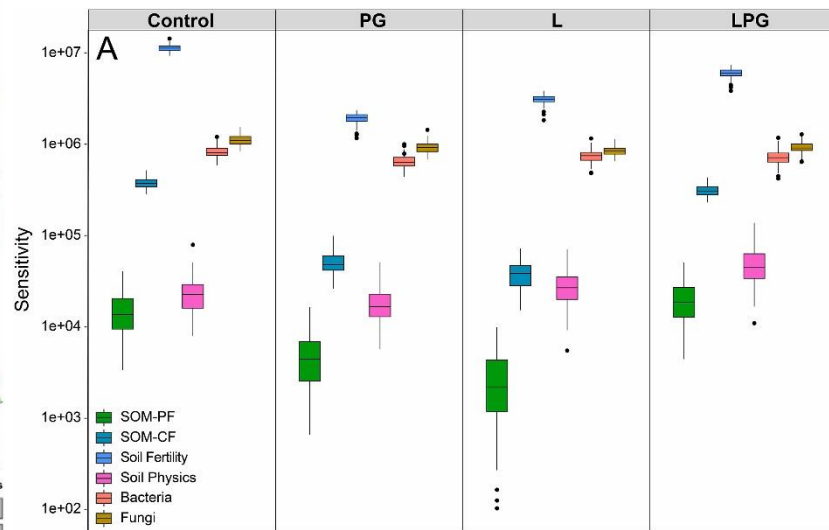
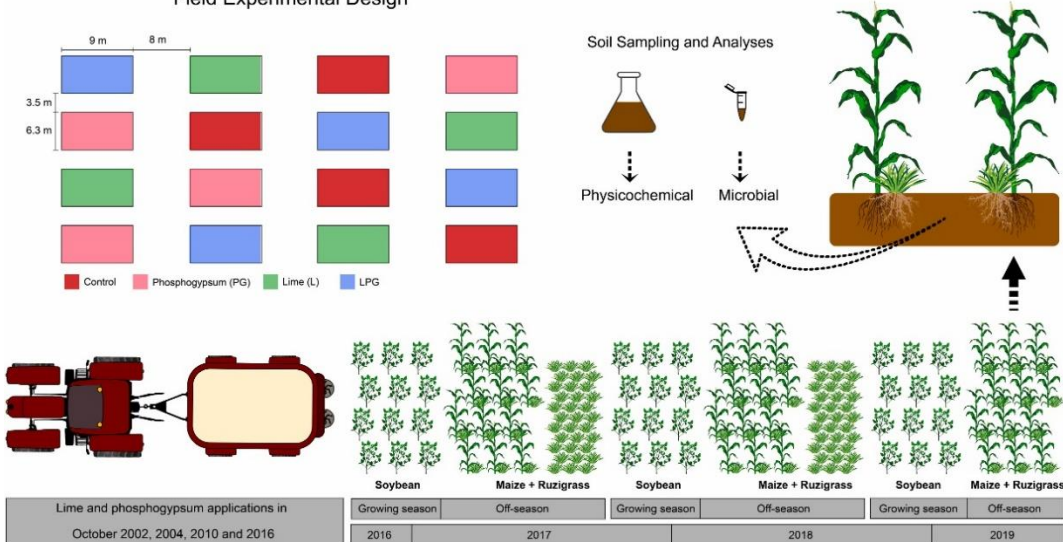


Generalized Joint
Attribute Model
(GJAM)

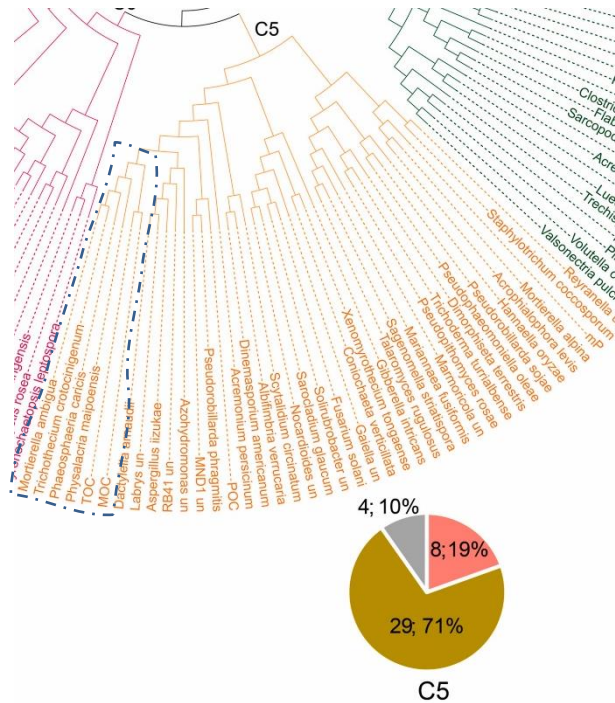
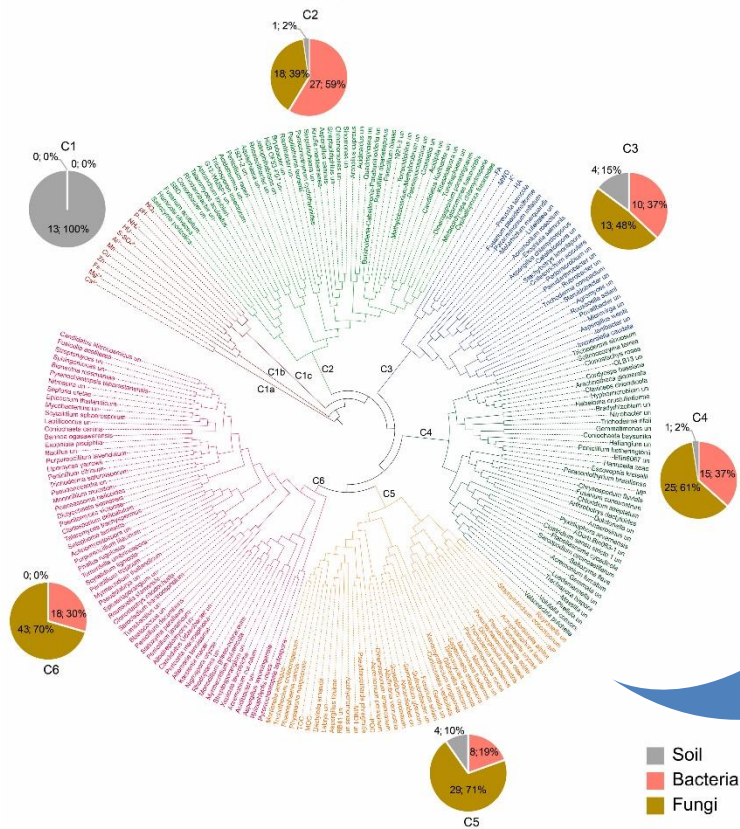


Data integration for agriculture

Field Experimental Design



Data integration for agriculture



Take-home message

- When facing a large arsenal of choices: assumption checking as the first step;
- Model-based approaches can explicitly account for key statistical properties of metagenomics data.

Thank you!



@MarcioLeiteEco



Leitemfa



M.Leite@nioo.knaw.nl

/DMG-symposium



NEDERLANDS INSTITUUT VOOR ECOLOGIE
NETHERLANDS INSTITUTE OF ECOLOGY



KURAMAE GROUP
SUSTAINABLE FOOD & FIBRE IN AFRICA