

Iterative Reference Selection for Normalizing Compositional Microbiome Data

Yiming Shi¹, Lili Liu¹, Jun Chen², Kristine M. Wylie³, Todd N. Wylie³, Sung Hee Park¹, Ruiwen Zhou¹, Yin Cao^{4, 5, 6},
Stephanie A. Fritz^{3, 7}, Molly J. Stout⁸, M. Cristina Vazquez Guillamet^{7, 9}, Lei Liu¹

¹ Institute for Informatics, Data Science & Biostatistics, Washington University in St. Louis; ² Mayo Clinic; ³ Department of Pediatrics, Washington University; ⁴ Department of Surgery, WUSM; ⁵ Siteman Cancer Center; ⁶ Gastroenterology, WUSM; ⁷ Infectious Diseases, WUSM; ⁸ University of Michigan; ⁹ Pulmonary and Critical Care Medicine, WUSM

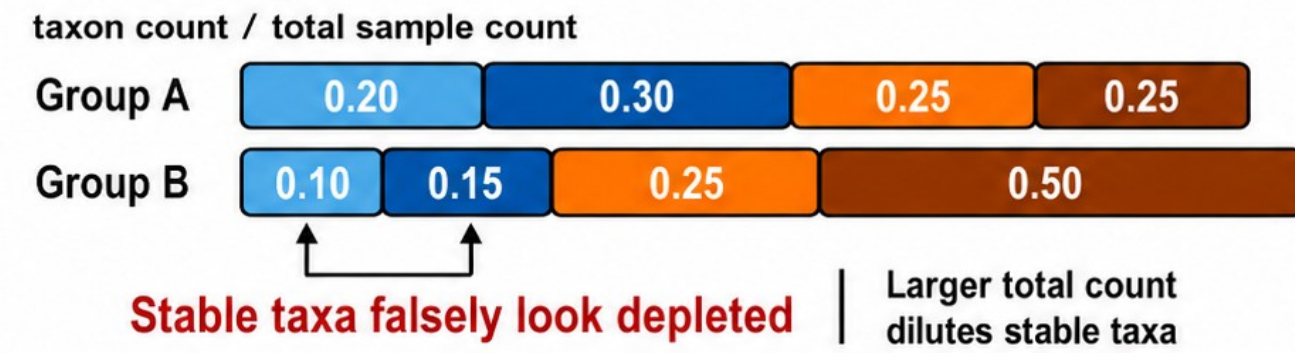
Background

Compositional nature is a fundamental characteristic of microbiome data: the abundances are often expressed as relative proportions due to unknown sampling fraction and variable sequencing depth.

Truth: absolute abundance



TSS artifact



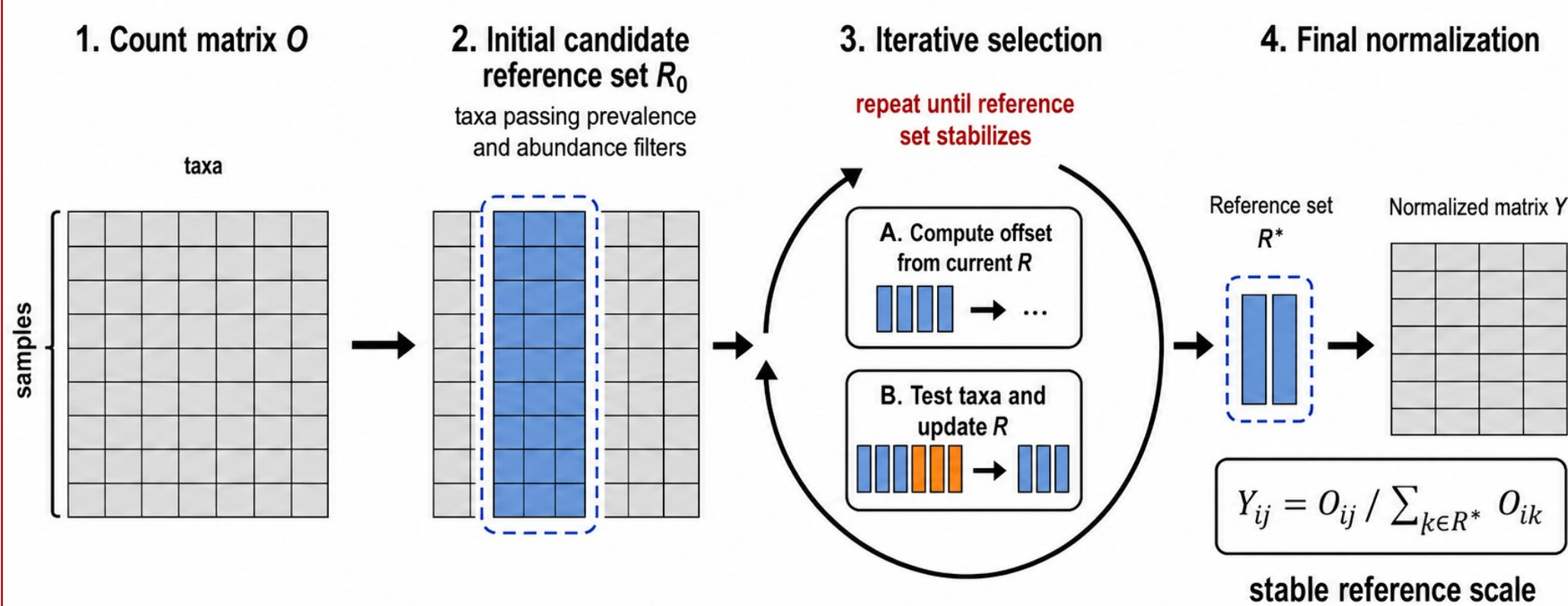
Reference normalization



Compositional normalization can make stable taxa appear differentially abundant. When total microbial load differs across samples, standard normalization can distort fold changes: stable taxa may appear depleted and differential taxa can induce spurious negative correlations.

IRS uses stable reference taxa to preserve true relative changes while reducing compositional artifacts.

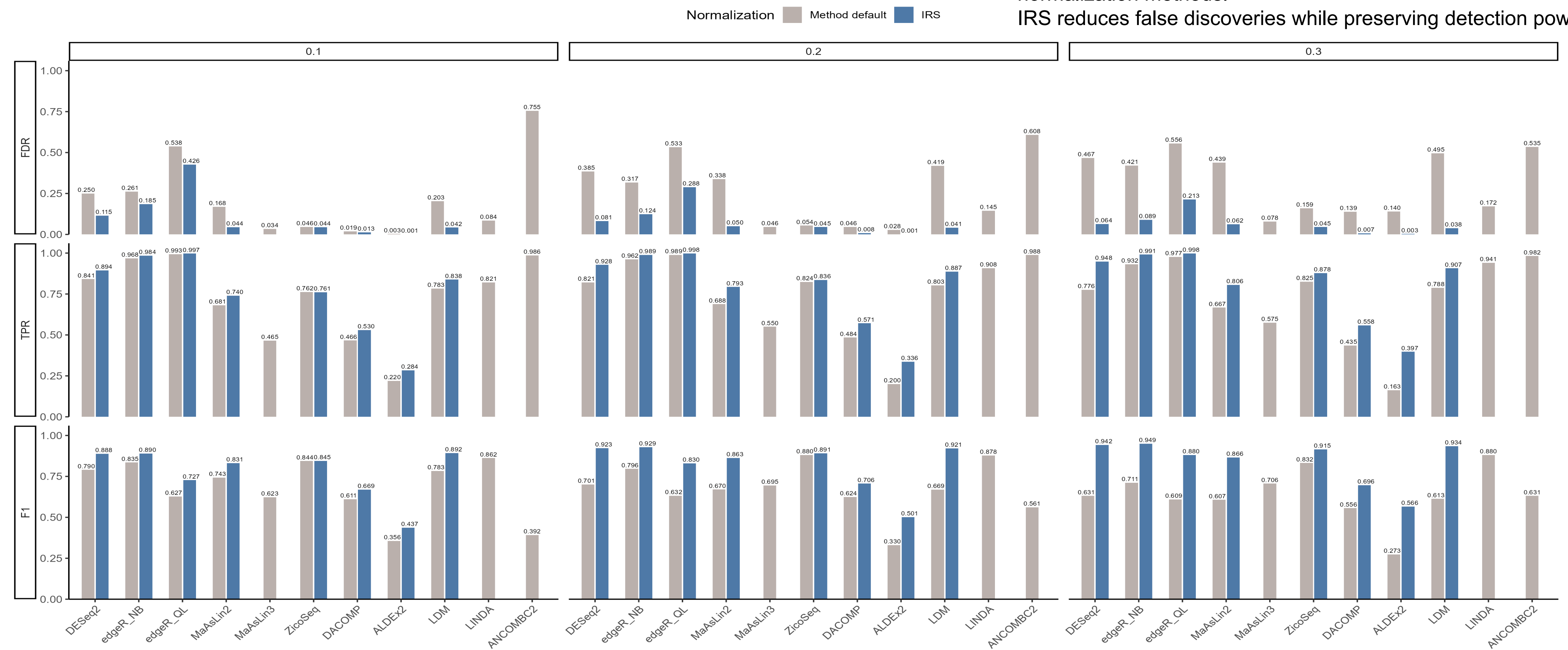
IRS: Iterative reference selection normalization



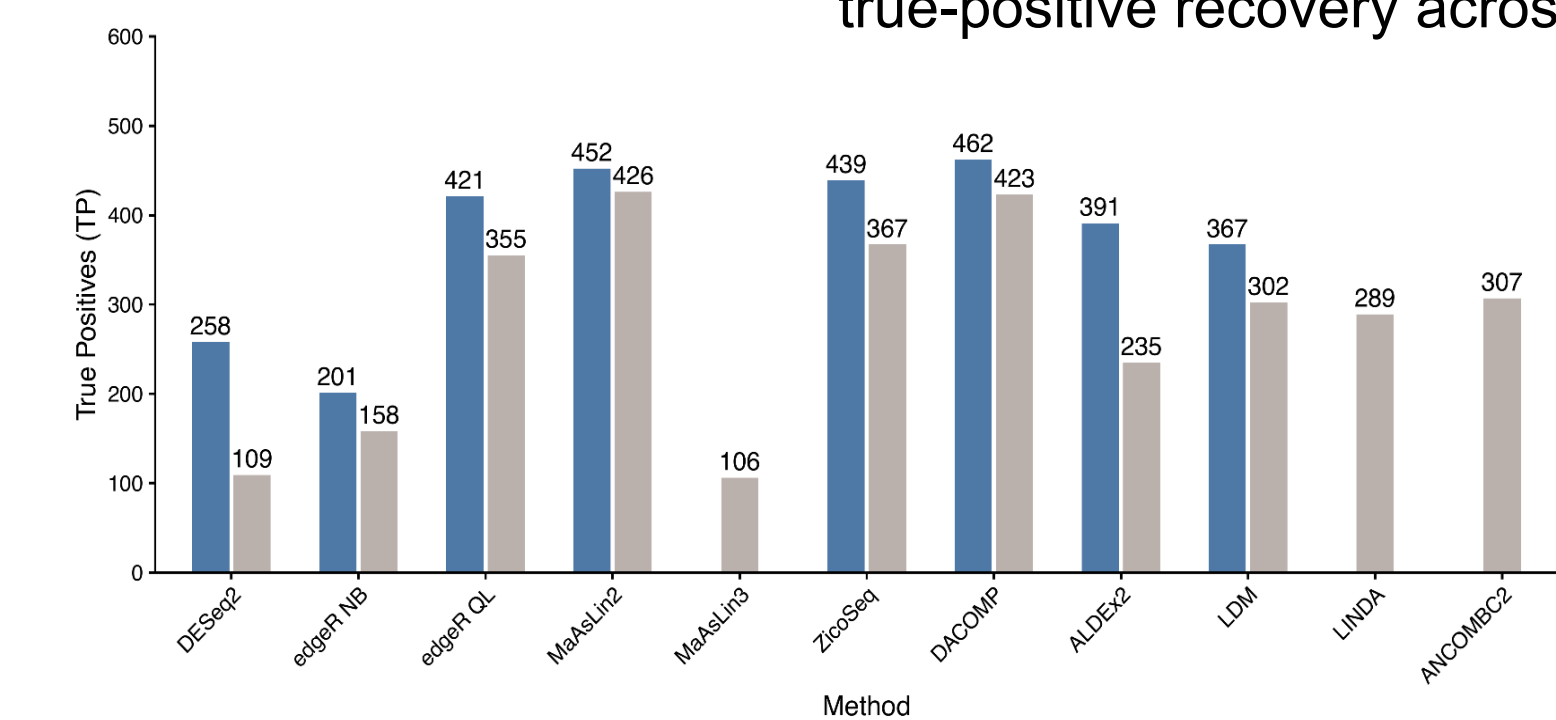
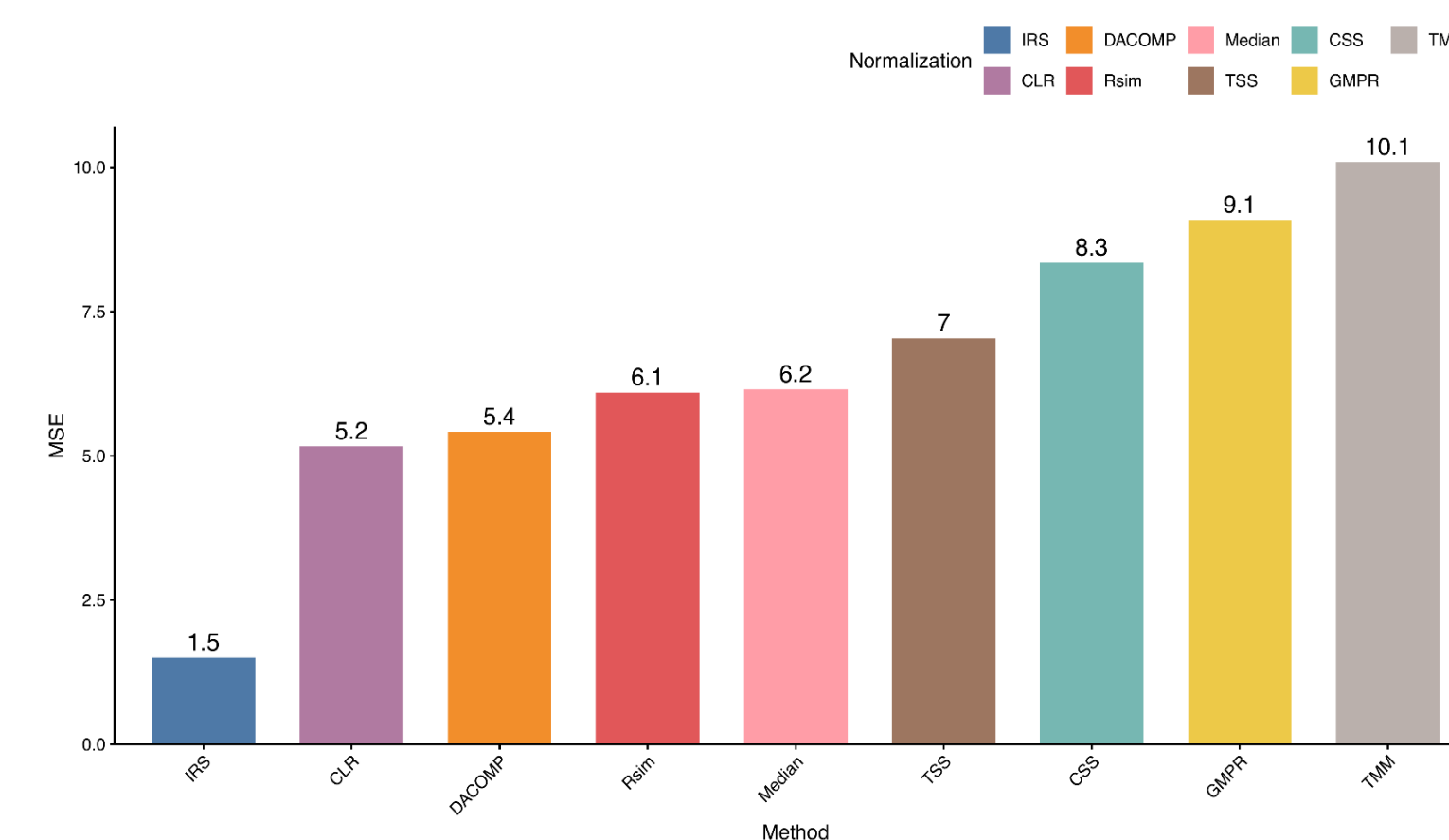
IRS iteratively selects stable taxa as an internal normalization reference.

Simulation studies

Simulated count data with known true signals benchmark normalization methods.
IRS reduces false discoveries while preserving detection power.



Real data validation



- The MetaCardis dataset with flow-cytometry microbial load was used to define an absolute-abundance truth set.
- IRS shows the lowest log₂ fold-change MSE and strong true-positive recovery across DA frameworks.

Conclusions & Future Directions

- Iterative Reference selection method efficiently identifies a large, clean set of **non-differential** taxa, enabling reference-based normalization that **mitigates negative correlation bias** in compositional microbiome analyses.
- Compared with existing normalization methods: More accurate size-factor estimates, smaller systematic bias, and a larger/cleaner reference set.**
- When used in differential abundance analysis models: **Lower FDR with equal or higher power** (higher TPR/F1) due to better compositional bias control.