

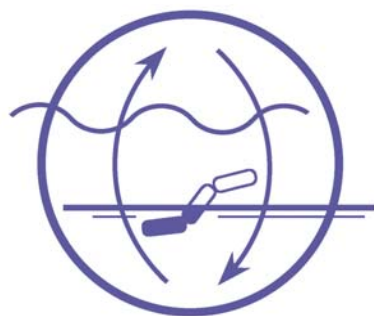
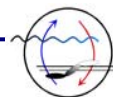
# Next Generation Sequencing Technologies in Microbial Ecology



**Frank Oliver Glöckner**



JACOBS  
UNIVERSITY



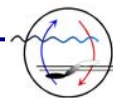
# Max Planck Institute for Marine Microbiology

Investigation of the role, diversity and features of microorganisms

Interactions with physical and chemical processes in marine and other aquatic habitats



Founded 1992 in Bremen, Germany



## Marine Microbiology at MPI – a Holistic Approach

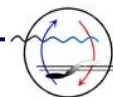
**Who** is out there and

**How much** of which kind?

**What** are they doing and

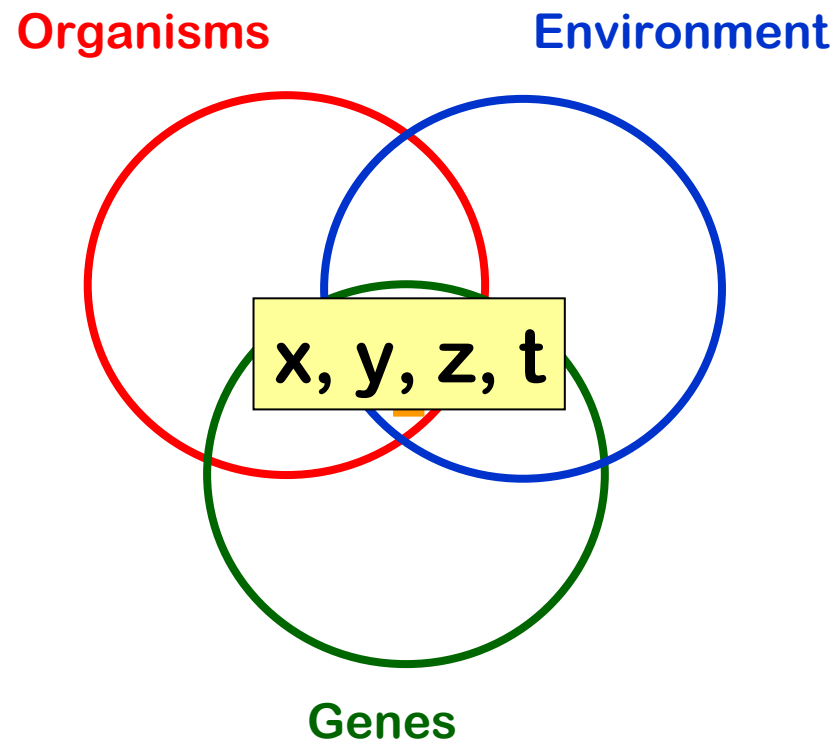
Under **which conditions** are they doing **what**?

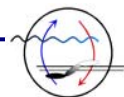




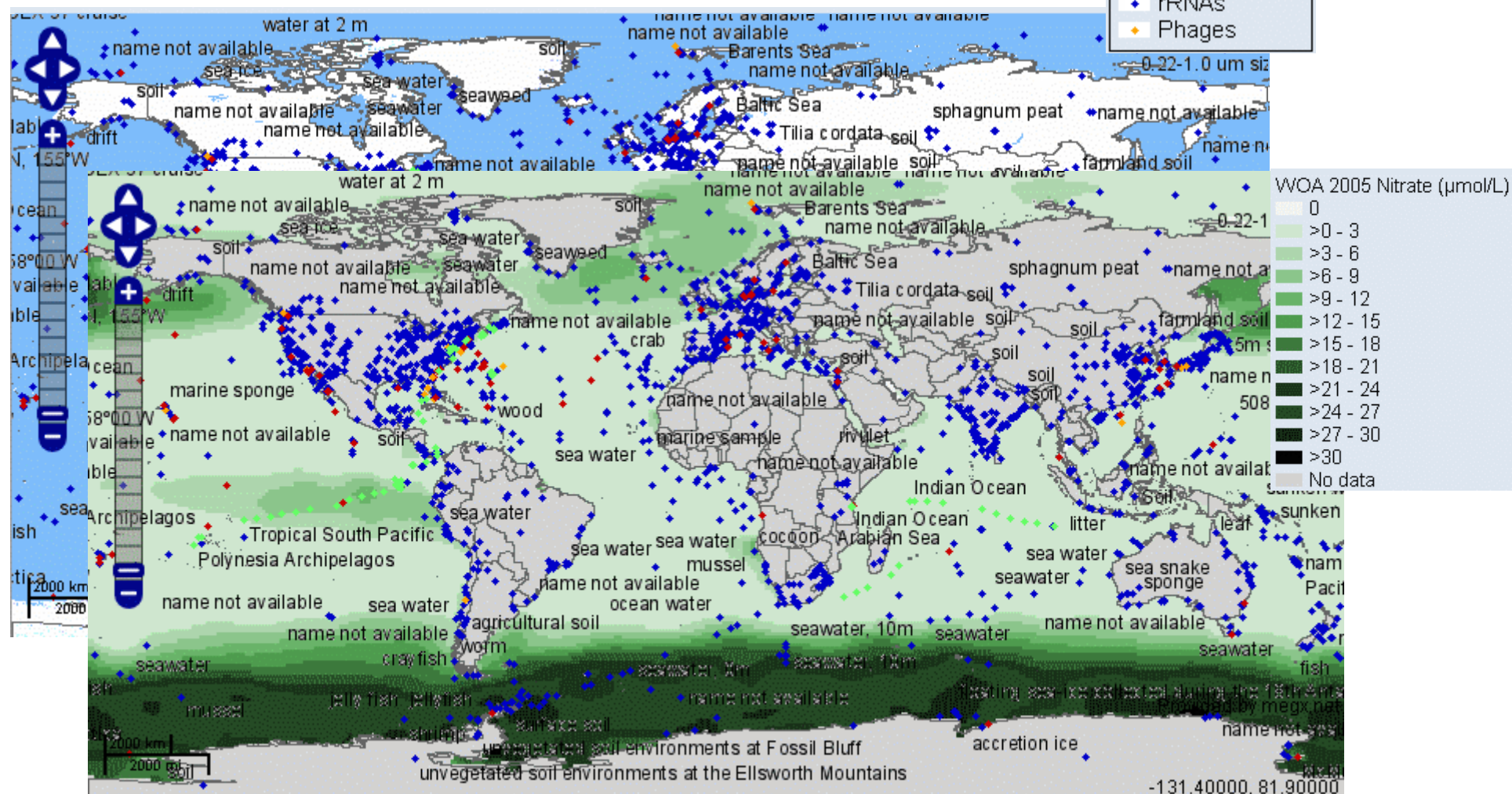
## Promise of NGS: Much Denser Network of Data

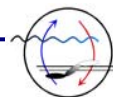
- ▶ **Phylogenetic diversity**
    - Qualitative data
    - Quantitative data
  
  - ▶ **Functional diversity**
    - Functional inventory
    - Operon structures
    - Expression profiles
  
  - ▶ **Environmental descriptors**
- > Integrated datasets



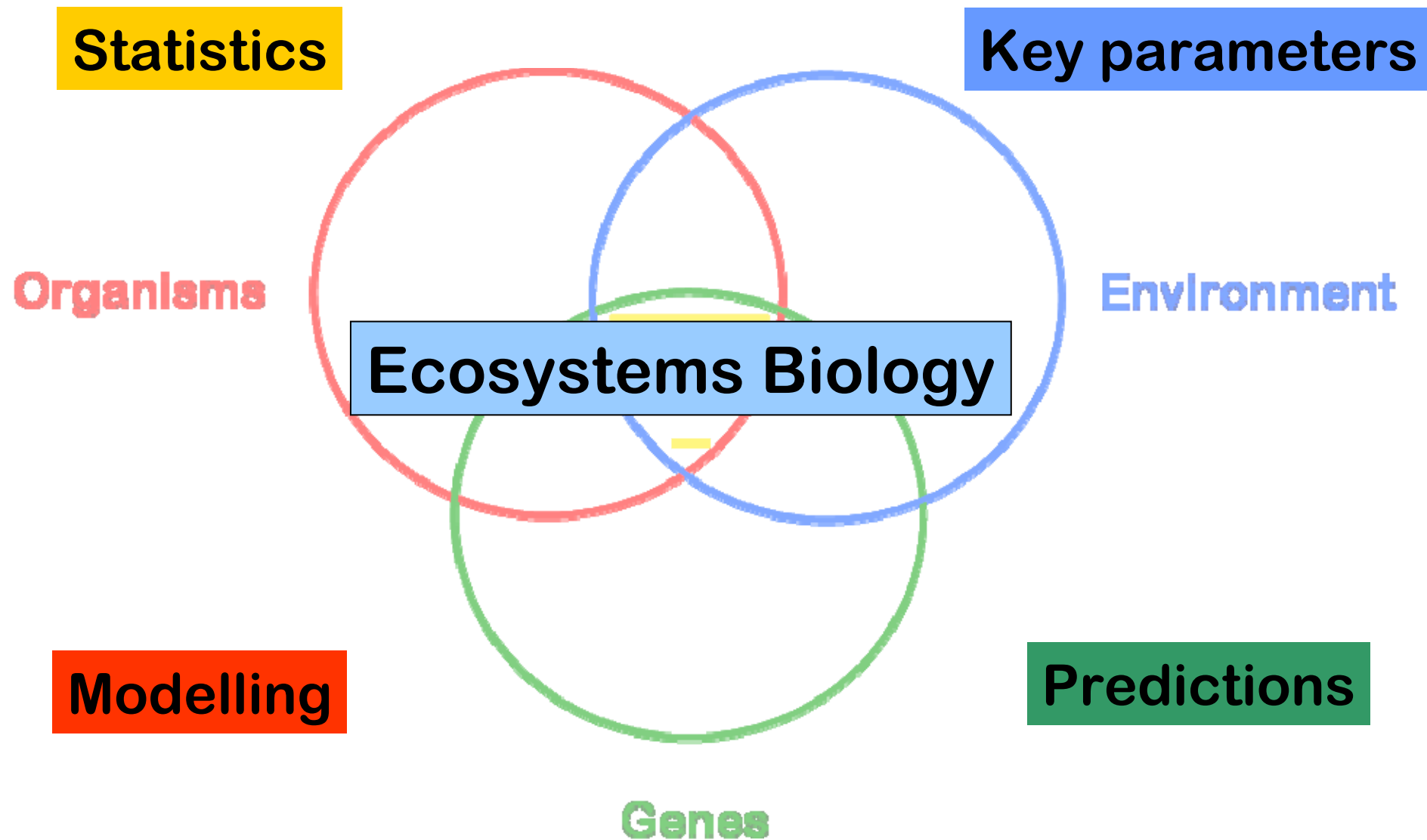


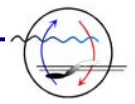
# Data Integration





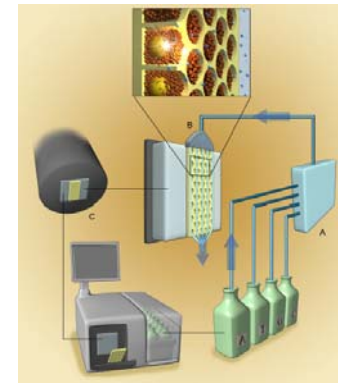
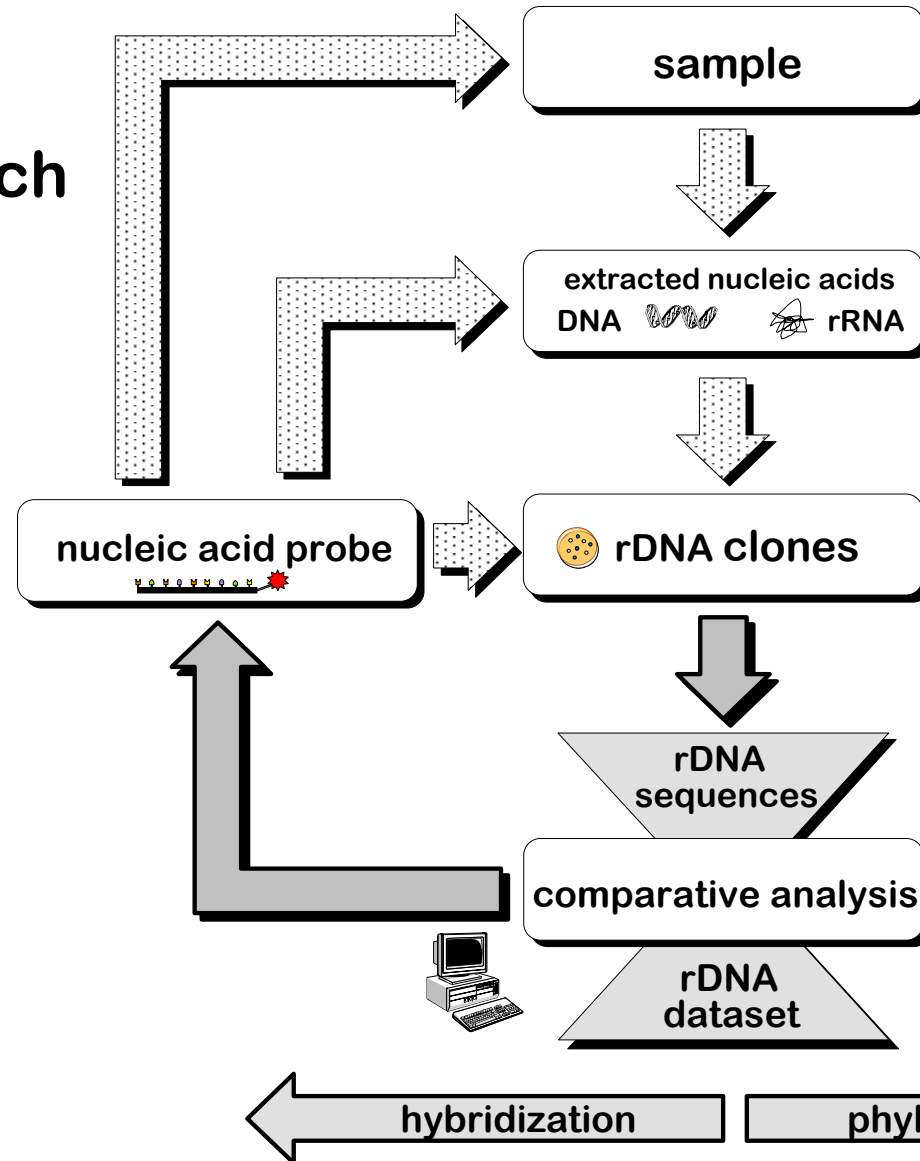
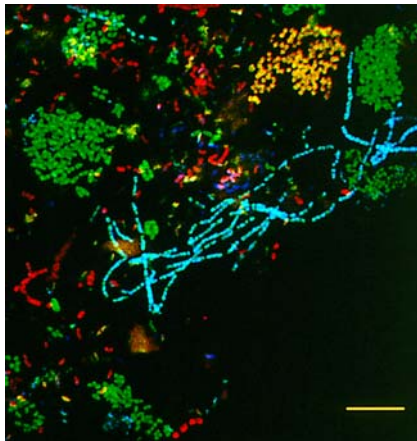
# Ecological Genomics – The Vision



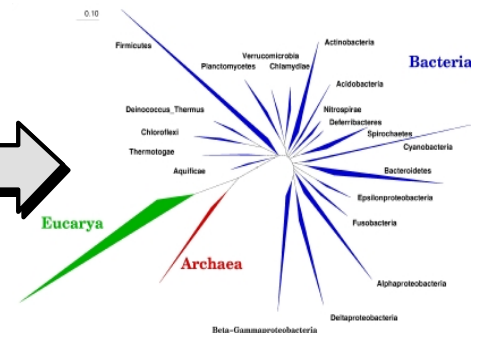


# Ribosomal RNA as a universal marker gene

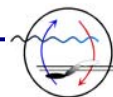
## Full cycle rRNA-approach



Pyrosequencing

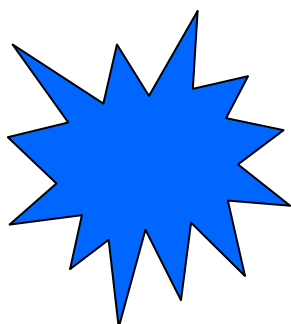


Amann, 1995

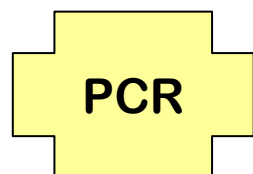


# Diversity Analysis

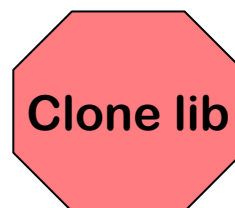
Sample



High diversity

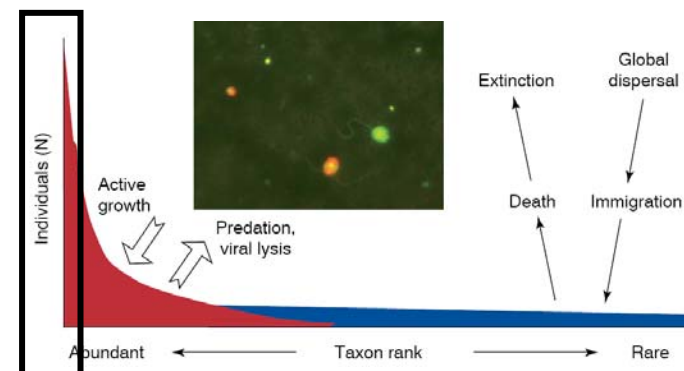


PCR

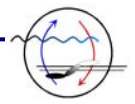


Clone lib

100-500  
2-3 month

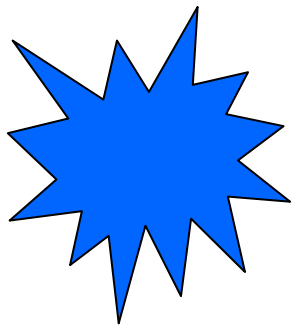






# Diversity Analysis

Sample



High diversity

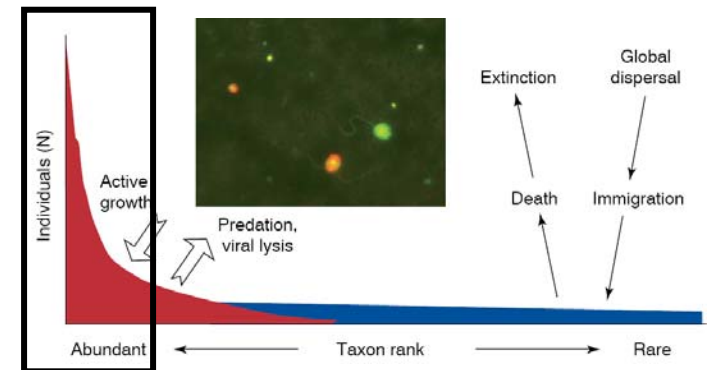
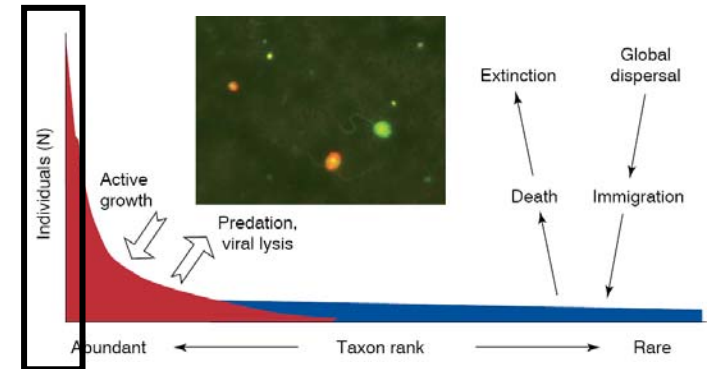
Clone lib

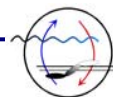
100-500  
2-3 month

PCR

Tags

10,000-50,000  
1 week





## Problems

### ▶ Processing the data

Your next-generation sequencing data just arrived...

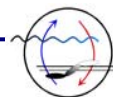
# Now What?



Richard was having a great day, until the arrival of his **next-generation data files**.

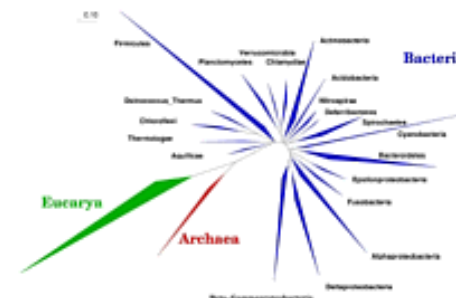
### ▶ Accuracy/Quantitative?

- DNA/RNA extraction
- Multiple Operons
- Technical replicates
- Noise (sequencing errors...)



# SILVA Databases Specifications

- ▶ **Comprehensive & Aligned**
  - *Bacteria, Archaea, Eukarya*
  - SSU, LSU
  - Regularly updated

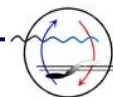


- ▶ **Quality first**
  - Quality management
  - Transparent process documentation

Length	633						
Quality	<table> <tr> <td>Sequence</td> <td><div style="width: 100%; height: 10px; background-color: green;"></div></td> </tr> <tr> <td>Alignment</td> <td><div style="width: 100%; height: 10px; background-color: yellow;"></div></td> </tr> <tr> <td>Pintail</td> <td><div style="width: 100%; height: 10px; background-color: red;"></div></td> </tr> </table>	Sequence	<div style="width: 100%; height: 10px; background-color: green;"></div>	Alignment	<div style="width: 100%; height: 10px; background-color: yellow;"></div>	Pintail	<div style="width: 100%; height: 10px; background-color: red;"></div>
Sequence	<div style="width: 100%; height: 10px; background-color: green;"></div>						
Alignment	<div style="width: 100%; height: 10px; background-color: yellow;"></div>						
Pintail	<div style="width: 100%; height: 10px; background-color: red;"></div>						

- ▶ **Integrative**
  - Nomenclature
  - Taxonomy
  - Cultured, Typestrains
  - Habitat (r100)

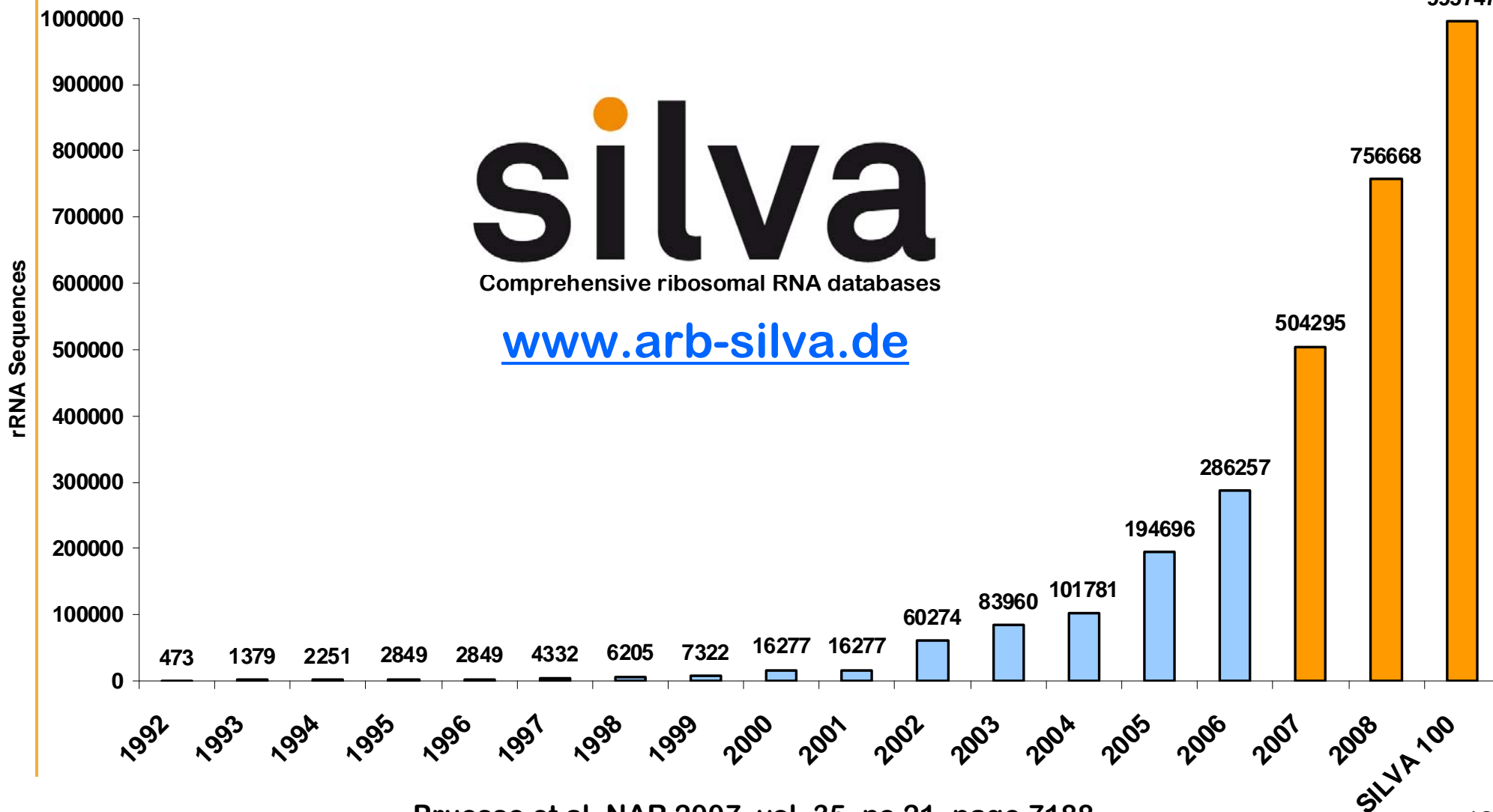




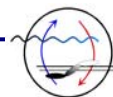
# Growth of rRNA databases (RDP & SILVA)

Growth of SSU ribosomal RNA databases (RDP II & SILVA)

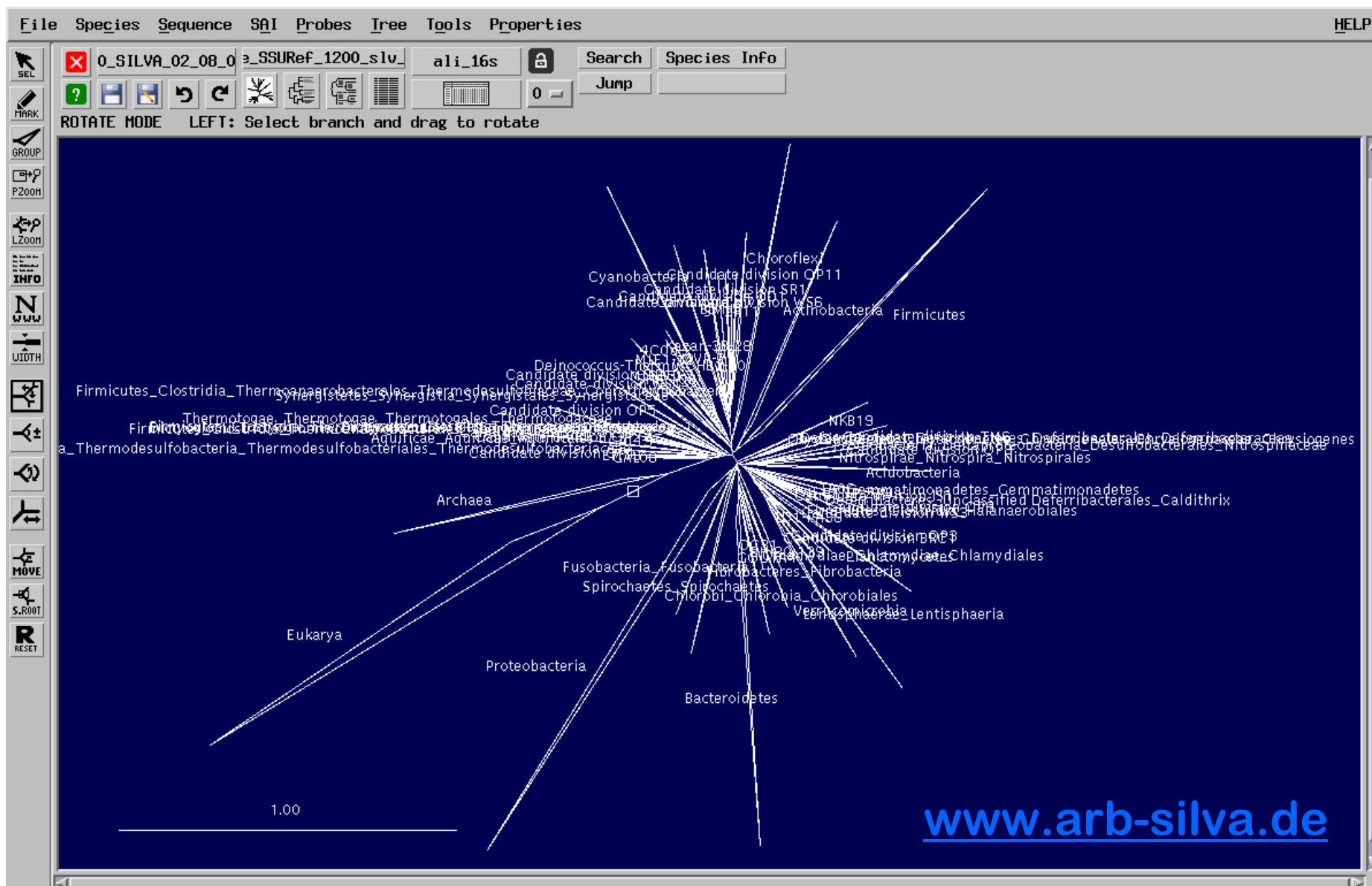
[www.arb-silva.de](http://www.arb-silva.de)

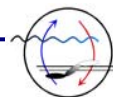


Pruesse et al. NAR 2007, vol. 35, no 21, page 7188

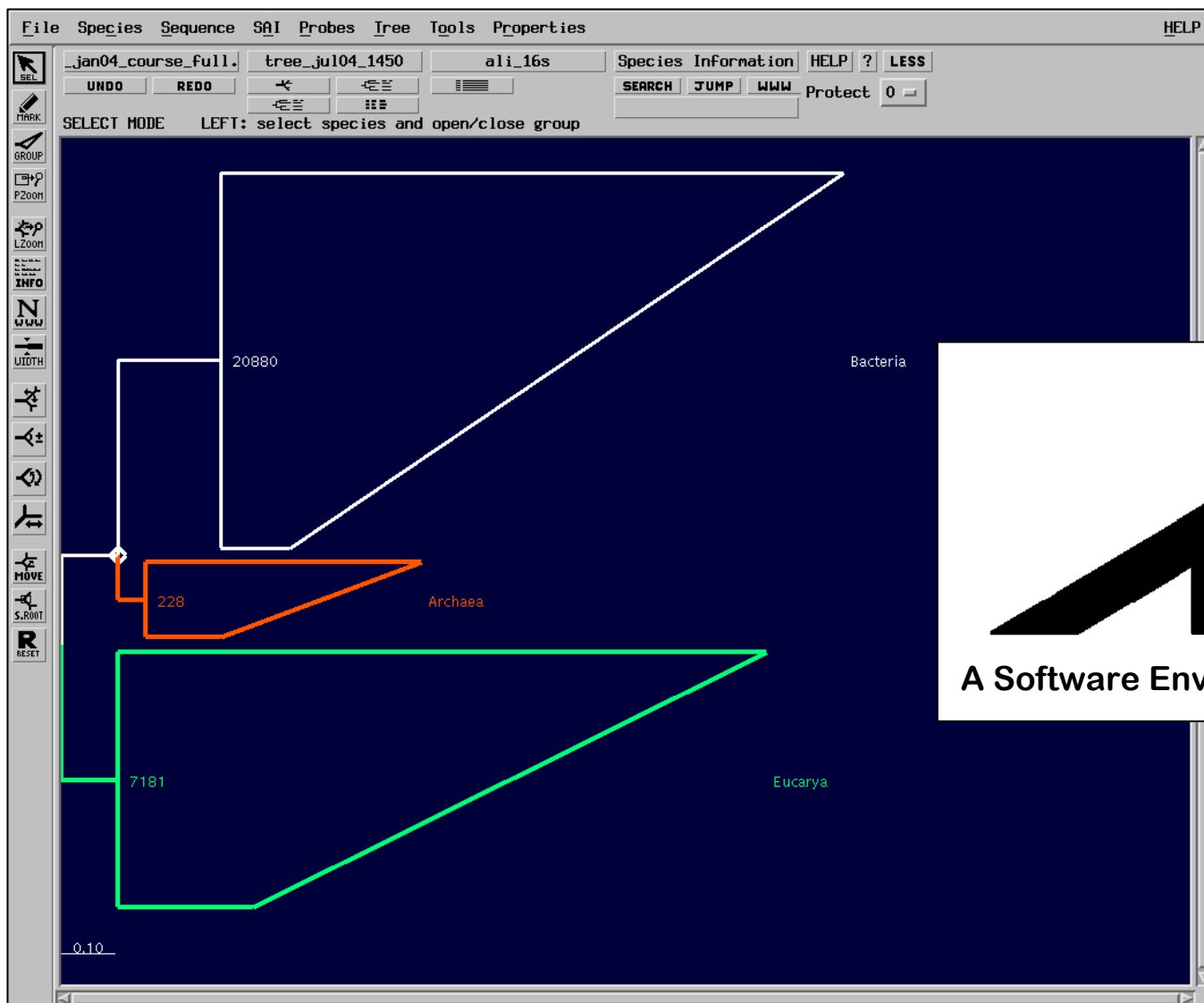


# SILVA SSURef 100: Fully classified guide tree





# ARB Software Suite

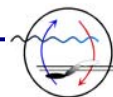


[www.arb-home.de](http://www.arb-home.de)



Ludwig et al. NAR, 2004

**ARB 5.0, 64 bit version released on 04. September 2009**



## Problems

### ▶ Processing the data

Your next-generation sequencing data just arrived...

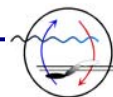
# Now What?



Richard was having a great day, until the arrival of his **next-generation data files**.

### ▶ Accuracy/Quantitative?

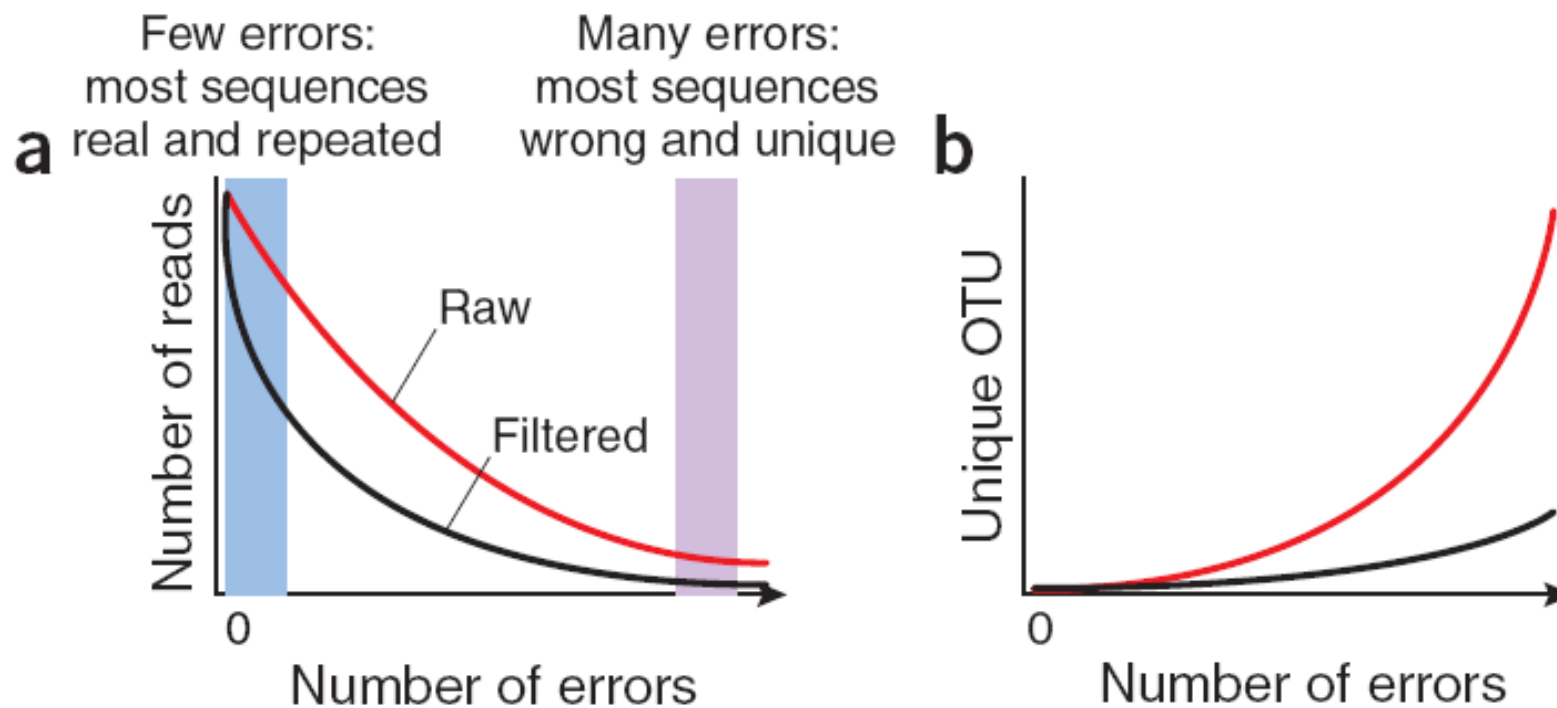
- DNA/RNA extraction
- Multiple Operons
- Noise (sequencing errors...)
- Technical replicates



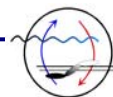
# Accuracy

## The 'rare biosphere': a reality check

Reeder and Knight, 2009, Nature Methods vol. 6, no. 9, p. 636

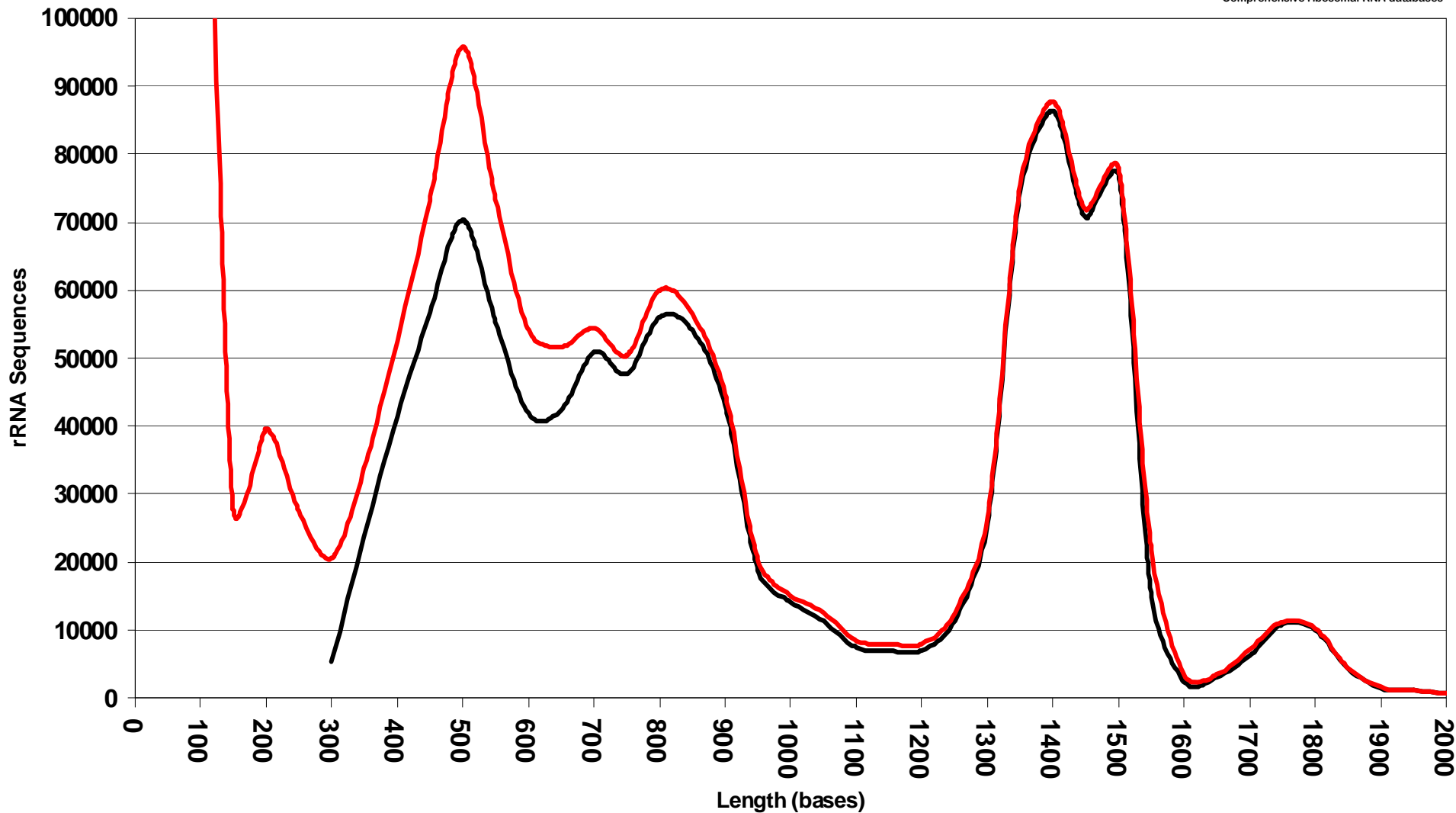


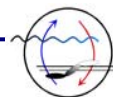




SILVA SSUParc 100, Sequence Length Distribution  
www.arb-silva.de

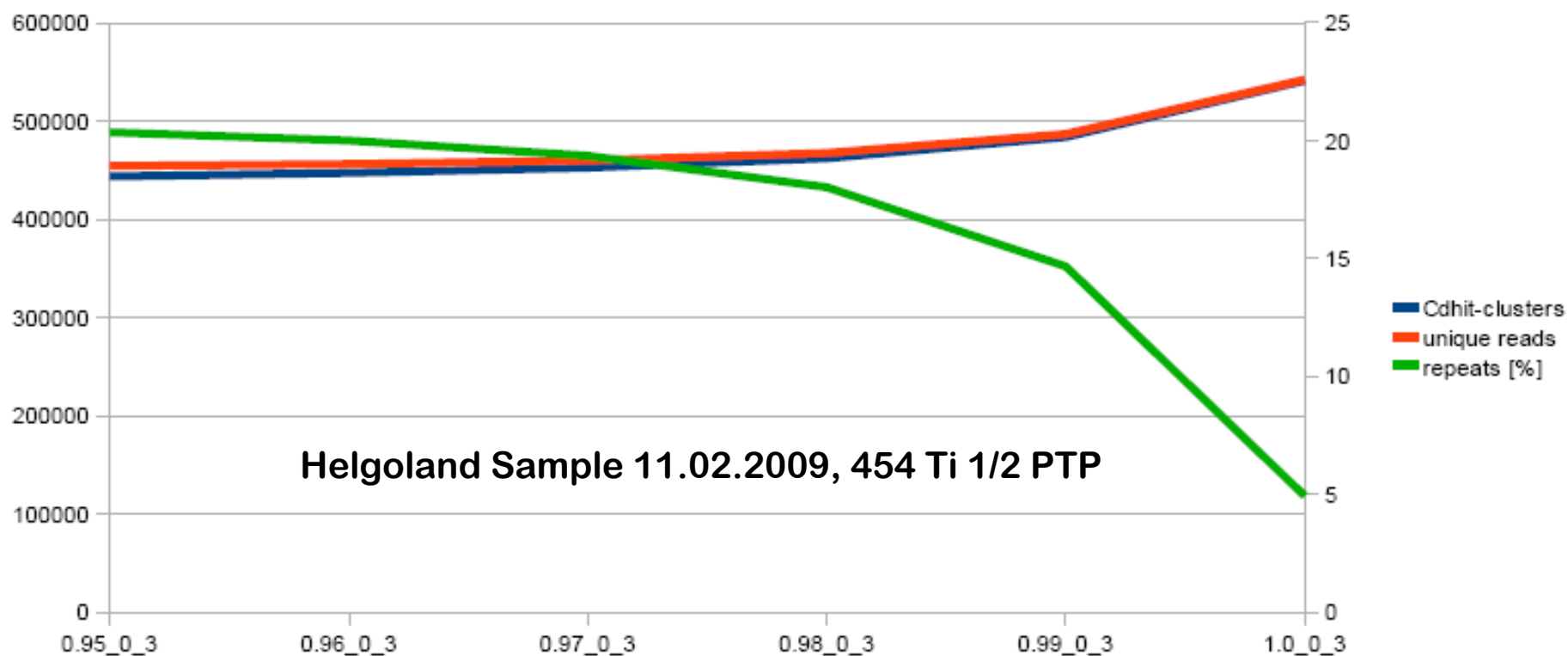
**silva**  
Comprehensive ribosomal RNA databases

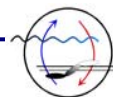




# Technical Replicates I

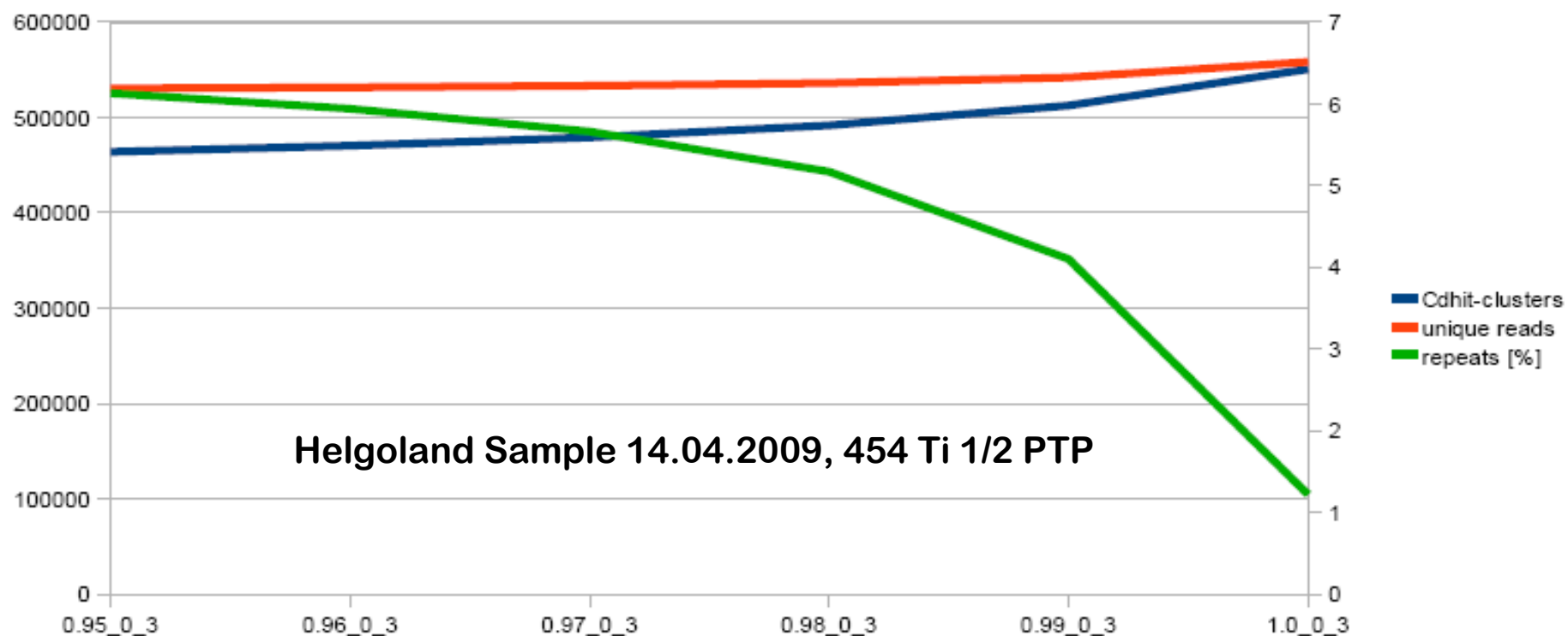
name	Cdhit-clusters	unique reads	repeats [%]	identity cutoff	length difference req.	initial base pair req.
0.95_0_3	443508	454154	20.35	0.95	0	3
0.96_0_3	447004	456150	20	0.96	0	3
0.97_0_3	452473	459824	19.35	0.97	0	3
0.98_0_3	461977	467396	18.02	0.98	0	3
0.99_0_3	483491	486594	14.66	0.99	0	3
1.0_0_3	540920	542038	4.93	1	0	3

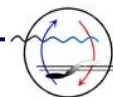




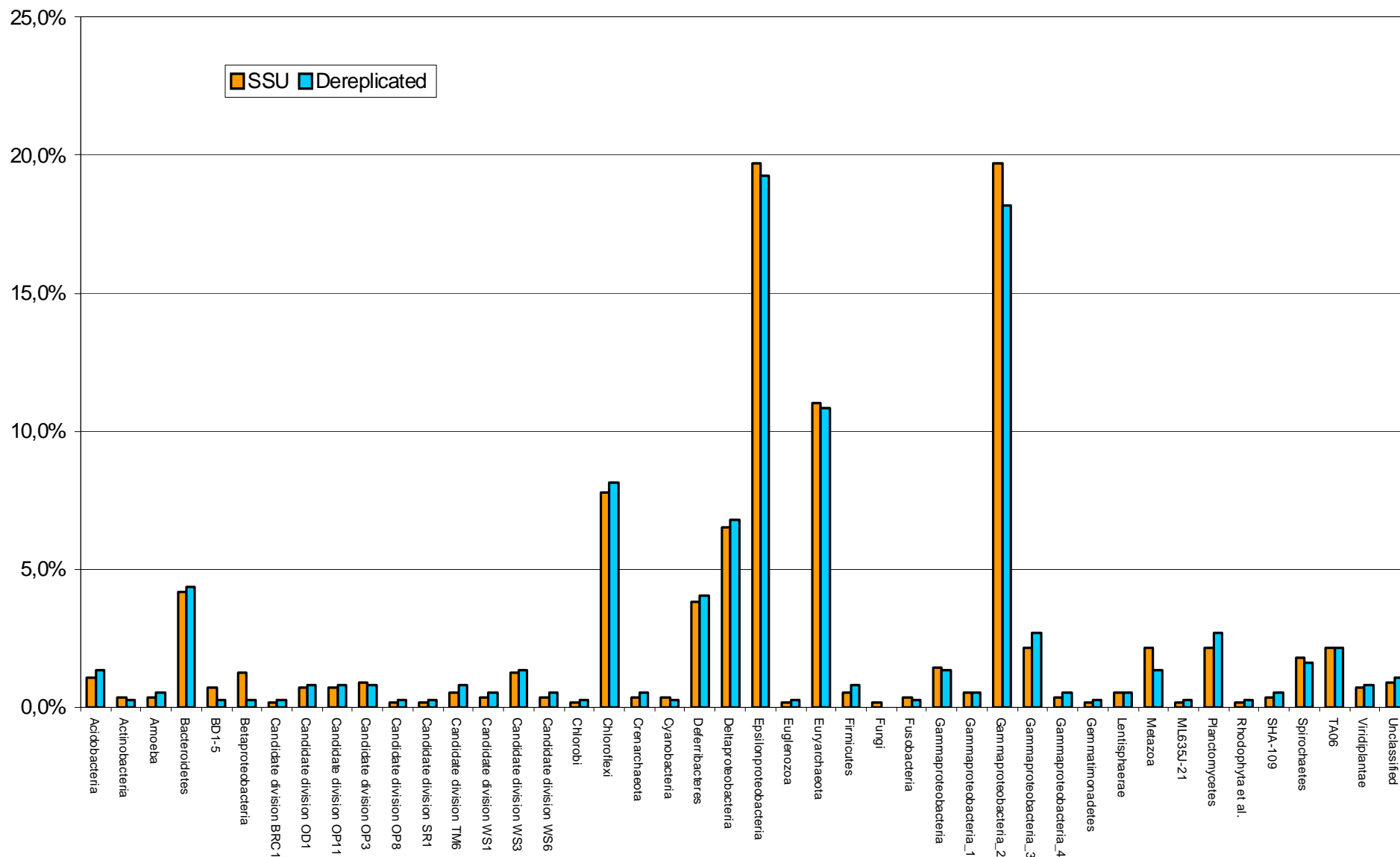
## Technical Replicates II

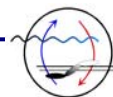
name	Cdhit-clusters	unique reads	repeats [%]	identity cutoff	length difference req.	initial base pair req.
0.95_0_3	463763	530516	6.13	0.95	0	3
0.96_0_3	470233	531581	5.94	0.96	0	3
0.97_0_3	479141	533202	5.66	0.97	0	3
0.98_0_3	491641	535978	5.17	0.98	0	3
0.99_0_3	512518	542010	4.1	0.99	0	3
1.0_0_3	550538	558287	1.22	1	0	3



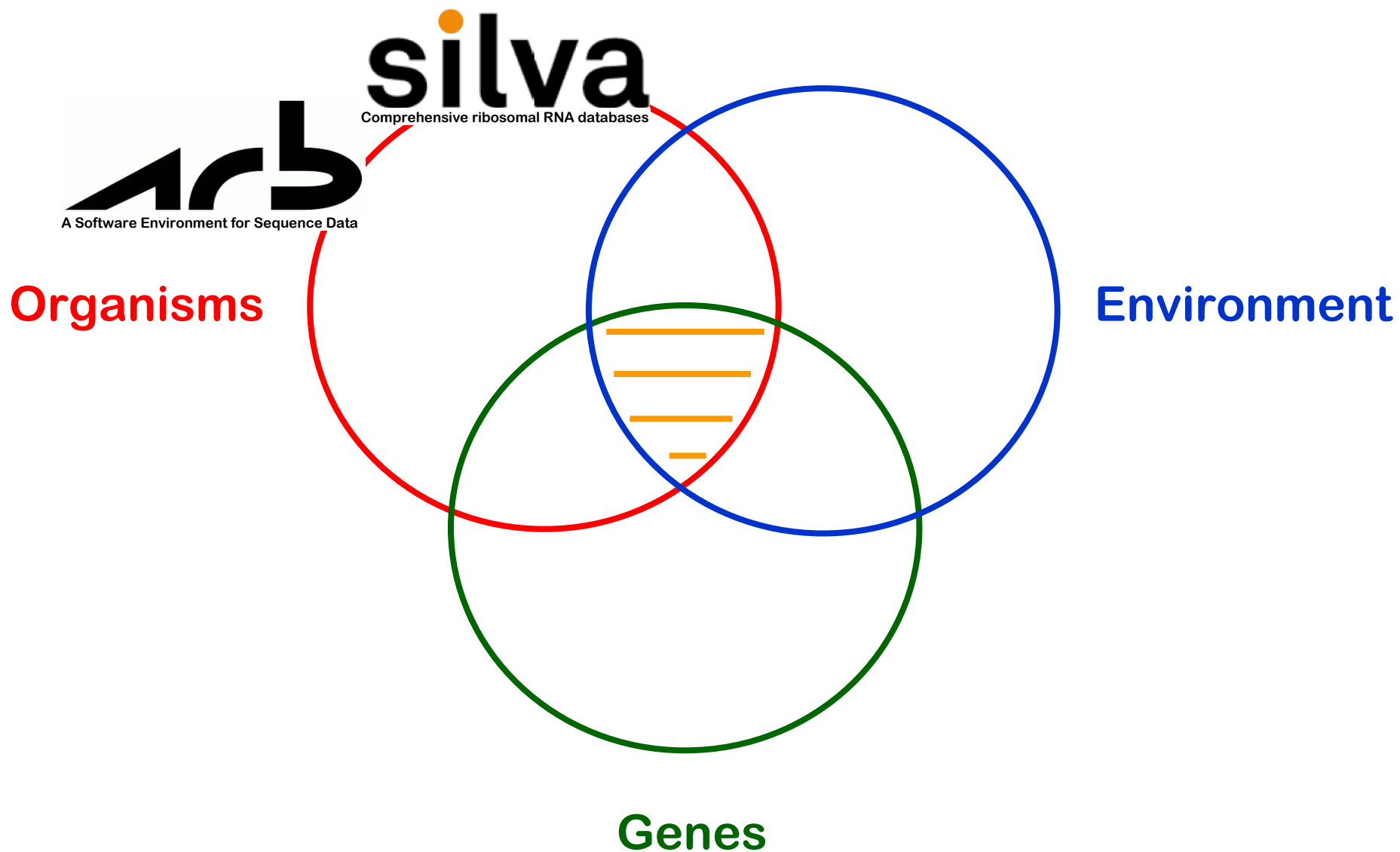


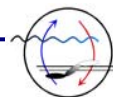
# Technical Replicates III - Dereplication



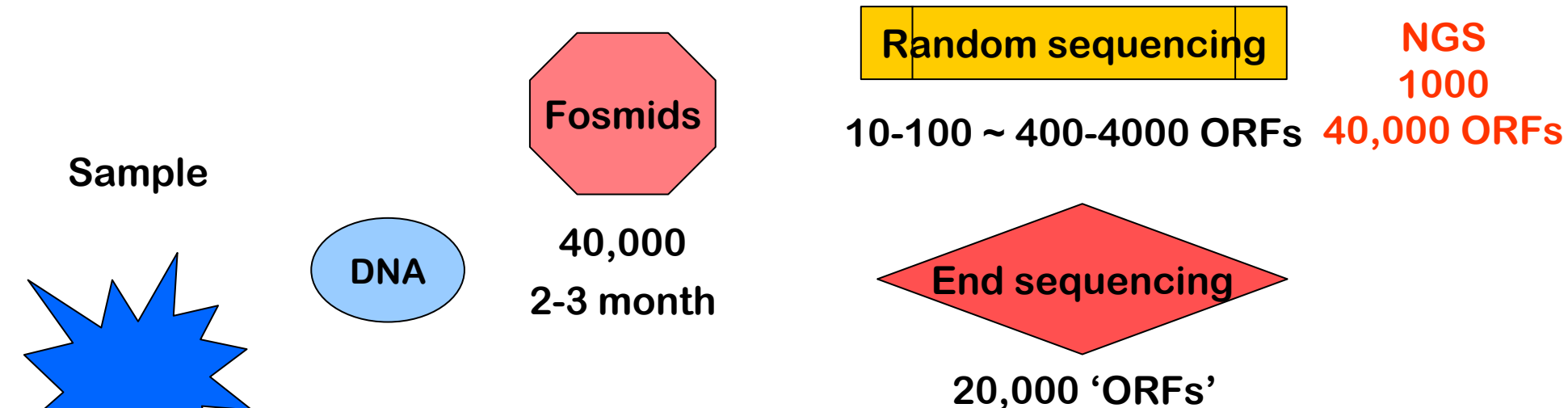


# A Bioinformatic Workbench for Ecological Genomics

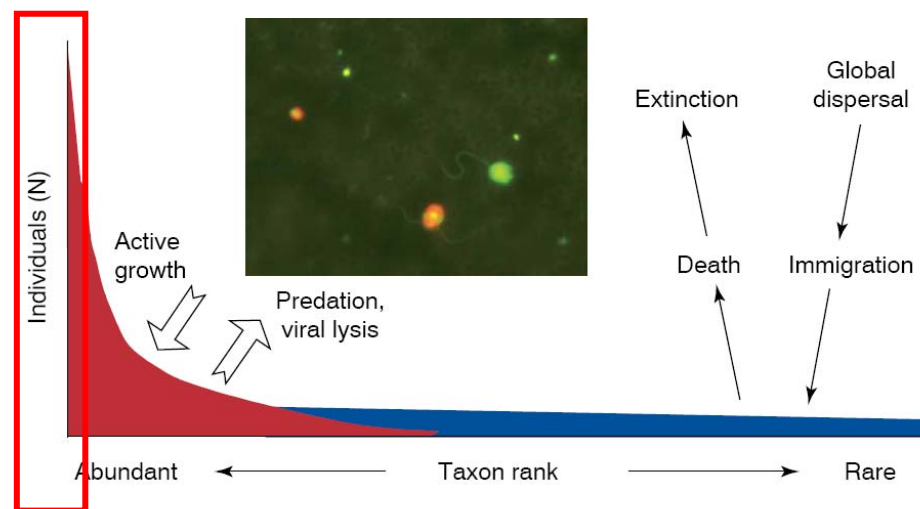


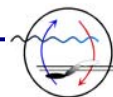


# Functional Diversity Analysis

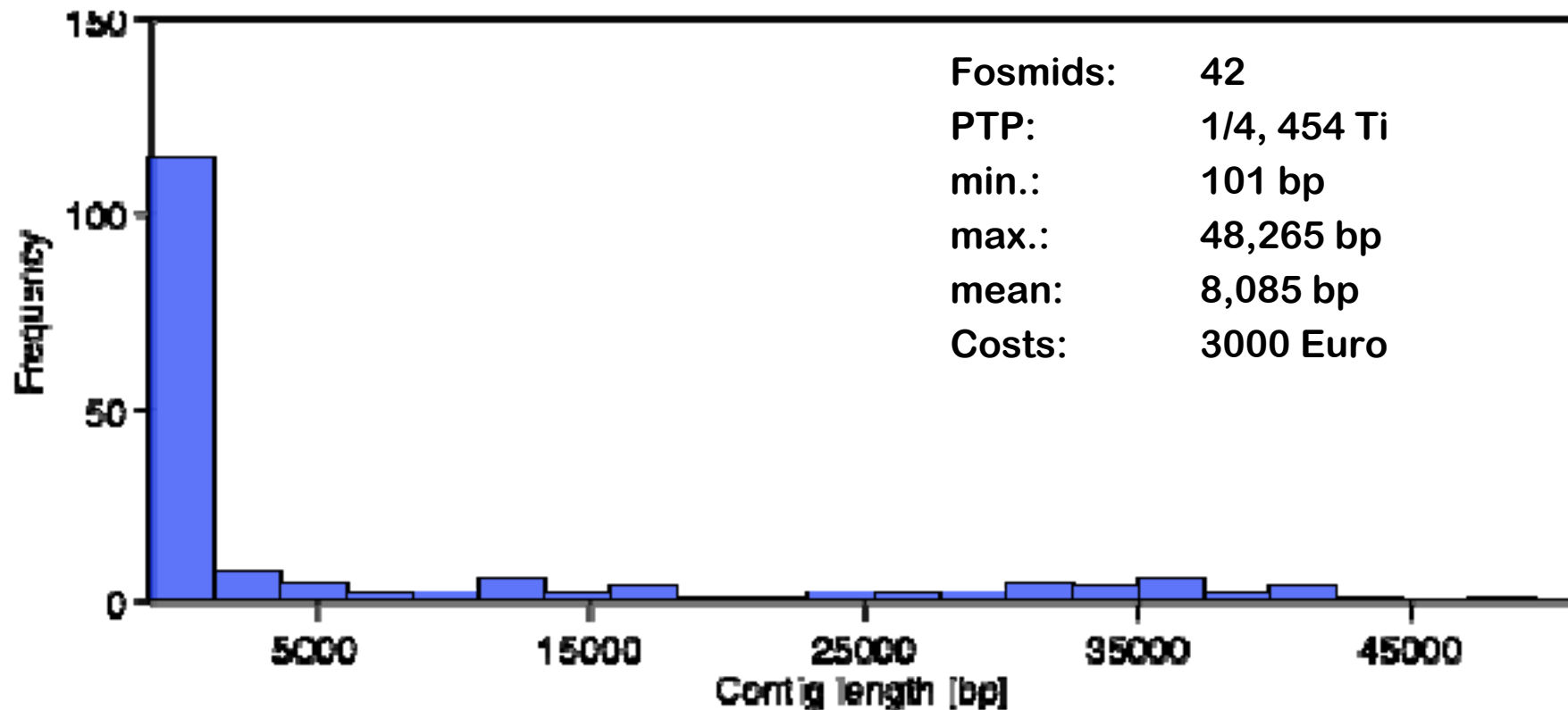


**High diversity**



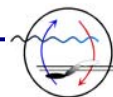


# Fosmid Sequencing



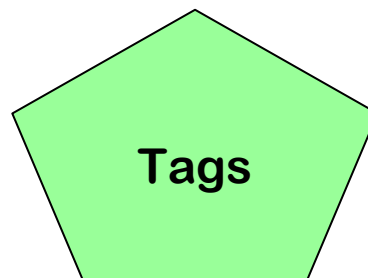
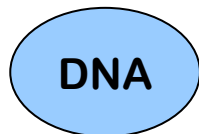
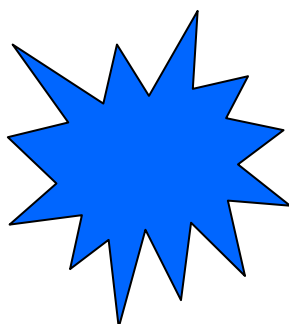
Fosmids: 42  
PTP: 1/4, 454 Ti  
min.: 101 bp  
max.: 48,265 bp  
mean: 8,085 bp  
Costs: 3000 Euro

48 contigs > 10 kb

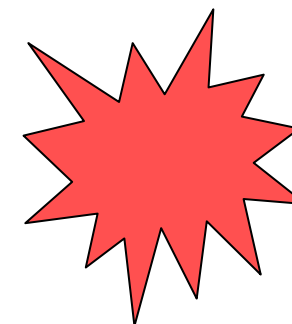


# Functional Diversity Analysis

Sample

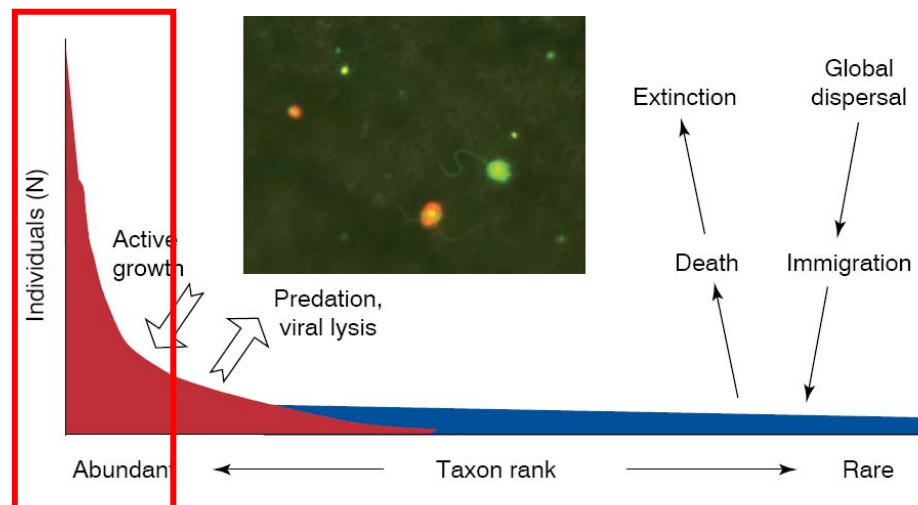


500,000 - 1 Mio  
1 week

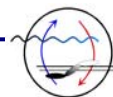


Statistics  
BLAST  
COGs  
Classify  
**Assembly?**

High diversity

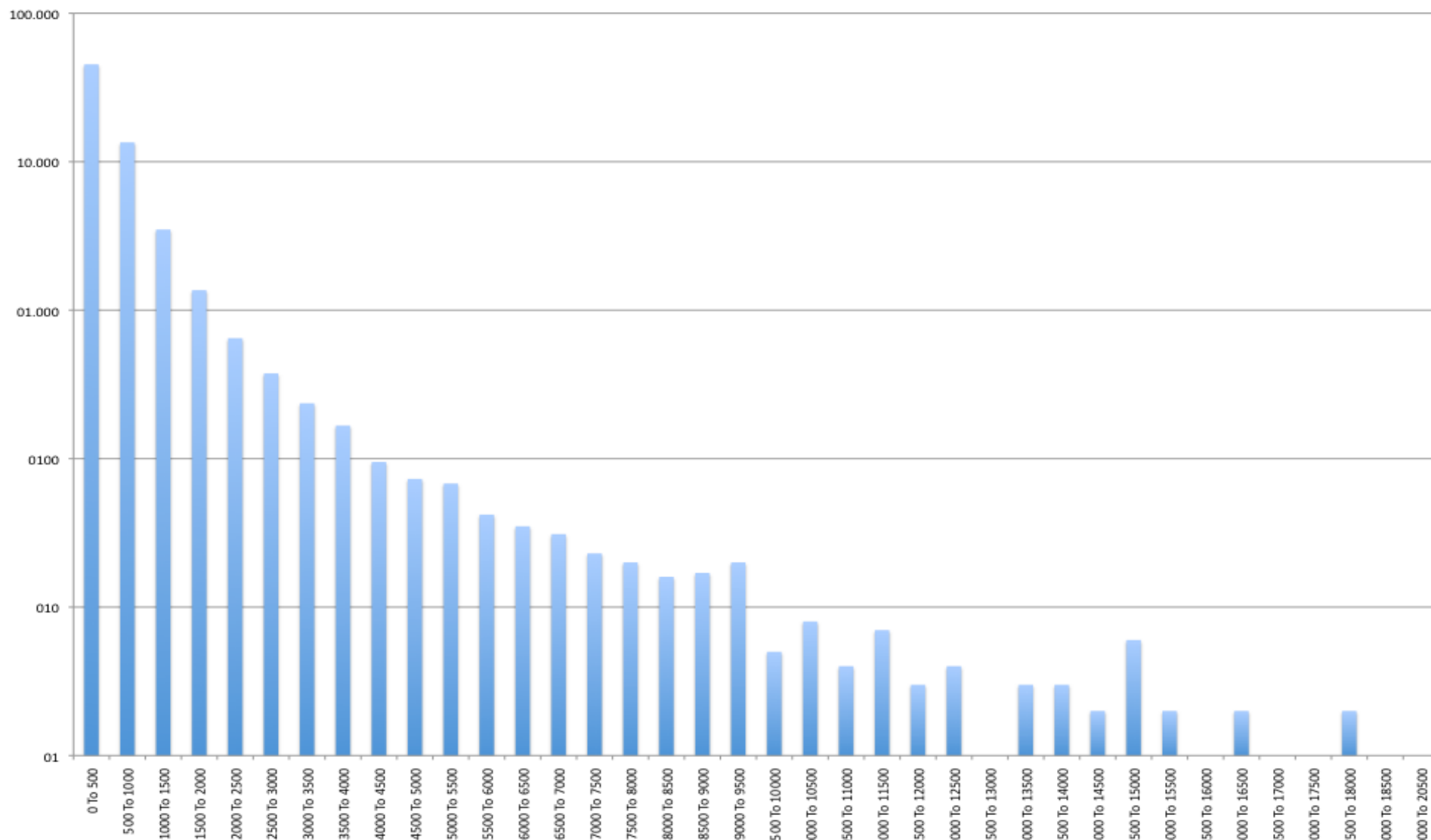


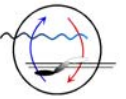




# Assembly

Histogramm of assembly lengths, 454 Flx Ti, 2 x 1/2 PTP, Helgoland Sample, bulk DNA, 04/14/2009





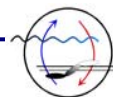
# Are we prepared for the Data Flood?

Your next-generation sequencing data just arrived...

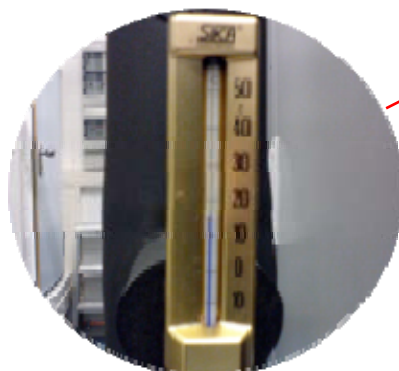
## Now What?



Richard was having a great day,  
until the arrival of his **next-generation data files**.



# Computing Infrastructure



## Cooling

- Liquid cooling with 5,000 L/h at 8° C



## Power

- 25 - 30 kW at peak
- installed: 40 kVA



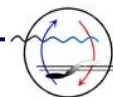
## Storage

- 8 TByte RAID file server
- 4 TByte RAID database server

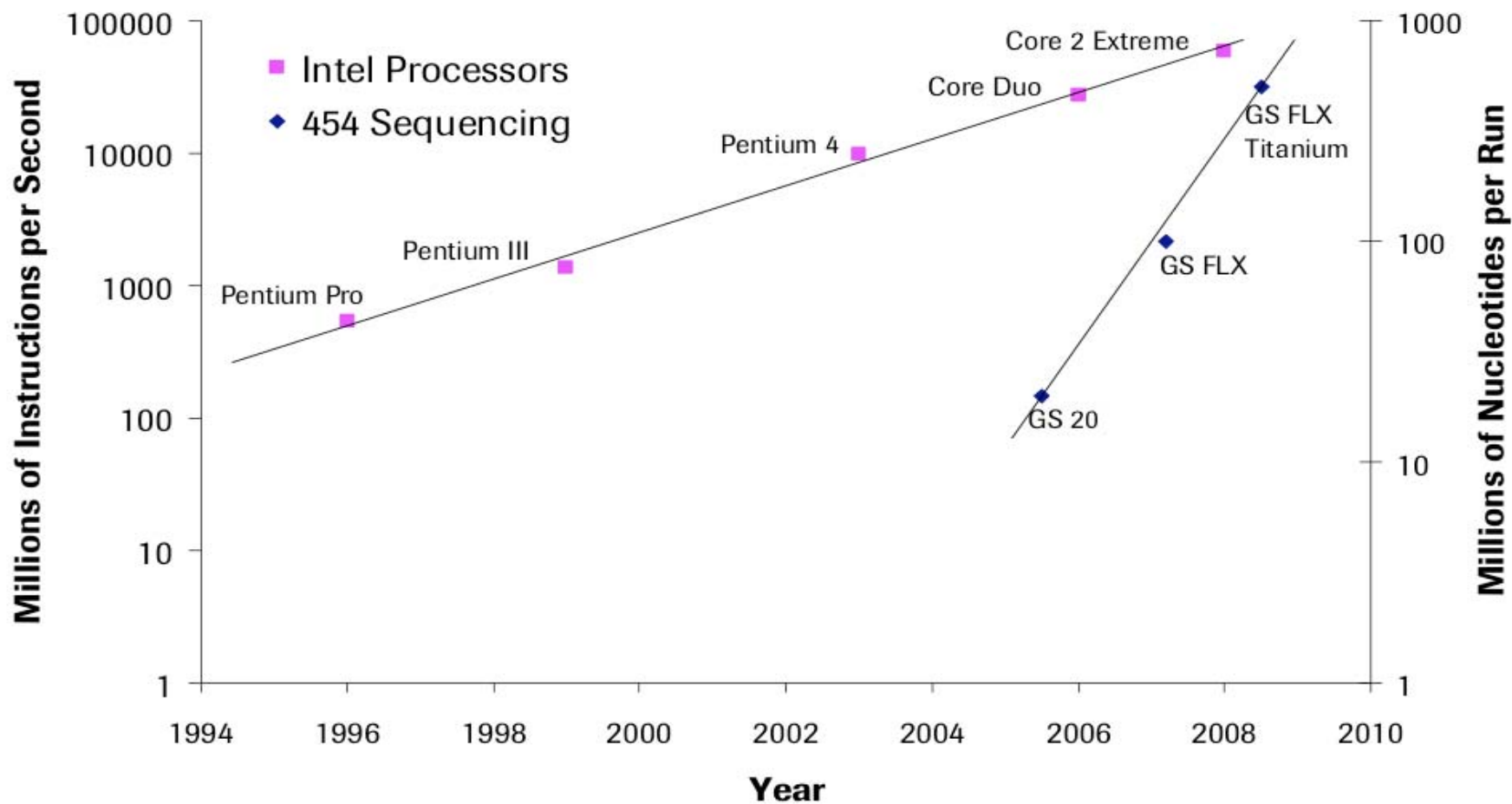


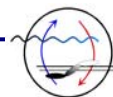
## Computers

- 43 cluster nodes
- several larger servers
- 400 CPU cores



# Moore's Law - Outcompeted

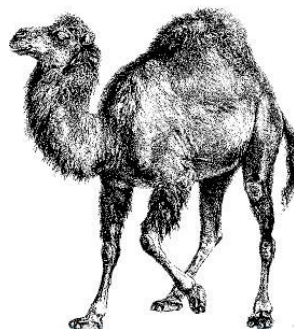




## Take Home Message for the Next Generation Biologists

### ▶ Three languages!

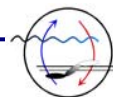
- Mother tongue
- English
- Perl, Python...



- Data management
  - ◆ Garbage in -> garbage out!
- Standardisation



[www.gensc.org](http://www.gensc.org)



# The Group

<http://www.microbial-genomics.de>



**Thanks for your attention**

