

EVE 161: Microbial Phylogenomics

Class #16: Predicting Metagenomes

UC Davis, Winter 2018

Instructor: Jonathan Eisen

Teaching Assistant: Cassie Ettinger

Profiling phylogenetic marker genes, such as the 16S rRNA gene, is a key tool for studies of microbial communities but does not provide direct evidence of a community's functional capabilities. Here we describe PICRUSt (phylogenetic investigation of communities by reconstruction of unobserved states), a computational approach to predict the functional composition of a metagenome using marker gene data and a database of reference genomes. PICRUSt uses an extended ancestral-state reconstruction algorithm to predict which gene families are present and then combines gene families to estimate the composite metagenome. Using 16S information, PICRUSt recaptures key findings from the Human Microbiome Project and accurately predicts the abundance of gene families in host-associated and environmental communities, with quantifiable uncertainty. Our results demonstrate that phylogeny and function are sufficiently linked that this 'predictive metagenomic' approach should provide useful insights into the thousands of uncultivated microbial communities for which only marker gene surveys are currently available.

Profiling phylogenetic marker genes, such as the 16S rRNA gene, is a key tool for studies of microbial communities but does not provide direct evidence of a community's functional capabilities. Here we describe PICRUSt (phylogenetic investigation of communities by reconstruction of unobserved states), a computational approach to predict the functional composition of a metagenome using marker gene data and a database of reference genomes. PICRUSt uses an extended ancestral-state reconstruction algorithm to predict which gene families are present and then combines gene families to estimate the composite metagenome. Using 16S information, PICRUSt recaptures key findings from the Human Microbiome Project and accurately predicts the abundance of gene families in host-associated and environmental communities, with quantifiable uncertainty. Our results demonstrate that phylogeny and function are sufficiently linked that this 'predictive metagenomic' approach should provide useful insights into the thousands of uncultivated microbial communities for which only marker gene surveys are currently available.

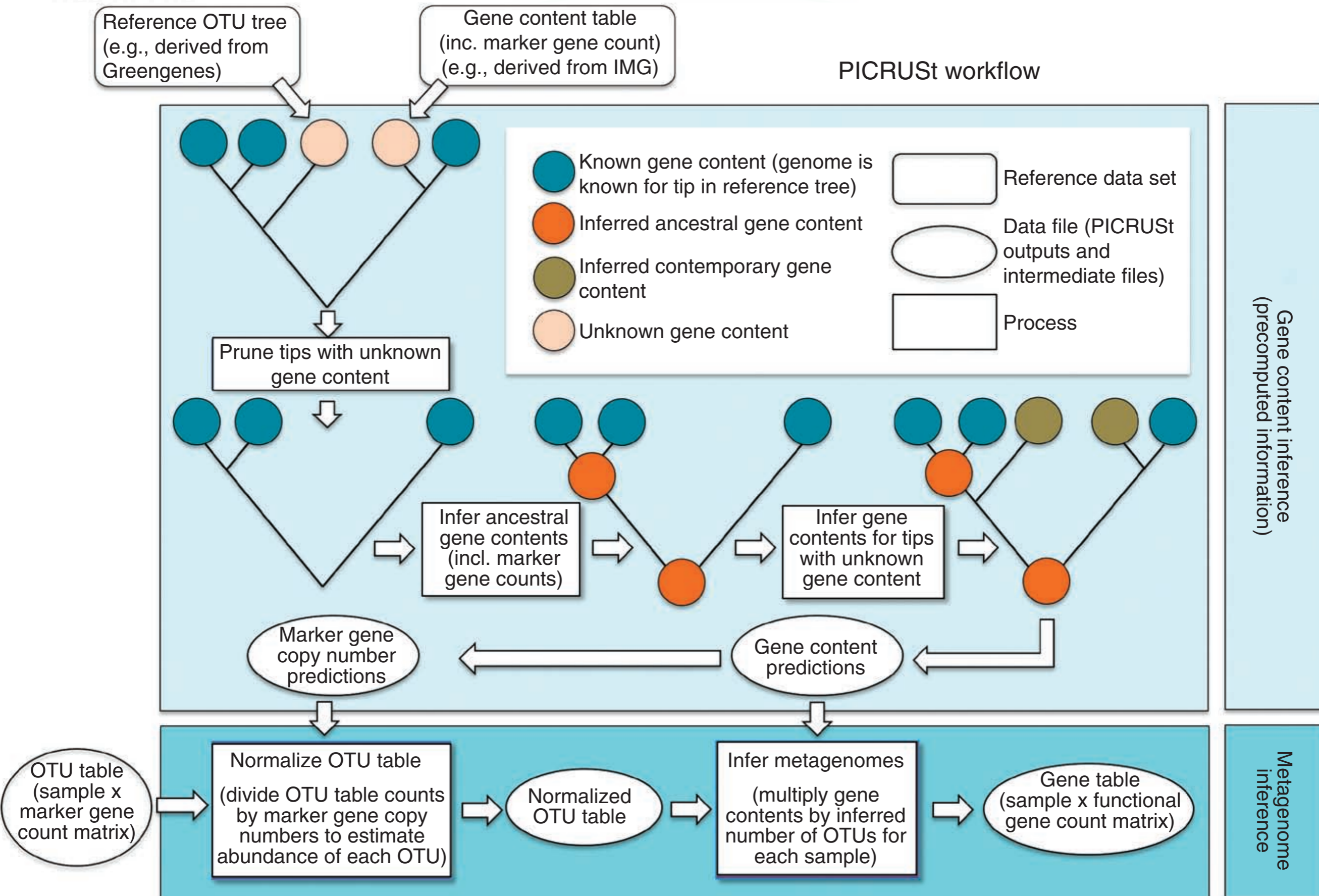


Figure 1 The PICRUSSt workflow. PICRUSSt is composed of two high-level workflows: gene content inference (top box) and metagenome inference (bottom box). Beginning with a reference OTU tree and a gene content table (i.e., counts of genes for reference OTUs with known gene content), the gene content inference workflow predicts gene content for each OTU with unknown gene content, including predictions of marker gene copy number. This information is precomputed for 16S based on Greengenes and IMG, but all functionality is accessible in PICRUSSt for use with other marker genes and reference genomes. The metagenome inference workflow takes an OTU table (i.e., counts of OTUs on a per sample basis), where OTU identifiers correspond to tips in the reference OTU tree, as well as the copy number of the marker gene in each OTU and the gene content of each OTU (as generated by the gene content inference workflow), and outputs a metagenome table (i.e., counts of gene families on a per-sample basis).

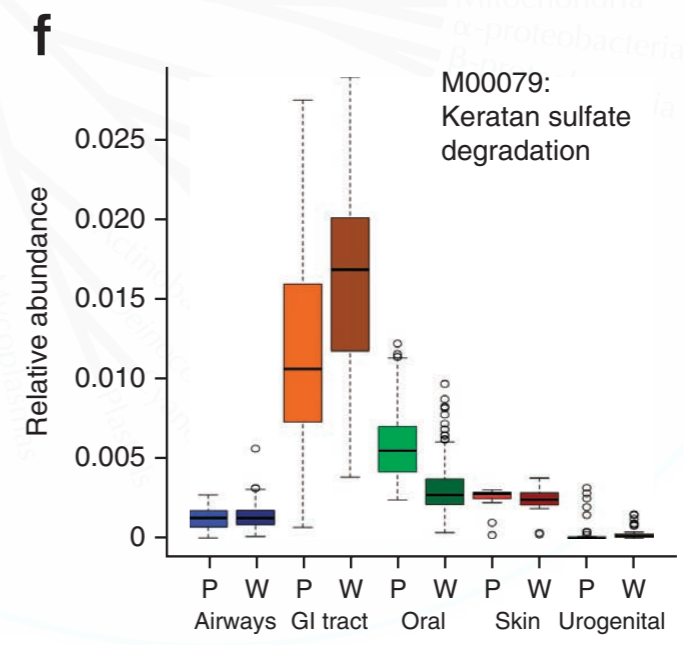
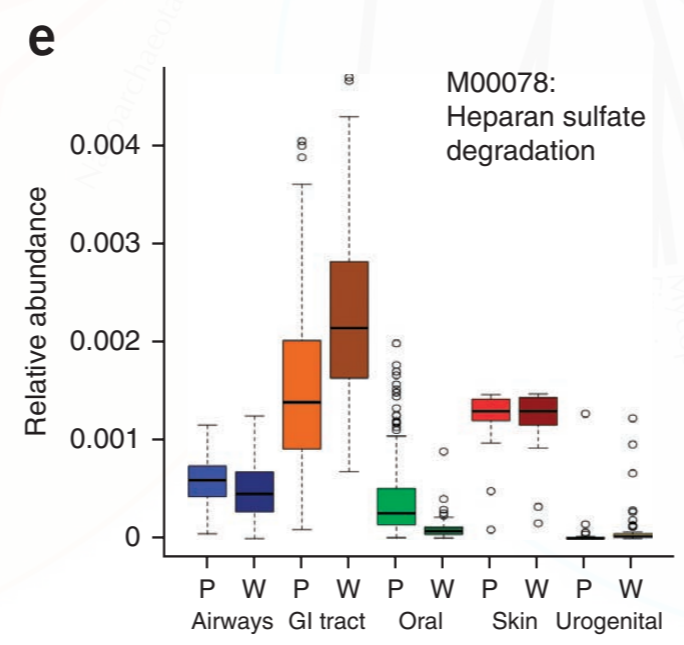
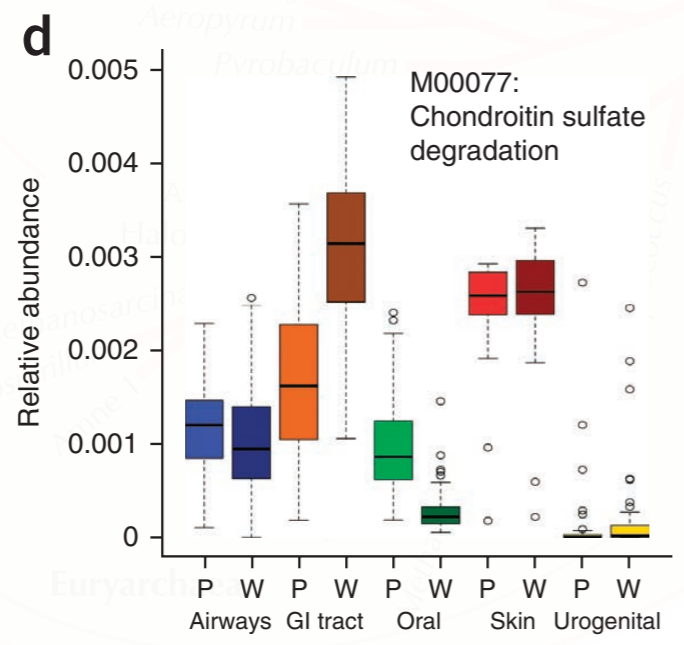
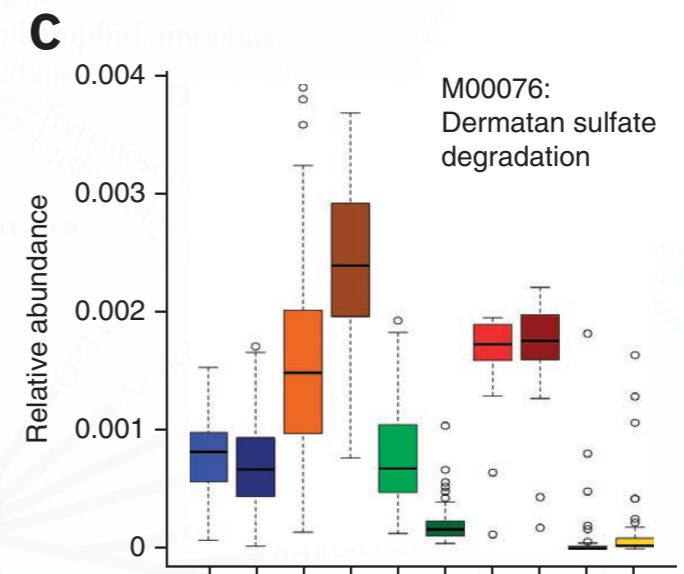
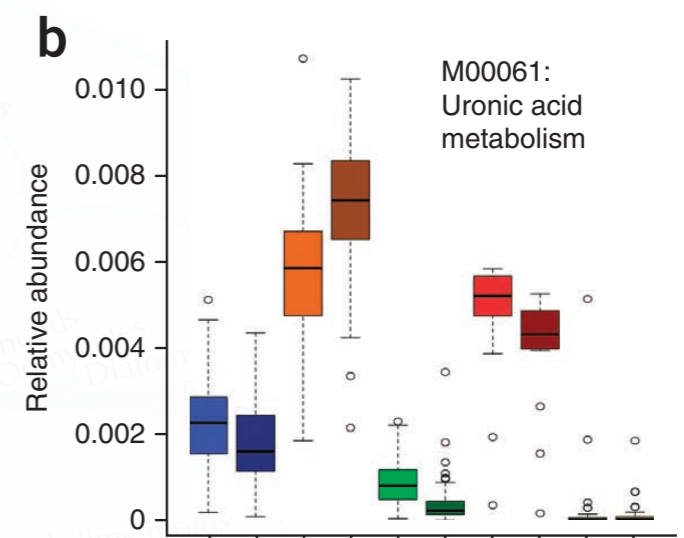
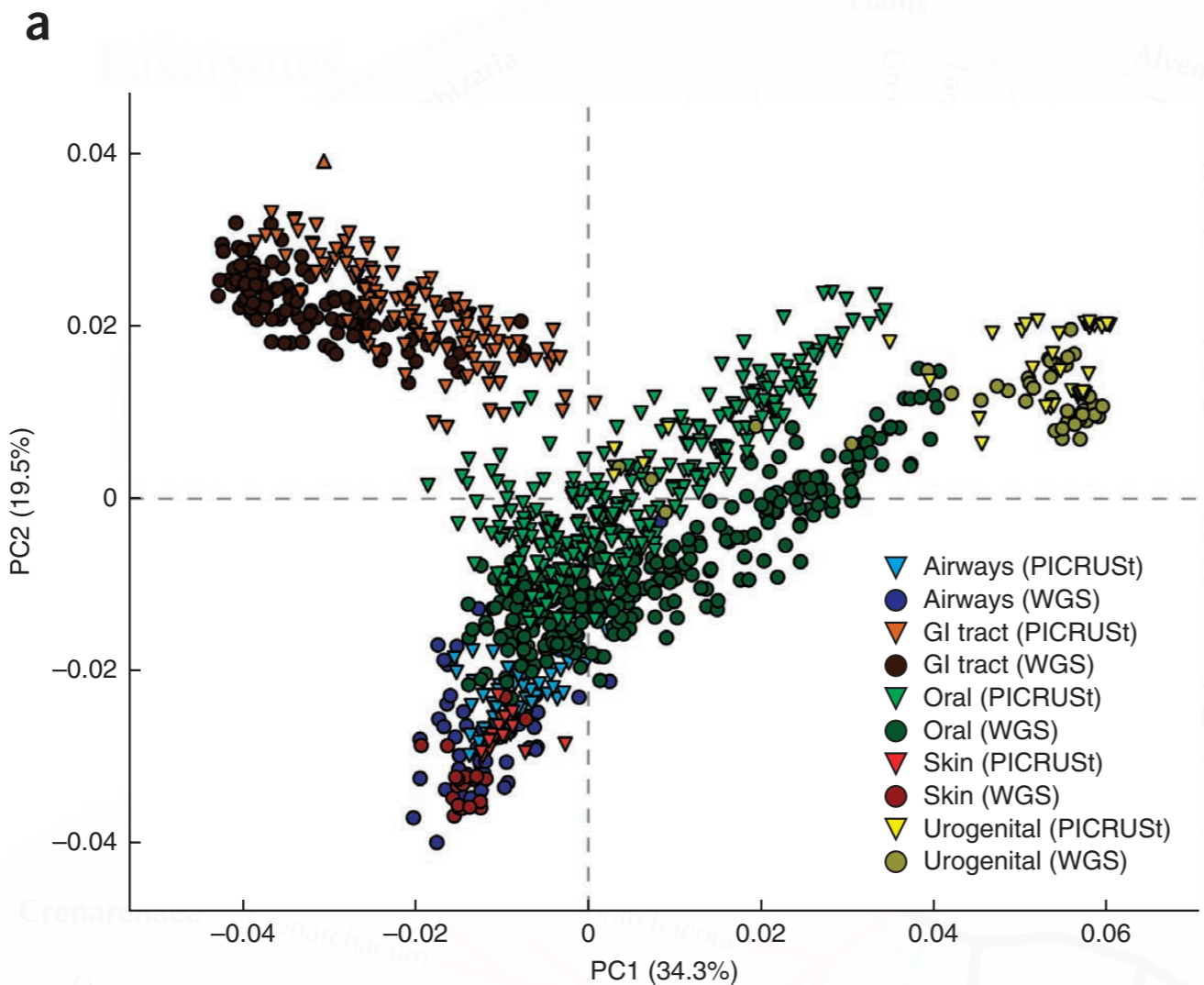


Figure 2 PICRUSt recapitulates biological findings from the Human Microbiome Project. **(a)** Principal component analysis (PCA) plot comparing KEGG module predictions using 16S data with PICRUSt (lighter colored triangles) and sequenced shotgun metagenome (darker colored circles) along with relative abundances for five specific KEGG modules: **(b)** [M00061](#): Uronic acid metabolism. **(c)** [M00076](#): Dermatan sulfate degradation. **(d)** [M00077](#): Chondroitin sulfate degradation. **(e)** [M00078](#): Heparan sulfate degradation. **(f)** [M00079](#): Keratan sulfate degradation. All KEGG modules are involved in glycosaminoglycan degradation (KEGG pathway [ko00531](#)) using 16S with PICRUSt (P) and whole genome sequencing (W) across human body sites. Color key is the same as in **a**.

Tree of Life

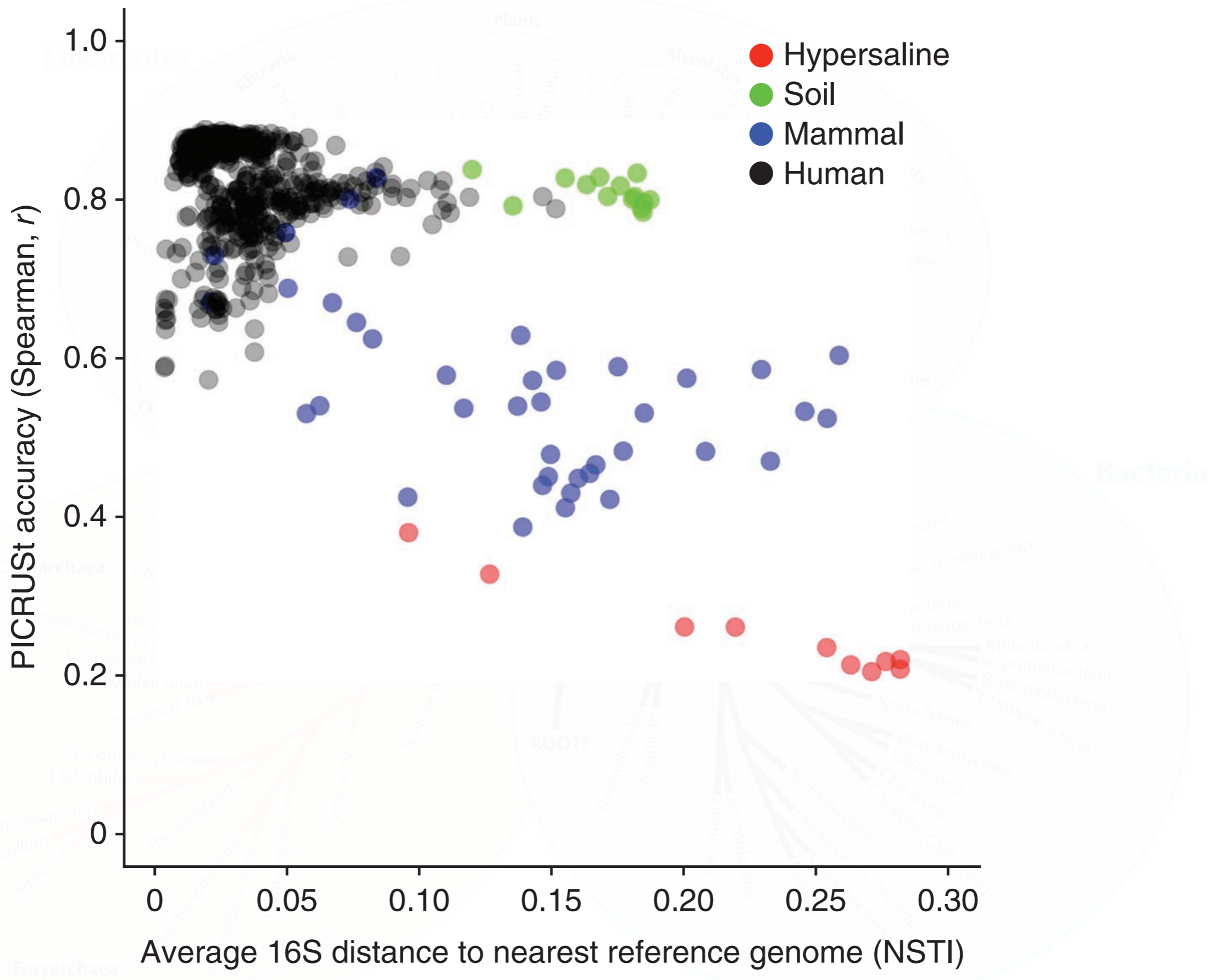
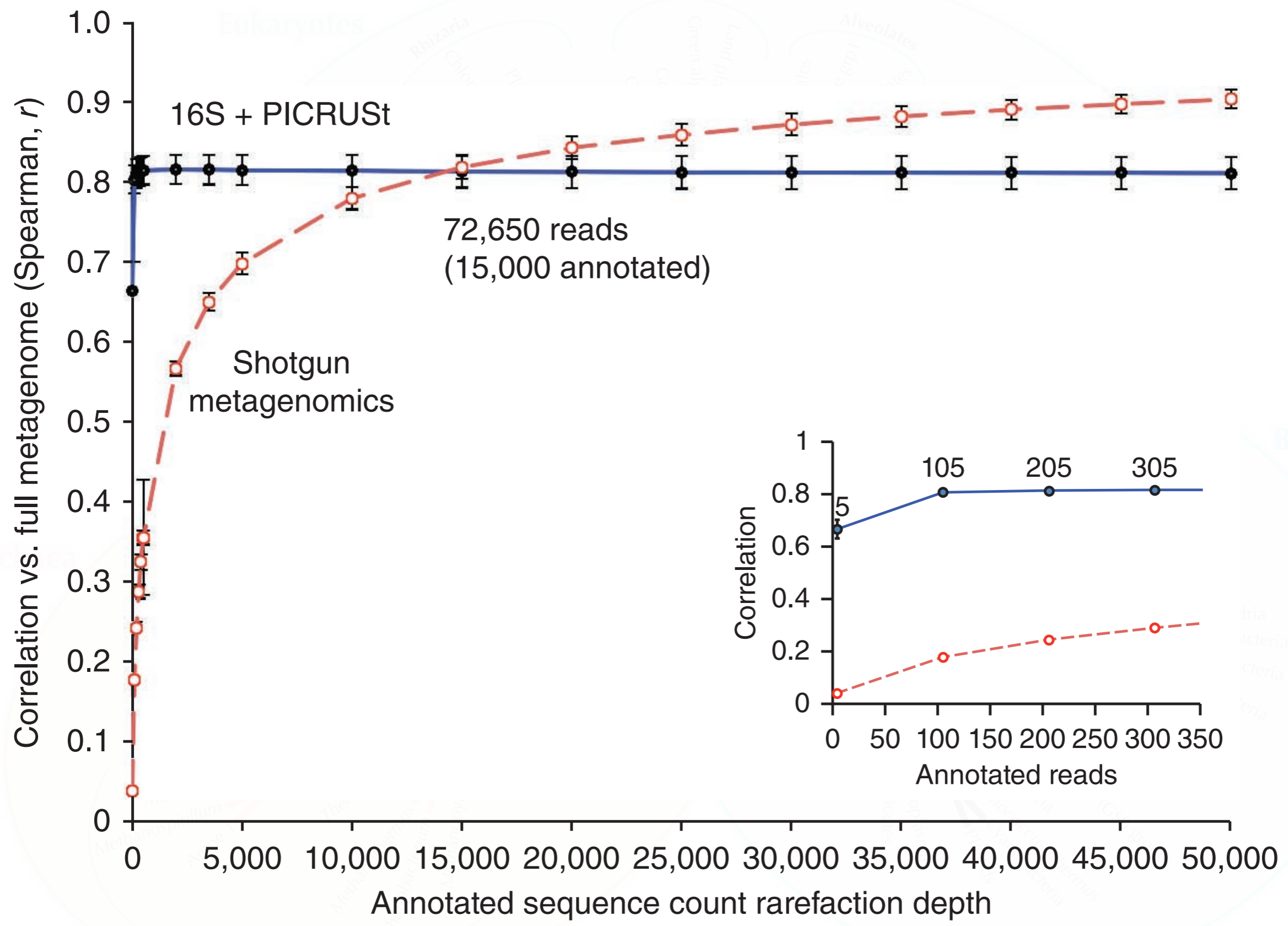


Figure 3 PICRUSt accuracy across various environmental microbiomes. Prediction accuracy for paired 16S rRNA marker gene surveys and shotgun metagenomes are plotted against the availability of reference genomes as summarized by NSTI. Accuracy is summarized using the Spearman correlation between the relative abundance of gene copy number predicted from 16S data using PICRUSt versus the relative abundance observed in the sequenced shotgun metagenome. In the absence of large differences in metagenomic sequencing depth, relatively well-characterized environments, such as the human gut, had low NSTI values and can be predicted accurately from 16S surveys. Conversely, environments containing much unexplored diversity (e.g., phyla with few or no sequenced genomes), such as the Guerrero Negro hypersaline microbial mats, tended to have high NSTI values.

Tree of Life

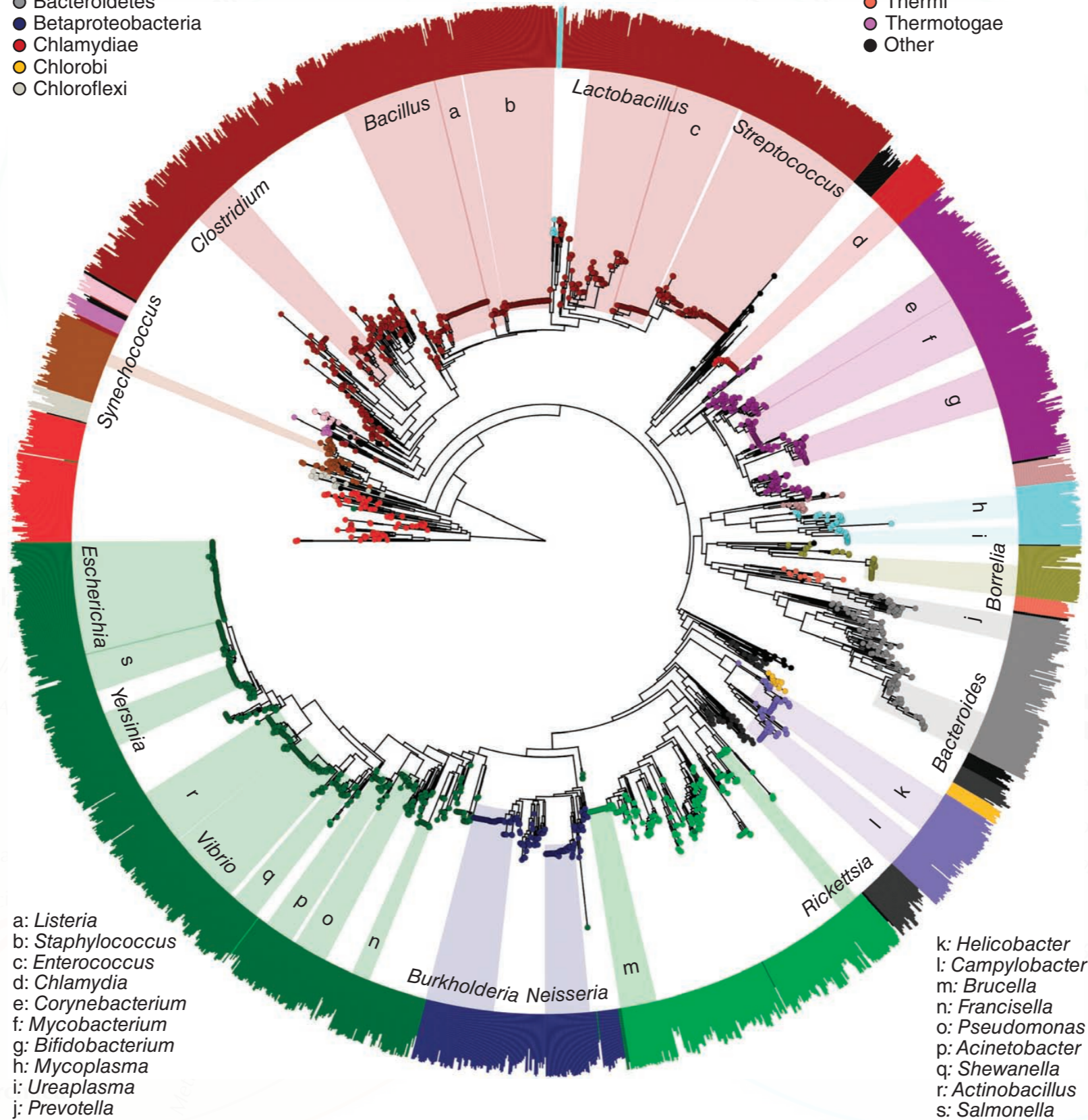


Bacteria

Figure 4 Accuracy of PICRUSt prediction compared with shotgun metagenomic sequencing at shallow sequencing depths. Spearman correlation between either PICRUSt-predicted metagenomes (blue lines) or shotgun metagenomes (dashed red lines) using 14 soil microbial communities subsampled to the specified number of annotated sequences. This rarefaction reflects random subsets of either the full 16S OTU table (blue) or the corresponding gene table for the sequenced metagenome (red). Ten randomly chosen rarefactions were performed at each depth to indicate the expected correlation obtained when assessing an underlying true metagenome using either shallow 16S rRNA gene sequencing with PICRUSt prediction or shallow shotgun metagenomic sequencing. The data label describes the number of annotated reads below which PICRUSt-prediction accuracy exceeds metagenome sequencing accuracy. Note that the plotted rarefaction depth reflects the number of 16S or metagenomic sequences remaining after standard quality control, dereplication and annotation (or OTU picking in the case of 16S sequences), not the raw number returned from the sequencing facility. The number of total metagenomic reads below which PICRUSt outperforms metagenomic sequencing (72,650) for this data set was calculated by adjusting the crossover point in annotated reads (above) using annotation rates for the soil data set (17.3%) and closed-reference OTU picking rates for the 16S rRNA data set (68.9%). The inset figure illustrates rapid convergence of PICRUSt predictions given low numbers of annotated reads (blue line).

Tree of Life

- Actinobacteria
- Alphaproteobacteria
- Epsilonproteobacteria
- Fusobacteria
- Aquificae
- Cyanobacteria
- Euryarchaeota
- Gammaproteobacteria
- Archaea
- Deltaproteobacteria
- Firmicutes
- Spirochaetes
- Bacteroidetes
- Betaproteobacteria
- Chlamydiae
- Tenericutes
- Thermi
- Chlorobi
- Chloroflexi
- Thermotogae
- Other



- a: *Listeria*
- b: *Staphylococcus*
- c: *Enterococcus*
- d: *Chlamydia*
- e: *Corynebacterium*
- f: *Mycobacterium*
- g: *Bifidobacterium*
- h: *Mycoplasma*
- i: *Ureaplasma*
- j: *Prevotella*

- k: *Helicobacter*
- l: *Campylobacter*
- m: *Brucella*
- n: *Francisella*
- o: *Pseudomonas*
- p: *Acinetobacter*
- q: *Shewanella*
- r: *Actinobacillus*
- s: *Salmonella*

Figure 5 PICRUSt prediction accuracy across the tree of bacterial and archaeal genomes. Phylogenetic tree produced by pruning the Greengenes 16S reference tree down to those tips representing sequenced genomes. Height of the bars in the outermost circle indicates the accuracy of PICRUSt for each genome (accuracy: 0.5–1.0) colored by phylum, with text labels for each genus with at least 15 strains. PICRUSt predictions were as accurate for archaeal (mean = 0.94 0.04 s.d., $n = 103$) as for bacterial genomes (mean = 0.95 0.05 s.d., $n = 2,487$).

Tree of Life

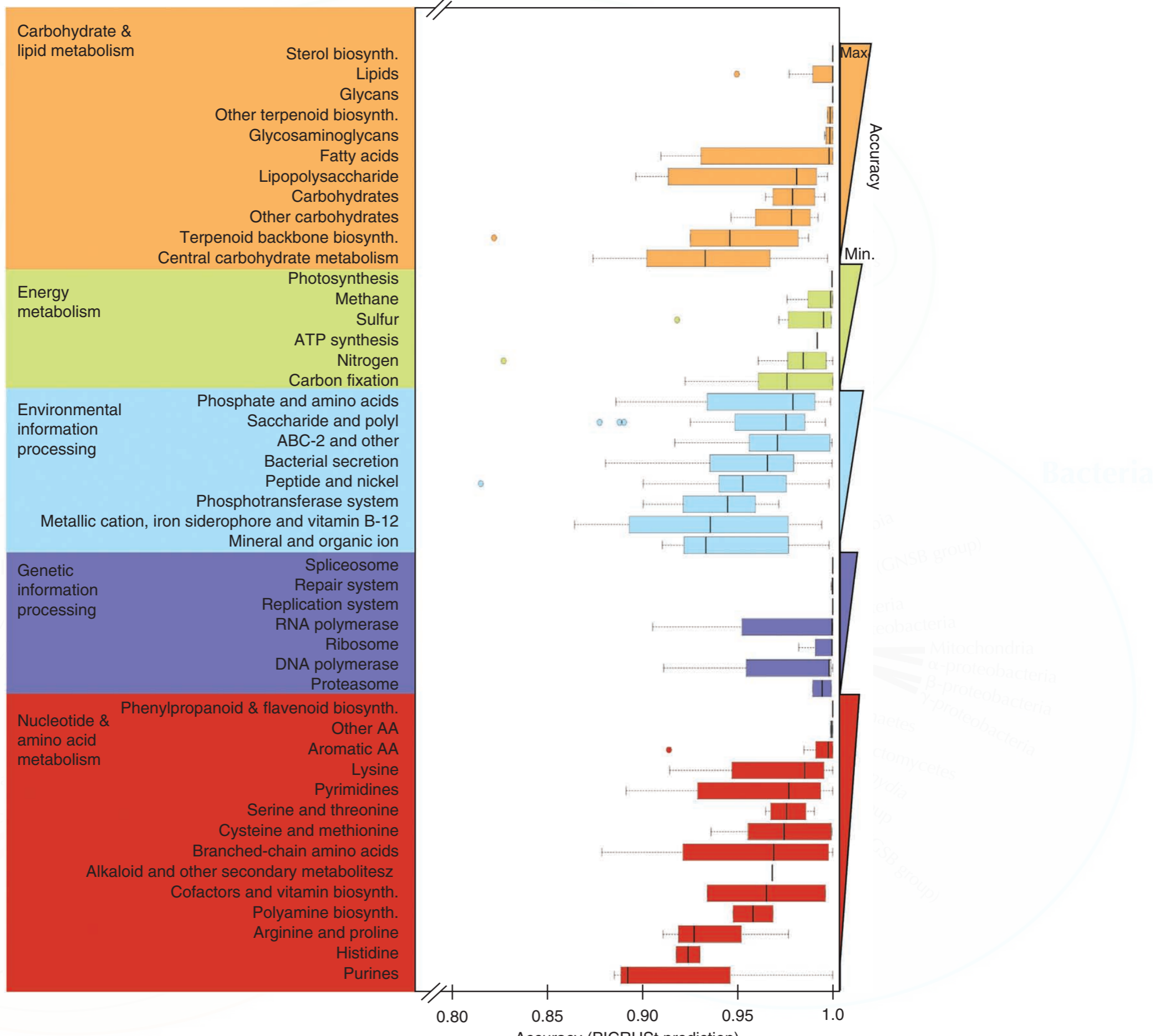


Figure 6 Variation in inference accuracy across functional modules within single genomes. Results are colored by functional category and sorted in decreasing order of accuracy within each category (indicated by triangular bars, right margin). Note that accuracy was >0.80 for all, and therefore the region 0.80–1.0 is displayed for clearer visualization of differences between modules.