Bacterial genomics

Dmitri Mavrodi

WSU Plant Pathology, dmavrodi@wsu.edu



PLP 514 • December 06, 2012

Bacterial genome features



Bacterial genome features

They are small and typically expressed in megabases $(1 \text{ Mb} = 10^6 \text{ bp})$

- Bacteria 0.8 9 Mb
- Yeast 18 Mb
- C. elegans 100 million
- Rice 420 Mb
- Arabidopsis 125 Mb
- Fugu 450 Mb
- Mouse 3,000 Mb
- Corn 2,500 Mb
- Human 3,200 Mb
- Wheat 16,000 Mb
- Loblolly pine 20,000+ Mb



Gene structure in prokaryotes and eukaryotes



Eukaryotes



Bacterial genome features



http://wishart.biology.ualberta.ca



In the 1960s and 1970s, **British scientists Frederick** Sanger and Alan Coulson, and Alan Maxam and Walter Gilbert in the US, developed DNA sequencing techniques. Gilbert and Sanger won the **1980 Nobel Prize in Chemistry** for their

Fred Sanger

DNA replication



Sanger's dideoxy sequencing technique



Cycle sequencing and fluorescent labels





Capillary sequencing





Leroy Hood



Capillary sequencers







ABI PRISM 3100 Genetic Analyzer (16 capillaries) ABI PRISM 3730xl DNA Analyzer (96 capillaries) MegaBACE 4000 DNA Analysis System (384 capillaries)

The greatest achievement...



Comparison of 1st, 2nd and 3rd generation sequencing technologies

	1st generation	2nd generation	3rd generation
Approach	Separation of labeled fragments produced by SBS	Wash-and-scan SBS	SBS or by degradation
Resolution	Many molecules are sequenced	Many molecules are sequenced	Single-molecule sequencing
Read accuracy	High	High	Moderate
Read length	Moderate (800–1000 bp)	Short	Long, 1 kb and longer
Throughput	Low	High	Moderate
Cost/base	High	Low	Low-to-moderate
Cost/run	Low	High	Low
Time	Hours	Days	Hours
Sample preparation	Moderately complex, no amplification	Complex, PCR amplification required	Ranges from complex to simple
Data analysis	Routine	Complex: large data volumes, short reads	Complex: large data volumes, signal processing

2nd generation sequencing includes **wash-and-scan** techniques

The DNA polymerase and other reagents are washed off after adding every base or an oligonucleotide; these steps are repeated many times, which consumes a lot of reagents and time









AB Applied Biosystems









a

round

mer 3

Pri

2

4

Universal seg primer (n-1)

3'11 11 11 11 11

Universal seq primer (n-2)

3 10 10 10 11 11 5 Universal seq primer (n-4)

Universal seq primer (n-3)

3'1 11 11 11

Bridge probe

Bridge probe

Bridge prob



Template 2nd base

CGT

2 3 4 A

11 March 1 Control 1 Contr Anna Anna Anna CAman Anna Anna Cana 2

st base C G

5 6 7 ... (n cycles)

Indicates positions of interrogation Ligation cycle 1 2 3 4 5 6 7

Comparison of 1st, 2nd and 3rd generation sequencing technologies

1st generation	2nd generation	3rd generation
Separation of labeled fragments produced by SBS	Wash-and-scan SBS	SBS or by degradation
Many molecules are sequenced	Many molecules are sequenced	Single-molecule sequencing
High	High	Moderate
Moderate (800–1000 bp)	Short	Long, 1 kb and longer
Low	High	Moderate
High	Low	Low-to-moderate
Low	High	Low
Hours	Days	Hours
Moderately complex, no amplification	Complex, PCR amplification required	Ranges from complex to simple
Routine	Complex: large data volumes, short reads	Complex: large data volumes, signal processing
	1st generationSeparation of labeled fragments produced by SBSMany molecules are sequencedHighModerate (800–1000 bp)LowHighLowModerately complex, no amplificationRoutine	1st generation2nd generationSeparation of labeled fragments produced by SBSWash-and-scan SBSMany molecules are sequencedMany molecules are sequencedHighHighModerate (800–1000 bp)ShortIowHighLowHighMany molecules are sequencedHighMany molecules are sequencedShortModerate (800–1000 bp)ShortHighDaysShortSampler, SamplificationMany moleculesSomplex, SamplificationShortSampler, SamplificationShortSampler, Samplification

 PacBio – a third generation platform that relies on real-time sequencing from a single DNA polymerase molecule





Ritz et al (2010) Bioinformatics



Cost of per raw Mb of DNA sequencing



Date	Cost per Mb	Cost per Genome	Date	Cost per Mb	Cost per Genome
Sep-01	\$5,292.39	\$95,263,072	Jul-07	\$495.96	\$8,927,342
Mar-02	\$3,898.64	\$70,175,437	Oct-07	\$397.09	\$7,147,571
Sep-02	\$3,413.80	\$61,448,422	Jan-08	\$102.13	\$3,063,820
Mar-03	\$2,986.20	\$53,751,684	Apr-08	\$15.03	\$1,352,982
Oct-03	\$2,230.98	\$40,157,554	Jul-08	\$8.36	\$752,080
Jan-04	\$1,598.91	\$28,780,376	Oct-08	\$3.81	\$342,502
Apr-04	\$1,135.70	\$20,442,576	Jan-09	\$2.59	\$232,735
Jul-04	\$1,107.46	\$19,934,346	Apr-09	\$1.72	\$154,714
Oct-04	\$1,028.85	\$18,519,312	Jul-09	\$1.20	\$108,065
Jan-05	\$974.16	\$17,534,970	Oct-09	\$0.78	\$70,333
Apr-05	\$897.76	\$16,159,699	Jan-10	\$0.52	\$46,774
Jul-05	\$898.90	\$16,180,224 <u></u>	Apr-10	\$0.35	\$31,512
Oct-05	\$766.73	\$13,801,124	Jul-10	\$0.35	\$31,125
Jan-06	\$699.20	\$12,585,659	Oct-10	\$0.32	\$29,092
Apr-06	\$651.81	\$11,732,535	Jan-11	\$0.23	\$20,963
Jul-06	\$636.41	\$11,455,315	Apr-11_	\$0.19	\$16,712
Oct-06	\$581.92	\$10,474,556	Jul-11	\$0.12	\$10,497
Jan-07	\$522.71	\$9,408,739	Oct-11	\$0.09	\$7,743
Apr-07	\$502.61	\$9,047,003	Jan-12	\$0.09	\$7,666



genome.gov/sequencingcosts

Increase in the number of genomes sequenced per year (1000 genomes in 2009)





Beijing Genomics Institute



Today...

As of December 2012, 17009 bacterial genome projects are listed in the Genomes Online Database



Proteobacteria 8166 (49.1%)

What do you do with all these data?

The main aim is to correlate the genome sequence of an organism with its phenotype

Genome sequences need to be annotated first searched for putative genes and analyzed for the function of the encoded proteins

The annotated genomes can be searched for genes responsible for a phenotype of interest (in the case of plant pathogens, any aspect of their interactions with plants)

What do you do with all these data?

Genome sequences can also be compared with those of closely related organisms

Transcriptomic, proteomic and molecular genetic screens can further elucidate gene function

Genome comparison data and experimental data can be dowloades in the internet databases to maximize utility to the research community

P. fluorescens group genome sequencing project

OPEN O ACCESS Freely available online

PLOS GENETICS

Comparative Genomics of Plant-Associated *Pseudomonas* spp.: Insights into Diversity and Inheritance of Traits Involved in Multitrophic Interactions

Joyce E. Loper^{1,2*}, Karl A. Hassan³, Dmitri V. Mavrodi⁴, Edward W. Davis II¹, Chee Kent Lim³, Brenda T. Shaffer¹, Liam D. H. Elbourne³, Virginia O. Stockwell², Sierra L. Hartney², Katy Breakwell³, Marcella D. Henkels¹, Sasha G. Tetu³, Lorena I. Rangel², Teresa A. Kidarsa¹, Neil L. Wilson³, Judith E. van de Mortel⁵, Chunxu Song⁵, Rachel Blumhagen¹, Diana Radune⁶, Jessica B. Hostetler⁶, Lauren M. Brinkac⁶, A. Scott Durkin⁶, Daniel A. Kluepfel⁷, W. Patrick Wechter⁸, Anne J. Anderson⁹, Young Cheol Kim¹⁰, Leland S. Pierson III¹¹, Elizabeth A. Pierson¹², Steven E. Lindow¹³, Donald Y. Kobayashi¹⁴, Jos M. Raaijmakers⁵, David M. Weller¹⁵, Linda S. Thomashow¹⁵, Andrew E. Allen¹⁶, Ian T. Paulsen³

1 Agricultural Research Service, U.S. Department of Agriculture, Corvallis, Oregon, United States of America, 2 Department of Botany and Plant Pathology, Oregon State University, Corvallis, Oregon, United States of America, 3 Department of Chemistry and Biomolecular Sciences, Macquarie University, Sydney, Australia, 4 Department of Plant Pathology, Washington State University, Pullman, Washington, United States of America, 5 Laboratory of Phytopathology, Wageningen University, Wageningen, The Netherlands, 6 The J. Craig Venter Institute, Rockville, Maryland, United States of America, 7 Agricultural Research Service, U.S. Department of Agriculture, Davis, California, United States of America, 8 Agricultural Research Service, U.S. Department of Agriculture, Charleston, South Carolina, United States of America, 9 Department of Biology, Utah State University, Logan, Utah, United States of America, 10 Institute of Environmentally-Friendly Agriculture, Chonnam National University, Gwangju, Korea, 11 Department of Plant Pathology and Microbiology, Texas A&M University, College Station, Texas, United States of America, 12 Department of Horticultural Sciences, Texas A&M University, College Station, Texas, United States of America, 13 Department of Plant and Microbial Biology, University of California Berkeley, Berkeley, California, United States of America, 14 Department of Plant Biology and Pathology, Rutgers, The State University of New Jersey, New Brunswick, New Jersey, United States of America, 15 Agricultural Research Service, U.S. Department of Agriculture, Pullman, Washington, United States of America, 16 The J. Craig Venter Institute, San Diego, California, United States of America

P. fluorescens group genome sequencing project

- Sequenced the genomes of 7 strains of the *P. fluorescens* group that colonize plant surfaces and function as biological control agents
- Compared strains to each other and to three other strains that were sequenced previously
- Looked for known and new genes that are likely to confer biocontrol traits and/or contribute to plantmicrobe interactions

P. fluorescens group genome sequencing project

Strain	Source	Target disease(s) for biological control				
P. chlororaphis s	ubsp. aureofaciens:					
30-84	Wheat rhizosphere, Washington, USA	Take-all of wheat				
06	Soil, Utah, USA	Wildfire of tobacco target spot of cucumber				
P. protegens:						
Pf-5	Soil, Texas, USA	Seedling emergence				
P. brassicacearu	m:					
Q8r1-96	Wheat rhizosphere, Washington, USA	Take-all of wheat				
P. fluorescens:						
Pf0-1	Soil, Massachusetts, USA					
Q2-87	Wheat rhizosphere, Washington, USA	Take-all of wheat				
SBW25	Sugar beet phyllosphere, Oxfordshire, UK	Seedling emergence				
A506	Pear phyllosphere, California, USA	Fire blight of pear and apple, frost injury, fruit russeting				
SS101	Wheat rhizosphere, The Netherlands	Diseases caused by Pythium spp. and Phytophthora spp.				
Pseudomonas sp).:					
BG33R	Peach rhizosphere, South Carolina, USA	The plant-parasitic nematode <i>Mesocriconema</i> xenoplax				

Steps in a typical genome sequencing project

Sample preparation and sequencing

Finishing Assembly Gap filling Conflict resolution Verification

Analysis Gene predictions Homology searches Annotation



Base calling: reading the chromatograms

RESEARCH

Base-Calling of Automated Sequencer Traces Using Phred. I. Accuracy Assessment

Brent Ewing,¹ LaDeana Hillier,² Michael C. Wendl,² and Phil Green^{1,3}

¹Department of Molecular Biotechnology, University of Washington, Seattle, Washington 98195-7730 USA; ²Genome Sequencing Center, Washington University School of Medicine, Saint Louis, Missouri 63108 USA

Genome Research (1998) 8:175-185

Developed by Phil Green (University of Washington, Seattle)

De facto industry standard for base calling

Phred assigns probability value to each base

Phred q	Probability that the base is called wrong	Accuracy of the base call
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1,000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%

Draft vs finished genomes

Finishing is the process of assembling and refining raw sequence data into a highly accurate final genomic sequence

Automated sequence editing Manual, interactive sequence inspection Directed sequencing Assembly verification

Disadvantages: cost and time

GC skew and scaffold order

%G+C composition (reflecting horizontal gene transfer) and GC skew (reflects strand-specific mutational bias)



GC skew and scaffold order



Gene synteny analyses

Finished genomes are useful for studying the preservation of the gene order (establishment of the orthology of genomic regions in different species; important functional relationships between genes; evolution)

Some whole-genome alignment tools:

MUAVE (http://gel.ahabs.wisc.edu/mauve/)

ACT (http://www.sanger.ac.uk/Software/ACT)



Plasmid gene synteny analysis

	50'00	10000	15000	20000	25000	30000	35000	40000	45000	50000	55000				
2			NA.		And the second lines		and a second	And the second sec	Chi Millionald	Manager and a second of the	wideling bild				
R'□⊂ ♥											/				
Pseu	domonas fluorescen	s A506	15000	20000	25000	30000	35000	40000	_						
			13000	2000	A Dependence of the second sec		33000								
R						1 00 0									
*															
Pseu	5000	10000	15000	20000	25000	30000	35000	40000	45000	50000					
_		nn			and Mathematic Photogram	Territory of Jacob	tal paolonya ar anto	Martin Contractor			_				
R'œ ≫				1 N H H H		1 100	11 11		0						
Pseu	domonas syringae p	v. phaseolicola 1448	A 1448A; BAA	-978	25000	30000	35000	40000	450						
		10000	13000	20000	august to the first of the	Source and the second	55000	A Alexa Materia (Mart							
R															
•															
rseu	5000	10000	15000	20000	25000	30000	35000	40000	45000						
_		<u>u</u> (and phone and a second		and down and the state of the								
R ⊨ ≫							11 11)								
Pseu	domonas syringae p	v. maculicola ES4326	15000	20000	25000	30000	35000	40000	45000	50000	55000	60000	65000	70000	75000
»	South Party		13000	20000		30000	0 0	40000	43000	30000		1 Million Day and A	Array Constants International	history history date and the	73000
R															
₩ []	U.						0.00		00						
Pseu	5000	10000	15000	20000	25000	30000	35000	40000	45000	50000	55000	60000	65000	70000	
~			1	-	M					- Dr. and the set of the	and the second balance	hand for the or a second		Contract of the	
R W			io — I											n '	
Pseu	domonas syringae p	v. syringae A2	15000	20000	25000	30000	35000	40000	45000	50000	55000	60000	65000		
~	1 Jan			20000	23000	50000	00066	10000	43000	50000	33000	00000			
R V III								סס סכ סכ סכ							
Salm	onella sp. 96A-2919 5000	10000	15000	20000	25000	3000									
~		1 and the state of the		1 and											
R						_									

Draft vs finished genomes

Sequence is experimental data, subject to experimental error

Virtually any accuracy can be obtained, but the cost of a project and the accuracy determined are directly related

It makes sense to use an accuracy appropriate for the results to be obtained

Finding the genes in a genome and their relative location does not require extremely high accuracy

The goal of annotation is to locate genes (aka coding sequences or CDSs) and assign their putative functions

Open reading frames are predicted by length of by using interpolated Markov models (Glimmer)

The proteins encoded by each putative gene are compared by BLAST searches against closely related proteomes and the NCBI databases

To improve gene function prediction, protein motifs are assigned (TIGRFAM, Pfam and COG databases)

BLAST scores and other evidence for function are compared to choose CDSs most likely to represent functional genes

Automated annotation

All analysis step can be assembled into annotation "pipeline" and performed without human intervention

Some free online automatic annotation engines:

- RAST (http://rast.nmpdr.org)
- BASYS (http://basys.ca)
- IGS Annotation Engine (http://ae.igs.umaryland.edu)

The Prokaryotic Genome Annotation Pipeline (PGAAP) (http://www.ncbi.nlm.nih.gov/genomes/static/Pipeline.html)

DOE-JGI Microbial Annotation Pipeline (https://img.jgi.doe.gov)











CloVR virtual machines



Community annotation

Automated annotation often followed by "community annotation", when a community of researchers integrates computer predictions with experimentally validated data

Some tools for integrating manual and automated annotations:

Artemis (http://www.sanger.ac.uk/resources/software/artemis)

Manatee (http://manatee.sourceforge.net/)

Apollo (http://apollo.berkeleybop.org)

Traits of potential importance for plant-microbe interactions and biological control



Mining individual genomes

Keyword searches

BLAST searches is the most common way of mining a genome (online vs local; nr vs specialized databases)

BLAST searches alone are often not sufficient

T3 effectors is a good example – we need to look at motifs (a stretch of protein with conserved amino acids in conserved positions), patterns (certain amino acids occur more frequently), regulatory elements (hrp boxes) and neighboring genes

Traits of potential importance for plant-microbe interactions and biological control (contd.)



Community annotation

Other interesting features: repeats, operons, islands

Some free online automatic annotation engines: RepeatScout (http://bix.ucsd.edu/repeatscout/) RepeatFinder (http://cbcb.umd.edu/software/RepeatFinder/)



REP elements



Anomalous regions: genomic islands and mobile elements



Genomic islands differentially distributed in the genomes of related strains

Often integrate into tRNAs and, in addition to cargo genes, carry mobility genes (integrases, insertion elements) and have atypical G+C content

Core and lineage-specific regions in the sequenced genomes

Phl⁺ strain *Pseudomonas fluorescens* Q8r1-96



- core *Pseudomonas* genes
- lineage-specific genes
- strain-specific genes
- REPa elements
- REPb elements

Cross-genome comparisons

All-against-all BLAST analyses to look at the presence/absence of genes

Helps to define core and pan genomes

Helps to find differences in gene content between organisms and identify candidate genes that may determine the observed phenotypic variations

Presence/absence of genomic regions:

MUAVE (http://gel.ahabs.wisc.edu/mauve/) ACT (http://www.sanger.ac.uk/Software/ACT)

Proteome conservation between strains from the *P. fluorescens* group



Loper et al (2012) PLoS Genetics 8: e1002784

Core and signature genome analyses

Core and signature genome comparisons of three Phz⁺ and three Phl⁺ *Pseudomonas* spp.



Conclusions I

- Thousands of bacterial genomes including those of plant pathogenic and biocontrol bacteria have been sequenced
- These genomes represent an invaluable resource for molecular plant pathologists and plant microbiologists
- Many different approaches can be used for mining bacterial genomes - one by one or through powerful comparative analyses
- There are numerous free and easy-to-use web databases and tools that can be utilized for genome mining