



Murdoch
UNIVERSITY

MURDOCH RESEARCH REPOSITORY

This is the author's final version of the work, as accepted for publication following peer review but without the publisher's layout or pagination.

The definitive version is available at

<http://dx.doi.org/10.1016/j.ijpara.2011.11.006>

Bellgard, M.I., Moolhuijzen, P.M., Guerrero, F.D., Schibeci, D., Rodriguez-Valle, M., Peterson, D.G., Dowd, S.E., Barrero, R., Hunter, A., Miller, R.J. and Lew-Tabor, A.E. (2011)
CattleTickBase: An integrated Internet-based bioinformatics resource for Rhipicephalus (Boophilus) microplus. *International Journal for Parasitology*, 42 (2). pp. 161-169.

<http://researchrepository.murdoch.edu.au/6300/>

Copyright: © 2011 Elsevier Ltd.

It is posted here for your personal use. No further distribution is permitted.

Accepted Manuscript

CattleTickBase: An integrated Internet-based bioinformatics resource for *Rhipicephalus (Boophilus) microplus*

Matthew I. Bellgard, Paula M. Moolhuijzen, Felix D. Guerrero, David Schibeci, Manuel Rodriguez-Valle, Daniel G. Peterson, Scot E. Dowd, Roberto Barrero, Adam Hunter, Robert J. Miller, Ala E. Lew-Tabor

PII: S0020-7519(11)00281-5
DOI: [10.1016/j.ijpara.2011.11.006](https://doi.org/10.1016/j.ijpara.2011.11.006)
Reference: PARA 3344



To appear in: *International Journal for Parasitology*

Received Date: 9 September 2011
Revised Date: 16 November 2011
Accepted Date: 17 November 2011

Please cite this article as: Bellgard, M.I., Moolhuijzen, P.M., Guerrero, F.D., Schibeci, D., Rodriguez-Valle, M., Peterson, D.G., Dowd, S.E., Barrero, R., Hunter, A., Miller, R.J., Lew-Tabor, A.E., CattleTickBase: An integrated Internet-based bioinformatics resource for *Rhipicephalus (Boophilus) microplus*, *International Journal for Parasitology* (2011), doi: [10.1016/j.ijpara.2011.11.006](https://doi.org/10.1016/j.ijpara.2011.11.006)

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.

1 **CattleTickBase: An integrated Internet-based**
2 **bioinformatics resource for *Rhipicephalus (Boophilus)***
3 ***microplus* ★**

4

5 Matthew I. Bellgard^{a, b, 1}, Paula M. Moolhuijzen^{a, b, 1}, Felix D. Guerrero^{c, *}, David
6 Schibeci^a, Manuel Rodriguez-Valle^{b, d}, Daniel G. Peterson^e, Scot E. Dowd^f, Roberto
7 Barrero^a, Adam Hunter^a, Robert J. Miller^g, Ala E. Lew-Tabor^{a, b, d}

8 ^a *Centre for Comparative Genomics, Murdoch University, Perth, WA 6150, Australia*

9 ^b *Cooperative Research Centre for Beef Genetic Technologies, Armidale, NSW,*

10 *Australia*

11 ^c *USDA-ARS Knippling Bushland US Livestock Insect Research Laboratory, 2700*

12 *Fredericksburg Rd., Kerrville, TX 78028, USA*

13 ^d *Queensland Alliance for Agriculture, Food & Innovation, The University of*

14 *Queensland and Dept. of Employment, Economic Development & Innovation, P. O.*

15 *Box 6097, St. Lucia 4067 QLD, Australia*

16 ^e *Department of Plant & Soil Sciences and Life Sciences & Biotechnology Institute,*

17 *Mississippi State University, 117 Dorman Hall, Box 9555, Mississippi State, MS*

18 *39762, USA*

19 ^f *Molecular Research, 503 Clovis Road, Shallowater, TX, 79363, USA*

20 ^g *USDA-ARS Cattle Fever Tick Research Laboratory, 22675 North Moorefield Road,*

21 *Building 6419, Edinburg, TX 78541, USA*

22 ¹These authors contributed equally.

23 *Corresponding author. Felix D. Guerrero, USDA-ARS Knippling Bushland US

24 Livestock Insect Research Laboratory, 2700 Fredericksburg Rd., Kerrville, TX 78028,

25 USA. Tel.: +1-830-792-0327; fax: +1-830-792-0314.

26 *E-mail address:* Felix.Guerrero@ars.usda.gov

27

28 ★Note: Nucleotide sequence data reported in this paper are available in GenBank

29 under accession numbers [HN108288-HN118367](#), [HM748958-HM748967](#),

30 [HN108288-HN118367](#).

31

32 Note: Supplementary data associated with this article.

33

34 **ABSTRACT**

35 The *Rhipicephalus microplus* genome is large and complex in structure, making it
36 difficult to assemble a genome sequence and costly to resource the required
37 bioinformatics. In light of this, a consortium of international collaborators was formed
38 to pool resources to begin sequencing this genome. We have acquired and assembled
39 genomic DNA into contigs that represent over 1.8 Gigabase pairs (Gbp) of DNA from
40 gene-enriched regions of the *R. microplus* genome. We also have several datasets
41 containing transcript sequences from a number of gene expression experiments
42 conducted by the consortium. A web-based resource was developed to enable the
43 scientific community to access our datasets and conduct analysis through a web-based
44 bioinformatics environment called YABI. The collective bioinformatics resource is
45 termed CattleTickBase. Our consortium has acquired genomic and transcriptomic
46 sequence data at approximately 0.9X coverage of the gene-coding regions of the *R.*
47 *microplus* genome. The YABI tool will facilitate access and manipulation of cattle
48 tick genome sequence data as the genome sequencing of *R. microplus* proceeds.
49 During this process the CattleTickBase resource will continue to be updated.

50

51 *Keywords:* Cattle tick; Transcriptome; Genome project; Bioinformatics

52

53

54 **1. Introduction**

55 The global cattle population is estimated at approximately 1 billion. Of this
56 population, 80% inhabit areas that have been considered suitable habitat for ticks and
57 tick-borne diseases (Snelson, 1975). The cattle tick *Rhipicephalus (Boophilus)*
58 *microplus* is considered the most significant cattle parasite in the world, having
59 established populations in most of the world's tropical and subtropical countries. This
60 tick causes blood loss and physical damage to hides of infested animals. In addition,
61 *R. microplus* is the vector for several bovine diseases, including babesiosis (caused by
62 protozoan species *Babesia bovis* and *Babesia bigemina*) and anaplasmosis (caused by
63 the rickettsia *Anaplasma marginale*), with severe impact on agricultural systems
64 globally (de Castro, 1997). Economic losses to cattle producers from ticks and tick-
65 borne diseases are US\$13-18 billion globally on an annual basis (de Castro, 1997).
66 Annual losses attributable to *R. microplus* in Brazil and Australia alone are estimated
67 at US\$2 billion (Grisi et al., 2002) and AUS\$175 million (Playford et al., 2005),
68 respectively.

69 Ticks are believed to be among the most ancient terrestrial arachnids and
70 possibly the earliest organisms to have evolved blood-feeding capabilities (Mans and
71 Neitz, 2004). *Rhipicephalus microplus* is a single-host species and has evolved such
72 that it must maintain sustained contact with its host during the life stages, from the
73 attached and feeding larva through to the fully engorged female. This period of
74 attachment typically lasts approximately 3 weeks with some variation depending on
75 environmental conditions. The species has developed a unique means of avoiding the
76 host animal's immune responses during infestation (Wikel, 1999) and *R. microplus*
77 salivary gland extracts have been shown to have an immunosuppressive effect on the
78 bovine host (Turni et al., 2004, 2007). The tick must also respond to many

79 microorganisms, both symbiotic and parasitic, from the external environment or those
80 ingested through feeding-associated activities (Andreotti et al., 2011).

81 With this interplay between bovine host, tick and microbiota, determining the
82 whole genome sequence of *R. microplus* will greatly advance tick gene discovery,
83 enable a better understanding of tick-host-pathogen immunology and provide insight
84 on how the cattle tick responds to environmental perturbations, including pressures
85 from moisture and temperature extremes and acaricidal applications. The *R. microplus*
86 genome size is estimated to be 7.1 Gigabase pairs (Gbp), more than twice the size of
87 the human genome, and consists of greater than 70% repetitive DNA (Ullmann et al.,
88 2005). It is therefore a challenge for de novo assembly, even with contemporary DNA
89 sequencing technologies. A 4X shotgun coverage genome sequence for the
90 blacklegged tick, *Ixodes scapularis*, is available (Lawson et al., 2009) and is the only
91 reported tick genome sequence to date. The version 1.1 sequence assembly consists of
92 369,492 supercontigs, totalling 1.76 Gbp with a supercontig N50 size of 72kb.

93 From a taxonomic perspective, both *R. microplus* and *I. scapularis* are
94 classified as hard ticks. There are two lineages of hard ticks, the Prostriata, which
95 consists of the single genus *Ixodes* containing approximately 250 species, and the
96 Metastricata, which consists of approximately 464 species from several genera
97 including *R. microplus* (Barker and Murrell, 2004). Given the sequence divergence
98 between *R. microplus* and *I. scapularis* (Guerrero et al., 2006), gene discovery efforts
99 solely using *I. scapularis* as the model tick genome would prove limiting for *R.*
100 *microplus* research efforts.

101 Towards a goal of generating a genetic resource for this economically
102 important tick species, efforts have focused on a combination of sequencing

103 strategies. The goal of this project was to maximize the utility of the data that could
104 be generated with the resources available. To date, these strategies include Cot-
105 filtered genomic DNA sequencing, bacterial artificial chromosome (BAC)-end
106 sequencing (BES), targeted whole BAC sequencing, whole transcriptome sequencing
107 and small RNA sequencing. Presently, we have acquired, assembled and annotated
108 over 2 Gb of sequence data. This is comprised of 1.7 Gb of assembled contigs from
109 three Cot reassociation experiments. These Cot experiments utilised methodologies to
110 select randomly sheared genomic DNA for fractions depleted in highly repetitive
111 sequences and enriched for putative gene coding regions (Guerrero et al., 2010). Also
112 available are three transcriptome library assemblies (21 Mb) representing over 33,000
113 transcripts. Integrating data generated from these various approaches is already
114 starting to provide new insights into the very large and complex cattle tick genome.
115 This paper provides an overview of the coordinated cattle tick genome sequence
116 resource as well as a new internet-based bioinformatics resource that is designed to
117 integrate our various genomic and transcriptomic datasets. This enables the cattle tick
118 research community to access and analyse genomic and transcriptomic data at a single
119 on-line resource, similar to approaches used by insect and worm researchers e.g.
120 FlyBase, WormBASE (Harris and Stein, 2006; Drysdale, 2008). Examples are
121 provided of these new genomic analysis tools and how they can be utilised by the
122 research community to understand the genomic structure, organisation and content of
123 the cattle tick genome.

124

125 **2. Materials and methods**126 *2.1. Source of tick materials*

127 For the USA ticks, genomic and Cot DNA were extracted from eggs of the *R.*
128 *microplus* Deutsch strain, f7, f10, f11 and f12 generations. These were pooled and a
129 total of 10 g was used to purify very high molecular weight genomic DNA (Guerrero
130 et al., 2010). This strain was started from only a few individual engorged females
131 collected from a 2001 tick outbreak in South Texas. Although the strain has been
132 inbred since its creation in 2001, it is not genetically homogeneous. For the Australian
133 ticks, the larvae and fully engorged N strain of Australian *R. microplus* were utilized
134 in these analyses. The N strain is maintained by the Biosecurity laboratories at the
135 Department of Employment, Economic Development and Innovation (DEEDI),
136 Queensland, under controlled conditions of 28 °C and 80% relative humidity prior to
137 bovine infestation (Stewart et al., 1982).

138

139 *2.2. Sequencing and assembly*

140 For the BAC library synthesis, approximately 2 g of larvae from the f8
141 generation of the Deutsch strain were used by Amplicon Express Inc. (Pullman, WA,
142 USA) to isolate genomic DNA partially digested with *Mbo*I to synthesize a BAC
143 library of approximately 0.8X coverage (Guerrero et al., 2010). Subsequently, a
144 second library of 2.4X coverage was synthesized from genomic DNA partially
145 digested with *Hind*III. Five BAC assemblies, BM-074-Random-F12, BM-077-
146 Random-J09, BM-129-CzEst9-N14, BM-066-M07, BM-077-G20, are as described by
147 Guerrero et al. (2010). The remaining 10 BAC sequences were trimmed for vector
148 and bacterial contamination by phred-phrap software (Ewing and Green, 1998)

149 cross_match, with options set at minmatch 12 and minscore 20. Contig order and
150 orientation were based on Phrapview paired end reads.

151 Total *R. microplus* genomic DNA was prepared and processed by three Cot
152 filtration experiments to enrich for single/low-copy and moderately repetitive DNAs.
153 Cot-filtered DNA was sequenced using 454 FLX and Titanium pyrosequencing
154 (Research and Testing Laboratory, Lubbock, TX, USA). Methods are as described in
155 Guerrero et al. (2010).

156 The filtered genomic DNA (total number of reads 7,289,230 and total number
157 of bases 1,798,400,445) was de novo assembled using the Newbler assembler for 454
158 reads (Margulies et al., 2005) with default settings. All contigs (745,975 sequences)
159 and BAC end sequences (BES) (GenBank Accession Numbers HN108288-
160 HN118367) were then assembled with Cap3 (Huang and Madan, 1999) default
161 settings. Whole Genome Shotgun (WGS) project ADMZ02000000 is the result of this
162 two-step assembly.

163

164 2.3. BAC and Cot read alignment

165 BAC and Cot read alignments were carried out with BWA-SW (Li and
166 Durbin, 2010) for long reads as our average read length was 245 bp. A mapping
167 accuracy >99% was expected with a mapping quality (MapQ) of 10 and sensitivity
168 (Z) of 100.

169

170 2.4. Gene predictions

171 Gene predictions for the BAC sequences and WGS were made with GenScan
172 version 1.0 (Burge and Karlin, 1997) for default 'optimal exons', parameters for
173 human/vertebrates and coding sequences (CDS) option. BAC predicted gene and
174 AutoFACT (Koski et al., 2005) annotation can be found in Supplementary Table S1
175 and Supplementary Fig. S1.

176

177 2.5. Repeat analysis

178 Repeat sequences and rRNA were identified using RepeatMasker version
179 3.2.6 Smit, A.F.A., Hubley, R., Green, P., 2004. RepeatMasker Open-3.0. 1996-2010
180 <<http://www.repeatmasker.org>>.) with parameters set up for the arthropod clade of
181 input sequences.

182

183 2.6. RNA Searches: tRNA and rRNA

184 Transfer RNA (tRNA) searches were conducted with trnscan version 1.23
185 (Schattner et al., 2005) and rRNA using rnammer version 1.2 (Lagesen et al., 2007).

186

187 2.7. Sequence comparative analysis

188 The genomic data set *Ixodes scapularis* SUPERCONTIGS-
189 Wikel.IscaW1.fa.gz (dataset downloaded from VectorBase (Lawson et al., 2009)
190 Date: Sept. 21, 2010) was aligned to *R. microplus* assembled Cot DNA (GenBank
191 WGS division second version, ADMZ02000000) using BLASTn (Altschul et al.,
192 1990). Homologous regions of interest were selected at an expected value < 1e-50.

193 The *R. microplus* BAC sequences submitted to GenBank: **HM748958-**
194 **HM748967** and five BACs as described in Guerrero et al. (2010) were aligned to the
195 BES submitted to the Genome Survey Sequence (GSS) division of GenBank:
196 **HN108288-HN118367** using BLAT with 70% identity, a length greater than 100 bp,
197 and the option 'fastMap'.

198 The *R. microplus* BES, predicted BAC gene content and new transcript
199 sequences were comparatively aligned to Dana Farber Cancer Institute (DFCI) gene
200 indices, IsGI version 3.0 and *BmiGI* version 2.1 (Quackenbush et al., 2001), NCBI
201 RefSeq mRNA, and the Subtraction Library clones as described previously (Lew-
202 Tabor et al., 2009) using BLASTn at an expected value $<1e-10$, and to NCBI RefSeq
203 Protein and GenPeptide, *iscapularis*.PEPTIDES-IscaW1.1.fa (VectorBase Date: Sep
204 21, 2010, Lawson et al., 2009) using BLASTx at an expected value $<1e-10$. The
205 collective *R. microplus* transcriptome (*RmiTr* Version 1.0, Table 1) and *I. scapularis*
206 (DFCI IscGI) sequences were searched using tBLASTx at an expected value of $1e-05$
207 to uniref100, and orthologous sequences were determined for those alignments that
208 had greater than or equal to 60% *R. microplus* sequence coverage and an amino acid
209 conservation greater than or equal to 30%. Transcript alignments were also made
210 using BLAT (Kent, 2002) to a comparative species with 70% identity and a length
211 greater than 100 bp.

212

213 2.8. Transcriptome sequencing

214 For transcriptome sequencing, female adult dissected gut and 'frustrated'
215 larvae were prepared as described previously from Australian N strain ticks (Lew-
216 Tabor et al., 2009). The frustrated larval sample contains larvae placed in a gas-

217 permeable bag taped directly to the host animal feed source. Thus, the larvae are able
218 to sense the presence of the host but the bag presents a barrier that prevents
219 attachment and feeding. Approximately 30 mg of total RNA from each of these
220 samples were collected for high-throughput sequencing of the tick transcriptome
221 using the Illumina/GA single-end reads format as described previously (Mortazavi et
222 al., 2008).

223

224 2.9. Transcriptome assembly and clustering

225 The de novo transcriptome assembly of the adult female gut and frustrated
226 larvae transcriptomes using 60 bp single-end Illumina/GA reads was conducted using
227 Abyss (Birol et al., 2009) with *k*-mer sizes ranging from 36 to 64. The assembled
228 contigs were then clustered using cap3 (Huang and Madan, 1999) with a 98%
229 sequence identity threshold and an overlap region of at least 30 bases to remove
230 transcript redundancy. Non-redundant sets of transcripts for the two libraries can be
231 found on the CattleTickBase website.

232 The *RmiTr* version 1.0 data set contains sequences from DFCI_bmigi.V2.1,
233 adult female gut transcriptome, frustrated larvae transcriptome and *R. microplus*
234 subtraction library (Lew-Tabor et al., 2009), which were clustered into contigs using
235 cap3 (Huang and Madan, 1999) with following the options: -p 99.99999 -m 1 -n -10
236 -g 1 -b 16 -y 6.

237

238 2.10. YABI

239 The YABI application consists of a front-end web application responsible for
240 the user interface. Users create a secure account and are free to access the datasets and
241 analysis tools available within the system. Users can create bioinformatics pipelines
242 from the available tools. The tools currently available for sequence analysis include:
243 similarity/homology searches, feature prediction, high throughput downstream
244 analysis, assembly and annotation. Datasets (including other tick-related GenBank
245 bioprojects) in YABI are updated from GenBank/EMBL/DDBJ and VectorBase at
246 regular intervals. Individual dataset contributors can deposit and update data sets by
247 contacting yabi@ccg.murdoch.edu.au with an option to have a secure account to
248 conduct their analysis within their own teams. The Centre for Comparative Genomics
249 (CCG), Murdoch University, Australia is committed to supporting the bioinformatics
250 aspects of the *R. microplus* project. The CCG houses a supercomputer (currently
251 ranked number 87 in the world - <http://www.top500.org/list/2010/11/100>) and
252 provides support to national and international bioinformatics and other high-end
253 science-based activities. The CCG develops and deploys sophisticated software
254 solutions, supports and conducts a diverse range of bioinformatics analysis. Requests
255 and suggestions can be made by contacting info@ccg.murdoch.edu.au. The *R.*
256 *microplus* datasets currently available for sequence similarity searching and other
257 bioinformatics analyses are summarised in Table 1.

258

259 3. Results

260 3.1. *Rhipicephalus microplus* datasets currently available

261 The nine datasets that are available for access and further analysis on the
262 CattleTickBase website are summarised in Table 1 and described below. For Dataset

263 1, total *R. microplus* genomic DNA was prepared and processed by three Cot
264 filtration experiments to enrich for single/low-copy and moderately repetitive DNA. It
265 was anticipated that the DNA obtained via this process would remove a significant
266 portion of highly repetitive DNA and contain predominantly gene rich regions. The
267 resulting DNA fragments ranged in size from 250 to 600 bp and were sequenced in
268 three separate sequencing experiments with one experiment using six runs of 454
269 FLX pyrosequencing (Guerrero et al., 2010) and two experiments each using three
270 Titanium 454 runs. The data from the latter two experiments have been deposited in
271 GenBank SRA, submission: [SRA012677.4/SID00001](#). Approximately 1.8 Gbp of
272 sequence were generated and assembled. This assembly (Dataset 2) was submitted to
273 GenBank Whole Genome Sequencing Project; under the United States Department of
274 Agriculture, Agricultural Research Service (USDA-ARS) *R. microplus* Project ID
275 46685 assigned Project accession [ADMZ0000000](#). The most recent version for this
276 project reported in this paper has the accession number [ADMZ0200000](#); this
277 submission consists of 175,208 contig sequences of average contig size 825 bp and a
278 maximum contig length of 9,681 bp. These sequences have been submitted under
279 GenBank Accession Numbers [ADMZ02000001](#)-[ADMZ02175208](#). Based on the
280 estimate that the *R. microplus* genome size is approximately 7.1 Gbp, the assembled
281 Cot DNA dataset represents approximately 2% of the *R. microplus* genome. We
282 undertook a comparative analysis of this Cot DNA dataset with the *BmiGI* Version
283 2.1 Gene Index ([http://compbio.dfci.harvard.edu/cgi-](http://compbio.dfci.harvard.edu/cgi-bin/tgi/gimain.pl?gudb=b_microplus)
284 [bin/tgi/gimain.pl?gudb=b_microplus](http://compbio.dfci.harvard.edu/cgi-bin/tgi/gimain.pl?gudb=b_microplus); Quackenbush et al., 2001; Guerrero et al., 2005)
285 to confirm that the Cot DNA filtration process filtered for gene rich regions. When we
286 aligned the 175,208 Cot DNA contigs to the 14,586 *BmiGI* transcript contigs, 52% of
287 the *BmiGI* entries had a match (>85% identity) to at least one contig in the Cot

288 dataset. Thus, even though the coverage of the Cot DNA over the entire 7.1 Gbp *R.*
289 *microplus* genome was very low (2%), the coverage of the *BmiGI* transcriptome
290 dataset was high (52%). This also indicates that deeper sequencing of the Cot-selected
291 DNA would be warranted to obtain fuller coverage of the gene-encoding regions, as
292 48% of the *BmiGI* entries did not have a match in the Cot DNA dataset.

293 For Dataset 3, a 3.2x coverage BAC library was used to sequence 10,582 BES
294 resulting in 7,290,530 bp of sequence. The amount of BES data generated represents
295 approximately 0.1% of the *R. microplus* genome. BES greater than 500 bp in size
296 have been deposited in GenBank GSS under Accession Numbers **HN108288-**
297 **HN118367**. Approximately 70% matched with the Cot DNA Dataset 2. Similarity
298 sequence analysis between BES and the *BmiGI* Version 2.0 reveals that a total of 502
299 BES (4.7%) aligned, at a length greater than 100 nucleotides and at greater than 90%
300 identity (percent identity, PID), to 224 *BmiGI* entries or approximately 1.6% of the
301 *BmiGI*.

302 The 10,582 BES were compared with the *I. scapularis* genome sequence, the
303 nearest taxonomic whole genome sequencing project, comprising 369,492 scaffolds
304 (supercontigs). Only 58 BES aligned with greater than 80% BES coverage and 80%
305 PID. The number of BES found in comparative searches at specified thresholds were:
306 2,418 (23%) at an expected alignment value (e-value) of 1e-05 to the NCBI protein
307 non-redundant (nr) database, 416 (4%) at an e-value 1e-20 to *I. scapularis* proteins
308 (Lawson et al., 2009; Megy et al., 2009), 2,559 (24%) BES at an e-value 1e-20 to *R.*
309 *microplus* gene indices (Guerrero et al., 2005), and 134 (1.2%) to the *I. scapularis*
310 gene index (Ribeiro et al., 2006) at an e-value of 1e-20. The protein functional
311 analyses by Gene Ontology (GO) classification are presented in Supplementary Table
312 S2 and Supplementary Fig. S2.

313 For Dataset 4, 15 BAC clones were selected and sequenced to completion; 13
314 were based on hybridization of BAC library filter arrays to probes from known
315 transcripts of interest involved in tick feeding and acaricide resistance and two were
316 randomly selected BACs. These 15 BACs were subjected to sequencing and de-novo
317 assembly. The sequencing of five of the BACs has been reported (Guerrero et al.,
318 2010), and the sequences from the remaining 10 BACs have been deposited in
319 GenBank (Accession Numbers HM748958-HM748967). Gene prediction analysis
320 found 180 predicted genes with 919 predicted exons with an average exon size of 311
321 bp (Supplementary Fig. S2). From these predicted genes, there are 166 full-length
322 genes comprised of 145 multiple exon genes containing both the initial and terminal
323 predicted exons and 21 are single exon genes. Genes of particular interest to our
324 research group that were identified in these BACs include cytochrome P450
325 (Guerrero et al., 2010), the permethrin-degrading carboxylesterase CzEst9 (Guerrero
326 et al., 2010), papilin (Moolhuijzen et al., 2011), transmembrane protein 215-like,
327 transpanin and serpin (Moolhuijzen et al., 2011). The Cot DNA dataset was also
328 mapped to the BACs. Table 2 provides a summary of all 15 BACs, including the
329 length of each assembled BAC, the number of Cot-selected DNA reads from Dataset
330 1 that mapped to each BAC using the next generation alignment tool BWA (Li and
331 Durbin, 2010) with an alignment score > 10 , the percentage of coverage of BAC
332 sequence provided by the Cot reads, and the number of predicted genes and exons. As
333 shown in Table 2, the percentage of Cot DNA coverage over these BACs ranged from
334 6% (BM-012-E08) to as high as 82% (BM-005-B21). Thirteen BACs had an average
335 GC content between 44-48% (Table 2). One BAC (BM-012-E08) had a GC content of
336 58% and another BAC (BM-074-Random-F12) had a GC content of 40%. BAC BM-
337 012-E08 contains rDNA and intergenic spacer and a significant proportion of highly

338 repetitive DNA. Present in the 15 BACs are the Ruka SINE elements identified
339 previously in *Rhipicephalus appendiculatus* (Sunter et al., 2008) and a range of other
340 interspersed repeats such as LINEs L2 and R1, LTR Gypsy and DNA transposons as
341 found by RepeatMasker (Smit, A.F.A., Hubley, R., Green, P., 2004. RepeatMasker
342 Open-3.0. 1996-2010 <<http://www.repeatmasker.org>>.). BAC BM-005-G14
343 contained a number of Ruka elements, which are a novel SINE element that was
344 reported to occur frequently in both genomic and transcribed ixodid ticks (Sunter et
345 al., 2008). A single 195 bp length Ruka element found in BM-004-G14 had 91%
346 sequence identity to the 138 bp *R. appendiculatus* (*Rap*) Ruka (Genbank:
347 **EU018139.1** complement (9,947-10,084 bp)). A total of 70 Ruka (> 100 bp) elements
348 were found in the *Rmi*Tr Version 1.0 transcriptome (Table 1 Dataset 8), while 409
349 were found in the assembled Cot DNA.

350 The remaining datasets are expression related. Dataset 5 consists of the *R.*
351 *microplus* Gene Index *Bmi*GI Version 2.1 (Quackenbush et al., 2001) entries that
352 were extended by comparative analysis with the assembled Cot-selected genomic
353 DNA (Guerrero et al., 2010). Datasets 6 and 7 consist of transcriptome sequence from
354 adult female tick gut and ‘frustrated’ larvae libraries, respectively. These datasets
355 were de novo assembled using Illumina short reads. We wanted to compare the
356 transcriptome data contained in our datasets 6 and 7 with the *Bmi*GI Version 2.1 to
357 determine whether our data presented new transcripts that could be added to the
358 current *Bmi*GI to create a more comprehensive resource. Fig. 1 provides an overview
359 of the CAP3 sequence clustering between datasets 6 and 7 and *Bmi*GI Version 2.1 at
360 greater than 98% identity and default settings. The Venn diagram shows that datasets 6
361 and 7 (Aus *Rmi* as noted in Fig. 1) contain 9,652 contigs not found in the *Bmi*GI
362 version 2.1. Dataset 8 is an updated *R. microplus* transcript data set made up of the

363 combination of the original *BmiGI* version 2.1 and Datasets 6 and 7, plus the
364 Subtraction Library data (GenBank: GO253189.1- GO253184.1, GE650059.1-
365 GE650181.1) and is referred to as *RmiTr* version 1.0, containing 28,893 sequences.
366 Finally, Dataset 9 consists of data from microarray experiments using Nimblegen
367 arrays (Rodriguez-Valle et al., 2010; Saldivar et al., 2008) to characterize the *R.*
368 *microplus* transcriptome responses to attaching and feeding upon *Bos indicus* and *Bos*
369 *taurus* cattle (Rodriguez-Valle et al., 2010).

370

371 3.2. *CattleTickBase* web resource

372 The *CattleTickBase* web resource provides a central point for the cattle tick
373 research community to access genome and transcriptome information on *R.*
374 *microplus*. The home page is shown in Fig. 2. This website contains a summary of the
375 datasets that have been generated to date and are available for download. It also
376 contains links to precomputed results, which are either figures from supplementary
377 material from publications or the datasets precomputed into the GBrowse genome
378 browser (Stein et al., 2002). To date, precomputed results include a summary of the
379 protein hits from similarity searches contained within the sequenced BACs (Dataset
380 4) as well as genes that were extended by the Cot DNA contigs that can be viewed
381 within GBrowse (Dataset 5).

382 A unique feature of our genome project web resource is the ability for the
383 research community to conduct their own sophisticated bioinformatics analysis
384 online. The open source YABI system (<http://ccg.murdoch.edu.au/yabi>) consists of a
385 front-end web application responsible for the user interface. Users create a secure
386 account and are free to access the datasets and analysis tools available within the

387 system. Users can create customized bioinformatics pipelines from the available tools
388 and scripts that capture provenance information of the tools used, such as parameters
389 used for each tool, and outputs of tools generated at each step. Fig. 3 shows the layout
390 of the Design tab within YABI for constructing analysis workflows. The tools
391 available are presented in the menu on the left. In the centre, dragging and dropping
392 tools from the menu onto the workflow panel will construct a workflow. The example
393 in Fig. 3 shows a screen shot of the creation of an automated workflow to search for
394 G-Protein Coupled Receptors (GPCRs) in frustrated larvae transcripts that make up
395 Dataset 7 (Table 1). GPCRs are an interesting family of membrane proteins that
396 perform a range of critical biological functions in eukaryotes. Approximately 40% of
397 the prescription pharmaceuticals target GPCRs (Filmore, 2004) and they are attractive
398 as potential targets for acaricidal product development. Predicted open reading frames
399 (ORFs) from the 6,082 assembled contig transcriptome from Dataset 7 were
400 determined and the resulting sequences were searched using GPCRHMM (Wistrand
401 et al., 2006). From a total of 17, 230 predicted ORFs (≥ 200 bp), six ORFs were
402 predicted to encode GPCRs (data not shown). A screen cast of this workflow can be
403 found at (<http://ccg.murdoch.edu.au/yabi>). In CattleTickBase, the web-based analysis
404 workflow environment (YABI) enables researchers to create sophisticated
405 bioinformatics analysis pipelines from a diverse range of tools. For example, tools are
406 available for multiple sequence alignment and assembly (including data from next-
407 generation sequencing technologies), sequence searches, predictions for genes,
408 rRNAs, tRNAs and small RNAs, repeat masking and various annotations. Also, the
409 calculation of the Codon Adaptation Index for a given nucleotide sequence, given a
410 reference codon usage table, can be used for predicting the level of expression of a

411 given gene and for making comparisons among codon usage in different organisms.

412 *R. microplus* codon usage tables will be available for use in YABI.

413

414 **4. Discussion**

415 This paper provides an overview of the status of our collaborative efforts
416 towards sequencing the *R. microplus* genome. Currently there is a scarcity of tick
417 genome sequence and that limits progress in projects designed to address the
418 significant worldwide impacts that *R. microplus* presents to cattle producers, both
419 large and small. With the difficulties associated with sequencing the large and
420 complex genome of *R. microplus*, an international collaboration has been formed to
421 commence a coordinated effort of targeted sequencing of both genome and
422 transcriptome with the longer-term view to obtain resources to sequence the complete
423 genome of *R. microplus*.

424 As described in the results section, nine datasets have been generated and
425 these contain the cumulative data our consortium has obtained which is relevant to *R.*
426 *microplus* genome and transcriptome sequences. In broad terms, these are the Cot-
427 filtration genomic DNA sequences, the BAC-associated sequences and the
428 transcriptome data. The Cot-filtration process filters out highly repetitive DNA and
429 allows the focus of resources on regions of the genome that are enriched for gene
430 coding sequences. The result from our mapping of the Cot-selected DNA to each of
431 the BACs is consistent with the Cot DNA experiments filtering out the highly
432 repetitive DNA fraction. Our analysis demonstrated that Cot-filtration is an effective
433 strategy to focus on gene-rich regions, because our Cot-filtered assembled contigs
434 found matches with 50% of the *BmiGI* Version 2.1 Gene Index despite the assembled

435 Cot sequences (Table 1 Dataset 2) only representing 2% of the *R. microplus* genome.
436 The availability of contigs from the Cot genomic DNA data set to align to
437 corresponding transcript sequences in *BmiGI* will enable the determination of intronic
438 and promoter regions associated with the exonic regions contained in *BmiGI*.

439 The acquisition of BES data from *R. microplus* represents the first project of
440 its type for this species. The 502 genes that were matched to the BES data represent a
441 16-fold increase over the expected number of matches based on the genome coverage
442 of the BES. The BES represents 0.1% of the genome, yet matches to 1.6% of the
443 *BmiGI* version 2.1 sequences. Further, the repetitive nature of the *R. microplus*
444 genome has been determined by reassociation kinetics analysis and was shown to
445 consist of 0.8% foldback, 30% unique DNA, 38% moderately repetitive DNA and
446 31% highly repetitive DNA (Ullmann et al., 2005). As the Cot DNA was selected to
447 remove the highly repetitive sequence fraction (32% of the genome), this is consistent
448 with 70% of BES alignment to the assembled Cot DNA contigs. This implies that the
449 BES is high in gene content and accounted for a high percentage alignment to the
450 gene-enriched Cot-selected DNA contigs. One possible reason for this is that the BAC
451 libraries were made with *MboI* partially digested genomic DNA. *MboI* has GATC as
452 the recognition site, and as GC-rich areas tend to be in gene-rich areas of genomes,
453 *MboI* digestion which forms the BES may lend bias to these gene-rich areas.

454 The 15 BACs provide an insight into the genomic structure and organisation
455 of *R. microplus*. The 15 BAC sequences are composed of 4.9% low complexity
456 sequence and known retrotransposable elements, for example GYPSY and LTR.
457 There is an average of 12 genes per BAC and this was a higher than expected ratio,
458 given the estimated gene space. It is possible that these BACs, 13 of which were
459 selected to be sequenced based on hybridization to known genes, occur in regions of

460 gene clusters. One particular BAC (BM-012-E08) had 6% of the BAC sequence
461 length mapped by Cot DNA and the remaining unmapped BAC sequence contained
462 highly repetitive intergenic regions nested between rDNA genes. In these selected 15
463 BACs a total of 97 genes out of 180 had significant alignments to known proteins.
464 These included 28 hypothetical proteins of which 20 were found in *Ixodes scapularis*,
465 five alignments to cytochrome P450 (*I. scapularis*), a receptor for egg jelly protein, a
466 papilin (*Pediculus humanus corporis*), a zinc finger, five unknown *Tribolium*
467 proteins, a serpin, an esterase and a transpanin. The remainder were comprised of
468 polyproteins, helicases and transcriptases (transposable element-like), (gene
469 predictions and annotation can be found in Supplementary Table S1 and
470 Supplementary Fig. S1).

471 The content of the newly assembled *R. microplus* transcriptome for the cattle
472 tick (Table 1 Dataset 8: *RmiTr* Version 1.0) contains 2,379 full-length sequences in
473 common with the previous *R. microplus* gene index and an additional 9,652 novel
474 full-length transcript sequences (Fig. 1). The *BmiGI* version 2.1 dataset consists of
475 9,851 contigs and 4,735 singletons. In our experience working with *BmiGI*,
476 approximately one-third to one-half of the contigs encode full-length proteins while a
477 very low percentage of the singletons encode full-length proteins. Thus, Dataset 8
478 should contain approximately 14,000 full-length transcripts. Further, the newly
479 assembled transcripts for the cattle tick contained 3,829 sequences with orthologues to
480 *I. scapularis* transcript sequences version 3.0 while there were 4,948 sequences in *I.*
481 *scapularis* Gene Index Version 3.0 that had orthologues to sequences in *RmiTr*
482 Version 1.0 (Table 3). As an example of the utility of CattleTickBase and *RmiTr*
483 Version 1.0, we searched both *RmiTr* Version 1.0 and *I. scapularis* Gene Index
484 Version 3.0 for sequences with similarity to esterases, cytochrome P450, and

485 glutathione S-transferases, metabolic proteins that often play roles in tick resistance to
486 acaricides. The results are shown in Table 3 and Supplementary Data S1-S3.

487 CattleTickBase is designed to become the comprehensive bioinformatics web
488 resource for *R. microplus*. CattleTickBase currently contains over 1.8 Gbp of genomic
489 sequence from *R. microplus*. Most of this sequence was derived from the Cot-selected
490 genomic DNA. This DNA was selected to yield the unique gene-enriched fraction of
491 the genomic DNA from *R. microplus* (Guerrero et al., 2010). Thus, if we assume the
492 30% unique fraction value from Ullmann et al. (2005), and also assume that our 1.8
493 Gbp of Cot-selected data primarily sources from the unique fraction (Table 1 Dataset
494 1), then with the *R. microplus* genome size of 7.1 Gbp, 30% of unique fraction would
495 represent 2.1 Gbp. Thus continuing the assumptions, CattleTickBase can be
496 considered to represent 0.86X coverage of the unique fraction (gene-enriched
497 fraction) of the *R. microplus* genome. However our Cot filtration dataset, although
498 effective in the recovery of cattle tick genomic DNA regions enriched in gene coding
499 sequences, has a drawback in that the low coverage has resulted in a somewhat
500 fragmented sequence assembly. This needs to be rectified by deeper sequencing of the
501 Cot-selected DNA or whole genome sequencing. Additionally, we expect genetic
502 variations amongst *R. microplus* populations from different countries or ecological
503 regions and this is necessarily an important consideration when defining consensus
504 reference data sets for *R. microplus*. Most of the data in CattleTickBase is derived
505 from North American *R. microplus*, although datasets 6 and 7 in Table 1 are from
506 Australian ticks. An example of a well-studied *R. microplus* transcript with sequence
507 that varies among different population is *Bm86*. Freeman et al. (2010) found 8.3%
508 sequence variation between *Bm86* isolated from *R. microplus* ticks from South Texas,

509 USA and from the Australian Yeerongpilly strain. Andreotti et al. (2008) reported
510 similar geographical differences in *Bm86* isolated from South American populations.

511 Nevertheless, we believe CattleTickBase is a significant resource for the tick
512 research community that should facilitate tick research in a number of areas. Not only
513 does CattleTickBase provide standard features of a genome sequencing project
514 browser and a BLAST interface; it also contains an online web-based analysis
515 environment (YABI) that enables the scientific community to conduct sophisticated
516 bioinformatics analysis. YABI's ease of use and facility for results to be easily
517 downloaded and used in other bioinformatics analysis environments makes it a
518 flexible web-based bioinformatics environment. To our knowledge, this type of
519 environment is not offered elsewhere. With the advent of third generation sequencing
520 technologies that promise longer reads than current 454 and Illumina technologies, we
521 expect to focus upon whole genome sequencing and targeted transcriptomic analysis
522 to further advance the *R. microplus* genome sequencing project. Naturally, progress
523 will be dependent on obtaining the necessary resources (both financial and scientific)
524 to further sequence these more difficult regions of the cattle tick genome. As
525 obtained, this information will be integrated into CattleTickBase to provide the
526 scientific community with access to the latest information from the *R. microplus*
527 genome-sequencing project.

528 In summary, we have acquired and assembled genomic DNA and
529 transcriptomic sequence data into contigs which represent over 1.8 Gbp of DNA from
530 gene-enriched regions of the *R. microplus* genome. The data was compiled into a
531 resource called CattleTickBase and a web-based YABI resource was developed to
532 enable the scientific community to access our databases. The YABI tool will facilitate

533 access and manipulation of cattle tick genome sequence data and, as the genome of *R.*

534 *microplus* proceeds to completion, the CattleTickBase resource will be updated.

535

536 **Acknowledgements**

537 The authors acknowledge funding support from the Beef CRC and the
538 National and International Research Alliances Program (DEEDI, Queensland
539 Government Smart Futures Funds, Australia) and U.S. Department of Agriculture
540 (USDA)-ARS Knipling Bushland US Livestock Insect Research Laboratory CRIS
541 Project Number 6205-32000-031-00 for the tick BAC and BES. We further
542 acknowledge technical assistance from Ms. Kylie G. Bendele, USDA-ARS Kerrville,
543 TX, USA and Dr. Zenaida Magbanua, Mississippi State University, USA and Dr.
544 Louise Jackson from Biosecurity, DEEDI, Australia for providing Australian ticks.
545 Thanks to Dr. Rudi Appels for helpful discussions during data analysis and
546 manuscript preparation. USDA is an equal opportunity employer.
547

548 **References**

- 549 Altschul, S.F., Gish, W., Miller, W., Myers, E.W., Lipman, D.J., 1990. Basic local
550 alignment search tool. *J. Mol. Biol.* 215, 403-410.
- 551 Andreotti, R., Perez de Leon, A.A., Dowd, S.E., Guerrero, F.D., Bendele, K.G.,
552 Scoles, G.A., 2011. Assessment of bacterial diversity in the cattle tick
553 *Rhipicephalus (Boophilus) microplus* through tag-encoded pyrosequencing.
554 *BMC Microbiol.* 11, 6.
- 555 Andreotti, R., Pedroso, M.S., Caetano, A.R., Martins, N.F., 2008. Comparison of
556 predicted binders in *Rhipicephalus (Boophilus) microplus* intestine protein
557 variants Bm86 Campo Grande strain, Bm86 and Bm95. *Rev. Brasileira*
558 *Parasitol. Vet.* 17, 93-98.
- 559 Barker, S.C., Murrell, A., 2004. Systematics and evolution of ticks with a list of valid
560 genus and species names. *Parasitol.* 129 Suppl., S15-S36.
- 561 Birol, I., Jackman, S.D., Nielsen, C.B., Qian, J.Q., Varhol, R., Stazyk, G., Morin,
562 R.D., Zhao, Y., Hirst, M., Schein, J.E., Horsman, D.E., Connors, J.M.,
563 Gascoyne, R.D., Marra, M.A., Jones, S.J., 2009. De novo transcriptome
564 assembly with ABySS. *Bioinformatics* 25, 2872-2877.
- 565 Burge, C., Karlin, S., 1997. Prediction of complete gene structures in human genomic
566 DNA. *J. Mol. Biol.* 268, 78-94.
- 567 de Castro, J.J., 1997. Sustainable tick and tickborne disease control in livestock
568 improvement in developing countries. *Vet. Parasitol.* 71, 77-97.
- 569 Drysdale, R., 2008. FlyBase : a database for the *Drosophila* research community.
570 *Methods Mol. Biol.* 420, 45-59.
- 571 Ewing, B., Green, P., 1998. Base-calling of automated sequencer traces using phred.
572 II. Error probabilities. *Genome Res.* 8, 186-194.

- 573 Filmore, D., 2004. It's a GPCR world. *Modern Drug Disc.* 7, 24-28.
- 574 Freeman, J.M., Davey, R.B., Kappmeyer, L.S., Kammlah, D.M., Olafson, P.U., 2010.
- 575 Bm86midgut protein sequence variation in South Texas cattle fever ticks.
- 576 *Parasites Vect.* 3, 101.
- 577 Grisi, L., Massard, C.L., Moya Borja, G.E., Pereira, J.B., 2002. Impacto, economico
- 578 das principais ectoparasitoses em bovinos no Brasil. *Hora Vet.* 21, 8-10.
- 579 Guerrero, F.D., Miller, R.J., Rousseau, M.E., Sunkara, S., Quackenbush, J., Lee, Y.,
- 580 Nene, V., 2005. BmiGI: a database of cDNAs expressed in *Boophilus*
- 581 *microplus*, the tropical/southern cattle tick. *Insect Biochem. Mol. Biol.* 35,
- 582 585-595.
- 583 Guerrero, F.D., Nene, V.M., George, J.E., Barker, S.C., Willadsen, P., 2006.
- 584 Sequencing a new target genome: the *Boophilus microplus* (Acari: Ixodidae)
- 585 genome project. *J. Med. Entomol.* 43, 9-16.
- 586 Guerrero, F.D., Moolhuijzen, P.M., Peterson, D.G., Bidwell, S., Caler, E., Appels, R.,
- 587 Bellgard, M., Nene, V.M., Djikeng, A., 2010. Reassociation kinetics-based
- 588 approach for partial genome sequencing of the cattle tick, *Rhipicephalus*
- 589 (*Boophilus*) *microplus*. *BMC Genomics* 11, 374.
- 590 Harris, T.W., Stein, L.D., 2006. WormBase: methods for data mining and
- 591 comparative genomics. *Methods Mol. Biol.* 351, 31-50.
- 592 Huang, X., Madan, A., 1999. CAP3: A DNA sequence assembly program. *Genome*
- 593 *Res.* 9, 868-877.
- 594 Kent, W.J., 2002. BLAT--the BLAST-like alignment tool. *Genome Res.* 12, 656-664.
- 595 Koski, L.B., Gray, M.W., Lang, B.F., Burger, G., 2005. AutoFACT: an automatic
- 596 functional annotation and classification tool. *BMC Bioinformatics* 6, 151.

- 597 Lagesen, K., Hallin, P., Rodland, E.A., Staerfeldt, H.H., Rognes, T., Ussery, D.W.,
598 2007. RNAmmer: consistent and rapid annotation of ribosomal RNA genes.
599 Nucleic Acids Res. 35, 3100-3108.
- 600 Lawson, D., Arensburger, P., Atkinson, P., Besansky, N.J., Bruggner, R.V., Butler,
601 R., Campbell, K.S., Christophides, G.K., Christley, S., Dialynas, E., Emmert,
602 D., Hammond, M., Hill, C.A., Kennedy, R.C., Lobo, N.F., MacCallum, M.R.,
603 Madey, G., Megy, K., Redmond, S., Russo, S., Severson, D.W., Stinson, E.O.,
604 Topalis, P., Zdobnov, E.M., Birney, E., Gelbart, W.M., Kafatos, F.C., Louis,
605 C., Collins, F.H., 2007. VectorBase: a home for invertebrate vectors of human
606 pathogens. Nucleic Acids Res. 35, D503-D505.
- 607 Lawson, D., Arensburger, P., Atkinson, P., Besansky, N.J., Bruggner, R.V., Butler,
608 R., Campbell, K.S., Christophides, G.K., Christley, S., Dialynas, E.,
609 Hammond, M., Hill, C.A., Konopinski, N., Lobo, N.F., MacCallum, R.M.,
610 Madey, G., Megy, K., Meyer, J., Redmond, S., Severson, D.W., Stinson, E.O.,
611 Topalis, P., Birney, E., Gelbart, W.M., Kafatos, F.C., Louis, C., Collins, F.H.,
612 2009. VectorBase: a data resource for invertebrate vector genomics. Nucleic
613 Acids Res. 37, D583-D587.
- 614 Lew-Tabor, A.E., Moolhuijzen, P.M., Vance, M.E., Kurscheid, S., Valle, M.R.,
615 Jarrett, S., Minchin, C.M., Jackson, L.A., Jonsson, N.N., Bellgard, M.I.,
616 Guerrero, F.D., 2009. Suppressive subtractive hybridization analysis of
617 *Rhipicephalus (Boophilus) microplus* larval and adult transcript expression
618 during attachment and feeding. Vet. Parasitol. 167, 304-320.
- 619 Li, H., Durbin, R., 2010. Fast and accurate long-read alignment with Burrows-
620 Wheeler transform. Bioinformatics 26, 589-595.

- 621 Mans, B.J., Neitz, A.W., 2004. Adaptation of ticks to a blood-feeding environment:
622 evolution from a functional perspective. *Insect Biochem. Mol. Biol.* 34, 1-17.
- 623 Margulies, M., Egholm, M., Altman, W.E., Attiya, S., Bader, J.S., Bembien, L.A.,
624 Berka, J., Braverman, M.S., Chen, Y.J., Chen, Z., Dewell, S.B., Du, L., Fierro,
625 J.M., Gomes, X.V., Godwin, B.C., He, W., Helgesen, S., Ho, C.H., Irzyk,
626 G.P., Jando, S.C., Alenquer, M.L., Jarvie, T.P., Jirage, K.B., Kim, J.B.,
627 Knight, J.R., Lanza, J.R., Leamon, J.H., Lefkowitz, S.M., Lei, M., Li, J.,
628 Lohman, K.L., Lu, H., Makhijani, V.B., McDade, K.E., McKenna, M.P.,
629 Myers, E.W., Nickerson, E., Nobile, J.R., Plant, R., Puc, B.P., Ronan, M.T.,
630 Roth, G.T., Sarkis, G.J., Simons, J.F., Simpson, J.W., Srinivasan, M., Tartaro,
631 K.R., Tomasz, A., Vogt, K.A., Volkmer, G.A., Wang, S.H., Wang, Y.,
632 Weiner, M.P., Yu, P., Begley, R.F., Rothberg, J.M., 2005. Genome sequencing
633 in microfabricated high-density picolitre reactors. *Nature* 437, 376-380.
- 634 Megy, K., Hammond, M., Lawson, D., Bruggner, R.V., Birney, E., Collins, F.H.,
635 2009. Genomic resources for invertebrate vectors of human pathogens, and the
636 role of VectorBase. *Infect. Genet. Evol.* 9, 308-313.
- 637 Moolhuijzen, P., Lew-Tabor, A., Morgan, A.T.J., Rodriguez Valle, M., Peterson,
638 D.G., Dowd, S.E., Guerrero, F.D., Bellgard, M.I., Appels, R., 2011. The
639 complexity of *Rhipicephalus (Boophilus) microplus* genome characterised
640 through detailed analysis of two BAC clones. *BMC Res. Notes* 4, 254.
- 641 Mortazavi, A., Williams, B.A., McCue, K., Schaeffer, L., Wold, B., 2008. Mapping
642 and quantifying mammalian transcriptomes by RNA-Seq. *Nat. Methods* 5,
643 621-628.

- 644 Playford, M., Rabiee, A.R., Lean, I.J., Ritchie, M., 2005. Review of research needs
645 for cattle tick control, Phases I and II. In, Meat & Livestock Australia Ltd,
646 Sydney.
- 647 Quackenbush, J., Cho, J., Lee, D., Liang, F., Holt, I., Karamycheva, S., Parvizi, B.,
648 Pertea, G., Sultana, R., White, J., 2001. The TIGR Gene Indices: analysis of
649 gene transcript sequences in highly sampled eukaryotic species. *Nucleic Acids*
650 *Res.* 29, 159-164.
- 651 Ribeiro, J.M., Alarcon-Chaidez, F., Francischetti, I.M., Mans, B.J., Mather, T.N.,
652 Valenzuela, J.G., Wikel, S.K., 2006. An annotated catalog of salivary gland
653 transcripts from *Ixodes scapularis* ticks. *Insect Biochem. Mol. Biol.* 36, 111-
654 129.
- 655 Rodriguez-Valle, M., Lew-Tabor, A., Gondro, C., Moolhuijzen, P., Vance, M.,
656 Guerrero, F.D., Bellgard, M., Jorgensen, W., 2010. Comparative microarray
657 analysis of *Rhipicephalus (Boophilus) microplus* expression profiles of larvae
658 pre-attachment and feeding adult female stages on *Bos indicus* and *Bos taurus*
659 cattle. *BMC Genomics* 11, 437.
- 660 Saldivar, L., Guerrero, F.D., Miller, R.J., Bendele, K.G., Gondro, C., Brayton, K.A.,
661 2008. Microarray analysis of acaricide-inducible gene expression in the
662 southern cattle tick, *Rhipicephalus (Boophilus) microplus*. *Insect Mol. Biol.*
663 17, 597-606.
- 664 Schattner, P., Brooks, A.N., Lowe, T.M., 2005. The tRNAscan-SE, snoscan and
665 snoGPS web servers for the detection of tRNAs and snoRNAs. *Nucleic Acids*
666 *Res.* 33, W686-689.
- 667 Snelson, J.T., 1975. Animal ectoparasites and disease vectors causing major
668 reductions in world food supplies. *FAO Plant Protect. Bull.* 23, 103-117.

- 669 Stein, L.D., Mungall, C., Shu, S., Caudy, M., Mangone, M., Day, A., Nickerson, E.,
670 Stajich, J.E., Harris, T.W., Arva, A., Lewis, S., 2002. The generic genome
671 browser: a building block for a model organism system database. *Genome*
672 *Res.* 12, 1599-1610.
- 673 Stewart, N.P., Callow, L.L., Duncalfe, F., 1982. Biological comparisons between a
674 laboratory-maintained and a recently isolated field strain of *Boophilus*
675 *microplus*. *Parasitology* 68, 691-694.
- 676 Sunter, J.D., Patel, S.P., Skilton, R.A., Githaka, N., Knowles, D.P., Scoles, G.A.,
677 Nene, V., de Villiers, E., Bishop, R.P., 2008. A novel SINE family occurs
678 frequently in both genomic DNA and transcribed sequences in ixodid ticks of
679 the arthropod sub-phylum Chelicerata. *Gene* 415, 13-22.
- 680 Turni, C., Lee, R.P., Jackson, L.A., 2004. A comparison of the immunosuppressive
681 effects of salivary gland extracts from two laboratory strains of *Boophilus*
682 *microplus*. *Int. J. Parasitol.* 34, 833-838.
- 683 Turni, C., Lee, R.P., Jackson, L.A., 2007. The effects of salivary gland extracts from
684 *Boophilus microplus* ticks on mitogen-stimulated bovine lymphocytes. *Vet.*
685 *Res. Commun.* 31, 545-552.
- 686 Ullmann, A.J., Lima, C.M., Guerrero, F.D., Piesman, J., Black IV, W.C., 2005.
687 Genome size and organization in the blacklegged tick, *Ixodes scapularis* and
688 the Southern cattle tick, *Boophilus microplus*. *Insect Mol. Biol.* 14, 217-222.
- 689 Wikel, S.K., 1999. Tick modulation of host immunity: an important factor in
690 pathogen transmission. *Int. J. Parasitol.* 29, 851-859.
- 691 Wistrand, M., Kall, L., Sonnhammer, E.L.L., 2006. A general model of G protein-
692 coupled receptor sequences and its application to detect remote homologs.
693 *Protein Sci.* 15, 509-521.
- 694

695 **Legends to Figures**

696 **Fig. 1.** Venn diagram summarizing the overlap between tick transcripts (S1 and S2)
697 and the *Rhipicephalus microplus* (*Rmi*) Gene Index (*BmiGI* Version 2.1). Transcript
698 sequences were clustered between Australian datasets 6 and 7 (Table 1) and the USA
699 *BmiGI* Version 2.1 gene index ([http://compbio.dfci.harvard.edu/cgi-](http://compbio.dfci.harvard.edu/cgi-bin/tgi/gimain.pl?gudb=b_microplus)
700 [bin/tgi/gimain.pl?gudb=b_microplus](http://compbio.dfci.harvard.edu/cgi-bin/tgi/gimain.pl?gudb=b_microplus)) at greater than 98 percent identity (PID). The
701 Venn diagram shows that the Australian datasets 6 and 7 contain 9,652 transcripts not
702 found in the *BmiGI*.

703 **Fig. 2.** Screen shot of the CattleTickBase home page. The CattleTickBase main web
704 page is a major resource for *Rhipicephalus microplus* data and analysis. On the
705 CattleTickBase home page pre-computed analyses can be viewed under the "View
706 Analyses" box; this includes genome browsers and full transcript annotations.
707 *Rhipicephalus microplus* sequence datasets are freely available for download from the
708 "Download Data" box. In the "Analysis Tools" box *R. microplus* datasets are
709 available for researchers to conduct their own analysis through a number of listed
710 bioinformatics applications and create customized bioinformatics pipelines. Users are
711 free to access the datasets and analysis tools available within the system, or can create
712 a secure account.

713 **Fig. 3.** YABI workflow analysis example for the identification of G protein-coupled
714 receptors (GPCRs) from transcript sequences in Dataset 7 of Table 1. The flow
715 diagram on the left hand side shows a simplified three-step bioinformatics workflow,
716 Step 1. Select a file of nucleotide sequences, Step 2. Translate all available open
717 reading frames (ORF) to amino acid sequences, Step 3. Search for GPCRs in the
718 translated ORF. On the right hand side of the figure is a screen capture of the YABI

719 online application; in the "Jobs" tab the automated three-step analysis pipeline
720 labelled "GPCR Rmi FL" is running. Steps 1 and 2 have completed successfully
721 (green circle with white mark) and step 3 is running, waiting for completion indicated
722 by the green bar. The remaining tabs include a "design" tab to create workflows and a
723 "files" tab where all files and results are stored.

ACCEPTED MANUSCRIPT

724 **Table 1.** Cattle Tick genome sequence and other relevant datasets available on CattleTickBase.

725	Dataset	Description	No. of sequences	No. of bp	References
726	1	Cot-selected genomic DNA seqs.	7,289,230	1,798,400,445	Guerrero et al., 2010
727	2	Cot-selected genomic DNA assembly	175,226	144,709,321	
728	3	BAC ^a end sequences	10,582	7,290,530	
729	4	Full-length sequenced BACs (15 BAC clones)	15	1,502,117	Guerrero et al., 2010, Moolhuijzen et al., 2011
730	5	BmiGI ^b contigs extended by Cot-selected seqs.	3,913	4,240,351	Guerrero et al., 2010
731	6	Tick gut transcriptome	11,333	9,228,737	
732	7	Frustrated larval transcriptome	6,082	3,617,080	
733	8	Assembled transcriptome <i>RmiTr</i> ^c Version 1.0	28,893	24,673,517	
734	9	NimbleGen microarray analyses GEO dataset GSE20605 (10 arrays)			Rodriguez-Valle et al., 2010

735 ^aBacterial Artificial Chromosome

736 ^b*Boophilus microplus* Gene Index at http://compbio.dfci.harvard.edu/cgi-bin/tgi/gimain.pl?gudb=b_microplus

737 ^c*Rhipicephalus microplus* transcriptome

738

ACCEPTED MANUSCRIPT

739 **Table 2.** Statistics of the 15 sequenced Bacterial Artificial Chromosomes (BACs) compared with Cot-selected genomic DNA.

BAC	BAC length (bp)	# Mapped Cot DNA reads ^a	BAC Cot % coverage ^b	No. of Genes/ No. of Exons ^c	GC content %
BM-005-B21	92,305	4,068	82	12/42	45
BM-005-G14	135,319	3,138	58	5/76	44
BM-013-M17	90,249	2,286	64	9/44	46
BM-026-P08	108,580	2,649	70	11/62	45
BM-031-L02	125,915	3,115	65	14/64	46
BM-118-H10	92,057	3,658	72	10/50	46
BM-004-A11	103,837	5,085	79	15/66	45
BM-006-B07	102,433	4,416	79	12/53	46
BM-010-J12	172,065	3,014	53	24/124	48
BM-012-E08	51,489	146	6	3/14	58
BM-074-Random-F12	95,687	1,534	61	15/40	40
BM-077-Random-J09	103,645	2,777	59	12/52	47
BM-129-CzEst9-N14	126,498	3,919	64	14/71	46
BM-066-M07	151,523	3,608	51	16/70	46
BM-077-G20	94,838	1,692	54	8/64	46

740

741 ^aThe number of Cot DNA reads mapped onto BAC sequence742 ^bThe percentage coverage of BAC sequence length by aligned Cot DNA reads743 ^cCount of predicted genes and total number of exons for each BAC sequence

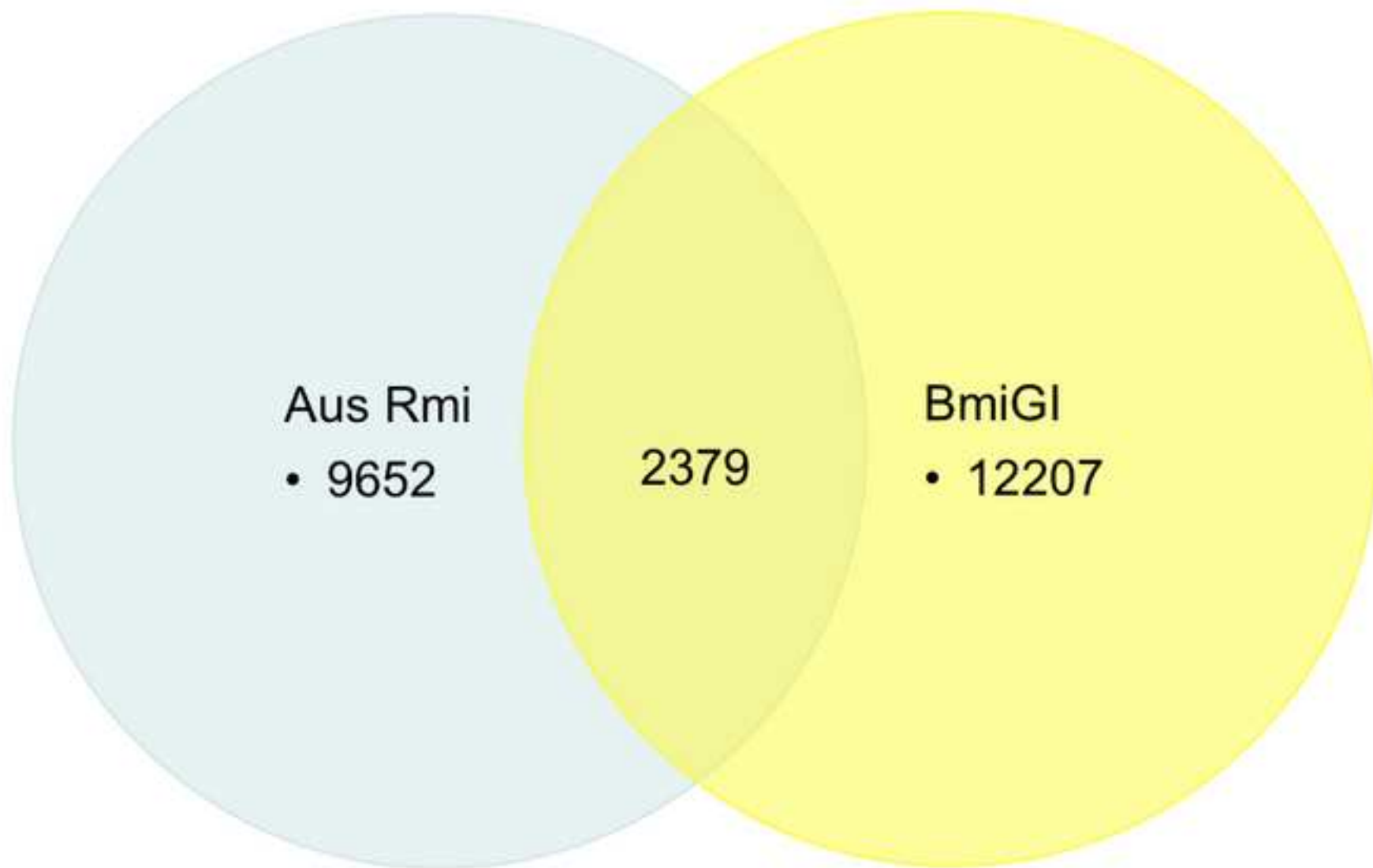
744

745 **Table 3.** *Ixodes scapularis* (*Isc*) transcript orthologous sequences in the *Rhipicephalus microplus* (*Rmi*)Tr Version 1.0.

746

747

Data	<i>Rmi</i> Tr1.0		<i>Isc</i> Gene Index Version 3.0	
	Total Seqs.	Orthologues with <i>Isc</i>	Total Seqs.	Orthologues with <i>Rmi</i> Tr
Overall Statistics	28,893	3,829	38,392	4,948
Cytochrome P450	115	55	357	235
Esterases	81	14	133	33
Glutathione S-transferases	39	28	96	60



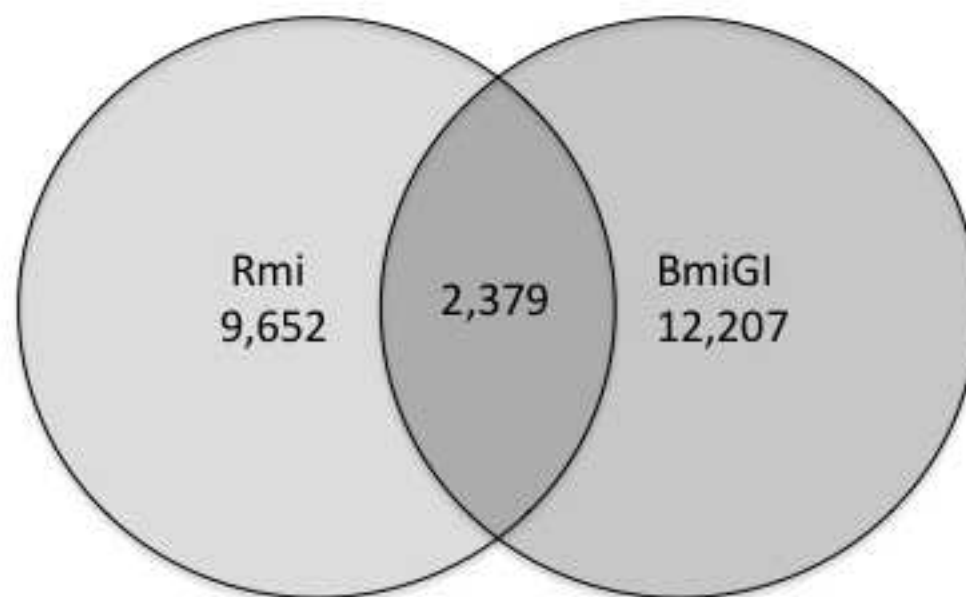


Figure 1

CCG Organism Portal

http://157.140.2.100/murdoch.edu.au/bioinformatics/rhipicephalus/

CCG Organism Portal

CENTRE FOR COMPARATIVE GENOMICS
Western Australia

Murdoch UNIVERSITY

Rhipicephalus microplus

View Analyses

- Browse Genome Annotations GBrowse
- Select a BAC sequence to view in GBrowse
- Select a Brugia Col extended transcripts to view in GBrowse
- View gene annotations for larval transcript library
- View gene annotations for tick gut transcript library
- Go back to organism portal

Analysis Tools

- Create your own bioinformatics pipelines to analyse data - YAS!
- Browse Genome sequence and add your own tracks - GBrowse
- Blat your data against our data sets

Download Data

- [Link to download cattle tick data](#)

Publications

- [Link to search for cattle tick publications](#)

Queensland Government **ARS** **Agricultural Research Service**

Introduction to CCG Rhipicephalus microplus Bioinformatics

A collaborative project between the State of Queensland (Department of Employment Economic Development and Innovation), Murdoch University's Centre for Comparative Genomics, and the United States Department of Agriculture Agricultural Research Service has generated a large amount of novel data and led to the creation of this online *R. microplus* genome resource. *Rhipicephalus microplus*, The cattle tick (formerly known as *Boophilus microplus*) is the most significant ectoparasite of cattle worldwide causing production losses and hide damage. Ticks also rapidly develop resistance to the chemicals (acaricides) used to treat cattle during heavy tick burden. Ticks are also vectors for highly pathogenic organisms including protozoan parasites *Babesia bovis* and *B. bigemina* causing bovine babesiosis and rickettsia *Anaplasma marginale* causing bovine anaplasmosis.

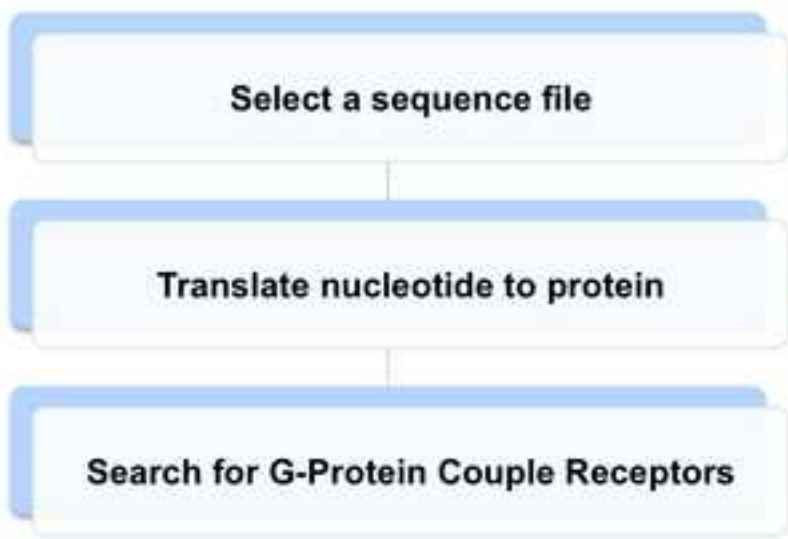
Research Aims

Integration of novel vector sequences and analyses to identify the genomic mechanisms to effectively control vector *R. microplus* livestock infestations.

Project site contains

This site contains access to data, analysis and to tools specific to the analysis of *Rhipicephalus microplus*. Data types include EST, BAC, and Cot DNA sequences, microarray data and miRNA studies direct link table can be found under "Download Data". Access to tools specific to the analysis of *Rhipicephalus microplus*, this includes interactive Bioinformatics pipeline analysis construction (YAS!). Get started using our online tutorial (link) and online help. Or go direct to View Analysis tools links.

© 2010 Centre for Comparative Genomics, Murdoch University



yabi

jobs design files

Find:

Date range: past 7 days

Status: **All** Ready Complete

GPCR Rmi FL	2011-08-05
-------------	------------

GPCR Rmi FL

Tags:

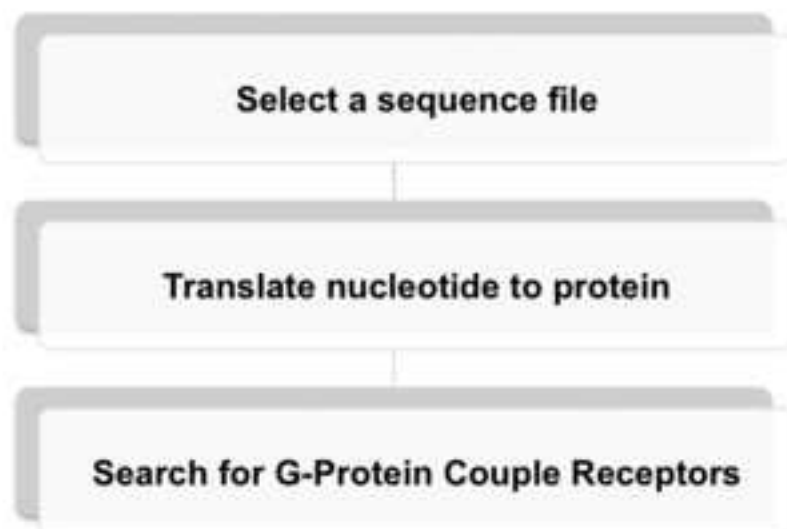
re-use

start

- 1 - select file
- 2 - getorf
- 3 - gpcrhmm

end

workflow running, waiting for completion...



The screenshot shows the yabi web interface. At the top left is the yabi logo with a stylized plant icon. To the right are navigation tabs for 'jobs', 'design', and 'files'. Below the logo is a search bar with 'Find:' and a date range selector set to 'past 7 days'. The status filter is set to 'All', with 'Ready' and 'Complete' options also visible. A table lists a job named 'GPCR Rmi FL' with a date of '2011-08-05'. To the right of the table is a detailed view of the job 'GPCR Rmi FL'. It includes a 'top' button, a 're-use' button, and a 'start' button. Below these are three steps in a workflow: '1 - select file', '2 - getorf', and '3 - gpcrhm', each with a refresh icon. At the bottom, there is an 'end' button and a status indicator: 'workflow running, waiting for completion...'.



>Consortium formed to begin sequencing cattle tick genome. >Acquired and assembled over 1.8 Gbp from gene-enriched regions of genome. >Web resource developed for access to data through web-based environment.

ACCEPTED MANUSCRIPT