

KH-domain proteins: another family of bacterial RNA matchmakers?

Mikołaj Olejniczak^{1,*}, Xiaofang Jiang², Maciej M. Basczok¹, Gisela Storz^{3,}**

¹*Institute of Molecular Biology and Biotechnology, Faculty of Biology, Adam Mickiewicz University, Uniwersytetu Poznańskiego 6, 61-614 Poznań, Poland.*

²*Intramural Research Program, National Library of Medicine, National Institutes of Health, Bethesda, MD 20894, USA.*

³*Division of Molecular and Cellular Biology, Eunice Kennedy Shriver National Institutes of Child Health and Human Development, Bethesda, MD 20892-4417, USA.*

For correspondences. *E-mail mol@amu.edu.pl and **E-mail storzg@mail.nih.gov

Supplementary Information

Table S1. Sequences and taxonomic information on KhpA and KhpB homologs identified in Genome Taxonomy Database (GTDB) (release 202).

Fig. S1. The taxonomic distribution of KhpA and KhpB based on the Genome Taxonomy Database (GTDB) (release 202).

Fig. S2. Gene synteny surrounding *khpA*.

Fig. S3. Gene synteny surrounding *khpB*.

Table S1. Sequences and taxonomic information on KhpA and KhpB homologs identified in Genome Taxonomy Database (GTDB) (release 202).

Protein families and domains were identified using InterProScan with default parameters.

Proteins that contain only the conserved KH domain (PTHR34654 or MF_00088) were identified as KhpA. Three conserved domains were used to identify KhpB genes: the Jag-N domain (PF14804, SM01245 or G3DSA:3.30.30.80), the KH domain (cd02414, G3DSA:3.30.300.20 or PF13083), and the R3H domain (PS51061, SM00393, PF01424, cd02644, G3DSA:3.30.1370.50 or SSF82708). Manual curation together with synteny information was used to improve the predictions. Some possible contaminants are noted in the .xlsx table.

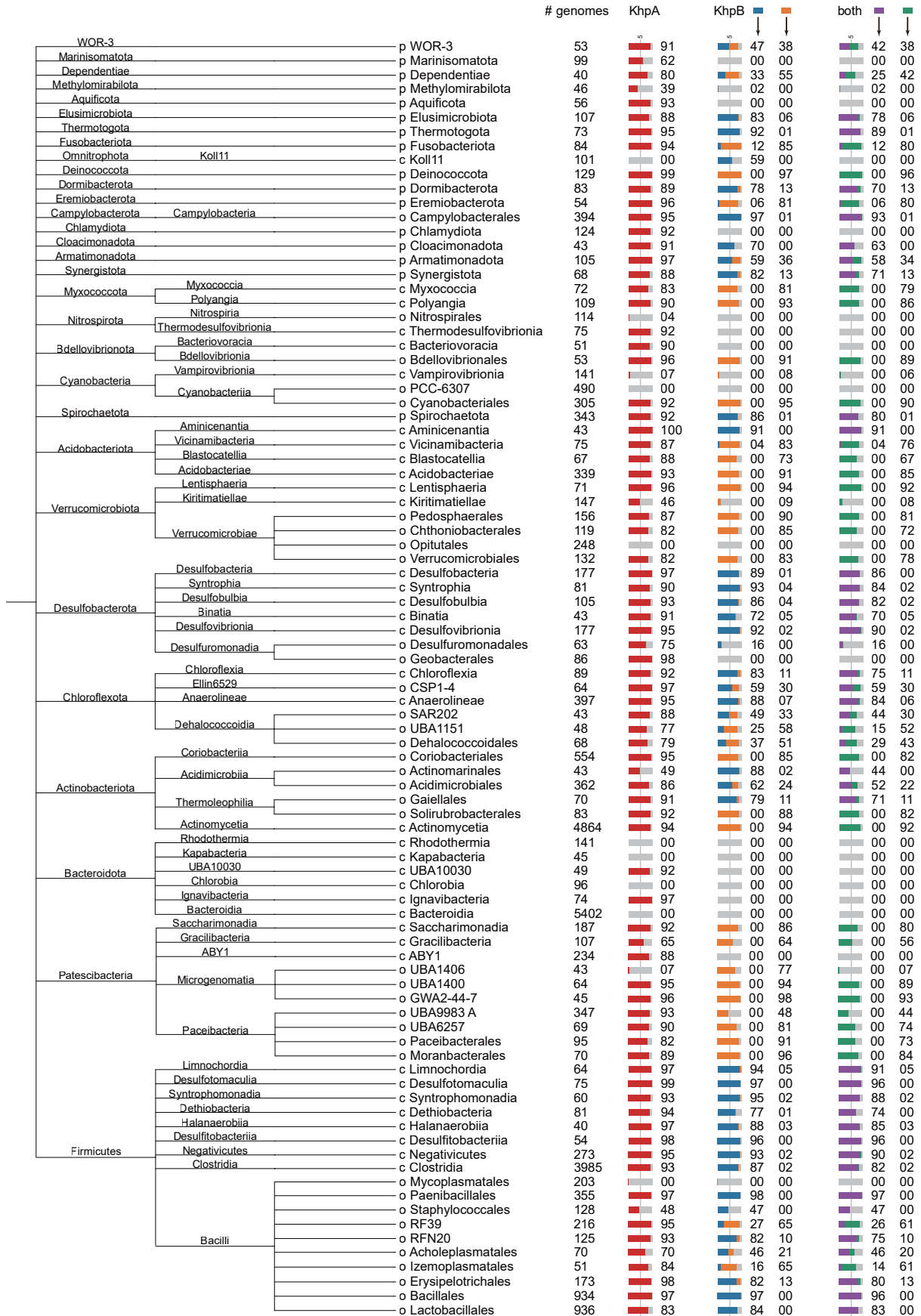


Fig. S1. The taxonomic distribution of KhpA and KhpB based on the Genome Taxonomy Database (GTDB) (release 202).

Based on Table S1. The numbers to the right of the cladogram represent the number of genomes examined in the corresponding taxa. Only clades with 40 or more genomes were included in the diagram. The red bar represents the percentage of genomes with KhpA genes in the taxa; the blue and orange bars represent the percentage of genomes with KhpB gene with and without the Jag-N domain, respectively; the purple bar represents the percentage of genomes with both KhpA and KhpB with the Jag-N domain; the green bar represents the percentage of genomes with both KhpA and KhpB without the Jag-N domain.



Fig. S2. Gene synteny surrounding *khpA*.

Orthologous gene families surrounding *khpA* genes for organisms shown in Fig. 1 were identified using OrthoFinder (version 2.3.8). Genome synteny was visualized with a custom script. The large triangle indicates the position of a gene encoding the 16S rRNA.



Fig. S3. Gene synteny surrounding *khpB*.

Orthologous gene families surrounding *khpB* genes for organisms shown in Fig. 1 were identified using OrthoFinder (version 2.3.8). Genome synteny was visualized with a custom script. The large triangle indicates the position of a gene encoding the 16S rRNA.