# EVE 161: Microbial Phylogenomics

# Era IV: Shotgun Metagenomics

UC Davis, Winter 2016 Instructor: Jonathan Eisen • How do they connect "who" and "what"?

# Metagenomic Binning



Classification of sequences by taxa

• Method 1: align reads to reference genome



Best binning method: reference genomes



#### Fragment Recruitment





#### Figure 5. Fragment Recruitment at Sites of Rearrangements

Environmental sequences recruited near breaks in synteny have characteristic patterns of recruitment metadata. Indeed, each of five basic rearrangements (i.e., insertion, deletion, translocation, inversion, and inverted translocation) produced a distinct pattern when examining the recruitment metadata. Here, example recruitment plots for each type of rearrangement have been artificially generated. The "good" and "no mate" categories have been suppressed. In each case, breaks in synteny are marked by the presence of stacks of "missing" mate reads. The presence or absence of other categories distinguishes each type of rearrangement from the others. doi:10.1371/journal.pbio.0050077.g005



Position (kb)

Figure 7. Fragment Recruitment Plots to 20-kb Segments of SAR11-Like Contigs Show That Many SAR11 Subtypes, with Distinct Distributions, Can Be Separated by Extreme Assembly

Each segment is constructed of a unique set of GOS sequencing reads (i.e., no read was used in more than one segment). Segments are arbitrarily labeled (A=X) for reference in Figure 8. doi:10.1371/journal.pbio.0050077.g007

• Method 1: align reads to reference genome

- How would this work?
- Advantages and disadvantages?

• Method 2: Assembly

# Four main challenges for *de novo* sequencing.

- Repeats.
- Low coverage.
- Errors

These introduce breaks in the construction of contigs.

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

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

Slide from talk by Titus Brown

# More assembly issues

- Many parameters to optimize!
- Metagenomes have variation in copy number; naïve assemblers can treat this as repetitive and eliminate it.
- Assembly requires gobs of memory (4 lanes, 60m reads => ~ 150gb RAM)
- How do we evaluate assemblies?
  What's the best assembler?

Slide from talk by Titus Brown

### Wolbachia Metagenomic Sequencing







 $\mathbf{X}$ 



# Wolbachia pipientis wMel



Endosymbiosis and the Wolbachia genome Refuting Neandertal ancestry Dioxin receptor network Amphetamine and dopamine transporters

in the state of the second sec

BIOLOGY

Without Intern ULVIII of

Wu et al., 2004. Collaboration between Jonathan Eisen and Scott O'Neill (Yale, U. Queensland).

#### Pierce's Disease

UC Statewide IPM Project © 2000 Regents, University of California

#### Glassy winged sharpshooter



# GLASSY-WINGED SHARPSHOOTER

A Serious Threat to California Agriculture





FROM THE UNIVERSITY OF CALIFORNIA'S Pierce's Disease Research and Emergency Response Task Force

- Obligate xylem feeder
- Transmits *Xylella* between plants
- Much like mosquitoes transmit malarial pathogen
- Only animal listed as possible "bioterror" agent by US DHS



#### PCR and phylogenetic analysis of rRNA genes



rRNA1 5' ...TACAGTATAGGTG GAGCTAGCGATCGAT CGA... 3'

#### Phylogenetic tree



Sequence alignment = Data matrix

rRNA1	A	С	А	С	A	С
Yeast	Т	А	С	А	G	Т
E. coli	Α	G	А	С	Α	G
Humans	Т	А	Т	А	G	Т

#### Baumania is close relative of Buchnera symbionts of aphids



Wu et al. 2006 PLoS Biology 4: e188.

#### Baumania is close relative of Buchnera symbionts of aphids



Wu et al. 2006 PLoS Biology 4: e188.

# Sharpshooter Shotgun Sequencing





Collaboration with Nancy Moran's lab

Wu et al. 2006 PLoS Biology 4: e188





Amino Asid

Serve

Candidatus Sulaia musifer

Arginine

Lysine

Outempta

Aspertatio



Wu et al. 2006 <u>PLoS</u> <u>Biology 4:</u> e188. • Method 2: Assembly in More Complex Ecosystems

- How would this work?
- Advantages and disadvantages?

Fig. 2. Gene conservation among closely related Prochlorococ*cus*. The outermost concentric circle of the diagram depicts the competed genomic sequence of Prochlorococcus marinus MED4 (11). Fragments from environmental sequencing were compared to this completed Prochlorococcus genome and are shown in the inner concentric circles and were given boxed outlines. Genes for the outermost circle have been assigned psuedospectrum colors based on the position of those genes along the chromosome, where genes nearer to the start of the genome are colored in red, and genes



nearer to the end of the genome are colored in blue. Fragments from environmental sequencing were subjected to an analysis that identifies conserved gene order between those fragments and the completed *Prochlorococcus* MED4 genome. Genes on the environmental genome segments that exhibited conserved gene order are colored with the same color assignments as the *Prochlorococcus* MED4 chromosome. Colored regions on the environmental segments exhibiting color differences from the adjacent outermost concentric circle are the result of conserved gene order with other MED4 regions and probably represent chromosomal rearrangements. Genes that did not exhibit conserved gene order are colored in black.

Fig. 3. Comparison of Sargasso Sea scaffolds to Crenarchaeal clone 4B7. Predicted proteins from 4B7 scaffolds and the showing significant homology to 4B7 by tBLASTx are arrayed in positional order along the x and yaxes. Colored boxes represent **BLASTp** matches scoring at least 25% similarity and with an e value of better than 1e-5. Black vertical and horizontal lines delineate scaffold borders.



- Studying Sar86 and other marine plankton
- Note published one of first genomic studies of uncultured microbes - in 1996

JOURNAL OF BACTERIOLOGY, Feb. 1996, p. 591–599 0021-9193/96/\$04.00+0 Copyright © 1996, American Society for Microbiology Vol. 178, No. 3

#### Characterization of Uncultivated Prokaryotes: Isolation and Analysis of a 40-Kilobase-Pair Genome Fragment from a Planktonic Marine Archaeon

JEFFEREY L. STEIN,<sup>1</sup>\* TERENCE L. MARSH,<sup>2</sup> KE YING WU,<sup>3</sup> HIROAKI SHIZUYA,<sup>4</sup> AND EDWARD F. DELONG<sup>3</sup>\*



**Fig. 4.** Circular diagrams of nine complete megaplasmids. Genes encoded in the forward direction are shown in the outer concentric circle; reverse coding genes are shown in the inner concentric circle. The genes have been given role category assignment and colored accordingly: amino acid biosynthesis, violet; biosynthesis of cofactors, prosthetic groups, and carriers, light blue; cell envelope, light green; cellular processes, red; central intermediary metabolism, brown; DNA metabolism, gold; energy metabolism, light gray; fatty acid and phospholipid metabolism, magenta; protein fate and protein synthesis, pink; purines, pyrimidines, nucleosides, and nucleotides, orange; regulatory functions and signal transduction, olive; transcription, dark green; transport and binding proteins, blue-green; genes with no known homology to other proteins and genes with homology to genes with no known function, white; genes of unknown function, gray; Tick marks are placed on 10-kb intervals.

Α							
	TTREE						
в	ear loving elected, opper begins sents one						
		Linute ter					
	International Contractor Contractor Contractor Contractor Contractor	(MILE (MILE) (MILE (MILE)					
		INTERACTORY INTERACTION					
	1	(2012), 2010) (2012), 2010) (2012), 2010) (2012), 2010)					
	Philipped and the second		•				
	hornbeedronetaan adamaa kariineerin markeerinee a	Soluble Same	A				
		10.075, 00.001 10.075, 00.001 10.075, 00.001					
		CARDING, CONTRACT CARDING, CONTRACT CARDING, CONTRACT		=		=	
		(A480*20.148308*6 (A10486_A04048) (A10486_A04048)	_	_ <i>Ξ</i>		<u>·</u>	
	······································	DATES A. DOMINIC				·	
		1+4840, 1049475 1+9840, 104975	0	10000	20000	30000	40000
		(ARROY, INSIGE (ARROY, INSIGE)	<u> </u>			<u> </u>	
	ex hota 6000, and hot. 107 re-		_	Ξ	= = = = = = = = = = = = = = = = = = = =	- =	
	2001 created			-			
	mental and the second s	Default: Net			= = =		
		DAMES, DATE: NO DAMES, DATE: NO	2		—		
		(autor Canada)					
		CARDINA NO. LOCKED IN CARDINAL INCOMES					
		10180.0100					
	1	(Artistin, Musting) (Artistin, 1874 art)					
		12-100 C. 10/7 C.					
		Canada TT, Salas of S					
		Darry, G. Darry, C.					
	and particular strained sectors and a first						
	IN pages						
	ADALINA DALIMININA DI MANAGAMPANA DALIMININA DALIM	Selecter type					
		INTERNAL DISC NO.					
		Destriction of the					
		124407123, 242210770 144407723, 242210770					
		TELEPISE INCOMES					
		Conception, and the second					
		(LA 108 0. 1000 975 (LA 108 0. 1000 975					
		CARRY F. INCOME.					
		14490.6. 1849900 (1499.0. 184090.0					
		2400.00,00100.00					

#### MS 1093857: Environmental Genome Shotgun Sequencing of the Sargasso Sea Venter et al., revised



**Figure S10.** Scaffold 2217664, containing the gene encoding Proteorhodopsin. Genes are colored using color assignments described in Fig. 2, and contig boundaries are indicated with red vertical lines. In this scaffold, rhodopsin is associated with a DNA-directed RNA polymerase, sigma subunit (rpoD) originating in the CFB group.

• Method 2: Assembly in More Complex Ecosystems

- How would this work?
- Advantages and disadvantages?

• Method 3: Composition of Reads / Contigs

- How would this work?
- Advantages and disadvantages?

# PCA separates species



Gluconobacter oxydans[65% GC] and Rhodospirillum rubrum[61% GC]

### Effect of Skewed Relative Abundance





Abundance 1:1

Abundance 20:1

B. anthracis and L. monogocytes

• Method 4: Coverage of Reads / Contigs

- How would this work?
- Advantages and disadvantages?



• Method 5: Phylogenetic Analysis of Reads / Contigs

- How would this work?
- Advantages and disadvantages?



# NO PATHWAYS FOR ESSENTIAL AMINO ACID SYNTHESIS



Wu et al. 2006 <u>PLoS</u> <u>Biology 4:</u> e188.







# Binning challenge



No reference genome? What do you do?

Phylogeny ....



# CFB Phyla



### Sulcia

#### Sulcia makes vitamins and cofactors

#### Baumannia makes amino acids

3μm

### Baumannia

Wu et al. 2006 PLoS Biology 4: e188

• Method 6: Long Reads

- How would this work?
- Advantages and disadvantages?

• Method 7: Cross-linking

- How would this work?
- Advantages and disadvantages?

### HiC



From Belton JM1, McCord RP, Gibcus JH, Naumova N, Zhan Y, Dekker J. Methods. 2012 Nov;58(3):268-76. doi: 10.1016/j.ymeth.2012.05.001. Hi-C: a comprehensive technique to capture the conformation of genomes.

• Method 8: Experiments ...

- How would this work?
- Advantages and disadvantages?