

Evaluating Genome Sequencing and Assembly Strategies for Diverse Microbial Pathogens

Xinyue Liu, Qi Su, Erin Hine, Naomi Sengamalay, Lisa D. Sadzewicz, Ivette Santana-Cruz, Sushma Parankush Das, Alvaro Godinez, Mark Eppinger, Jacques Ravel, Hervé Tettelin, Emmanuel F. Mongodin, W. Florian Fricke, Claire M. Fraser, Luke J. Tallon
Institute for Genome Sciences, University of Maryland School of Medicine, Baltimore, MD

Abstract

As new high-throughput sequencing technologies are rapidly evolving, it is becoming increasingly inexpensive to generate large numbers of draft microbial genome sequences. However, the quality and completeness of these genome sequences is often highly variable and limits comparative analyses and conclusions. Selecting the most appropriate sequencing and assembly strategy is a challenge facing most large-scale microbial pathogen genome projects. Project designers must balance the competing interests of higher quality and cost for more complete genome sequences versus larger numbers of samples to enable more comprehensive comparative analyses. In order to evaluate the optimal balance of sequencing platforms and assembly strategies for a wide range of microbial species, we conducted a comprehensive study using a series of samples from five bacterial pathogens that range in genome size and %GC content. Each genome was sequenced using three complementary platforms (454 FLX, Illumina HiSeq2000, and Pacific Biosciences RS) that offer a wide range of read length, depth of coverage, and accuracy. These diverse data were assembled in multiple combinations and at varying depths of coverage using a suite of six genome assemblers. The results of this study show that optimal quality genome sequences are obtained using different strategies for each of the analyzed species, including different combinations of sequencing and assembly techniques. These conclusions will inform future large-scale microbial pathogen genome studies, leading to more efficient and improved project design and, ultimately, the availability of more comprehensive genomic data set resources to the scientific community.

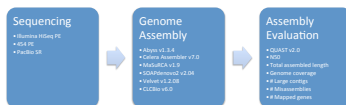
Genomes

Organism/Strain	Strain	GC%	Genome Size (Mb)	Reference
<i>Staphylococcus aureus</i>	CGS185	32.9	2.9	NC_003921.1
<i>Staphylococcus aureus</i>	CGS1410	32.9	2.9	NC_002952.2
<i>Helicobacter pylori</i>	CPH621	38.8	1.7	NC_012973.1
<i>Helicobacter pylori</i>	R037c	38.8	1.7	NC_012973.1
<i>Vibrio cholerae</i>	CPH632_5	47.5	4.0	NC_002505.1
<i>Vibrio cholerae</i>	CPH638_21	47.5	4.0	NC_002505.1
<i>Escherichia coli</i>	TW0319	50.3	5.5	NC_011933.1
<i>Escherichia coli</i>	PH0209	50.3	5.5	NC_011933.1
<i>Mycobacterium abscessus</i>	3A_0810_R	64.1	5.1	CJ458896.1
<i>Mycobacterium abscessus</i>	60_0125_R	64.1	5.1	CJ458896.1

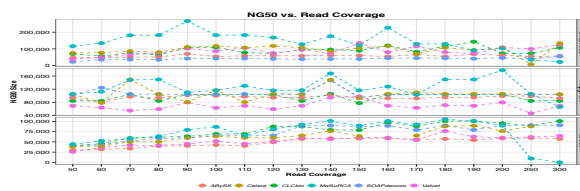
Data

Genome	Strain	Platform	Library Size	# Reads	Mean Read Length	N50
<i>H. pylori</i>	CPH621	454	2,634	113,447	36,437,744	339
<i>H. pylori</i>	287	Illumina	2,607	1,268,611	156,961,100	100
<i>H. pylori</i>	PacBio	7,224	147,490	30,283,138	2,389	6,420
<i>H. pylori</i>	R037c	454	2,947	132,225	79,383,793	342
<i>H. pylori</i>	Illumina	365	4,238,047	426,543,747	101	
<i>H. pylori</i>	PacBio	7,881	74,521	142,453,112	1,510	5,015
<i>M. abscessus</i>	3A_0810_R	454	3,377	105,634	113,635,548	363
<i>M. abscessus</i>	60_0125_R	454	3,422	5,813,087	581,309,700	100
<i>M. abscessus</i>	PacBio	6,566	393,389	265,508,977	1,500	3,854
<i>M. abscessus</i>	60_0125_R	454	2,760	340,896	115,688,196	339
<i>M. abscessus</i>	PacBio	333	2,843,005	284,100,500	100	
<i>M. abscessus</i>	PacBio	3,985	228,117	292,137,962	1,481	3,103
<i>S. aureus</i>	CGS185	454	2,734	258,335	71,883,367	285
<i>S. aureus</i>	381	Illumina	3,962	5,362,806	103,917,796	99
<i>S. aureus</i>	PacBio	9,247	117,319	277,628,745	2,366	6,133
<i>S. aureus</i>	CGS1410	454	2,455	348,656	115,963,700	332
<i>S. aureus</i>	393	Illumina	4,536,445	647,127,895	99	
<i>S. aureus</i>	PacBio	8,179	257,689	729,454,435	2,431	7,404
<i>E. coli</i>	PH0209	454	1,135	438,395	140,677,133	336
<i>E. coli</i>	TW0319	454	375	1,177,807	54,173,607	101
<i>E. coli</i>	PacBio	9,187	133,192	1,107,191,248	2,077	6,105
<i>E. coli</i>	TW0319	454	3,357	387,748	124,632,236	331
<i>E. coli</i>	Illumina	300	9,950,375	914,087,875	101	
<i>E. coli</i>	PacBio	8,058	229,432	313,780,077	2,454	6,238
<i>V. cholerae</i>	CPH632_5	454	2,655	851,750	188,946,036	305
<i>V. cholerae</i>	Illumina	346	1,968,045	196,004,500	100	
<i>V. cholerae</i>	PacBio	5,479	10,265	124,508,413	1,036	4,541
<i>V. cholerae</i>	CPH638_21	454	2,911	307,084	94,048,892	308
<i>V. cholerae</i>	Illumina	321	2,846,188	284,518,800	100	
<i>V. cholerae</i>	PacBio	5,132	288,640	443,238,191	1,176	4,183

Method

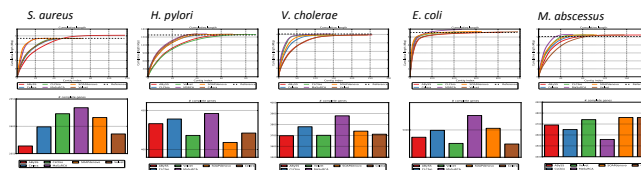


Coverage Analysis

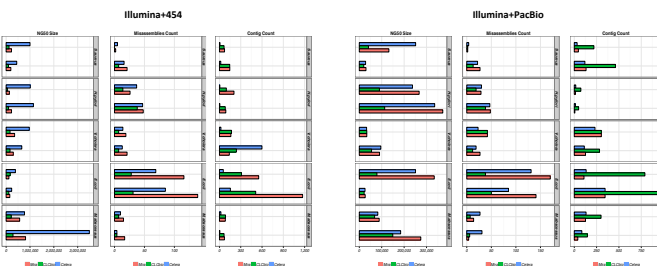
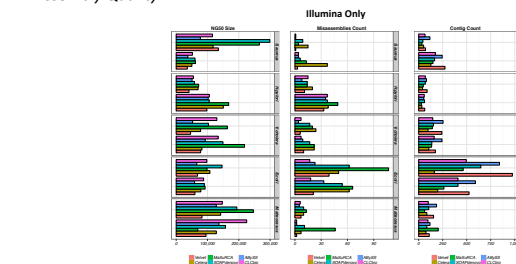


Assembler Comparison

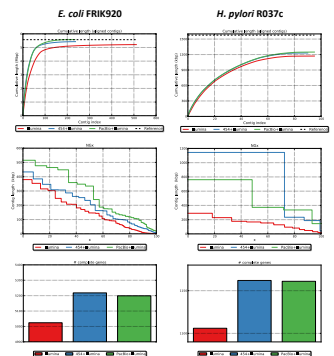
Assembly Completeness



Assembly Quality



Platform Comparison



E. coli FRK920 Reference: E. coli O157:H7 str. EC4155 5.57Mbp, 58.52% GC, 5153 genes									
Platform	Total Length	Total Contigs	GC (%)	N50	Largest Contig	Unaligned Ctg. Length	Misassembly Count	# Complete	
Illumina	5,705,734	648	58.24	146,441	379,877	397,560	62	5,024	
Illumina+454	5,783,554	154	58.20	212,022	432,132	266,599	45	5,218	
Illumina+PacBio	5,782,894	56	58.33	334,411	515,106	44,487	102	5,199	

H. pylori R037c Reference: H. pylori 838 1.58Mbp, 39.10% GC, 1382 genes									
Platform	Total Length	Total Contigs	GC (%)	N50	Largest Contig	Unaligned Ctg. Length	Misassembly Count	# Complete	
Illumina	1,622,138	28	39.18	167,666	285,435	16,380	49	1,106	
Illumina+454	1,618,432	6	39.16	1,543,918	1,543,918	4,075	47	1,163	
Illumina+PacBio	1,641,904	7	39.16	372,484	761,218	5,900	48	1,161	

Discussion

Our coverage analysis demonstrates that overlap-layout-consensus assemblers are more sensitive to read coverage than De Bruijn graph assemblers. However, optimal Illumina-only assemblies for each assembler and genome type were achieved between 100x and 200x coverage. Deeper coverage did not result in improved assembly for any of our samples. Additional coverage analysis (not shown) indicates that optimal hybrid assemblies are achieved using 20-25x coverage of either 454 or PacBio data in combination with Illumina short reads.

When measuring assembly completeness, N50, and contig count for Illumina-only assemblies, MaSuRCA and SOAPdenovo outperformed the other assemblers tested. However, MaSuRCA and SOAPdenovo also generated more assembly errors on average. ABySS and Velvet produced fewer misassemblies, but a larger number of smaller contigs.

Three of the assemblers tested were capable of hybrid assembly of Illumina data with both 454 and PacBio data. In all cases, hybrid assembly using PacBio data was preceded by error correction of the PacBio data using Illumina reads and the Celera Assembler pacBioToCA module. Celera Assembler and Mira achieved similar assembly statistics in hybrid assembly of Illumina with either 454 or PacBio data.

While high-quality draft genome assembly is possible using an Illumina-only approach, significant improvement can be achieved when combining data from multiple platforms. An Illumina+PacBio approach often achieves comparable or better results than an Illumina+454 strategy, with much lower cost and faster turnaround. Recent improvements in both PacBio sequencing and assembly methods have resulted in continued improvement in microbial genome quality and future studies will continue to evaluate these strategies.

References

Benjamin, S. and Fraser, C. (2010). "Using the metatranscriptome for reliable and automated microbial transcript assembly and SNP detection in sequenced DNA." *Genome Research* 20(10):1467-1475.
 Chen, S., Peng, Y., and Xu, J. (2010). "Hybrid error correction and de novo assembly of single-molecule sequencing reads." *Genome Research* 20(10):1476-1484.
 Li, H., et al. (2010). "The noise assembly of human genomes with massively parallel short read sequencing." *Genome Research* 20(10):1485-1493.
 Miller, J. R., Decker, A. L., et al. (2008). "Aggressive assembly of pyrosequencing reads with velvet." *Genome Research* 18(12):1818-1824.
 Quackenbush, J. (2004). "Quality Assessment Tool for genome assembly (QUAST), version 2.0." <http://www.bioinformatics.org/Quast/>.
 Quackenbush, J. (2008). "Quast: a quality assessment tool for genome assembly." *Genome Research* 18(12):1818-1824.
 Quackenbush, J. (2010). "Quast: a quality assessment tool for genome assembly." *Genome Research* 20(10):1467-1475.
 Quackenbush, J. (2011). "Quast: a quality assessment tool for genome assembly." *Genome Research* 21(10):1818-1824.

Acknowledgements

This project has been funded in whole or part with federal funds from the National Institute of Allergy and Infectious Diseases, National Institutes of Health, Department of Health and Human Services under contract number HHSN272200900003C.