# Software for Metagenomics and Metadata Standards

Iddo Friedberg

CAMERA / CALIT2

University of California San Diego

# What is Metagenomics?

- Culture-free approach to study microbial communities
  - < 1% of microbes can be cultured
  - DNA directly isolated from environmental sample and sequenced
- Examining genomic content of organisms in community/environment to better understand:
  - Diversity of organisms
  - Their roles and interactions in the ecosystem

# Metagenomics: Appy Genomics to Populations & Communities



| | Individual | Population | Community |
|---|---|---|---|
| **Ecology** | Physiology: Differential gene expression in response to change | Demographics: Birth, death, immigration, emigration | Community ecology: Interspecific interactions that shape community structure and function |
| **Genomics** | Fine-scale mapping of individual genomes | Population genomics: Comparative genomic analyses to assess variation | Metagenomics: Genetic potential of collective members of community |
| **Genetics** | Bacterial genetics: Role of genes under various conditions | Population genetics: Allele frequency distribution | Community genetics: Interplay between genetic composition of community and ecological community properties |

Little AEF, et al. 2008.
Annu. Rev. Microbiol. 62:375–401

# What Can We Learn?

- **Taxonomic content**: Taxon diversity in a habitat (using taxonomic markers)

- **Functional content**: biological functions, qualitative and quantitative profiles

- **Coping with the environment**: differences in functional content between habitats

- **Decompose the biotic / abiotic elements in a habitat**: metadata analysis

- …

# Some Metagenomic Studies

- **Bacterial rhodopsins** (Beja *et al.* 2000)
- **Acid mine drainage study** (Tyson *et al.* 2004)
- **Sargasso sea study** (Venter *et al.* 2004)
- **Wisconsin soil study** (DeLong *et al.* 2006)
- **Termite Hindgut** (Warnecke *et al.* 2007)
- **Human Obesity** (Turnbaugh *et al.* 2006)

… and many, many others

# What is CAMERA?

- <u>C</u>ommunity <u>C</u>yberinfrastructure for <u>A</u>dvanced <u>M</u>arine Microbial <u>R</u>esearch and <u>A</u>nalysis

- 7 year $24.5 mil grant from the Moore foundation

- Goal: build a community computational resource for researchers in metagenomics

- "Cyberinfrastrcture": hardware, software & data

# What is CAMERA? Hardware

512 CPU, 200 TB, 5 TFlops

# What is CAMERA? Data

- Metagenomic sequence data are:
  - Voluminous
  - Noisy
  - Partial
- At the very least they should be:
  - Standardized for processing
  - Associated with Metadata

# Why Metadata?

- Microbial communities are affected by and affect their habitats

- Therefore habitat data, in addition to sequence data, is crucial for an **environmental** genomic picture

- Also, sample condition data is needed for reproducibility

# Why Metadata?

- **Microbial communities are affected by and affect their habitat**

- **Sequence information + metadata = whole picture**

- Habitat Type
- Geographic Location (large area)
- Sample Location (smaller area)
- Country
- Filter Size
- Latitude (exact location)
- Longitude (exact location)
- Depth
- Wat. Dep.
- Chlorophyll
- Oxygen
- Fluor.
- Salin.
- Temp
- Trans.
- BioMass
- Inorg. Carbon
- Inorg. Phospate
- Org. Carbon
- Nitr.
- # Pooled
- # Sampled
- Volume
- Coll. Date

# Why Metadata?

- Microbial communities are affected by and affect their habitat
- Sequence information + metadata = whole picture

Habitat

Altitude

pH

Country

Temperature

Filter size

# Data Standards and Data Acquisition

- Minimal Information for (Meta)Genomic Sequences: MIGS/MIMS

- A Metadata standard, developed by the Genomics Standards Consortium

  - Controlled vocabularies e.g. EnvO, PATO, CABRI
  - Common language: GCDML

- Submissions shall comply with a MIMS/MIGS core, but any metadata can be entered via keywords and free text

- Different metadata submission forms for different habitats: (water, soil, air, hosts)

# Standards Compliance: MIMS/MIGS

# Standards Compliance: MIMS/MIGS

# Standards Compliance: PATO

# Standards Compliance: Controlled Vocabularies

# Things to Think About

- Visualization:
    - How do we look at "disembodied" sequence data?
    - "Fragment recruitment" track
    - Visualization of sequence data <--> metadata associations
- Database: association of metadata and sequence data; queries by metadata
-

# Would you Like to Know More?

- The Genomic Standards Consortium.
  MIMS & GCDML:
  http://gensc.org
- BioMIST (soon):
  http://sourceforge.net/projects/biomist/
- CAMERA: http://camera.calit2.net
- Me: http://iddo-friedberg.net

# Would you Like to Know More?

- The Geno... ...nsortium. MIMS & G... http://gens...
- BioMIST ( http://sour... ...ts/biomist/
- CAMERA: ...t2.net
- Me: http://...

# Thanks

- CAMERA data acquisition team:
    - Brian Fox
    - Shulei Sun
    - Jing Chen
    - Laurence Bohannan
    - Jeffrey Grethe
- Genomic Standards Consortium
    - Dawn Field
    - George Garrity
    - Renzo Kottman