

Supplementary Materials

A regulatory hydrogenase gene cluster observed in the thioautotrophic symbiont of *Bathymodiolus* mussel in the East Pacific Rise

Ajit Kumar Patra¹, Maëva Perez², Sook-Jin Jang³ and Yong-Jin Won^{1*}

¹Division of Ecoscience, Ewha Womans University, Seoul, Republic of Korea

²Department of Biological Sciences, Université de Montréal, Montreal, Canada

³Ocean Science and Technology Institute, Inha University, Republic of Korea

*Correspondence: won@ewha.ac.kr

Symbiont genome assembly and plasmid prediction.

Filtered PacBio subreads were used for assembly by all three assemblers. De novo assembly was conducted using the hierarchical genome-assembly process (HGAP3) pipeline of the SMRT Analysis v2.3.032. After gap closing, the newly generated contig served as a reference to which raw PacBio reads were mapped using the resequencing module of the SMRT protocol. Assembly by Flye v2.9 (--pacbio-raw) and Unicycler v3.0 (-l) was used with command for PacBio reads as mentioned by the developers. Then, we used these assemblies to cluster contigs to groups similar contigs together and then reconciled manually validated contig clusters by Tricycler. Assembly statistics generated by three assemblers are presented in Table S10. Plasmids from the assemblies are predicted by PlasFlow are presented in Table S11.

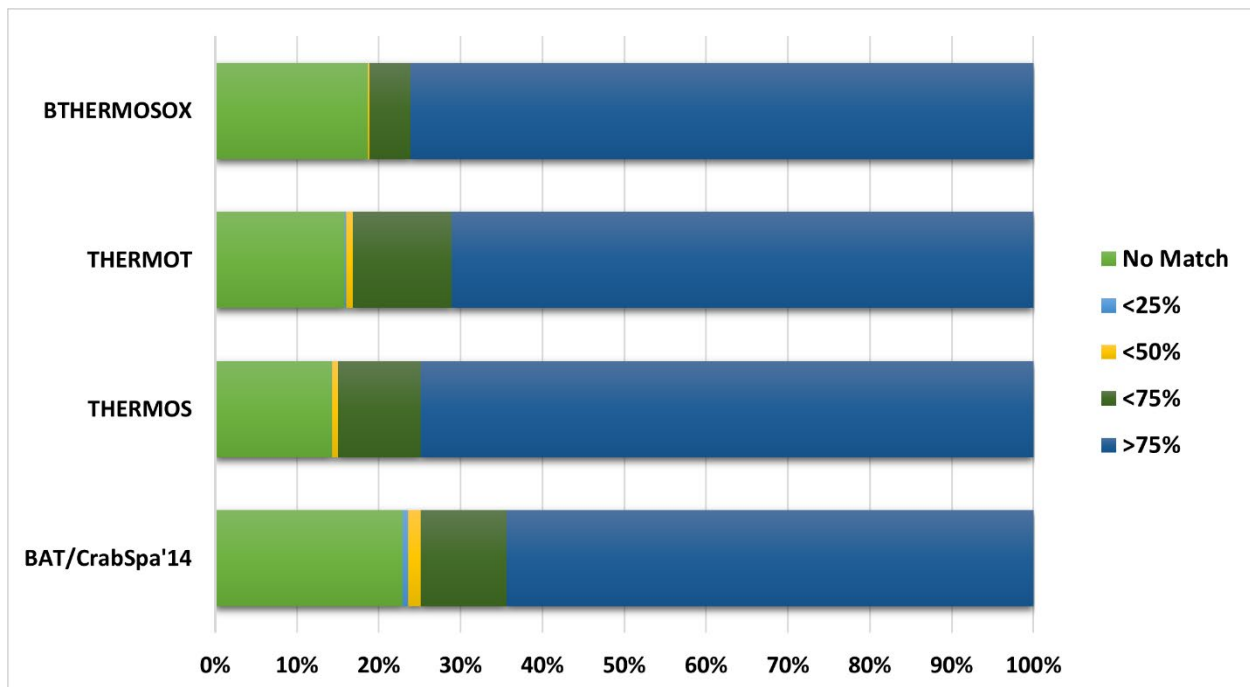


Figure S1 Sequence alignment similarity of *B. thermophilus* thiotrophic symbiont genomes to the EPR9N genome.

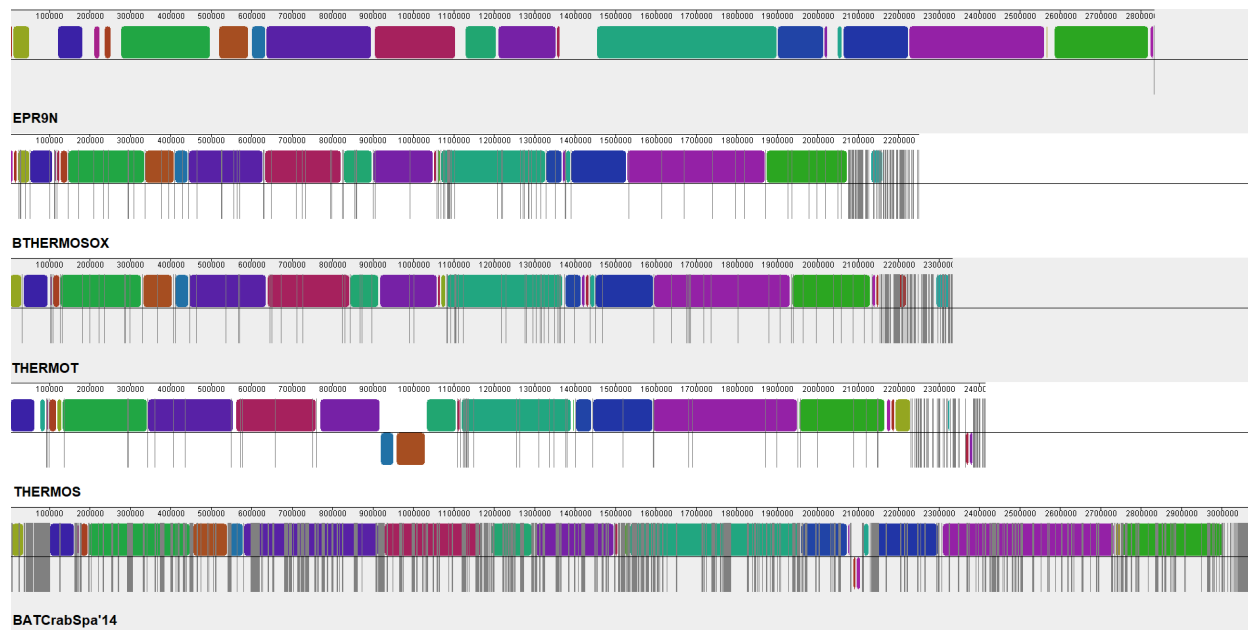
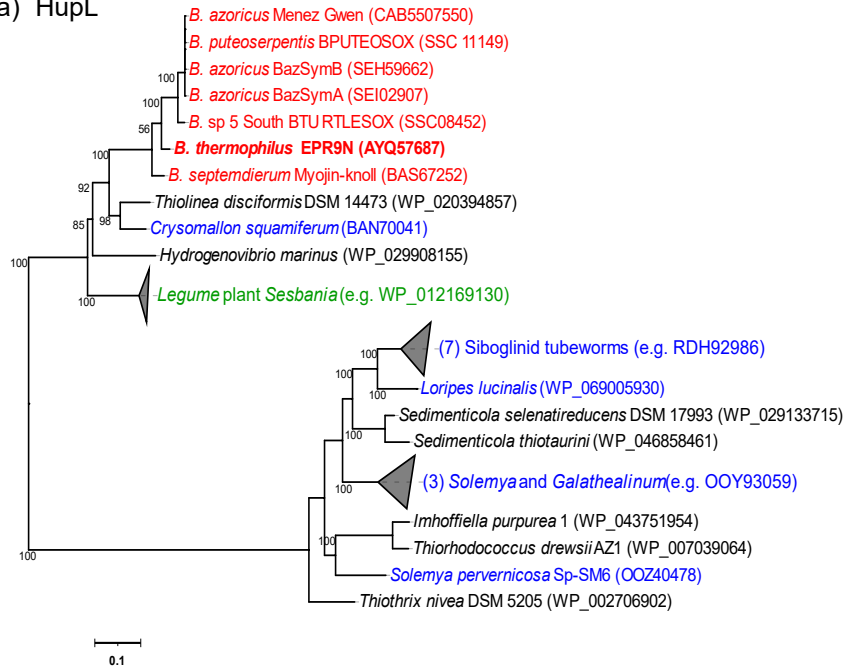


Figure S2. Structural rearrangements between *B. thermophilus* thiotrophic symbiont genome assemblies. Contigs of the fragmented assemblies were sorted and reoriented to match as best as possible the reference EPR9N assembly. Colored fragments indicate locally collinear blocks (LCBs) where sequences are homologous and show no putative structural variation. Fragments under the central horizontal line indicate inversions compared to EPR9N. Contig boundaries of the assemblies are marked with grey bars.

(a) HupL



(b) HupS

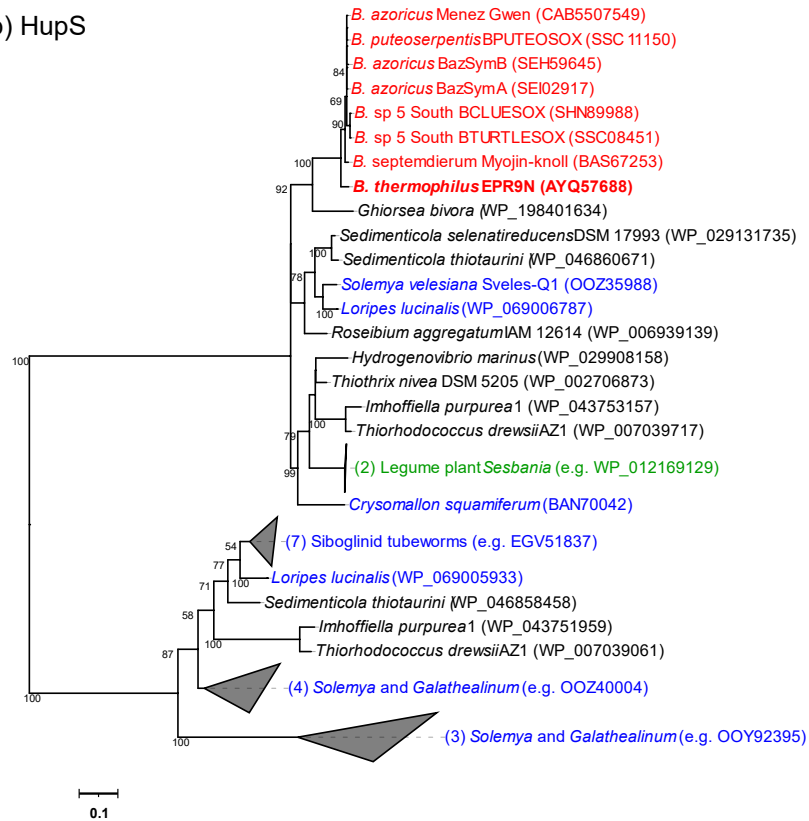


Figure S3. Phylogenies of hydrogenases (a) HupL and (b) HupS genes. Phylogenetic trees were inferred by using the maximum likelihood method with bootstrap values on the tree branches (only over 50 are shown). The taxa in green are of the symbionts of legume host *Sesbania*, in blue are the symbionts of marine invertebrates, and in bold red are the symbiont of *B. thermophilus* EPR9N. GenBank accession numbers of proteins are presented in parentheses.

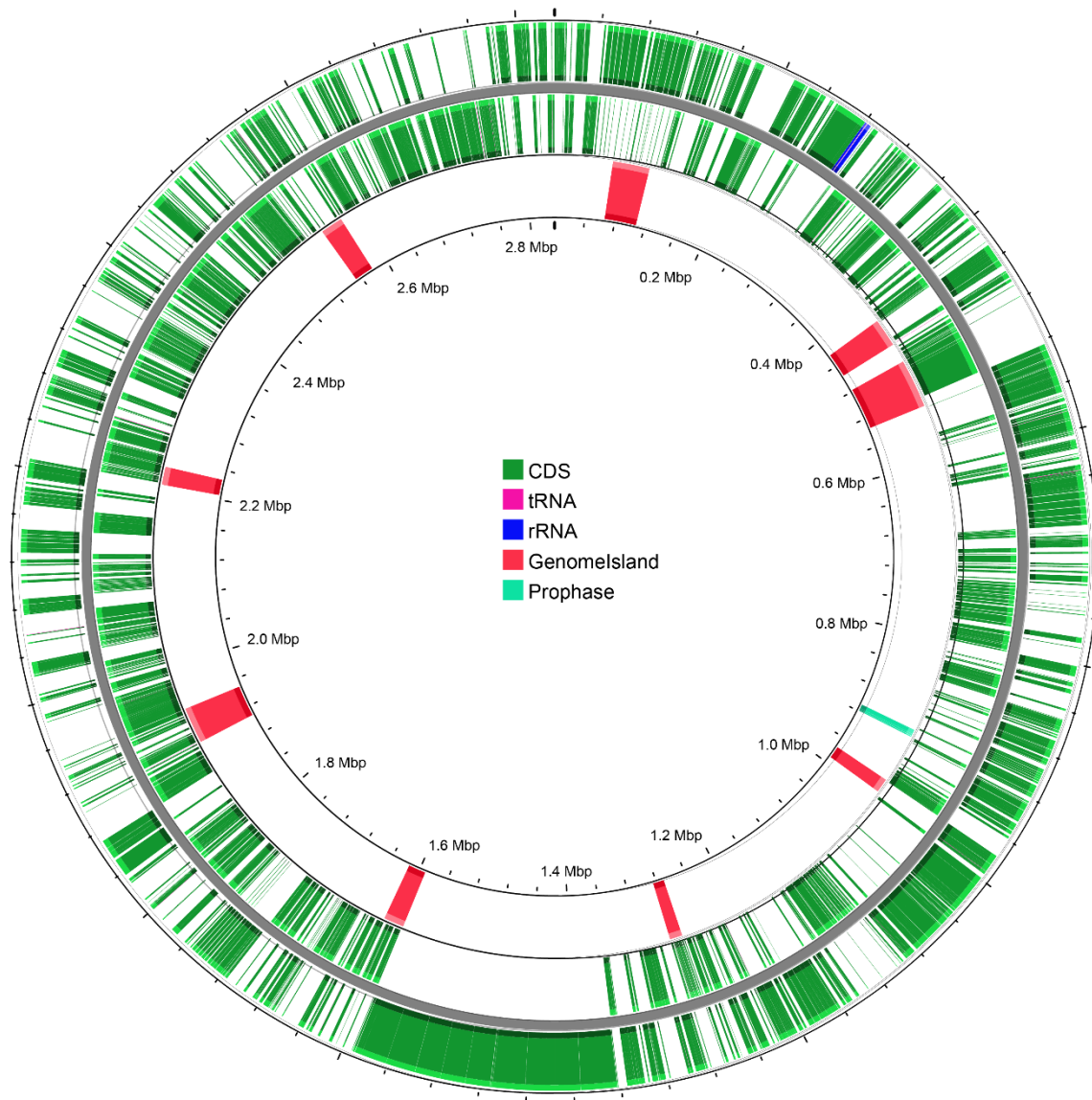


Figure S4. Genomic islands, and prophage regions in *B. thermophilus* EPR9N genome.