EVE 161: Microbial Phylogenomics

Class #10: Genome Sequencing

UC Davis, Winter 2018 Instructor: Jonathan Eisen Teaching Assistant: Cassie Ettinger

Whole-Genome Random Sequencing and Assembly of *Haemophilus influenzae* Rd

Robert D. Fleischmann, Mark D. Adams, Owen White, Rebecca A. Clayton, Ewen F. Kirkness, Anthony R. Kerlavage, Carol J. Bult, Jean-Francois Tomb, Brian A. Dougherty, Joseph M. Merrick, Keith McKenney, Granger Sutton,
Will FitzHugh, Chris Fields,* Jeannine D. Gocayne, John Scott, Robert Shirley,
Li-Ing Liu, Anna Glodek, Jenny M. Kelley, Janice F. Weidman, Cheryl A. Phillips, Tracy Spriggs, Eva Hedblom, Matthew D. Cotton, Teresa R. Utterback,
Michael C. Hanna, David T. Nguyen, Deborah M. Saudek, Rhonda C. Brandon,
Leah D. Fine, Janice L. Fritchman, Joyce L. Fuhrmann, N. S. M. Geoghagen,
Cheryl L. Gnehm, Lisa A. McDonald, Keith V. Small, Claire M. Fraser,
Hamilton O. Smith, J. Craig Venter[‡]





Whole-Genome Random Sequencing and Assembly of *Haemophilus influenzae* Rd

Robert D. Fleischmann, Mark D. Adams, Owen White, Rebecca A. Clayton, Ewen F. Kirkness, Anthony R. Kerlavage, Carol J. Bult, Jean-Francois Tomb, Brian A. Dougherty, Joseph M. Merrick, Keith McKenney, Granger Sutton,
Will FitzHugh, Chris Fields,* Jeannine D. Gocayne, John Scott, Robert Shirley,
Li-Ing Liu, Anna Glodek, Jenny M. Kelley, Janice F. Weidman, Cheryl A. Phillips, Tracy Spriggs, Eva Hedblom, Matthew D. Cotton, Teresa R. Utterback,
Michael C. Hanna, David T. Nguyen, Deborah M. Saudek, Rhonda C. Brandon,
Leah D. Fine, Janice L. Fritchman, Joyce L. Fuhrmann, N. S. M. Geoghagen,
Cheryl L. Gnehm, Lisa A. McDonald, Keith V. Small, Claire M. Fraser,
Hamilton O. Smith, J. Craig Venter[‡]



Figure 1-17 Brock Biology of Microorganisms 11/e © 2006 Pearson Prentice Hall, Inc.



An approach for genome analysis based on sequencing and assembly of unselected pieces of DNA from the whole chromosome has been applied to obtain the complete nucleotide sequence (1,830,137 base pairs) of the genome from the bacterium *Haemophilus influenzae* Rd. This approach eliminates the need for initial mapping efforts and is therefore applicable to the vast array of microbial species for which genome maps are unavailable. The *H. influenzae* Rd genome sequence (Genome Sequence DataBase accession number L42023) represents the only complete genome sequence from a free-living organism.



What Genomes Done Before?

Tree of Life

• What Genomes Had Been Sequenced?



A prerequisite to understanding the complete biology of an organism is the determination of its entire genome sequence. Several viral and organellar genomes have been completely sequenced. Bacteriophage $\phi X174$ [5386 base pairs (bp)] was the first to be sequenced, by Fred Sanger and colleagues in 1977 (1). Sanger et al.

How Were Genomes Being Sequenced?

Tree of Life

• How Were Genomes Being Sequenced?





Homo sapiens (11). These projects, as well as viral genome sequencing, have been based primarily on the sequencing of clones usually derived from extensively mapped restriction fragments, or λ or cosmid clones. Despite advances in DNA sequencing technology

What is new here?

Tree of Life

• What did they do that was new in terms of sequencing?



The computational methods developed to create assemblies from hundreds of thousands of 300- to 500-bp complementary DNA (cDNA) sequences (13) led us to test the hypothesis that segments of DNA several megabases in size, including entire microbial chromosomes, could be sequenced rapidly, accurately, and cost-effectively by applying a shotgun sequencing strategy to whole genomes. With this strategy, a single random DNA fragment library may be prepared, and the ends of a sufficient number of randomly selected fragments may be sequenced and assembled to produce the complete genome. We chose the free-living organism Haemophilus influenzae Rd as a pilot project because its genome size (1.8 Mb) is typical among bacteria, its G+Cbase composition (38 percent) is close to that of human, and a physical clone map did not exist.

Plant

Table 1. Whole-genome sequencing strategy.

Stage	Description
Random small insert and large insert library construction	Shear genomic DNA randomly to ~2 kb and 15 to 20 kb, respectively
Library plating	Verify random nature of library and maximize random selection of small insert and large insert clones for template production
High-throughput DNA sequencing	Sequence sufficient number of sequence fragments from both ends for 6× coverage
Assembly	Assemble random sequence fragments and identify repeat regions
Gap closure	
Physical gaps	Order all contigs (fingerprints, peptide links, λ clones, PCR) and provide templates for closure
Sequence gaps	Complete the genome sequence by primer walking
Editing	Inspect the sequence visually and resolve sequence ambiguities, including frameshifts
Annotation	Identify and describe all predicted coding regions (putative identifications, starts and stops, role assignments, operons, regulatory regions)

Organisms Chosen

Tree of Life

- What is H. Influenzae?
- Why did they choose it?

complete genome. We chose the free-living organism Haemophilus influenzae Rd as a pilot project because its genome size (1.8 Mb) is typical among bacteria, its G+Cbase composition (38 percent) is close to that of human, and a physical clone map did not exist.

Sequencing Details?

Tree of Life

• Outline of sequencing approach?



Shotgun Sequencing

Tree of Life

- What is shotgun sequencing?
- How does it work?

Shotgun Sequencing

Tree of Life

- What is shotgun sequencing?
- How does it work?
- How can you complete a genome this way?

Genome sequencing. The strategy for a shotgun approach to whole genome sequencing is outlined in Table 1. The theory follows from the Lander and Waterman (14) application of the equation for the Poisson distribution. The probability that a base is not sequenced is $P_0 = e^{-m}$, where m is the sequence coverage. Thus after 1.83 Mb of sequence has been randomly generated for the H. influenzae genome (m = 1, 1)× coverage), $P_0 = e^{-1} = 0.37$ and approximately 37 percent of the genome is unsequenced. Fivefold coverage (approximately 9500 clones sequenced from both insert ends and an average sequence read length of 460 bp) yields $P_0 = e^{-5} = 0.0067$, or 0.67 percent unsequenced. If L is genome length and n is the number of random sequence segments done, the total gap length is Le^{-m} , and the average gap size is L/n. Fivefold coverage would leave about 128 gaps averaging about 100 bp in size.

Lander Waterman

Tree of Life

$P_0 = e^{-m}$ $P_0 = probability a base is$

not sequenced

m = coverage

$P_0 = e^{-m}$ $P_0 = probability a base is$

not sequenced

MATTER?

WHY DOES RANDOMNESS

Show Curve

Lander Waterman Curve



Other issues

Tree of Life

- Why do they talk about ESTS?
- Why do they need a database?

Assembly

TIGR ASSEMBLER is the software component that enabled us to assemble the H. influenzae genome. It simultaneously clusters and assembles fragments of the genome. In order to obtain the speed necessary to assemble more than 10⁴ fragments, the algorithm builds a table of all 10-bp oligonucleotide subsequences to generate a list of potential sequence fragment overlaps. When TIGR ASSEMBLER is used, a single fragment begins the initial contig; to extend the contig, a candidate fragment is chosen with the best overlap based on oligonucleotide content. The current contig and candidate fragment are aligned by a modified version of the Smith-Waterman (23) algoAssembly

Tree of Life

rithm, which provides for optimal gapped alignments. The contig is extended by the fragment only if strict criteria for the quality of the match are met. The match criteria include the minimum length of overlap, the maximum length of an unmatched end, and the minimum percentage match. The algorithm automatically lowers these criteria in regions of minimal coverage and raises them in regions with a possible repetitive element. The number of potential overlaps for each fragment determines which fragments are likely to fall into repetitive elements. Fragments representing the boundaries of repetitive elements and potentially chimeric fragments are often rejected on the basis of partial mismatches at the needs of alignments and excluded from the contig.

Assembly Method

Tree of Life

- TIGR Assembler
- Smith Waterman
- Paired Ends



Assembly Results

• Contigs Sequencing Gaps Physical Gaps

Gap Filling



Gap Filling

Tree of Life

- DNA Hybridization
- Peptide Links
- Phage lambda libraries
- PCR

Accuracy Checks?



Accuracy Checks?

Tree of Life

- Frameshift in protein sequences vs. Homologs
- Coverage > 1x
- Few ambiguities
- Comparison to H. Influenza sequence in DBs

Cost

• 0.48\$ / finished bp

Annotation

Annotation

Table 2. Summary of features of whole-genome sequencing of *H. influenzae* Rd.

Description	Number	
Double-stranded templates	19,687	
Forward-sequencing reactions (M13-21 primer)	19,346	
Successful (%)	16,240 (84)	
Average edited read length (bp)	485	
Reverse sequencing reactions (M13RP1 primer)	9,297	
Successful (%)	7,744 (83)	
Average edited read length (bp)	444	
Sequence fragments in random assembly	24,304	
Total base pairs	11,631,485	
Contigs	140	
Physical gap closure	42	
PCR	37	
Southern analysis	15	
λ clones	23	
Peptide links	2	
Terminator sequencing reactions*	3,530	
Successful (%)	2,404 (68	
Average edited read length (bp)	375	
Genome size (bp)	1,830,137	
G+C content (%)	38	
rRNA operons	6	
rrnA, rrnC, rrnD (spacer region) (bp)	723	
rrnB, rrnE, rrnF (spacer region) (bp)	478	
tRNA genes identified	54	
Number of predicted coding regions	1,743	
Unassigned role (%)	736 (42	
No database match	389	
Match hypothetical proteins	347	
Assigned role (%)	1,007 (58	
Amino acid metabolism	68 (6.8	
Biosynthesis of cofactors, prosthetic groups, and carriers	54 (5.4	
Cell envelope	84 (8.3	
Cellular processes	53 (5.3	
Central intermediary metabolism	30 (3.0	
Energy metabolism	105 (10.4	
Fatty acid and phospholipid metabolism	25 (2.5	
Purines, pyrimidines, nucleosides and nucleotides	53 (5.3	
Regulatory functions	64 (6.3	
Replication	87 (8.6	
Transcription	27 (2.7	
Translation	141 (14.0	
Transport and binding proteins	123 (12.2	
Other	93 (9.2	

*Includes gap closure, walks on rRNA repeats, random end-sequencing of λ clones for assembly confirmation, and alternative reactions for ambiguity resolution.

Euryarchaea

Plant

Eukaryotes /

HI0077 HI0080 HI0083 Oys Lys HI0075 and HI0076 tess HI0079 HIS HI0083 Oys Lys 3 HI0078 cyss HI0078 trad HI0086 melb HI0087 thrC HI0089 thrA HI0078 cyss HI0086 ddh HI0088 thrB	IIC099 sful IIC095 sful IIC095 sful IIC095 sful IIC095 sful IIC095 sful IIC095 store IIC095 sec2 IIC092 IIC092 dagE IIC0104 htpp IIC095 sec2 IIC095 store IIC095 store IIC09	HIGHEN HI	E10115 fina Len E10113 dot 150,000 nt Net E10113 tiols data HT0112 Q/V F HT0125 V/V F HT0125 V/V F Net E10114 chas HT0120 Q/V F HT0125 V/V F HT0125 V/V F Ité HT0120 Q/V F HT0126 HT0127 HT0126 HT0127 Ité HT0120 Q/V F HT0126 HT0127 HT0126 HT0127 Ité HT0120 Q/V F HT0126 HT0127 HT0126 HT0127
NO213 IO214 pelo HIO215 hadd HIO216 hadd HIO218 prrD Qlu HIO219 HIO219 HIO220 arcB rHMLF 5s-33s-16s	HI0221 guns HI0222 guns HI0224 lrp HI0226 brug HI0229 pap HI0230 HI0231 danb HI0223 rard HI0225 mbas HI0227 HI0228	R10213 R10233 R10233 R10237 R10242 R10236 RFC R10239 sec7 R10246 tet R10246 tet R10247 igs1 R10249 w R10216 RFC R10239 sec7 R10246 sec R10246 tet R10246 t	300,000 nt E10356 s.cb E10357 Ser E10351 opox E10353 hous E10357 Revo E10357 Revo FX E10351 toob E10353 Ser E10358 E10359 E10358 E10359 E10358 E10359 E10358 E10359 FX E10351 toob E10353 hous E10358 E10359 E10358 E10359 E10358 FX E10353 toob E10354 houp E10359 E10359 E10359
HI0355 tenk HI0356 tenk HI0356 tenk HI0366 tenk HI0376 tenk HI0376 tenk HI0376 tenk HI0377 tended tenk HI0377	LOYS LOJSO TAG	XI0335 XI0402 dati XI0402 mill XI0403 mill XI0403 mill XI0335 XI0337 maak XI0395 (or XI0401 mega) XI0403 mill XI0403 mill XI0403 mill XI0403 mill XI0404 mill X	Notes Notes Notes 450,000 nt Notes Notes Notes Notes Notes Notes Notes Notes Notes Notes Notes
HIOSDA LacIDM HIOSDA FPAC HIOSDS FPAB HIOSDS FPAL HIOSDA Gamma HIOSDA HIOSDS FPAB HIOSDS FPAL HIOSDA GAMMA HIOSDA HIOSDA FPALL CONTACTIONS	EIC0522 dawd EIC0527 fdw EIC0511 pp621 EIC0513 ppc0 523 HIC0524 fbm EIC0526 type EIC0526 type EIC0526 type 523 HIC0524 fbm EIC0526 type EIC0526 type EIC0526 type EIC0526 type 523 HIC0526 fbm EIC0526 type EIC0526 type EIC0526 type EIC0526 type EIC0526 type	EX0543 maph EX0553 EX0543 maph EX0554 EX0545 maph EX0555 EX0555 maph EX0555 maph E	Eitöff sama Eitöff gres
RIO643 PLAD Ser Arg Arg MIT RIO640 rpl.10 Arg HIO645 bisc RIO644 torc RIO646 and RIO645 min65 RIO650 RIO653 RIO645 plan RIO647 RIO645 rep RIO651 Ant	E10664 cyde F E10653 STORE E10656 Etof5 E10659 E10663 E10663 E10665 E1065 E10665 E10665 E1065 E10665 E10665 E10665 E10665 E10665 E1065	18469 MLGC 646 TLGC70 MLG675 Pappo MLG675 Papto MLG675 P	TIGGES and RIGES alar TIGGES RIGES AL
III0769 fts# III0769 fts# III0769 fts# III0769 fts# III0769 fts# III0769 fts# III0775	11/757 ppik 11/755 mpik 11/7555 mpik 11/75555 mpik 11/755555 mpik 11/755555 mpik 11/755555 mpik 11/755555 mpik 11/755555 mpik 11/75555555 mpik 11/755555555555555555555555555555555555	NICOLI aglis NICOL	EIG433 EIG43 EIG43 <t< td=""></t<>
R10933 holA 922 r1gs R10934 g1ys R10923 R10927 g1yg R10939 R10930 R10932 eeo R10933 R10933 R10933 R10933 R10934 g1ys R10925 R10927 g1yg R10930 R10930 R10932 eeo R10933 R10933 R10934 g1ys R10936 R10931 R10931 R10931 R10931 R10935 R10933 R10935 R10935 R10933 R10935 R10955 R1095 R109	NICOJAS NICOJAS <t< th=""><th>K dat 300556 K10955 E00577 erp K10955 E00577 erp</th><th>RIG974 RIG975 prak 1,050,00 nt NA RIG975 RIG976 RIG977 RIG975 RIG976 RIG977 RIG975 RIG976 RIG977 NA RIG975 RIG976 RIG977 RIG975 RIG976 RIG977 RIG975 RIG976 RIG977 NA RIG976 RIG977 RIG977 RIG976 RIG977 RIG977 NA RIG977 RIG976 RIG977 RIG977 RIG977 RIG977 RIG977 RIG977 RIG977</th></t<>	K dat 300556 K10955 E00577 erp K10955 E00577 erp	RIG974 RIG975 prak 1,050,00 nt NA RIG975 RIG976 RIG977 RIG975 RIG976 RIG977 RIG975 RIG976 RIG977 NA RIG975 RIG976 RIG977 RIG975 RIG976 RIG977 RIG975 RIG976 RIG977 NA RIG976 RIG977 RIG977 RIG976 RIG977 RIG977 NA RIG977 RIG976 RIG977 RIG977 RIG977 RIG977 RIG977 RIG977 RIG977
ALA HI1077 pyr0 1066 mrfD H11068 mrfB H11070 H11072 H11074 H11075 cyd3 H11078 sfpo H1 H11067 mrfC H11069 mrfA H11078 tr1075 cyd3 H11078 cyd3 H11079 cyd3 H11079 cyd3	HIG90 HIG91 HIG92 HIG93	RIILDS open RIILDS open RIILS RIILDS open RIILS RII	1,200,00 nt HIL156 400C HIL153 1000 T HIL155 HIL153 GOOD HIL153 HIL153 GOOD HIL153 HIL153 GOOD HIL153 HIL153 GOOD HIL154 GOOD HIL155 HIL155 GOOD HIL155 GOOD HIL154 GOOD HIL155 HIL155 GOOD HIL155 GO
HI1214 rec7 HI1220 rp51 HI1222 HI1223 HI1225 HI1213 spra HI1225 HI1214 sprr HI1219 cak HI1221 Has HI1224 sprr 12 prf8 HI1217 tbp1 HI1218 letp HI1224 dnak H HI1227 U	H1236 upp H1230 apt H1232 scor H1233 acor H1235 faar H1226 upp H1230 apt H1232 scor H1233 acor H1235 H1235 proj rak H1229 daax H1231 lpda H1235	D HI351 HI354 HI354 M12211 HI325 HI326 HI326<	HI1258 MI1258 mfd HI1261 fold S MI1257 MI1258 mfd HI1261 fold HI1269 mcob HI1261 met2 HI1267 for HI1269 mi1271 for HI1260 m
HIIJSS malq HIIJSS malq HIIJSS malk HIIJSS glgs HIIJSS	RIIJ63 pata RIIJ65 pata RIIJ65 pata RIIJ66 RIIJ66 R	UIU77 gefes HI1376 HI1376 HI1377 HI1377 HI1377 docts HI1376 HI1376 HI1377 docts HI1376 HI1376 HI1377 docts UIU77 gefes HI1376 HI1377 HI1377 HI1377 docts HI1376 HI1376 HI1377 docts HI1376 HI1376 HI1377 docts UU77 gefes HI1376 HI1376 HI1377 HI1377 docts HI1376 HI1376 HI1376 HI1377 docts HI1376 HI13776 HI1376 HI13766 HI13766 HI1376 HI1376 HI1376 HI1376 HI1376 HI1376 HI1376 HI13	2 reak HI1397 tags HI1398 tago HI1398 bago HI1398 tago
HI3508 HI3508 HI3508 HI350 HI3512 HI3514 HI3516 HI3520 HI3522 HI3509 HI3513 HI353 HI353 HI HI353 HI3524 HI3509 HI3513 HI353 HI353 HI353 HI3537 HI353 HI3521 HI3525 HOAS HI3525 HOAS HI3526 HI3527 HI3528 HI3527 HI3528 HI3527 HI3528 HI3527 HI3528 HI3528 HI3527 HI3528 HI358 HI3588 HI358 HI3588 HI358 HI3588 HI35888 HI3588 HI3588 HI3588 HI3588 HI3588 HI358	HIIS20 part HIIS20 part HIIS20 part HIIS30 gits HIIS32 grx HIIS35 HIIS37 Hea HIIS30 lies H	CO ELISAS	HI1546 HI1574 daab HI1576 ppt HI1570 1gtb 1,650,000 nt HI1576 taab HI1575 taak val HI1575 ppA HI1575 alx HI1577 1gtb 1,650,000 nt HI1575 tab HI1575 taak val HI1575 alx HI1577 1gtb 1,650,000 nt HI1575 HI1565 HI1572 reb HI1572 alx HI1577 1gtb 1,550,000 nt HI1575 HI1565 HI1572 reb HI1572 alx HI1577 1gtb 1,550,000 nt HI1575 HI1565 HI1572 reb H
#11673 #11673 #11673 #11673 #11663 #11665 #11677 #11689 #11673 #11676 #11662 #11676 #11676 #11670 #11670 #11675 #11662 #11676 #11670 #11670 #11675 #11675 #11672 #11676 #11670 #11670 #11675 #11675	NII694 KII695 rafk KII695 11678 kpeF _{RII679} RII662 sola KII663 tola KII667 KII688 KII699 sch KII660 KII691 KII693 and KII693 sch	HI1694 HI1701 HI1705 pepA 1693 HI1695 HI1697 HI1698 HI1700 HI1702 metz HI1703 HI1706 betr HI1707 HI1696 HI1699 HI1702 HI1702 HI1704 HI1704 HI1704 HI1696 HI1699 HI1702 HI1704 HI1704 HI1704	1,800,000 nt E13735 Gåy E11726 E11728 E11728 E11728 E11728 E11728 E11728 E11728 E11728 E11728 E11728 E2220 E11728 E11728 E11728 E11728 E11728 E11726 E11728

Euryarchaea

Fig. 1. A circular representation of the *H. influenzae* Rd chromosome illustrating the location of each predicted coding region containing a database match as well as selected global features of the genome. Outer perimeter: The location of the unique Not I restriction site (designated as nucleotide 1), the Rsr II sites, and the Sma I sites. Outer concentric circle: Coding regions for which a gene identification was made. Each coding region location is classified as to role according to the color code in Fig. 2. Second concentric circle: Regions of high G+C content (>42 percent, red; >40 percent, blue) and high A+T content (>66 percent, black; >64 percent, green). Third concentric circle: Coverage by λ clones (blue). More than 300 λ clones were sequenced from each end to confirm the overall structure of the genome and identify the six ribosomal operons. Fourth concentric circle: The locations of the six ribosomal operons (green), the tRNAs (black) and the cryptic mu-like prophage (blue). Fifth concentric circle: Simple tandem repeats. The locations of the following repeats are shown: CTGGCT, GTCT, ATT, AATGGC, TTGA, TTGG, TTTA, TTATC ,TGAC, TCGTC, AACC, TTGC, CAAT, CCAA. The putative origin of replication is illustrated by the outward pointing arrows (green) originating near base 603,000. Two potential termination sequences are shown near the opposite midpoint of the circle (red).

Sim Length ld Identification Best Location match* (%) (%) (bp) number Sensors 63.9 200 39.5 HI0220 239,378 arcB 38.1 68.0 562 HI0267 299,541 narQ 250 27.7 51.5 basS HI1707 1,781,143 280 38.1 61.6 1,475,017 phoR HI1378 Regulators 59.3 77.0 209 narP 777,934 HI0726 73.0 229 51.9 HI0837 887,011 cpxR 236 77.2 87.8 936,624 arcA HI0884 52.9 71.4 228 phoB 1,475,502 HI1379 219 43.5 59.3 HI1708 1,781,799 basR

Table 4. Two-component systems in H. influenzae Rd. ID, identity; Sim, similarity.

*In all cases, the best match was to a gene of E. coli.

Haemophilus influenzae type b

Fig. 3. A comparison of the region of the *H. influenzae* chromosome containing the eight genes of the fimbrial gene cluster present in *H. influenzae* type b and the same region in *H. influenzae* Rd. The region is flanked by *pepN* and *purE* in both organisms. However, in the noninfectious Rd strain the eight genes of the fimbrial gene cluster have been excised. A 172-bp spacer region is located in this region in the Rd strain and continues to be flanked by the *pepN* and *purE* genes.

Fig. 4. Hydrophobicity analysis of five potential channel proteins. The amino acid sequences of five predicted coding regions that do not display similarity with known peptide sequences (GenBank release 87), each exhibit multiple hydrophobic domains that are characteristic of channel-forming proteins. The predicted coding region sequences were analyzed by the Kyte-Doolittle algorithm (46) (with a range of 11 residues) with the GENE-WORKS software package (Intelligenetics).

SCIENCE • VOL. 269 • 28 JULY 1995

Plant

Table 1. Whole-genome sequencing strategy.

Stage	Description
Random small insert and large insert library construction	Shear genomic DNA randomly to ~2 kb and 15 to 20 kb, respectively
Library plating	Verify random nature of library and maximize random selection of small insert and large insert clones for template production
High-throughput DNA sequencing	Sequence sufficient number of sequence fragments from both ends for 6× coverage
Assembly	Assemble random sequence fragments and identify repeat regions
Gap closure	
Physical gaps	Order all contigs (fingerprints, peptide links, λ clones, PCR) and provide templates for closure
Sequence gaps	Complete the genome sequence by primer walking
Editing	Inspect the sequence visually and resolve sequence ambiguities, including frameshifts
Annotation	Identify and describe all predicted coding regions (putative identifications, starts and stops, role assignments, operons, regulatory regions)

