

**Fig. S2.** Additional information supporting the use case “Application to Comparative Genomics: Erythritol Utilization in *Brucella*”. (A) Genes encoding enzymes involved in erythritol transport and catabolism in *Brucella* spp., with emphasis on candidate pseudogenes across 16 genomes [1, 2]. (B) Whole genome-based phylogeny estimation for 107 Rhizobiales taxa. One outgroup (Sphingomonadales) is used to root the tree. Arrow points to inset showing cladogram of the 41 *Brucella* genomes. Red bar marks the origin of the *abortus* clade, and red asterisk denotes the monophyly of *B. abortus* strains S19 and NCTC 8038 (see text for details). The following pipeline was implemented to estimate Rhizobiales phylogeny: BLAT (refined BLAST algorithm) [3] searches were performed to identify similar protein sequences between all genomes, including the outgroup taxon. To predict initial homologous protein sets, mcl [4] was used to cluster BLAT results, with subsequent refinement of these sets using in-house hidden Markov models [5]. These protein families were then filtered to include only those with membership in >80% of the analyzed genomes (85 or more taxa included per protein family). Multiple sequence alignment of each protein family was performed using MUSCLE (default parameters) [6, 7], with masking of regions of poor alignment (length heterogeneous regions) done using Gblocks (default parameters) [8, 9]. All modified alignments were then concatenated into one dataset. Tree-building was performed using FastTree [10]. Support for generated lineages was estimated using a modified bootstrapping procedure, with 100 pseudoreplications sampling only half of the aligned protein sets per replication (NOTE: standard bootstrapping tends to produce inflated support values for very large alignments). Local refinements to tree topology were attempted in instances where highly supported nodes have subnodes with low support. This refinement is executed by running the entire pipeline on only those genomes represented by the node being refined (with additional sister taxa for rooting purposes). The refined subtree was then spliced back into the full tree. More information pertaining to this phylogeny pipeline is available at PATRIC (see “Phylogeny FAQs” at <http://enews.patricbrc.org/faqs/>). (C) Application of the Multiple Sequence Alignment Viewer tool to evaluate the origin and diversification of the ATP-binding (*eryE*), permease (*eryF*), and substrate-binding (*eryG*) components of erythritol ABC transporters 1 and 2. Using the ‘Protein Family Sorter’ tool, orthologous proteins for each *ery* protein (three FIGfams for EryE, five FIGfams for EryF and four FIGfams for EryG) were extracted from the interactive heatmap. Specifically, once a set of FIGfams was captured, the “show proteins” option was selected. From this table, all proteins were selected (65 for EryE, 72 for EryF and 56 for EryG). Next, the “Integrated Protein Tree and Alignment” option was selected, resulting in the display of the full length multiple sequence alignment coupled with an estimated phylogeny (the Multiple Sequence Alignment Viewer tool). For more information regarding the Multiple Sequence Alignment Viewer tool, see the “Alignment FAQs”.

1. Crasta OR, Folkerts O, Fei Z, Mane SP, Evans C, Martino-Catt S, Bricker B, Yu G, Du L, Sobral BW: **Genome sequence of *Brucella abortus* vaccine strain S19 compared to virulent strains yields candidate virulence genes.** *PLoS One* 2008, **3**(5):e2193.
2. Tsolis RM, Seshadri R, Santos RL, Sangari FJ, Lobo JM, de Jong MF, Ren Q, Myers G, Brinkac LM, Nelson WC *et al*: **Genome degradation in *Brucella ovis* corresponds with narrowing of its host range and tissue tropism.** *PLoS One* 2009, **4**(5):e5519.
3. Kent WJ: **BLAT--the BLAST-like alignment tool.** *Genome Res* 2002, **12**(4):656-664.
4. Van Dongen S: **Graph Clustering Via a Discrete Uncoupling Process.** *SIAM J Matrix Anal & Appl* 2008, **30**:121-141.
5. Durbin R, Eddy S, Krogh A, Mitchison G: **Biological sequence analysis: probabilistic models of proteins and nucleic acids:** Cambridge University Press; 1998.
6. Edgar RC: **MUSCLE: a multiple sequence alignment method with reduced time and space complexity.** *BMC Bioinformatics* 2004, **5**:113.

7. Edgar RC: **MUSCLE: multiple sequence alignment with high accuracy and high throughput.** *Nucleic Acids Res* 2004, **32**(5):1792-1797.
8. Castresana J: **Selection of conserved blocks from multiple alignments for their use in phylogenetic analysis.** *Mol Biol Evol* 2000, **17**(4):540-552.
9. Talavera G, Castresana J: **Improvement of phylogenies after removing divergent and ambiguously aligned blocks from protein sequence alignments.** *Syst Biol* 2007, **56**(4):564-577.
10. Price MN, Dehal PS, Arkin AP: **FastTree 2--approximately maximum-likelihood trees for large alignments.** *PLoS One* 2010, **5**(3):e9490.

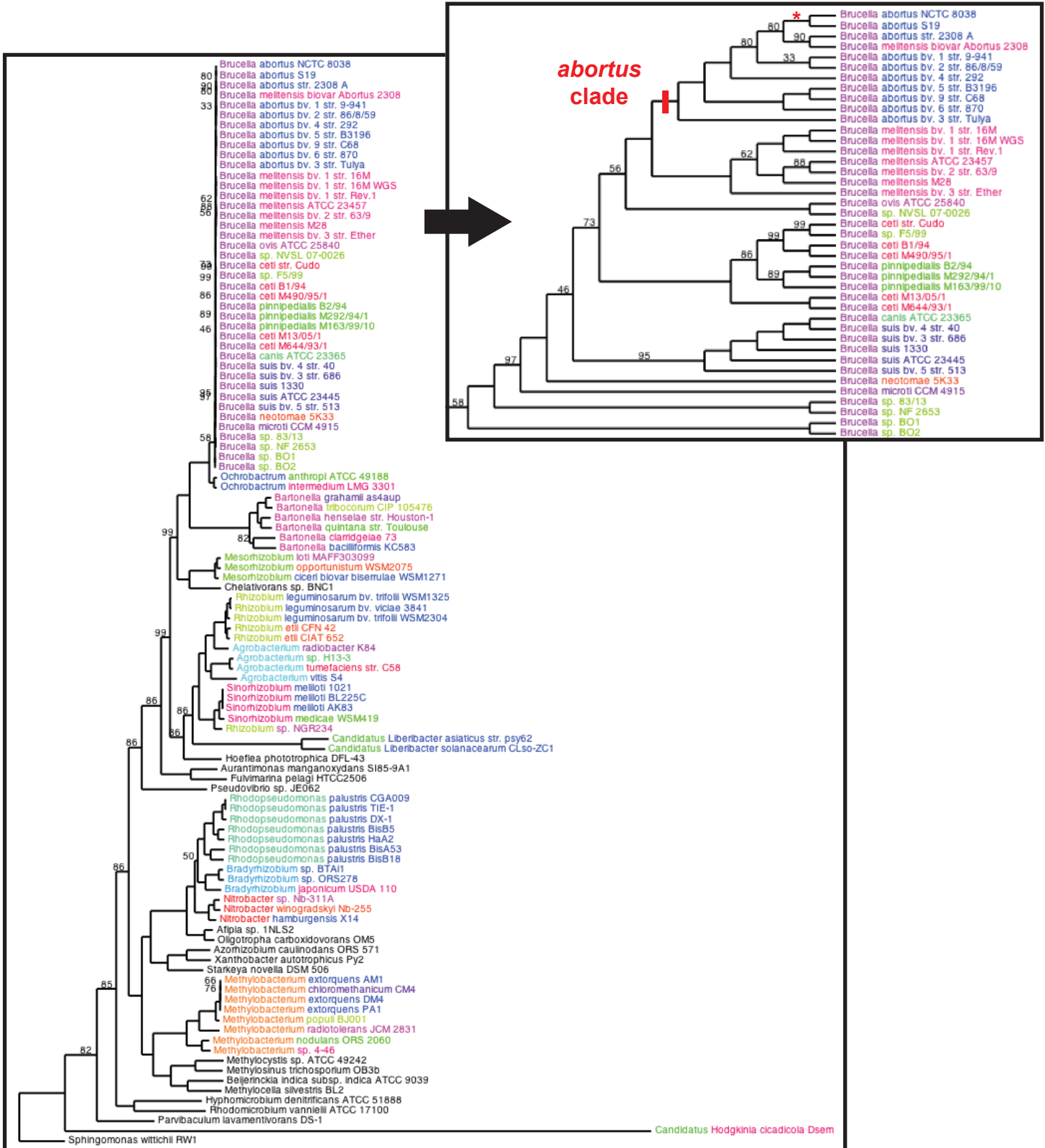
# A

Gene	PATRIC Annotation	Putative Pseudogenes			Ref. <sup>3</sup>
		Genome <sup>1</sup>	Mutation <sup>2</sup>	Effect	
eryA	Erythritol kinase EryA	C	D (1)	frameshift	--
		E,G	D (1)	nonsense/premature stop	--
		I	PM	altered start site; short	T
		N	I (1)	nonsense/premature stop	--
		P	amb.s	disrupted reading frame	--
eryB	Erythritol phosphate dehydrogenase EryB	--	--	--	--
eryC	Possible D-erythrulose 4-phosphate dehydrogenase EryC	B	D (703)	not annotated	C
eryD	Erythritol transcriptional regulator EryD	B	D (703)	altered start site; short	C
		I	D (7)	nonsense/premature stop	T
--	Predicted erythritol ABC transporter 2, hypothetical lipoprotein	E,G	D (1)	nonsense/premature stop	--
eryE	Predicted erythritol ABC transporter 2, ATP-binding component	O	D (1)	frameshift	--
eryF	Predicted erythritol ABC transporter 2, permease component	A,B	D (67)	altered start site; short	--
		D-H	D (1)	nonsense/premature stop	--
		I	D (2)	altered start site; short	T
		J-M	D (1)	nonsense/premature stop	--
eryG	Predicted erythritol ABC transporter 2, substrate-binding component	I,O	D (1)	frameshift	T

<sup>1</sup> A, *Brucella abortus* NCTC 8038; B, *Brucella abortus* S19; C, *Brucella canis* ATCC 23365; D, *Brucella ceti* B1/94; E, *Brucella ceti* M13/05/1; F, *Brucella ceti* M490/95/1; G, *Brucella ceti* M644/93/1; H, *Brucella ceti* str. Cudo; I, *Brucella ovis* ATCC 25840; J, *Brucella pinnipedialis* B2/94; K, *Brucella pinnipedialis* M163/99/10; L, *Brucella pinnipedialis* M292/94/1; M, *Brucella* sp. F5/99; N, *Brucella* sp. NF 2653; O, *Brucella* sp. NVSL 07-0026; P, *Brucella suis* bv. 3 str. 686.

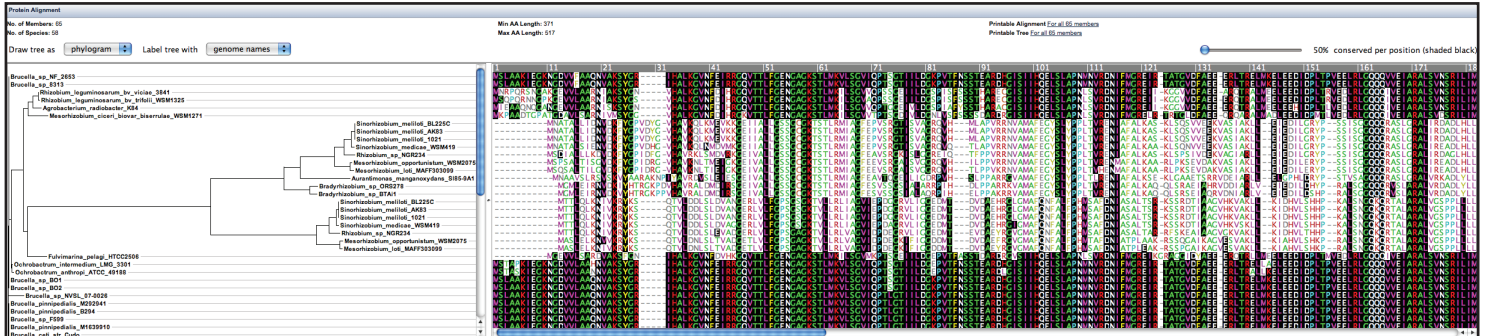
<sup>2</sup> D, deletion; PM, point mutation; I, insertion. Numbers in parentheses denote nucleotides.

<sup>3</sup> T, Tsolis et al. [1]; C, Crasta et al. [2].

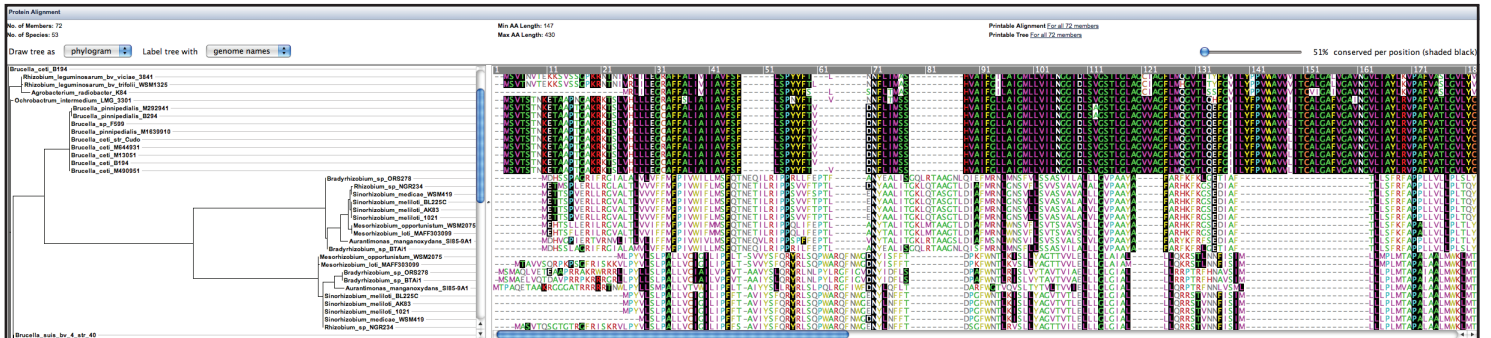
**B**

C

## ATP-binding protein (*eryE*)



## Permease (*eryF*)



## Substrate binding protein (*eryG*)

