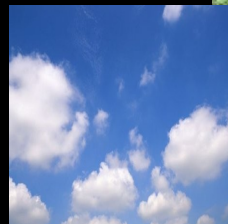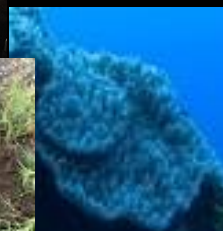# Characterization of environmental viral diversity using metagenomics

*Florent Angly*

RohwerLab

CSRC

# Ecological importance of viruses
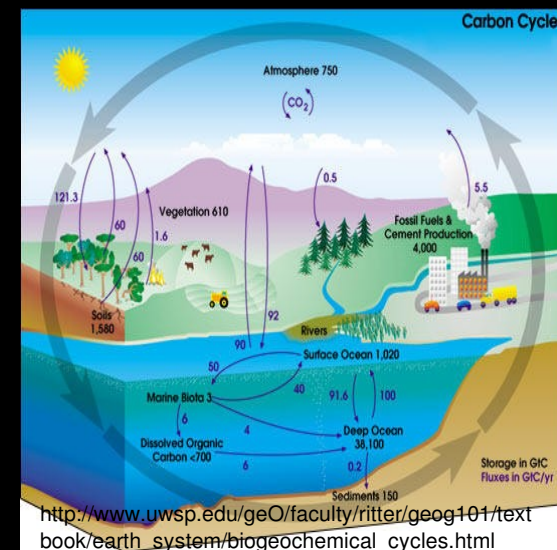
Ubiquity



Large abundance



http://www.virusecology.org/MOVE/Method%206.html

Global impact



http://www.uwsp.edu/geO/faculty/ritter/geog101/text book/earth_system/biogeochemical_cycles.html

*Culturable*

# Metagenomics



Environmental sample: uncultured phages



Environmental phage DNA

Who is there?

What do they do?

How many are they?

~~16S rDNA~~

Taxonomy



Function



Diversity

# Diversity

Diversity types: $\alpha, \beta, \gamma$

$\alpha$-diversity: richness, evenness, indices

Shannon-Wiener index:

$$H' = -\sum_{i=1}^{S} r_i \ln r_i$$

$S$: richness
$r_i$: relative abundance
of the $i^{th}$ species

# Hypothesis and plan

*Is it possible to estimate viral diversity from metagenomes? If so, how?*

1. $\alpha$-diversity (PHACCS)
2. $\beta$-diversity (MaxiPhi)
3. Average genome length (GAAS)
4. Diversity workflow
5. Environmental viral diversity

# - 1 -

# $\alpha$-diversity

# PHACCS

# $\alpha$-diversity from metagenomic data

$\alpha$-diversity: local diversity (one sample)

1. Assemble metagenomic sequences

2. Count the number of contigs

3. Assume that only sequences from the same species form contigs

4. Model diversity: the more abundant a species, the larger the number of its sequences forming contigs

Circonspect



Metagenome

1-contig {

2-contig {

3-contig {

Contig spectrum [11 4 1 0 0 ...]

$\alpha$-diversity estimates

# Characterization of environmental viral diversity using metagenomics

## φHACCS
### Phage communities from Contig Spectrum



Assumed community structure

Generalized Lander-Waterman equations

$$c_q = \sum_{i=1}^{M} n_i w_{qi}$$

$n_i$ : expected number of sequences for species $i$

$w_{qi}$ : probability that a sequence of species $i$ is in a contig of size $q$
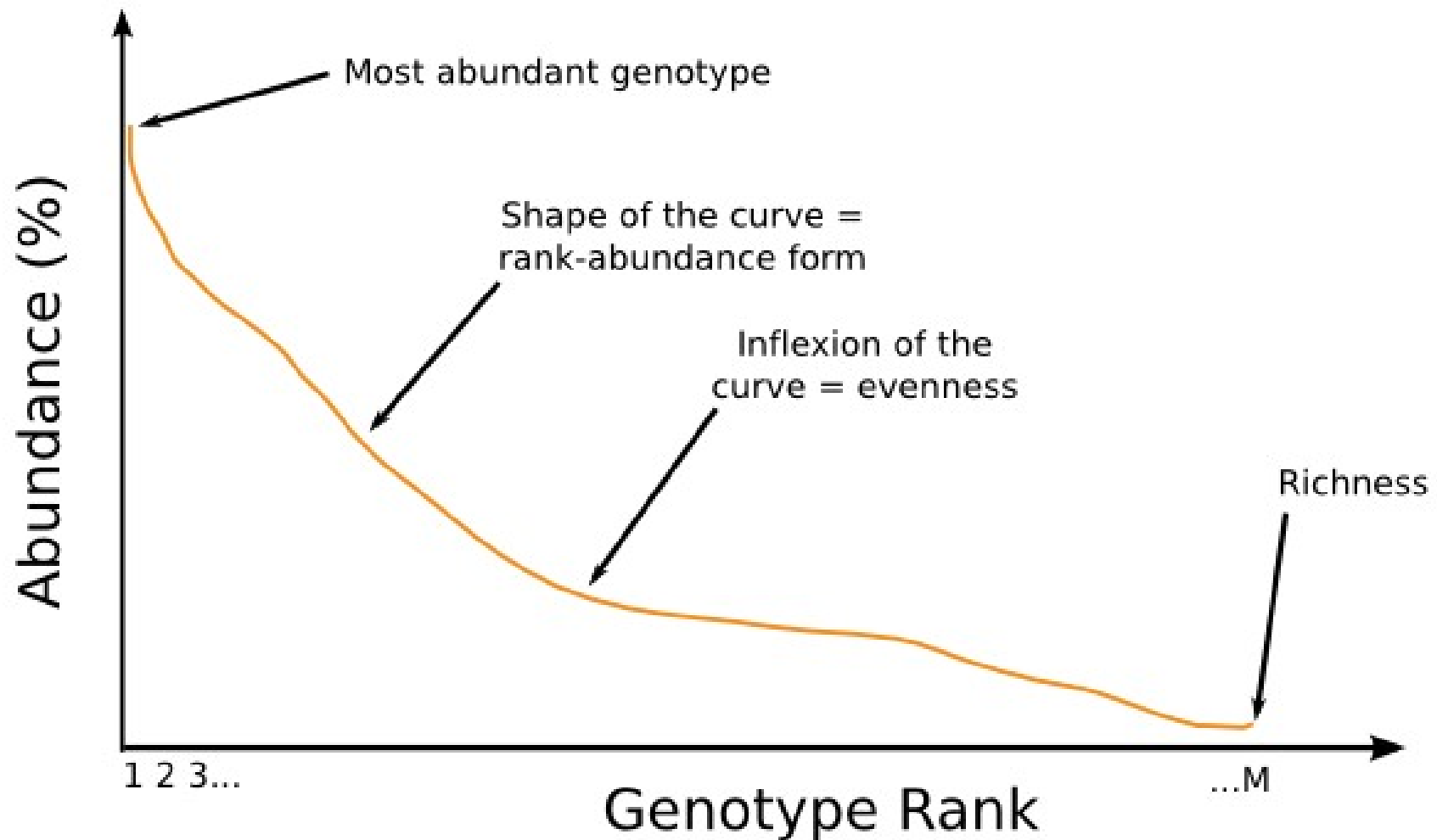
Assumptions:

- No chimeric contigs

- Rank-abundance distribution

- Genomes have an average length

*Angly et al. 2005*

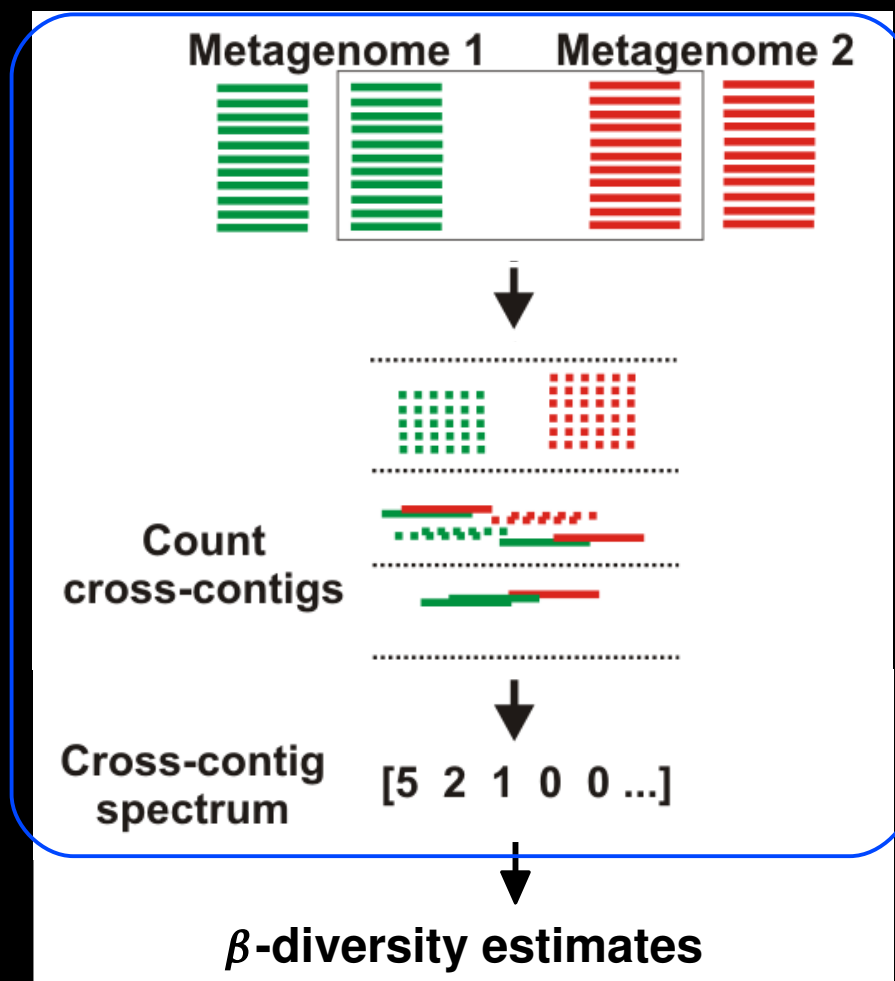# Community structure and diversity

# - 2 -

# $\beta$-diversity

# MaxiPhi

# $\beta$-diversity from metagenomic data

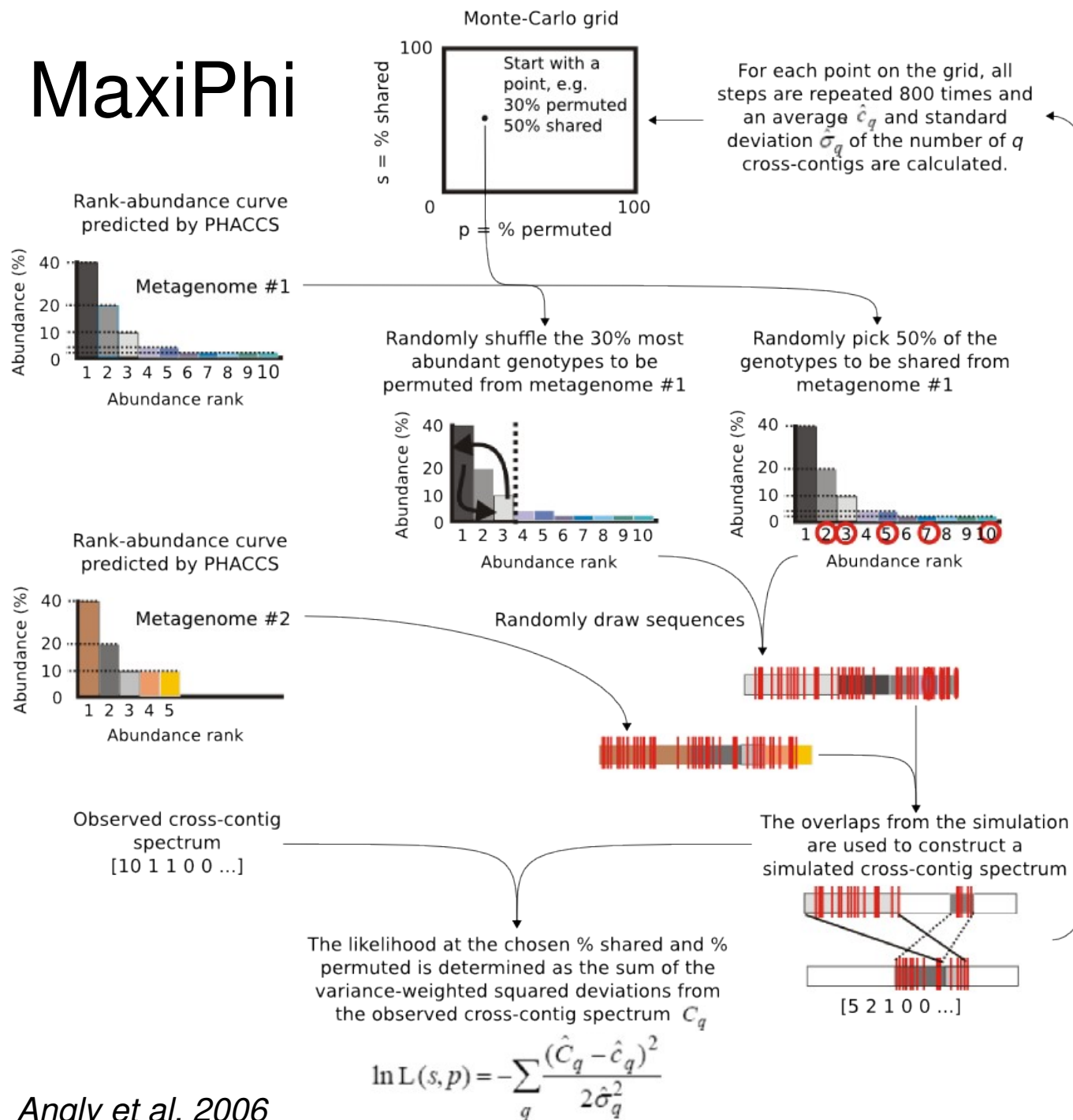- $\beta$-diversity: difference in diversity between several samples

Cross-contig spectrum

# MaxiPhi



Monte-Carlo grid

s = % shared

Start with a point, e.g. 30% permuted 50% shared

p = % permuted

For each point on the grid, all steps are repeated 800 times and an average $\hat{c}_q$ and standard deviation $\hat{\sigma}_q$ of the number of $q$ cross-contigs are calculated.

Rank-abundance curve predicted by PHACCS

Metagenome #1

Abundance (%)

Abundance rank

Randomly shuffle the 30% most abundant genotypes to be permuted from metagenome #1

Randomly pick 50% of the genotypes to be shared from metagenome #1

Rank-abundance curve predicted by PHACCS

Metagenome #2

Abundance (%)

Abundance rank

Randomly draw sequences

Observed cross-contig spectrum
[10 1 1 0 0 ...]

The overlaps from the simulation are used to construct a simulated cross-contig spectrum

[5 2 1 0 0 ...]

The likelihood at the chosen % shared and % permuted is determined as the sum of the variance-weighted squared deviations from the observed cross-contig spectrum $C_q$

$$\ln L(s,p) = -\sum_q \frac{(\hat{C}_q - \hat{c}_q)^2}{2\hat{\sigma}_q^2}$$

*Angly et al. 2006*

# Modeling $\beta$-diversity

- Percent shared and percent permuted

- 3 -

Average genome length

GAAS

# Genome length of different organisms



Genome length (bp)

| | | | |
|---|---|---|---|
| 10³ | 10⁶ | 10⁹ | 10¹² |

Animals — *Pratylenchus coffeae* Plant-parasitic nematode **20 Mb** — **130 Gb** *Protopterus aethiopicus* Marbled lungfish

Plants — *Ostreococcus tauri* Green alga **13 Mb** — **250 Gb** *Psilotum nudum* Fern

Fungi — *Pneumocystis pneumonia* Yeast-like fungus **6.5 Mb** — **800 Mb** *Scutellospora castanea* Mycorrhizal fungi
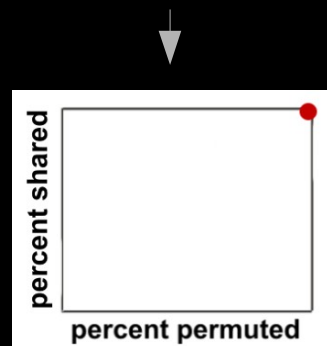
Protists — *Bigelowiella natans* Rhizaria **440 kb** — **1.4 Tb** *Chaos chaos* Amoeba

Bacteria — *Carsonella rudii* Endosymbiotic γ-Proteobacteria **160 kb** — **13 Mb** *Sorangium cellulosum* Soil Myxobacteria

Archaea — *Nanoarchaeum equitans* Hyperthermophilic archaea **490 kb** — **4.1 Mb** *Methanosarcina acetivorans* Colony-forming archaea

Phage * — *Leuconostoc phage L5* **2.4 kb** — **500 kb** *Bacillus phage G* Tailed phage

Eukaryotic viruses — *Porcine circovirus* Small circular virus **1.7 kb** — **1.7 Mb** *Mamavirus* Amoeba virus

Mimivirus
1,700 kb
750 nm

http://www.microbiologybytes.com/virology/Mimivirus.html

Circovirus
1.7 kb
17 nm

http://www.pcvd.org/

Microbes
and
viruses

Data compiled from the ICTVdb, NCBI RefSeq, Microbe Wiki, Fungal Genome Size Database, Plant DNA C-values Database and Animal Genome Size Database

# GAAS

Metagenome    Complete genome database    ★ Cutoff alignment E-value, identity percent, and relative length

→ BLAST ←

target sequence ▬▬▬▬▬▬▬▬▬▬

50% relative size

query sequence ▬▬▬▬

95% relative size

target sequence ▬▬▬▬▬▬▬▬▬▬

Filter out weak similarities ★

★★ Normalize by genome length

Keep & weight all similarities ★★

★★ For each query sequence *i*, keep all the similarities to genomes *j* but give similarities a weight that has a statistical meaning:

**Environment**
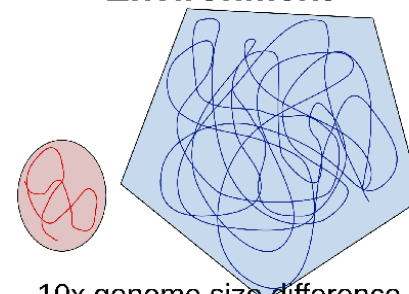
★★ Relative abundance

$$w_{ij} \propto k_i / E_{ij}$$

10x genome size difference
Same relative abundance

**Metagenome**

Average genome length
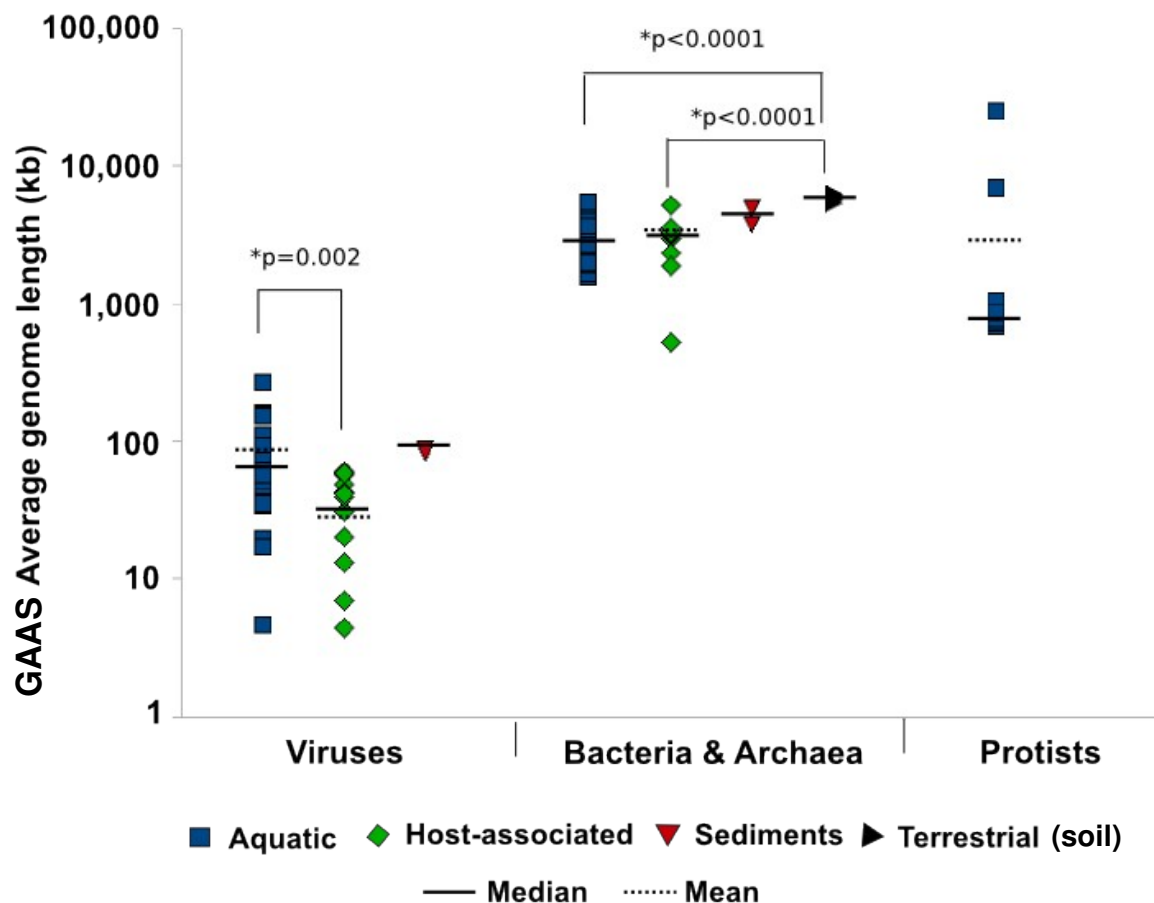
$w_{ij}$ : weight
$E_{ij}$ : E-value
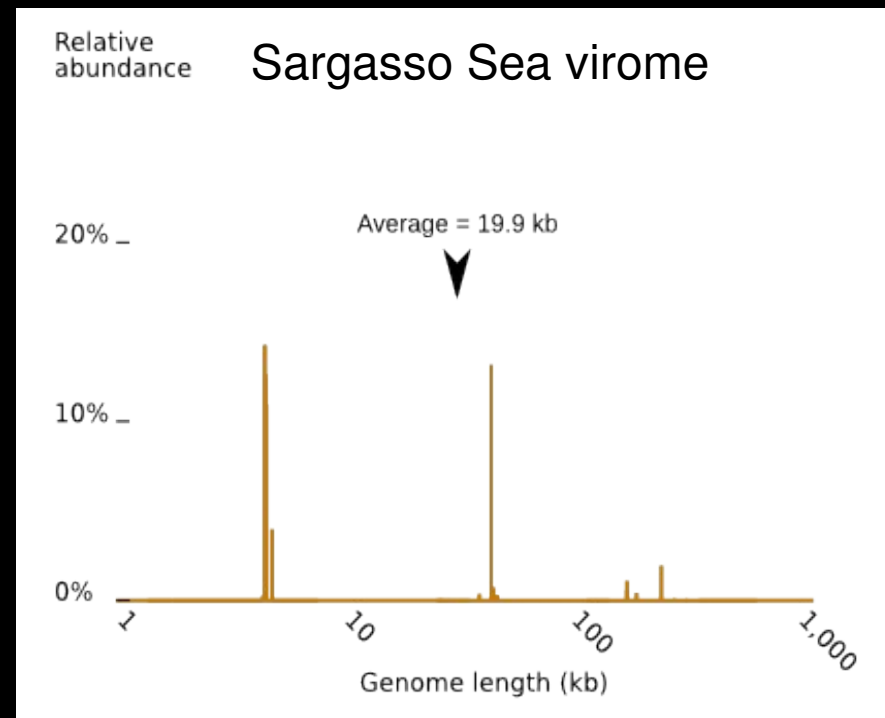$k_i$ : constant

10x more sequences from the larger genome

# Genome length in the environment

Meta-analysis of 174 metagenomes

- Variability between biomes
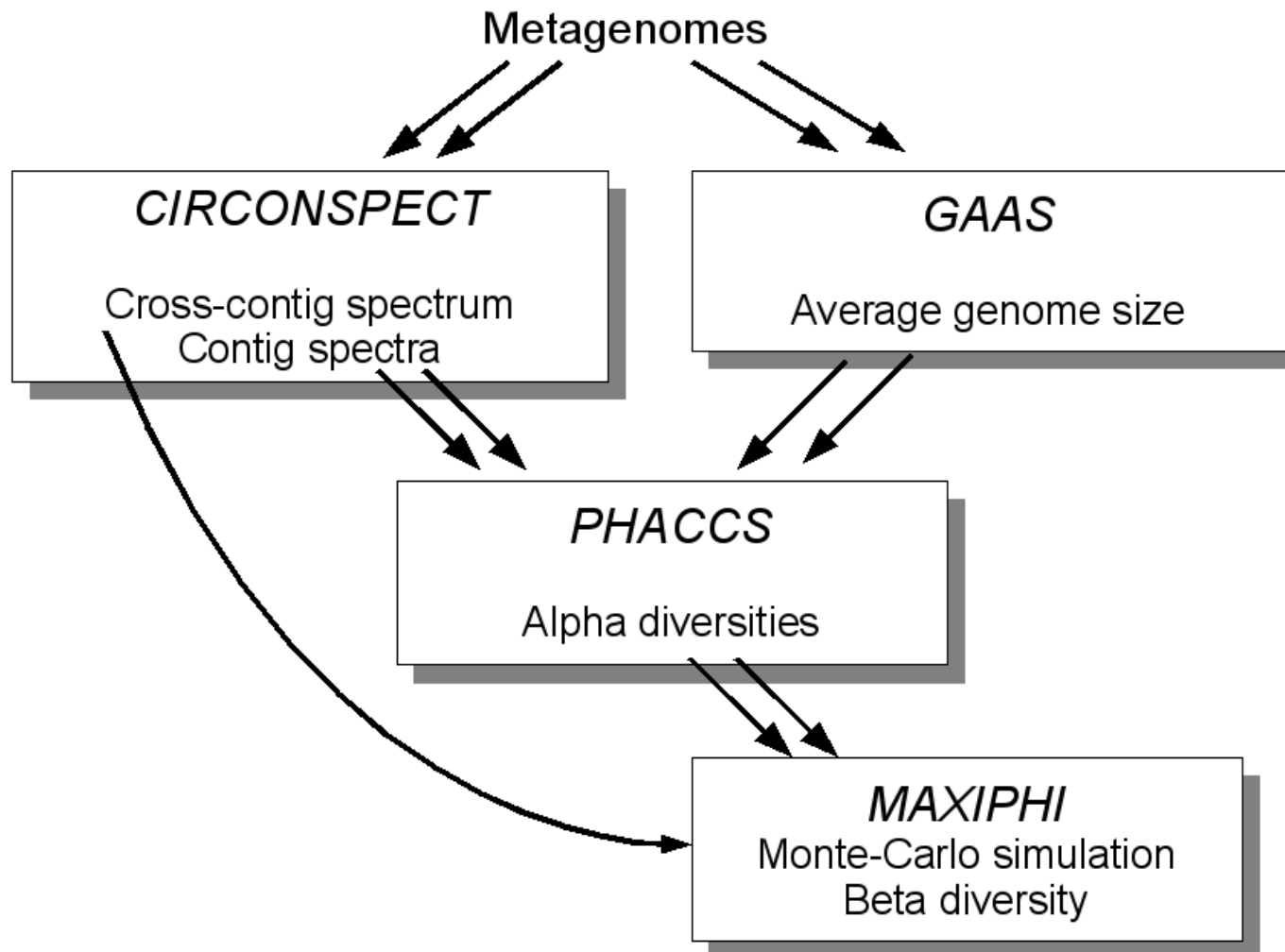- Variability within biomes
- Variability within sample



*Angly et al. In revision*

- 4 -

# Diversity workflow

# Diversity workflow

# α-diversity

CAMERA
Marine Microbial Ecology

Logout
Welcome, Florent Angly

Home | Browse Data | **Data Analysis** | Submit Data | Get Help

Activities   User Projects   All Projects   Create Project   **Workflows**   BLAST Wizard   Expert(Advanced) BLAST

## Execute Workflow: Alpha Diversity (Rohwer)

**Download documentation**

Alpha diversity is the biodiversity within a particular area, community or ecosystem, and is usually expressed as the Species richness of the area. This can be measured by counting the number of taxa (distinct groups of organisms) within the ecosystem (eg. families, genera, species). However, such estimates of species richness are strongly influenced by sample size, so a number of statistical techniques can be used to correct for sample size to get comparable values

| **Default Parameters** | Advanced Parameters |

JobName                        My Workflow 07/11/2009 0

**Circonspect**

Trim Size        100

Min Coverage     1

Repetitions      7

Size             1000

Seed             644715020

Fasta File 1        Select sequence

**Parameters**

type             power

Submit Workflow!

July 11, 2009

Workflows Menu:

- Home
- New Workflow
- Current Jobs
- Provenance Browser

Browse:

- Project
- Public
- CAMERA supported

Workflows Help

### Kepler window

File  Edit  View  Workflow  Tools  Window  Help

Components \ Data

Search

☐ Search repository

Search    Reset

Components
- Data Input
- Data Operation
- Data Output
- Director
  - DE Director
  - CT Director
  - DDF Director
  - PN Director
  - SDF Director
- File System
- General Purpose
- Workflow
Projects
Disciplines
Statistics

0 results found.

PN Director

Circ

AtachedFile
f1= $SequenceFileName       fileInput

s
Size
v
Trimsize
k
Min_metaG_coverage
z
Seed

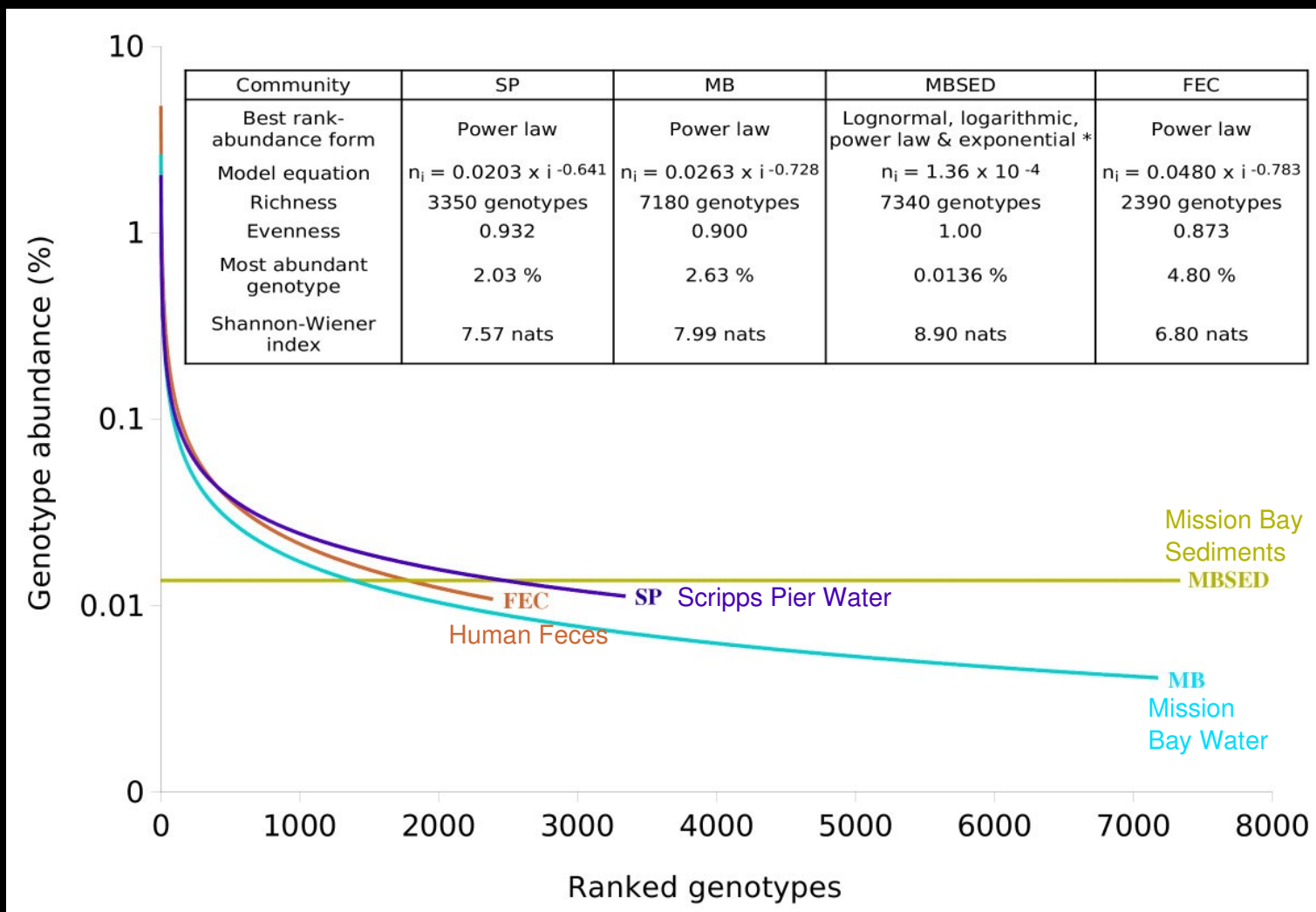- 5 -

# Environmental viral diversity

# $\alpha$-diversity of four viromes

- PHACCS estimate of community structure and diversity



| Community | SP | MB | MBSED | FEC |
|---|---|---|---|---|
| Best rank-abundance form | Power law | Power law | Lognormal, logarithmic, power law & exponential * | Power law |
| Model equation | $n_i = 0.0203 \times i^{-0.641}$ | $n_i = 0.0263 \times i^{-0.728}$ | $n_i = 1.36 \times 10^{-4}$ | $n_i = 0.0480 \times i^{-0.783}$ |
| Richness | 3350 genotypes | 7180 genotypes | 7340 genotypes | 2390 genotypes |
| Evenness | 0.932 | 0.900 | 1.00 | 0.873 |
| Most abundant genotype | 2.03 % | 2.63 % | 0.0136 % | 4.80 % |
| Shannon-Wiener index | 7.57 nats | 7.99 nats | 8.90 nats | 6.80 nats |

*Angly et al. 2005*

# Marine latitudinal gradient of diversity

- The latitudinal richness gradient is the most documented pattern: Higher diversity close to the tropics

- Affects macroorganisms (*Hillebrand et al. 2004*), microorganisms (*Pommier et al. 2007, Fuhrman et al. 2008*)

- Affects viruses?

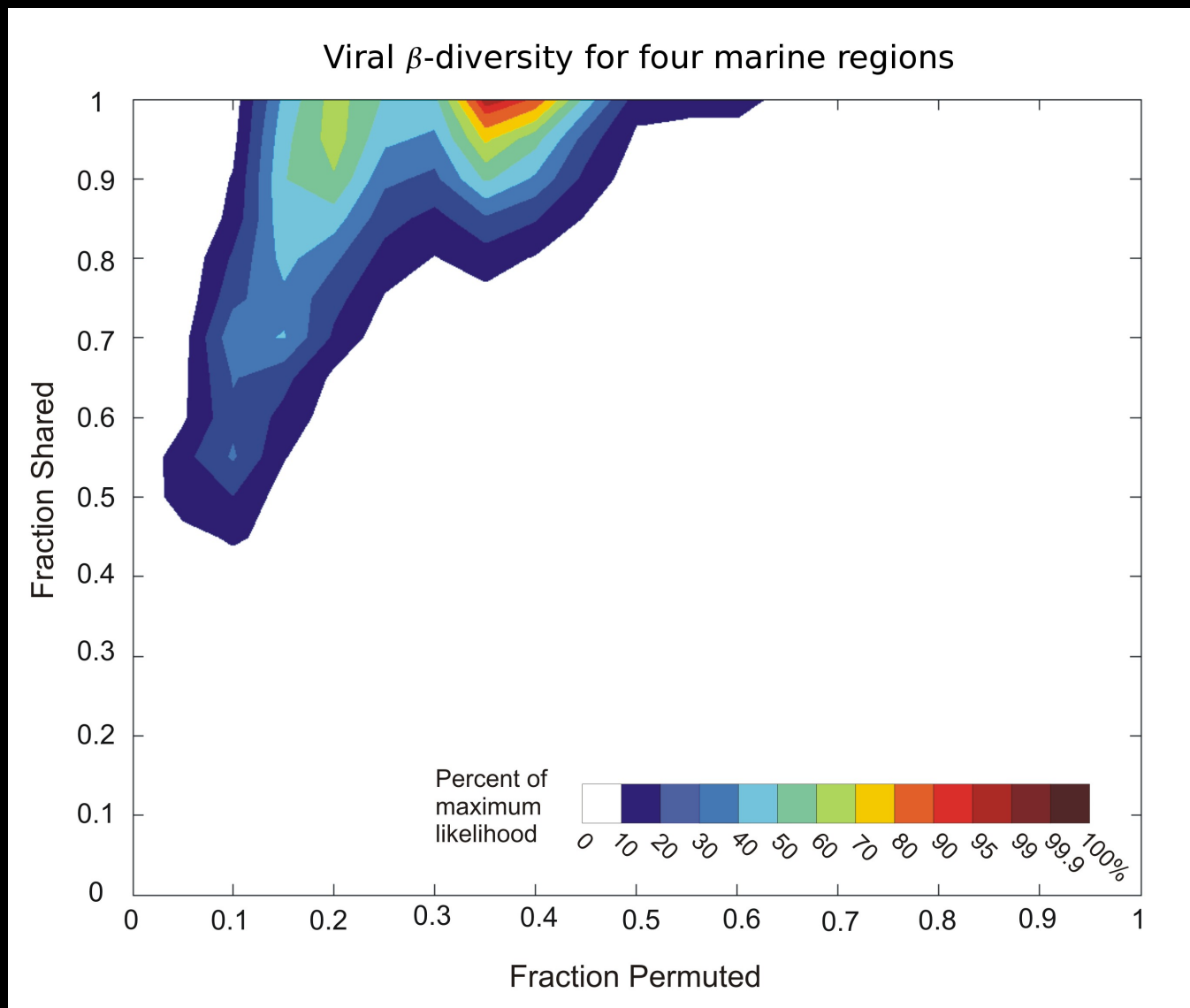- How universal are the rules that govern the distribution of life on Earth?

Viral richness in the oceans



**Arctic** 532

**SAR** 5,140

**GOM** 15,400

**North-South latitudinal gradient**

*Angly et al. 2006*

# Marine $\beta$-diversity

- Viruses are dispersed world-wide

- "Everything is everywhere"

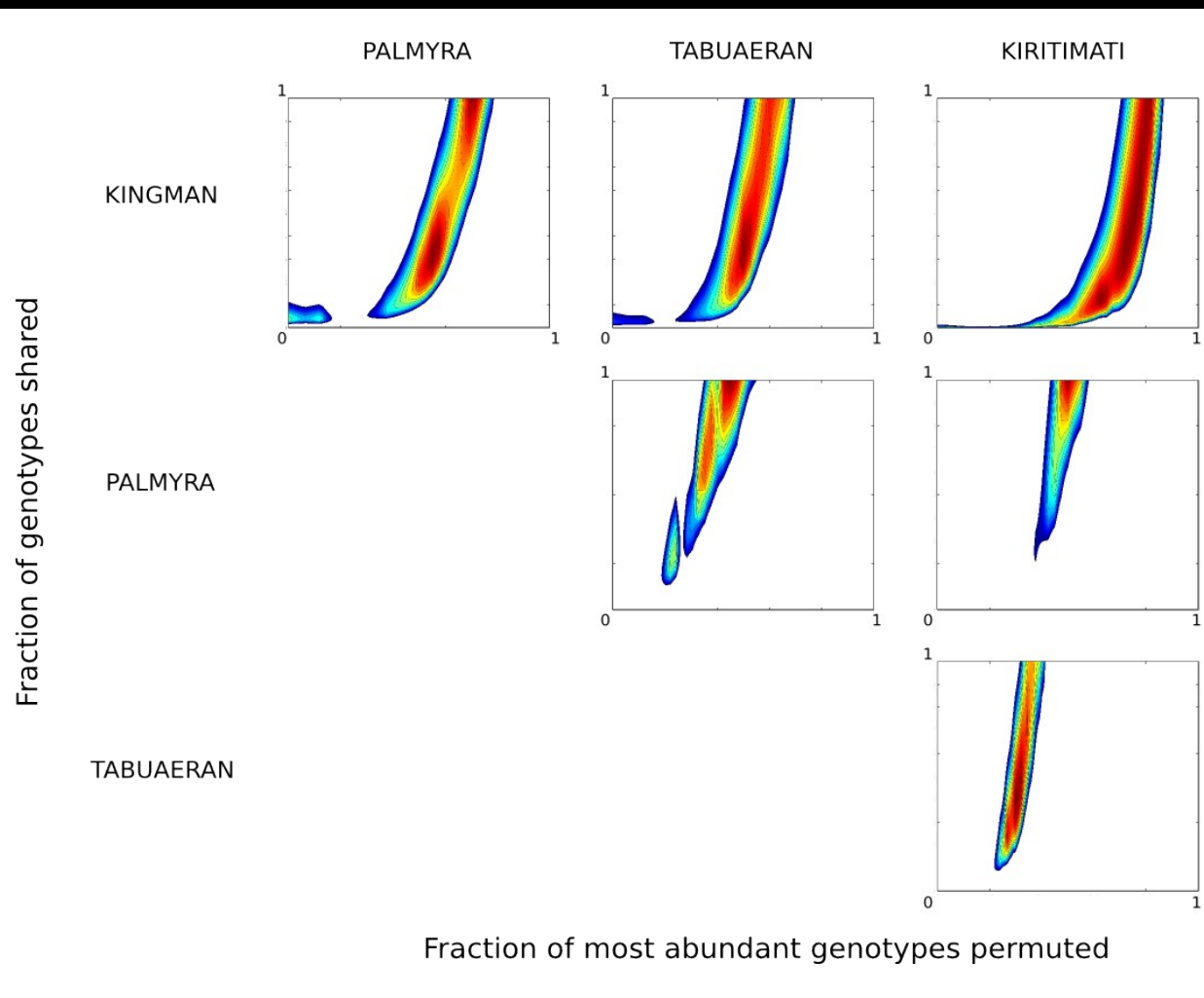- Marine viruses are cosmopolitan but the environment selects!

Viral $\beta$-diversity for four marine regions



*Angly et al. 2006*

# Diversity in the Line islands coral reefs

| Atoll | Human population | Viral richness |
|---|---|---|
| Kingman | 0 | 8,380 |
| Palmyra | 20 | 17,100 |
| Tabuaeran | 1,000 | 24,800 |
| Kiritimati | 5,100 | 102,000 |

Viral $\beta$-diversity



*Sandin et al 2008*          *Plot by Steve Rayhawk*

# Conclusions

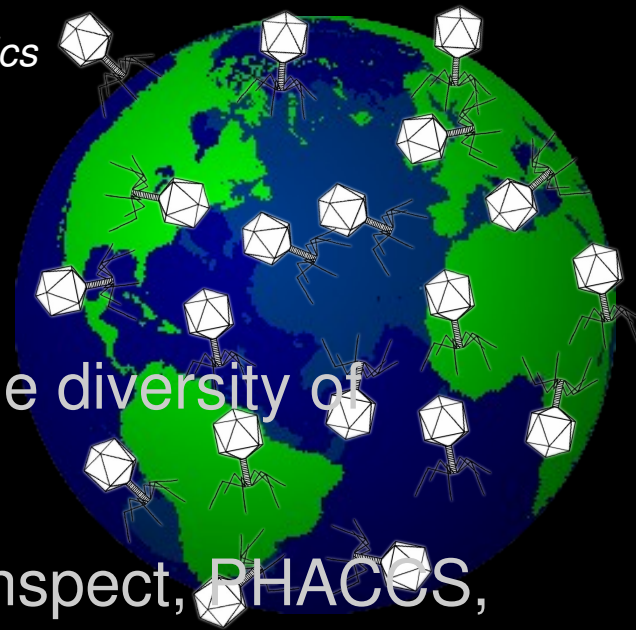Metagenomics is well suited to investigate the diversity of uncultured viral communities

Novel computational methods: GAAS, Circonspect, PHACCS, MaxiPhi

Publicly available tools and integration into easy-to-use software workflow

Diversity methodology does not require similarities to databases

Viral diversity may follow the same patterns of diversity as microorganisms and macroorganisms

As more and more viral metagenomes are sequenced, the metagenomic diversity workflow will be used to analyze the global virome and estimate Earth's total viral richness

# Viral metagenome locations

# Acknowledgments

San Diego State University

Rob Edwards
Peter Salamon
Ben Felts
Jim Nulton
Joseph Mahaffy
Robert Schmieder
Liz Dinsdale
Alejandra Prieto-Davo

NSF Biocomplexity

Gordon and Betty Moore Foundation

RohwerLab

Forest
Bahador
Beltran
Dana
John
Katie
Linda
Liz
Mark
Matt
Mike
Yan Wei

Anna
Becky
Bethany
Betty
Christelle
Cynthia
David
Danielle
Elysa
Emiko
Fairoz
Karin
Jennifer
John
Marina
Megan
Morrigan
Mya
Neilan
Olga
Priscila
Selina
Steve
Veronica
Yanan
...

*Characterization of environmental viral diversity using metagenomics*

# Characterization of environmental viral diversity using metagenomics