

# Using the Concept of Informative Genomic Segment to Investigate Microbial Diversity of Metagenomics Sample

Qingpeng Zhang, C. Titus Brown

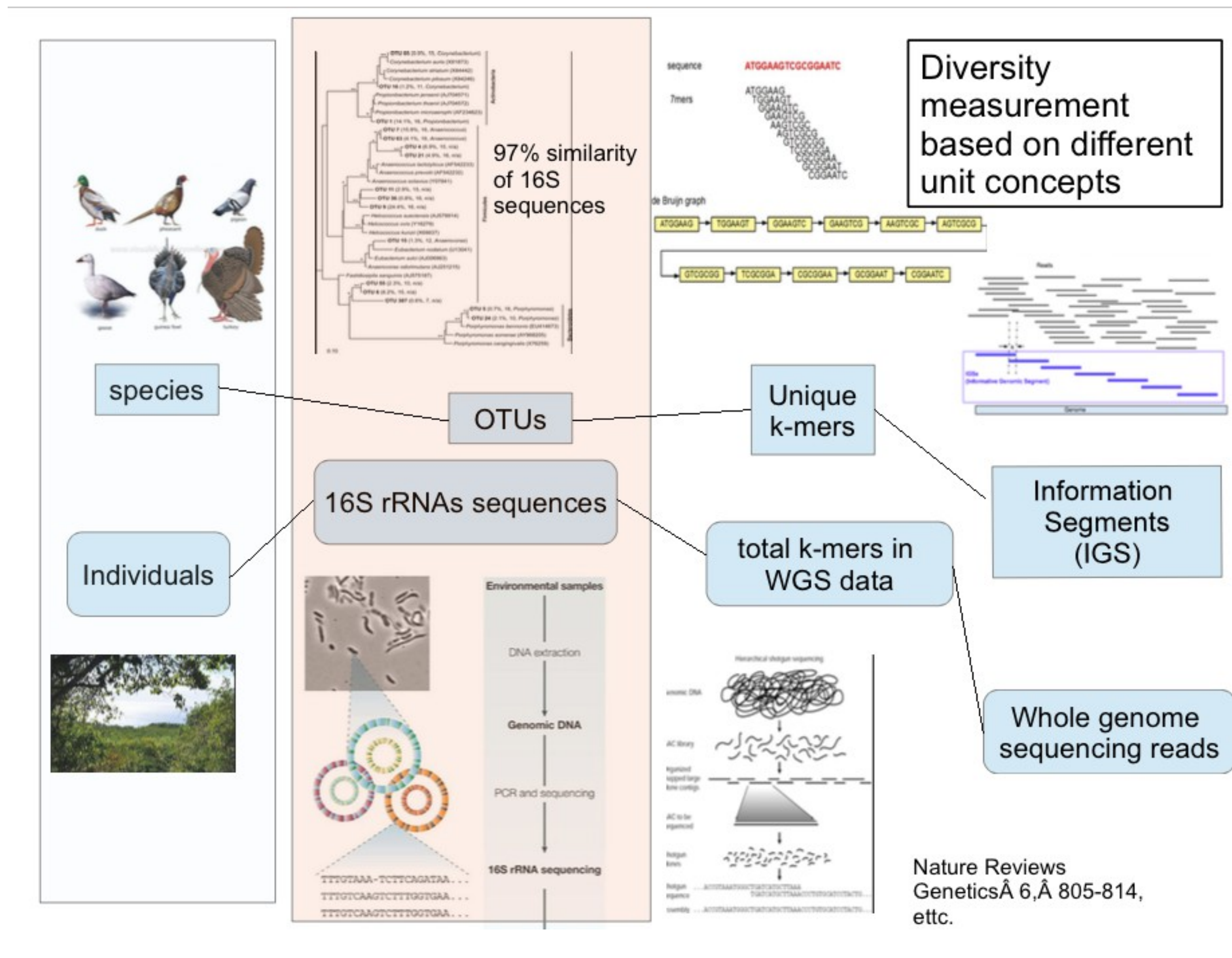
Department of Computer Science and Engineering  
Michigan State University, East Lansing, MI, United States

## Introduction

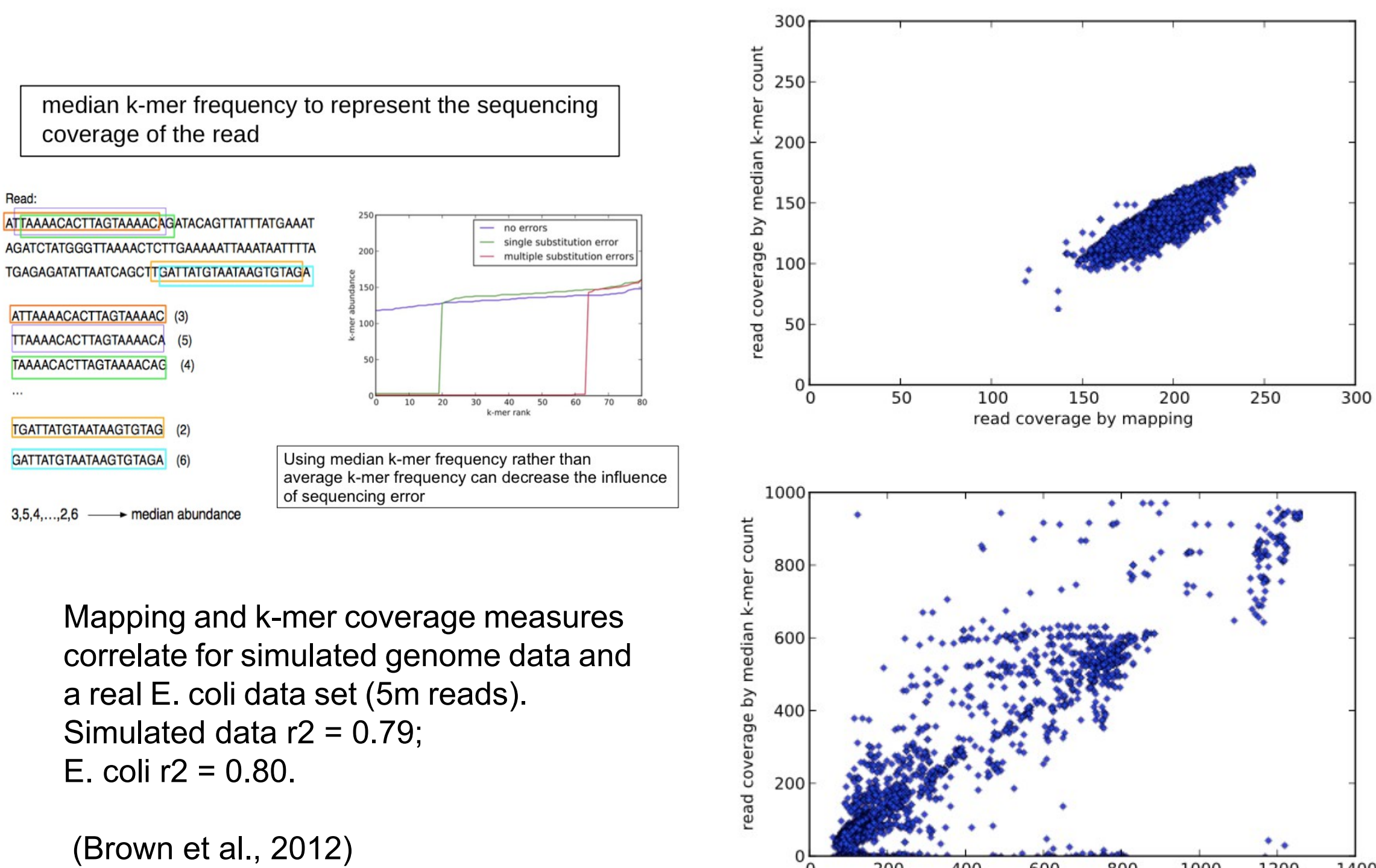
In almost all the metagenomics projects, diversity analysis plays an important role to supply information about the richness of species, the species abundance distribution in a sample or the similarity and difference between different samples, all of which are crucial to draw insightful and reliable conclusion. Traditionally OTUs(Operational Taxonomic Units) are used as the cornerstone for diversity analysis. Here we propose a novel concept - IGS (informative genomic segment) and use IGS as a replacement of OTUs to be the cornerstone for diversity analysis of whole shotgun metagenomics data sets. IGSs represent the unique information in a metagenomics data set and the abundance of IGSs in different samples can be retrieved by the reads coverage through an efficient k-mer counting method. This samples-by-IGS abundance data matrix is a promising replacement of samples-by-OTU data matrix used in 16S rRNA based analysis and all existing statistical methods can be borrowed to work on the samples-by-IGS data matrix to investigate the diversity. We applied the IGS-based method to Global Ocean Sampling Expedition (GOS) dataset and the samples were clustered more accurately than existing alignment-based method. We also tried this novel method to MetaHIT data sets. Since this method is totally binning-free, assembly-free, annotation-free, reference-free, it is specifically promising to deal with the highly diverse samples, while we are facing large amount of “dark matters” in it, like soil.

## Background

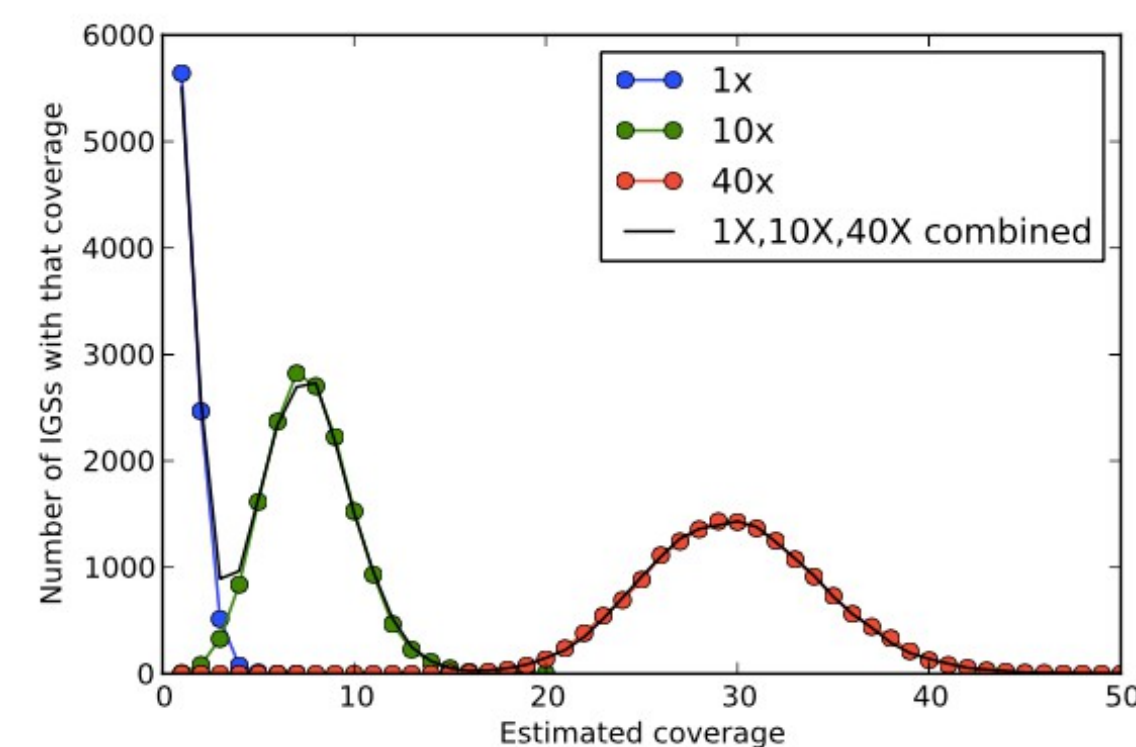
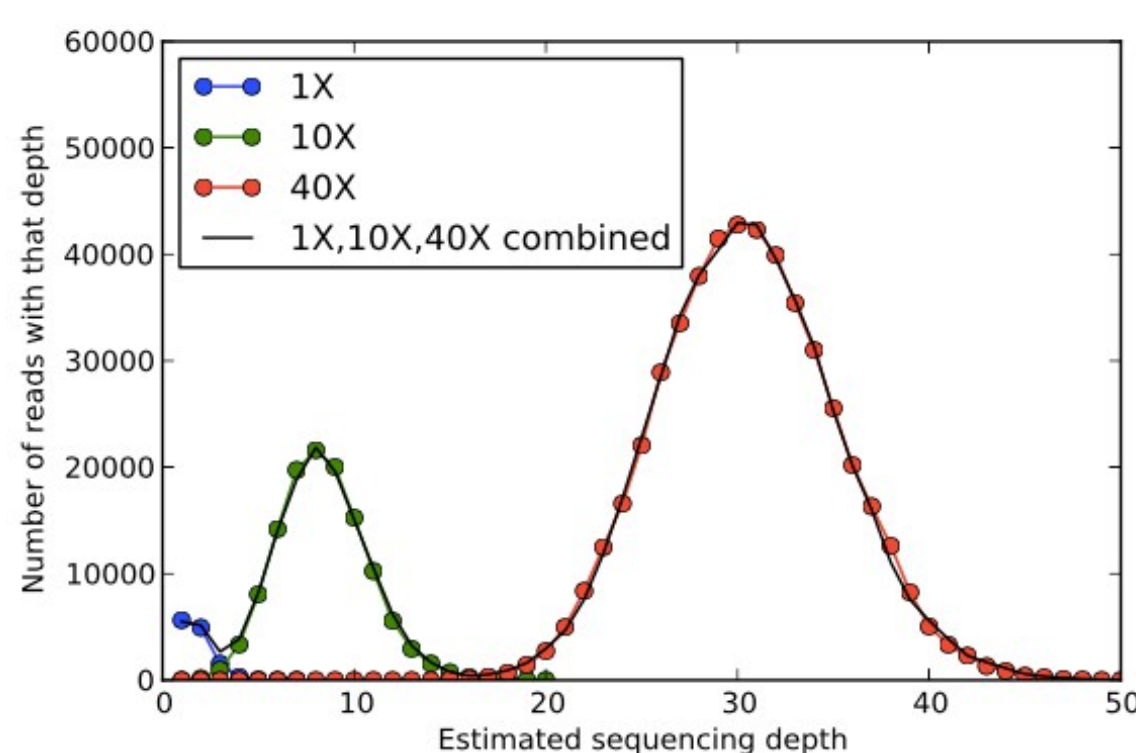
In traditional ecology, the concept of species is used to investigate diversity. In microbial ecology, the concept of OTU is used to investigate microbial diversity. OTU is mostly used for 16S data sets. And binning reads into OTU is typically required for OTU-based diversity analysis. Here we propose a new concept - IGS to replace the concept of OTU in 16S based diversity analysis and the concept of species in traditional ecology.



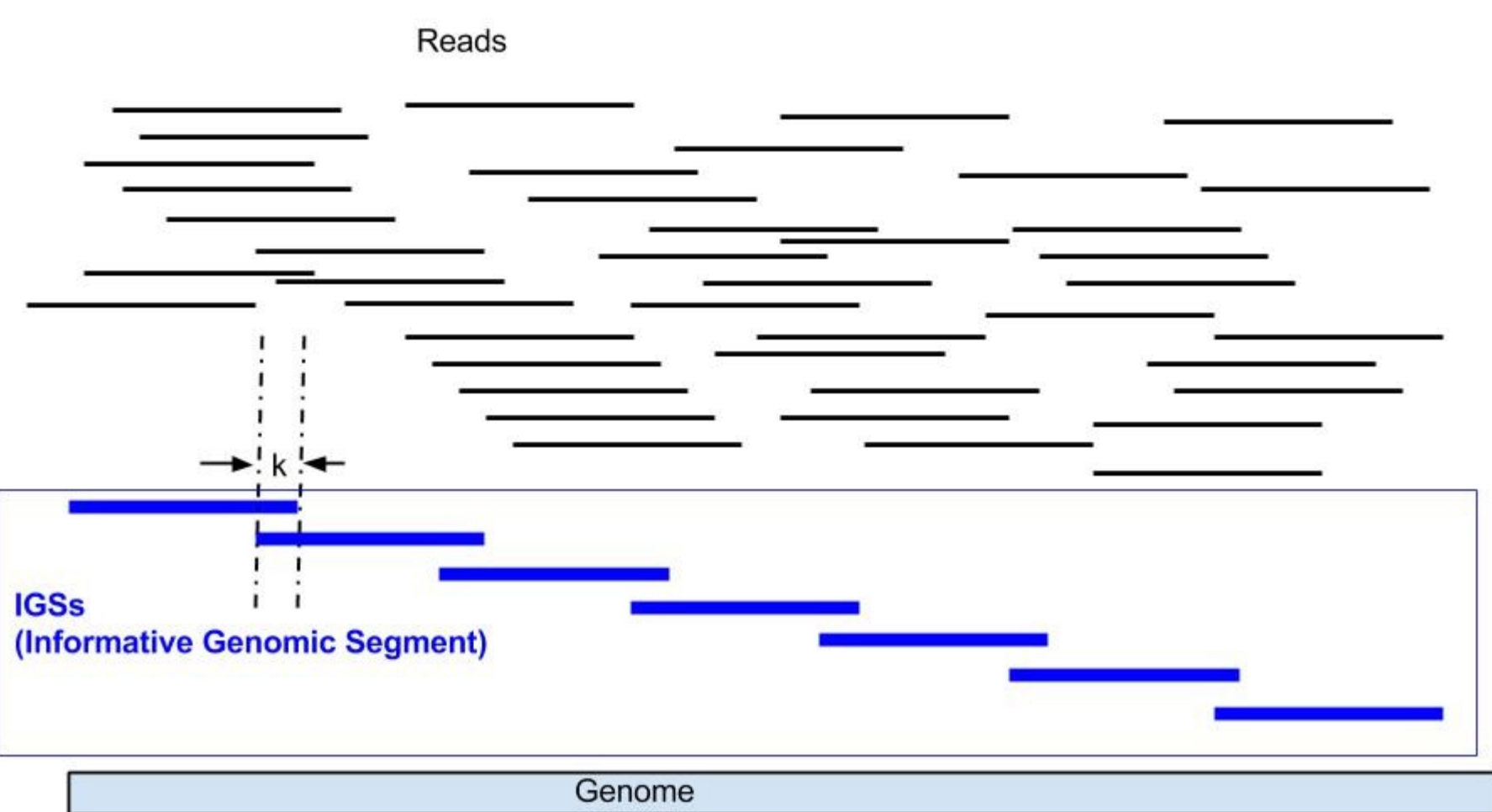
## The concept of IGS(informative genomic segment)



## IGS(informative genomic segment) can represent the novel information of a genome



Data set	total number of IGSs
1X depth	8714
10X depth	16321
40X depth	16794



The concept of IGS for single genome

## Using IGS to do diversity analysis

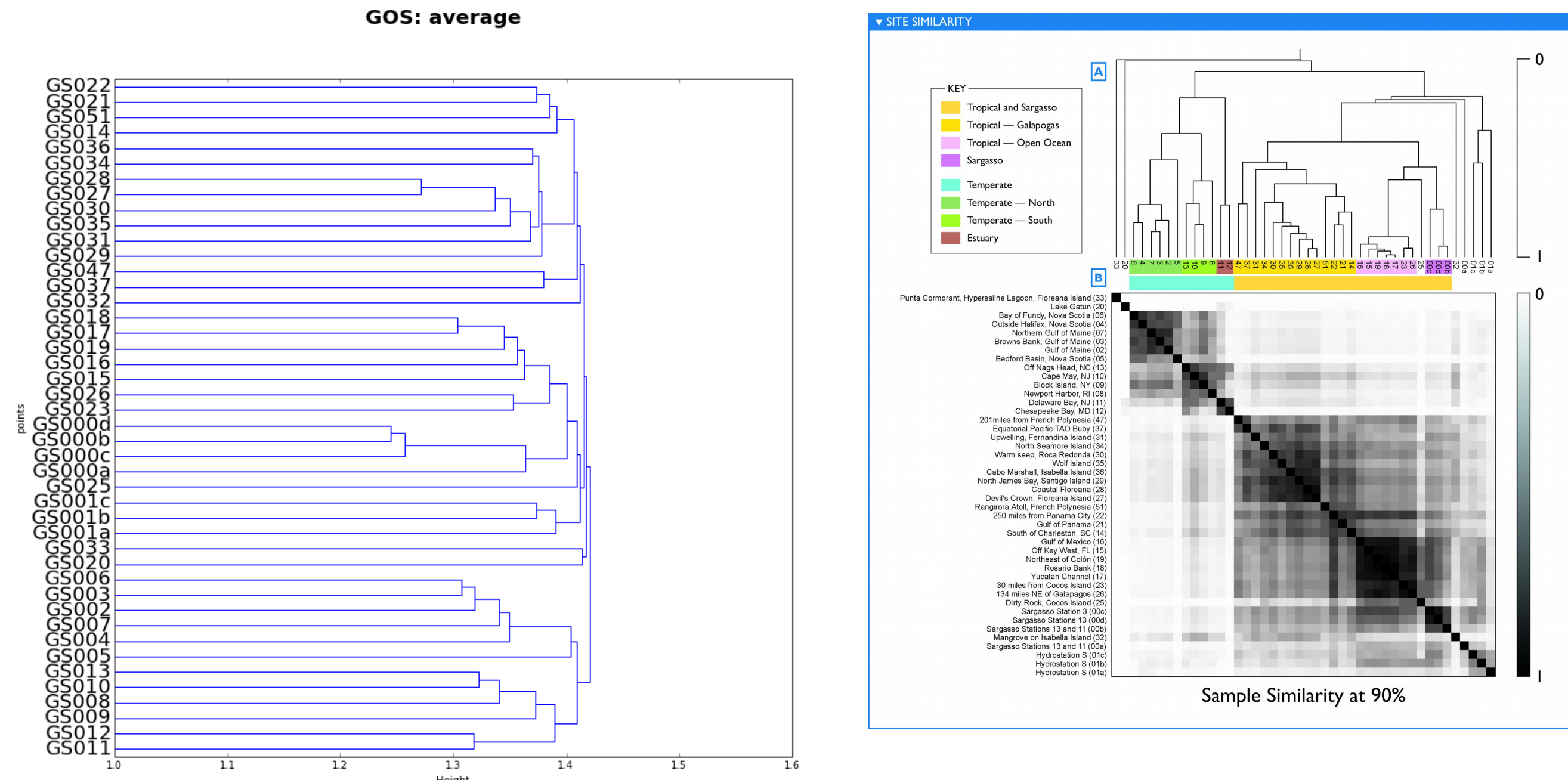
### Build samples-by-IGS matrix to replace samples-by-OUT matrix



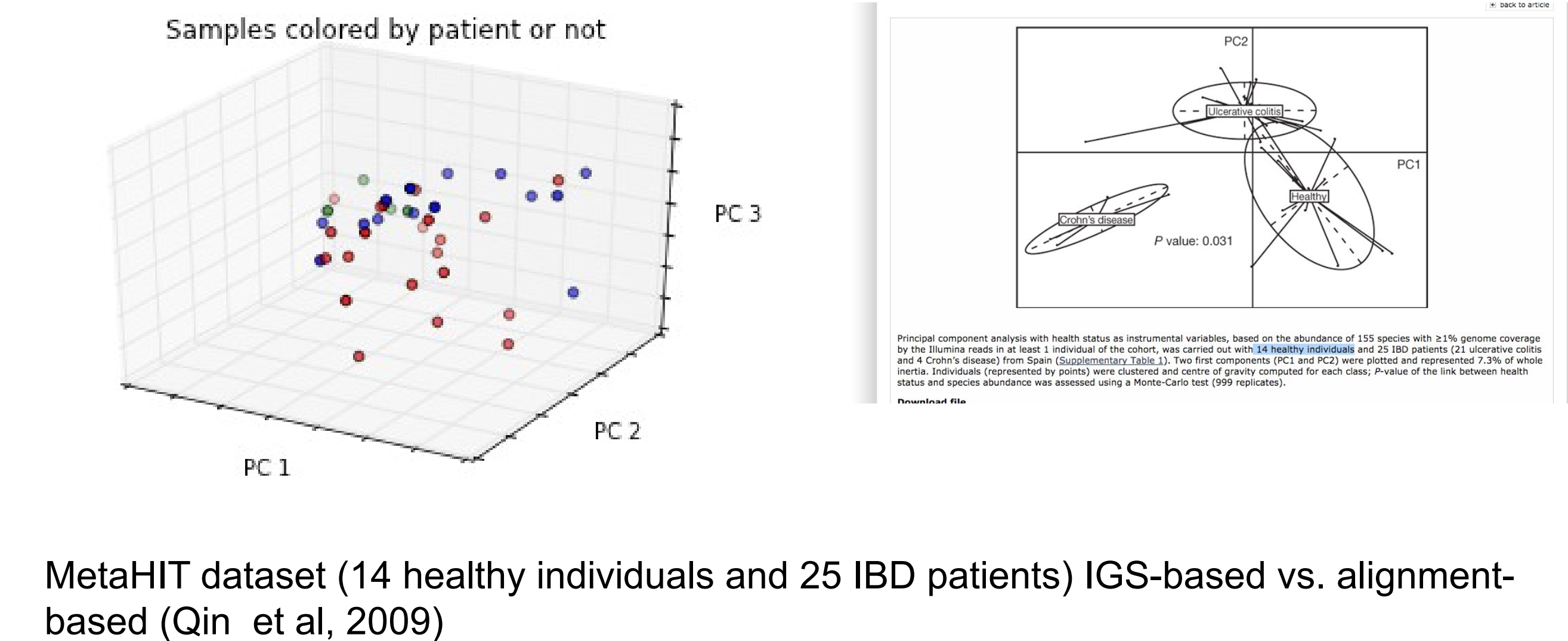
We can generate a sample-by-IGS data matrix as the counterpart of samples-by-OTU data matrix so many of the existing tools/methods used for OTU-based diversity can be borrowed for this kind of IGS-based analysis.

### Apply IGS-based method on real data sets

IGS-based method can get comparable if not better beta-diversity result than traditional methods based on reference/alignment.



Global Ocean Sampling Expedition (GOS) IGS-based vs. alignment-based (Rusch et al, 2007)



This IGS-based method to do microbial diversity analysis is totally binning-free, assembly-free, annotation-free, reference-free and using this method to do alpha diversity analysis is under investigation.