

Analysis and assembly methods for microbiome sequencing data

Marcus Fedarko



UCSDCSE
Computer Science and Engineering

UC San Diego

Analysis and assembly methods for **microbiome** sequencing data

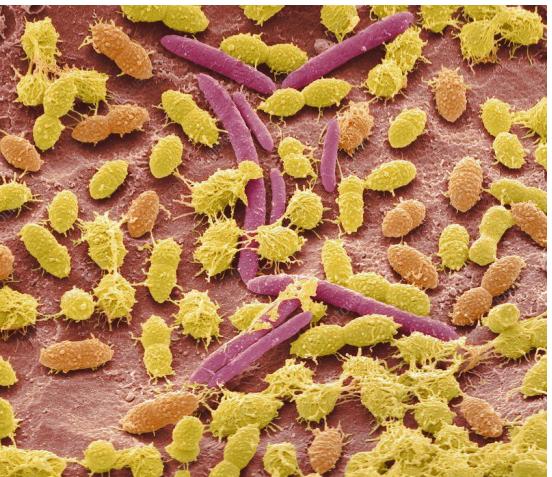
Marcus Fedarko



UCSDCSE
Computer Science and Engineering

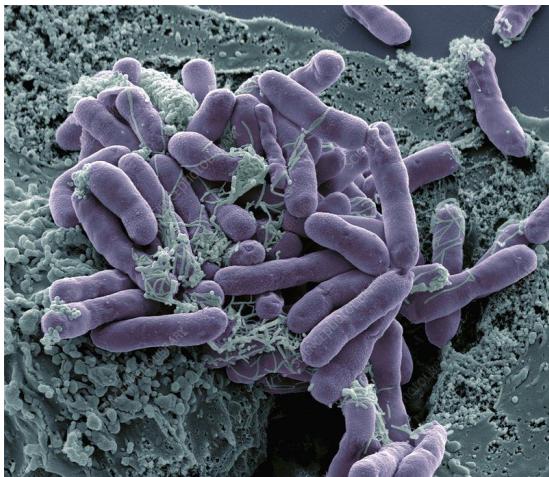
UC San Diego

Microbiomes



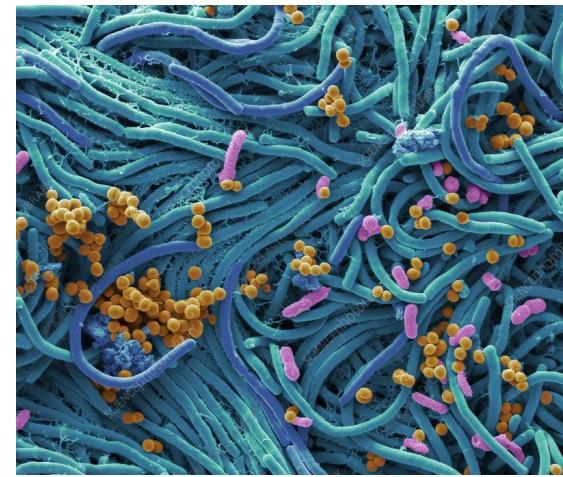
[Steve Gschmeissner, Science Photo Library](#)

“Scanning electron micrograph (SEM) of bacteria cultured from a sample of human faeces.”



[Steve Gschmeissner, Science Photo Library](#)

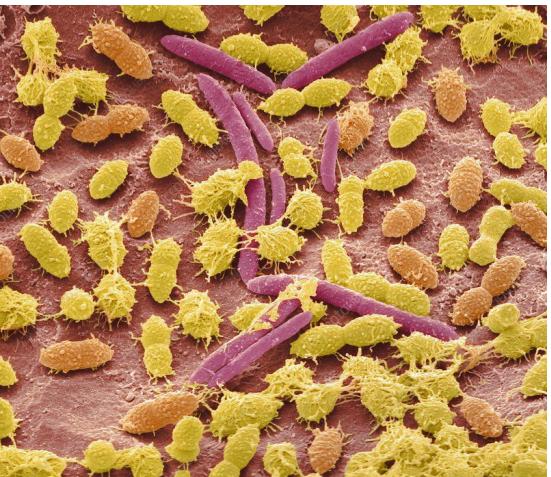
“Bacterial contamination, coloured scanning electron micrograph (SEM). *Escherichia coli* bacteria in a cell culture. This contamination has come from an unclean water source.”



[Steve Gschmeissner, Science Photo Library](#)

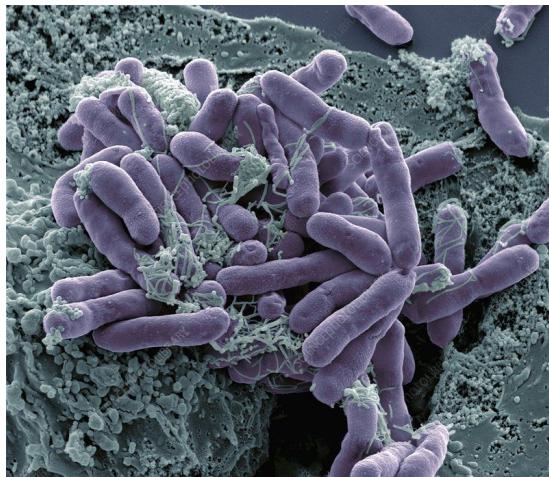
“Coloured scanning electron micrograph (SEM) of bacteria cultured from a mobile phone.”

Microbiomes



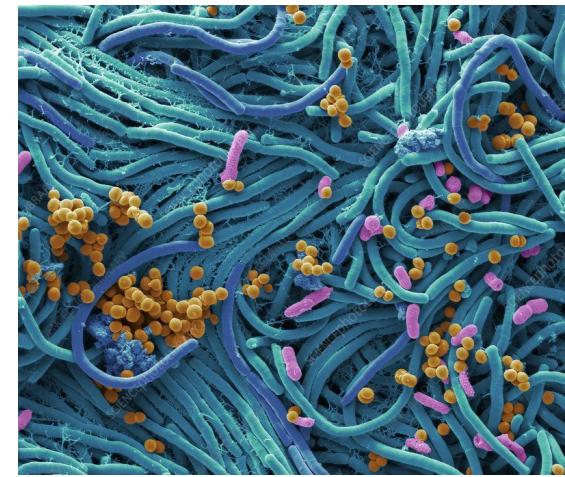
[Steve Gschmeissner, Science Photo Library](#)

“Scanning electron micrograph (SEM) of bacteria cultured from a sample of human faeces.”



[Steve Gschmeissner, Science Photo Library](#)

“Bacterial contamination, coloured scanning electron micrograph (SEM). *Escherichia coli* bacteria in a cell culture. This contamination has come from an unclean water source.”



[Steve Gschmeissner, Science Photo Library](#)

“Coloured scanning electron micrograph (SEM) of bacteria cultured from a mobile phone.”



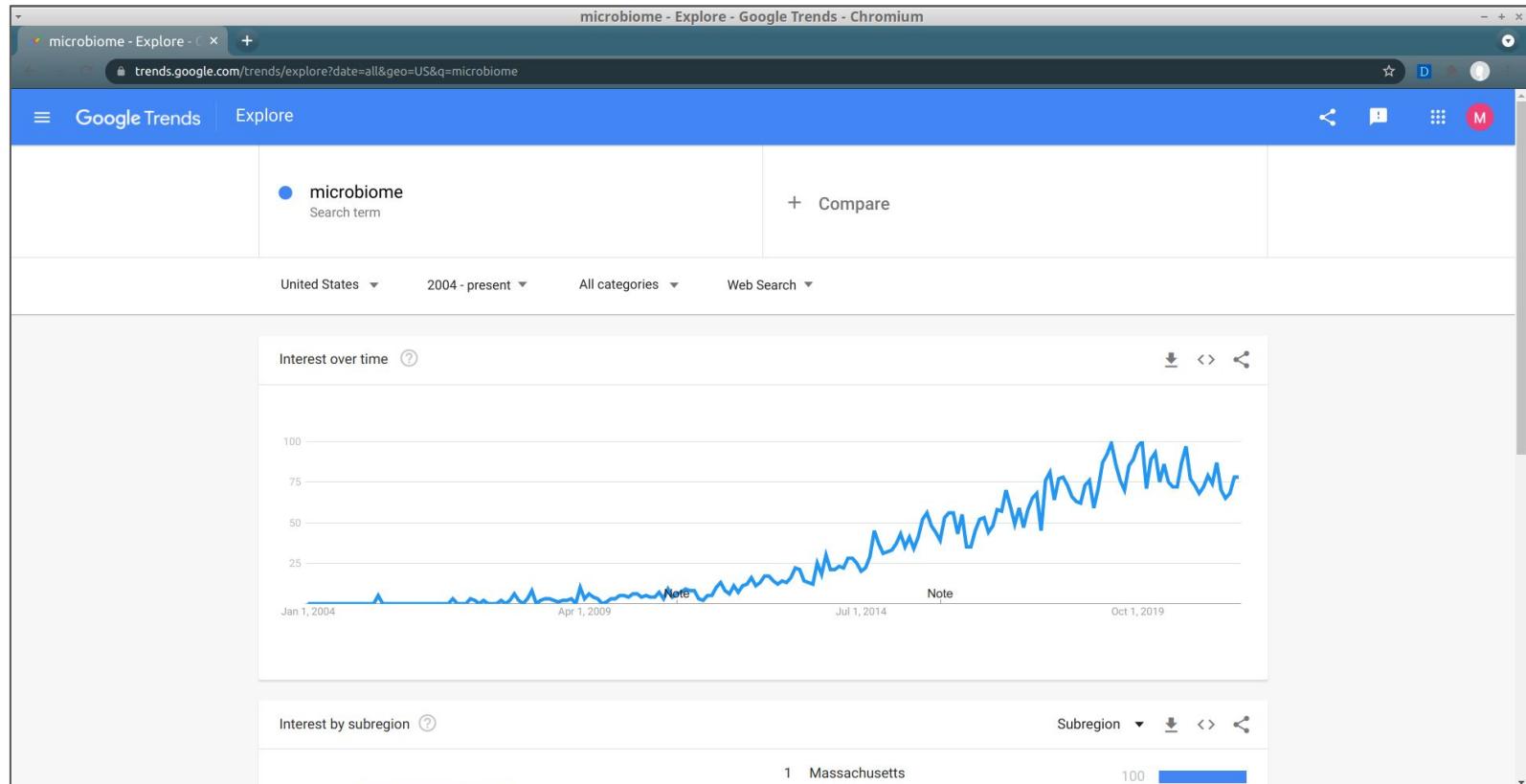
This talk

1. Introduction: Studying microbiomes
 - a. Why bother?
 - b. Why hasn't this research been more useful?
 - c. Defining a goal for this talk
2. Culture-independent (a.k.a. sequencing-based) methods
 - a. Marker gene sequencing
 - b. Metagenome sequencing
3. Metagenome assembly
 - a. Input (*reads*)
 - b. Outputs (*contigs*, an *assembly graph*, ...)
 - c. Methods (*de novo* vs. reference-based, overlap graph vs. de Bruijn graph, ...)
4. Future work: Solving the *strain separation problem*

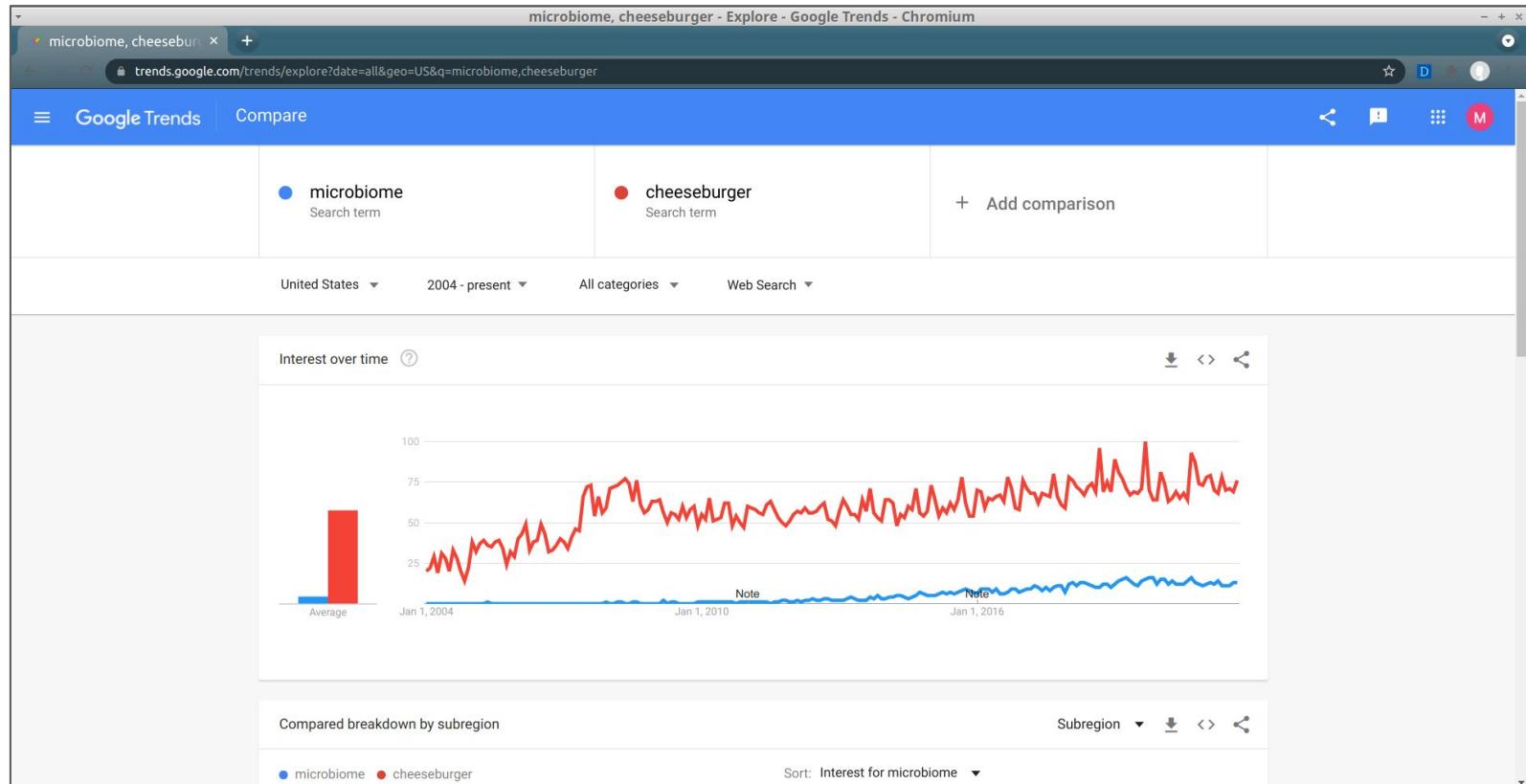
This talk

1. **Introduction: Studying microbiomes**
 - a. Why bother?
 - b. Why hasn't this research been more useful?
 - c. Defining a goal for this talk
2. **Culture-independent (a.k.a. sequencing-based) methods**
 - a. Marker gene sequencing
 - b. Metagenome sequencing
3. **Metagenome assembly**
 - a. Input (*reads*)
 - b. Outputs (*contigs*, an *assembly graph*, ...)
 - c. Methods (*de novo* vs. reference-based, overlap graph vs. de Bruijn graph, ...)
4. **Future work: Solving the *strain separation problem***

Introduction: Why bother studying microbiomes?



Introduction: Why bother studying microbiomes?

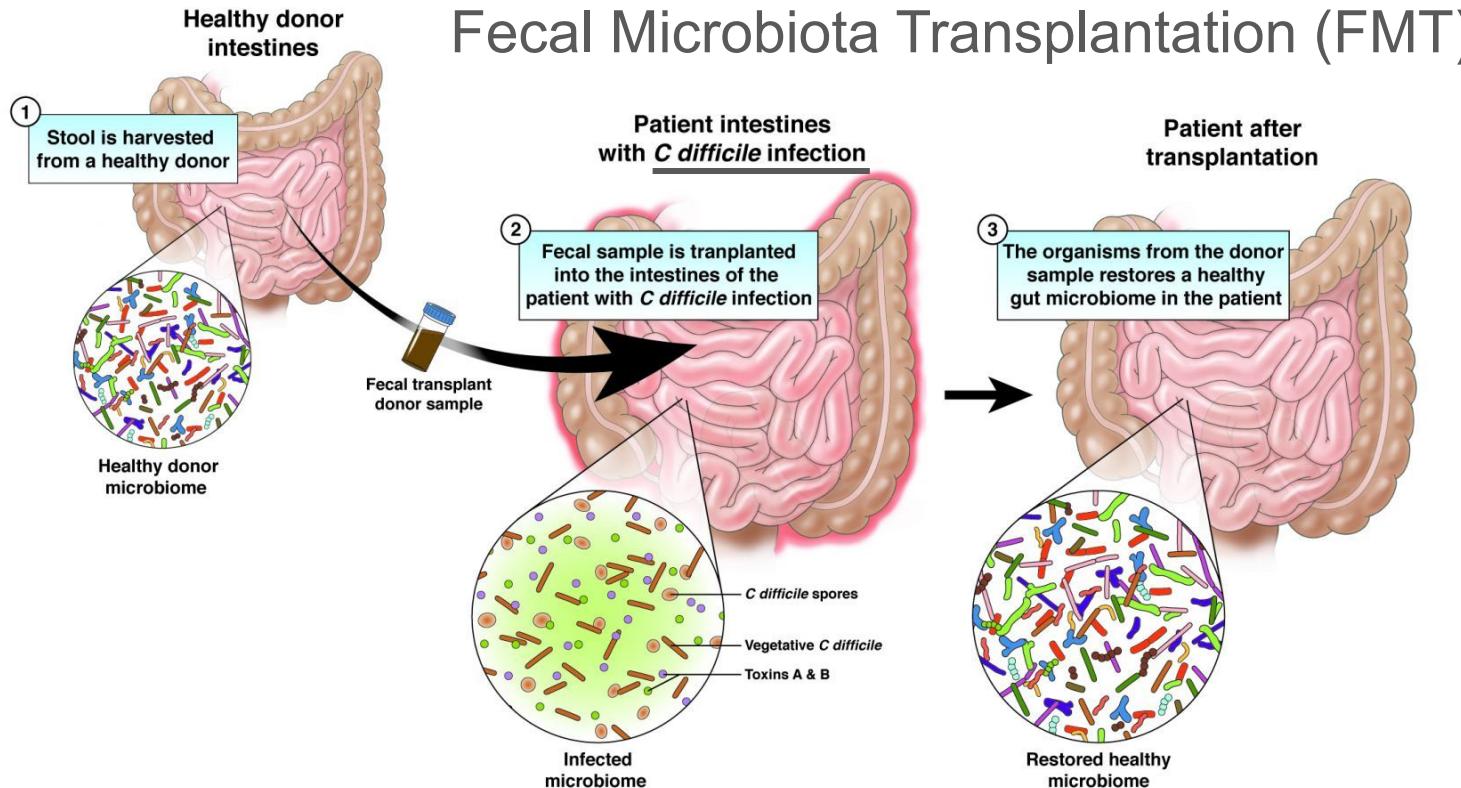


Introduction: Why bother studying microbiomes?

“During the Eastern Jin dynasty (AD 300–400 years), ‘*Zhou Hou Bei Ji Fang*’, a well-known monograph of traditional Chinese medicine (TCM) written by Hong Ge, recorded a case of **treating patients with food poisoning or severe diarrhea by ingesting human fecal suspension** (known as yellow soup or Huang-Long decoction).”

H. Du, T.-t. Kuang, S. Qiu, T. Xu, C.-L. G. Huan, G. Fan, and Y. Zhang. Fecal medicines used in traditional medical system of China: a systematic review of their names, original species, traditional uses, and modern investigations. *Chinese medicine*, 14(1):1–16, 2019.

Introduction: Why bother studying microbiomes?



Introduction: Why bother studying microbiomes?



Clinical Gastroenterology and Hepatology

Volume 9, Issue 12, December 2011, Pages 1044-1049



Perspective

Treating *Clostridium difficile* Infection With Fecal Microbiota Transplantation

Johan S. Bakken *, Thomas Borody ‡, Lawrence J. Brandt §, Joel V. Brill ||, Daniel C. Demarco ¶, Marc Alaric Franzos #,
Colleen Kelly **, Alexander Khoruts #‡, Thomas Louie §§, Lawrence P. Martinelli ||||, Thomas A. Moore ¶¶, George
Russell ##, Christina Surawicz ***, Fecal Microbiota Transplantation Workgroup

J Med Life. 2016 Apr-Jun; 9(2): 160–162.

PMCID: PMC4863507

PMID: [27453747](#)

Fecal transplantation – the new, inexpensive, safe, and rapidly effective approach in the treatment of gastrointestinal tract diseases

R Oprita, * M Bratu, * B Oprita, * and B Diaconescu *

ALIMENTARY TRACT: ORIGINAL ARTICLES

Fecal Microbiota Transplantation for the Treatment of *Clostridium difficile* Infection A Systematic Review

Cammarota, Giovanni MD; Ianiro, Gianluca MD; Gasbarrini, Antonio MD

Fecal Microbiota Transplant for Treatment of *Clostridium difficile* Infection in Immunocompromised Patients

Colleen R. Kelly, MD, FACP,¹ Chioma Ihunna, MD, MPH,¹ Monika Fischer, MD, MSCR,² Alexander Khoruts, MD,³ Christina Surawicz, MD, MACG,⁴ Anita Afzali, MD, MPH,⁴ Olga Aroniadis, MD,⁵ Amy Barto, MD,⁶ Thomas Borody, MD, PhD, FACP,⁷ Andrea Giovannelli, BS,⁸ Shelley Gordon, MD, PhD,⁹ Michael Gluck, MD,¹⁰ Elizabeth L. Hohmann, MD,¹¹ Dina Kao, MD,¹² John Y. Kao, MD,¹³ Daniel P. McQuillen, MD,⁶ Mark Mellow, MD, FACP,¹⁴ Kevin M. Rank, MD,³ Krishna Rao, MD,¹³ Armap Ray, MD,¹⁵ Margot A. Schwartz, MD, MPH,¹⁰ Namita Singh, MD,¹⁶ Neil Stollman, MD, FACP,⁸ David L. Suskind, MD,¹⁶ Stephen M. Vindigni, MD, MPH,⁴ Ilan Youngster, MD,¹¹ and Lawrence Brandt, MD, MACG⁵

Preliminary Communication

Oral, Capsulized, Frozen Fecal Microbiota Transplantation for Relapsing *Clostridium difficile* Infection

Ilan Youngster, MD, MMSc; George H. Russell, MD, MSc; Christina Pinder, BA; Tomer Ziv-Baran, PhD; Jenny Sauk, MD; Elizabeth L. Hohmann, MD

Introduction: Why bother studying microbiomes?

ORIGINAL ARTICLE

Helicobacter pylori Infection and the Risk of Gastric Carcinoma

Julie Parsonnet, M.D., Gary D. Friedman, M.D., M.S., Daniel P. Vandersteen, Yuan Chang, M.D., Joseph H. Vogelman, D.E.E., Norman Orentreich, M.D., and Richard K. Sibley, M.D.

Infection with *Helicobacter pylori* Strains Possessing *cagA* Is Associated with an Increased Risk of Developing Adenocarcinoma of the Stomach¹

Martin J. Blaser,² Guillermo I. Perez-Perez, Harry Kleanthous, Timothy L. Cover, Richard M. Peek, P. H. Chyou, Grant N. Stemmermann, and Abraham Nomura

Helicobacter pylori Persistence: an Overview of Interactions between *H. pylori* and Host Immune Defenses

Holly M. Scott Algood¹ and Timothy L. Cover^{1,2,3,*}

The gastric microbiome, its interaction with *Helicobacter pylori*, and its potential role in the progression to stomach cancer

Jennifer M. Noto^{1,*}, Richard M. Peek, Jr.^{1,2,3}

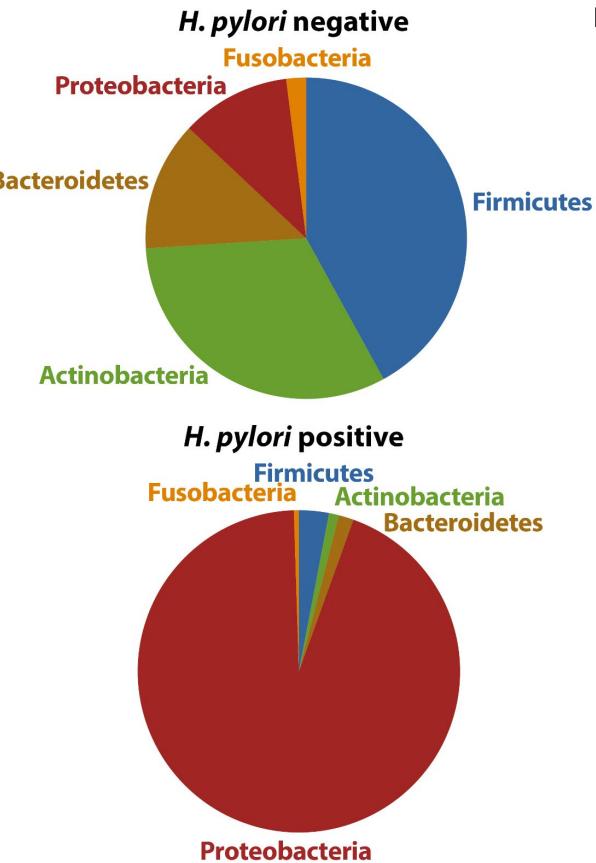
Introduction: Why bother studying microbiomes?

The gastric microbiome, its interaction with
Helicobacter pylori, and its potential role in the
progression to stomach cancer

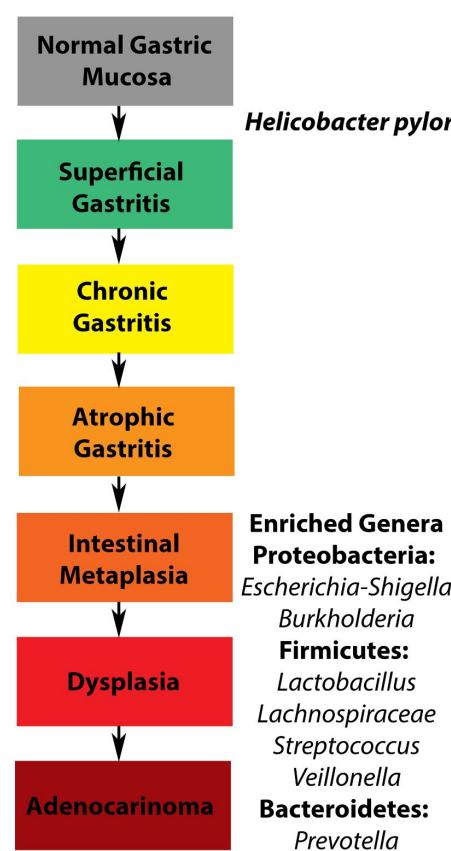
Jennifer M. Noto^{1*}, Richard M. Peek, Jr.^{1,2,3}

Introduction: Why bother studying microbiomes?

A



B



“*H. pylori*-negative individuals harbor a microbiota that is more complex and highly diverse compared to *H. pylori*-positive individuals. [...] Following infection with *H. pylori*, Proteobacteria and specifically *H. pylori* dominate the gastric microbiota. This leads to the development of chronic gastritis.”

The gastric microbiome, its interaction with *Helicobacter pylori*, and its potential role in the progression to stomach cancer

Introduction: Why is it always *C. diff* and *H. pylori*?

“There are two well-documented diseases in the microbiome field that link a microbial biomarker with causation in disease: ***Helicobacter pylori*-associated peptic ulceration and gastric cancer** (Parsonnet et al., 1991) and ***Clostridium* (or *Clostridioides*) difficile infection-associated diarrhea** (Gupta et al., 2016).”

“However, causal inferences between complex microbiomes and other inflammatory, metabolic, neoplastic, and neuro-behavioral disorders have been **neither compelling nor conclusive** [...]”

Introduction: Obesity and the gut microbiome

An obesity-associated gut microbiome with increased capacity for energy harvest

Peter J. Turnbaugh¹, Ruth E. Ley¹, Michael A. Mahowald¹, Vincent Magrini², Elaine R. Mardis^{1,2} & Jeffrey I. Gordon¹

Introduction: Obesity and the gut microbiome

An obesity-associated gut microbiome with increased capacity for energy harvest

Peter J. Turnbaugh¹, Ruth E. Ley¹, Michael A. Mahowald¹, Vincent Magrini², Elaine R. Mardis^{1,2} & Jeffrey I. Gordon¹

“Adult germ-free [wild-type] C57BL/ 6J mice were colonized (by gavage) with a microbiota harvested from the caecum of **obese (ob/ob)** or **lean (+/+)** donors (1 donor and 4–5 germ-free recipients per treatment group per experiment; two independent experiments). [...]”

(Context: **ob/ob mice** have a specific mutation “[...] that produces a stereotyped, fully penetrant obesity phenotype”; **+/+ mice** lack this mutation.)

Introduction: Obesity and the gut microbiome

An obesity-associated gut microbiome with increased capacity for energy harvest

Peter J. Turnbaugh¹, Ruth E. Ley¹, Michael A. Mahowald¹, Vincent Magrini², Elaine R. Mardis^{1,2} & Jeffrey I. Gordon¹

“Adult germ-free [wild-type] C57BL/ 6J mice were colonized (by gavage) with a microbiota harvested from the caecum of **obese (ob/ob)** or **lean (+/+)** donors (1 donor and 4–5 germ-free recipients per treatment group per experiment; two independent experiments). [...]”

Results: “Strikingly, mice colonized with an **ob/ob** microbiota exhibited a **significantly greater percentage increase in body fat over two weeks** than mice colonized with a **+/+ microbiota** [...]”

Introduction: Obesity and the gut microbiome

An obesity-associated gut microbiome with increased capacity for energy harvest

Peter J. Turnbaugh¹, Ruth E. Ley¹, Michael A. Mahowald¹, Vincent Magrini², Elaine R. Mardis^{1,2} & Jeffrey I. Gordon¹

Gut microbiome, obesity, and metabolic dysfunction

Herbert Tilg¹ and Arthur Kaser²

Looking for a Signal in the Noise: Revisiting Obesity and the Microbiome

Marc A. Sze, Patrick D. Schloss

The Human Microbiome and Obesity: Moving beyond Associations

Padma Maruvada¹, Vanessa Leone², Lee M. Kaplan³, Eugene B. Chang²✉

GASTROINTESTINAL INFECTIONS: EDITED BY MITCHELL COHEN

Obesity and the human microbiome

Ley, Ruth E

A Taxonomic Signature of Obesity in the Microbiome?
Getting to the Guts of the Matter

Mariel M. Finucane, Thomas J. Sharpton, Timothy J. Laurent, Katherine S. Pollard ✉

Introduction:

An obesity-associated gut microbiome with increased

Peter J. Turnbaugh¹, Ruth E. Ley¹

Gut microbiome

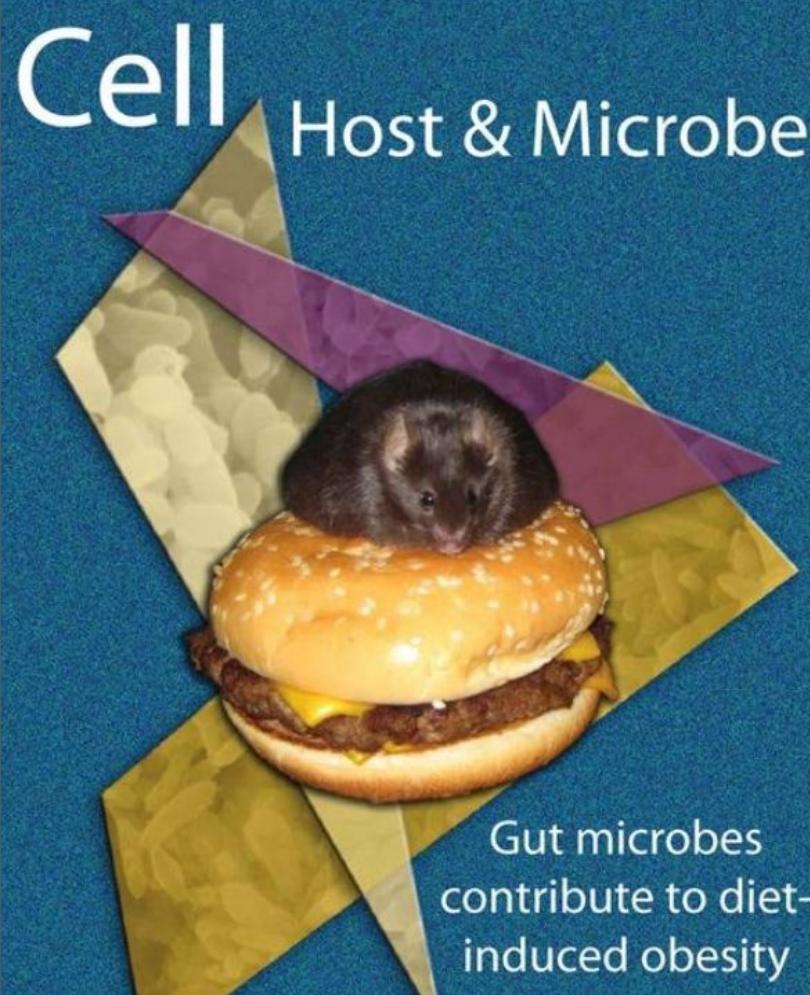
Herbert Tilg¹ and Arthur

Looking for a Signal in the Noise: Revisiting the Human Microbiome

Marc A. Sze, Patrick D. Schloss

The Human Microbiome and Beyond Associations

Padma Maruvada¹, Vanessa Leone², Lee M. Kaplan³, Eugene B.



Turnbaugh, P. J. (2017). Microbes and diet-induced obesity: fast, cheap, and out of control. *Cell Host & Microbe*, 21(3), 278-281.

obiome
robiome
ergy harvest
R. Mardis^{1,2} & Jeffrey I. Gordon¹
lic dysfunction

ITED BY MITCHELL COHEN

human microbiome

of Obesity in the Microbiome? The Matter

. Laurent, Katherine S. Pollard

Introduction: Obesity and the gut microbiome?

Looking for a Signal in the Noise: Revisiting Obesity and the Microbiome

 Marc A. Sze,  Patrick D. Schloss

Department of Microbiology and Immunology, University of Michigan, Ann Arbor, Michigan, USA

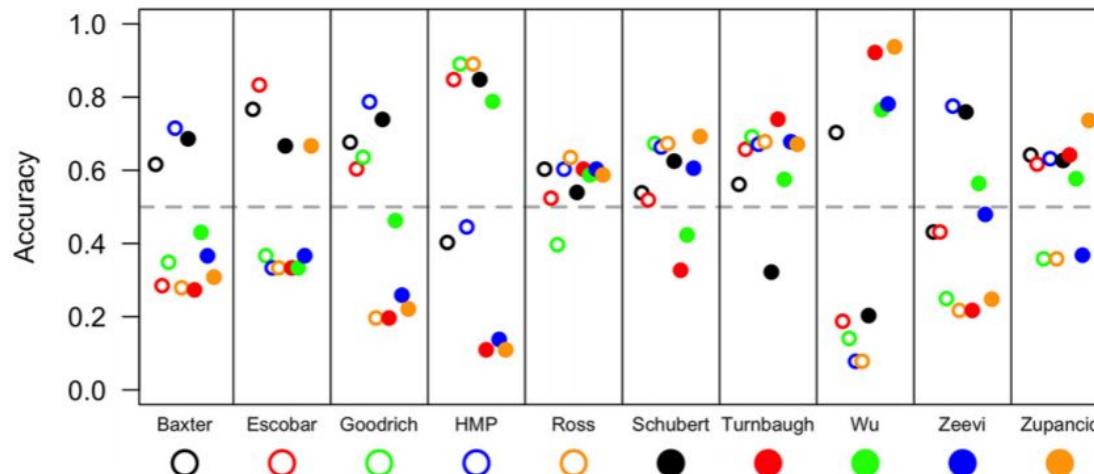


FIG 5 Overall accuracy with which each study predicted nonobese and obese individuals on the basis of that study's random forest machine learning model applied to each of the other studies. HMP, Human Microbiome Project.

Introduction: Why is it always *C. diff* and *H. pylori*?

Microbiome science needs a healthy dose of scepticism

The human microbiome in health and disease: hype or hope

Gwen Falony ^{a,b}, Doris Vandepitte ^{a,b}, Clara Caenepeel^c, Sara Vieira-Silva ^{a,b}, Tanine Daryoush^{a,b}, Séverine Vermeire ^c and Jeroen Raes^{a,b}

Separating the microbiome from the hyperbolome

Fergus Shanahan

Establishing or Exaggerating Causality for the Gut Microbiome: Lessons from Human Microbiota-Associated Rodents

Jens Walter,^{1,2,3,4,8,*} Anissa M. Armet,^{1,8} B. Brett Finlay,^{5,6,7} and Fergus Shanahan³

Introduction: Why is this so difficult?

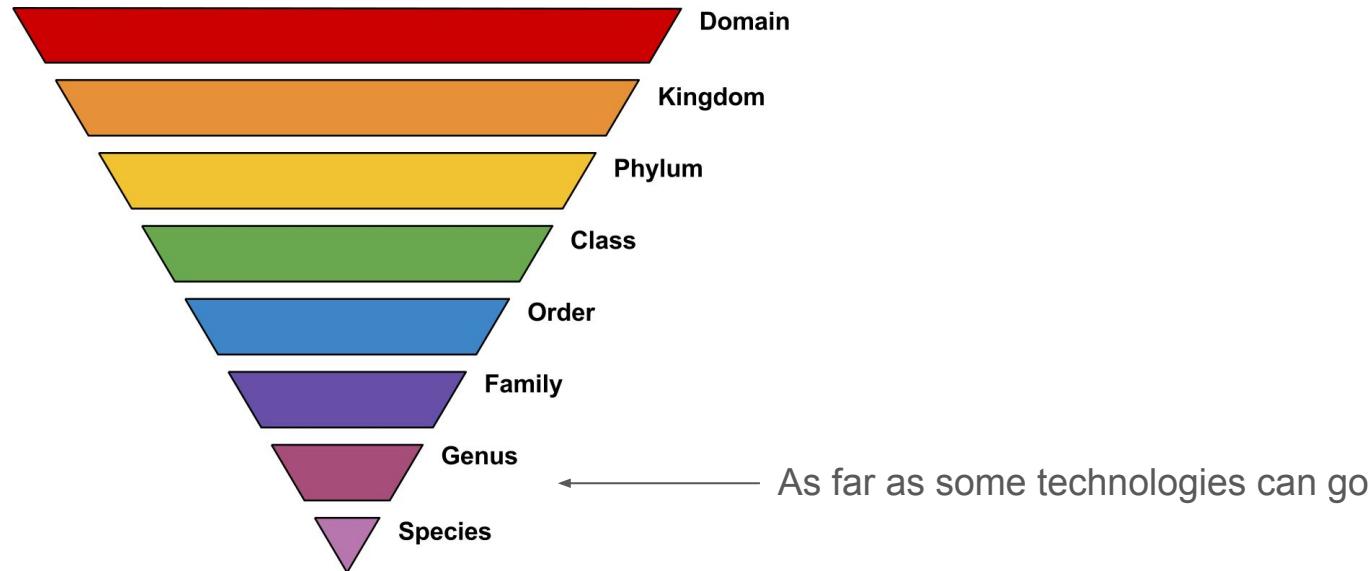
- Challenges common to many areas of science
 - Limited sample sizes
 - Institutional and personal biases against the publication of null results
 - Failure to pre-register studies
 - Hypothesizing after the collection of data without reporting the study as exploratory
- Challenges more specific to microbiome / bioinformatics research
 - Limitations of using mice (or other non-human organisms) as models
 - The “curse of dimensionality”
 - Sparsity
 - Compositionality
 - Uneven sampling depths
 - Methods that only provide limited resolution about the types of microbes in a sample

Introduction: Why is this so difficult?

- Challenges common to many areas of science
 - Limited sample sizes
 - Institutional and personal biases against the publication of null results
 - Failure to pre-register studies
 - Hypothesizing after the collection of data without reporting the study as exploratory
- Challenges more specific to microbiome / bioinformatics research
 - Limitations of using mice (or other non-human organisms) as models
 - The “curse of dimensionality”
 - Sparsity
 - Compositionality
 - Uneven sampling depths
 - **Methods that only provide limited resolution about the types of microbes in a sample**

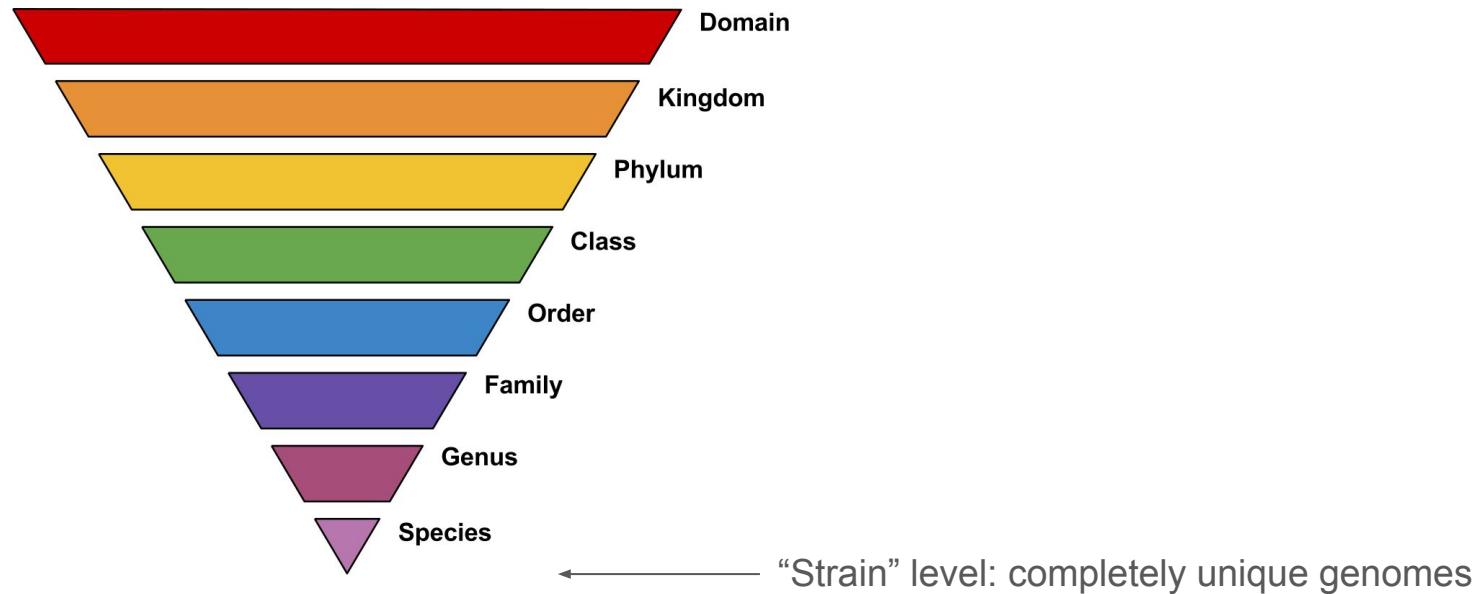
Introduction: Our goal for this talk

Improving the **resolution** with which we can study microbes.



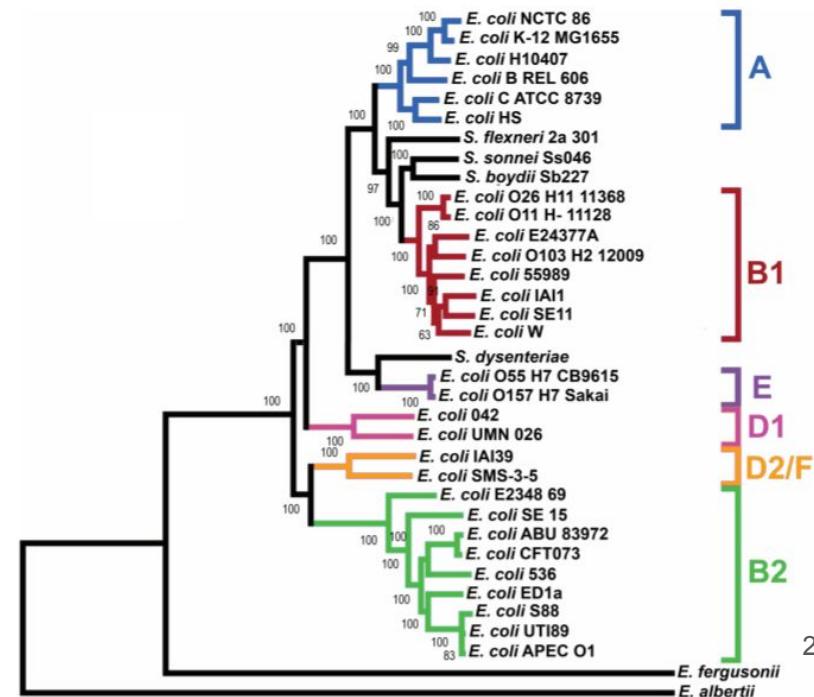
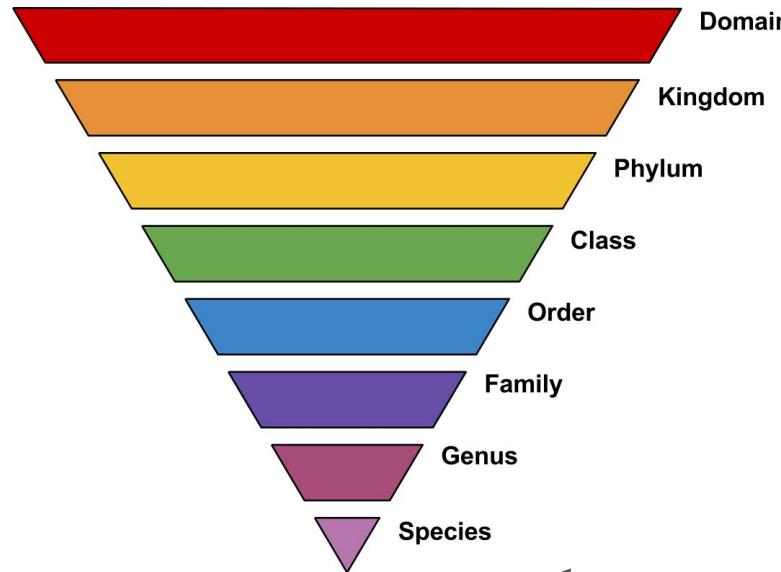
Introduction: Our goal for this talk

Improving the **resolution** with which we can study microbes.



Introduction: Our goal for this talk

Improving the **resolution** with which we can study microbes.



Introduction: Our goal for this talk

Improving the **resolution** with which we can study microbes.
Small strain-level differences can make a big difference!

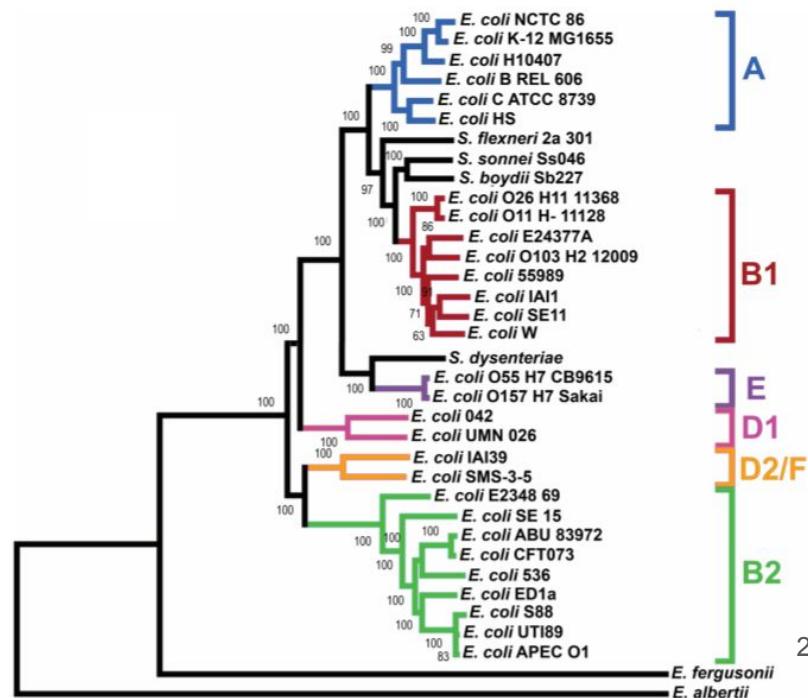
Epidemic Profile of Shiga-Toxin–Producing *Escherichia coli* O104:H4 Outbreak in Germany

Christina Frank, Ph.D., Dirk Werber, D.V.M., Jakob P. Cramer, M.D., Mona Askar, M.D.,
Mirko Faber, M.D., Matthias an der Heiden, Ph.D., Helen Bernard, M.D., Angelika Fruth,
Ph.D., Rita Prager, Ph.D., Anke Spode, M.D., Maria Wadl, D.V.M., Alexander Zoufaly,
M.D., Sabine Jordan, M.D., Markus J. Kemper, M.D., Per Follin, M.D., Ph.D., Luise
Müller, M.Sc., Lisa A. King, M.P.H., Bettina Rosner, Ph.D., Udo Buchholz, M.D., M.P.H.,
Klaus Stark, M.D., Ph.D., and Gérard Krause, M.D., Ph.D. for the HUS Investigation

Team*

Escherichia coli induces DNA damage in vivo and triggers genomic instability in mammalian cells

Gabriel Cuevas-Ramos^{a,b}, Claude R. Petit^{a,b}, Ingrid Marcq^{a,b}, Michèle Boury^{a,b}, Eric Oswald^{a,b,c,d}, and
Jean-Philippe Nougayrède^{a,b,1}



Introduction: Our goal for this talk

Improving the **resolution** with which we can study microbes.

Small strain-level differences can make a big difference!

Our goal, then, is **reconstructing the full genomes of all microbes in a sample**.

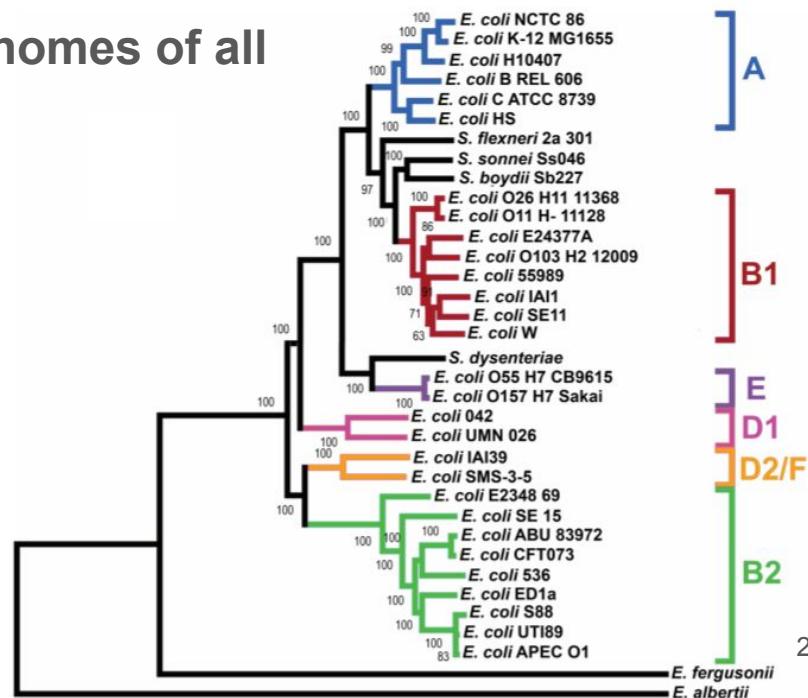
Epidemic Profile of Shiga-Toxin–Producing *Escherichia coli* O104:H4 Outbreak in Germany

Christina Frank, Ph.D., Dirk Werber, D.V.M., Jakob P. Cramer, M.D., Mona Askar, M.D.,
Mirko Faber, M.D., Matthias an der Heiden, Ph.D., Helen Bernard, M.D., Angelika Fruth,
Ph.D., Rita Prager, Ph.D., Anke Spode, M.D., Maria Wadl, D.V.M., Alexander Zoufaly,
M.D., Sabine Jordan, M.D., Markus J. Kemper, M.D., Per Follin, M.D., Ph.D., Luise
Müller, M.Sc., Lisa A. King, M.P.H., Bettina Rosner, Ph.D., Udo Buchholz, M.D., M.P.H.,
Klaus Stark, M.D., Ph.D., and Gérard Krause, M.D., Ph.D. for the HUS Investigation

Team*

***Escherichia coli* induces DNA damage in vivo and triggers genomic instability in mammalian cells**

Gabriel Cuevas-Ramos^{a,b}, Claude R. Petit^{a,b}, Ingrid Marcq^{a,b}, Michèle Boury^{a,b}, Eric Oswald^{a,b,c,d}, and Jean-Philippe Nougayrède^{a,b,1}



Introduction: Our goal for this talk

Improving the **resolution** with which we can study microbes.

Small strain-level differences can make a big difference!

Our goal, then, is **reconstructing the full genomes of all microbes in a sample**.

Introduction: Our goal for this talk

Improving the **resolution** with which we can study microbes.

Small strain-level differences can make a big difference!

Our goal, then, is **reconstructing the full genomes of all microbes in a sample**.

... This isn't really possible right now, but we'll see what we can do!

Introduction: Our goal for this talk

Improving the **resolution** with which we can study microbes.

Small strain-level differences can make a big difference!

Our goal, then, is **reconstructing the full genomes of all microbes in a sample**.

... This isn't really possible right now, but we'll see what we can do!

One final note for the introduction, though:

Introduction: Why is this so difficult?

- Challenges common to many areas of science
 - Limited sample sizes
 - Institutional and personal biases against the publication of null results
 - Failure to pre-register studies
 - Hypothesizing after the collection of data without reporting the study as exploratory
- Challenges more specific to microbiome / bioinformatics research
 - Limitations of using mice (or other non-human organisms) as models
 - The “curse of dimensionality”
 - Sparsity
 - Compositionality
 - Uneven sampling depths
 - **Methods that only provide limited resolution about the types of microbes in a sample**

Introduction: Why is this so difficult?

- Challenges common to many areas of science
 - Limited sample sizes
 - Institutional and personal biases against the publication of null results
 - Failure to pre-register studies
 - Hypothesizing after the collection of data without reporting the study as exploratory
- Challenges more specific to microbiome / bioinformatics research
 - Limitations of using mice (or other non-human organisms) as models
 - The “curse of dimensionality”
 - Sparsity
 - Compositionality
 - Uneven sampling depths
 - Methods that only provide limited resolution about the types of microbes in a sample

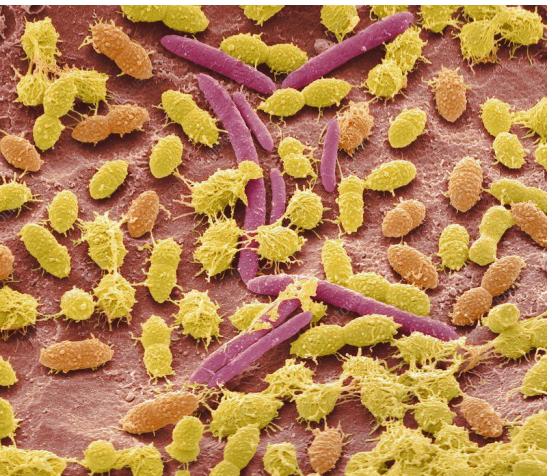


This talk

1. Introduction: Studying microbiomes
 - a. Why bother?
 - b. Why hasn't this research been more useful?
 - c. Defining a goal for this talk
2. **Culture-independent (a.k.a. sequencing-based) methods**
 - a. Marker gene sequencing
 - b. Metagenome sequencing
3. Metagenome assembly
 - a. Input (*reads*)
 - b. Outputs (*contigs*, an *assembly graph*, ...)
 - c. Methods (*de novo* vs. reference-based, overlap graph vs. de Bruijn graph, ...)
4. Future work: Solving the *strain separation problem*

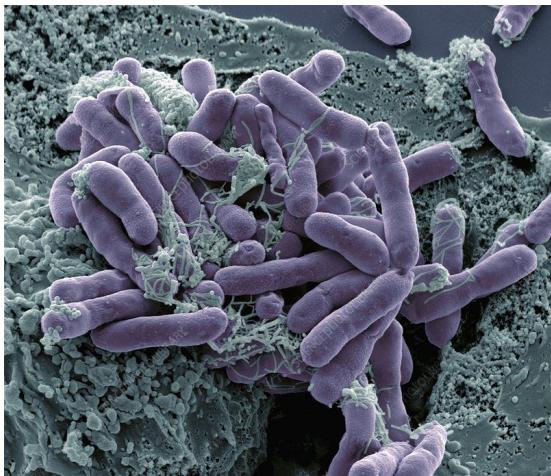
Culture-Independent Methods

Culture-Independent Methods



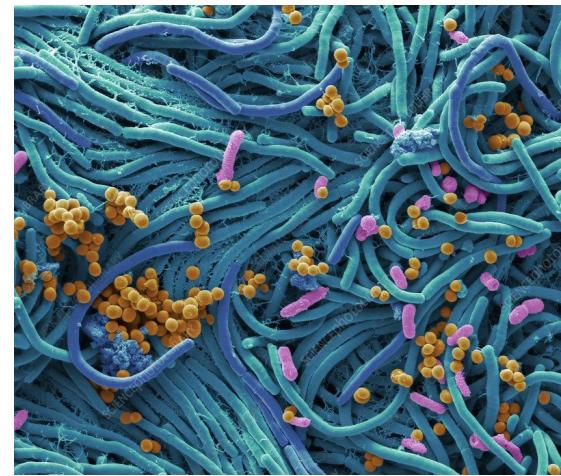
[Steve Gschmeissner, Science Photo Library](#)

“Scanning electron micrograph (SEM) of bacteria cultured from a sample of human faeces.”



[Steve Gschmeissner, Science Photo Library](#)

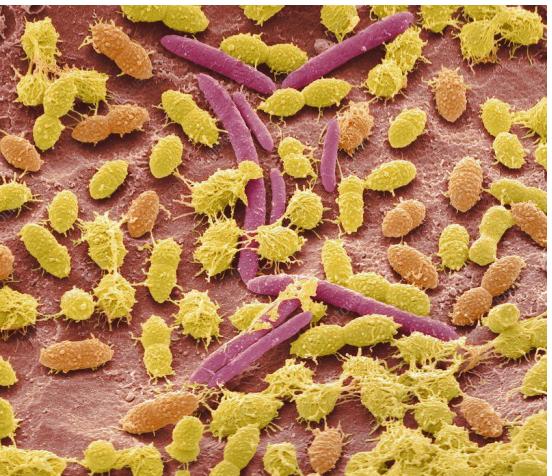
“Bacterial contamination, coloured scanning electron micrograph (SEM). *Escherichia coli* bacteria in a cell culture. This contamination has come from an unclean water source.”



[Steve Gschmeissner, Science Photo Library](#)

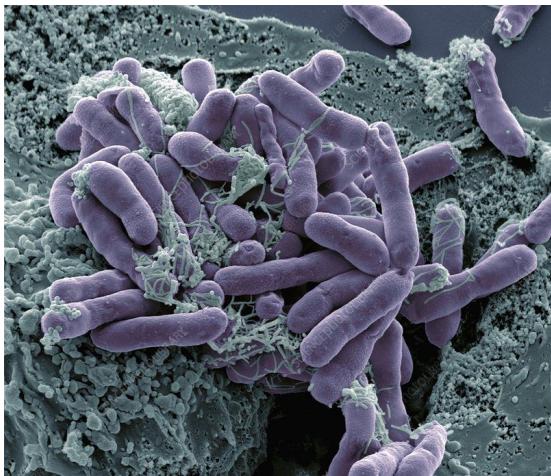
“Coloured scanning electron micrograph (SEM) of bacteria cultured from a mobile phone.”

Culture-Independent Methods



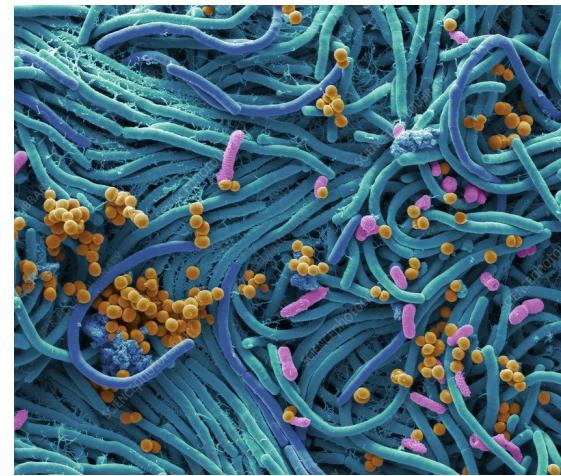
[Steve Gschmeissner, Science Photo Library](#)

“Scanning electron micrograph (SEM) of bacteria **cultured** from a sample of human faeces.”



[Steve Gschmeissner, Science Photo Library](#)

“Bacterial contamination, coloured scanning electron micrograph (SEM). *Escherichia coli* bacteria in a **cell culture**. This contamination has come from an unclean water source.”



[Steve Gschmeissner, Science Photo Library](#)

“Coloured scanning electron micrograph (SEM) of bacteria **cultured** from a mobile phone.”

Culture-Independent Methods

“It is estimated that **>99% of microorganisms observable in nature** typically are not cultivated by using standard techniques.”

Hugenholtz, P., Goebel, B. M., & Pace, N. R. (1998). Impact of culture-independent studies on the emerging phylogenetic view of bacterial diversity. *Journal of bacteriology*, 180(18), 4765-4774.

Culture-Independent Methods

“It is estimated that **>99% of microorganisms observable in nature** typically are not cultivated by using standard techniques.”

Hugenholtz, P., Goebel, B. M., & Pace, N. R. (1998). Impact of culture-independent studies on the emerging phylogenetic view of bacterial diversity. *Journal of bacteriology*, 180(18), 4765-4774.

Although not all microbes are easily culturable, all have **genomes**.

Idea: look at the DNA in a sample and use that to study the microbes there!

For a nice history of these and other methods, see: M. H. Fraher, P. W. O'Toole, and E. M. Quigley. Techniques used to characterize the gut microbiota: a guide for the clinician. *Nature Reviews Gastroenterology & Hepatology*, 9(6):312–322, 2012.

Culture-Independent Methods

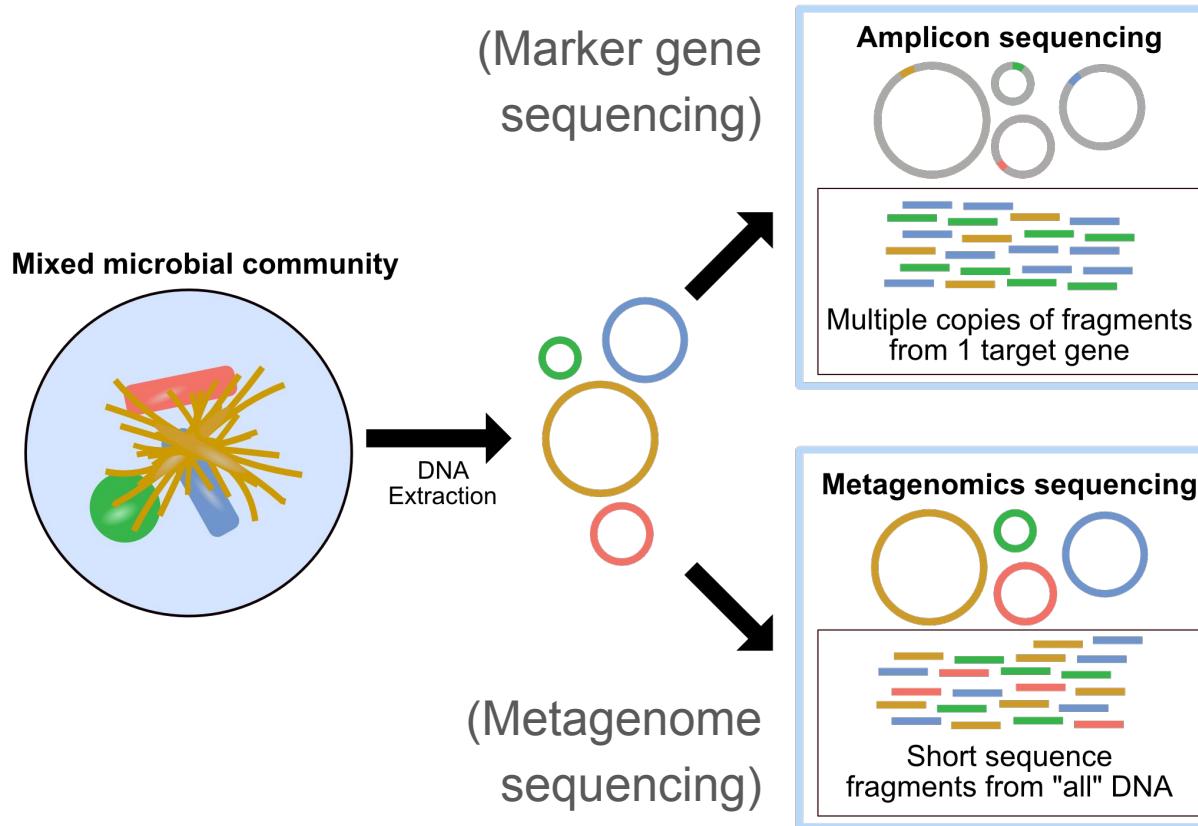
1. Terminal restriction fragment length polymorphism
a.k.a. *T-RFLP*
2. Denaturing / temperature gradient gel electrophoresis
a.k.a. *DGGE / TGGE*
3. Fluorescence *in situ* hybridization
a.k.a. *FISH*
4. Marker gene sequencing
a.k.a. *amplicon sequencing, metabarcoding, metataxonomics, ...*
5. Metagenome sequencing
a.k.a. *metagenomics, shotgun metagenome sequencing, whole metagenome sequencing, ...*

For a nice history of these and other methods, see: M. H. Fraher, P. W. O'Toole, and E. M. Quigley. Techniques used to characterize the gut microbiota: a guide for the clinician. *Nature Reviews Gastroenterology & Hepatology*, 9(6):312–322, 2012.

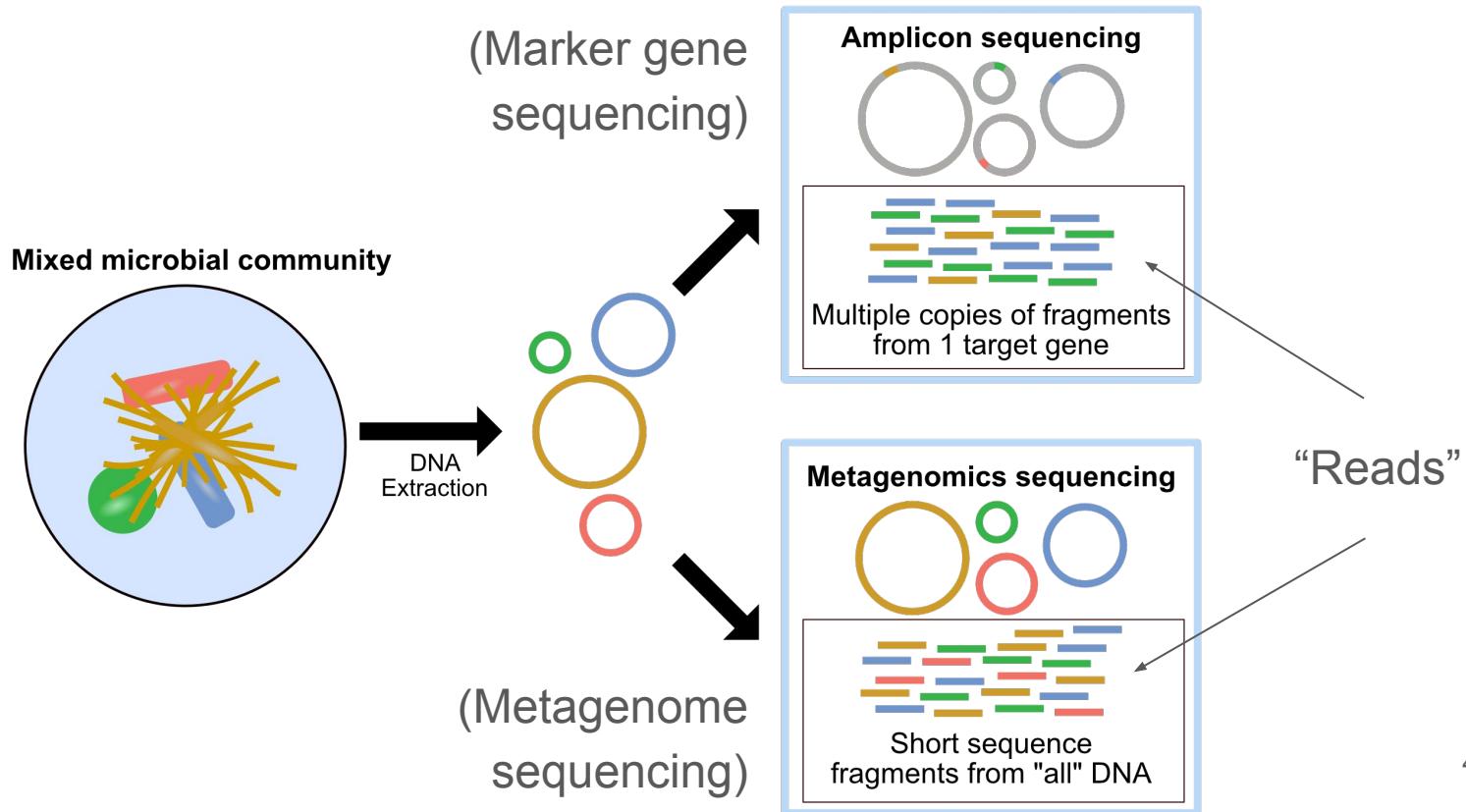
Culture-Independent Methods

1. Terminal restriction fragment length polymorphism
 - a.k.a. *T-RFLP*
2. Denaturing / temperature gradient gel electrophoresis
 - a.k.a. *DGGE / TGGE*
3. Fluorescence *in situ* hybridization
 - a.k.a. *FISH*
4. Marker gene sequencing
 - a.k.a. *amplicon sequencing, metabarcoding, metataxonomics, ...*
5. Metagenome sequencing
 - a.k.a. *metagenomics, shotgun metagenome sequencing, whole metagenome sequencing, ...*

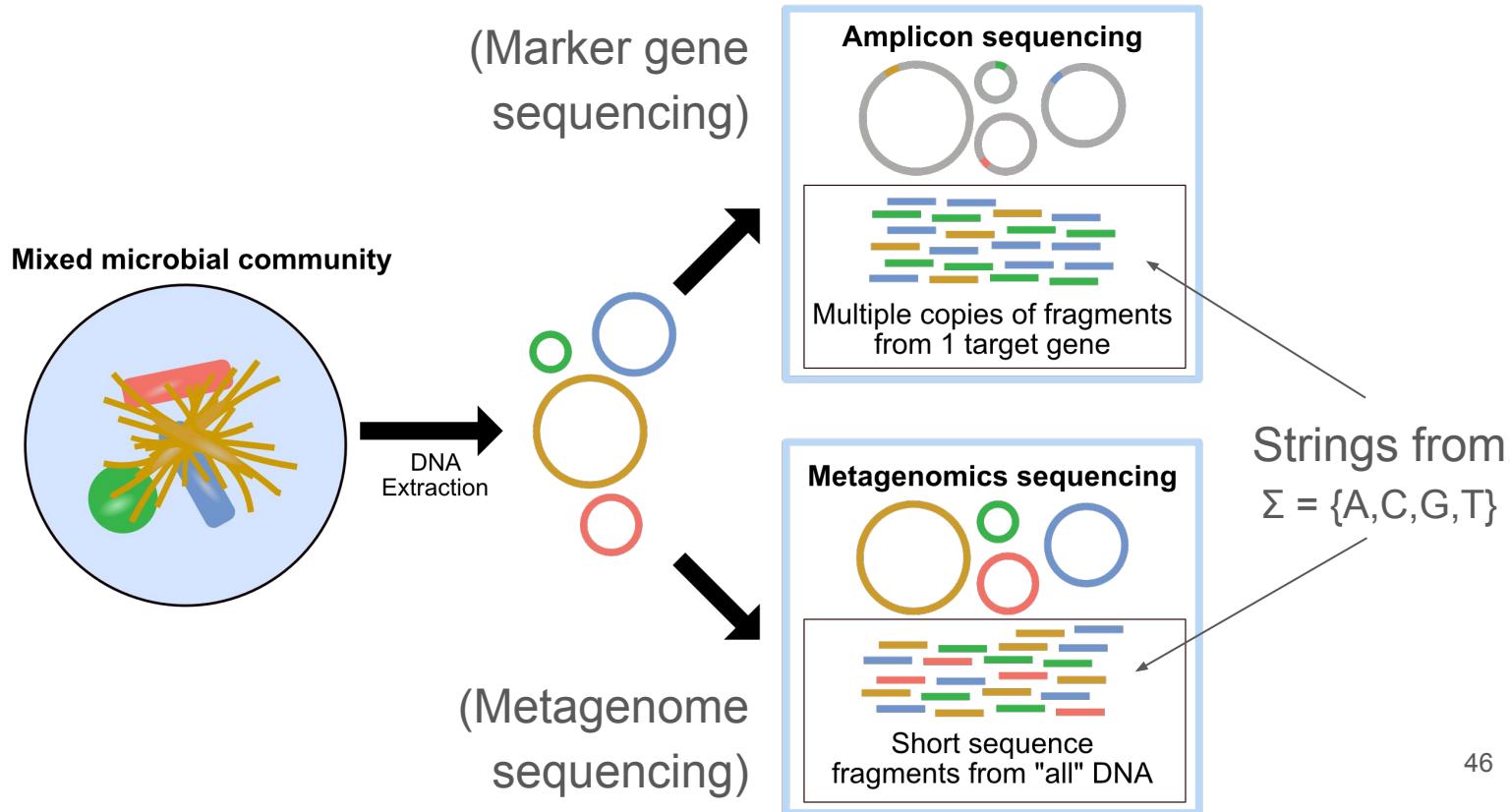
Culture-Independent Methods



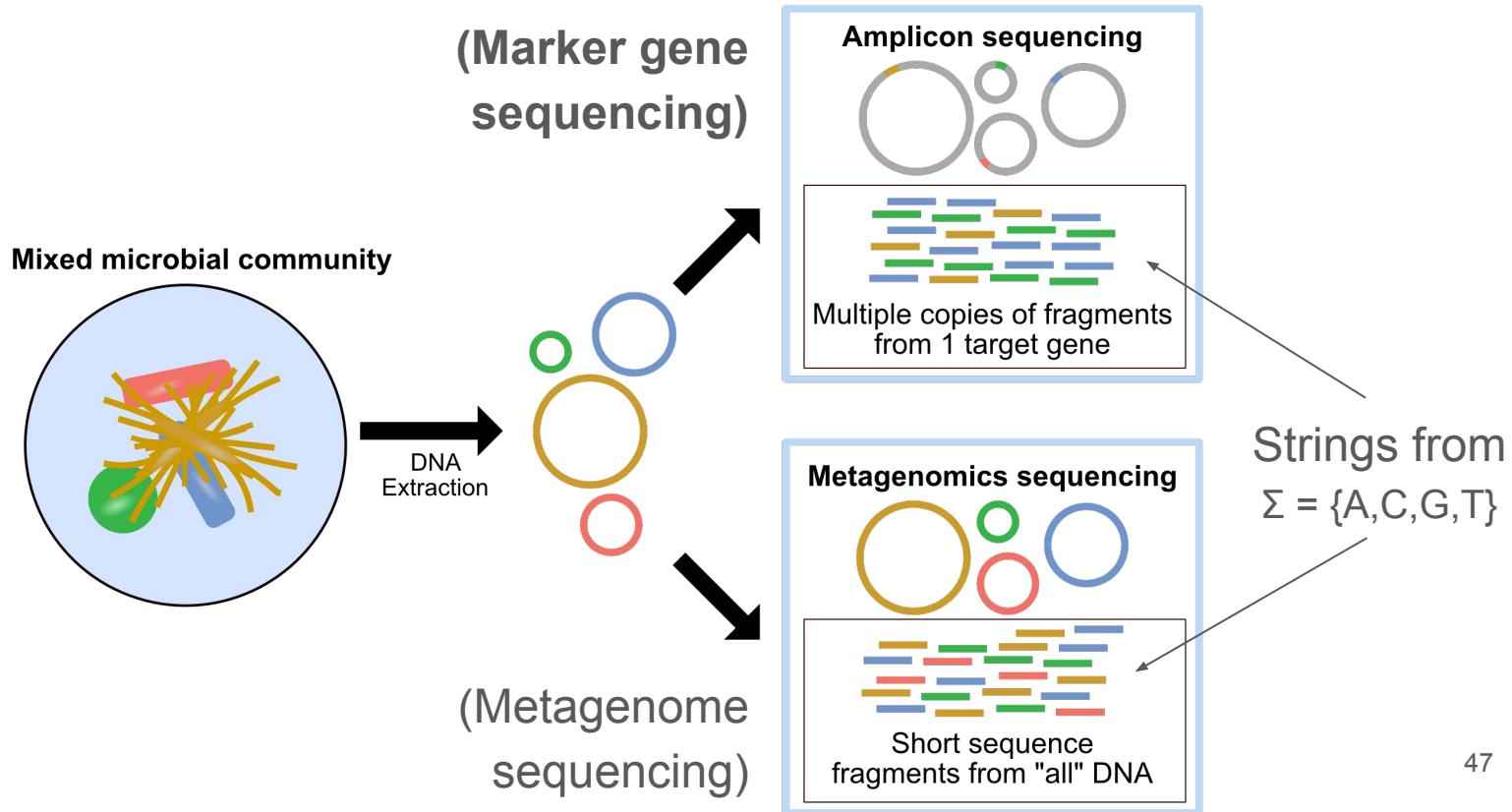
Culture-Independent Methods



Culture-Independent Methods



Culture-Independent Methods



C. I. Methods: Marker gene sequencing

C. I. Methods: Carl Woese and rRNA genes

Proc. Natl. Acad. Sci. USA
Vol. 74, No. 11, pp. 5088–5090, November 1977
Evolution

Phylogenetic structure of the prokaryotic domain: The primary kingdoms

(archaeabacteria/eubacteria/urkaryote/16S ribosomal RNA/molecular phylogeny)

CARL R. WOESE AND GEORGE E. FOX*

Department of Genetics and Development, University of Illinois, Urbana, Illinois 61801

ABSTRACT A phylogenetic analysis based upon ribosomal RNA sequence characterization reveals that living systems represent one of three aboriginal lines of descent:

- (i) the eubacteria, comprising all typical bacteria;
- (ii) the archaeabacteria, containing methanogenic bacteria; and
- (iii) the urkaryotes, now represented in the cytoplasmic component of eukaryotic cells.

C. I. Methods: Carl Woese and rRNA genes

Proc. Natl. Acad. Sci. USA
Vol. 74, No. 11, pp. 5088–5090, November 1977
Evolution

Phylogenetic structure of the prokaryotic domain: The primary kingdoms

(archaeabacteria/eubacteria/urkaryote/16S ribosomal RNA/molecular phylogeny)

CARL R. WOESE AND GEORGE E. FOX*

Department of Genetics and Development, University of Illinois, Urbana, Illinois 61801

ABSTRACT A phylogenetic analysis based upon ribosomal RNA sequence characterization reveals that living systems represent one of three aboriginal lines of descent:

- (i) the eubacteria, comprising all typical bacteria;
- (ii) **the archaeabacteria, containing methanogenic bacteria;** and
- (iii) the urkaryotes, now represented in the cytoplasmic component of eukaryotic cells.

C. I. Methods: Carl Woese and rRNA genes

Chicago Tribune
PAGE 43
📍 Chicago, Illinois
📅 Thursday, November 03, 1977
1 of 1 matches on this page.

Martianlike bugs may be oldest life

By Ronald Kotulak
Source: science

BACTERIA that live in the dark recesses of hot-spring vents, where no other organisms can survive, may be the oldest form of life on earth, according to new evidence.

These bacteria have been found because at the time they evolved and went into hiding, no rocks had yet formed.

Discovery of the bacteria may provide the missing link to understanding how life evolved on earth and how it may develop on other planets, said Dr. Carl

ABSTRACT A phylogenetic analysis based upon ribosomal RNA sequence characterization reveals that living systems represent one of three aboriginal lines of descent:

- (i) the eubacteria, comprising all typical bacteria;
- (ii) the archaeabacteria, containing methanogenic bacteria; and**
- (iii) the urkaryotes, now represented in the cytoplasmic component of eukaryotic cells.

For a more thorough history of Carl Woese's life and work, see:
D. Quammen. The Scientist Who Scrambled Darwin's Tree of Life. *The New York Times*, 2018.

C. I. Methods: Carl Woese and rRNA genes

The New York Times Magazine

PLAY THE CROSSWORD Account ▾

FEATURE

The Scientist Who Scrambled Darwin's Tree of Life

How the microbiologist Carl Woese fundamentally changed the way we think about evolution and the origins of life.

52

C. I. Methods: What's in a marker gene?

Proc. Natl. Acad. Sci. USA
Vol. 74, No. 11, pp. 5088–5090, November 1977
Evolution

Phylogenetic structure of the prokaryotic domain: The primary kingdoms

(archaeabacteria/eubacteria/urkaryote/16S ribosomal RNA/molecular phylogeny)

CARL R. WOESE AND GEORGE E. FOX*

Department of Genetics and Development, University of Illinois, Urbana, Illinois 61801

ABSTRACT A phylogenetic analysis based upon ribosomal RNA sequence characterization reveals that living systems represent one of three aboriginal lines of descent:

- (i) the eubacteria, comprising all typical bacteria;
- (ii) the archaeabacteria, containing methanogenic bacteria; and
- (iii) the urkaryotes, now represented in the cytoplasmic component of eukaryotic cells.

C. I. Methods: What's in a marker gene?

Proc. Natl. Acad. Sci. USA
Vol. 74, No. 11, pp. 5088–5090, November 1977
Evolution

Phylogenetic structure of the prokaryotic domain: The primary kingdoms

(archaeabacteria/eubacteria/urkaryote/16S ribosomal RNA/molecular phylogeny)

CARL R. WOESE AND GEORGE E. FOX*

Department of Genetics and Development, University of Illinois, Urbana, Illinois 61801

ABSTRACT A phylogenetic analysis based upon **ribosomal RNA sequence** characterization reveals that living systems represent one of three aboriginal lines of descent:

- (i) the eubacteria, comprising all typical bacteria;
- (ii) the archaeabacteria, containing methanogenic bacteria; and
- (iii) the urkaryotes, now represented in the cytoplasmic component of eukaryotic cells.

C. I. Methods: What's in a marker gene?

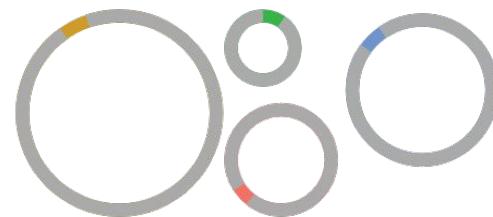
“To determine relationships covering the entire spectrum of extant living systems, one optimally needs a molecule of appropriately broad distribution. None of the readily characterized proteins fits this requirement. However, ribosomal RNA does. It is

- a component of all self-replicating systems;
- it is readily isolated;
- and **its sequence changes but slowly with time—**

permitting the detection of relatedness among very distant species.”

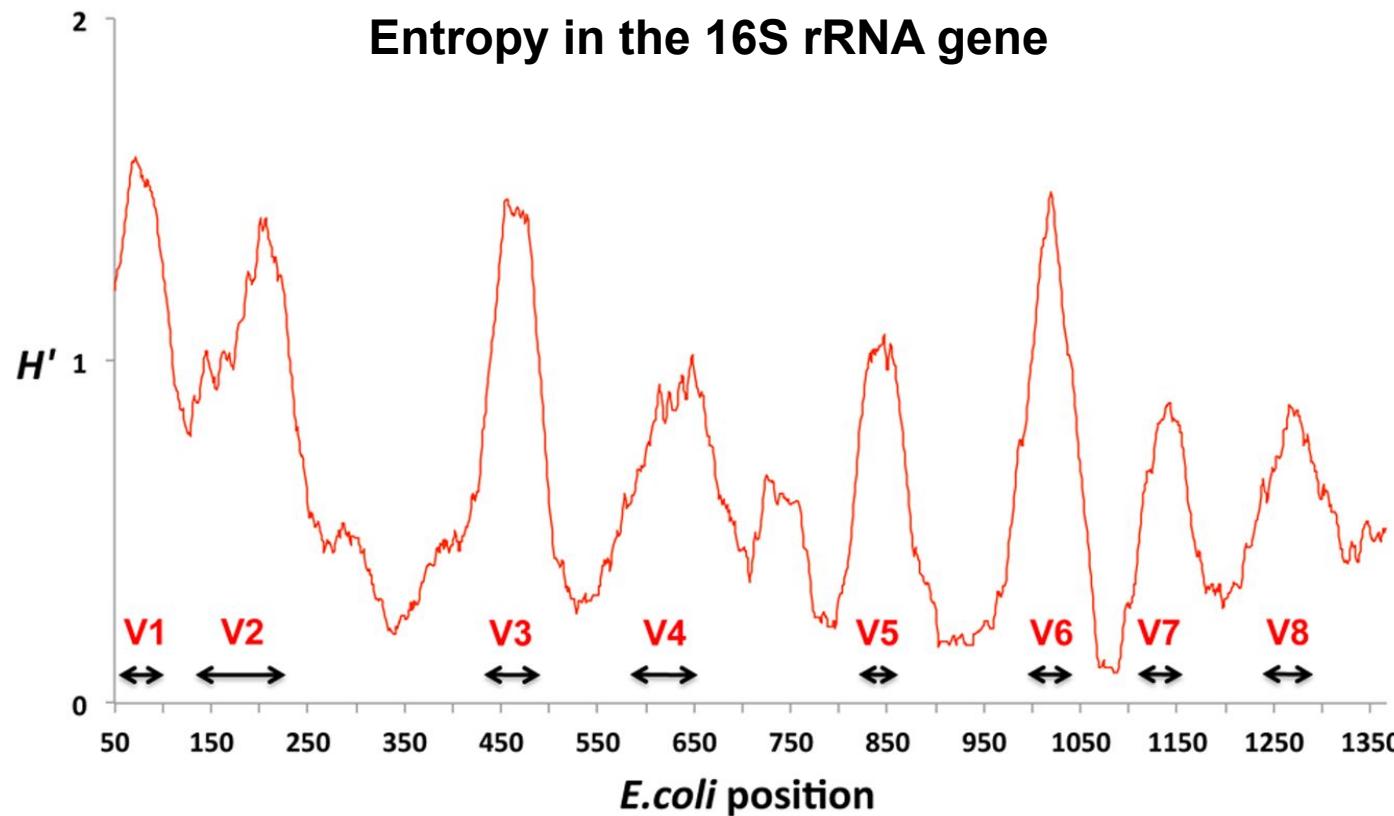
C. I. Methods: What's in a marker gene?

Amplicon sequencing

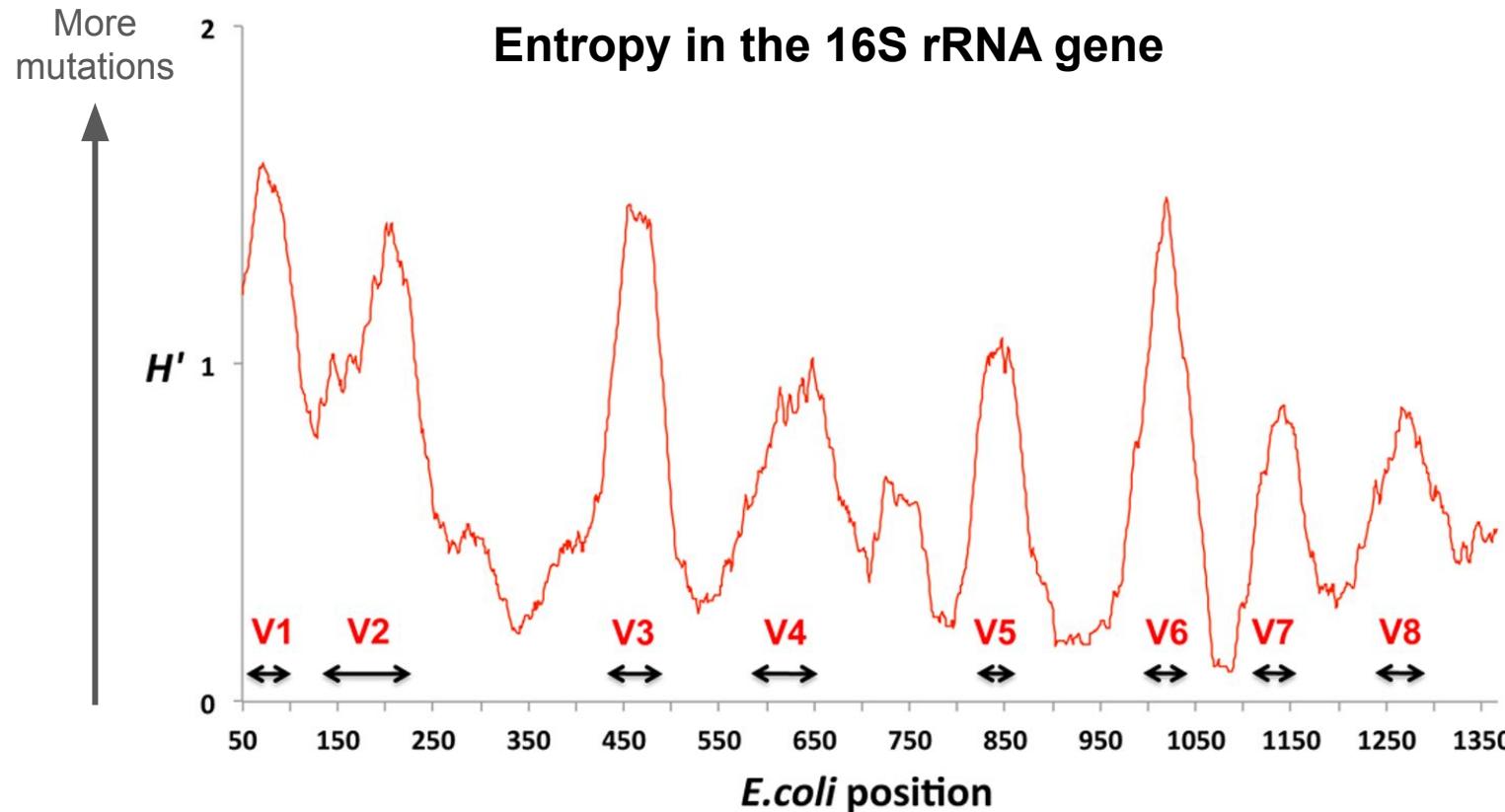


Multiple copies of fragments
from 1 target gene

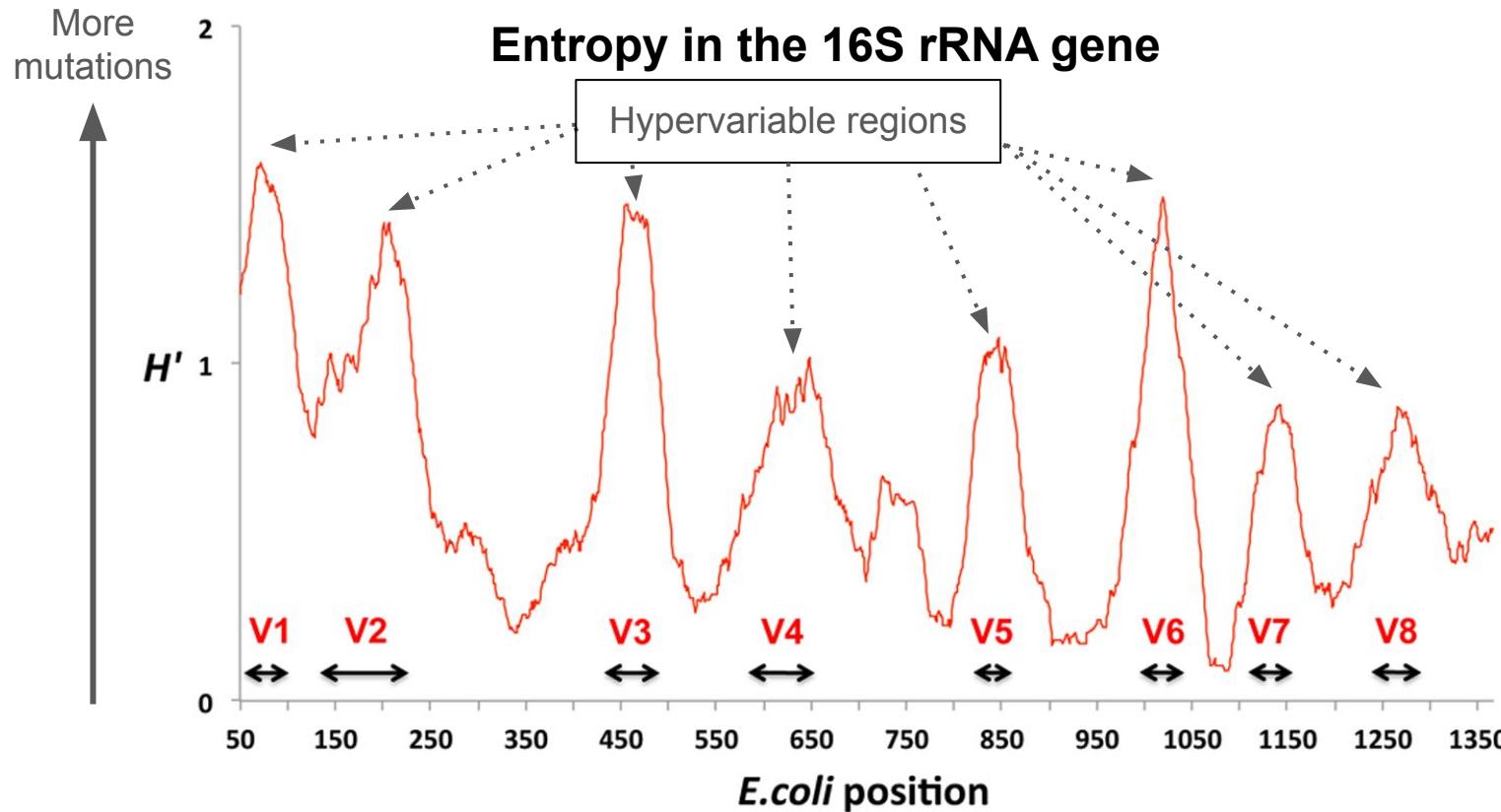
C. I. Methods: What's in a marker gene?



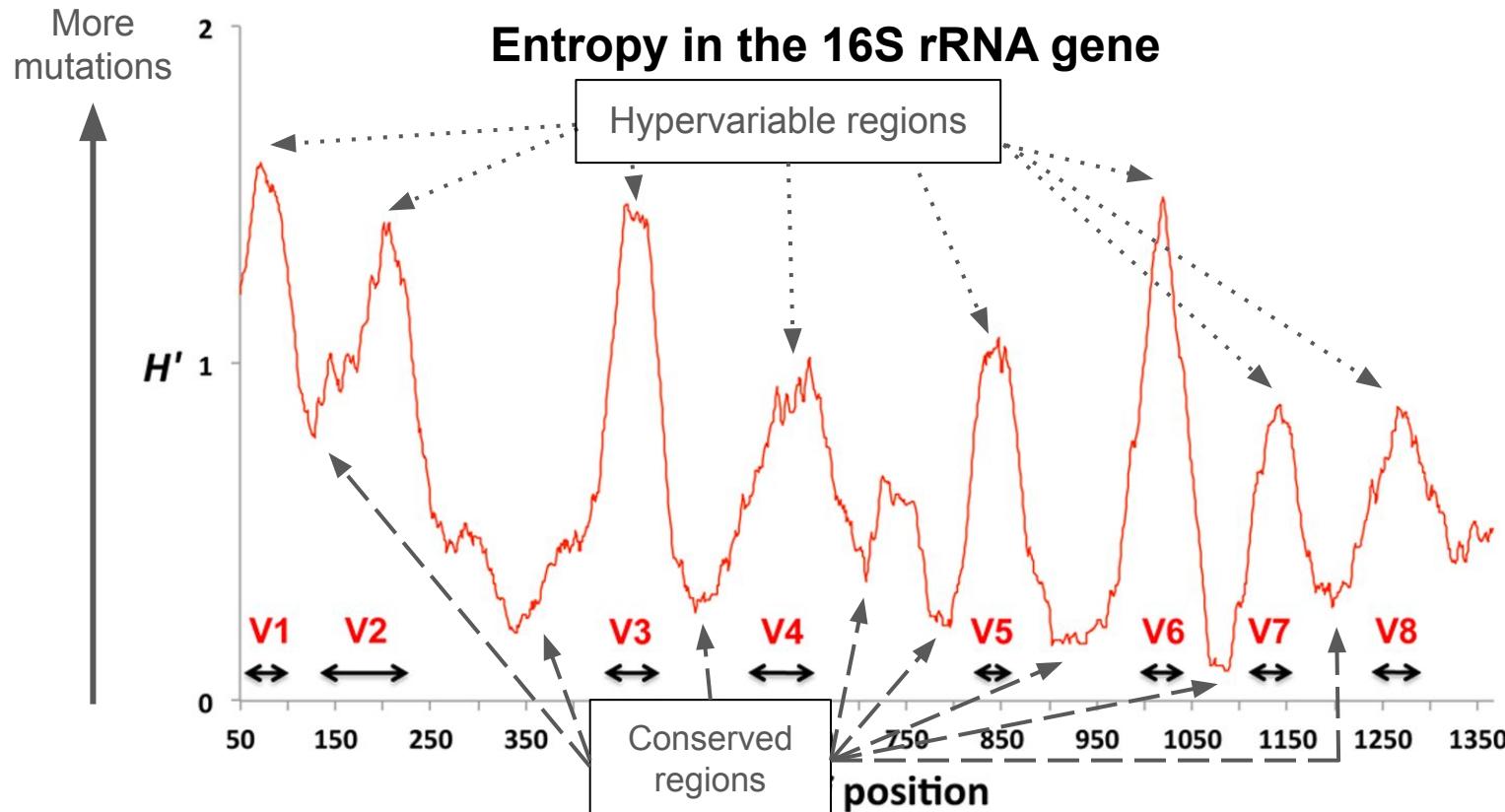
C. I. Methods: What's in a marker gene?



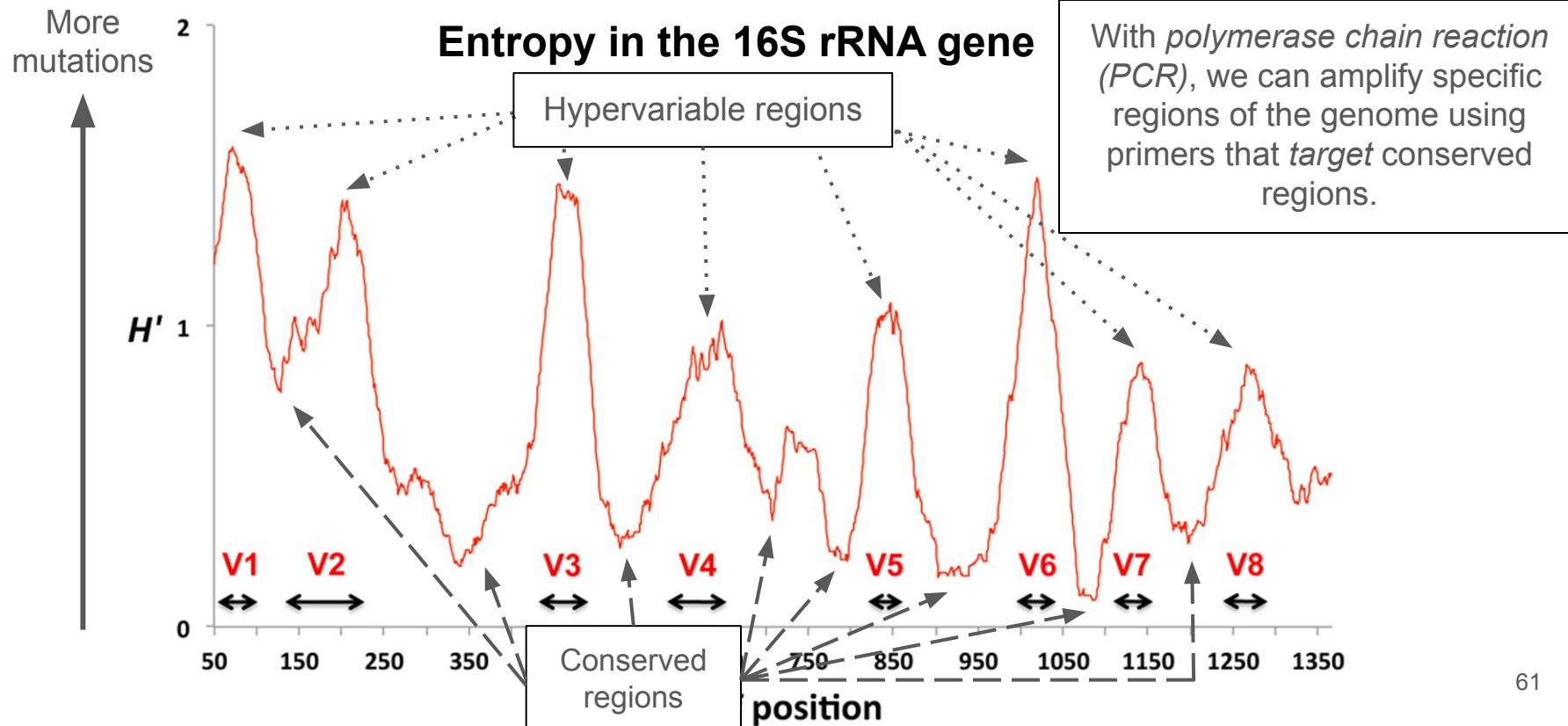
C. I. Methods: What's in a marker gene?



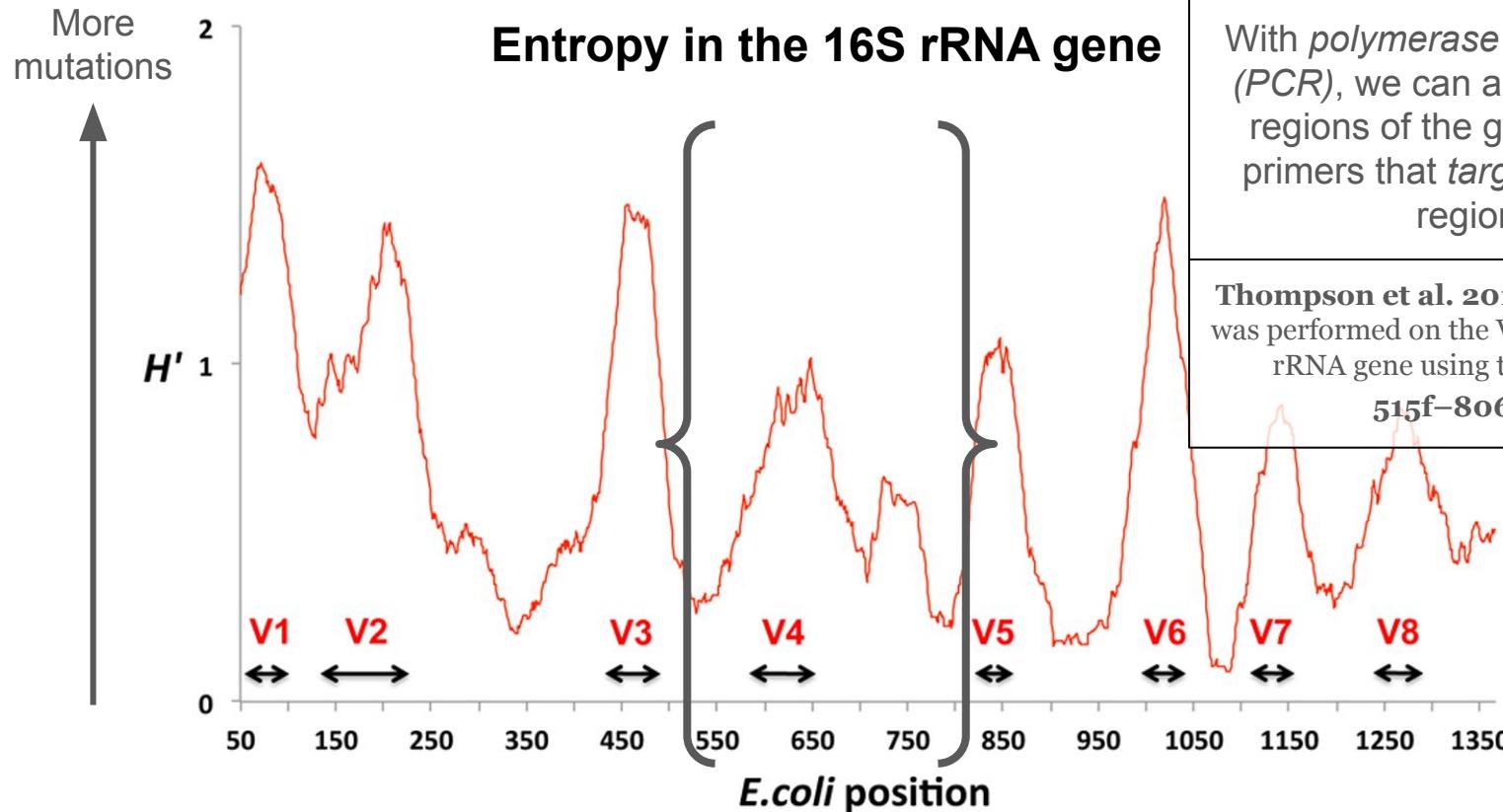
C. I. Methods: What's in a marker gene?



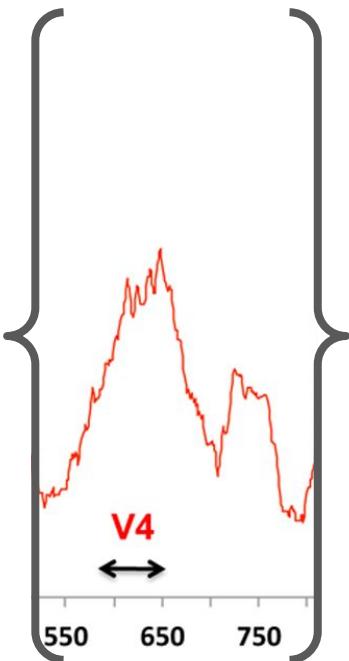
C. I. Methods: What's in a marker gene?



C. I. Methods: What's in a marker gene?



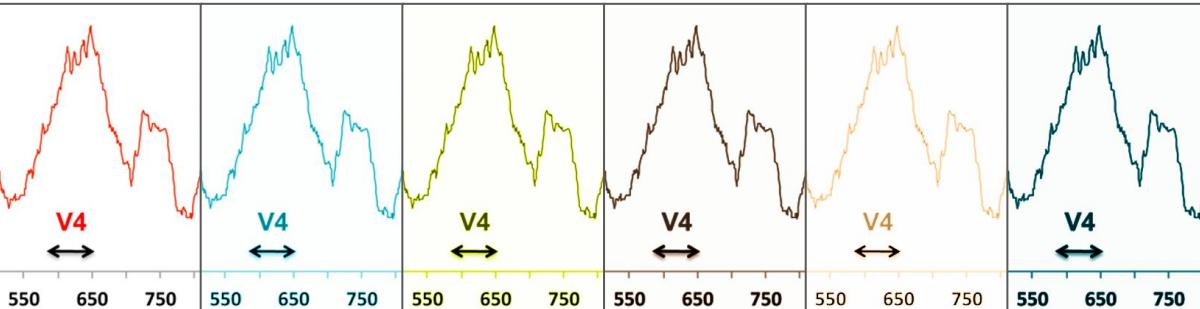
C. I. Methods: What's in a marker gene?



With *polymerase chain reaction* (PCR), we can amplify specific regions of the genome using primers that *target* conserved regions.

Thompson et al. 2017: “Amplicon PCR was performed on the V4 region of the 16S rRNA gene using the primer pair **515f–806r** [...]”

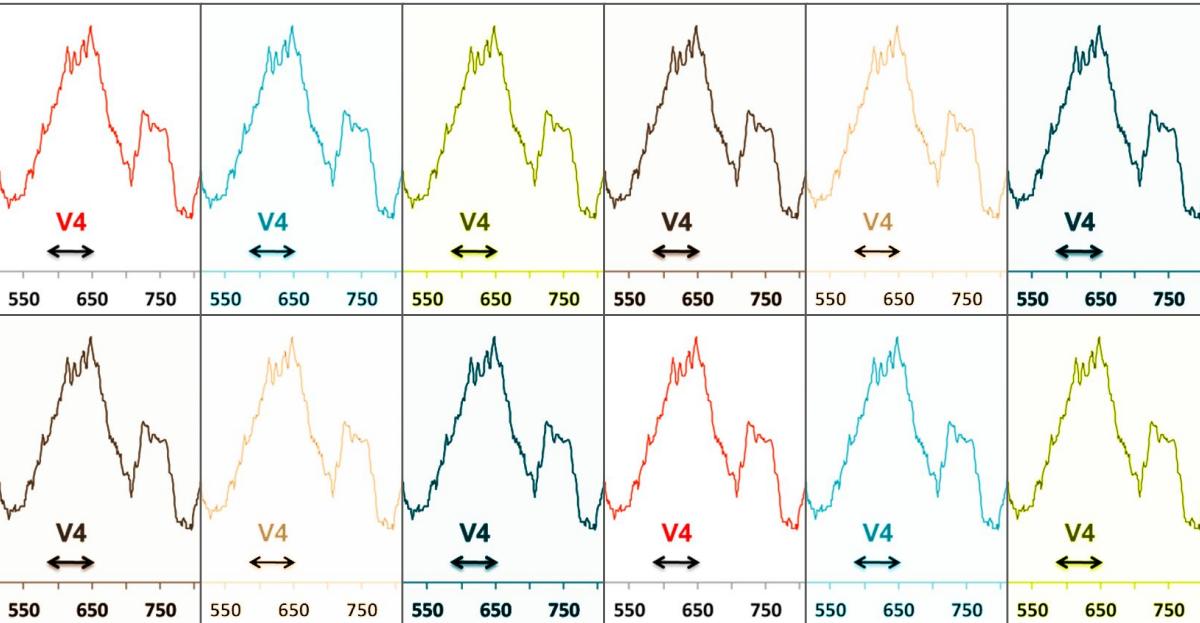
C. I. Methods: What's in a marker gene?



With *polymerase chain reaction* (PCR), we can amplify specific regions of the genome using primers that *target* conserved regions.

Thompson et al. 2017: “Amplicon PCR was performed on the V4 region of the 16S rRNA gene using the primer pair **515f–806r** [...]”

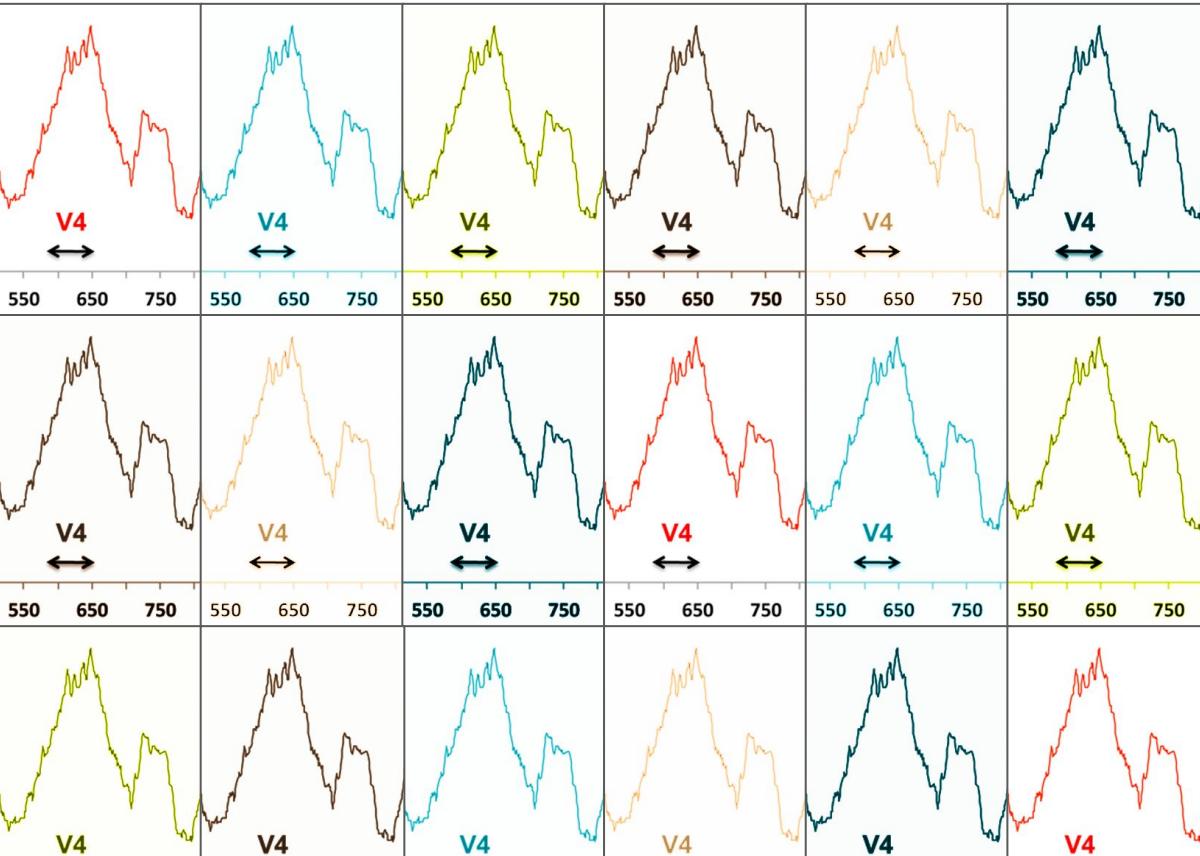
C. I. Methods: What's in a marker gene?



With *polymerase chain reaction* (PCR), we can amplify specific regions of the genome using primers that *target* conserved regions.

Thompson et al. 2017: “Amplicon PCR was performed on the V4 region of the 16S rRNA gene using the primer pair **515f–806r** [...]”

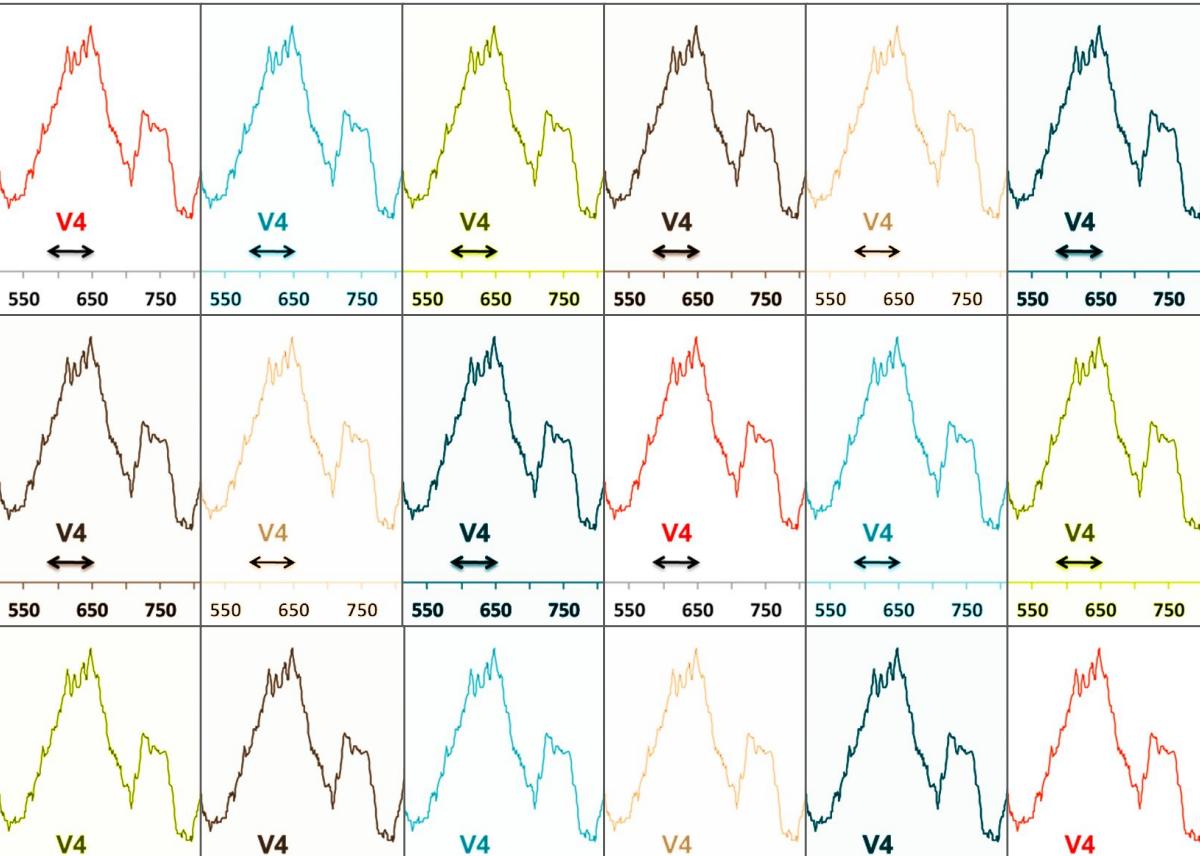
C. I. Methods: What's in a marker gene?



With *polymerase chain reaction* (PCR), we can amplify specific regions of the genome using primers that *target* conserved regions.

Thompson et al. 2017: “Amplicon PCR was performed on the V4 region of the 16S rRNA gene using the primer pair **515f–806r** [...]”

C. I. Methods: What's in a marker gene?

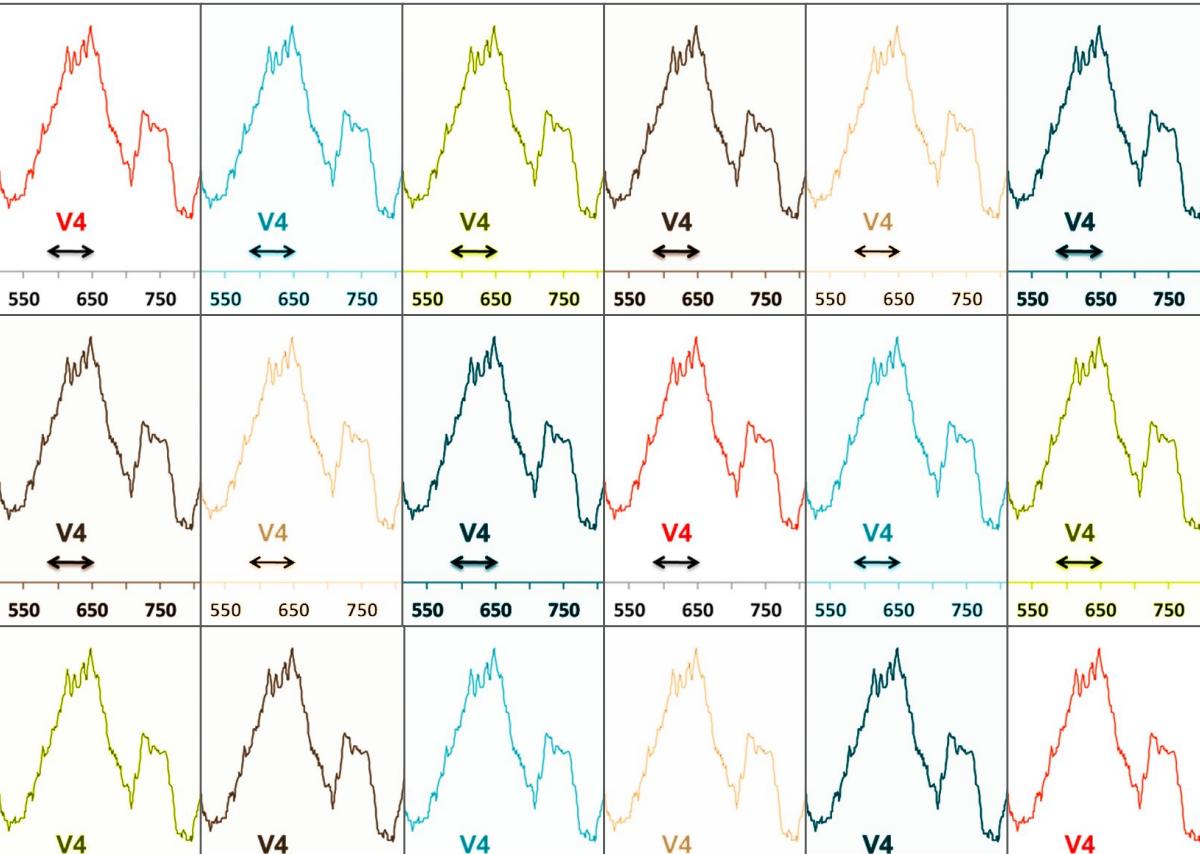


With *polymerase chain reaction* (PCR), we can amplify specific regions of the genome using primers that *target* conserved regions.

Thompson et al. 2017: “Amplicon PCR was performed on the V4 region of the 16S rRNA gene using the primer pair **515f–806r** [...]”

• • •

C. I. Methods: What's in a marker gene?



With *polymerase chain reaction* (PCR), we can amplify specific regions of the genome using primers that *target* conserved regions.

Thompson et al. 2017: “Amplicon PCR was performed on the V4 region of the 16S rRNA gene using the primer pair **515f–806r** [...]”

...

Since these amplified sequences contain hypervariable region(s), these regions help us determine which sequences came from which microbe.

C. I. Methods: What's in a marker gene?

ATCGACTGA CTGACGTAC TGTACGGAT ACCGGGGAC ATACTACTA CTACTACTA CTTTTCCCA	ATCGACTGA CTGACGTAC TGTACGGAT ACCGGGGAC ATACTACTA CTACTACTA CTTTTCCCA	ATCGACTGA CTGACGTAC TGTACGGAT ACCGGGGAC ATACTACTA CTACTACTA CTTTTCCCA	ATCGACTGA CTGACGTAC TGTACGGAT ACCGGGGAC ATACTACTA CTACTACTA CTTTTCCCA	ATCGACTGA CTGACGTAC TGTACGGAT ACCGGGGAC ATACTACTA CTACTACTA CTTTTCCCA	ATCGACTGA CTGACGTAC TGTACGGAT ACCGGGGAC ATACTACTA CTACTACTA CTTTTCCCA
---	---	---	---	---	---

ATCGACTGA CTGACGTAC TGTACGGAT ACCGGGGAC ATACTACTA CTACTACTA CTTTTCCCA	ATCGACTGA CTGACGTAC TGTACGGAT ACCGGGGAC ATACTACTA CTACTACTA CTTTTCCCA	ATCGACTGA CTGACGTAC TGTACGGAT ACCGGGGAC ATACTACTA CTACTACTA CTTTTCCCA	ATCGACTGA CTGACGTAC TGTACGGAT ACCGGGGAC ATACTACTA CTACTACTA CTTTTCCCA	ATCGACTGA CTGACGTAC TGTACGGAT ACCGGGGAC ATACTACTA CTACTACTA CTTTTCCCA	ATCGACTGA CTGACGTAC TGTACGGAT ACCGGGGAC ATACTACTA CTACTACTA CTTTTCCCA
---	---	---	---	---	---

ATCGACTGA CTGACGTAC TGTACGGAT ACCGGGGAC ATACTACTA CTACTACTA CTTTTCCCA	ATCGACTGA CTGACGTAC TGTACGGAT ACCGGGGAC ATACTACTA CTACTACTA CTTTTCCCA	ATCGACTGA CTGACGTAC TGTACGGAT ACCGGGGAC ATACTACTA CTACTACTA CTTTTCCCA	ATCGACTGA CTGACGTAC TGTACGGAT ACCGGGGAC ATACTACTA CTACTACTA CTTTTCCCA	ATCGACTGA CTGACGTAC TGTACGGAT ACCGGGGAC ATACTACTA CTACTACTA CTTTTCCCA	ATCGACTGA CTGACGTAC TGTACGGAT ACCGGGGAC ATACTACTA CTACTACTA CTTTTCCCA
---	---	---	---	---	---

With *polymerase chain reaction* (PCR), we can amplify specific regions of the genome using primers that *target* conserved regions.

Thompson et al. 2017: “Amplicon PCR was performed on the V4 region of the 16S rRNA gene using the primer pair **515f–806r** [...]”

...

Since these amplified sequences contain hypervariable region(s), these regions help us determine which sequences came from which microbe.

C. I. Methods: What's in a marker gene?

70

k_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae; g_Bacteroides; s_	k_Bacteria; p_Proteobacteria; c_Betaproteobacteria; o_Neisseriales; f_Neisseriaceae; g_Neisseria	k_Bacteria; p_Firmicutes; c_Bacilli; o_Lactobacillales; f_Streptococcaceae; g_Streptococcus	k_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae; g_Bacteroides; s_	k_Bacteria; p_Firmicutes; c_Bacilli	k_Bacteria; p_Proteobacteria; c_Gammaproteobacteria; o_Pasteurellales; f_Pasteurellaceae; g_Haemophilus; s_parainfluenzae
k_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae; g_Bacteroides; s_	k_Bacteria; p_Firmicutes; c_Bacilli	k_Bacteria; p_Proteobacteria; c_Gammaproteobacteria; o_Pasteurellales; f_Pasteurellaceae; g_Haemophilus; s_parainfluenzae	k_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae; g_Bacteroides; s_	k_Bacteria; p_Proteobacteria; c_Bacilli; o_Lactobacillales; f_Streptococcaceae; g_Streptococcus	k_Bacteria; p_Firmicutes; c_Bacilli; o_Lactobacillales; f_Streptococcaceae; g_Streptococcus
k_Bacteria; p_Firmicutes; c_Bacilli; o_Lactobacillales; f_Streptococcus	k_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae	k_Bacteria; p_Proteobacteria; c_Betaproteobacteria; o_Neisseriales; f_Neisseriaceae	k_Bacteria; p_Firmicutes; c_Bacilli	k_Bacteria; p_Proteobacteria; c_Gammaproteobacteria; o_Pasteurellales; f_Pasteurellaceae; g_Haemophilus	k_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae

...

With *polymerase chain reaction* (PCR), we can amplify specific regions of the genome using primers that target conserved regions.

Thompson et al. 2017: “Amplicon PCR was performed on the V4 region of the 16S rRNA gene using the primer pair **515f–806r** [...]”

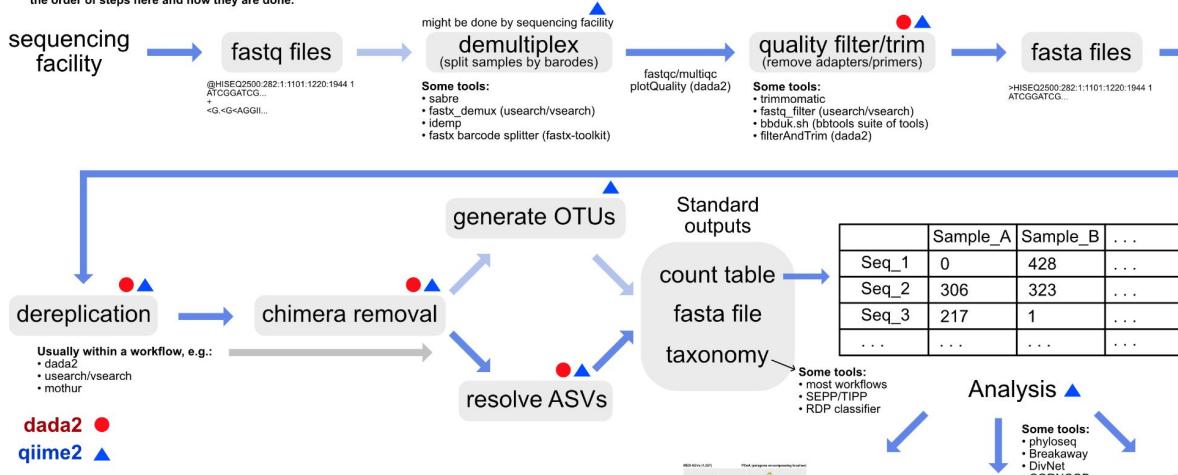
Since these amplified sequences contain hypervariable region(s), these regions help us determine which sequences came from which microbe.

Example taxonomic annotations from the QIIME 2 “Moving Pictures” tutorial: <https://docs.qiime2.org>

C. I. Methods: Marker gene sequencing

Overview of generic* amplicon workflow

*This is generic; specific workflows can vary on the order of steps here and how they are done.



Some tools that provide whole workflows:

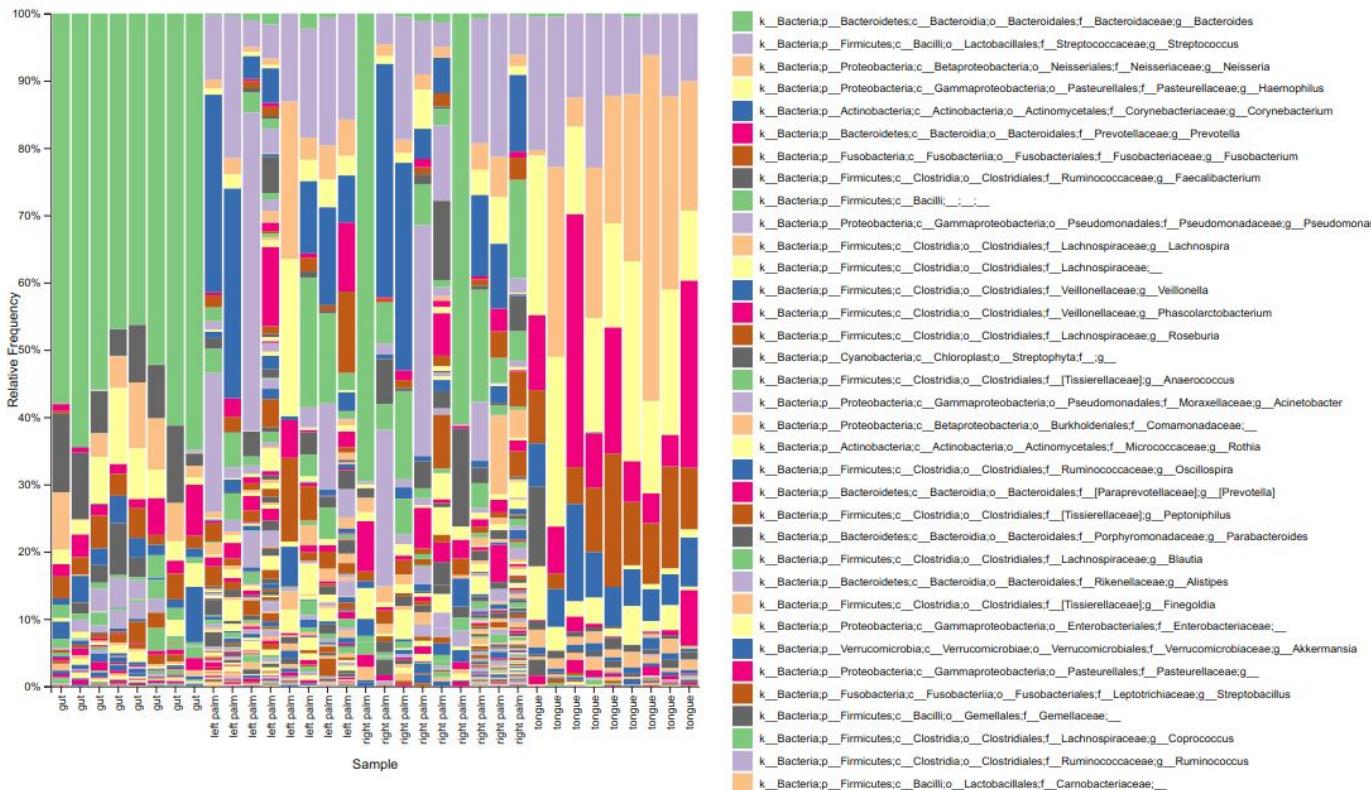
dada2 runs within R (ASVs)

usearch/vsearch runs at the command line (ASVs and OTUs)
mothur runs at the command line (OTUs only currently)

qiime2 provides a multi-interface environment that employs processing tools like those above, infrastructure for easily documenting all processing performed, and interactive visualizations

C. I. Methods: Marker gene sequencing

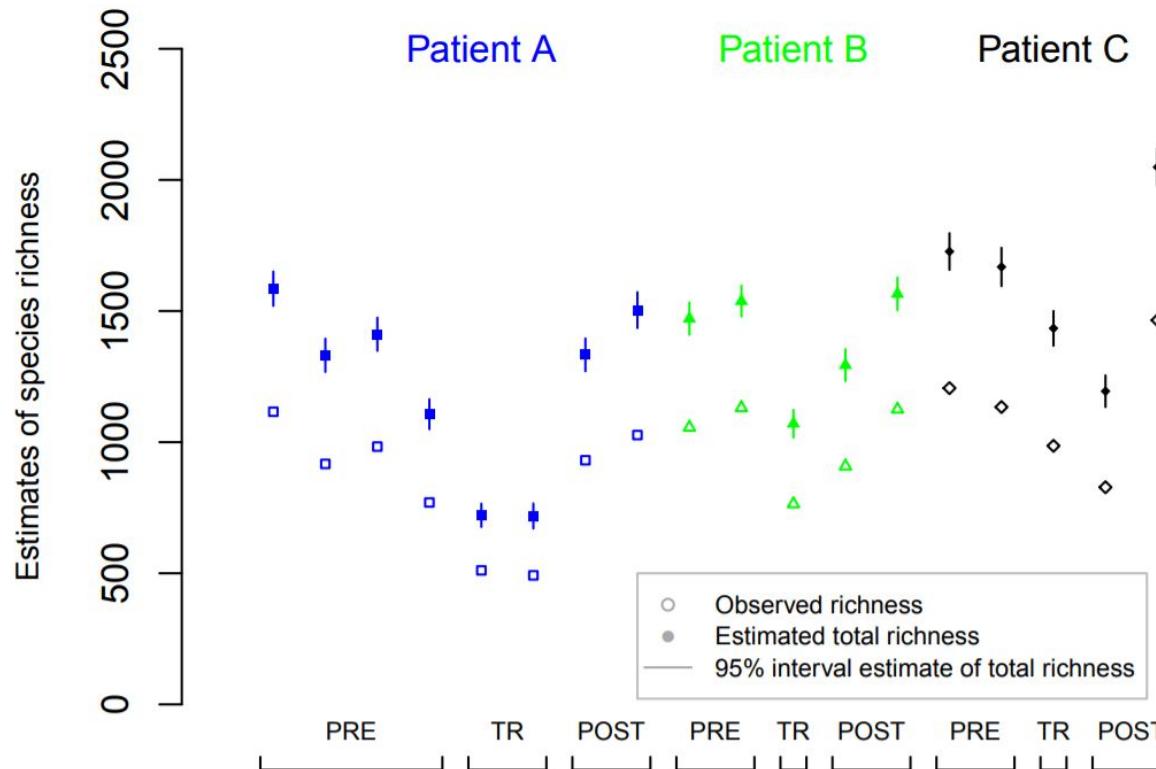
Broad
taxonomic
comparisons



C. I. Methods: Marker gene sequencing

Estimating and comparing diversity within samples

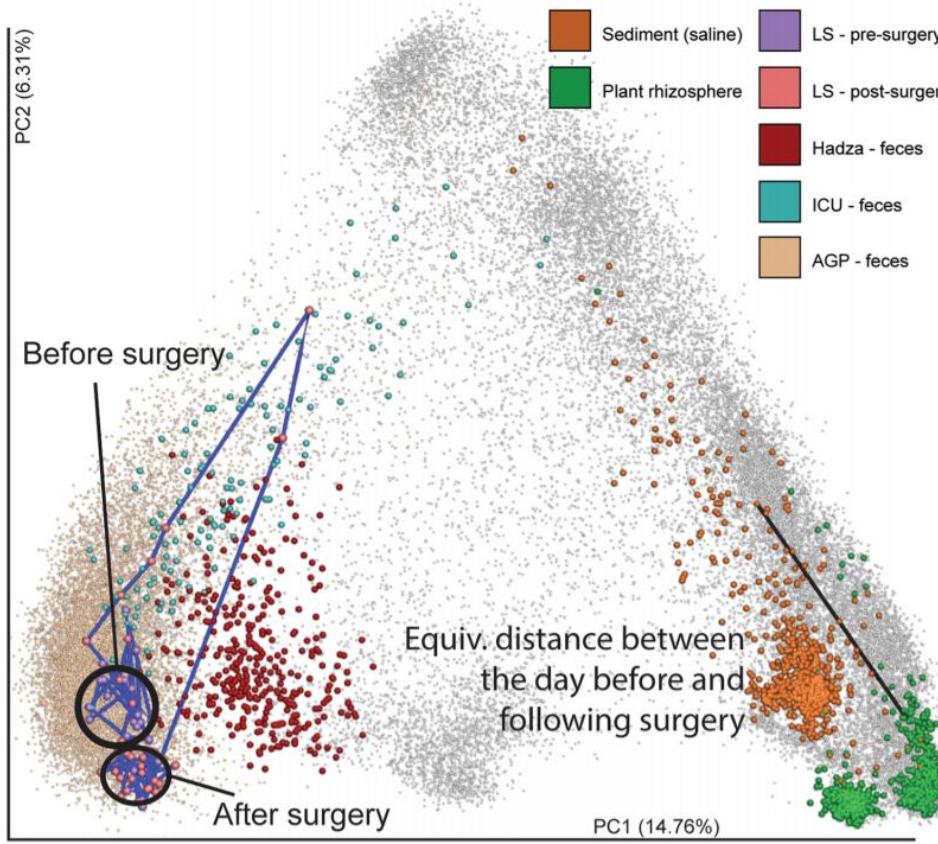
Also referred to as α -diversity ("alpha diversity")



C. I. Methods: Marker gene sequencing

Unsupervised dimensionality reduction (e.g. PCA / PCoA)

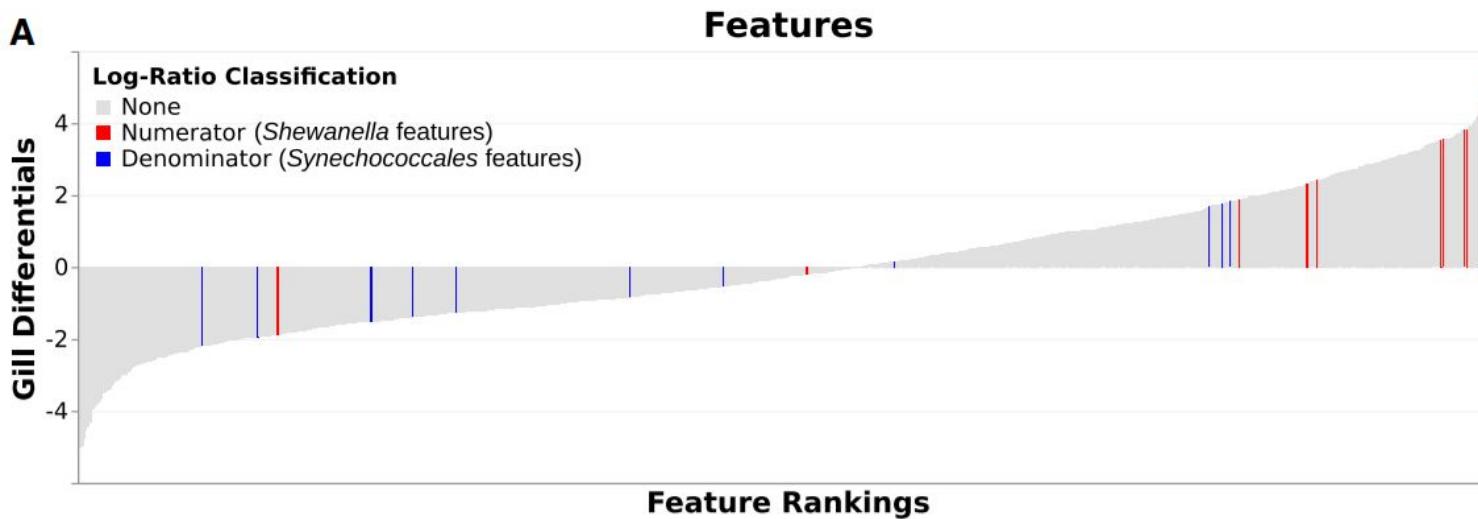
Also referred to as β -diversity ("beta diversity")



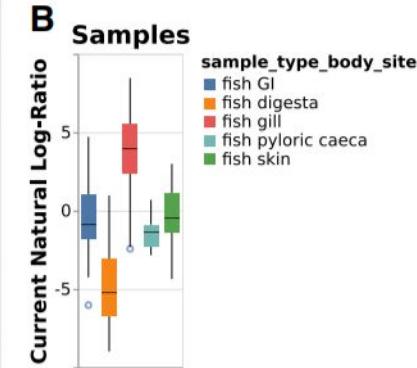
C. I. Methods: Marker gene sequencing

Identifying
differentially abundant
features (or ratios
of features) in
groups of samples

A



B Samples

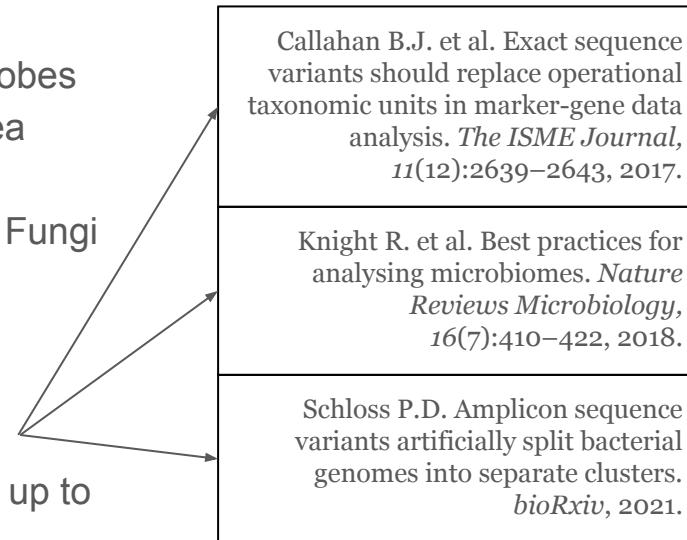


C. I. Methods: Marker gene sequencing

- **Pros**
 - Relatively cheap: much less sequencing needed to profile a community at a given depth than metagenome sequencing
 - Very well-studied, so many established pipelines have been created (mothur, QIIME / QIIME 2, ...)
- **Cons**
 - Marker genes are usually specific to certain types of microbes
 - 16S rRNA gene: only identifies Bacteria and Archaea
 - 18S rRNA gene: only identifies certain Eukaryotes
 - Internal transcribed spacer region: mostly identifies Fungi
 - ...
 - Copy number variation
 - PCR errors can result in “chimeric” gene sequences
 - Still disagreements on how to correct errors in raw reads
 - Limited resolution: marker genes are usually only reliable up to the genus level

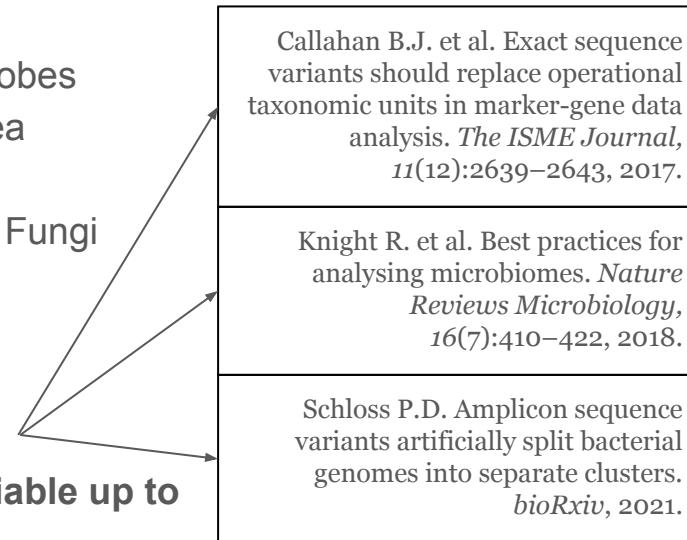
C. I. Methods: Marker gene sequencing

- **Pros**
 - Relatively cheap: much less sequencing needed to profile a community at a given depth than metagenome sequencing
 - Very well-studied, so many established pipelines have been created (mothur, QIIME / QIIME 2, ...)
- **Cons**
 - Marker genes are usually specific to certain types of microbes
 - 16S rRNA gene: only identifies Bacteria and Archaea
 - 18S rRNA gene: only identifies certain Eukaryotes
 - Internal transcribed spacer region: mostly identifies Fungi
 - ...
 - Copy number variation
 - PCR errors can result in “chimeric” gene sequences
 - Still disagreements on how to correct errors in raw reads
 - Limited resolution: marker genes are usually only reliable up to the genus level



C. I. Methods: Marker gene sequencing

- **Pros**
 - Relatively cheap: much less sequencing needed to profile a community at a given depth than metagenome sequencing
 - Very well-studied, so many established pipelines have been created (mothur, QIIME / QIIME 2, ...)
- **Cons**
 - Marker genes are usually specific to certain types of microbes
 - 16S rRNA gene: only identifies Bacteria and Archaea
 - 18S rRNA gene: only identifies certain Eukaryotes
 - Internal transcribed spacer region: mostly identifies Fungi
 - ...
 - Copy number variation
 - PCR errors can result in “chimeric” gene sequences
 - Still disagreements on how to correct errors in raw reads
 - **Limited resolution: marker genes are usually only reliable up to the genus level**

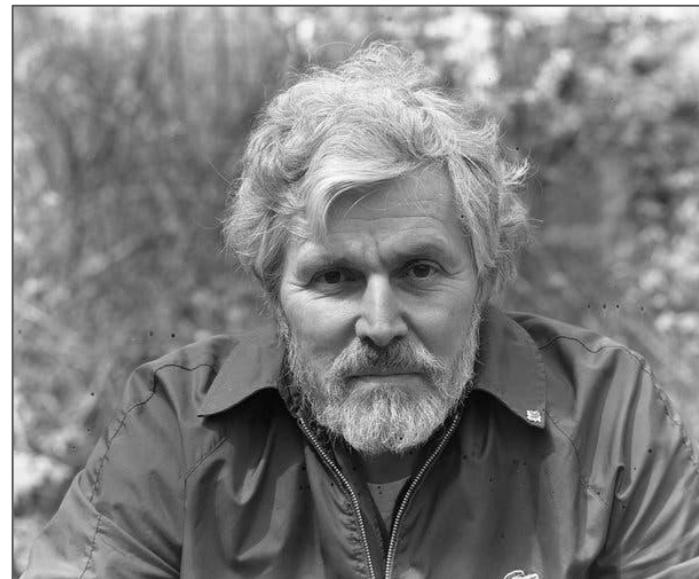


For a more thorough history of Carl Woese's life and work, see:
D. Quammen. The Scientist Who Scrambled Darwin's Tree of Life. *The New York Times*, 2018.

C. I. Methods: back to Carl Woese and rRNA genes

The original paper on PCR was published in **1986**.

The first automatic sequencer (the “AB370”) was developed in **1987**.



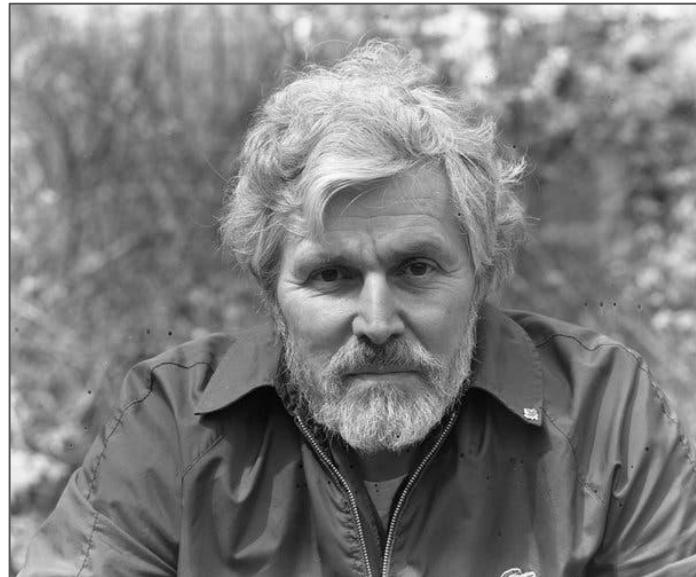
For a more thorough history of Carl Woese's life and work, see:
D. Quammen. The Scientist Who Scrambled Darwin's Tree of Life. *The New York Times*, 2018.

C. I. Methods: back to Carl Woese and rRNA genes

The original paper on PCR was published in **1986**.

The first automatic sequencer (the “AB370”) was developed in **1987**.

Woese and Fox’s three-domain paper was published in **1977**!



For a more thorough history of Carl Woese's life and work, see:
D. Quammen. The Scientist Who Scrambled Darwin's Tree of Life. *The New York Times*, 2018.

C. I. Methods: back to Carl Woese and rRNA genes

The original paper on PCR was published in **1986**.

The first automatic sequencer (the “AB370”) was developed in **1987**.

Woese and Fox’s three-domain paper was published in **1977**!

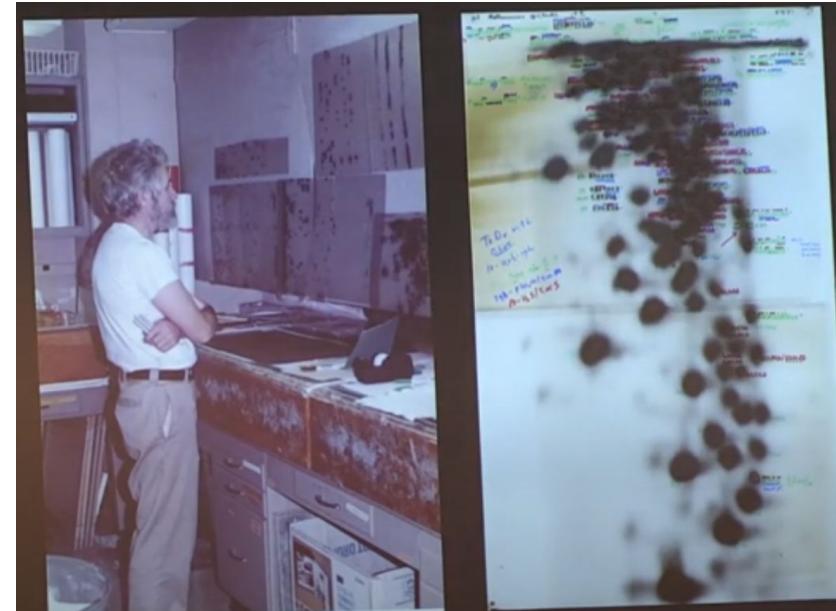


Photo credit: N. Pace. Carl Woese and the Beginnings of Metagenomics.
Looking in the Right Direction: Carl Woese and the New Biology, 2015.
<https://www.youtube.com/watch?v=h3K5oDD9kIM> (timestamp: 2:56)

C. I. Methods: back to Carl Woese and rRNA genes

The original paper on PCR was published in **1986**.

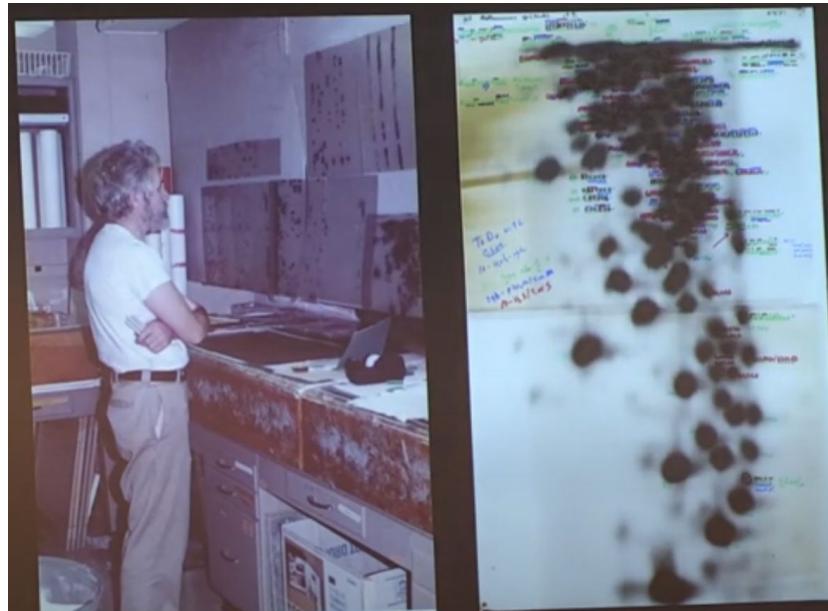
The first automatic sequencer (the “AB370”) was developed in **1987**.

Woese and Fox’s three-domain paper was published in **1977**!

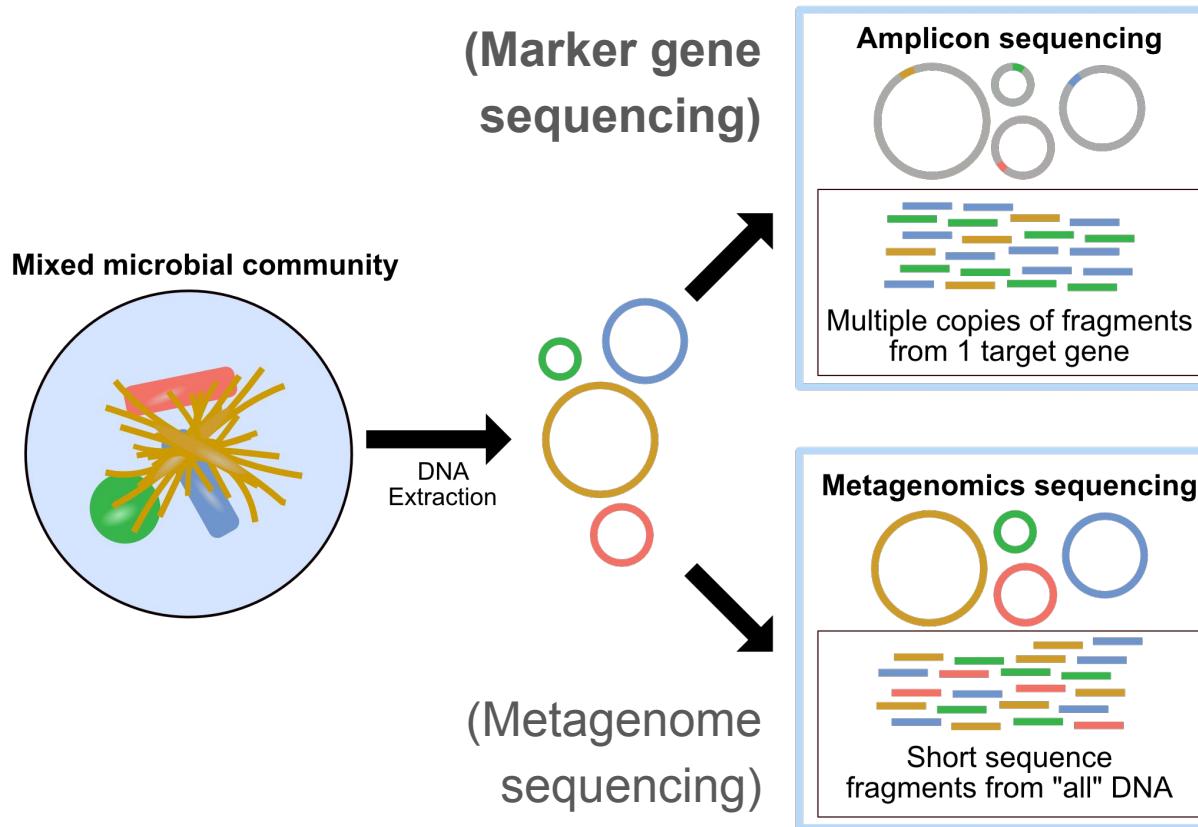
“While the grad students and technicians produced fingerprints, Woese spent his time staring at the spots. Was this effort tedious in practice as well as profound in its potential results? Yes.

‘There were days,’ he wrote later, ‘when I’d walk home from work saying to myself, ‘Woese, you have destroyed your mind again today.’ ’ ”

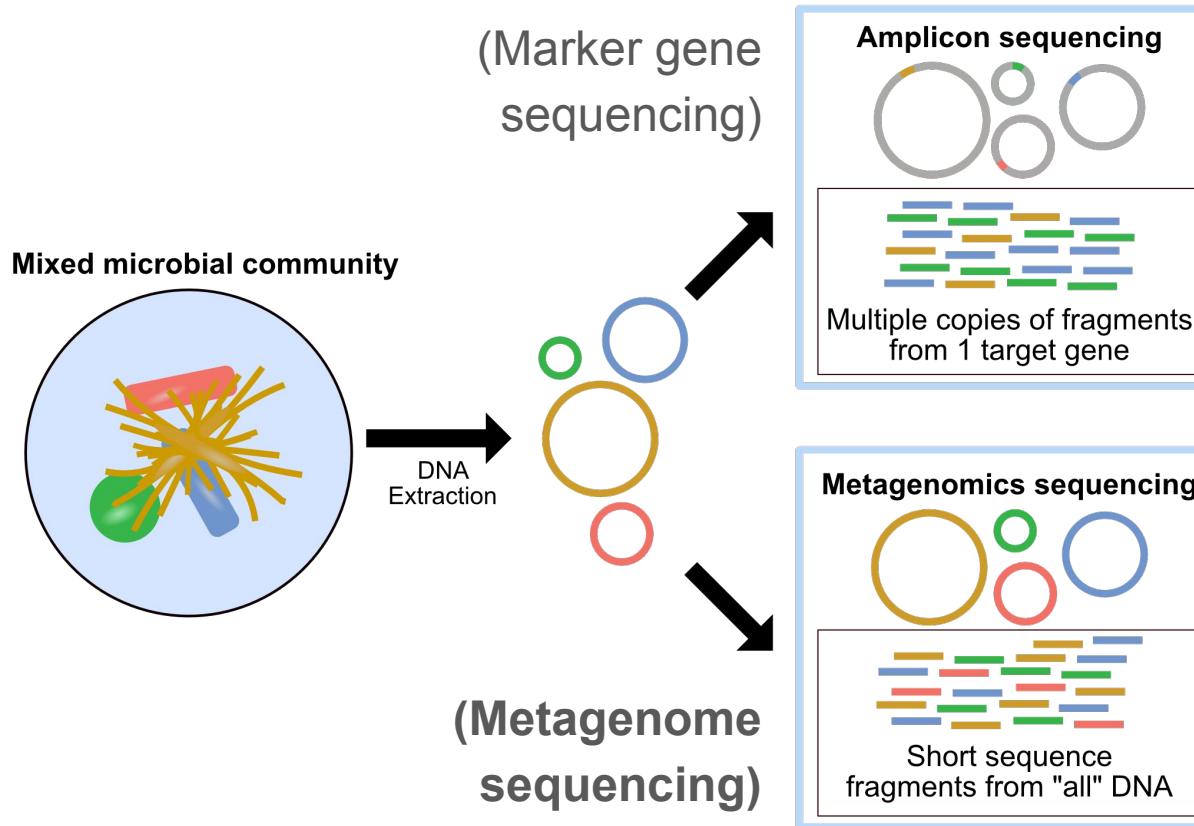
Photo credit: N. Pace. Carl Woese and the Beginnings of Metagenomics.
Looking in the Right Direction: Carl Woese and the New Biology, 2015.
<https://www.youtube.com/watch?v=h3K5oDD9kIM> (timestamp: 2:56)



C. I. Methods: Marker gene sequencing



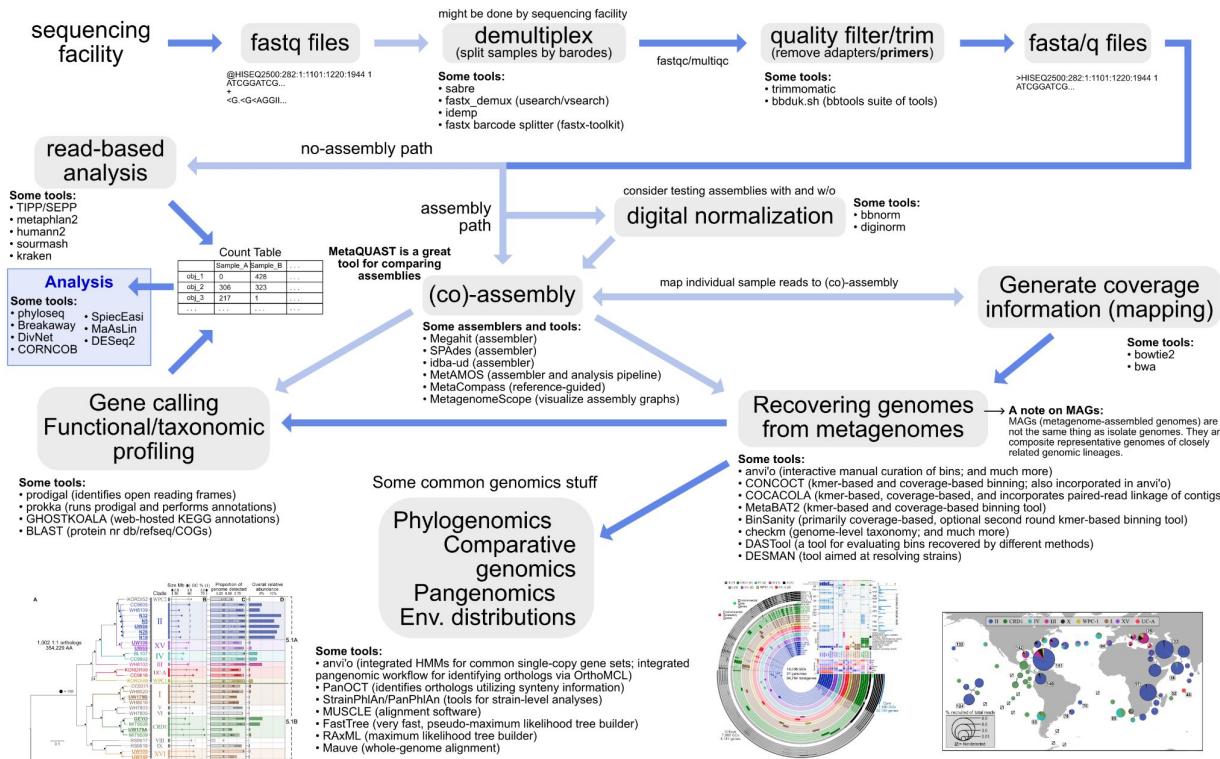
C. I. Methods: Metagenome sequencing



C. I. Methods: Metagenome sequencing

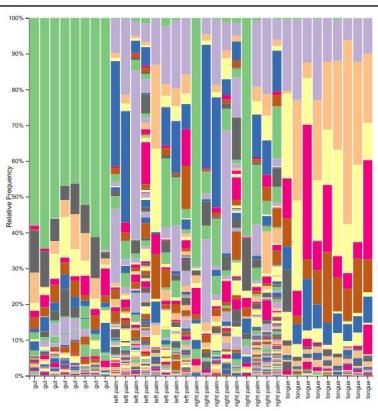
Overview of generic* metagenomics workflow

*This is generic; specific workflows can vary on the order of steps here and how they are done.

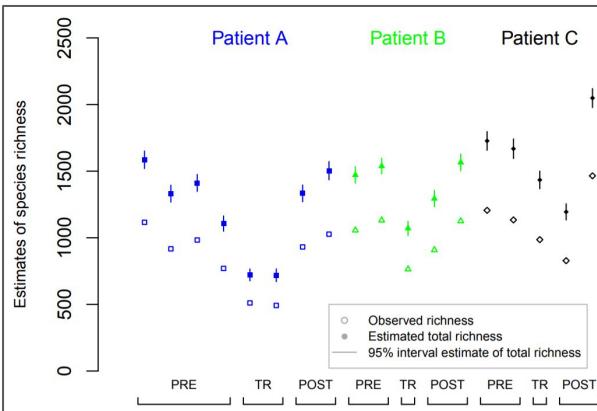


C. I. Methods: Metagenome sequencing

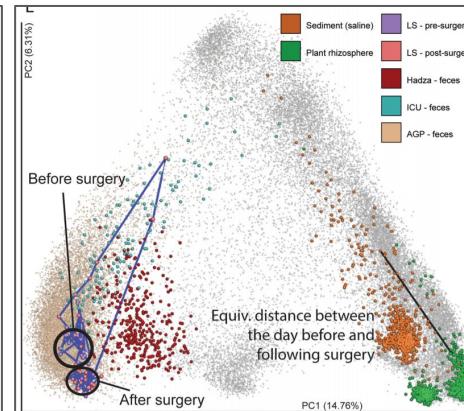
Taxonomy



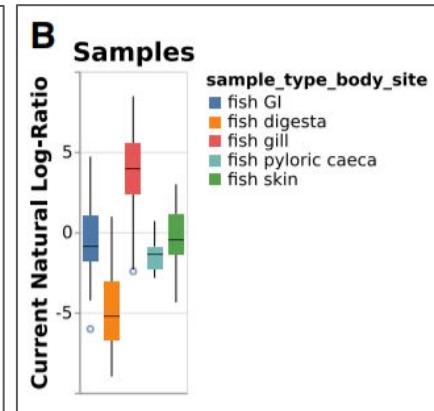
α -diversity



β -diversity



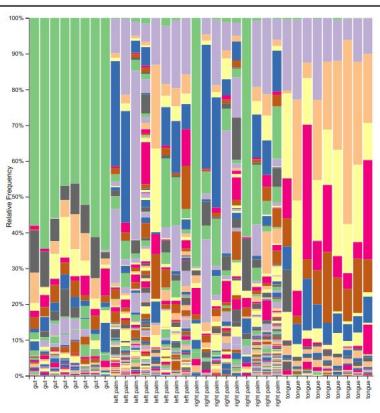
Differential abundance



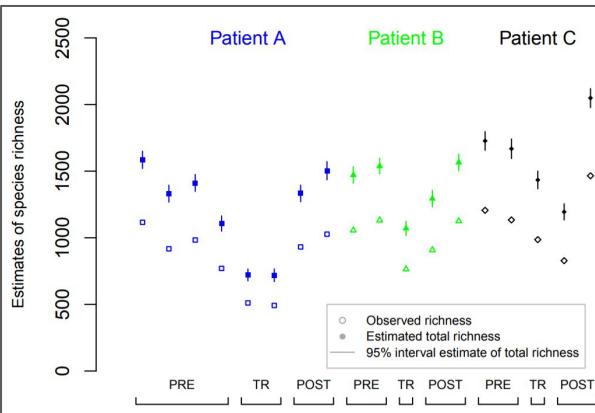
Many of the “standard analyses” for marker gene sequencing data are also applicable to metagenome sequencing data.

C. I. Methods: Metagenome sequencing

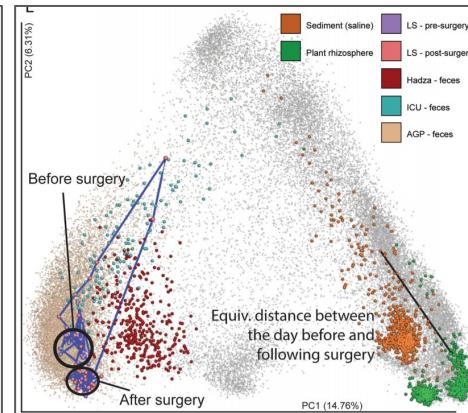
Taxonomy



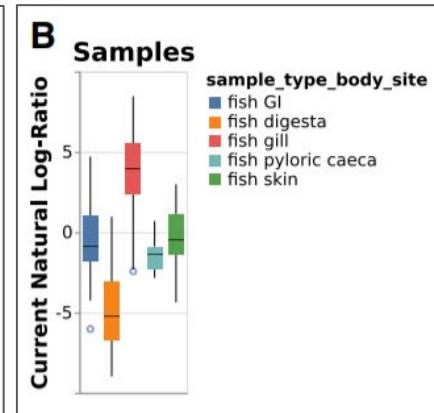
α -diversity



β -diversity



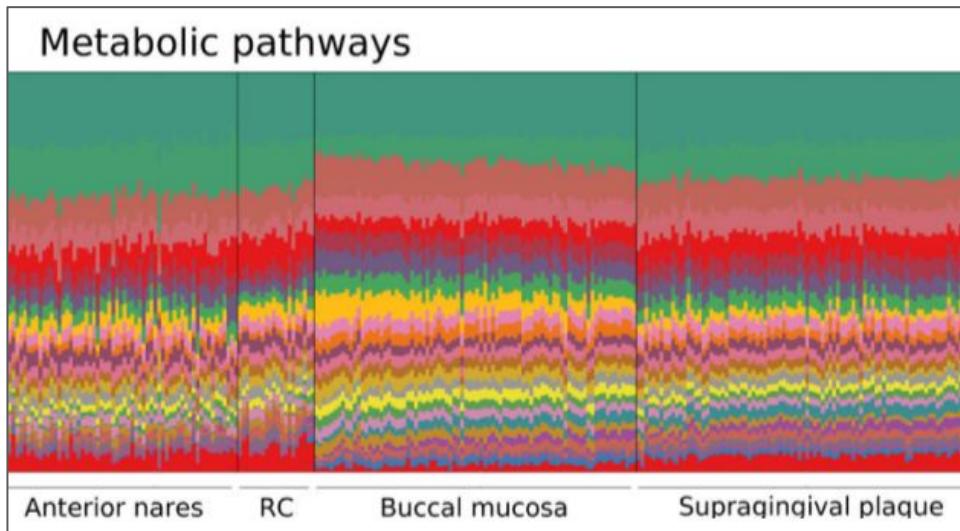
Differential abundance



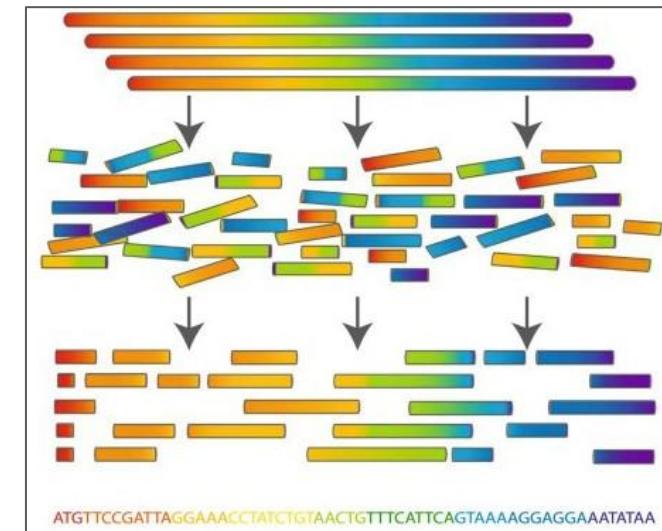
Metagenome sequencing enables two main additional types of analyses, compared to marker gene sequencing.

C. I. Methods: Metagenome sequencing

Functional annotation



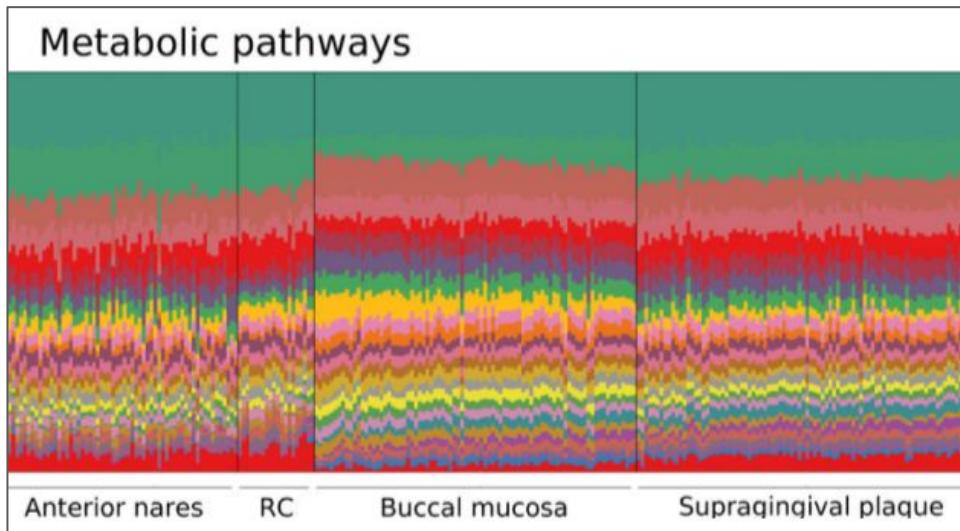
Sequence assembly



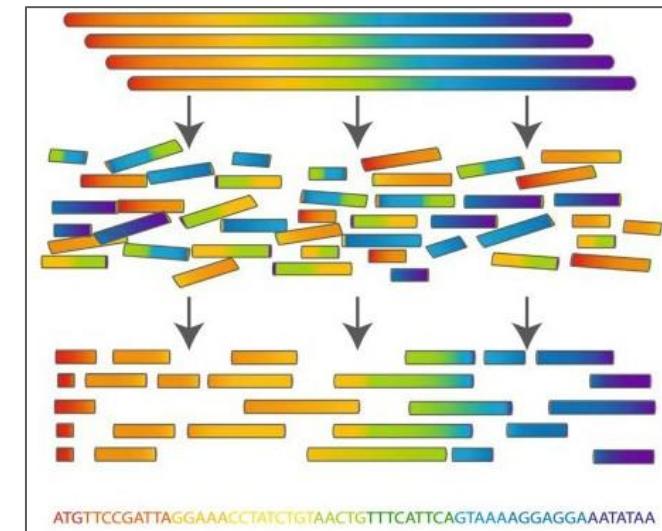
Metagenome sequencing enables two main additional types of analyses, compared to marker gene sequencing.

C. I. Methods: Metagenome sequencing

Functional annotation

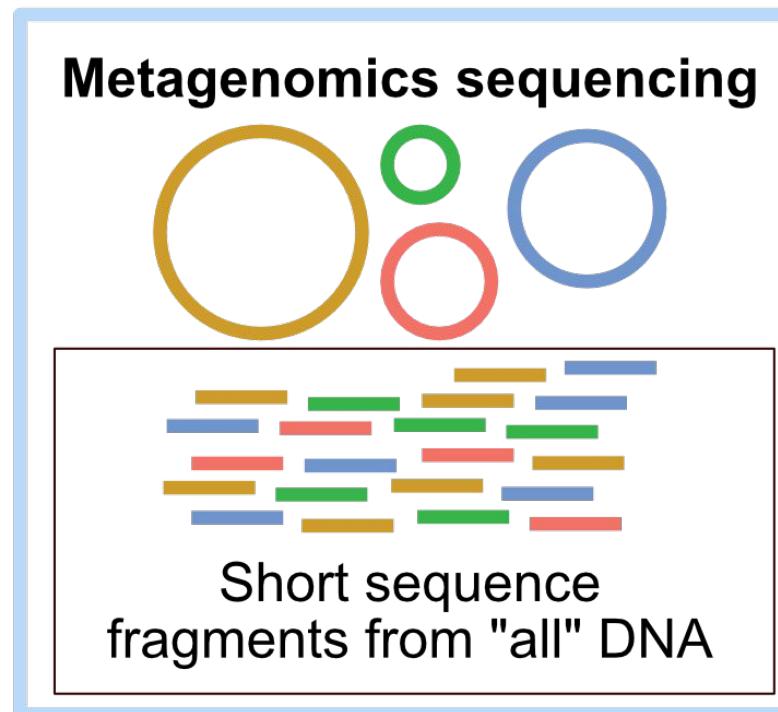


Sequence assembly

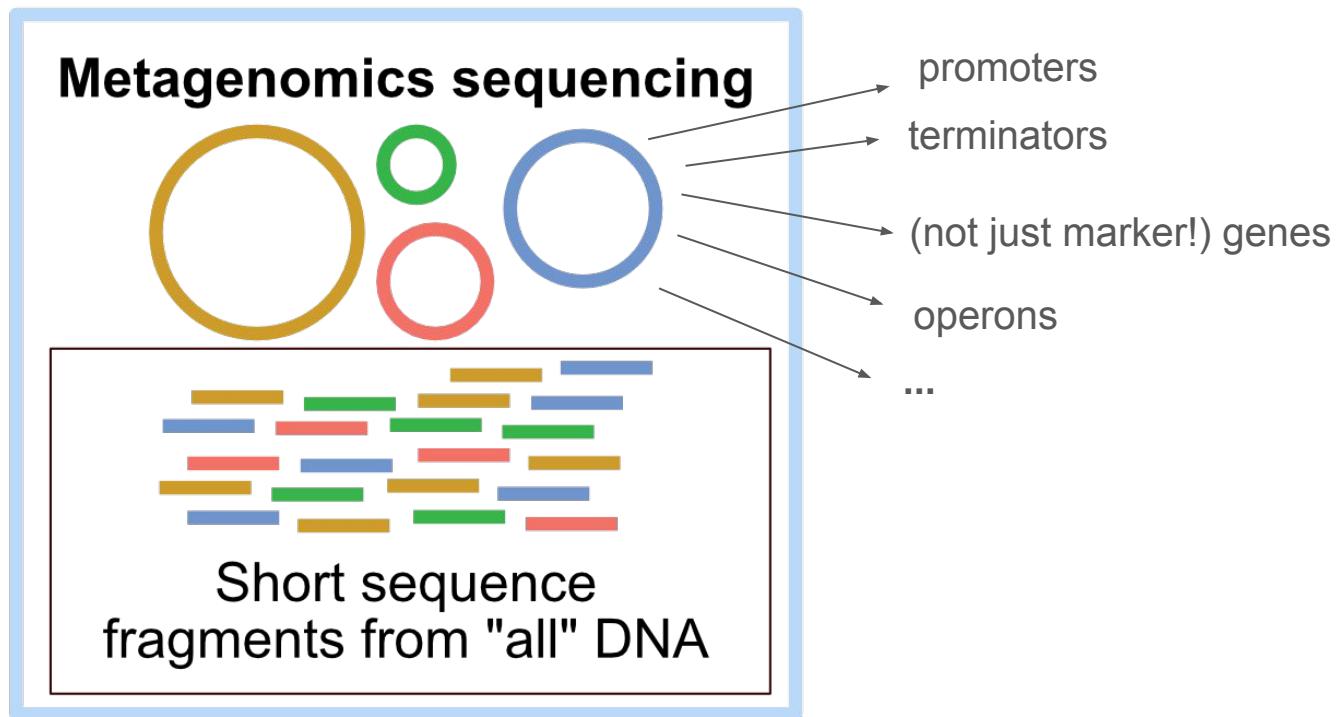


Metagenome sequencing enables two main additional types of analyses, compared to marker gene sequencing.

C. I. Methods: Functional annotation



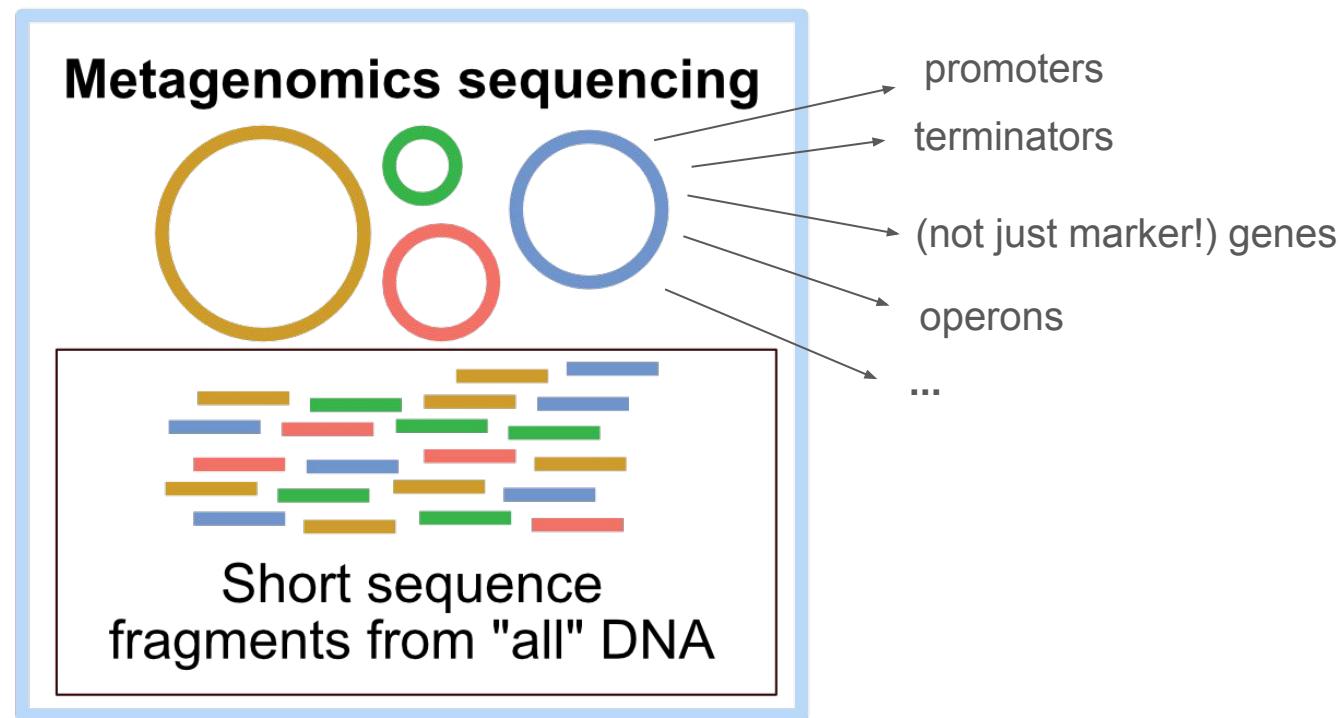
C. I. Methods: Functional annotation



C. I. Methods: Functional annotation, in practice

Usually, F.A. involves aligning (“mapping”) sequences to a reference database with information about “function” in well-studied organisms.

(But there are fancier approaches, e.g. metabolic modelling methods.)



C. I. Methods: Functional annotation, in practice?

An Integrated Encyclopedia of DNA Elements in the Human Genome

The ENCODE Project Consortium

“The Encyclopedia of DNA Elements (ENCODE) project has systematically mapped regions of transcription, transcription factor association, chromatin structure, and histone modification. These data enabled us to **assign biochemical functions for 80% of the [human] genome**, in particular outside of the well-studied protein-coding regions.”

C. I. Methods: Functional annotation, in practice??

An Integrated Encyclopedia of DNA Elements in the Human Genome

The ENCODE Project Consortium

“The Encyclopedia of DNA Elements (ENCODE) project has systematically mapped regions of transcription, transcription factor association, chromatin structure, and histone modification. These data enabled us to **assign biochemical functions for 80% of the [human] genome**, in particular outside of the well-studied protein-coding regions.”

On the Immortality of Television Sets: “Function” in the Human Genome According to the Evolution-Free Gospel of ENCODE

Dan Graur^{1,*}, Yichen Zheng¹, Nicholas Price¹, Ricardo B.R. Azevedo¹, Rebecca A. Zufall¹, and Eran Elhaik²

A recent slew of ENCYclopedia Of DNA Elements (ENCODE) Consortium publications, specifically the article signed by all Consortium members, put forward the idea that more than 80% of the human genome is functional. **This claim flies in the face of current estimates [...]**

C. I. Methods: Functional annotation, in practice???

An Integrated Encyclopedia of DNA Elements in the Human Genome

The ENCODE Project Consortium

“The Encyclopedia of DNA Elements (ENCODE) project has systematically mapped regions of transcription, transcription factor association, chromatin structure, and histone modification. These data enabled us to **assign biochemical functions for 80% of the [human] genome**, in particular outside of the well-studied protein-coding regions.”

On the Immortality of Television Sets: “Function” in the Human Genome According to the Evolution-Free Gospel of ENCODE

Dan Graur^{1,*}, Yichen Zheng¹, Nicholas Price¹, Ricardo B.R. Azevedo¹, Rebecca A. Zufall¹, and Eran Elhaik²

A recent slew of ENCyclopedia Of DNA Elements (ENCODE) Consortium publications, specifically the article signed by all Consortium members, put forward the idea that more than 80% of the human genome is functional. **This claim flies in the face of current estimates [...]**

Is junk DNA bunk? A critique of ENCODE

W. Ford Doolittle¹

Junk or functional DNA? ENCODE and the function controversy

Pierre-Luc Germain · Emanuele Ratti · Federico Boem

The ENCODE project: Missteps overshadowing a success

A farewell to bioinformatics

by Fred Ross

Last updated: March 26, 2012

genome being functionally annotated. It says something of our ability to annotate genomes, that the proportion of a genome functionally annotated is often correlated to the genetic distance to the very well researched *Escherichia coli* (anecdotal observation). How-

C. I. Methods: Functional annotation, in practice???

An Integrated Encyclopedia of DNA Elements in the Human Genome

The ENCODE Project Consortium

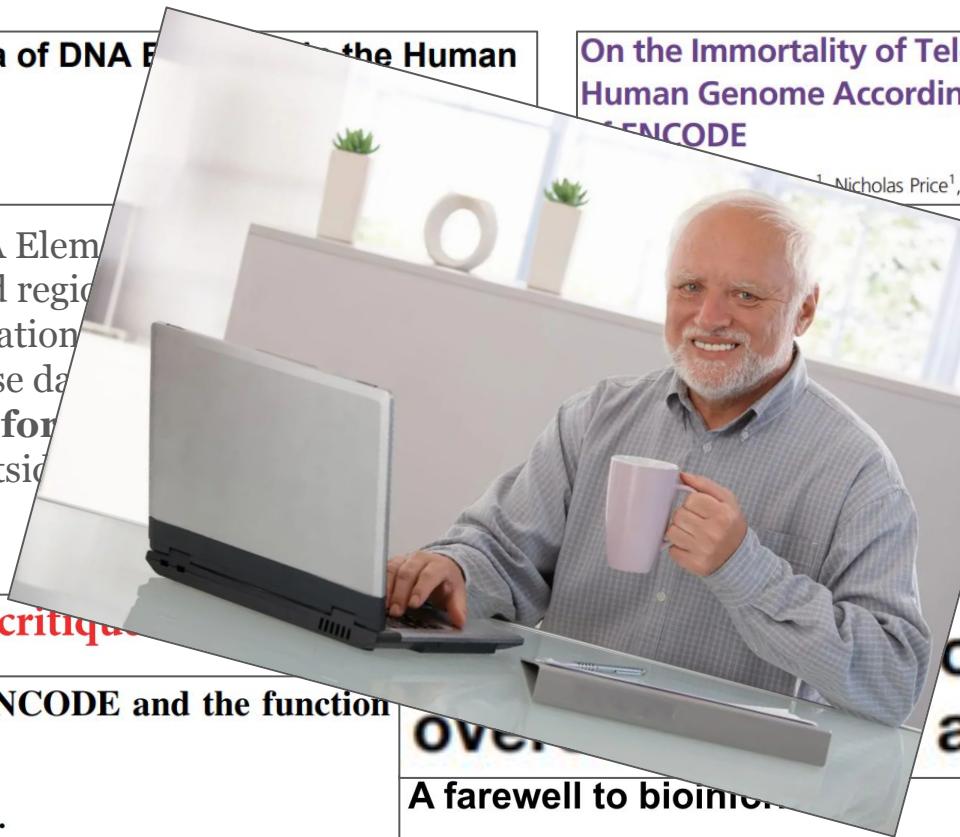
"The Encyclopedia of DNA Elements has systematically mapped regions of transcription factor association and histone modification. These data define the **biochemical functions for the genome**, in particular outside protein-coding regions."

Is junk DNA bunk? A critique

W. Ford Doolittle¹

Junk or functional DNA? ENCODE and the function controversy

Pierre-Luc Germain · Emanuele Ratti · Federico Boem



On the Immortality of Television Sets: "Function" in the Human Genome According to the Evolution-Free Gospel of ENCODE

¹ Nicholas Price¹, Ricardo B.R. Azevedo¹, Rebecca A. Zufall¹, and Eran Elhaik²

dia Of DNA Elements publications, specifically nsortium members, put e than 80% of the human is claim flies in the face [...]

Project: Missteps a success

A farewell to bioinformatics

by Fred Ross

Last updated: March 26, 2012

genome being functionally annotated. It says something of our ability to annotate genomes, that the proportion of a genome functionally annotated is often correlated to the genetic distance to the very well researched *Escherichia coli* (anecdotal observation). How-

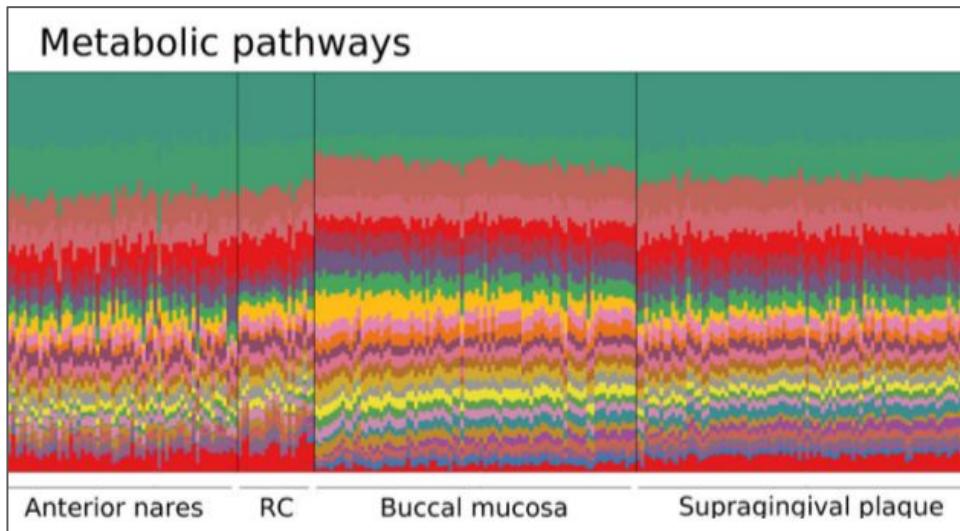
C. I. Methods: Functional annotation, in practice

Table 2. The information contained in different lengths of genomic DNA.

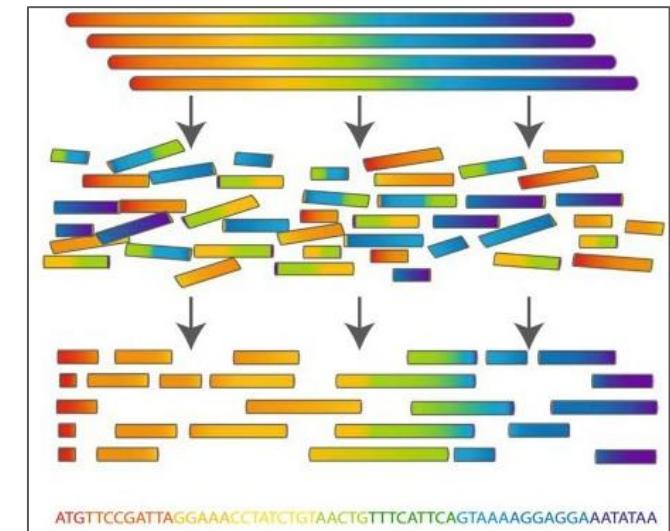
Sequence Length (bp)	Genome Element
25–75	SNPs, short frameshift mutations
100–400	Short functional signatures
500–1,000	Whole domains, single domain genes
1,000–5,000	Short operons, multidomain genes
5,000–10,000	Longer operons, some <i>cis</i> -control elements
>100,000	Prophages, pathogenicity islands, various mobile insertion elements
>1,000,000	Whole prokaryotic chromosome organization

C. I. Methods: Metagenome sequencing

Functional annotation

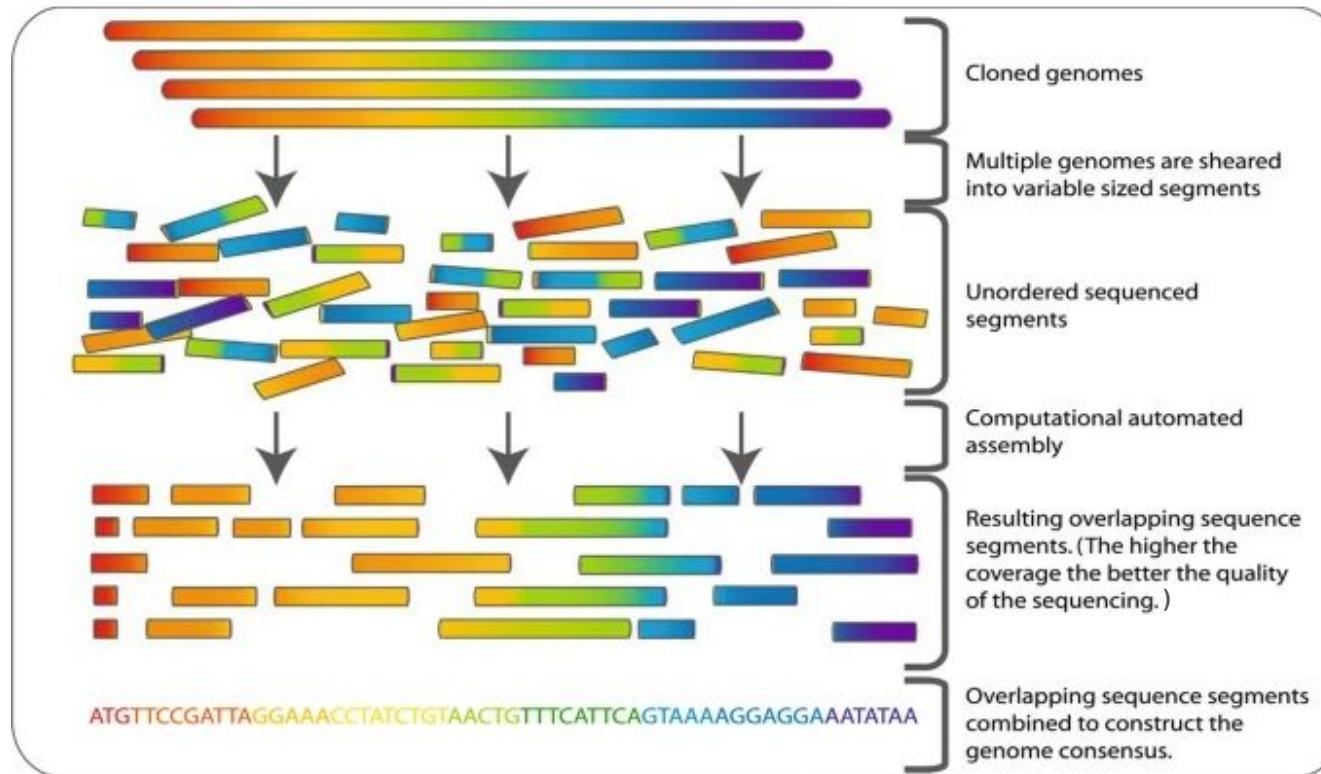


Sequence assembly

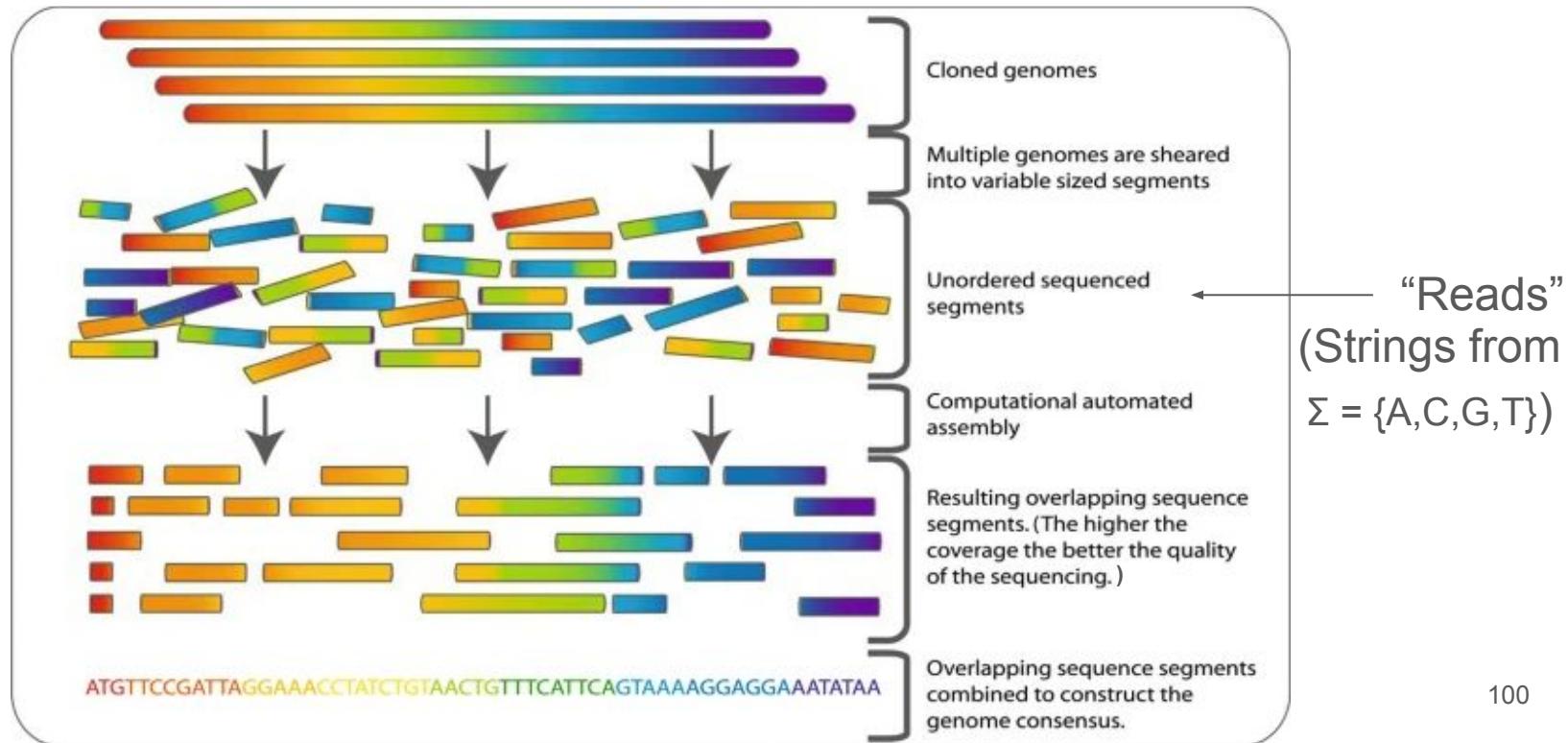


Metagenome sequencing enables two main additional types of analyses, compared to marker gene sequencing.

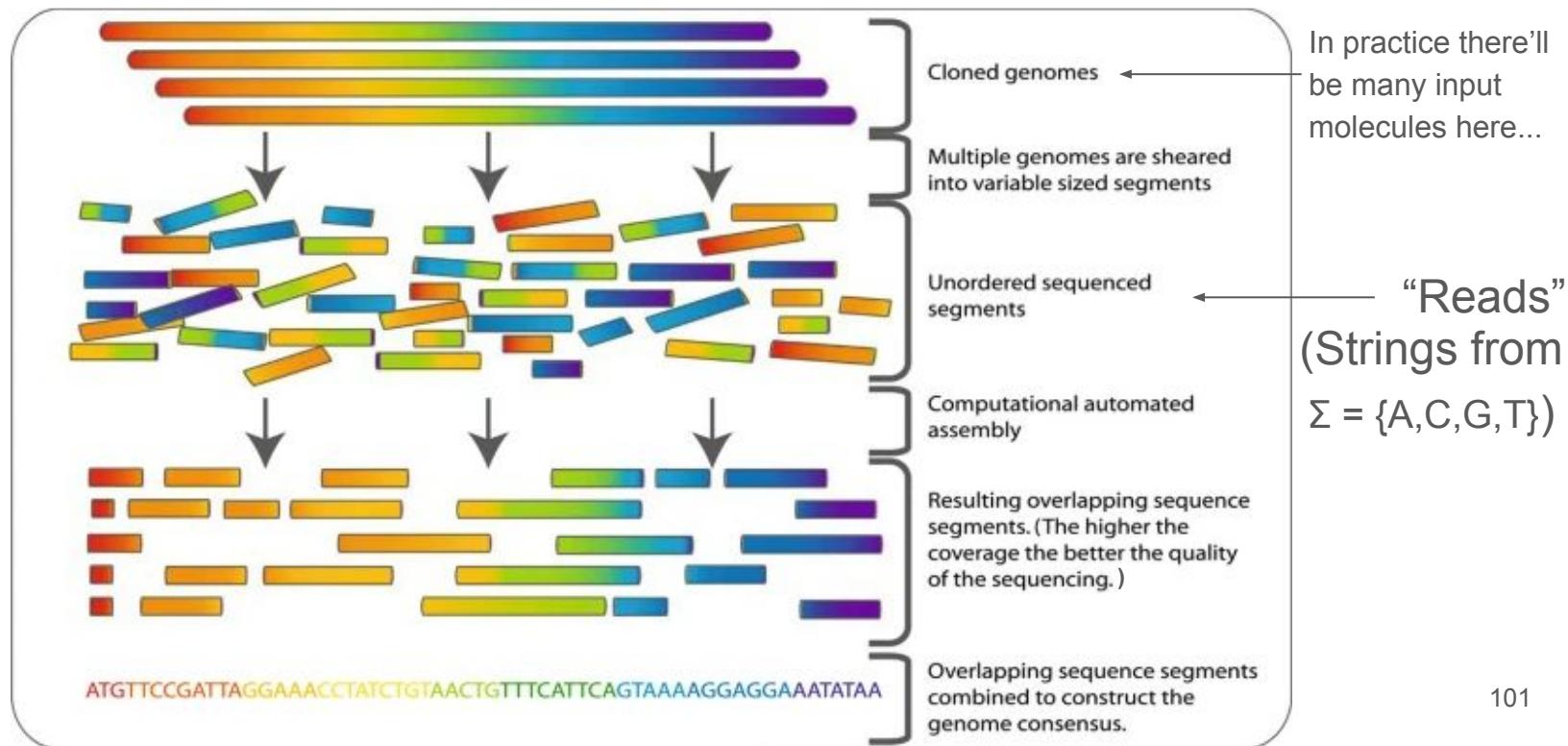
C. I. Methods: Metagenome sequence assembly



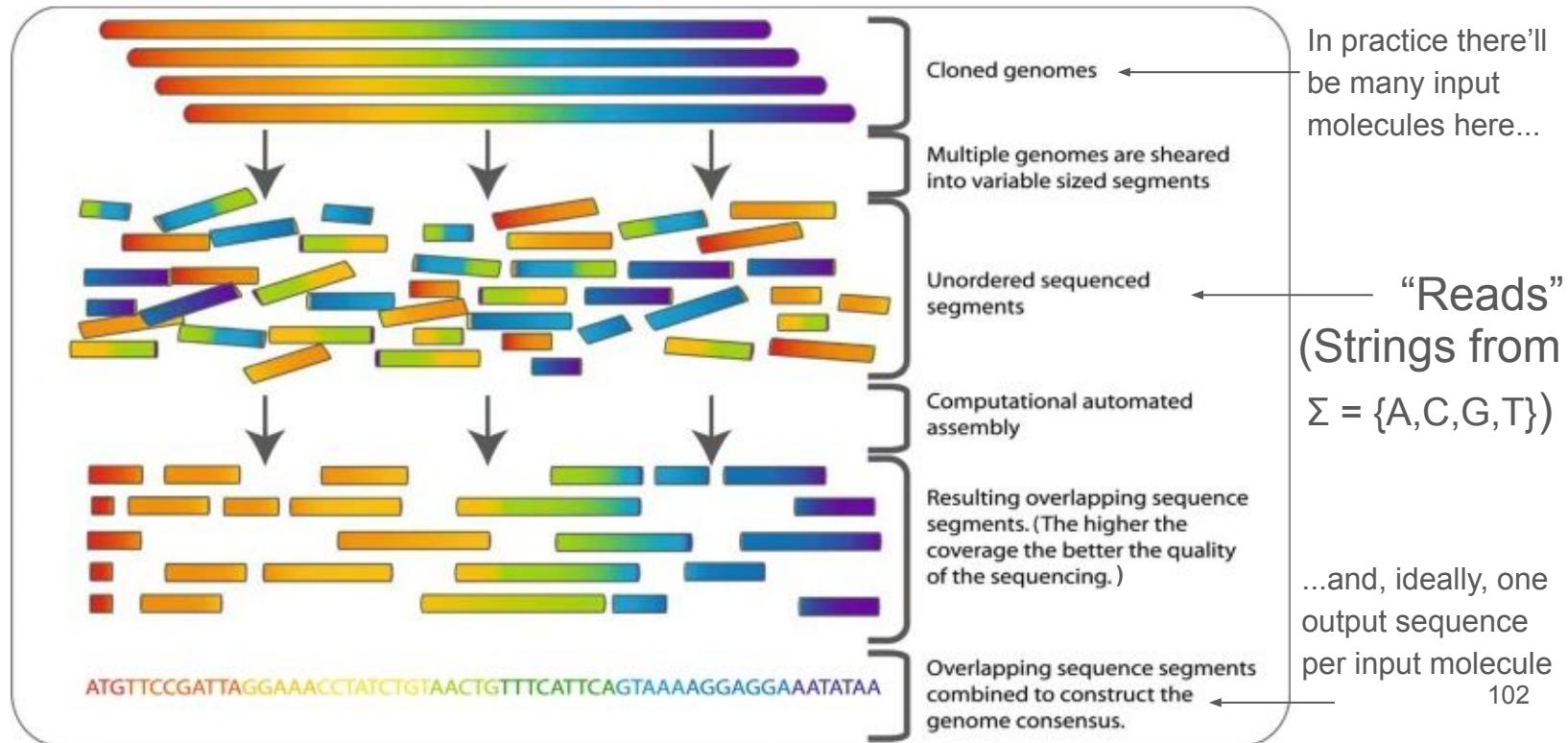
C. I. Methods: Metagenome sequence assembly



C. I. Methods: Metagenome sequence assembly



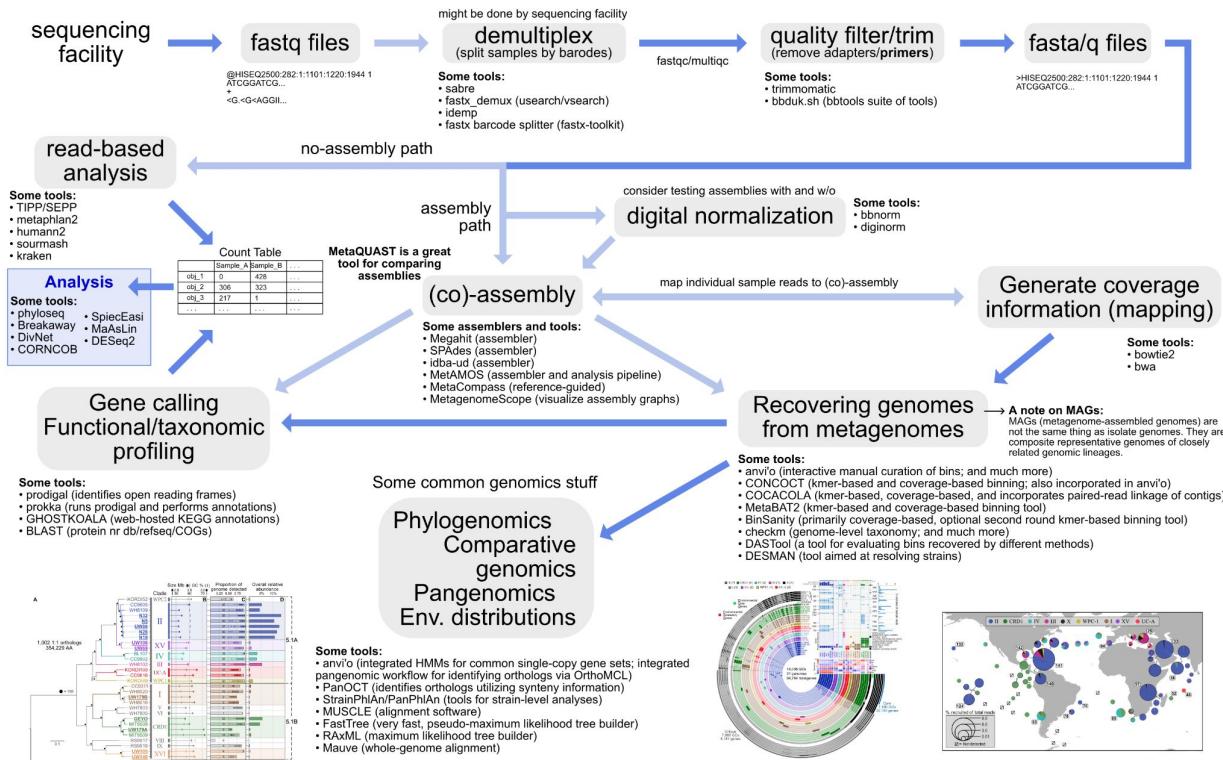
C. I. Methods: Metagenome sequence assembly



C. I. Methods: Metagenome sequence assembly

Overview of generic* metagenomics workflow

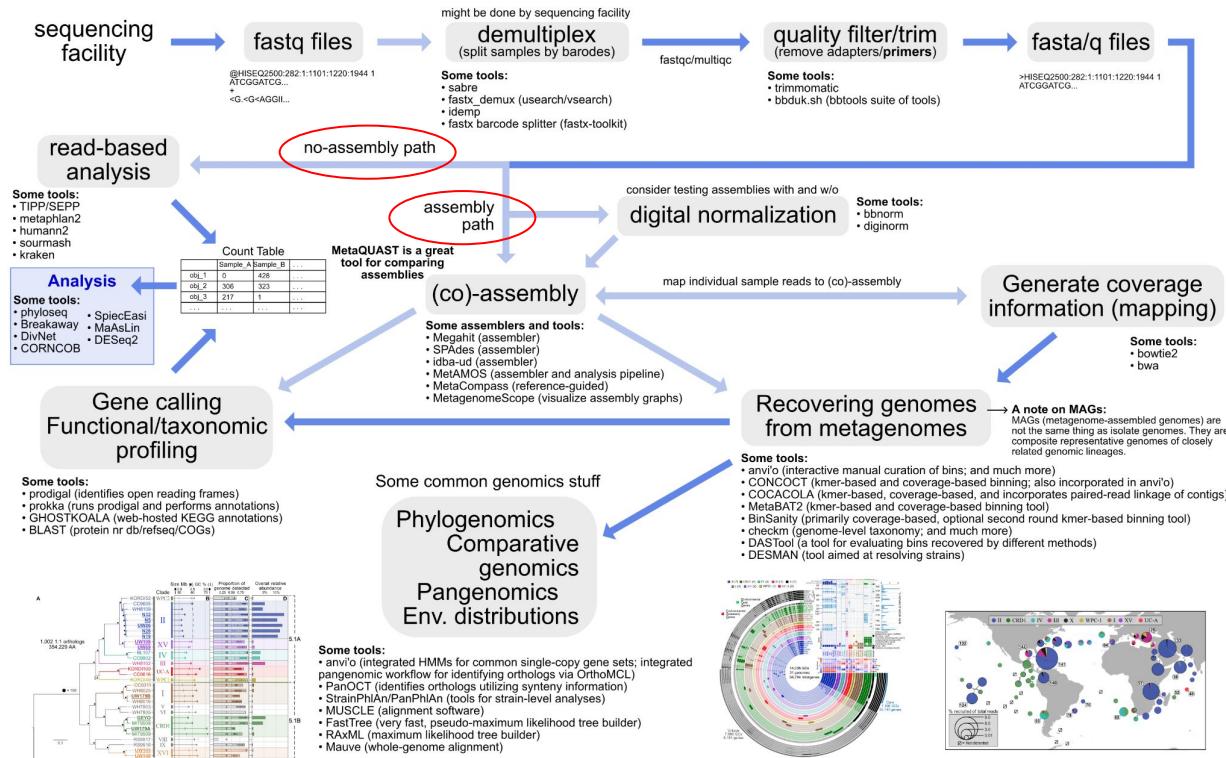
*This is generic; specific workflows can vary on the order of steps here and how they are done.



C. I. Methods: Metagenome sequence assembly

Overview of generic* metagenomics workflow

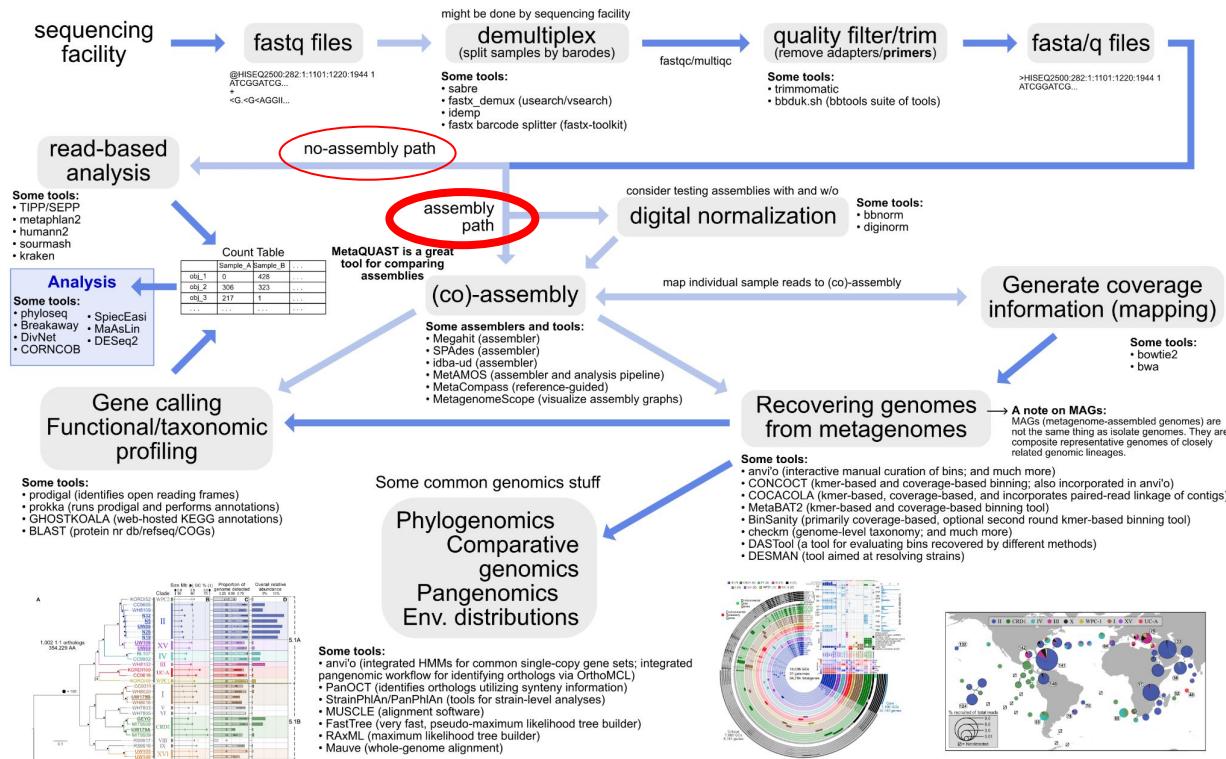
*This is generic; specific workflows can vary on the order of steps here and how they are done.



C. I. Methods: Metagenome sequence assembly

Overview of generic* metagenomics workflow

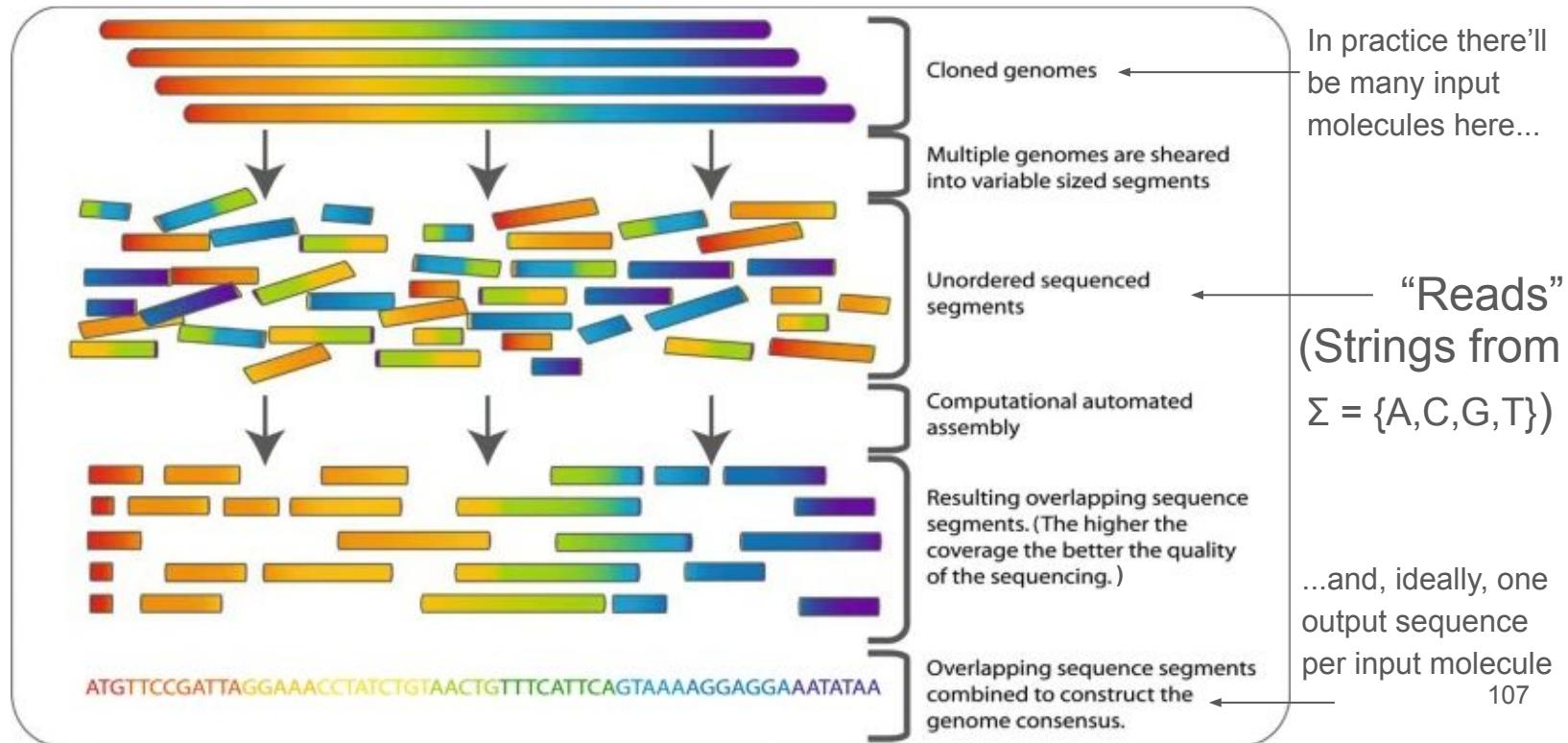
*This is generic; specific workflows can vary on the order of steps here and how they are done.



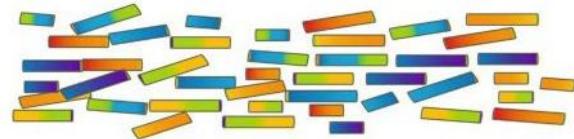
This talk

1. Introduction: Studying microbiomes
 - a. Why bother?
 - b. Why hasn't this research been more useful?
 - c. Defining a goal for this talk
2. Culture-independent (a.k.a. sequencing-based) methods
 - a. Marker gene sequencing
 - b. Metagenome sequencing
3. **Metagenome assembly**
 - a. Input (*reads*)
 - b. Outputs (*contigs*, an *assembly graph*, ...)
 - c. Methods (*de novo* vs. reference-based, overlap graph vs. de Bruijn graph, ...)
4. Future work: Solving the *strain separation problem*

Assembly (in the context of a metagenome)



Assembly: inputs

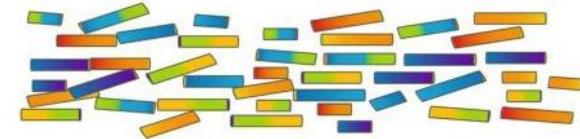


Reads: strings from the alphabet $\Sigma = \{\mathbf{A}, \mathbf{C}, \mathbf{G}, \mathbf{T}\}$

Occasionally this alphabet is extended if the base at a position is ambiguous:

e.g. $\mathbf{W} = (\mathbf{A} \text{ or } \mathbf{T})$, $\mathbf{S} = (\mathbf{C} \text{ or } \mathbf{G})$, $\mathbf{N} = (\mathbf{A} \text{ or } \mathbf{C} \text{ or } \mathbf{G} \text{ or } \mathbf{T})$, ...

Assembly: inputs



Reads: strings from the alphabet $\Sigma = \{\mathbf{A}, \mathbf{C}, \mathbf{G}, \mathbf{T}\}$

Occasionally this alphabet is extended if the base at a position is ambiguous:

e.g. $\mathbf{W} = (\mathbf{A} \text{ or } \mathbf{T})$, $\mathbf{S} = (\mathbf{C} \text{ or } \mathbf{G})$, $\mathbf{N} = (\mathbf{A} \text{ or } \mathbf{C} \text{ or } \mathbf{G} \text{ or } \mathbf{T})$, ...

Things that can vary based on the **sequencing technology** being used:

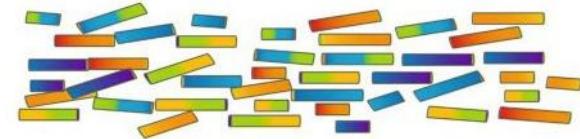
Read length

Read error rate

Number of reads

Read structure (e.g. single vs. *paired-end* reads)

Assembly: inputs



Reads: strings from the alphabet $\Sigma = \{\mathbf{A}, \mathbf{C}, \mathbf{G}, \mathbf{T}\}$

Occasionally this alphabet is extended if the base at a position is ambiguous:

e.g. $\mathbf{W} = (\mathbf{A} \text{ or } \mathbf{T})$, $\mathbf{S} = (\mathbf{C} \text{ or } \mathbf{G})$, $\mathbf{N} = (\mathbf{A} \text{ or } \mathbf{C} \text{ or } \mathbf{G} \text{ or } \mathbf{T})$, ...

Things that can vary based on the **sequencing technology** being used:

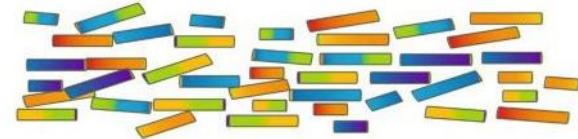
Read length

Read error rate

Number of reads

Read structure (e.g. single vs. *paired-end* reads)

Assembly: inputs



Reads: strings from the alphabet $\Sigma = \{\mathbf{A}, \mathbf{C}, \mathbf{G}, \mathbf{T}\}$

Occasionally this alphabet is extended if the base at a position is ambiguous:

e.g. $\mathbf{W} = (\mathbf{A} \text{ or } \mathbf{T})$, $\mathbf{S} = (\mathbf{C} \text{ or } \mathbf{G})$, $\mathbf{N} = (\mathbf{A} \text{ or } \mathbf{C} \text{ or } \mathbf{G} \text{ or } \mathbf{T})$, ...

Things that can vary based on the **sequencing technology** being used:

Read length

Read error rate

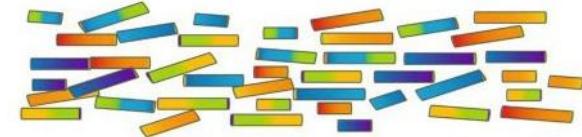
Number of reads

Read structure (e.g. single vs. *paired-end* reads)

Three (modern) sequencing technologies we'll focus on:

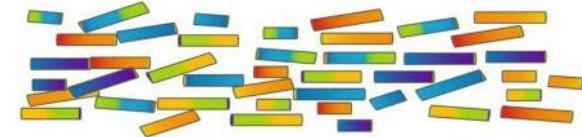
(1) short-read, (2) long, error-prone read, (3) HiFi

Assembly: sequencing technologies



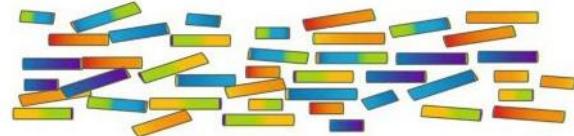
- Short-read technologies
 - e.g. Illumina
 - Read lengths: up to 150 nt
 - Error rates: less than 1%
- Long, error-prone read technologies
 - e.g. Oxford Nanopore, Pacific Biosciences
 - Read length: over 10,000 nt
 - Error rates: 10–25%
- Long, accurate read technologies (“HiFi”)
 - Pacific Biosciences circular consensus sequencing
 - Read length: over 10,000 nt, but generally shorter than long, error-prone reads
 - Error rates:
 - less than 1% for “point” mutations
 - somewhat higher (~5%) for insertions/deletions

Assembly: sequencing technologies



- Short-read technologies
 - e.g. Illumina
 - Read lengths: up to 150 nt
 - Error rates: **less than 1%**
- Long, error-prone read technologies
 - e.g. Oxford Nanopore, Pacific Biosciences
 - Read length: **over 10,000 nt**
 - Error rates: 10–25%
- Long, accurate read technologies (“HiFi”)
 - Pacific Biosciences circular consensus sequencing
 - Read length: over 10,000 nt, but generally shorter than long, error-prone reads
 - Error rates:
 - less than 1% for “point” mutations
 - somewhat higher (~5%) for insertions/deletions

Assembly: sequencing technologies

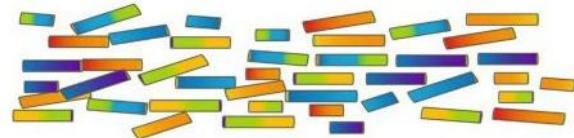


- Short-read technologies
 - e.g. Illumina
 - Read lengths: up to 150 nt
 - Error rates: **less than 1%**
- Long, error-prone read technologies
 - e.g. Oxford Nanopore, Pacific Biosciences
 - Read length: **over 10,000 nt**
 - Error rates: 10–25%
- Long, accurate read technologies (“HiFi”)
 - Pacific Biosciences circular consensus sequencing
 - Read length: over 10,000 nt, but generally shorter than long, error-prone reads
 - Error rates:
 - less than 1% for “point” mutations
 - somewhat higher (~5%) for insertions/deletions

This simple categorization ignores *a lot* of finicky details!

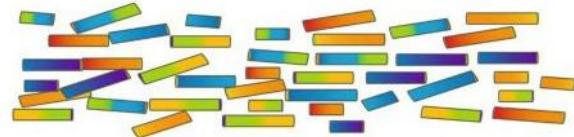
But this should be enough to understand how these technologies can complicate or simplify assembly.

Assembly: impacts of technologies

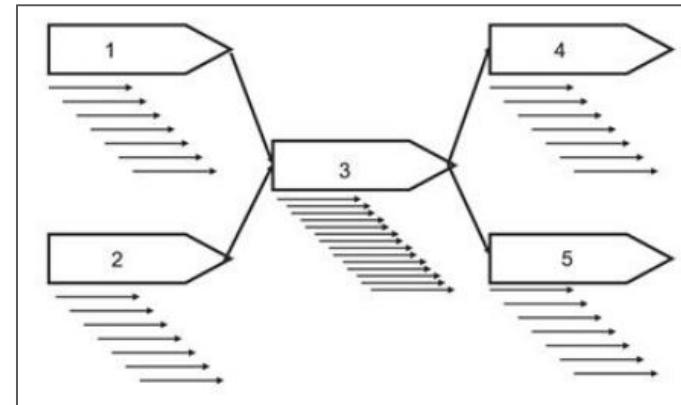


- Read lengths
 - Longer reads can span **repeats** (“identical, or nearly identical, stretches of DNA”).
 - Repeats that are longer than reads are especially challenging.
 - Accounting for repeats is the most important challenge in (metagenome) assembly!
 - So longer reads are generally better than short ones.

Assembly: impacts of technologies



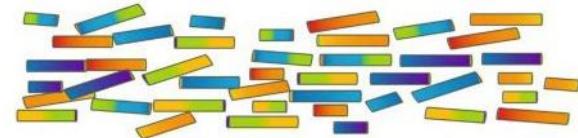
- Read lengths
 - Longer reads can span **repeats** (“identical, or nearly identical, stretches of DNA”).
 - Repeats that are longer than reads are especially challenging.
 - Accounting for repeats is the most important challenge in (metagenome) assembly!
 - So longer reads are generally better than short ones.



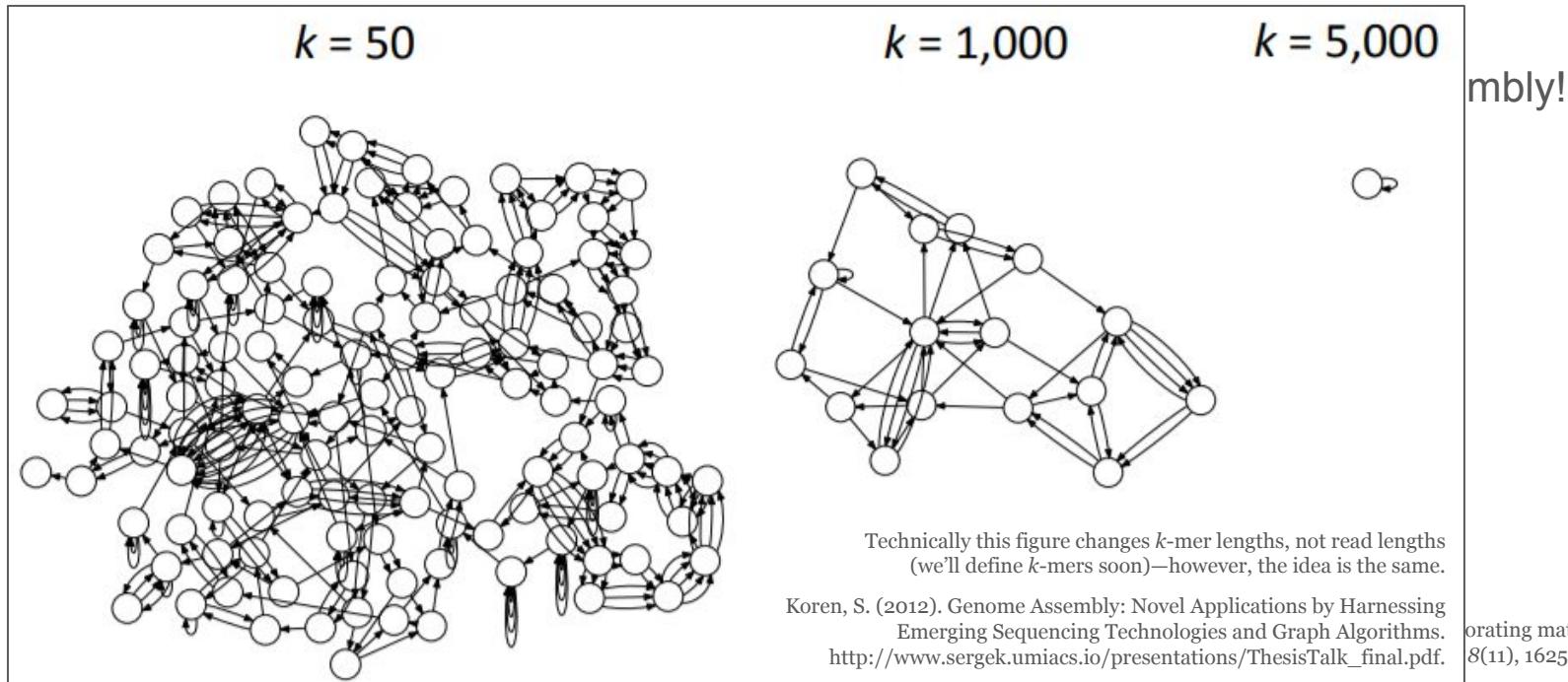
Koren, S., Treangen, T. J., & Pop, M. (2011). Bambus 2: scaffolding metagenomes. *Bioinformatics*, 27(21), 2964-2971.

Medvedev, P., Pham, S., Chaisson, M., Tesler, G., & Pevzner, P. (2011). Paired de bruijn graphs: a novel approach for incorporating mate pair information into genome assemblers. *Journal of Computational Biology*, 18(11), 1625-1634.

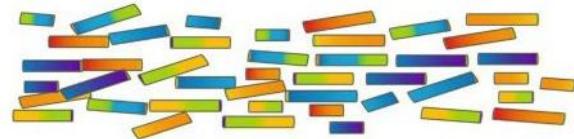
Assembly: impacts of technologies



- Read lengths
 - Longer reads can span **repeats** (“identical, or nearly identical, stretches of DNA”).

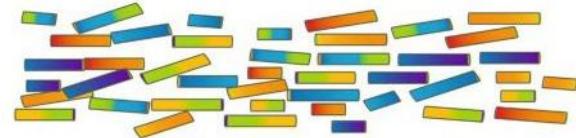


Assembly: impacts of technologies



- Read lengths
 - Longer reads can span **repeats** (“identical, or nearly identical, stretches of DNA”).
 - Repeats that are longer than reads are especially challenging.
 - Accounting for repeats is the most important challenge in (metagenome) assembly!
 - So longer reads are generally better than short ones.

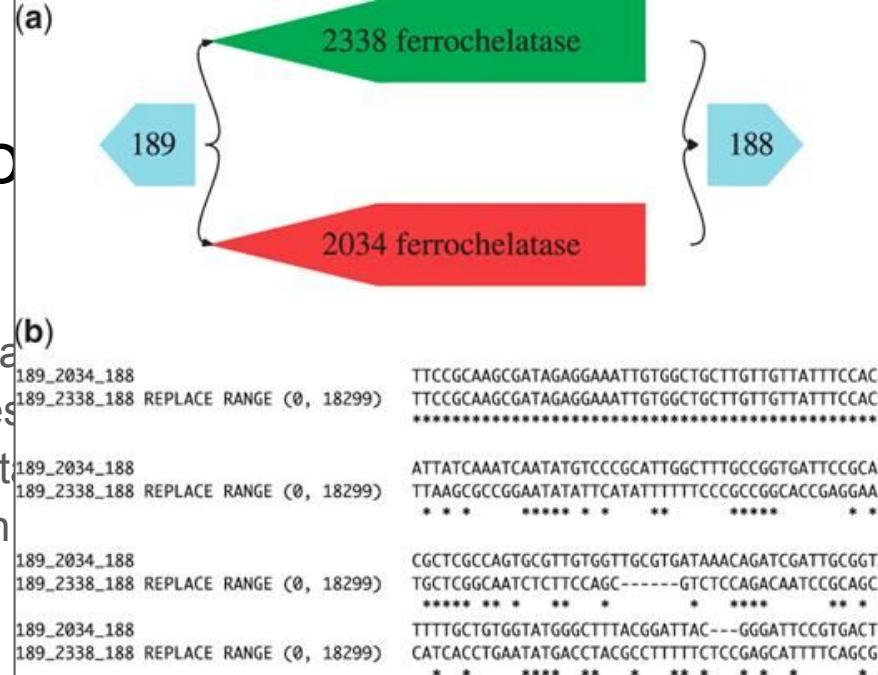
Assembly: impacts of technologies



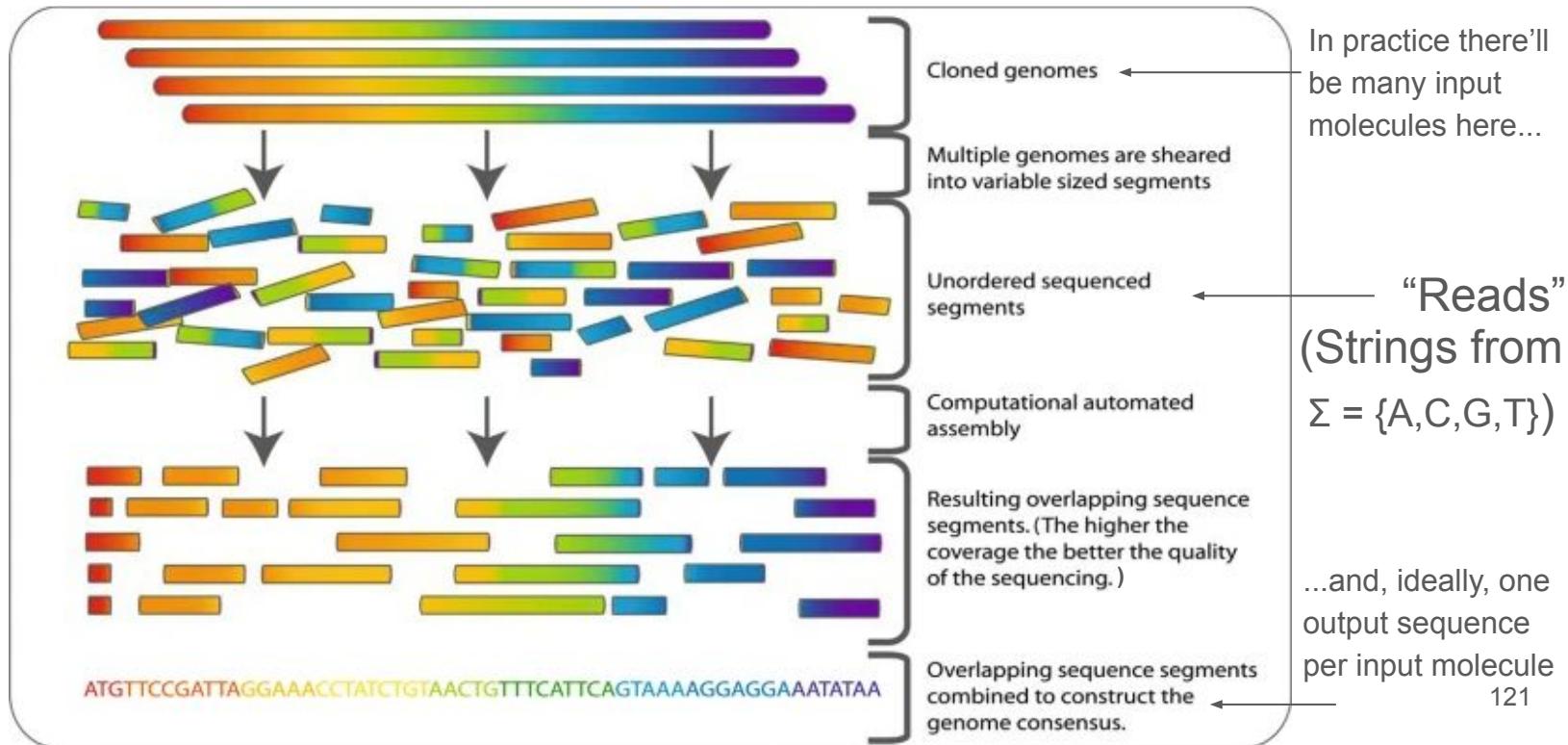
- Read lengths
 - Longer reads can span **repeats** (“identical, or nearly identical, stretches of DNA”).
 - Repeats that are longer than reads are especially challenging.
 - Accounting for repeats is the most important challenge in (metagenome) assembly!
 - So longer reads are generally better than short ones.
- Error rates
 - Lower error rates simplify **variant calling**, in which we attempt to distinguish real variations in the data from sequencing errors.
 - This is especially common for metagenome assembly: is a variation at some position the result of error, or is it indicative of a rare strain with this variation in its genome?

Assembly: impacts of technolo

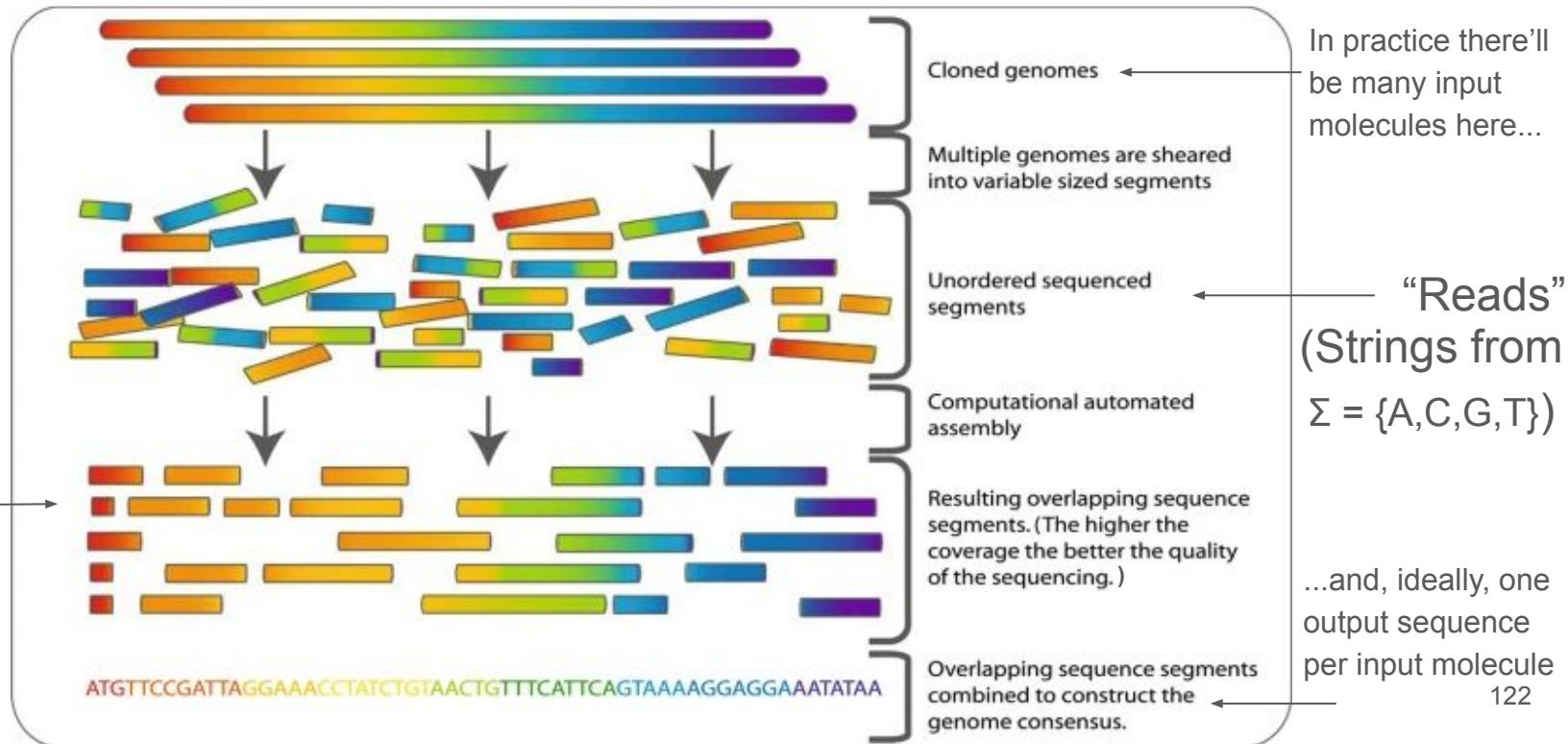
- Read lengths
 - Longer reads can span **repeats** (“identical DNA sequence”) → (b)
 - Repeats that are longer than reads are especially problematic
 - Accounting for repeats is the most important part of genome assembly
 - So longer reads are generally better than shorter ones
 - Error rates
 - Lower error rates simplify **variant calling**, in which we attempt to distinguish real variations in the data from sequencing errors.
 - This is especially common for metagenome assembly: is a variation at some position the result of error, or is it indicative of a rare strain with this variation in its genome?



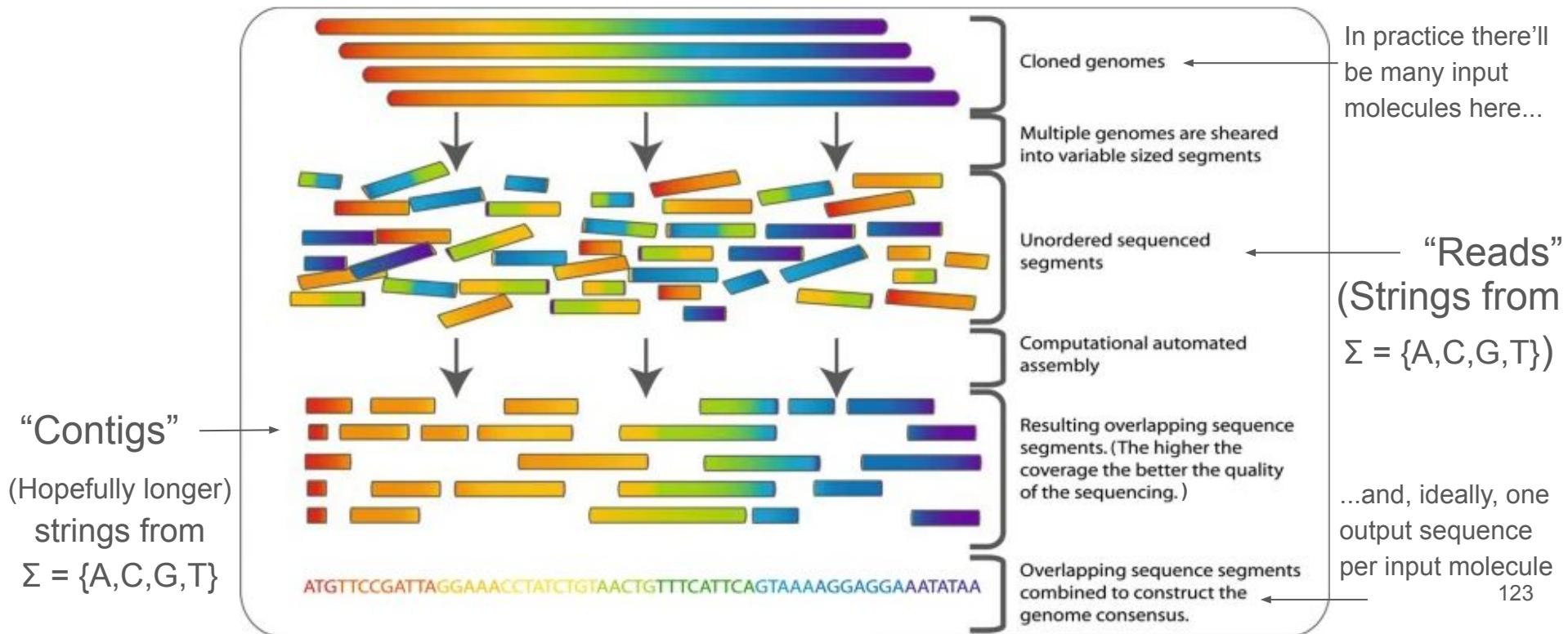
Assembly: outputs



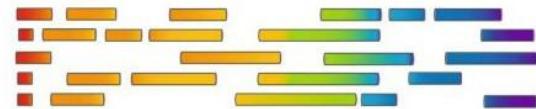
Assembly: outputs



Assembly: outputs



Assembly: outputs (contigs)



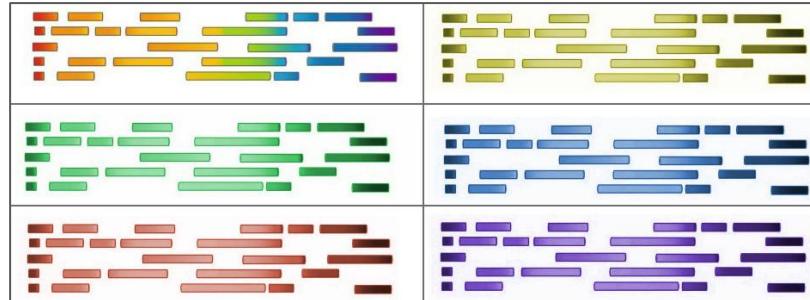
Ideally: one contig per input molecule of DNA
(e.g. each chromosome, plasmid, ...)

In practice: usually more contigs than that

Assembly: outputs (contigs)

Ideally: one contig per input molecule of DNA
(e.g. each chromosome, plasmid, ...)

In practice: usually more contigs than that

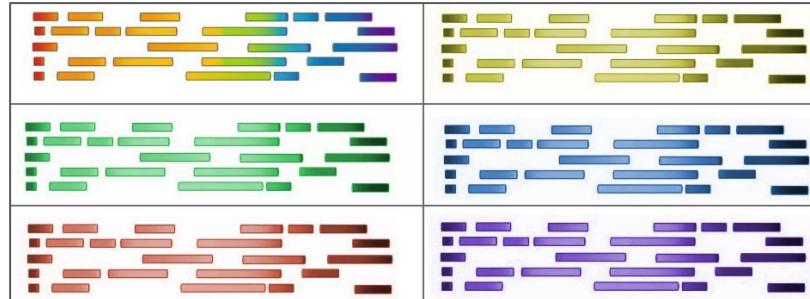


Some projects attempt to group contigs together into **bins** that likely originate from the same “genome”.

Assembly: outputs (contigs)

Ideally: one contig per input molecule of DNA
(e.g. each chromosome, plasmid, ...)

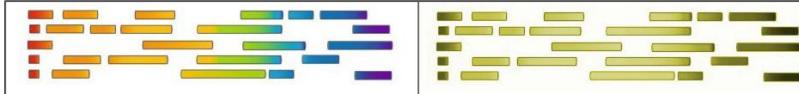
In practice: usually more contigs than that



Some projects attempt to group contigs together into **bins** that likely originate from the same “genome”.

Bins of contigs, or especially high-quality individual contigs, are referred to as **metagenome-assembled genomes (MAGs)**.

Assembly: outputs (contigs)



“We present a metagenomic HiFi assembly of a complex microbial community from sheep fecal material that resulted in **428 high-quality MAGs** from a single sample, the highest resolution achieved with metagenomic deconvolution to date.”

Bickhart, D. M., Kolmogorov, M., Tseng, E., Portik, D., Korobeynikov, A., Tolstoganov, I., ... & Smith, T. P. (2021). Generation of lineage-resolved complete metagenome-assembled genomes by precision phasing. *bioRxiv*.

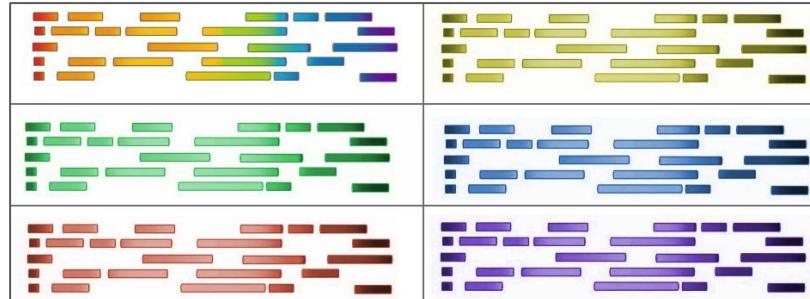
Bins of contigs, or especially high-quality individual contigs, are referred to as **metagenome-assembled genomes (MAGs)**.

Bowers, R. M., Kyrpides, N. C., Stepanauskas, R., Harmon-Smith, M., Doud, D., Reddy, T. B. K., ... & Woyke, T. (2017). Minimum information about a single amplified genome (MISAG) and a metagenome-assembled genome (MIMAG) of bacteria and archaea. *Nature Biotechnology*, 35(8), 725-731.

Assembly: outputs (contigs)

Ideally: one contig per input molecule of DNA
(e.g. each chromosome, plasmid, ...)

In practice: usually more contigs than that



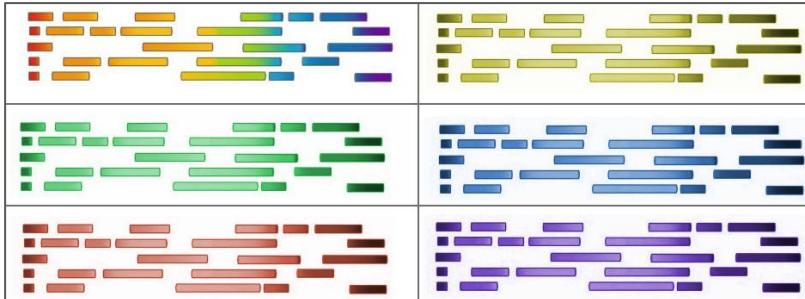
Some projects attempt to group contigs together into **bins** that likely originate from the same “genome”.

Bins of contigs, or especially high-quality individual contigs, are referred to as **metagenome-assembled genomes (MAGs)**.

Assembly: outputs (contigs)

Ideally: one contig per input molecule of DNA
(e.g. each chromosome, plasmid, ...)

In practice: usually more contigs than that



Assembly: outputs (assembly graph)

Most assembly algorithms model the problem as some sort of **graph traversal**.

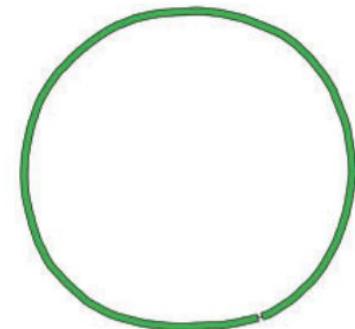
Assemblers usually output **assembly graphs**, which (generally) show overlaps between contigs.

Assembly: outputs (assembly graph)

Most assembly algorithms model the problem as some sort of **graph traversal**.

Assemblers usually output **assembly graphs**, which (generally) show overlaps between contigs.

Ideally: one connected component per input molecule of DNA



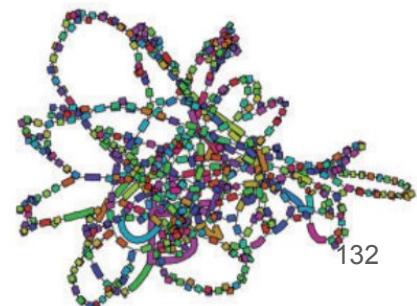
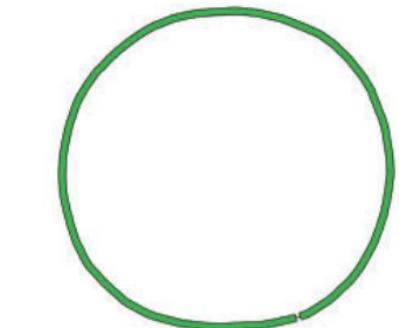
Assembly: outputs (assembly graph)

Most assembly algorithms model the problem as some sort of **graph traversal**.

Assemblers usually output **assembly graphs**, which (generally) show overlaps between contigs.

Ideally: one connected component per input molecule of DNA

In practice: the graph is usually tangled, fragmented, ...



Assembly: outputs (assembly graph)

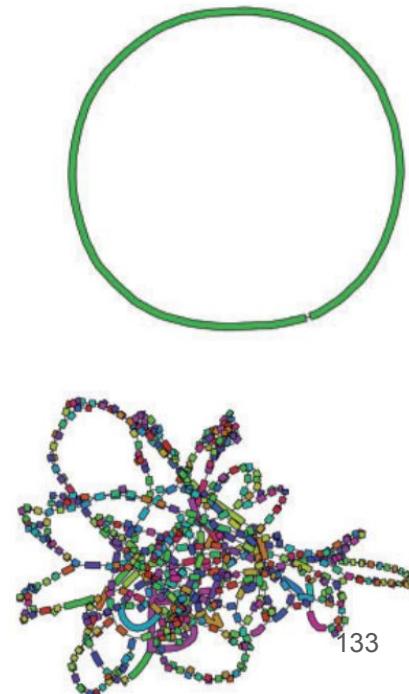
Most assembly algorithms model the problem as some sort of **graph traversal**.

Assemblers usually output **assembly graphs**, which (generally) show overlaps between contigs.

Ideally: one connected component per input molecule of DNA

In practice: the graph is usually tangled, fragmented, ...

These can be useful when “finishing” assemblies, or looking at subtle variations.



Assembly: outputs (assembly graph)

Most assembly algorithms model the problem as some sort of **graph traversal**.

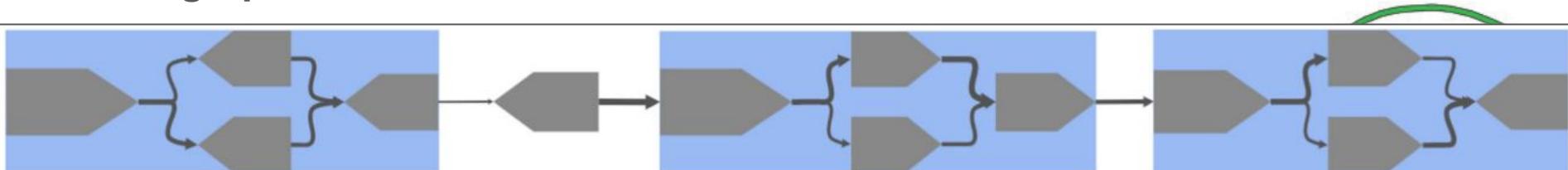
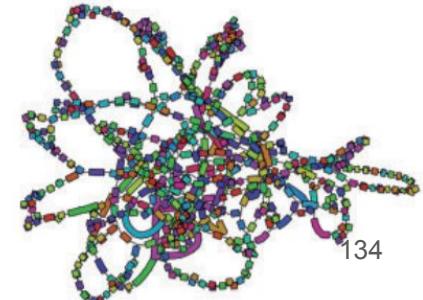


Fig. 2 Variants detected in one of the components of *Acinetobacter baumanii* scaffold graph. In this component, we find all the non-terminal nodes in a bubble are more than 97% identical to each other and originate from two different strains of *Acinetobacter baumannii* genome

Ghurye, J., Treangen, T., Fedarko, M., Hervey, W. J., & Pop, M. (2019). MetaCarvel: linking assembly graph motifs to biological variants. *Genome Biology*, 20(1), 1-14.

These can be useful when “finishing” assemblies, or looking at subtle variations.



This talk

1. Introduction: Studying microbiomes
 - a. Why bother?
 - b. Why hasn't this research been more useful?
 - c. Defining a goal for this talk
2. Culture-independent (a.k.a. sequencing-based) methods
 - a. Marker gene sequencing
 - b. Metagenome sequencing
3. **Metagenome assembly**
 - a. Input (*reads*)
 - b. Outputs (*contigs*, an *assembly graph*, ...)
 - c. **Methods** (*de novo* vs. reference-based, overlap graph vs. de Bruijn graph, ...)
4. Future work: Solving the *strain separation problem*

Assembly: methods

“A good genome assembler is like a good sausage:
you would rather not know what is inside.”

Assembly: methods (*de novo* vs. reference-based)

de novo assembly: Use only the read data available

“[...] reconstruction in its pure form, without consultation to previously resolved sequence including from genomes, transcripts, and proteins.”

Reference-based assembly: Also use available reference sequence(s)

“For some applications, sufficient information can be extracted from the mapping of reads to a reference sequence, such as a finished genome from a related individual.”

Assembly: methods (*de novo* vs. reference-based)

de novo assembly: Use only the read data available

Far more commonly used when working with metagenome sequencing data.

Reference-based assembly: Also use available reference sequence(s)

Some reference-based assemblers have been developed for metagenome sequencing data, but they have not (yet) seen widespread use in the field.

Assembly: methods (*de novo* vs. reference-based)

de novo assembly: Use only the read data available

Far more commonly used when working with metagenome sequencing data.

Reference-based assembly: Also use available reference sequence(s)

Some reference-based assemblers have been developed for metagenome sequencing data, but they have not (yet) seen widespread use in the field.

Assembly: methods (overlap vs. de Bruijn graphs)

OLC

DBG

Overlap graph (directed graph)

Nodes: input reads

Edges: an edge is created from $n_1 \rightarrow n_2$ if read n_1 overlaps with read n_2

de Bruijn graph (also a directed graph; takes a parameter k)

Nodes: unique $(k - 1)$ -mers (strings of length $k - 1$) in the reads

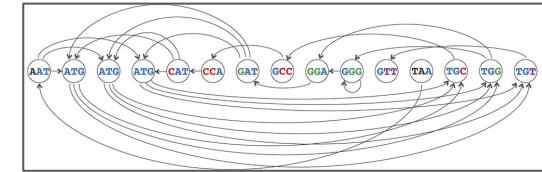
Edges: an edge is created from $n_1 \rightarrow n_2$ if there is a k -mer whose prefix is n_1 and whose suffix is n_2

Assembly: methods (overlap vs. de Bruijn graphs) OLC

Overlap graph (directed graph)

Nodes: input reads

Edges: an edge is created from $n_1 \rightarrow n_2$ if read n_1 overlaps with read n_2

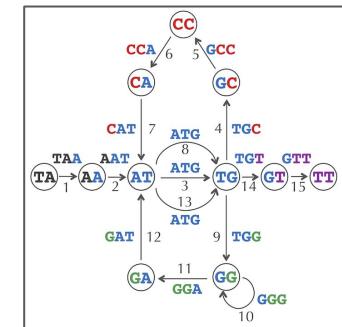


de Bruijn graph (also a directed graph; takes a parameter k)

Nodes: unique $(k - 1)$ -mers (strings of length $k - 1$) in the reads

Edges: an edge is created from $n_1 \rightarrow n_2$ if there is a k -mer whose prefix is n_1 and whose suffix is n_2

DBG



141

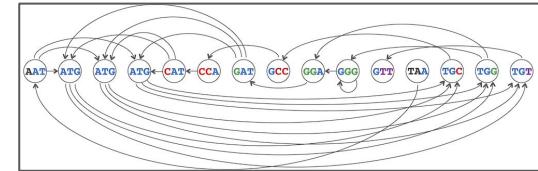
Assembly: methods (overlap vs. de Bruijn graphs) OLC

Overlap graph (directed graph)

Nodes: input reads

Edges: an edge is created from $n_1 \rightarrow n_2$ if read n_1 overlaps with read n_2

Goal: Find **Hamiltonian** Paths (or cycles) in this graph **NP-Complete!**



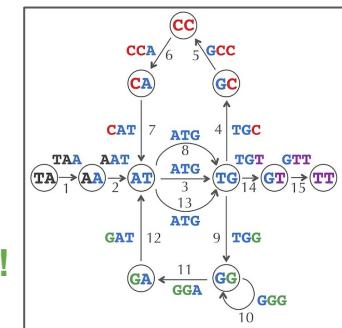
de Bruijn graph (also a directed graph; takes a parameter k)

Nodes: unique $(k - 1)$ -mers (strings of length $k - 1$) in the reads

Edges: an edge is created from $n_1 \rightarrow n_2$ if there is a k -mer whose prefix is n_1 and whose suffix is n_2

Goal: Find **Eulerian** Paths (or cycles) in this graph **Polynomial time!**

DBG



Assembly: methods (overlap vs. de Bruijn graphs)

OLC

- Newbler
- Celera
- ARACHNE
- Edena
- SHORTY
- HINGE
- BAUM
- Canu
- hifiasm
- ...

DBG

- EULER
- AllPaths
- SOAP
- Velvet
- ABySS
- E + V-SC
- SPAdes
- ABruijn
- Flye
- ...

Neither of these lists are comprehensive! Many assemblers (of both the OLC or DBG varieties) are still actively being developed and in use today. Also, most assemblers implement their own “twist” on how they use these graph structures, so cleanly categorizing assemblers in this way ignores many details.

Dida, F., & Yi, G. (2021). Empirical evaluation of methods for de novo genome assembly. *PeerJ Computer Science*, 7, e636.

Wang, A., Wang, Z., Li, Z., & Li, L. M. (2018). BAUM: improving genome assembly by adaptive unique mapping and local overlap-layout-consensus approach. *Bioinformatics*, 34(12), 2019-2028.

Miller, J. R., Koren, S., & Sutton, G. (2010). Assembly algorithms for next-generation sequencing data. *Genomics*, 95(6), 315-327.

Assembly: methods (overlap vs. de Bruijn graphs)

OLC

DBG

What do Eulerian and Hamiltonian cycles have to do with genome assembly?

Paul Medvedev^{1,2,3*}, Mihai Pop^{4,5}

- | | |
|---|---|
| <ul style="list-style-type: none">● BAUM● Canu● hifiasm● ... | <ul style="list-style-type: none">● SPAdes● ABruijn● Flye● ... |
|---|---|

Neither of these lists are comprehensive! Many assemblers (of both the OLC or DBG varieties) are still actively being developed and in use today. Also, most assemblers implement their own “twist” on how they use these graph structures, so cleanly categorizing assemblers in this way ignores many details.

Dida, F., & Yi, G. (2021). Empirical evaluation of methods for de novo genome assembly. *PeerJ Computer Science*, 7, e636.

Wang, A., Wang, Z., Li, Z., & Li, L. M. (2018). BAUM: improving genome assembly by adaptive unique mapping and local overlap-layout-consensus approach. *Bioinformatics*, 34(12), 2019–2028.

Miller, J. R., Koren, S., & Sutton, G. (2010). Assembly algorithms for next-generation sequencing data. *Genomics*, 95(6), 315–327.

Assembly: methods (overlap vs. de Bruijn graphs)

OLC

What do Euler
cycles have
to do with geno-

Paul Medvedev^{1,2,3*}, Mihai

- BAUM
- Canu
- hifiasm
- ...

DBG



Neither of these lists are comprehensive! Many assemblers (of both the OLC or DBG varieties) are still actively being developed and in use today. Also, most assemblers implement their own “twist” on how they use these graph structures, so cleanly categorizing assemblers in this way ignores many details.

Dida, F., & Yi, G. (2021). Empirical evaluation of methods for de novo genome assembly. *PeerJ Computer Science*, 7, e636.

Wang, A., Wang, Z., Li, Z., & Li, L. M. (2018). BAUM: improving genome assembly by adaptive unique mapping and local overlap-layout-consensus approach. *Bioinformatics*, 34(12), 2019–2028.

Miller, J. R., Koren, S., & Sutton, G. (2010). Assembly algorithms for next-generation sequencing data. *Genomics*, 95(6), 315–327.

Assembly: methods (overlap vs. de Bruijn graphs)

What do Euler
cycles have
to do with geno-

Paul Medvedev^{1,2,3*}, Mihai

OLC

DBG



- BAUM
- Canu
- hifiasm
- ...

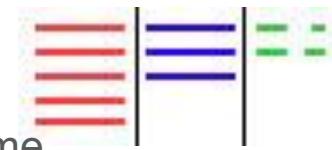
Neither of these lists are comprehensive! Many assemblers (of both the OLC or DBG varieties) are still actively being developed and in use today. Also, most assemblers implement their own “twist” on how they use these graph structures, so cleanly categorizing assemblers in this way ignores many details.

Kolmogorov, M., Bickhart, D. M., Behsaz, B., Gurevich, A., Rayko, M., Shin, S. B., ... & Pevzner, P. A. (2020). metaFlye: scalable long-read metagenome assembly using repeat graphs. *Nature Methods*, 17(11), 1103-1110.

Nurk, S., Meleshko, D., Korobeynikov, A., & Pevzner, P. A. (2017). metaSPAdes: a new versatile metagenomic assembler. *Genome Research*, 27(5), 824-834.

Namiki, T., Hachiya, T., Tanaka, H., & Sakakibara, Y. (2012). MetaVelvet: an extension of Velvet assembler to *de novo* metagenome assembly from short sequence reads. *Nucleic Acids Research*, 40(20), e155-e155.

Assembly: methods (single-genome vs. metagenome)



- **Problem:** Uneven coverage
 - Different microbes are present at different abundances in a microbiome.
 - The *coverage*, or number of reads “supporting” a position in a genome, varies a lot!
 - Worst case scenario: some genome(s) are only partially covered in the reads.

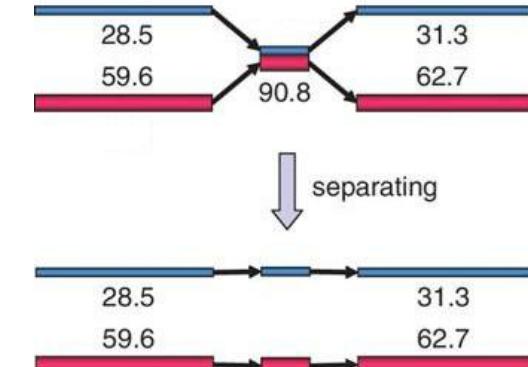
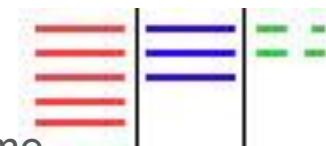
Kolmogorov, M., Bickhart, D. M., Behsaz, B., Gurevich, A., Rayko, M., Shin, S. B., ... & Pevzner, P. A. (2020). metaFlye: scalable long-read metagenome assembly using repeat graphs. *Nature Methods*, 17(11), 1103-1110.

Nurk, S., Meleshko, D., Korobeynikov, A., & Pevzner, P. A. (2017). metaSPAdes: a new versatile metagenomic assembler. *Genome Research*, 27(5), 824-834.

Namiki, T., Hachiya, T., Tanaka, H., & Sakakibara, Y. (2012). MetaVelvet: an extension of Velvet assembler to *de novo* metagenome assembly from short sequence reads. *Nucleic Acids Research*, 40(20), e155-e155.

Assembly: methods (single-genome vs. metagenome)

- **Problem:** Uneven coverage
 - Different microbes are present at different abundances in a microbiome.
 - The *coverage*, or number of reads “supporting” a position in a genome, varies a lot!
 - Worst case scenario: some genome(s) are only partially covered in the reads.
- **Problem:** Intergenomic repeats
 - Stretches of DNA shared by many organisms are repeats!
 - For example, marker genes (e.g. rRNA genes).
 - Some approaches attempt to get around this by *exploiting* uneven coverages across genomes.

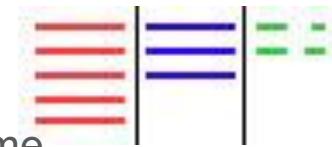


Kolmogorov, M., Bickhart, D. M., Behsaz, B., Gurevich, A., Rayko, M., Shin, S. B., ... & Pevzner, P. A. (2020). metaFlye: scalable long-read metagenome assembly using repeat graphs. *Nature Methods*, 17(11), 1103–1110.

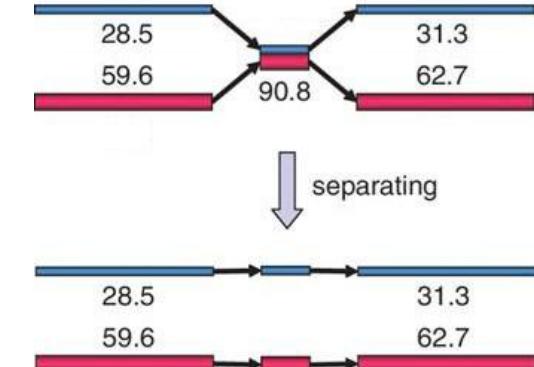
Nurk, S., Meleshko, D., Korobeynikov, A., & Pevzner, P. A. (2017). metaSPAdes: a new versatile metagenomic assembler. *Genome Research*, 27(5), 824–834.

Namiki, T., Hachiya, T., Tanaka, H., & Sakakibara, Y. (2012). MetaVelvet: an extension of Velvet assembler to *de novo* metagenome assembly from short sequence reads. *Nucleic Acids Research*, 40(20), e155–e155.

Assembly: methods (single-genome vs. metagenome)



- **Problem:** Uneven coverage
 - Different microbes are present at different abundances in a microbiome.
 - The *coverage*, or number of reads “supporting” a position in a genome, varies a lot!
 - Worst case scenario: some genome(s) are only partially covered in the reads.
- **Problem:** Intergenomic repeats
 - Stretches of DNA shared by many organisms are repeats!
 - For example, marker genes (e.g. rRNA genes).
 - Some approaches attempt to get around this by *exploiting* uneven coverages across genomes.
- **Problem:** Strain mixtures
 - Similar to intergenomic repeats: a microbiome can contain many strains of a microbe.
 - How can we distinguish between “real” variations and sequencing errors?
 - *Separating* these strains’ genomes is very challenging, especially with short reads.¹⁴⁹



Kolmogorov, M., Bickhart, D. M., Behsaz, B., Gurevich, A., Rayko, M., Shin, S. B., ... & Pevzner, P. A. (2020). metaFlye: scalable long-read metagenome assembly using repeat graphs. *Nature Methods*, 17(11), 1103-1110.

Nurk, S., Meleshko, D., Korobeynikov, A., & Pevzner, P. A. (2017). metaSPAdes: a new versatile metagenomic assembler. *Genome Research*, 27(5), 824-834.

Namiki, T., Hachiya, T., Tanaka, H., & Sakakibara, Y. (2012). MetaVelvet: an extension of Velvet assembler to *de novo* metagenome assembly from short sequence reads. *Nucleic Acids Research*, 40(20), e155-e155.

Assembly: methods (single-genome vs. metagenome)

- **Problem:** Strain mixtures
 - Similar to intergenic repeats: a microbiome can contain many strains of a microbe.
 - How can we distinguish between “real” variations and sequencing errors?
 - *Separating* these strains’ genomes is very challenging, especially with short reads.

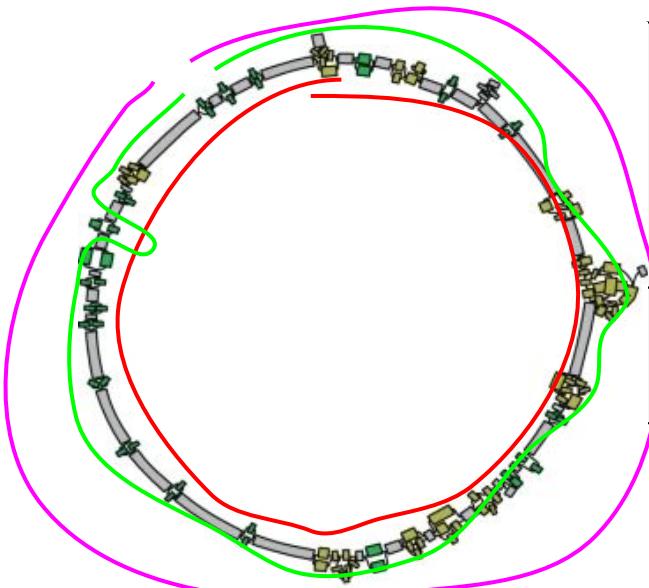
Kolmogorov, M., Bickhart, D. M., Behsaz, B., Gurevich, A., Rayko, M., Shin, S. B., ... & Pevzner, P. A. (2020). metaFly: scalable long-read metagenome assembly using repeat graphs. *Nature Methods*, 17(11), 1103-1110.

Nurk, S., Meleshko, D., Korobeynikov, A., & Pevzner, P. A. (2017). metaSPAdes: a new versatile metagenomic assembler. *Genome Research*, 27(5), 824-834.

Namiki, T., Hachiya, T., Tanaka, H., & Sakakibara, Y. (2012). MetaVelvet: an extension of Velvet assembler to *de novo* metagenome assembly from short sequence reads. *Nucleic Acids Research*, 40(20), e155-e155.

Assembly: methods (single-genome vs. metagenome)

- **Problem:** Strain mixtures
 - Similar to intergenic repeats: a microbiome can contain many strains of a microbe.
 - How can we distinguish between “real” variations and sequencing errors?
 - Separating these strains’ genomes is very challenging, especially with short reads.



“An assembly graph of a single connected component in the sheep microbiome dataset before strain collapsing [...] The component represents a bacterial genome of the *Clostridia* class [...] There are 20 simple bubbles (shown in green) and 10 superbubbles (shown in yellow) that account for 1.2 Mbp out of 2.4 Mbp long genome.”

Kolmogorov, M., Bickhart, D. M., Behsaz, B., Gurevich, A., Rayko, M., Shin, S. B., ... & Pevzner, P. A. (2020). metaFly: scalable long-read metagenome assembly using repeat graphs. *Nature Methods*, 17(11), 1103-1110.

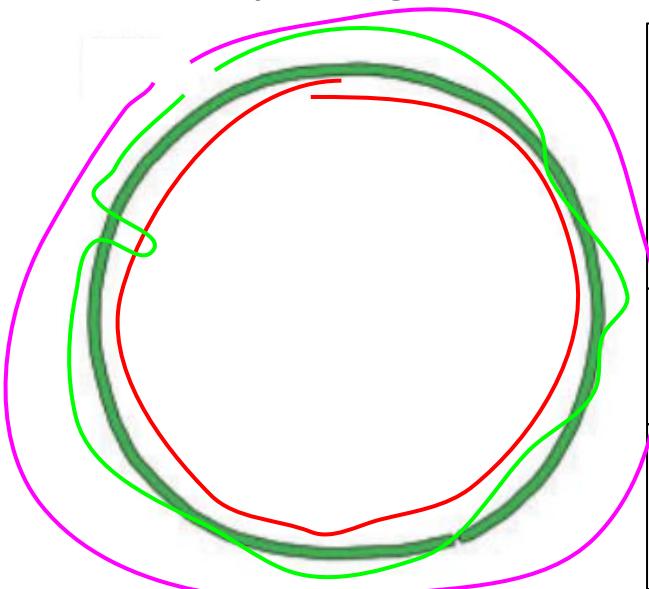
Kolmogorov, M., Bickhart, D. M., Behsaz, B., Gurevich, A., Rayko, M., Shin, S. B., ... & Pevzner, P. A. (2020). metaFly: scalable long-read metagenome assembly using repeat graphs. *Nature Methods*, 17(11), 1103–1110.

Nurk, S., Meleshko, D., Korobeynikov, A., & Pevzner, P. A. (2017). metaSPAdes: a new versatile metagenomic assembler. *Genome Research*, 27(5), 824–834.

Namiki, T., Hachiya, T., Tanaka, H., & Sakakibara, Y. (2012). MetaVelvet: an extension of Velvet assembler to *de novo* metagenome assembly from short sequence reads. *Nucleic Acids Research*, 40(20), e155–e155.

Assembly: methods (single-genome vs. metagenome)

- **Problem:** Strain mixtures
 - Similar to intergenic repeats: a microbiome can contain many strains of a microbe.
 - How can we distinguish between “real” variations and sequencing errors?
 - Separating these strains’ genomes is very challenging, especially with short reads.



“An assembly graph of a single connected component in the sheep microbiome dataset before strain collapsing [...] The component represents a bacterial genome of the *Clostridia* class [...] There are 20 simple bubbles (shown in green) and 10 superbubbles (shown in yellow) that account for 1.2 Mbp out of 2.4 Mbp long genome.”

Kolmogorov, M., Bickhart, D. M., Behsaz, B., Gurevich, A., Rayko, M., Shin, S. B., ... & Pevzner, P. A. (2020). metaFly: scalable long-read metagenome assembly using repeat graphs. *Nature Methods*, 17(11), 1103–1110.

Different assemblers will do different things to deal with these sorts of subtle variations: even “smooth” contigs can conceal a lot of variation. Sometimes this is desirable—sometimes not!

This talk

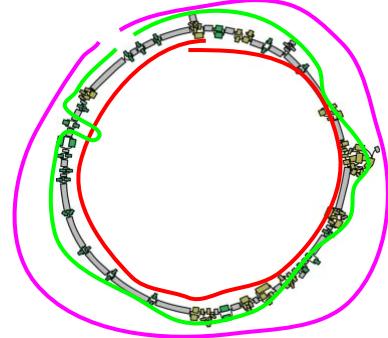
1. Introduction: Studying microbiomes
 - a. Why bother?
 - b. Why hasn't this research been more useful?
 - c. Defining a goal for this talk
2. Culture-independent (a.k.a. sequencing-based) methods
 - a. Marker gene sequencing
 - b. Metagenome sequencing
3. Metagenome assembly
 - a. Input (*reads*)
 - b. Outputs (*contigs*, an *assembly graph*, ...)
 - c. Methods (*de novo* vs. reference-based, overlap graph vs. de Bruijn graph, ...)
4. Future work: Solving the *strain separation problem*

(Vicedomini et al.'s definition of "strain" matches the one I defined ~45 minutes ago, i.e. a completely unique genome.)

Future work

Strain separation problem (Vicedomini et al., 2021)

"The reconstruction of partial or complete DNA sequences corresponding to strains, at the base level."

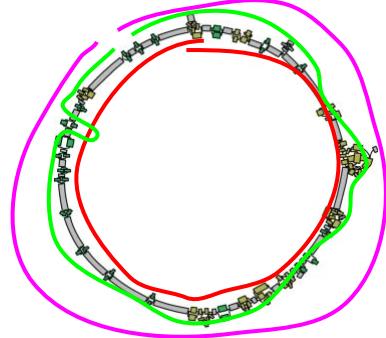


(Vicedomini et al.'s definition of "strain" matches the one I defined ~45 minutes ago, i.e. a completely unique genome.)

Future work

Strain separation problem (Vicedomini et al., 2021)

"The reconstruction of partial or complete DNA sequences corresponding to strains, at the base level."



(Local) Metagenome individual haplotyping problem (Nicholls et al., 2021)

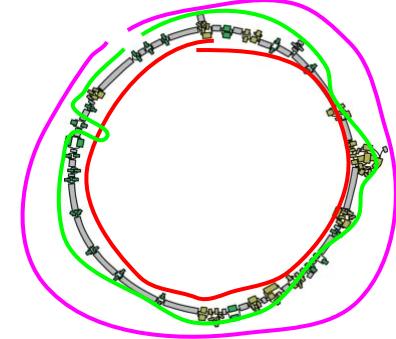
"The ideal output [...] is the collection of whole-genome sequences representing all the individual organisms in a microbial community."

(Vicedomini et al.'s definition of "strain" matches the one I defined ~45 minutes ago, i.e. a completely unique genome.)

Future work

Strain separation problem (Vicedomini et al., 2021)

"The reconstruction of partial or complete DNA sequences corresponding to strains, at the base level."



(Local) Metagenome individual haplotyping problem (Nicholls et al., 2021)

"The ideal output [...] is the collection of whole-genome sequences representing all the individual organisms in a microbial community."

Haplotype assembly problem (Lancia et al., 2001)

"Given a set of fragments obtained by DNA sequencing from the two copies of a chromosome, reconstruct two haplotypes that would be compatible with all the fragments observed."

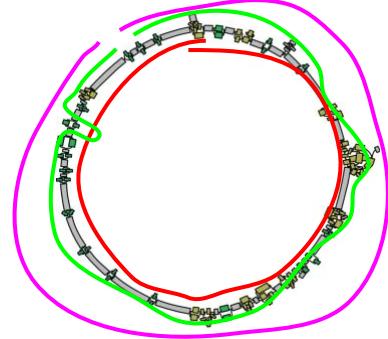
Vicedomini, R., Quince, C., Darling, A. E., & Chikhi, R. (2021). Strainberry: automated strain separation in low-complexity metagenomes using long reads. *Nature Communications*, 12(1), 1-14.

Nicholls, S. M., Aubrey, W., De Grave, K., Schietgat, L., Creevey, C. J., & Clare, A. (2021). On the complexity of haplotyping a microbial community. *Bioinformatics*, 37(10), 1360-1366.

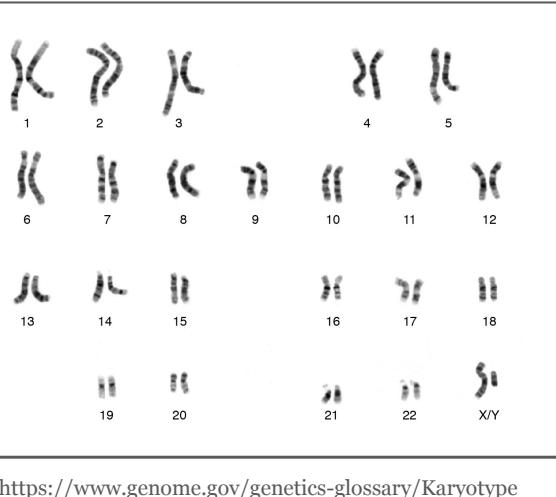
Lancia, G., Bafna, V., Istrail, S., Lippert, R., & Schwartz, R. (2001, August). SNPs problems, complexity, and algorithms. In *European Symposium on Algorithms* (pp. 182-193). Springer, Berlin, Heidelberg.

Future work: “haplotype”?

Humans (or other **diploid** organisms) usually have two copies of each chromosome.



Ordinary assemblers often “smooth out” differences between the chromosomes, creating contigs that are **chimeras** of both chromosomes’ sequences.



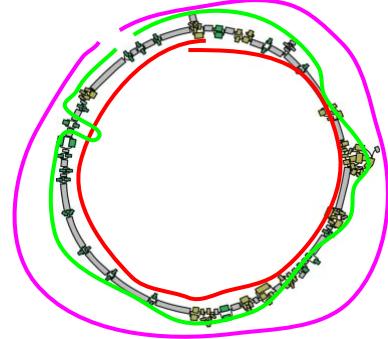
In **haplotype phasing**, we attempt to fix this. Usually, this involves looking for variations which occur on the same reads.

157

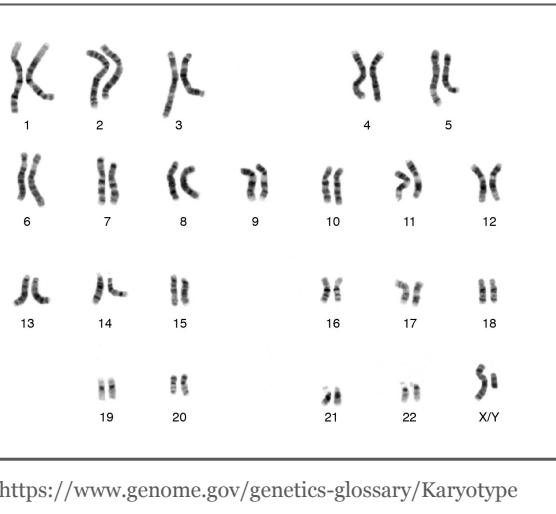
Lancia, G., Bafna, V., Istrail, S., Lippert, R., & Schwartz, R. (2001, August). SNPs problems, complexity, and algorithms. In *European Symposium on Algorithms* (pp. 182-193). Springer, Berlin, Heidelberg.

Future work: “haplotype”?

Humans (or other **diploid** organisms) usually have two copies of each chromosome.



Ordinary assemblers often “smooth out” differences between the chromosomes, creating contigs that are **chimeras** of both chromosomes’ sequences.



In **haplotype phasing**, we attempt to fix this. Usually, this involves looking for variations which occur on the same reads.

Even in the case of a single human genome, this problem is **NP-Hard**. It gets worse for metagenomes!

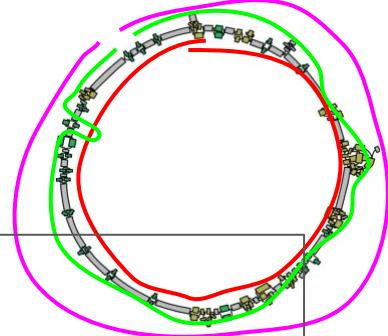
158

Lancia, G., Bafna, V., Istrail, S., Lippert, R., & Schwartz, R. (2001, August). SNPs problems, complexity, and algorithms. In *European Symposium on Algorithms* (pp. 182-193). Springer, Berlin, Heidelberg.

Future work: “metagenomic haplotyping”

A solution to [metagenomic haplotyping] is confounded by five problems:

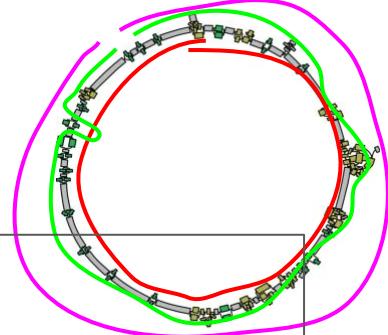
- (i) DNA from every genome needs to be extracted and sequenced to a depth sufficient for recovery,
- (ii) genomes share homologous regions that require disambiguation,
- (iii) reads may be of an insufficient length to disambiguate repeats or resolve bridges between variants,
- (iv) sequencing error can be indistinguishable from rare haplotypes and
- (v) the presence of an unknown number of haplotypes complicates the already computationally difficult (**NP-hard**) (Ciliberti et al., 2005) problem of haplotyping.



Future work: “metagenomic haplotyping”

A solution to [metagenomic haplotyping] is confounded by five problems:

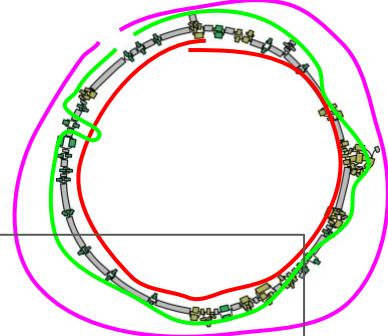
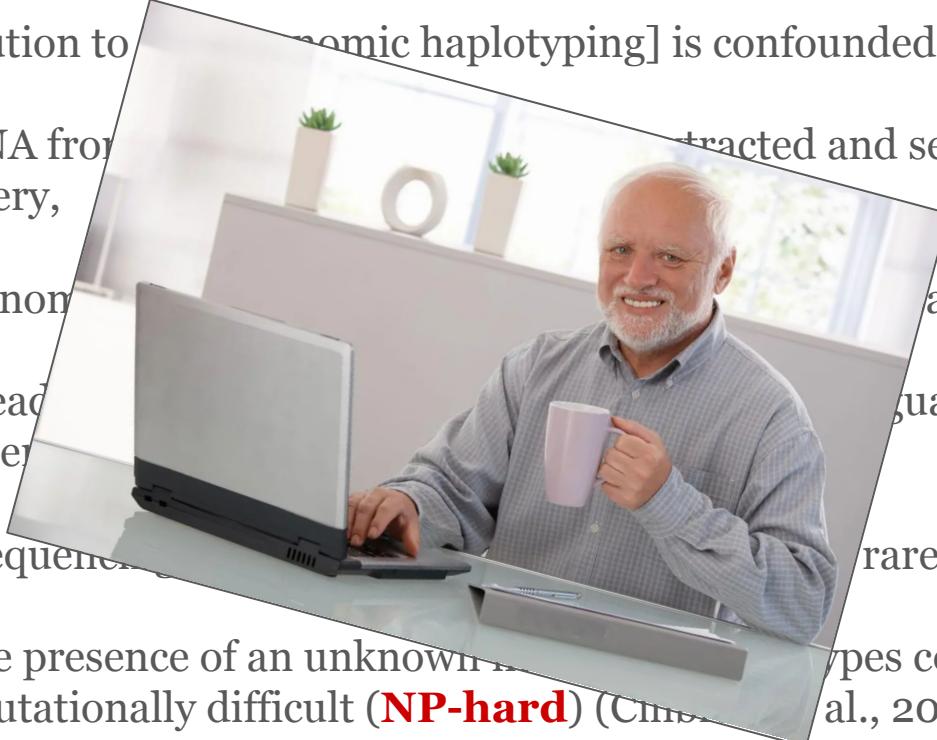
- (i) DNA from every genome needs to be extracted and sequenced to a depth sufficient for recovery,
So increase sequencing depth!
- (ii) genomes share homologous regions that require disambiguation,
So increase read lengths!
- (iii) reads may be of an insufficient length to disambiguate repeats or resolve bridges between variants,
So increase read lengths (again)!
- (iv) sequencing error can be indistinguishable from rare haplotypes and
So decrease error rates!
- (v) the presence of an unknown number of haplotypes complicates the already computationally difficult (**NP-hard**) (Cilibrasi et al., 2005) problem of haplotyping.
So use better algorithms?



Future work: “metagenomic haplotyping”

A solution to [metagenomic haplotyping] is confounded by five problems:

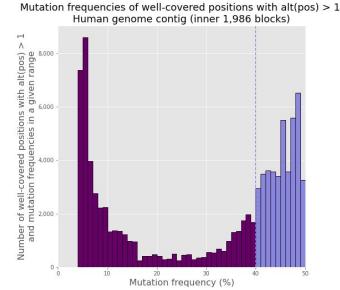
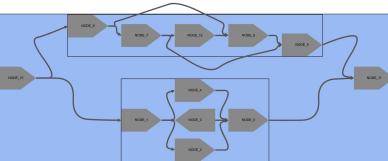
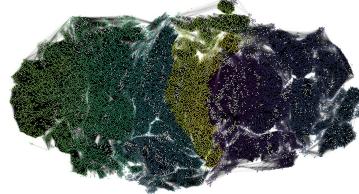
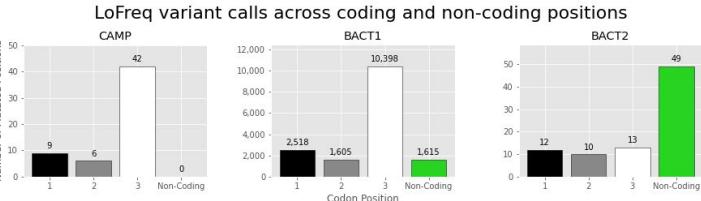
- (i) DNA from multiple sources must be extracted and sequenced to a depth sufficient for haplotype reconstruction.
So increase sequencing depth!
- (ii) genome assembly must resolve sequencing ambiguities.
So increase read lengths!
- (iii) reads from different individuals must be aligned to correctly delineate repeats or resolve bridges.
So increase read lengths (again)!
- (iv) sequencing errors must be corrected.
rare haplotypes and
So decrease error rates!
- (v) the presence of an unknown number of rare haplotypes complicates the already computationally difficult (**NP-hard**) (Clibanari et al., 2005) problem of haplotyping.



So use better algorithms?

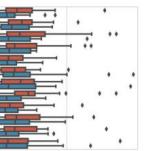
Future work: next steps forward

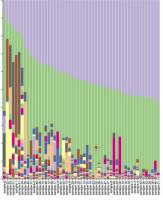
- Improving the detection (and phasing!) of rare mutations in HiFi metagenome sequencing data
 - with Misha Kolmogorov (UCSC) and Pavel Pevzner (UCSD)
- Improved methods for visualizing metagenome assembly graphs
 - with Jay Ghurye (Verily Life Sciences), Todd Treangen (Rice), Misha Kolmogorov (UCSC), Pavel Pevzner (UCSD), Jacquelyn Michaelis + Harihara Muralidharan + Mihai Pop (Maryland)
- Improved classification of prokaryotic/eukaryotic contigs
 - with Misha Kolmogorov (UCSC) and Pavel Pevzner (UCSD)



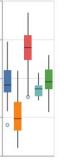
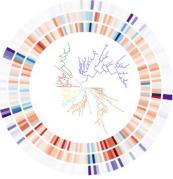
Some things I've been up to at UCSD

- Contributing to various analyses using marker gene / metagenome sequencing

- 
 - **"Co-assemblies" of metagenome sequencing data:** Sanders JG, Nurk S, Salido RA, Minich J, Xu ZZ, Martino C, **Fedarko M**, Arthur TD, Chen F, Boland BS, Humphrey GC, Brennan C, Sanders K, Gaffney J, Jepsen K, Khosroheidari M, Green C, Liyange M, Dang JW, Phelan VV, Quinn RA, Bankevich A, Chang JT, Rana TM, Conrad DJ, Sandborn WJ, Smarr L, Dorrestein PC, Pevzner PA, and Knight R (2019). "Optimizing sequencing protocols for leaderboard metagenomics by combining long and short reads." *Genome Biology*, 20(1):226.

- 
 - **Analyzing marker gene sequencing data:** Huey SL, Jiang L, **Fedarko MW**, McDonald D, Martino C, Ali F, Russell DG, Udiqi SA, Thorat A, Thakker V, Ghugre P, Potdar RD, Chopra H, Rajagopalan K, Haas JD, Finkelstein JL, Knight R, and Mehta S (2020). "Nutrition and the Gut Microbiota in 10- to 18-Month-Old Children Living in Urban Slums of Mumbai, India." *mSphere*, 5(5):e00731-20.

- Developing visualization tools for microbiome sequencing data

- 
 - **Differential abundance:** **Fedarko MW**, Martino C, Morton JT, González A, Rahman G, Marotz CA, Minich JJ, Allen EA, and Knight R (2020). "Visualizing 'omic feature rankings and log-ratios using Qurro." *NAR Genomics and Bioinformatics*, 2(2):lqaa023.
- 
 - **Phylogenetic trees (and associated data):** Cantrell K*, **Fedarko MW***, Rahman G, McDonald D, Yang Y, Zaw T, Gonzalez A, Janssen S, Estaki M, Haiminen N, Beck KL, Zhu Q, Sayyari E, Morton JT, Armstrong G, Tripathi A, Gauglitz JM, Marotz C, Matteson NL, Martino C, Sanders JG, Carrieri AP, Song SJ, Swafford AD, Dorrestein PC, Andersen KG, Parida L, Kim H-C, Vázquez-Baeza Y, and Knight R (2021). "EMPress Enables Tree-Guided, Interactive, and Exploratory Analyses of Multi-'omic Data Sets." *mSystems*, 6(2):eo1216-20. (* = contributed equally)

Thank you!

Gary Cottrell, Vineet Bafna, Melissa Gymrek, Julie Conner

Advice/support over the years

Pavel Pevzner	Kalen Cantrell	Jon Sanders	Greg Humphrey	Alison Vrbanc
Mikhail Kolmogorov	Daniel McDonald	Anna Paola Carrieri	Celeste Allaband	Bryn Taylor
Andrey Bzikadze	Yimeng Yang	Se Jin Song	Rodolfo Salido	Jerry Kennedy
Vikram Sirupurapu	Thant Zaw	Austin Swafford	Greg Poore	Yna Villanueva
Rob Knight	Stefan Janssen	Pieter Dorresteijn	Victor Cantu	Justine Debelius
Yoshiki Vázquez-Baeza	Mehrbod Estaki	Kristian Andersen	Jeffrey Chiu	Evan Bolyen
Lisa Marotz	Niina Haiminen	Laxmi Parida	Franck Lejzerowicz	Matthew Dillon
Cameron Martino	Kristen Beck	Ho-Cheol Kim	Shi Huang	Jay Ghurye
Jamie Morton	Qiyun Zhu	Larry Smarr	Sarah Adams	Jacquelyn Michaelis
Antonio González	Erfan Sayyari	Gail Ackermann	Tomasz Kosciolek	Harihara Muralidharan
Gibraan Rahman	George Armstrong	Jeff DeReus	Zech Xu	Nidhi Shah
Jake Minich	Priya Tripathi	Michiko Souza	Charles Cowart	Brian Brubach
Eric Allen	Julia Gauglitz	Justin Shaffer	Farhana Ali	Todd Treangen
Dan Hakim	Nate Matteson	Pedro Belda-Ferre	Robert Mills	Mihai Pop

Thank you!



Gary Cottrell, Vineet Bafna, Melissa Gymrek, Julie Conner

Advice/support over the years

Pavel Pevzner	Kalen Cantrell	Jon Sanders	Greg Humphrey	Alison Vrbanc
Mikhail Kolmogorov	Daniel McDonald	Anna Paola Carrieri	Celeste Allaband	Bryn Taylor
Andrey Bzikadze	Yimeng Yang	Se Jin Song	Rodolfo Salido	Jerry Kennedy
Vikram Sirupurapu	Thant Zaw	Austin Swafford	Greg Poore	Yna Villanueva
Rob Knight	Stefan Janssen	Pieter Dorresteijn	Victor Cantu	Justine Debelius
Yoshiki Vázquez-Baeza	Mehrbod Estaki	Kristian Andersen	Jeffrey Chiu	Evan Bolyen
Lisa Marotz	Niina Haiminen	Laxmi Parida	Franck Lejzerowicz	Matthew Dillon
Cameron Martino	Kristen Beck	Ho-Cheol Kim	Shi Huang	Jay Ghurye
Jamie Morton	Qiyun Zhu	Larry Smarr	Sarah Adams	Jacquelyn Michaelis
Antonio González	Erfan Sayyari	Gail Ackermann	Tomasz Kosciolek	Harihara Muralidharan
Gibraan Rahman	George Armstrong	Jeff DeReus	Zech Xu	Nidhi Shah
Jake Minich	Priya Tripathi	Michiko Souza	Charles Cowart	Brian Brubach
Eric Allen	Julia Gauglitz	Justin Shaffer	Farhana Ali	Todd Treangen
Dan Hakim	Nate Matteson	Pedro Belda-Ferre	Robert Mills	Mihai Pop

Misc. Acknowledgements

Emojis: Google emojis from emojiipedia.org: <https://emojiipedia.org/pile-of-poo/> (removed the emoji eyes manually in GIMP), <https://emojiipedia.org/non-potable-water/>, <https://emojiipedia.org/potable-water/>, <https://emojiipedia.org/mobile-phone/>

Citation of Blaser 1992 in the context of *H. pylori* based on the Strainberry paper's introduction: Vicedomini, R., Quince, C., Darling, A. E., & Chikhi, R. (2021). Strainberry: automated strain separation in low-complexity metagenomes using long reads. *Nature Communications*, 12(1), 1-14.

Taxonomic ranks figure modified from https://en.wikipedia.org/wiki/Taxonomic_rank#/media/File:Taxonomic_Rank_Graph.svg, % Annina Breen.

E. coli phylogeny from Dunne, K. A., Chaudhuri, R. R., Rossiter, A. E., Beriotto, I., Browning, D. F., Squire, D., ... & Henderson, I. R. (2017). Sequencing a piece of history: complete genome sequence of the original Escherichia coli strain. *Microbial Genomics*, 3(3).

Stock photos of someone in a suit kicking a can down the road all from Shutterstock.com % Jim Barber (all images marked as royalty-free). Why were there five separate images of this? That's a great question. I wish my job was literally just putting on a suit and kicking cans down a road. That would be so sick. I bet it pays better than grad school. Wait, why are you reading this? Seriously, there's nothing important here. It's just links. And this text.

<https://www.shutterstock.com/search/kick+the+can>,

<https://www.shutterstock.com/image-photo/politicians-shoe-stops-dented-can-rolling-85554970>,

<https://www.shutterstock.com/image-photo/close-politicians-shoe-kicking-dented-shiny-85554970>,

<https://www.shutterstock.com/image-photo/close-shiny-dented-can-sitting-on-85554964>,

<https://www.shutterstock.com/image-photo/close-politicians-shoe-kicking-dented-shiny-85554973>,

<https://www.shutterstock.com/image-photo/politician-kicks-shiny-dented-can-down-85554967>

Harold's face:

<https://www.independent.co.uk/arts-entertainment/interviews/hide-pain-harold-meme-gif-interview-model-real-name-arato-andras-thumbs-stock-photo-a7835076.html>

Quote about functional annotation and *E. coli* distance is from C. Frioux, D. Singh, T. Koresmaros, and F. Hildebrand. From bag-of-genes to bag-of-genomes: metabolic modelling of communities in the era of metagenome-assembled genomes. Computational and Structural Biotechnology Journal, 18:1722–1734, 2020.

Jigsaw puzzle photo: From VisitIndiana.com. Also, I acknowledge that I used this figure (and the Commins figure for assembly) in a talk I gave last December.

Original PCR paper: Mullis et al. 1986. Specific Enzymatic Amplification of DNA In Vitro: The Polymerase Chain Reaction.

Source about the AB370 being the first sequencer: <https://www.hindawi.com/journals/bmri/2012/251364/>

Use of the Bambus 2 “variant” figure in the context of variant calling based on Serge’s PhD defense from 9 years ago: Koren, S. (2012). Genome Assembly: Novel Applications by Harnessing Emerging Sequencing Technologies and Graph Algorithms. http://www.sergek.umiacs.io/presentations/ThesisTalk_final.pdf. I already cited this when using the “read lengths help” figure a few slides beforehand, but I figure I might as well make that clear here. Serge’s a cool dude.

Funding

Fall 2018–Winter 2019

Standard first-year CSE department fellowship

Spring 2019–Summer 2019

Joint University Microelectronics Program (JUMP)'s
Center for Research on Intelligent Storage and Processing-in-memory (CRISP)

Fall 2019–Fall 2020

IBM Research AI, via the AI Horizons Network and
the UCSD Center for Microbiome Innovation (CMI)

Winter 2021

Teaching assistantship (CSE 282)

Spring 2021–

Pevzner Lab grants

Bonus: So how many microbes are there?

Turnbaugh 2007: “The vast majority of the **10–100 trillion microbes** in the human gastrointestinal tract live in the colon.”

Locey and Lennon 2016: “[...] we predict that Earth is home to upward of **1 trillion microbial species.**”

Willis 2016: The method used by L&L 2016 isn’t statistically admissible!

∴ **Maybe the only answer right now that won’t anger any statistician or biologist:** “a lot, I guess”

An obesity-associated gut microbiome with increased capacity for energy harvest

Peter J. Turnbaugh¹, Ruth E. Ley¹, Michael A. Mahowald¹, Vincent Magrini², Elaine R. Mardis^{1,2} & Jeffrey I. Gordon¹

“We performed microbiota transplantation experiments to test directly the notion that the ob/ob microbiota has an increased capacity to harvest energy from the diet and to determine whether increased adiposity is a transmissible trait. **Adult germ-free C57BL/ 6J mice were colonized (by gavage) with a microbiota harvested from the caecum of obese (ob/ob) or lean (1/1) donors (1 donor and 4–5 germ-free recipients per treatment group per experiment; two independent experiments).** 16S-rRNA-gene-sequence-based surveys confirmed that the ob/ob donor microbiota had a greater relative abundance of Firmicutes compared with the lean donor microbiota (Supplementary Fig. 4 and Supplementary Table 7). Furthermore, the ob/ob recipient microbiota had a significantly higher relative abundance of Firmicutes compared with the lean recipient microbiota ($P < 0.05$, two-tailed Student’s t-test). UniFrac analysis of 16S rRNA gene sequences obtained from the recipients’ caecal microbiotas revealed that they cluster according to the input donor community (Supplementary Fig. 4): that is, the initial colonizing community structure did not exhibit marked changes by the end of the two-week experiment. **There was no statistically significant difference in (1) chow consumption over the 14-day period (55.4 6 2.5 g (ob/ob) versus 54.0 6 1.2 g (1/1); caloric density of chow, 3.7 kcal g⁻¹), (2) initial body fat (2.7 6 0.2 g for both groups as measured by dual-energy X-ray absorptiometry), or (3) initial weight between the recipients of lean and obese microbiotas.** Strikingly, mice colonized with an ob/ob microbiota exhibited a significantly greater percentage increase in body fat over two weeks than mice colonized with a 1/1 microbiota (Fig. 3c; 47 6 8.3 versus 27 6 3.6 percentage increase or 1.3 6 0.2 versus 0.86 6 0.1 g fat (dual-energy X-ray absorptiometry): at 9.3 kcal g⁻¹ fat, this corresponds to a difference of 4 kcal or 2% of total calories consumed.”

Bonus: Culture-Independent methods and dark matter

“It is estimated that **>99% of microorganisms observable in nature** typically are not cultivated by using standard techniques.”

Hugenholz, P., Goebel, B. M., & Pace, N. R. (1998). Impact of culture-independent studies on the emerging phylogenetic view of bacterial diversity. *Journal of bacteriology*, 180(18), 4765-4774.

Some folks have used the term “dark matter” to refer to these so-far-uncultured microbes, but it’s not a great analogy...

Microbial Dark Matter: The mullet of microbial ecology



A. Murat Eren (Meren)

[Web](#) [Email](#) [Twitter](#) [Github](#)

The Dark Matter Metaphor in Biology

By Iddo on November 27th, 2015

“Tempting as it may be, perhaps we should calm down on the use of the term dark matter in biology. Biology is confusing, complicated, and mysterious enough without it.”

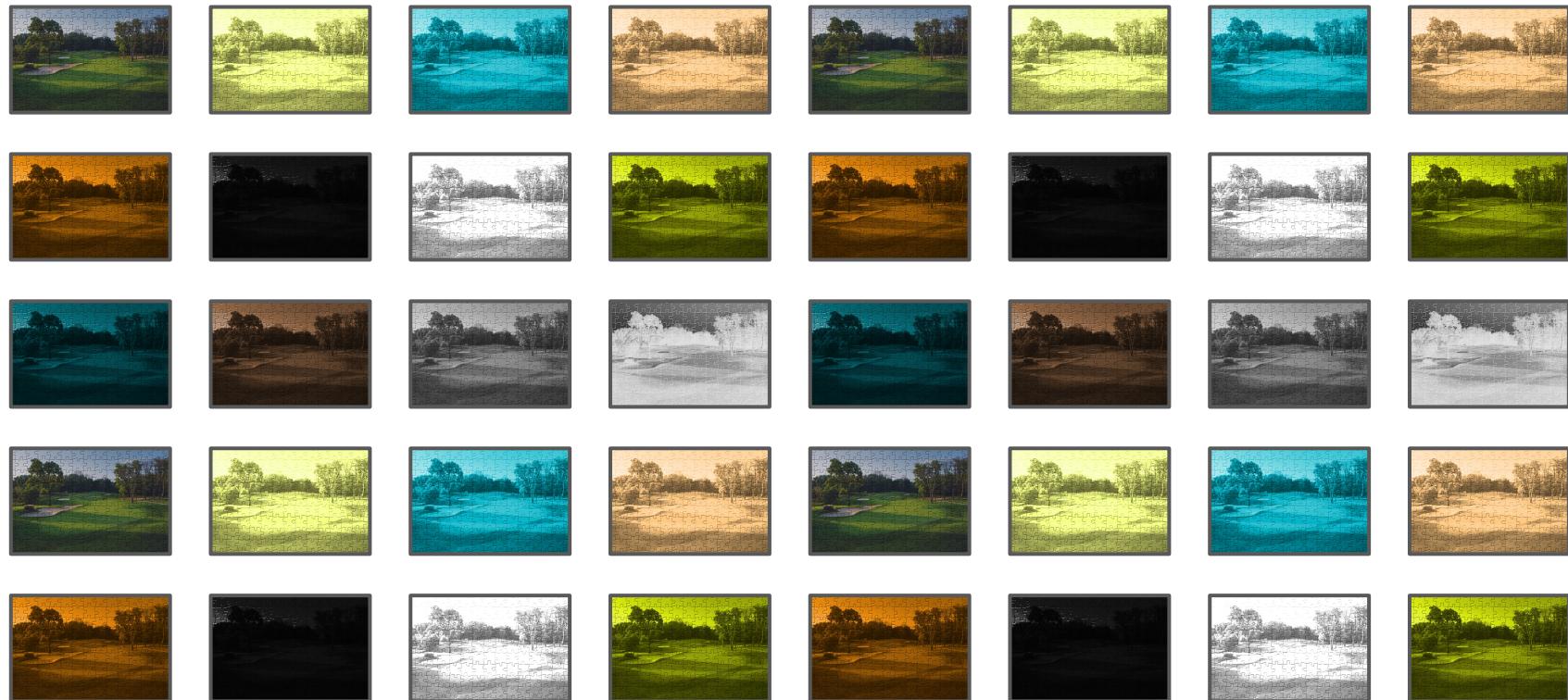
Bonus: About functional annotation...

“It says something of our ability to annotate genomes, that the proportion of a genome functionally annotated is often correlated to the genetic distance to the very well researched *Escherichia coli* (anecdotal observation).”

Bonus: Assembly (single genome)



Bonus: Assembly (metagenome)



Bonus: “metagenomic haplotyping”??????

A solution to [metagenomic] problems:

- (i) DNA from every genome in the sample must be recovered,
- (ii) genomes share homologous genes.
- (iii) reads may be of an insufficient length to distinguish between variants,
- (iv) sequencing error can be high.
- (v) the presence of an unknown number of variants makes it computationally difficult (NP-hard).



ve problems:

duced to a depth sufficient for

So increase sequencing depth!

solution,

So increase read lengths!

peats or resolve bridges

So increase read lengths (again)!

otypes and

So decrease error rates!

icates the already
problem of haplotyping.

So use better algorithms?

