

Finishing Strategies, Sequence Validation and Genomic Analyses: lessons from a GC- and repeat-rich bacterial genome



Warren RL, Morin RD, Yang G, Stott JM, Butterfield YB, Shin H, Smailus D, Schein JE, Siddiqui AS, Holt R, Jones SJM, Marra MA, McLeod MP, Mohn WW, Fukuda M, Davies JE, Eltis LD

Canada's Michael Smith
Genome Sciences Centre
www.bcgsc.ca

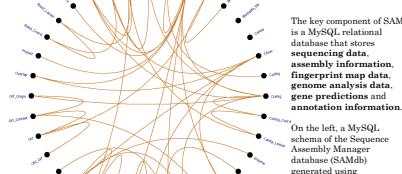
University of British Columbia
Department of Microbiology and Immunology

1. Introduction

Rhodococcus sp. RHA1 is a soil bacterium of environmental and industrial importance due to its ability to metabolize numerous organic substrates and its potential to produce secondary metabolites. Random Whole Genome Sequencing (WGS) of plasmid and fosmid clones yielded ~9.7 Mb coverage of the ~9.7 Mb genome.

In an effort to analyze the RHA1 genome, we have designed and implemented a system to manage whole genome shotgun sequences and whole genome sequence assembly (WGA) data flow. The Sequence Assembly Manager (SAM) enables prompt of a MySQL relational database and PERL applications designed to manage, analyze and coordinate the analysis of sequence information and to view and report genome assembly progress through its CGI web interface *sam.pl*.

2. SAM MySQL Relational Database

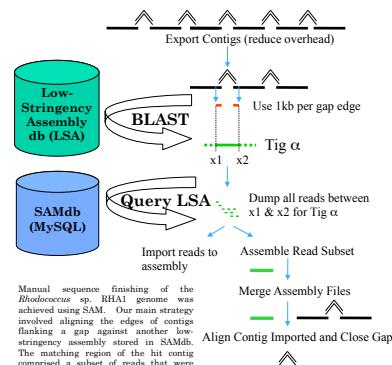


3. Visualization of WGA using sam.pl

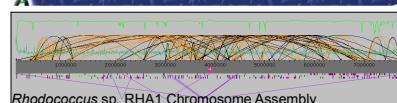
The main purpose of SAM is to organize WGA in a structured way such that users can assess WGA quality rapidly using an intuitive web interface that allows easy evaluation of the assembly quality and its progress. SAM renders WGA data generic, upon which genome analysis tools and viewers for different sequence assemblers were designed. The images shown below are generated dynamically on the web and accessed by running the PERL CGI application *sam.pl* in a web browser.



4. Strategy used to Finish the RHA1 Genome



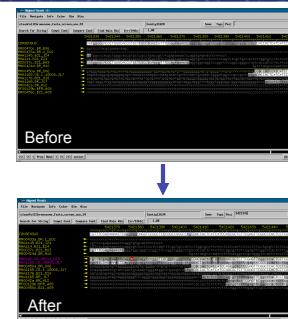
5. Resolving Repeats



The RHA1 genome contains numerous repeats that represented an enormous challenge for finishing. Most repeats were resolved during the final assembly: This was achieved using the *rhifinder* algorithm superimposed for ultimate use as mini-assembler bins for Phrap. Manual finishing is labor intensive, which is why we sought to automate the process.

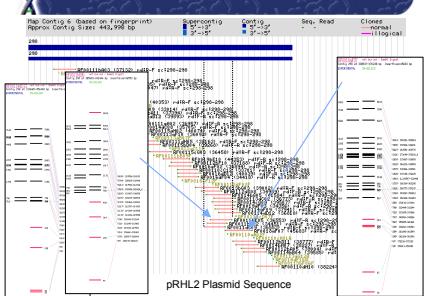
The above figure, taken from Canada's Assembly View, shows the entire 7.8 Mb RHA1 chromosome. Orange and black curves indicate direct and inverted repeats, respectively, of 0.50-2.0 kb and less than 10% mismatch.

6. Resolving Sequence Hardstops

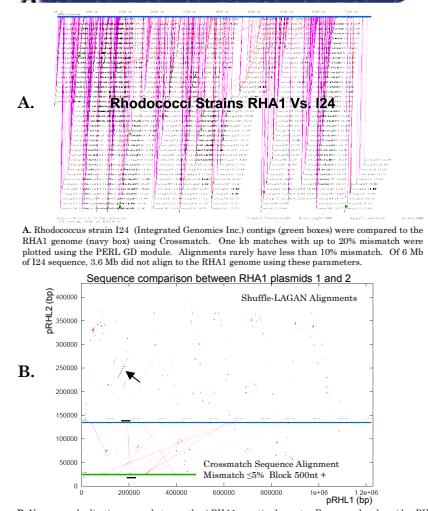


Due to the GC-rich nature of the RHA1 genome, several clones could not be assembled with standard assembly mix. Use of the 457TP kit in a 5 Mb DMSO solution helped considerably, but did not yield data for regions with unusually high GC-content (>75%). Ultimately, 72% of these recalcitrant clones were sequenced using the Amersham Sequence Finishing Kit, which was well suited. Analysis of selected traces revealed that this kit yielded long reads with no less GC content. The above screenshots from Consed show a sequence hardstop before and after merging sequences obtained using the Amersham Sequence Finishing Kit*.

7. Sequence Validation using Fingerprint Maps in SAM



8. Sequence Comparison



I thank Martin Krzywinski for helpful discussions about PERL programming & (schemaball) and Greg Taylor for his initial work on the large chromosome linear inversions.

References

Altshul SF et al. 1990. Basic local alignment search tool. *J. Mol. Biol.* 215:403-410.
*Amersham Sequence Finishing Kit: ©Amersham Biosciences Corp. (www.amsrhamsbio.com)
Batzoglou S et al. 2002. ARACHNE: a whole-genome shotgun assembler. *Genome Res.* 12:177-189.
Boetzel P et al. 2002. A new finishing strategy for bacterial genomes. *Genome Biol.* 3:0042.
Gordon D et al. 1998. Consed: A Graphical Tool for Sequence Finishing. *Genome Res.* 8:119-122.
Green P. 1996. <http://bozeman.mbt.washington.edu/phrap/docs/phrap.html>
Warren R et al. 2004. Functional Characterization of a Catabolic Plasmid from PCB-Degrading *Rhodococcus* sp. RHA1. *J. Bacteriology* (scheduled to appear in November).

Software
SAM is available for download at <http://www.bcgsc.ca/bioinfo/software/sam.pl>
Schemaball is available for download at <http://mkweb.bcgsc.ca/schemaball>