

# WHAT TO DO IN THE EVENT OF A DATA DELUGE

Adina Howe  
germslab.org (Genomics and  
Environmental Research in  
Microbial Systems)

Iowa State University, Ag &  
Biosystems Engr (January)

Slides available at  
[www.slideshare.com/  
adinachuanghowe](http://www.slideshare.com/adinachuanghowe)



**NGS SEQUENCING**

~~ZOMBIE~~ SURVIVAL TIP #3:

Panic fire only panics your allies and wastes ammunition.  
Stop. Breathe. Aim. Squeeze. Live.

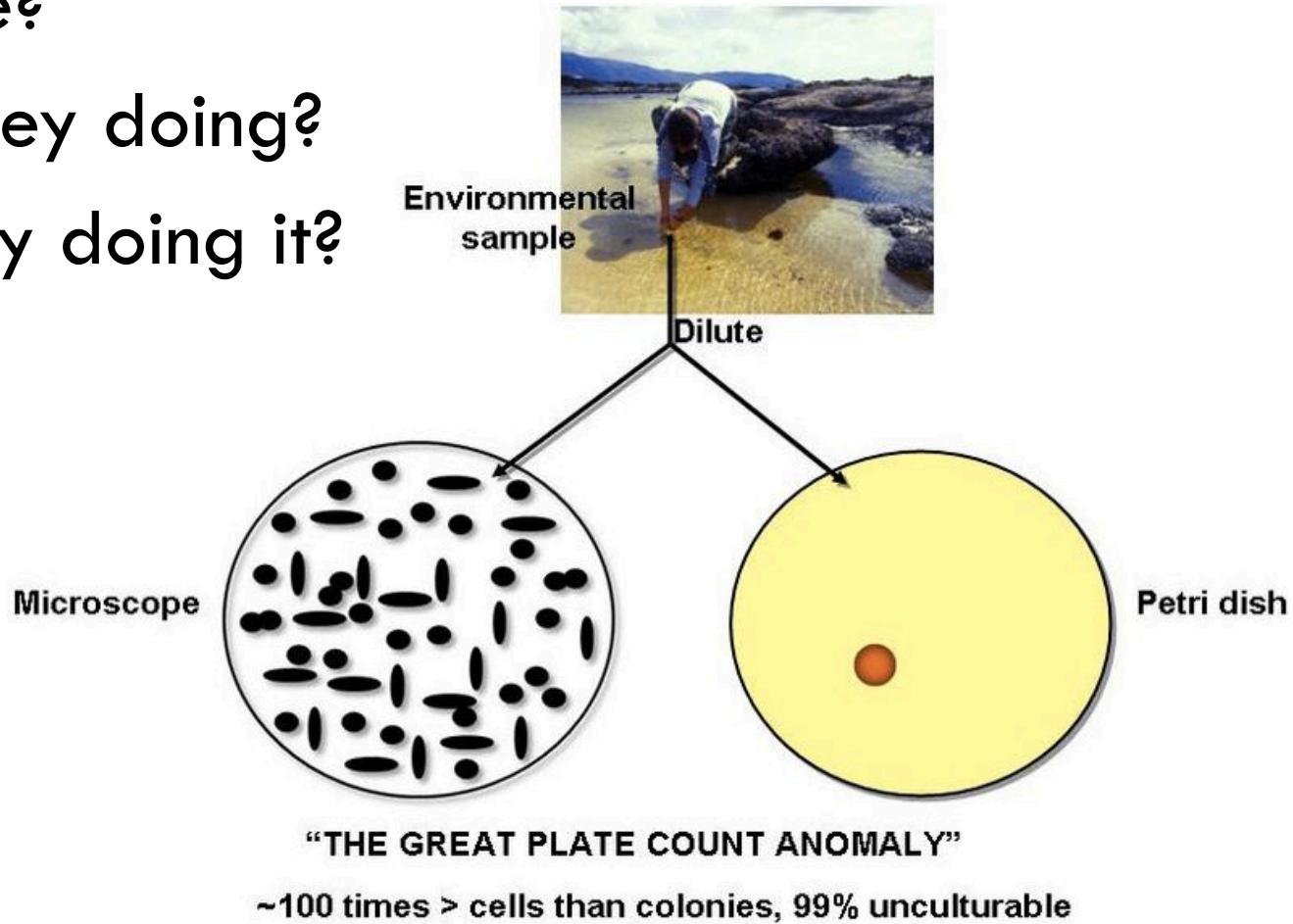
# THE PROMISE OF METAGENOMICS



HOW DID WE GET HERE

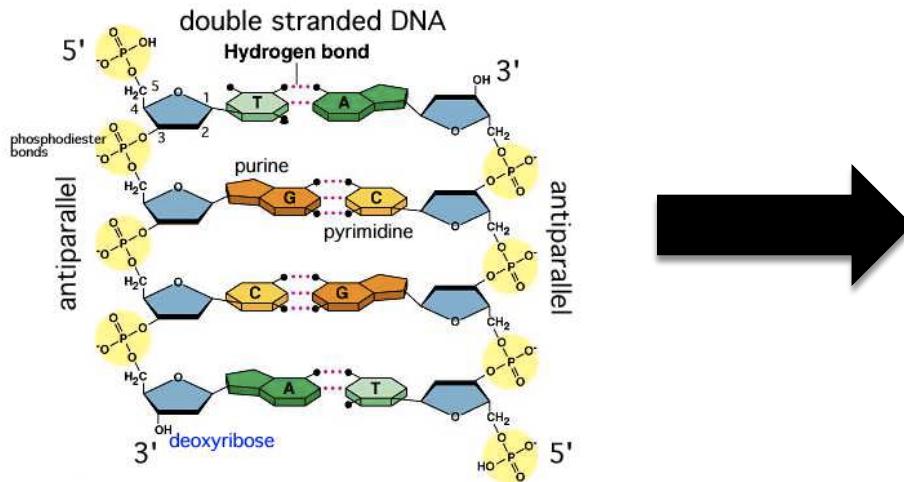
# Understanding community dynamics

- Who is there?
- What are they doing?
- How are they doing it?



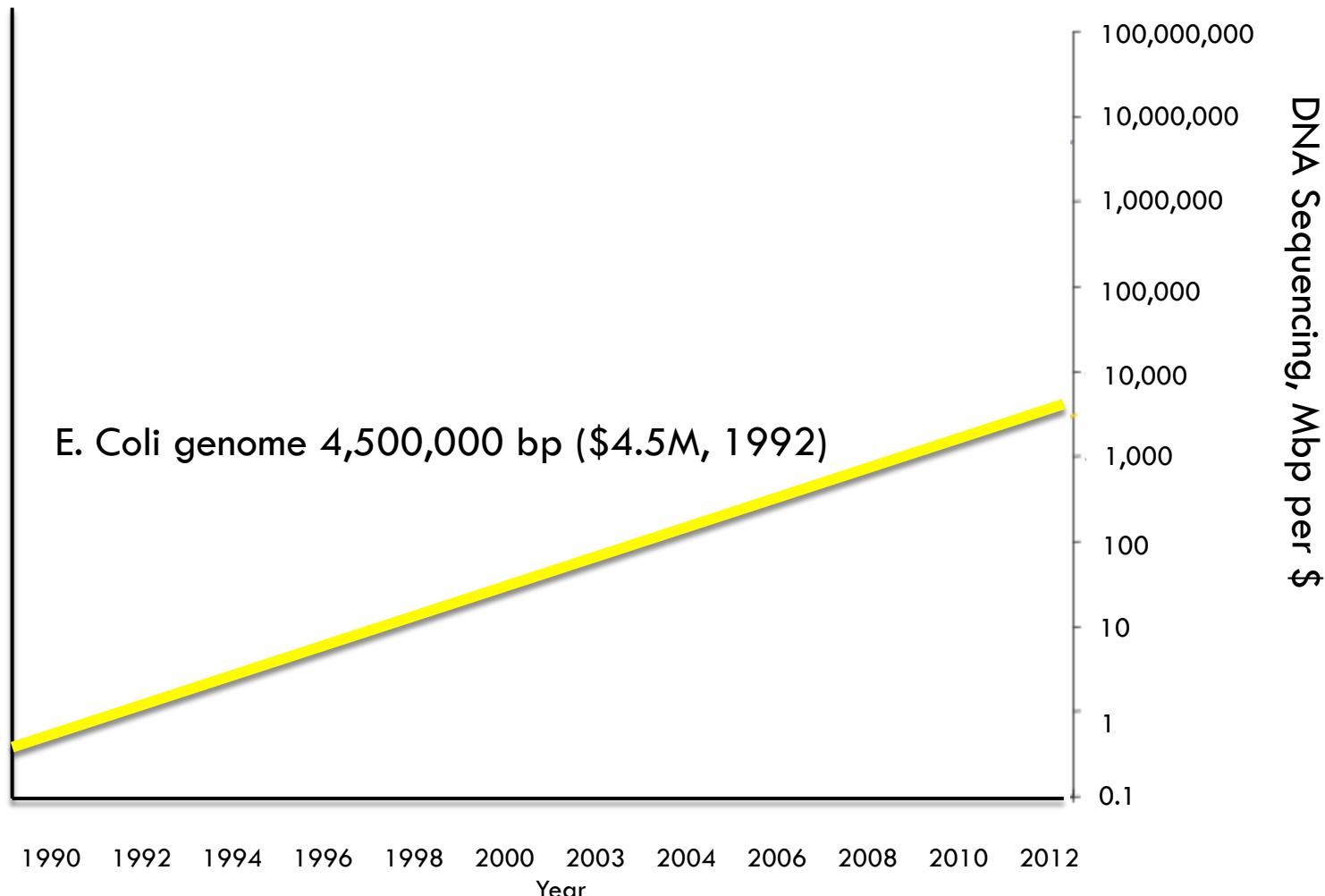
# Gene / Genome Sequencing

- Collect samples
- Extract DNA
- Sequence DNA
- “Analyze” DNA to identify its content and origin

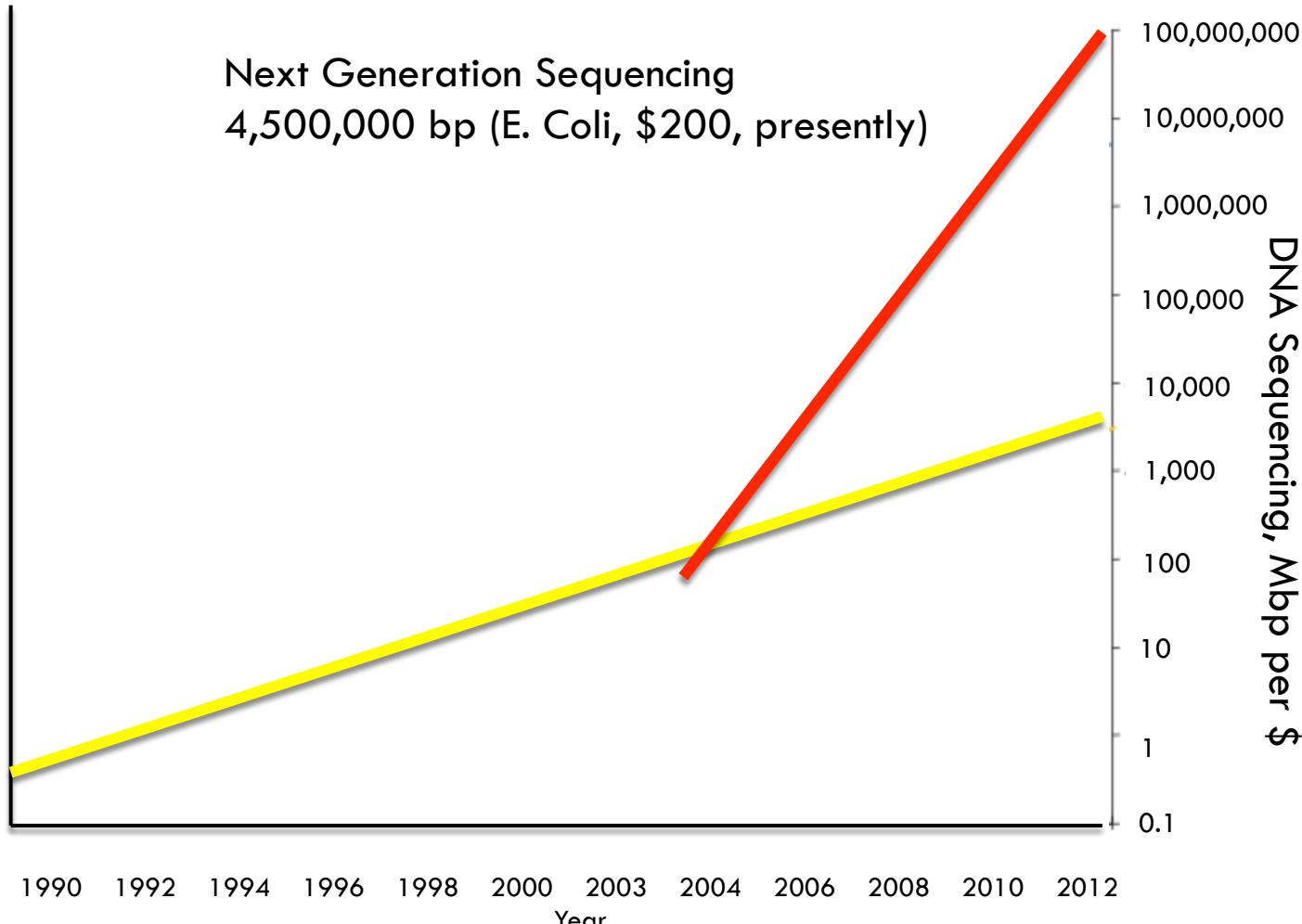


Taxonomy  
(e.g., pathogenic E. Coli)  
Function  
(e.g., degrades cellulose)

# Cost of Sequencing



# Rapidly decreasing costs with NGS Sequencing



# Effects of low cost sequencing...

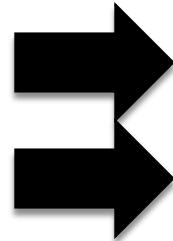
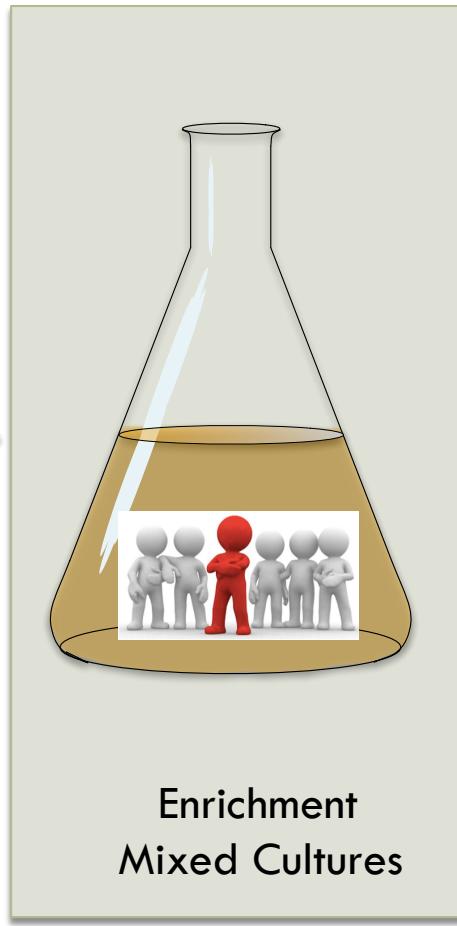
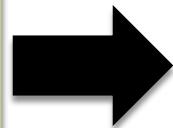
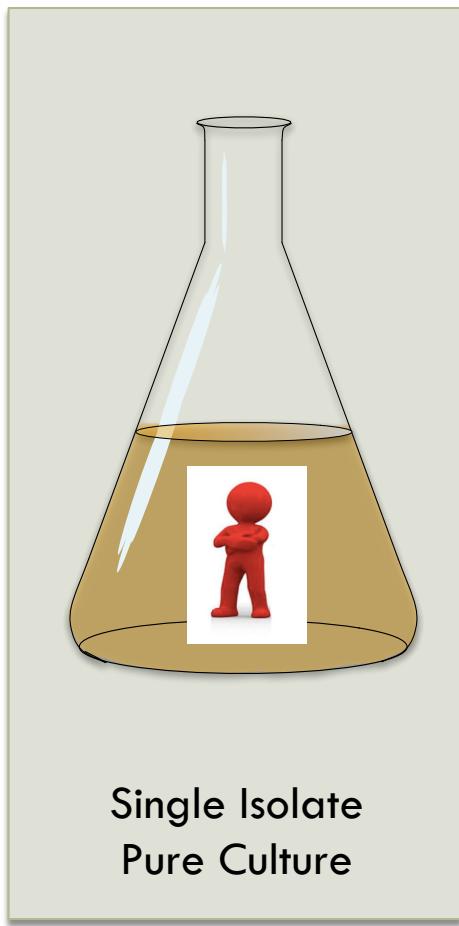


First free-living bacterium sequenced for  
billions of dollars and years of analysis

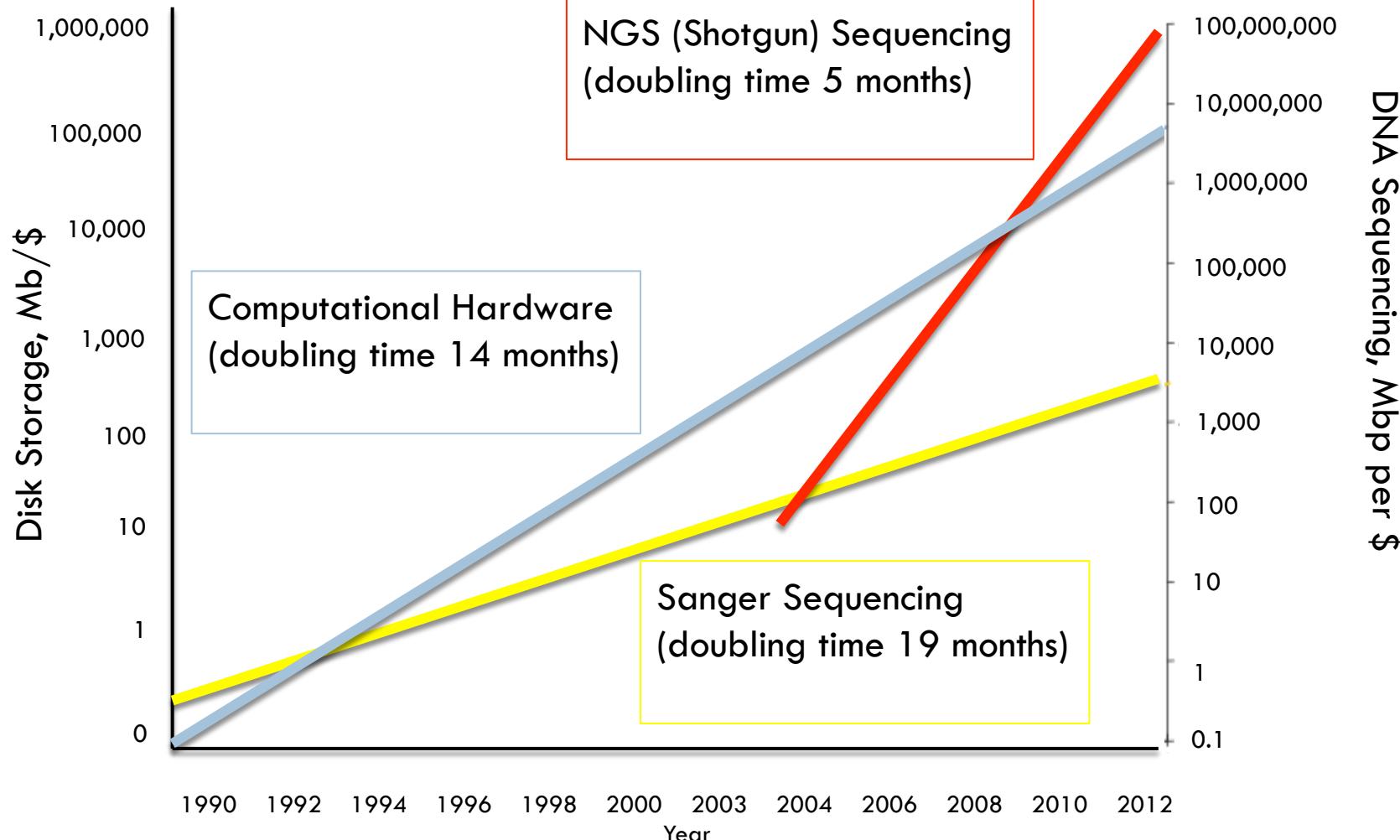
A screenshot of the 23andMe website homepage. The header includes the 23andMe logo, a search bar, and navigation links for welcome, ancestry, health, how it works, store, and help. A banner at the top states "23andMe can help you manage risk and make informed decisions...". The main content area features three main sections: "Ancestry" (with a tree icon), "Health" (with a person running icon), and "Research" (with a lightbulb icon). Below these are sub-sections: "Learn for the present." under Health and "Participate for the future." under Research. A large central callout reads "Learn valuable health &amp; ancestry information." Below it is a price of "\$99" and a pink "Order Now" button. At the bottom, there's an image of the 23andMe DNA spit kit in its green packaging.

Personal genome can be mapped in a few days and hundreds to few thousand dollars

# The experimental continuum



# The era of big data in biology



# Postdoc experience with data

2003-2008 Cumulative sequencing in PhD = 2000 bp

2008-2009 Postdoc Year 1 = 50 Gbp

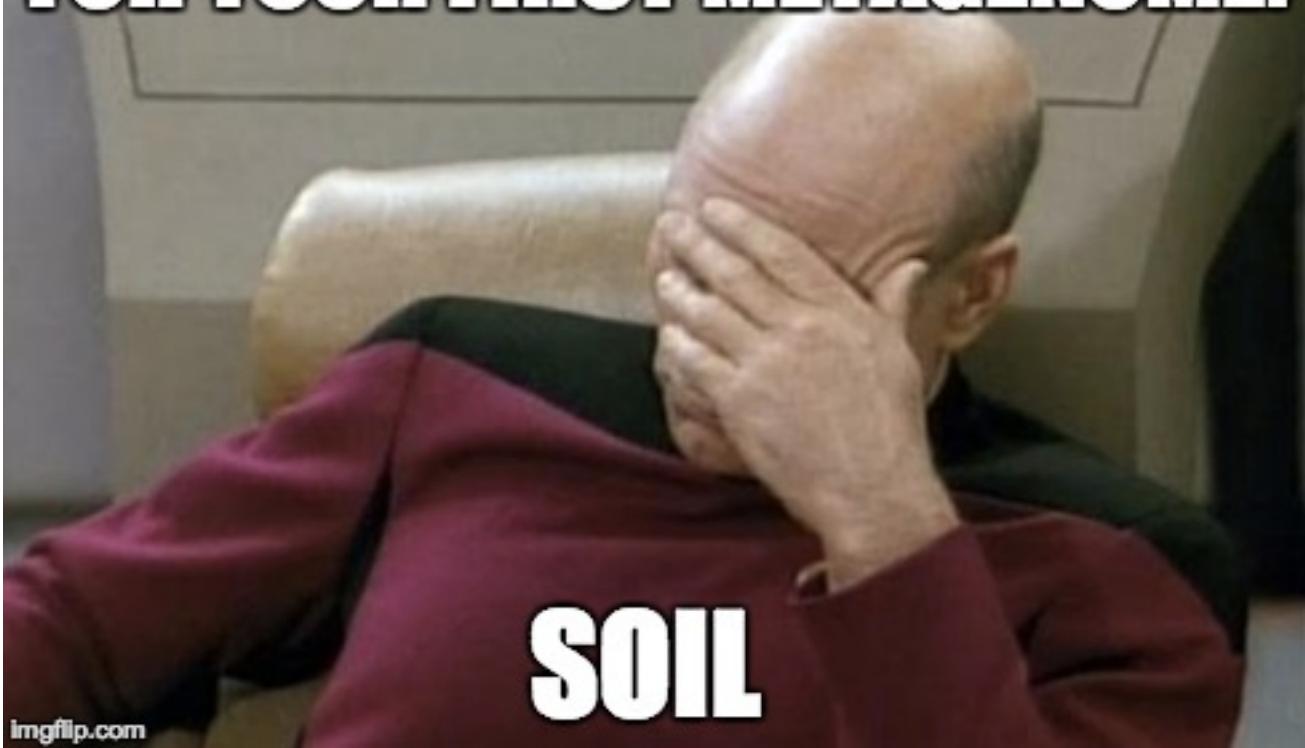
2009-2010 Postdoc Year 2 = 450 Gbp

2014 = 50 Tbp

2015 = 500 Tbp budgeted



**FOR YOUR FIRST METAGENOME:**

A photograph of Captain Jean-Luc Picard from Star Trek: The Next Generation. He is sitting in a light-colored armchair, wearing his signature maroon Starfleet uniform over a green shirt. His right hand is raised to his forehead, with his fingers resting against his temple, a gesture often used to indicate distress or deep thought. The background is a plain, light-colored wall.

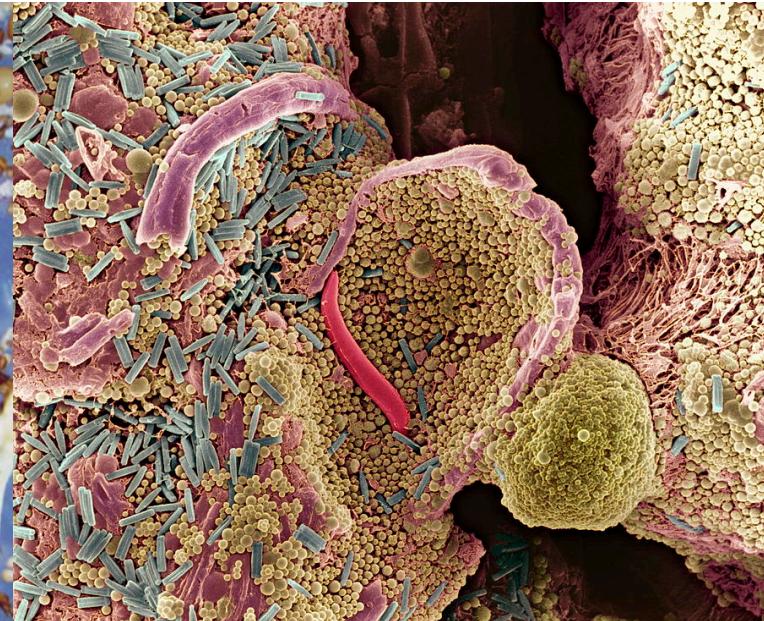
**SOIL**

# THE DIRT ON SOIL

## MAGNIFICENT BIODIVERSITY



Biodiversity in the dark, Wall et al., Nature Geoscience, 2010



Jeremy Burgess

# THE DIRT ON SOIL

## SPATIAL HETEROGENEITY

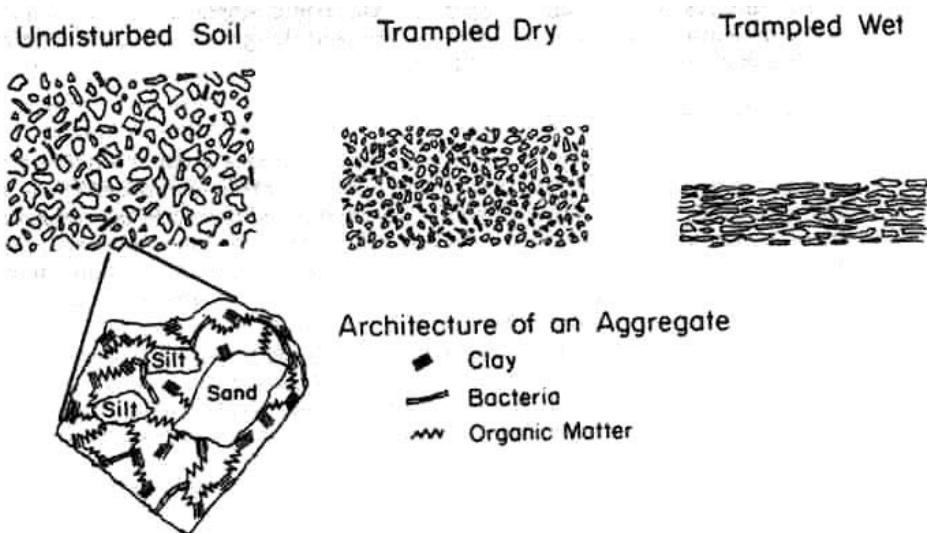
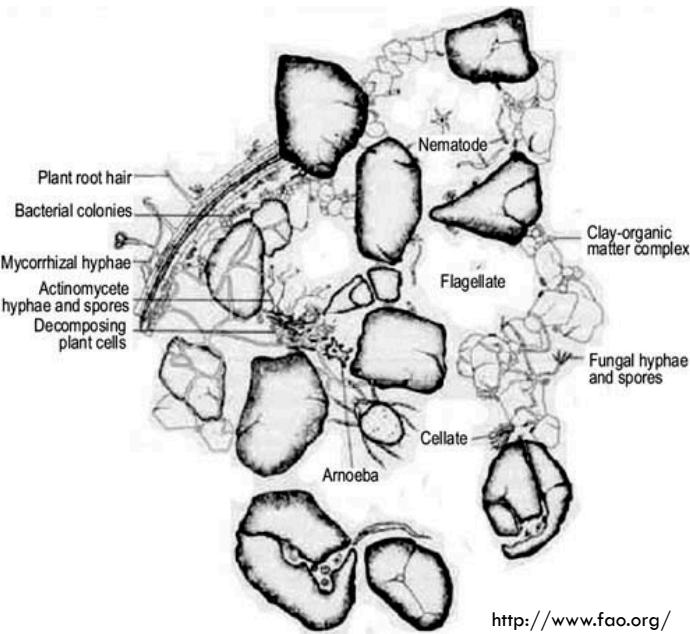


Figure 6.2. Conceptual architecture of a soil aggregate and the changes in soil aggregate structure caused by trampling under wet and dry conditions.

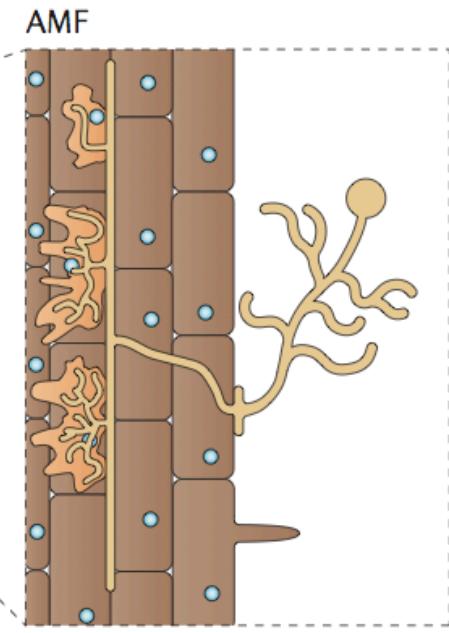
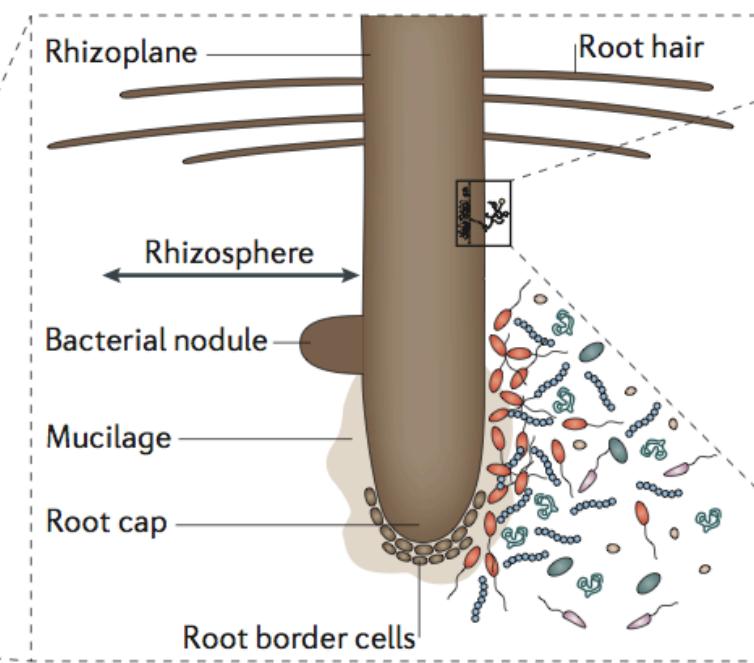
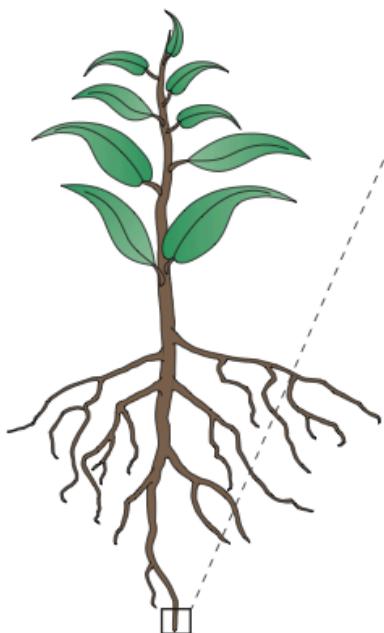
# THE DIRT ON SOIL

DYNAMIC



# THE DIRT ON SOIL

INTERACTIONS: BIOTIC, ABIOTIC, ABOVE, BELOW, SCALES



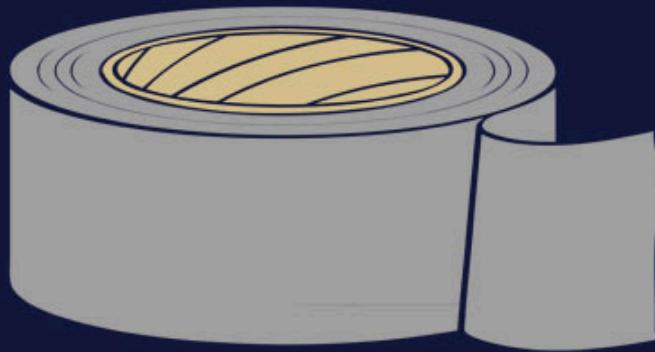
Philippot, 2013, *Nature Reviews Microbiology*

# Tackling Soil Biodiversity



C. Titus Brown, James Tiedje, Qingpeng Zhang, Jason Pell (MSU)  
Janet Jansson, Susannah Tringe (JGI)

I'LL GET  
MY TOOL KIT



A Slight Digression: Decisions for the new microbial ecologist

# Getting the most out of your data

*ID, Abundance, Function*

Complex  
Samples

16S rRNA amplicon  
sequencing

Pros:

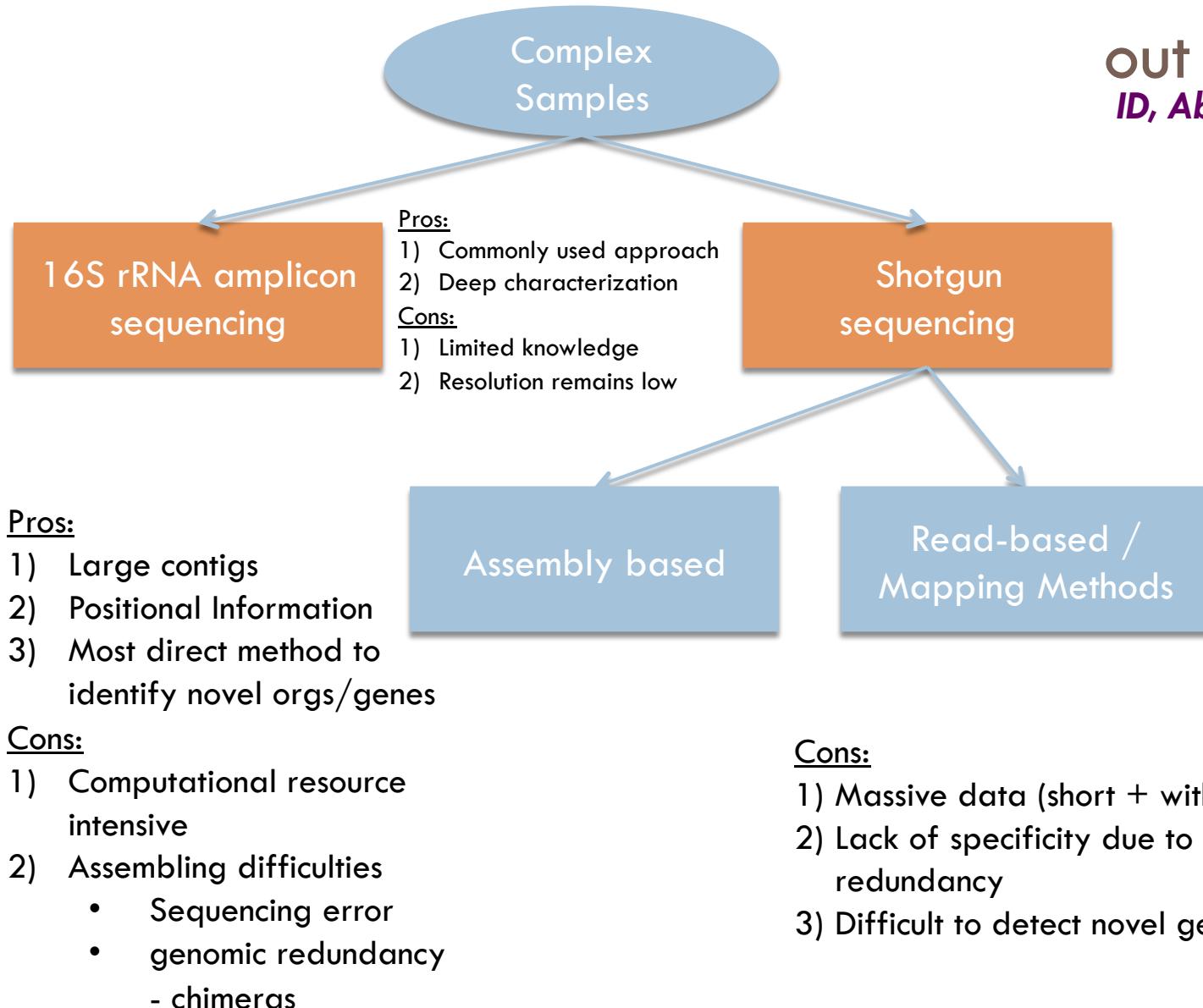
- 1) Commonly used approach
- 2) Deep characterization

Cons:

- 1) Limited knowledge
- 2) Resolution remains low

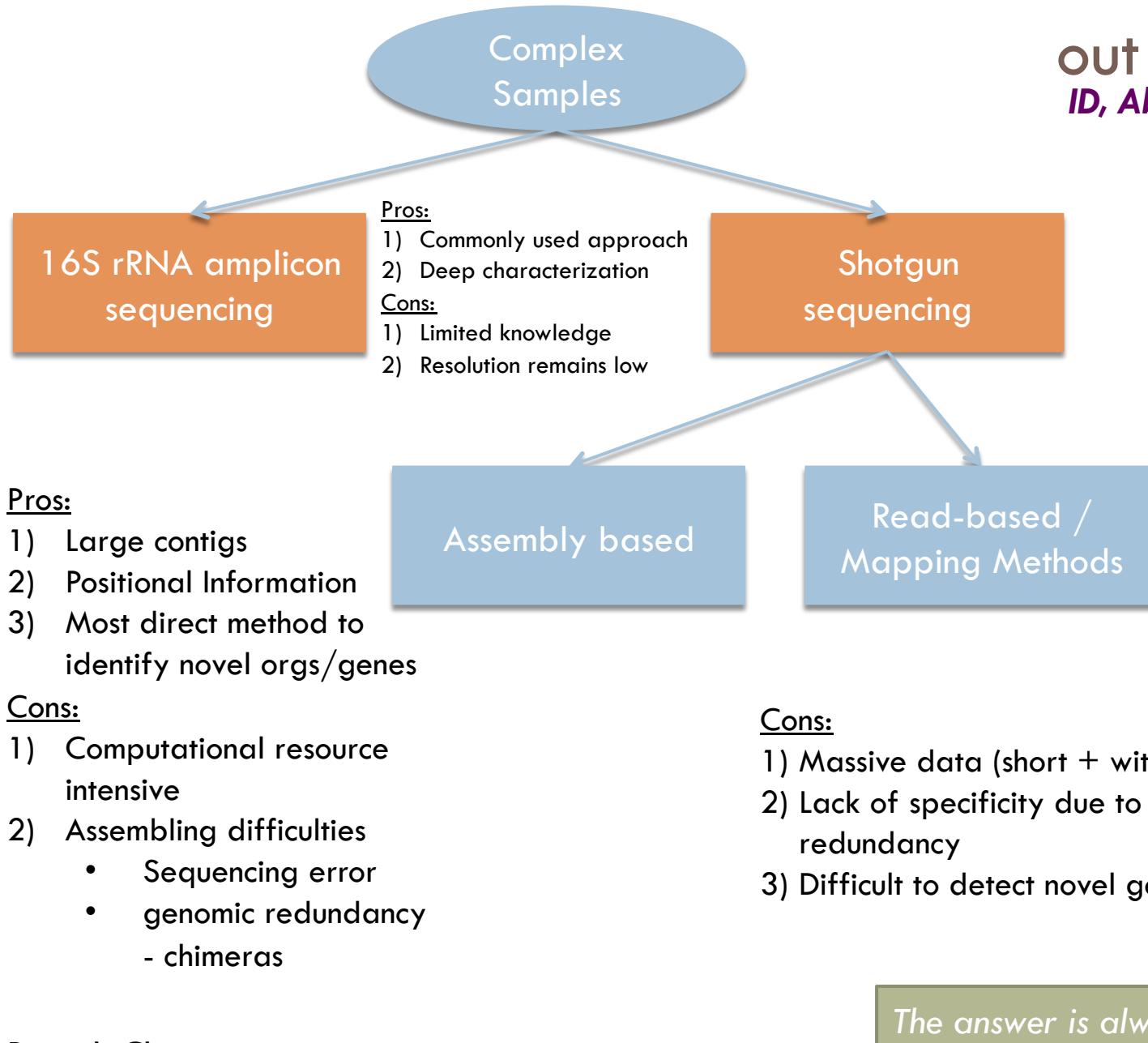
# Getting the most out of your data

*ID, Abundance, Function*

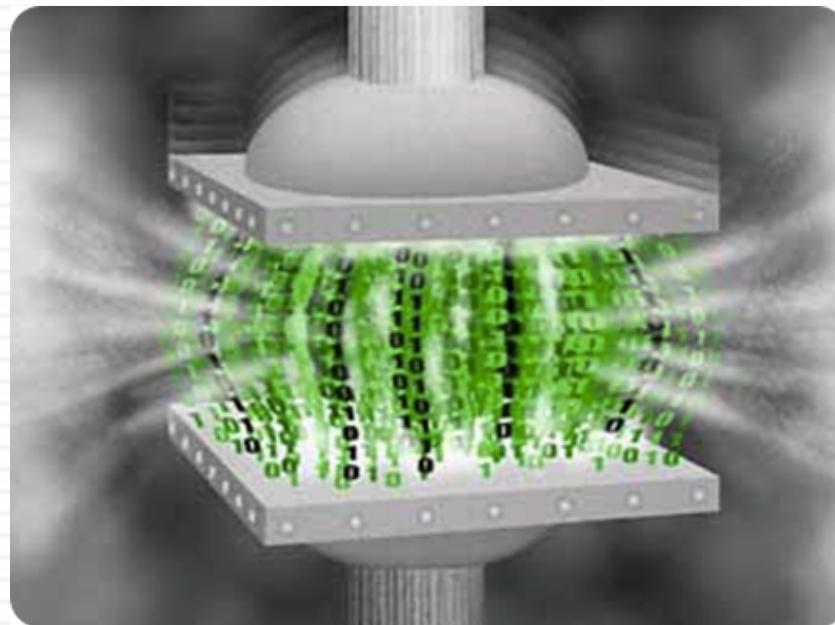


# Getting the most out of your data

*ID, Abundance, Function*



# Example #1: Data compression

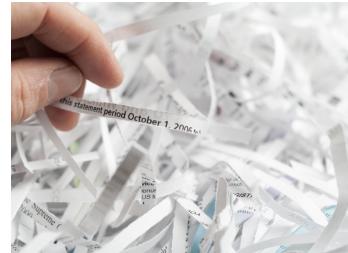


[http://siliconangle.com/files/2010/09/image\\_thumb69.png](http://siliconangle.com/files/2010/09/image_thumb69.png)

# *de novo* assembly



Raw sequencing data (“reads”)



Computational algorithms



Informative genes / genomes

- ❖ Compresses dataset size significantly
- ❖ Improved data quality (longer sequences, gene order)
- ❖ Reference not necessary (novelty)

# Metagenome assembly...a scaling problem.



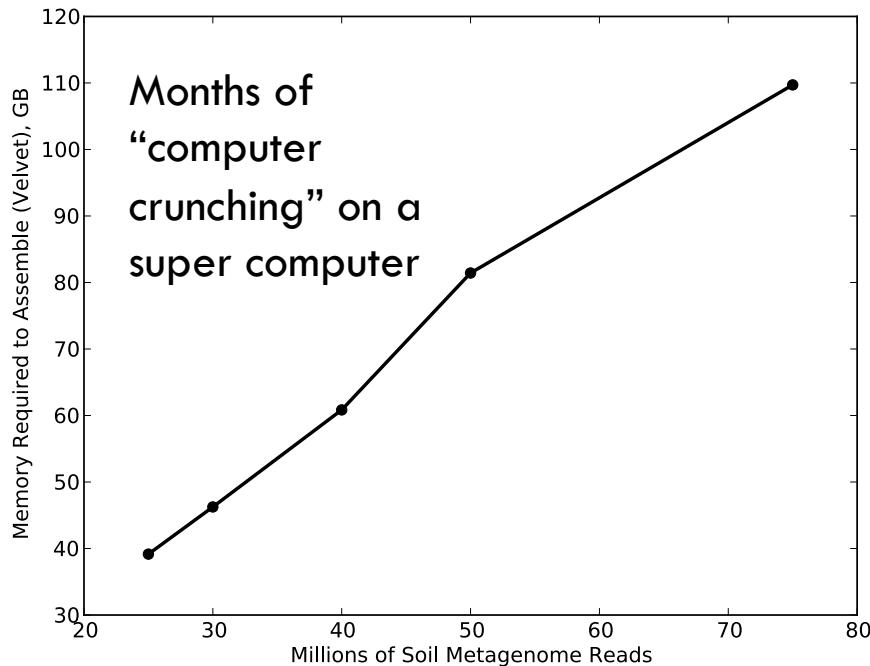
# Shotgun sequencing and de novo assembly

It was the best of times, it was the worst of times, it was the worst of timZs, it was the isdom, it was the age of foolisXness, it was the worVt of times, it was the mes, it was Ahe age of wisdom, it was th  
It was the best of times, it Gas the wor mes, it was the age of witsdom, it was th isdom, it was tle age of foolishness

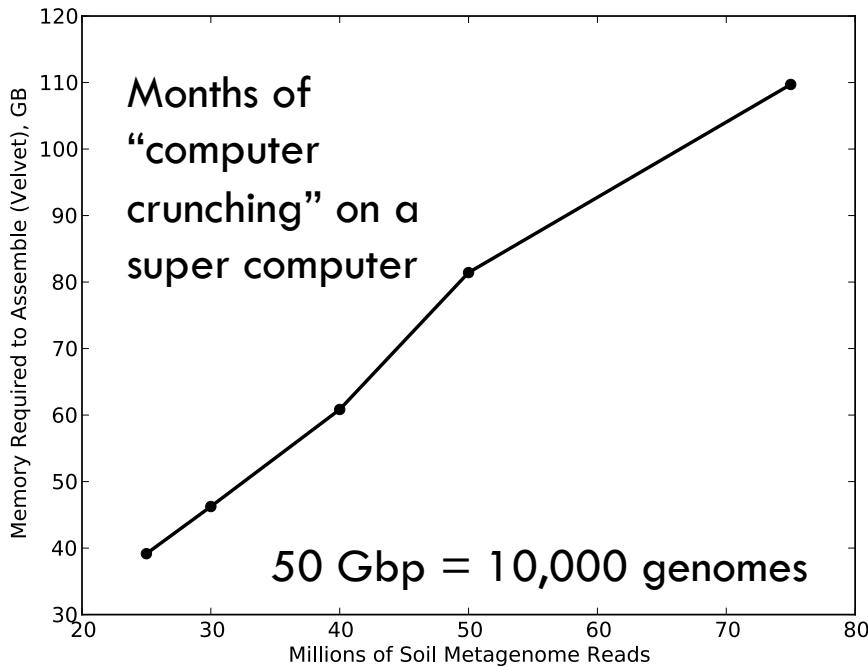


It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness

# Practical Challenges – Intensive computing



# Practical Challenges – Intensive computing



Assembly of 300 Gbp (70,000 genomes worth) can be done with any assembly program in less than 14 GB RAM and less than 24 hours.



# Four main challenges for *de novo* sequencing.

- Repeats.
- Low coverage.
- Errors

These introduce breaks in the construction of contigs.

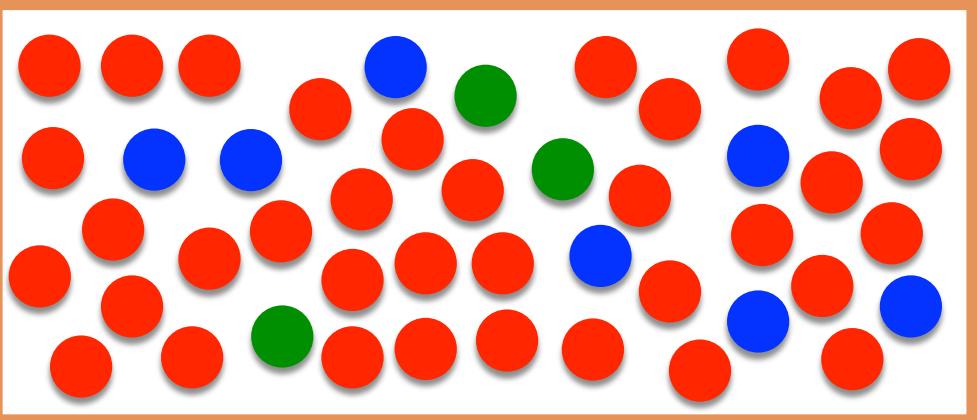
- Variation in coverage – transcriptomes and metagenomes, as well as amplified genomic.

This challenges the assembler to distinguish between erroneous connections (e.g. repeats) and real connections.

# Natural community characteristics

## ◆ Diverse

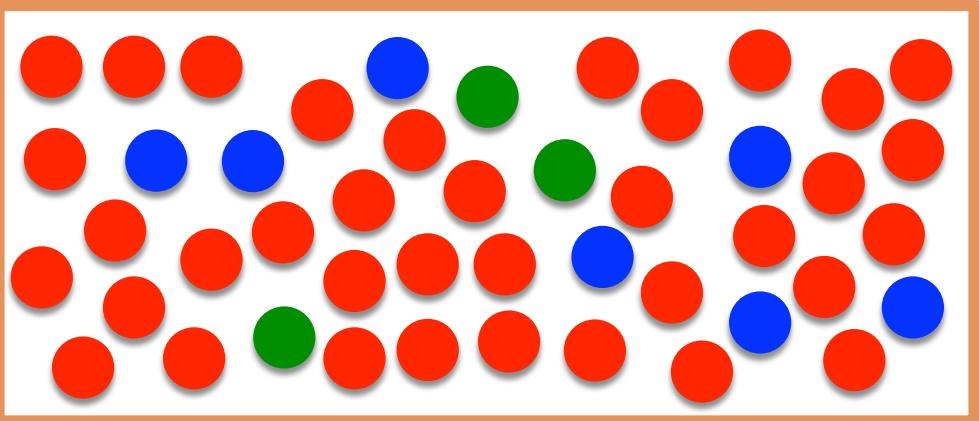
→ Many organisms  
(genomes)



# Natural community characteristics

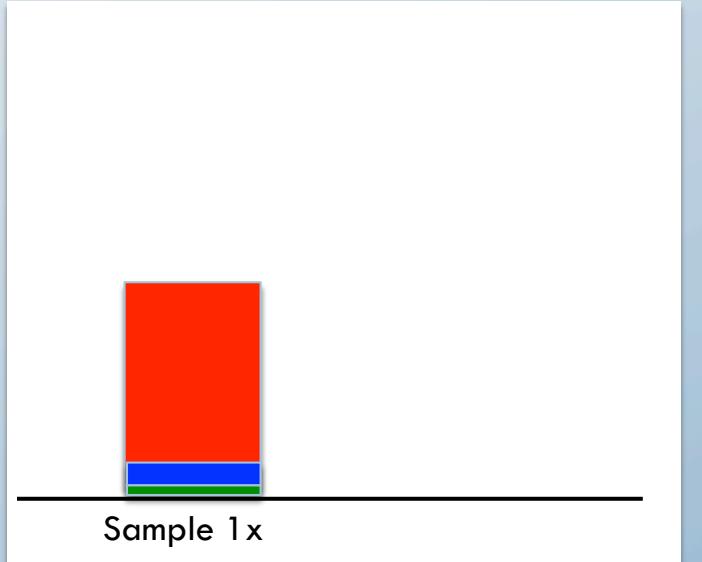
## ◆ Diverse

- Many organisms (genomes)



## ◆ Variable abundance

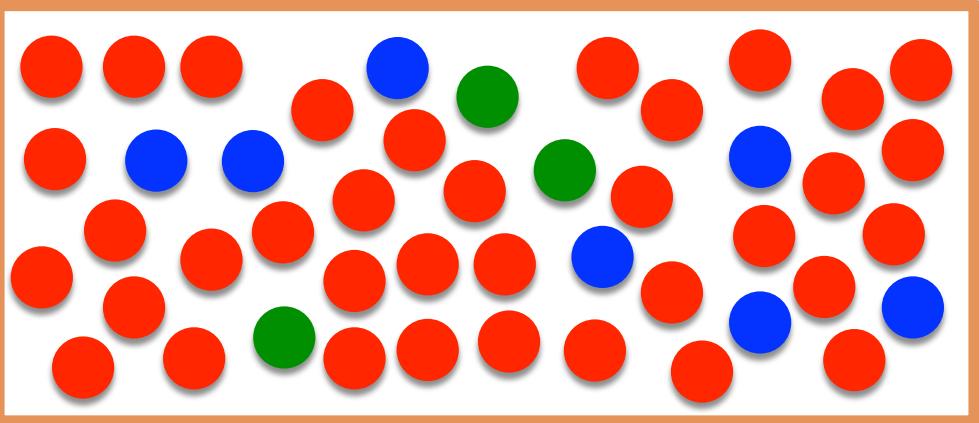
- Most abundant organisms, sampled more often
- Assembly requires a minimum amount of sampling
- More sequencing, more errors



# Natural community characteristics

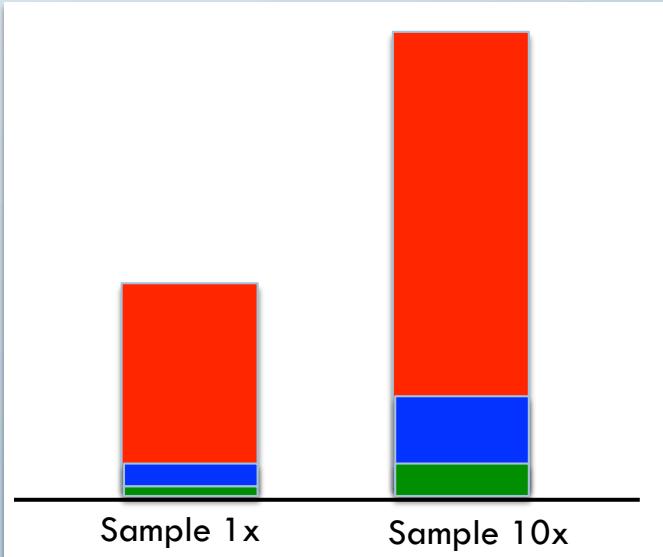
## ◆ Diverse

- Many organisms (genomes)



## ◆ Variable abundance

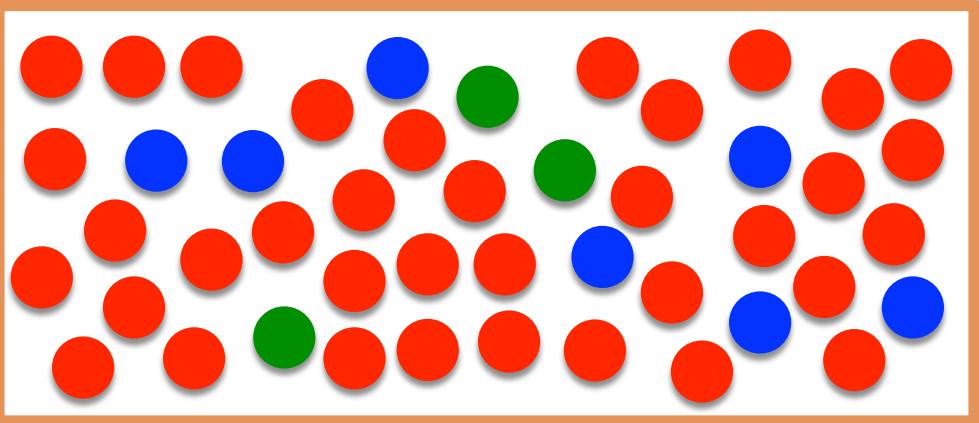
- Most abundant organisms, sampled more often
- Assembly requires a minimum amount of sampling
- More sequencing, more errors



# Natural community characteristics

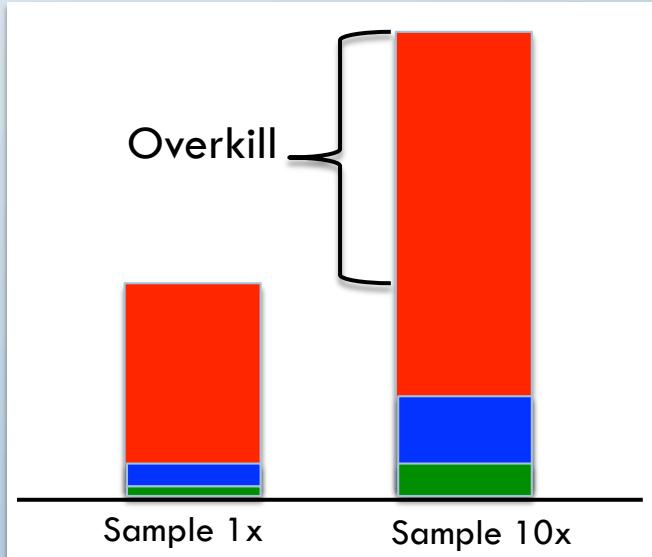
## ◆ Diverse

- Many organisms (genomes)



## ◆ Variable abundance

- Most abundant organisms, sampled more often
- Assembly requires a minimum amount of sampling
- More sequencing, more errors



# Digital normalization

Page 1 of 1

Reads  
(randomly sequenced)

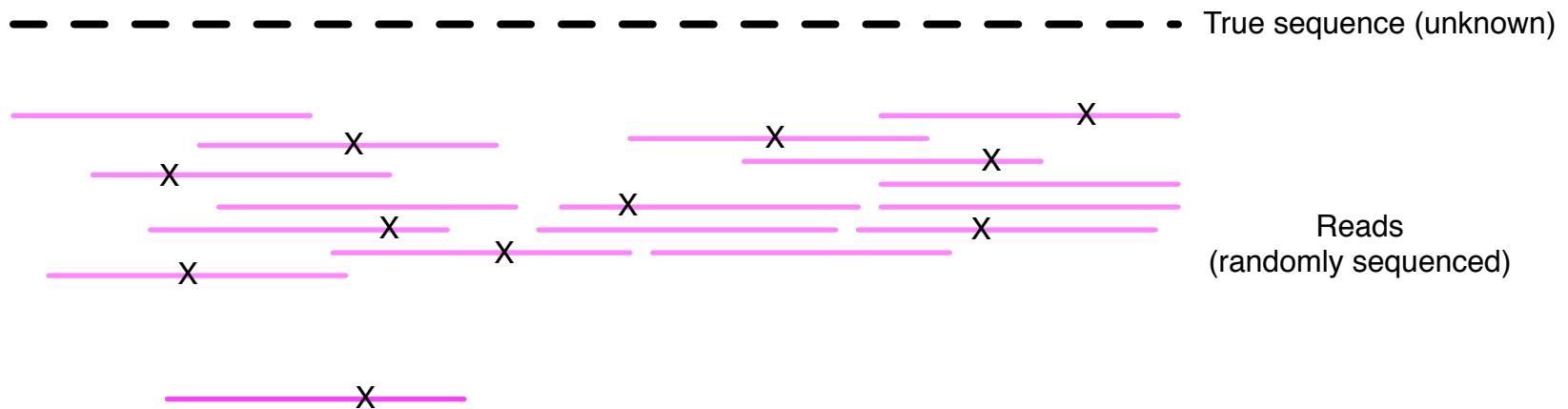
Brown et al., 2012, arXiv  
Howe et al., 2014, PNAS  
Zhang et al., 2014, PLOS One

# Digital normalization



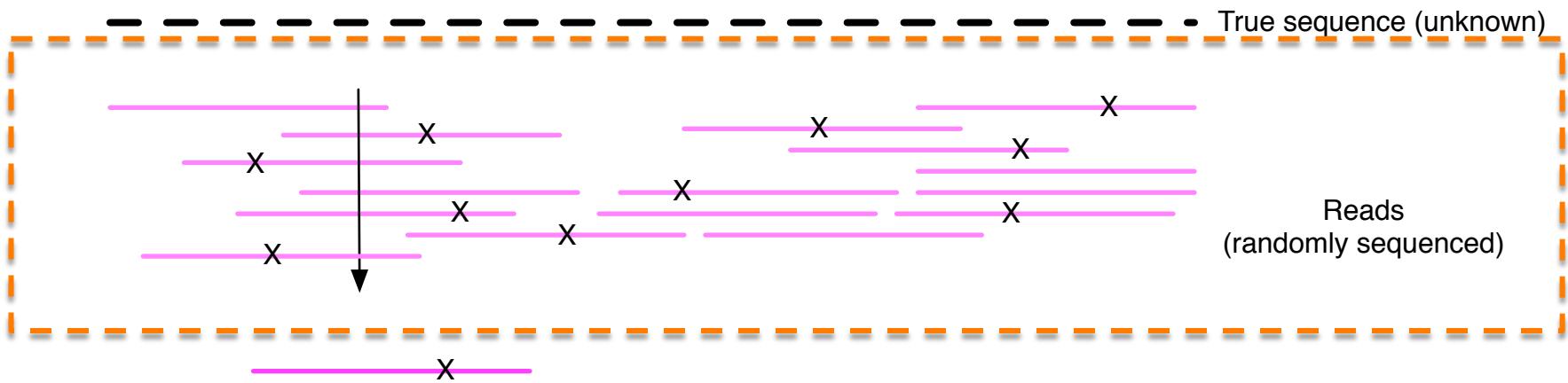
Brown et al., 2012, arXiv  
Howe et al., 2014, PNAS  
Zhang et al., 2014, PLOS One

# Digital normalization

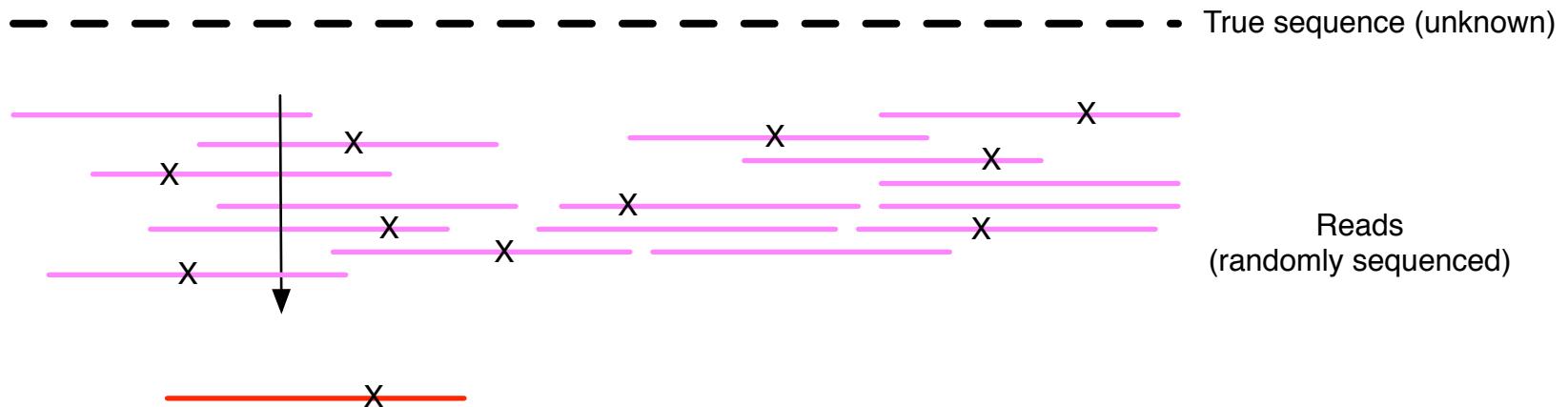


Brown et al., 2012, arXiv  
Howe et al., 2014, PNAS  
Zhang et al., 2014, PLOS One

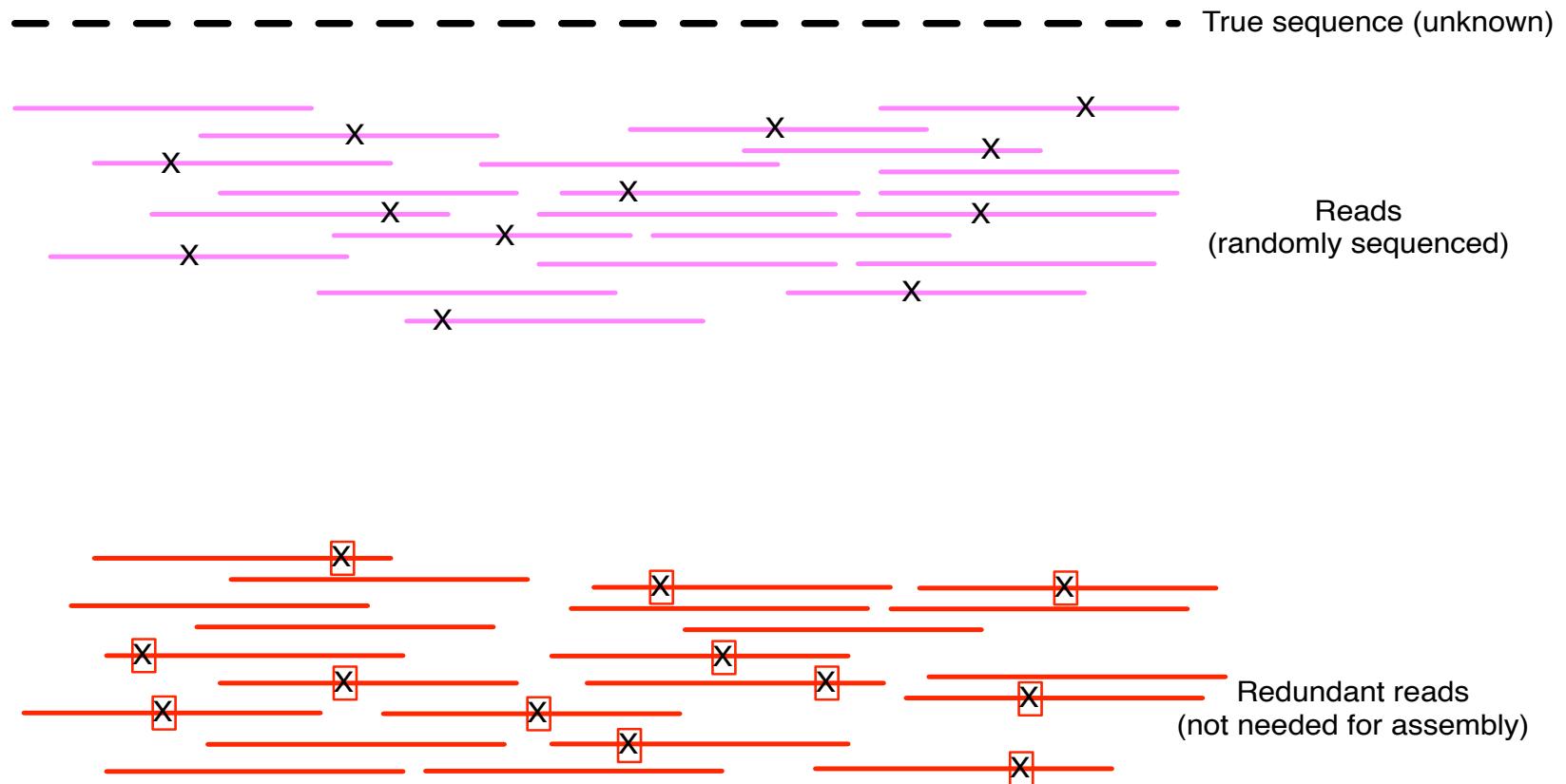
# Digital normalization



# Digital normalization



# Digital normalization



- ❖ Scales datasets for assembly up to 95% - same assembly outputs.
- ❖ Genomes, mRNA-seq, metagenomes (soils, gut, water)

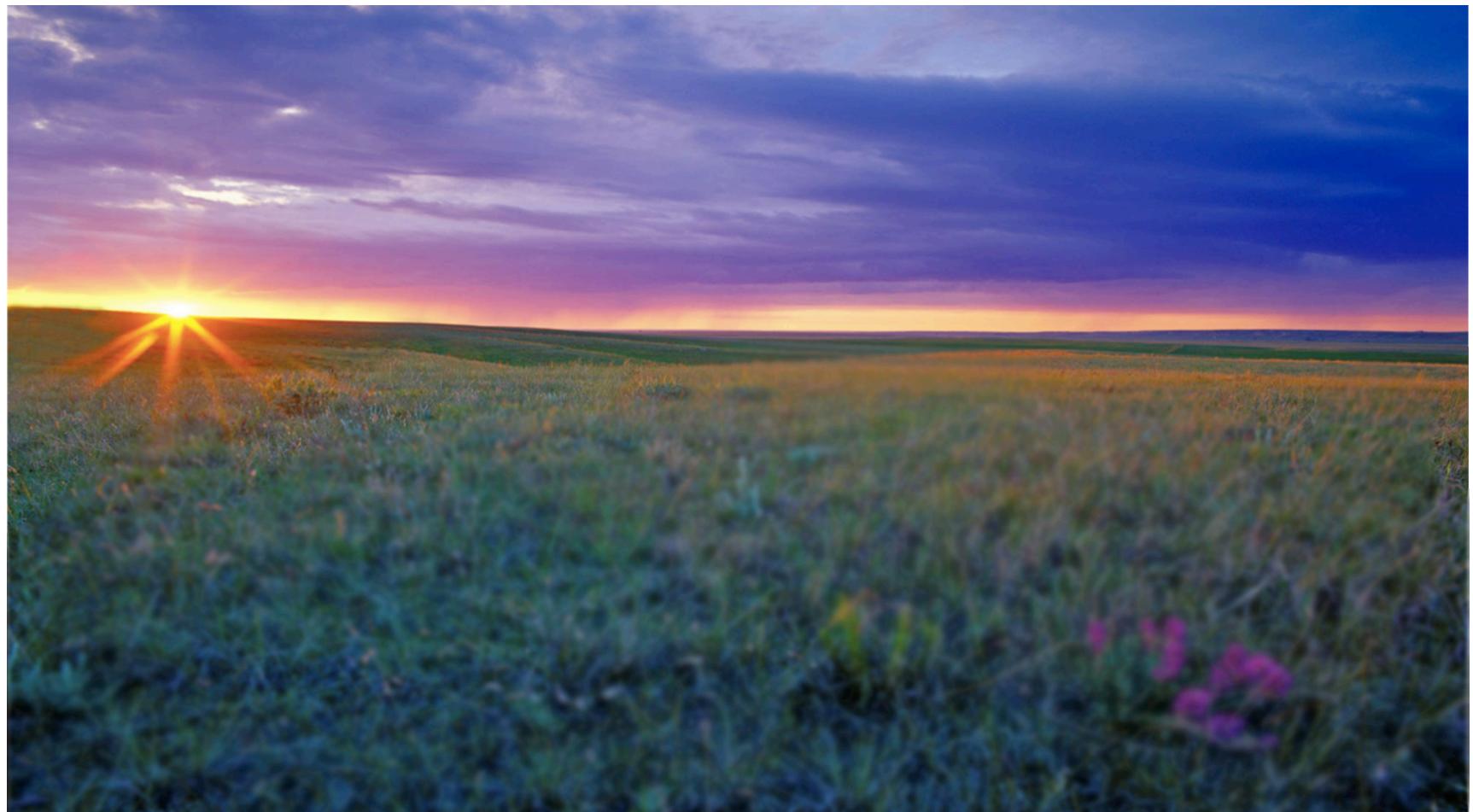
Brown et al., 2012, arXiv  
Howe et al., 2014, PNAS  
Zhang et al., 2014, PLOS One

# Tackling Soil Biodiversity



C. Titus Brown, James Tiedje, Qingpeng Zhang, Jason Pell (MSU)  
Janet Jansson, Susannah Tringe (JGI)

# The reality?



# More like...



# The Future

- More data, more samples, better references
- Expense will be in sampling – not sequencing or even data analysis
- All biologists will need to know how to use a pipette and write computer programs
- Large-scale, collaborative projects rather than single PI efforts