



HAL
open science

Design and implementation of a database for genome annotation

Benoît de Hertogh, Leïla Lahlimi, Christophe Lambert, Jean-Jacques Letesson, Eric Depiereux

► **To cite this version:**

Benoît de Hertogh, Leïla Lahlimi, Christophe Lambert, Jean-Jacques Letesson, Eric Depiereux. Design and implementation of a database for genome annotation. *Veterinary Microbiology*, 2008, 127 (3-4), pp.369. 10.1016/j.vetmic.2007.09.010 . hal-00532323

HAL Id: hal-00532323

<https://hal.science/hal-00532323>

Submitted on 4 Nov 2010

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Accepted Manuscript

Title: Design and implementation of a database for *Brucella melitensis* genome annotation

Authors: Benoît De Hertogh, Leïla Lahlimi, Christophe Lambert, Jean-Jacques Letesson, Eric Depiereux



PII: S0378-1135(07)00433-6
DOI: doi:10.1016/j.vetmic.2007.09.010
Reference: VETMIC 3810

To appear in: *VETMIC*

Received date: 13-3-2006
Revised date: 15-1-2007
Accepted date: 13-9-2007

Please cite this article as: De Hertogh, B., Lahlimi, L., Lambert, C., Letesson, J.-J., Depiereux, E., Design and implementation of a database for *Brucella melitensis* genome annotation, *Veterinary Microbiology* (2007), doi:10.1016/j.vetmic.2007.09.010

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

**Design and implementation of a database for *Brucella melitensis* genome
annotation.**

Benoît De Hertogh¹, Leïla Lahlimi¹, Christophe Lambert², Jean-Jacques Letesson¹, Eric
Depiereux^{1*}.

1. U.R.B.M., F.U.N.D.P., Rue de Bruxelles, 61, B-5000, Namur, Belgium.
2. BioXpr sa, Rue du séminaire, 22, B-5000, Namur, Belgium.

*To whom correspondence can be addressed:

Eric Depiereux,
Unité de Recherche en Biologie Moléculaire,
Facultés Universitaires Notre-Dame de la Paix,
Rue de Bruxelles, 61
B-5000 Namur,
Belgium

Tel: +32 81724415

Fax: +32 81724420

E-mail: eric.depiereux@fundp.ac.be

1 **Abstract**

- 2 The genome sequences of three *Brucella* biovars and of some species close to *Brucella* sp.

3 have become available, leading to new relationship analysis. Moreover, the automatic
4 genome annotation of the pathogenic bacteria *Brucella melitensis* has been manually
5 corrected by a consortium of experts, leading to 899 modifications of start sites predictions
6 among the 3198 open reading frames (ORFs) examined. This new annotation, coupled with
7 the results of automatic annotation tools of the complete genome sequences of the *B.*
8 *melitensis* genome (including BLASTs to 9 genomes close to *Brucella*), provides
9 numerous data sets related to predicted functions, biochemical properties and phylogenic
10 comparisons.

11 To made these results available, α PAGe, a functional auto-updatable database of the
12 corrected sequence genome of *B. melitensis*, has been built, using the entity-relationship
13 (ER) approach and a multi-purpose database structure. A friendly graphical user interface
14 has been designed, and users can carry out different kinds of information by three levels of
15 queries: (1) the basic search use the classical keywords or sequence identifiers; (2) the
16 original advanced search engine allows to combine (by using logical operators) numerous
17 criteria: (a) keywords (textual comparison) related to the pCDS's function, family domains
18 and cellular localization; (b) physico-chemical characteristics (numerical comparison) such
19 as isoelectric point or molecular weight and structural criteria such as the nucleic length or
20 the number of transmembrane helix (TMH); (c) similarity scores with *E. coli* and 10
21 species phylogenetically close to *B. melitensis*; (3) complex queries can be performed by
22 using a SQL field, which allows all queries respecting the database's structure.

23 The database is publicly available through a Web server at the following url:

24 <http://www.fundp.ac.be/urbm/bioinfo/aPAGe>

25

26 **Keywords:** *Brucella melitensis*, relational, database, annotation, db-main.

27 **Introduction**28 *Brucella melitensis*

29 Brucellosis is a worldwide distributed bacterial disease, affecting reindeers in Alaska and
30 Siberia, camels in the Middle East, cattles, pigs, goats, sheeps.... Sheeps and goats are the
31 most susceptible. In human, the disease is known as Malta fever, undulating fever,
32 mediterranean fever or Melitococcie. It causes abortion in females and orchitis in males.
33 Humans catch this disease from contaminated meat or milk, aborted fetuses and
34 slaughtering infected animals (Godfroid *et al.* 2005).

35 The six known species of *Brucella* (Proteobacteria; Alphaproteobacteria; Rhizobiales;
36 Brucellaceae; *Brucella*), of which four are pathogenic for human (*B. melitensis*, *B. abortus*,
37 *B. suis* and *B. canis*) are often regarded today as being biovars of one single species (for
38 contradictory informations, see Chain *et al.*, 2005).

39

40

Table 1

41

42 The development of many chronic diseases involves interactions between environmental
43 factors and genes that regulate important physiological processes. Several projects are
44 trying to understand those processes, and to develop methods to predict toxicity and
45 understand the genetic basis of differential susceptibility (Mattingly *et al.*, 2004). Genome
46 sequencing projects led to new disciplines in biology, and promise a better comprehension
47 of disease-associated genes (Miller and Kumar, 2001).

48 The availability of the genome sequence of three *Brucella* biovars and eight species
49 phylogenetically close to *Brucella sp.* (Table 1) opens the way to genomic comparisons
50 that may elucidate the molecular mechanisms of brucellosis, but also phylogenetic and
51 evolutionary relationships, speciation and emergence of new gene families.

52

53 The correction consortium

54 The 3198 pCDS of the *Brucella melitensis* genome were sequenced and automatically
55 annotated by Integrated Genomics Inc. (Delvecchio *et al.*, 2002). However, many errors
56 were notified in predictions of the position of the start codon of these pCDS. The Research
57 Unit in Molecular Biology (URBM) initiated a project of correction of the start position of
58 all the genome sequences. This project was carried out in collaboration with four research
59 teams, members of the European consortium of research COST845 (VLA Weybridge,
60 U.K.; U. Navarra, Spain; U. Cantabria, Spain; Inserm U. 431, France) (Dricot, 2004).

61

62 Objectives

63 Through this paper, we had five main objectives: (1) to made available the Consortium
64 corrections (2) and the results of the automatic annotation (3) in a polyvalent auto-
65 updatable and easily upgradable (i.e. to others genomes) database (4) using a performing
66 DataBase Management System (DBMS), in order to allow complex queries by the use of
67 (5) performing and user friendly search tools.

68

69 **Material and methods**

70 The correction protocol

71 Each pCDS was manually checked by at least two scientific teams, using the following
72 protocol:

- 73 1. Genome files of *Brucella melitensis* (NC_003318.gbk and NC_003317.gbk)
74 had been downloaded from the National Center for Biotechnology Information
75 (NCBI) ftp site. The free Artemis software was used to visualize annotations;
- 76 2. For each ORF, the predicted start was checked. The wrong starts were
77 detected by carrying out the following steps:
 - 78 a. Control of the start codon. TTG is usually a wrong start site;

- 79 b. Detection of overlapping ORFs, which usually hide a wrong start
80 codon;
- 81 c. Analyses of the corresponding protein sequence alignment against the
82 Non-Redundant (NR) database from NCBI. If similar sequences in
83 phylogenetically close organisms are longer or shorter, the start site is
84 probably wrong;
- 85 3. A new start site was checked using Artemis software with, in increasing
86 preferences corresponding to their observed frequency, start codons: ATG, GTG
87 and TTG;
- 88 4. If several potential start positions were found, without operon or Ribosome
89 Binding Site (RBS) (Salgado *et al.*, 2000) identified, the pCDS having the closest
90 size to homologous sequences detected in NR was selected.

91 Otherwise, the longer pCDS with an ATG or GTG was chosen.

92

93 Functional annotation

94 Functional annotations of proteins translated from the ORFeome library are done *in silico*
95 using the following programs:

- 96 1. BLASTP) against the NR protein database from GenBank and Swiss-Prot
97 databases;
- 98 2. BLASTN against the genomes listed in Table 1;
- 99 3. hmmpfam against the 8183 proteins of the Protein Family database (Pfam
100 19.0);
- 101 4. Prediction of the cellular localizations by an updated version of PSORTII
102 (Nakai *et al.*, 1999) for Gram-negative bacteria. PSORTII examines a given protein
103 sequence for amino acid composition, similarity to proteins of known localization,
104 presence of a signal peptide, transmembrane alpha-helices and motifs

105 corresponding to specific localizations. A probabilistic method integrates this
106 analysis, returning a list of five possible localization sites with associated
107 probability scores;

108 5. Prediction of the transmembrane segments are predicted by TMHMM v.2.0,
109 the most reliable transmembrane prediction program (Möller, 2001);

110 6. Prediction of the secondary structure of each pCDS using the PSIPRED2
111 protein structure prediction server (McGuffin *et al.*, 2000);

112 7. Prediction of the three-dimensional structure by ESyPred3D (Lambert *et al.*,
113 2002), with alignment of the peptidic sequences against the Protein data bank
114 (PDB) (Berman *et al.*, 2000). This program produces a multiple alignment of the
115 query sequence with several sequences from the PDB, and builds a consensus of
116 high reliability. This reliable alignment is subsequently used for the building of the
117 homology model.

118

119

120 The DB-MAIN CASE tool

121 Because a genome annotation generates numerous data sets that cannot be easily managed
122 with simple file management software, a powerful and complex knowledge-management
123 software, DB-MAIN, had been used to create a database using the relational model in
124 accordance with the ER model (Entity Relationship). DB-MAIN is a data-oriented
125 Computer Aided Software Engineering (CASE) environment, designed to support most
126 database engineering processes, including: (1) requirement analysis, conceptual design,
127 normalisation, schema integration, logical design, physical design, schema's optimisation
128 and code generation; (2) schema transformation, model transformation; (3) schema
129 analysis, code analysis, data reverse engineering; (4) database migration, database
130 evolution, database integration and federation, data wrapper design and generation; (5)

131 temporal database design, active database evaluation and generation.

132 In addition, the method-modelling component allows the user to define any of these sub-
133 models (Englebert and Hainaut, 1999; <http://www.db-main.be/>).

134

135 The database design

136 The database design follows globally the "Merise" method. It starts with the extraction of
137 the entities from the application domain (in this paper, the data sets provided by both
138 manual and automatic annotation of the *B. melitensis* genome). These entities will
139 constitute the database classes, which models the real-world organization and its important
140 data elements and relationships.

141 Although there are many conceptual models drawn for biological databases, the conceptual
142 schema remains specific to the application domains, which are various and depend on
143 projects subject. The conceptual model is an extension of the Entity-Relationship model. It
144 defines the logical relationships that link the database entities and the object attributes, in
145 order to build the data structure. This step involves the description of the entire information
146 content of the database.

147 The next step, the logical design, is the translation of the conceptual schema in logical
148 schema. These schemas represent the same information, but the second one (logical)
149 expresses the data through the construct of the DBMS. It shows how data are organised in
150 a relational way: the database is represented as a group of related tables that can be
151 managed by a relational database management system (RDBMS). Finally, the entity
152 relationship model is translated into a database structure that can be afterwards easily
153 expressed using structured query language (SQL).

154 Both conceptual and logical schemas are available on the database website.

155

156 The web interface

157 The data sets are stored in a MySQL (<http://www.mysql.com/>) relational database. The
158 database can be queried through a web interface built using PHP Hypertext Preprocessor
159 (PHP), a widely used general-purpose scripting language especially suited for web
160 development and more specifically to be connected to a database (<http://www.php.net/>).

161

162 Machines

163 The annotation of the genome and the construction of the database were carried out on
164 Silicon Graphics Octane duo and a cluster of PC (Table 2).

165

166 **Results**

167

168 Increasing of pCDS reliability

169 The correction project consortium has increased the database information reliability. Start
170 positions have been corrected for 899 pCDS. The mean difference in position is 69
171 nucleotides or 23 amino acids. 565 pCDS have been shortened and 334 have been
172 lengthened. The corrected pCDS are mentioned in the commentary field.

173

174 The α PAGe database

175 The sequence databases offer a great source of information for studies on biological
176 variation, evolutionary patterns and protein family characterization. In order to find all
177 pieces of information available, it is often necessary to search several databases or to click
178 on several links. Facing the complexity to obtain information, each α PAGe entry includes
179 the whole information about a sequence (functional and structural annotation, pCDS's
180 properties, links to other biological databases...).

181 The identification and the functional characterization of genes, that may be involved in
182 disease development or may be responsible for virulence, is a critical step. An advanced

183 search tool form—allows the combination of different criteria and may help to select a
184 subset of pCDS for further analysis. Most databases propose advanced search tools based
185 only on keywords searches. The α PAGe makes possible to perform searches combining
186 several criteria, with the possibility to impose choice constraints (and/or) at each level (see
187 below). The database is publicly available through a Web server at
188 <http://www.fundp.ac.be/urbm/bioinfo/aPAGe>.

189

190 ORFeome library

191 Entire libraries composed of all protein-encoding Open Reading Frames (ORF) cloned in
192 highly flexible vectors represent a new type of resource, which is needed to take full
193 advantage of informations generated by the sequencing efforts and to address the new type
194 of question and hypothesis generated in post-genomic area.

195 Thus, the complete genome sequence of *B. melitensis* had been used to generate a protein-
196 coding ORF database. This ORFeome library contains 3091 Gateway entry clones, each
197 one corresponding to a defined ORF. This strategy may help to validate the genome
198 annotation and to create a resource to functionally characterize the proteome (Dricot *et al.*,
199 2004). The cloning state field displays the information available for this experiment such
200 as the forward and reverse primers used for each cloned pCDS.

201 Moreover, the "Similarity (BLASTN) in close organisms" field allows a quick overview of
202 the distribution of the homologous sequences in the genomes of 10 alpha-proteobacteria,
203 phylogenetically close to each other, and *Escherichia coli K12*, as Gram-negative model.

204 In order to facilitate the databases mining, links are proposed to display the precomputed
205 best hits against the following databases: GenBank, Protein Information Resource (PIR),
206 Protein Data Bank (PDB), DNA Data Bank of Japan (dbj), the Protein Research
207 Foundation (PRF) (www4.prf.or.jp/en/), the European Molecular Biology Laboratory
208 (EMBL) and the Protein Family (Pfam) database.

209

210 Cellular localization and membrane topology

211 The cellular localization field contains three parts: the observed cellular localization (often
212 unknown), the predicted cellular localization (PSORT II) and the compatibility between
213 predicted and observed subcellular localization.

214 The membrane topology and the number of transmembrane helices (TMHs) are also
215 precomputed (TMHMM).

216

217 Secondary and three-dimensional structure

218 The predicted secondary structures are coloured in the following way: helices, beta sheets,
219 and coils are represented respectively in red, blue and grey. Transmembrane segments are
220 written with underlined bold characters.

221 The 3D structure field reports the homology modelling procedure if a protein with a
222 detectable similarity is found in the PDB. It includes two parts:

223 1. The "3D structure field", which proposes a link to download the 3D final
224 model of the treated pCDS. The model may be displayed using an external software
225 (i.e. the Swiss PDB viewer (<http://www.expasy.org>) or a browser plug in (Chime)).

226 2. The "Modelling characteristics" field summarizes information used during
227 the modelling step: the template (protein with known 3D structure) used for the
228 alignment, the model building, the template experimental method and its
229 parameters (resolution and R-value), the percentage of sequence modelled and a
230 link to find more information and coordinates of the template structure at the PDB
231 web site. Models are updated each month or at each new PDB release.

232

233 pCDS properties

234 The last sets of predictions concerns pCDS start and end positions, GC content, theoretical

235 physico-chemical properties, that includes molecular weight and isoelectric point and
236 nucleic and peptidic sequence in FASTA format.

237 Promoter and termination sites can easily be studied by displaying the sequence of a
238 chosen number of nucleotides located upstream or downstream of the pCDS. This tool may
239 be useful to identify the non-coding upstream and downstream of a pCDS. It also allows to
240 define the pCDS extremities for the cloning experiments.

241

242 References and cross-references

243 The cross-reference field allows users to be redirected towards KEGG and translated
244 EMBL (TrEMBL) public databases. In a reference field, users are invited to add references
245 related to the pCDS treated.

246

247 User interface

248 The database can be browsed using a powerful graphical user interface. One may consider
249 three levels of queries:

250

251 *Basic search*

252 The database can be browsed:

- 253 1. By the pCDS ID;
- 254 2. By keywords (text search). Various fields can be targeted: GenBank
255 annotation, Swiss-Prot annotation, observed function, comments, predicted or
256 observed localization, Pfam summary, templates used in the 3D modelling
257 procedure or a key word related to the cloning state;
- 258 3. By the display of the complete list of ORFs;
- 259 4. By similarity (using BLAST). Peptidic or nucleic sequences can be blasted
260 against the whole genomes of *B. melitensis*, *B. suis* and *B. abortus*;

261 5. By using regular expression to search for specific patterns (nucleic or
262 peptidic).

263

264 *Advanced search*

265 An advanced search tool allows users to limit queries to some defined fields or to combine
266 search terms with logical operators.

267

268 Following fields are available:

- 269 1. Keywords searches (text comparison) described in the basic search.
- 270 2. Quantitative data (numerical comparisons using "<" and ">" operators)
271 including:
 - 272 a. Physico-chemical properties, such as isoelectric point or molecular
273 weight;
 - 274 b. Structural properties such as:
 - 275 i. The number of trans-membrane helices (TMHs);
 - 276 ii. The target-template percentage identity (three dimensional structure
277 prediction);
 - 278 iii. The percentage identity of the deduced sequence modelled using the
279 ESyPred3D system.
 - 280 c. The pCDS properties such as:
 - 281 i. The percentage in GC;
 - 282 ii. The position of the first and the last nucleotide;
 - 283 iii. The nucleotidic and peptidic lengths.

284 It is also possible to add constraints (expected value range) related to similarities between
285 the *B. melitensis* pCDS considered and the pCDS of at least one of the 11 species described
286 in Table 1.

287

288

Figure 1

289

290 To allow combinations of quantitative data with the logical operators "and" and "or" and
291 with the comparisons operators "<" and ">" is probably the most interesting particularity of
292 the advanced search tool (Fig. 1). This unusual way to proceed permits to extract data in
293 order to perform statistical analysis. For example, the study of the frequency distribution of
294 the predicted pI shows a bimodal distribution (Fig. 2). It is well known that, in disruptive
295 selection, selection pressures act against individuals in the middle of the trait distribution,
296 resulting in a bimodal curve. Thus, this predicted pI bimodal distribution seems to suggest
297 an important evolutionary selection pressure applied against *B. melitensis*, which
298 unfavourise the *B. melitensis* ORFs around pI's of 8. However, for different authors, this
299 multimodal distribution is an effect of allowed combinations of the charged amino acids,
300 and not due to evolutionary causes (*i.e.* Nandy *et al.*, 2005; Schwartz *et al.*, 2001).
301 Similarly, data extraction for various statistical analysis, implying one or more fields, can
302 easily be performed.

303

304

Figure 2

305

306 SQL search

307 For advanced users, data may be retrieved from the Web interface using SQL queries.

308

309

310

311 Conclusion and Perspectives

312 We developed, under DBMS, a relational updatable database dedicated to the pathogenic
313 bacteria *Brucella melitensis*, from the sequenced genome manually corrected by a

314 specialists consortium and from the data generated *in silico* by several among the most
315 powerful prediction programs.

316 Search tools are efficient, friendly, but also original. In the advanced search, users may
317 combine, by using logical operators, several criteria related to the pCDS properties in order
318 to select one or more specific pCDS: (1) keywords (textual comparison) related to the
319 pCDS's function, family domains and cellular localization; (2) physico-chemical
320 characteristics (numerical comparison) such as isoelectric point, molecular weight or
321 structural criteria such as the nucleic length or the number of TMHs; (3) homology (or not)
322 in 10 species phylogenetically close to *B. melitensis* (expected value range).

323 The structure of the database allows its extension to other genomes, especially those
324 phylogenetically close to *Brucella melitensis*, but also to model organism like
325 *Saccharomyces cerevisiae*, *Drosophila melanogaster* or *Homo sapiens*. The functional and
326 structural annotation of these genomes could be incorporated automatically from
327 nucleotidic sequences.

328 The extension of the advanced search to new fields, like TMHs or motives is also on the
329 way. Such improvements will permit more complex queries.

330 New advanced search tool form that deal with all the data types and new links to biological
331 databases for complementary information will also be added.

332

333 Finally, a special attention must be paid in order to make sure that the alteration is
334 compatible with the existing database. The DB-MAIN CASE tool supports such meta-data
335 evolution. The approach relies on a generic database model and on the transformational
336 paradigm that states that database engineering processes can be modelled by schema
337 transformation. Indeed, a transformation provides both structural and instance mappings
338 that formally define how to modify database structures and contents.

339

340 Acknowledgements

341 The authors gratefully acknowledge the aMAZE team and BioXpr colleagues for their help
342 for the database conceptual schema design. They also address a special thank the DB-
343 MAIN team for their assistance and to Fabrice Berger and Xavier De Bolle for critical
344 discussions and comments.

345

346 **References**

347 Alsmark, C.M., Frank, A.C., Karlberg, E.O., Legault, B.A., Ardell, D.H., Canback, B.,
348 Eriksson, A.S., Naslund, A.K., Handley, S.A., Huvet, M., Scola, B.L., Holmberg, M. and
349 Andersson, S.G., 2004. The louse-borne human pathogen *Bartonella quintana* is a
350 genomic derivative of the zoonotic agent *Bartonella henselae*. Proc. Natl. Acad. Sci. USA.,
351 101, 9716-9721.

352 Barnett, M.J., Fisher, R.F., Jones, T., Komp, C., Abola, A.P., Barloy-Hubler, F., Bowser,
353 L., Capela, D., Galibert, F., Gouzy, J., Gurjal, M., Hong, A., Huizar, L., Hyman, R.W.,
354 Kahn, D., Kahn, M.L., Kalman, S., Keating, D.H., Palm, C., Peck, M.C., Surzycki, R.,
355 Wells, D.H., Yeh, K.C., Davis, R.W., Federspiel, N.A. and Long, S.R., 2001. Nucleotide
356 sequence and predicted functions of the entire *Sinorhizobium meliloti* pSymA
357 megaplasmid. Proc. Natl. Acad. Sci. USA., 98, 9883-9888.

358 Blattner, F.R., Plunkett, G., Bloch, C.A., Perna, N.T., Burland, V., Riley, M., Collado-
359 Vides, J., Glasner, J.D., Rode, C.K., Mayhew, G.F., Gregor, J., Davis, N.M., Kirkpatrick,
360 H.A., Goeden, M.A., Rose, D.J., Mau, B. and Shao, Y., 1997. The complete genome
361 sequence of *Escherichia coli* K-12. Science. 277, 1453-1474.

362 Chain PS, Commerci DJ, Tolmasky ME, Larimer FW, Malfatti SA, Vergez LM, Aguerro F,
363 Land ML, Ugalde RA, Garcia E., 2005. Whole-genome analyses of speciation events in
364 pathogenic *Brucellae*. Infect. Immun. 73, 8353-61.

365 Delvecchio, V.G., Kapatral, V., Redkar, R.J., Patra, G., Mujer, C., Los, T., Ivanova, N.,
366 Anderson, I., Bhattacharyya, A., Lykidis, A., Reznik, G., Jablonski, L., Larsen, N.,
367 D'souza, M., Bernal, A., Mazur, M., Goltsman, E., Selkov, E., Elzer, P.H., Hagijs, S.,
368 O'callaghan, D., Letesson, J.J., Haselkorn, R., Kyrpides, N. and Overbeek, R., 2002. The
369 genome sequence of the facultative intracellular pathogen *Brucella melitensis*. Proc. Natl.
370 Acad. Sci. USA. 99, 443-448.

371 Dricot, A., Rual, J.F., Lamesch, P., Bertin, N., Dupuy, D., Hao, T., Lambert, C., Hallez, R.,
372 Delroisse, J.M., Vandenhoute, J., Lopez-Goñi, I., Moriyon, I., Garcia-Lobo, J.M., Sangari,
373 F.J., Macmillan, A.P., Cutler, S.J., Whatmore, A.M., Bozak, S., Sequerra, R., Doucette-
374 Stamm, L., Vidal, M., Hill, D.E., Letesson, J.J. and Debolle, X., 2004. Generation of
375 *Brucella melitensis* ORFeome version 1.1. *Genome Research*. 14, 2201-2206.

376 Englebert, V. and Hainaut, J.L., 1999. DB-MAIN: A Next Generation Meta-CASE.
377 *Journal of Information Systems*. 24, 99-112.

378 Galibert F., Finan T.M., Long S.R., Puhler A., Abola P., Ampe F., Barloy-Hubler F.,
379 Barnett M.J., Becker A., Boistard P., Bothe G., Boutry M., Bowser L., Buhrmester J.,
380 Cadieu E., Capela D., Chain P., Cowie A., Davis R.W., Dreano S., Federspiel N.A., Fisher
381 R.F., Gloux S., Godrie T., Goffeau A., Golding B., Gouzy J., Gurjal M., Hernandez-Lucas
382 I., Hong A., Huizar L., Hyman R.W., Jones T., Kahn D., Kahn M.L., Kalman S., Keating
383 D.H., Kiss E., Komp C., Lelaure V., Masuy D., Palm C., Peck M.C., Pohl T.M., Portetelle
384 D., Purnelle B., Ramsperger U., Surzycki R., Thebault P., Vandenbol M., Vorholter F.J.,
385 Weidner S., Wells D.H., Wong K., Yeh K.C., Batut J., 2001. The composite genome of the
386 legume symbiont *Sinorhizobium meliloti*. *Science*. 293:668-72.

387 Godfroid J., Cloeckert A., Liautard J.P., Kohler S., Fretin D., Walravens K., Garin-Bastuji
388 B., Letesson J.J., (2005) From the discovery of the Malta fever's agent to the discovery of a
389 marine mammal reservoir, brucellosis has continuously been a re-emerging zoonosis. *Vet.*
390 *Res*. 36:313-26.

391 Goodner, B., Hinkle, G., Gattung, S., Miller, N., Blanchard, M., Quorollo, B., Goldman, B.
392 S., Cao, Y., Askenazi, M., Halling, C., Mullin, L., Houmiel, K., Gordon, J., Vaudin, M.,
393 Iartchouk, O., Epp, A., Liu, F., Wollam, C., Allinger, M., Doughty, D., Scott, C., Lappas,
394 C., Markelz, B., Flanagan, C., Crowell, C., Gurson, J., Lomo, C., Sear, C., Strub, G., Cielo,
395 C. and Slater, S., 2001. Genome sequence of the plant pathogen and biotechnology agent
396 *Agrobacterium tumefaciens* C58. *Science*. 294:2323-2328.

397 Halling S.M., Peterson-Burch B.D., Bricker B.J., Zuerner R.L., Qing Z., Li L.L., Kapur V.,
398 Alt D.P., Olsen S.C. (2005) Completion of the genome sequence of *Brucella abortus* and
399 comparison to the highly similar genomes of *Brucella melitensis* and *Brucella suis*. *J*
400 *Bacteriol.* 187:2715-26.

401 Kaneko, T., Nakamura, Y., Sato, S., Asamizu, E., Kato, T., Sasamoto, S., Watanabe, A.,
402 Idesawa, K., Ishikawa, A., Kawashima, K., Kimura, T., Kishida, Y., Kiyokawa, C.,
403 Kohara, M., Matsumoto, M., Matsuno, A., Mochizuki, Y., Nakayama, S., Nakazaki, N.,
404 Shimpo, S., Sugimoto, M., Takeuchi, C., Yamada, M. and Tabata, S., 2000. Complete
405 genome structure of the nitrogen-fixing symbiotic bacterium *Mesorhizobium loti*. *DNA*
406 *Res.* 7:331-338.

407 Lambert, C., Leonard, N., Xavier De Bolle and Depiereux, E., 2002. ESyPred3D:
408 Prediction of proteins 3D structures. *Bioinformatics*. 18:1250-1256.

409 Larimer, F.W., Chain, P., Hauser, L., Lamerdin, J., Malfatti, S., Do, L., Land, M.L.,
410 Pelletier, D.A., Beatty, J.T., Lang, A.S., Tabita, F.R., Gibson, J.L., Hanson, T.E., Bobst,
411 C., Torres, J.L., Peres, C., Harrison, F.H., Gibson, J. and Harwood, C.S. (2004) Complete
412 genome sequence of the metabolically versatile photosynthetic bacterium
413 *Rhodospseudomonas palustris*. *Nat. Biotechnol.* 22:55-61.

414 Mattingly, C.J., Colby, G.T., Rosenstein, M.C., Forrest, J.N. and Boyer, J.L., 2004.
415 Promoting comparative molecular studies in environmental health research: an overview of
416 the comparative toxicogenomics database (CTD). *Pharmacogenomics J.*, 4, 5-8.

417 McGuffin, L.J., Bryson, K. and Jones, D.T., 2000. The PSIPRED protein structure
418 prediction server. *Bioinformatics*. 16:404-405.

419 Miller, M. and Kumar, S., 2001. Understanding human disease mutations through the use
420 of interspecific genetic variation. *Hum. Mol. Genet.* 10:2319-2328.

421 Möller, S., Croning, M.D.R. and Apweiler, R., 2001. Evaluation of methods for the
422 prediction of membrane spanning regions. *Bioinformatics*. 17, 646-653.

423 Nandi S., Mehra N., Lynn A.M., Bhattacharya A., 2005. Comparison of theoretical
424 proteomes: Identification of COGs with conserved and variable pI within the multimodal
425 pI distribution. *BMC Genomics*. 6:116.

426 Nierman, W.C., Feldblyum, T.V., Laub, M.T., Paulsen, I.T., Nelson, K.E., Eisen, J.A.,
427 Heidelberg, J.F., Alley, M.R., Ohta, N., Maddock, J.R., Potocka, I., Nelson, W.C., Newton,
428 A., Stephens, C., Phadke, N.D., Ely, B., Deboy, R.T., Dodson, R.J., Durkin, A.S., Gwinn,
429 M.L., Haft, D.H., Kolonay, J.F., Smit, J., Craven, M.B., Khouri, H., Shetty, J., Berry, K.,
430 Utterback, T., Tran, K., Wolf, A., Vamathevan, J., Ermolaeva, M., White, O., Salzberg,
431 S.L., Venter, J.C., Shapiro, L., Fraser, C.M. and Eisen, J., 2001. Complete genome
432 sequence of *Caulobacter crescentus*. *Proc. Natl. Acad. Sci. USA.* 98, 4136-4141.

433

434 Paulsen, I.T., Seshadri, R., Nelson, K.E., Eisen, J.A., Heidelberg, J.F., Read, T.D., Dodson,
435 R.J., Umayam, L., Brinkac, L.M., Beanan, M.J., Daugherty, S.C., Deboy, R.T., Durkin,
436 A.S., Kolonay J.F., Madupu, R., Nelson, W.C., Ayodeji, B., Kraul, M., Shetty, J., Malek,
437 J., Aken, S.E.V., Riedmuller, S., Tettelin, H., Gill, S.R., White, O., Salzberg, S.L., Hoover,
438 D.L., Lindler, L.E., Halling, S.M., Boyle, S.M. and Fraser, C.M., 2002. The *Brucella suis*
439 genome reveals fundamental similarities between animal and plant pathogens and
440 symbionts. *Proc Natl. Acad. Sci. USA.* 99, 13148-13153.

441 Salgado, H., Moreno-Hagelsieb, G., Smith, T.F. and Collado-Vides, J., 2000. Operons in
442 *Escherichia coli*: genomic analyses and predictions. *Proc. Natl. Acad. Sci.*, 97, 6652-6657.

443 Sanchez, D.O., Zandomeni, R.O., Cravero, S., Verdun, R.E., Pierrou, E., Faccio, P., Diaz,
444 G., Lanzavecchia, S., Agüero, F., Frasch, A.C., Andersson, S.G., Rossetti, O.L., Grau, O.
445 and Ugalde, R.A., 2001. Gene discovery through genomic sequencing of *Brucella abortus*.
446 *Infect. Immun.* 69, 865-868.

447 Schwartz, R., Ting, C.S., King, J., 2001. Whole proteome pI values correlate with
448 subcellular localizations of proteins for organisms within the three domains of life.
449 *Genome Res.* 11:703-709.

450 Wood, D.W., Setubal, J.C., Kaul, R., Monks, D.E., Kitajima, J.P., Okura, V.K., Zhou, Y.,
451 Chen, L., Wood, G.E., Almeida, N.F., Woo, L., Chen, Y., Paulsen, I.T., Eisen, J.A., Karp,
452 P.D., Bovee, D., Chapman, P., Clendenning, J., Deatherage, G., Gillet, W., Grant, C.,
453 Kutyavin, T., Levy, R., Li, M.J., McClelland, E., Palmieri, A., Raymond, C., Rouse, G.,
454 Saenphimmachak, C., Wu, Z., Romero, P., Gordon, D., Zhang, S., Yoo, H., Tao, Y.,
455 Biddle, P., Jung, M., Krespan, W., Perry, M., Gordon-Kamm, B., Liao, L., Kim, S.,
456 Hendrick, C., Zhao, Z.Y., Dolan, M., Chumley, F., Tingey, S.V., Tomb, J.F., Gordon, M.
457 P., Olson, M.V. and Nester E.W., 2001. The genome of the natural genetic engineer
458 *Agrobacterium tumefaciens* C58. *Science*, 294, 2317-2323.

459 **Figure legends**

460

461 Fig. 1. Combinations of numerous fields and both logical and comparisons operators allow
462 to construct complex queries through a user-friendly interface.

463

464 . Fig. 2. The bimodal distribution of the predicted isoelectric point seems to suggest an
465 important evolutionary selection pressure to *B. melitensis*, but some authors invalid this
466 suggestion (see text) (each classes defines an half-unity of pI).

467

d Manuscript



Advanced search form (can take more than 1 minute)

Text comparison

Choose an Organism:

Search in: both chromosomes small chromosome large chromosome

and or

and or

 and or

Numeric comparison

| | | | | | | | |
|--|----------------------|------------------|----------------------|--|----------------------|-----------------|----------------------|
| <input type="radio"/> and <input type="radio"/> or | <input type="text"/> | < PI < | <input type="text"/> | <input type="radio"/> and <input type="radio"/> or | <input type="text"/> | < MW < | <input type="text"/> |
| <input type="radio"/> and <input type="radio"/> or | <input type="text"/> | < first nucl < | <input type="text"/> | <input type="radio"/> and <input type="radio"/> or | <input type="text"/> | < last nucl < | <input type="text"/> |
| <input type="radio"/> and <input type="radio"/> or | <input type="text"/> | < HTMs < | <input type="text"/> | <input type="radio"/> and <input type="radio"/> or | <input type="text"/> | < %GC < | <input type="text"/> |
| <input type="radio"/> and <input type="radio"/> or | <input type="text"/> | < nucl. length < | <input type="text"/> | <input type="radio"/> and <input type="radio"/> or | <input type="text"/> | < pep. length < | <input type="text"/> |
| <input type="radio"/> and <input type="radio"/> or | <input type="text"/> | < % id. simple < | <input type="text"/> | <input type="radio"/> and <input type="radio"/> or | <input type="text"/> | < % modelled < | <input type="text"/> |

E_value range in closed organisms

| | | | |
|--|--|--|---|
| <input type="radio"/> and <input type="radio"/> or | <input type="text" value="'Agrobacterium tumefaciens CS9' with e_value = 10"/> | <input type="radio"/> and <input type="radio"/> or | <input type="text" value="'Bartonella henselae' with e_value = 10"/> |
| <input type="radio"/> and <input type="radio"/> or | <input type="text" value="'Bartonella quintana' with e_value = 10"/> | <input type="radio"/> and <input type="radio"/> or | <input type="text" value="'Brucella abortus' with e_value = 10"/> |
| <input type="radio"/> and <input type="radio"/> or | <input type="text" value="'Brucella suis' with e_value = 10"/> | <input type="radio"/> and <input type="radio"/> or | <input type="text" value="'Caulobacter crescentus' with e_value = 10"/> |
| <input type="radio"/> and <input type="radio"/> or | <input type="text" value="'Escherichia coli K12' with e_value = 10"/> | <input type="radio"/> and <input type="radio"/> or | <input type="text" value="'Mycobacterium lei' with e_value = 10"/> |
| <input type="radio"/> and <input type="radio"/> or | <input type="text" value="'Rhizobium leguminosarum' with e_value = 10"/> | <input type="radio"/> and <input type="radio"/> or | <input type="text" value="'Stenotrophomonas maltophilia' with e_value = 10"/> |
| <input type="radio"/> and <input type="radio"/> or | <input type="text" value="'Rhopilema eysenhardti pulchra' with e_value = 10"/> | | |

Output format

Bimodal distribution of the predicted pI.

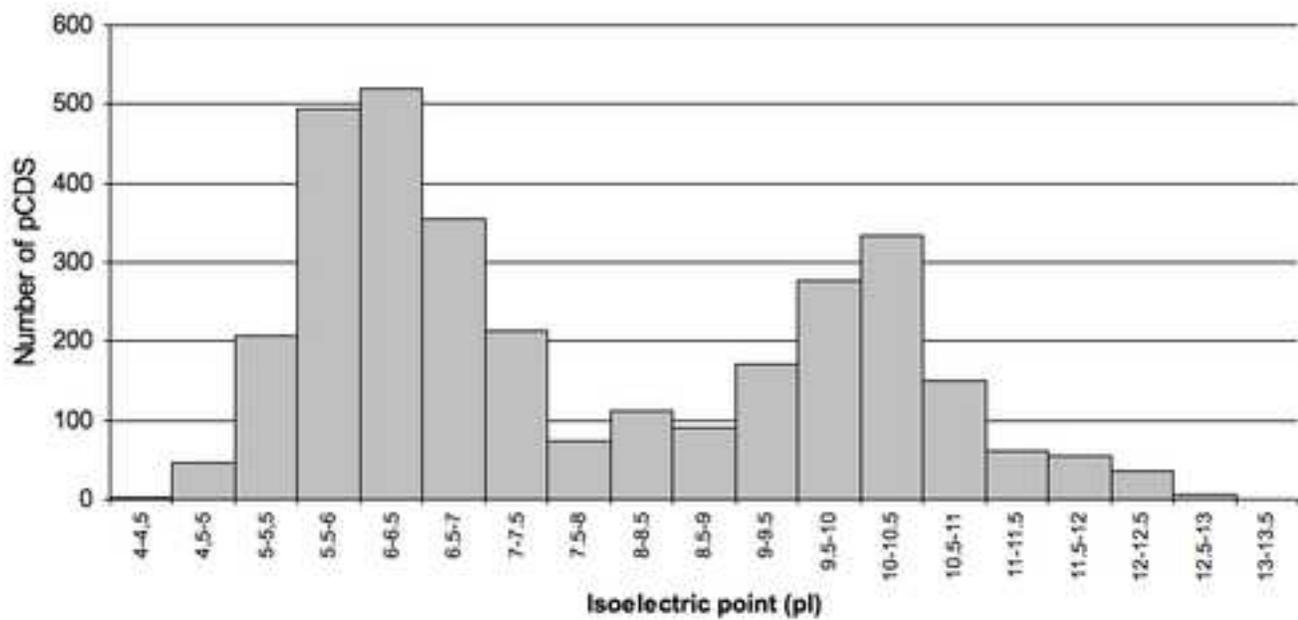


Table 1. Characteristics of the 10 species phylogenetically close to *Brucella sp.* present in the database.

Escherichia coli is used as a model Gram-negative bacteria (from EMBL-EBI and quoted authors).

| Brucella biovars | Reference | Life cycle | Host | Genome | Orfs |
|--|---|-------------------------|--|---|-------------|
| <i>Brucella abortus</i> | Sanchez <i>et al.</i> , 2001; Halling <i>et al.</i> , 2005 | pathogenic | human and livestock | 2,1 MB and 1,2 Mb circular chromosomes | 3,076 |
| <i>Brucella melitensis</i> | Delvecchio <i>et al.</i> , 2002 | pathogenic | human and livestock | 2,1 MB and 1,2 Mb circular chromosomes | 3,197 |
| <i>Brucella suis</i> | Paulsen <i>et al.</i> , 2002 | pathogenic | human and livestock | 2,1 MB and 1,2 Mb circular chromosomes | 3,256. |
| Close to Brucella | | | | | |
| <i>Agrobacterium tumefaciens</i> C58 | Wood <i>et al.</i> , 2001; Goodner <i>et al.</i> , 2001 | pathogenic | numerous plants (crown, roots and stems) | 2,8 MB circular chromosome; 0,54 Mb Plasmid pMLa; 0,2 Mb Plasmid pMLa | 5,304 |
| <i>Bartonella henselae</i> str. Houston-1, | Alsmark <i>et al.</i> , 2004 | pathogenic | human and cat | 1,9 Mb circular chromosome | 1,464 |
| <i>Bartonella quintana</i> str. Toulouse | Alsmark <i>et al.</i> , 2004 | pathogenic | human specific | 1,6 Mb circular chromosome | 1,137 |
| <i>Caulobacter crescentus</i> | Nierman <i>et al.</i> , 2001 | | dilute aquatic environment | 1,2 Mb circular chromosome | 3,718 |
| <i>Mesorhizobium loti</i> | Kaneko <i>et al.</i> , 2000 | symbiotic | nitrogen-fixing soil plant | 7 Mb circular chromosomes; 0,4 Mb Plasmid pMLa; 0,2 Mb Plasmid pMLa | 7,255 |
| <i>Sinorhizobium meliloti</i> | Barnett <i>et al.</i> , 2001; Galibert <i>et al.</i> , 2001 | symbiotic | nitrogen-fixing soil plant | 3,6 Mb chromosome and two megaplasmids, pSyma (1,3 Mb) and pSymb (1,7 Mb) | 6,148 |
| <i>Rhizobium leguminosarum</i> | unpublished data | symbiotic | nitrogen-fixing soil plant | 5 Mb circular chromosomes and six plasmid | - |
| <i>Rhodopseudomonas palustris</i> | Larimer <i>et al.</i> , 2004 | metabolically versatile | soils and water | 5,5 Mb circular chromosome; 8,427 bp plasmid pRPA | 4,798 |
| GRAM - | | | | | |
| <i>Escherichia coli</i> K12 | Blattner <i>et al.</i> , 1997 | pathogenic | human | 4,6 Mb circular chromosome; 1 MB plasmid F | 4,338 |

Table 2. Machines dedicated to the *B. melitensis* genome database.

| | Processors | CPU | RAM | OS |
|-----------------------------|--------------------------------|------------|------------|----------------------|
| Silicon Graphics Octane duo | 2 x R10000 MIPS | 225 Mhz | 512 MB | IRIX 6.5 |
| Priminfo Xeon Server | 2 x 32 bits Intel Pentium 4 | 2.2 Ghz | 2 GB | Red Hat Linux 7.3 |
| 6 nodes Priminfo Xeon | 2 x 32 bits Intel Pentium 4 | 2.2 Ghz | 2 GB | Red Hat Linux 7.3 |

Accepted Manuscript