

1 **BITACORA: A comprehensive tool for the identification and**  
2 **annotation of gene families in genome assemblies**

3

4 **Joel Vizqueta\*, Alejandro Sánchez-Gracia\* and Julio Rozas\***

5 Departament de Genètica, Microbiologia i Estadística and Institut de Recerca de la  
6 Biodiversitat (IRBio), Universitat de Barcelona, Barcelona, Spain

7 \*To whom correspondence should be addressed.

8

9 Corresponding authors: [jvizqueta@ub.edu](mailto:jvizqueta@ub.edu), [elsanchez@ub.edu](mailto:elsanchez@ub.edu) and [jrozas@ub.edu](mailto:jrozas@ub.edu)

10

11 Running head

12 BITACORA: A tool for gene family annotation

13

14 **Abstract**

15 Gene annotation is a critical bottleneck in genomic research, especially for the  
16 comprehensive study of very large gene families in the genomes of non-model  
17 organisms. Despite the recent progress in automatic methods, state-of-the-art tools used  
18 for this task often produce inaccurate annotations, such as fused, chimeric, partial or  
19 even completely absent gene models for many family copies, errors that require  
20 considerable extra efforts to be corrected. Here we present BITACORA, a  
21 bioinformatics solution that integrates popular sequence similarity-based search tools  
22 and Perl scripts to facilitate both the curation of these inaccurate annotations and the  
23 identification of previously undetected gene family copies directly in genomic DNA  
24 sequences. We tested the performance of BITACORA in annotating the members of  
25 two chemosensory gene families with different repertoire size in seven available  
26 genome sequences, and compared its performance with that of Augustus-PPX, a tool  
27 also designed to improve automatic annotations using a sequence similarity-based  
28 approach. Despite the relatively high fragmentation of some of these drafts,  
29 BITACORA was able to improve the annotation of many members of these families and  
30 detected thousands of new chemoreceptors encoded in genome sequences. The program  
31 creates general feature format (GFF) files, with both curated and newly identified gene  
32 models, and FASTA files with the predicted proteins. These outputs can be easily  
33 integrated in genomic annotation editors, greatly facilitating subsequent manual  
34 annotation and downstream evolutionary analyses.

35

## 36 **Introduction**

37 The falling cost of high-throughput sequencing (HTS) technologies made them  
38 accessible to small labs, promoting a large number of genome-sequencing projects even  
39 in non-model organisms. Nevertheless, genome assembly and annotation, especially in  
40 eukaryotic genomes, still represent major limitations (Dominguez Del Angel et al.,  
41 2018). The unique genomic characteristics of many non-model organisms, often lacking  
42 pre-existing gene models (Yandell & Ence, 2012), and the absence of closely related  
43 species with well-annotated genomes, means that the annotation process can be very  
44 challenging. State-of-the-art pipelines for *de novo* genome annotation, like BRAKER1  
45 (Hoff, Lange, Lomsadze, Borodovsky, & Stanke, 2016) or MAKER2 (Holt & Yandell,  
46 2011), allow integrating multiple evidences such as RNA-seq, EST data, gene models  
47 from other previously annotated species or *ab initio* gene predictions (using software  
48 such as GeneMark, (Lomsadze, Burns, & Borodovsky, 2014), Exonerate (Slater &  
49 Birney, 2005), GenomeThreader (Gremme, Brendel, Sparks, & Kurtz, 2005), Augustus  
50 (M. Stanke & Waack, 2003; Mario Stanke, Diekhans, Baertsch, & Haussler, 2008) or  
51 SNAP (Korf, 2004)). However, the gene models predicted by these automatic tools are  
52 often inaccurate, particularly for gene family members. Furthermore, these predictions  
53 can be especially inaccurate for medium or low-quality assemblies, which is a quite  
54 common situation in the increasing large number of genome drafts of non-model  
55 organisms used in molecular ecology studies. The correct annotation of gene families  
56 frequently requires additional programs, such as Augustus-PPX (Keller, Kollmar,  
57 Stanke, & Waack, 2011a), or semi-automatic, and even manual approaches, that  
58 evaluate the quality of supporting data. This latter task is usually performed in genomic  
59 annotation editors, such as Apollo, which give researchers the option to work  
60 simultaneously in the same annotation project (Lee et al., 2013).

61 There are a number of issues affecting the quality of gene family annotations, especially  
62 for either old or fast evolving families (Yohe et al., 2019). First, new duplicates within a  
63 family usually originate by unequal crossing-over and are found in tandem arrays in the  
64 genome, with the more recent duplicates also the physically closest (Clifton et al., 2020;  
65 Vieira, Sánchez-Gracia, & Rozas, 2007). This configuration often causes local  
66 misassemblies that result in the incorrect or failed identification of tandem duplicated  
67 copies (i.e., it produces artifact, incomplete, or chimeric genes along a genomic region).  
68 Secondly, the identification and characterization of gene copies in medium- to large-  
69 sized families tends to be laborious, requiring data from multiple sources, including  
70 well-annotated remote homologs and hidden Markov model (HMM) profiles. Certainly,  
71 the robust identification and annotation of the complete repertory of a gene family in a  
72 typical genome draft is a challenging task that requires important additional efforts,  
73 which are very tedious to perform manually.

74 In order to facilitate this curation task, we have developed BITACORA, a  
75 bioinformatics pipeline to assist the comprehensive annotation of gene families in  
76 genome assemblies. BITACORA requires a structurally annotated genome (GFF and  
77 FASTA format) or a draft assembly, and a curated database with well-annotated  
78 members of the focal gene families. The program will perform comprehensive BLAST  
79 and HMMER searches (Altschul, 1997; Eddy, 2011) to identify putative candidate gene  
80 regions (already annotated, or not), combine evidences from all searches and generate  
81 new gene models. The outcome of the pipeline consists of a new structural annotation  
82 (GFF) file along with their encoded sequences. These output sequences can be directly  
83 used to conduct downstream functional or evolutionary analyses or to facilitate a fine-  
84 scale re-annotation in genome browsers such as Apollo (Lee et al., 2013).

## 86 **Methods and implementation**

### 87 *Input data files*

88 BITACORA requires: i) a data file with the genome sequences (in FASTA format); ii)  
89 the associated GFF file with annotated features (either in GFF3 or GTF formats;  
90 features must include both transcript or mRNA, and CDS); iii) a data file with the  
91 predicted proteins included in the GFF (in FASTA format); and iv) a database (here  
92 referred as FPDB database) with the protein sequences of well annotated members of  
93 the gene family of interest (focal family; in FASTA format) along with its HMM profile  
94 (see Supplementary Material for a detailed description of FPDB construction). Since  
95 sequence similarity-based searches are very sensitive to the quality of the proteins in  
96 FPDB, it is important to include in this database highly curated proteins from closely  
97 related species. This is especially important for the annotation of very old or fast-  
98 evolving gene families. Also, the use of a HMM profile increases the likelihood of  
99 identifying sequences encoding new members; these profiles can be obtained from  
100 external databases (such as PFAM) or built using high quality protein alignments with  
101 the program *hmmbuild* (Finn *et al.*, 2014). Before starting the analysis, BITACORA  
102 checks whether input data files are correctly formatted; otherwise, it will suggest some  
103 format converters distributed with the program (see Troubleshooting section in  
104 Supplementary Material).

### 105 *Curating existing annotations*

106 The BITACORA workflow has three main steps (Fig. 1). The first step consists of the  
107 identification of all putative homologs of the FPDB sequences from the focal gene  
108 family that are already present in the input GFF file, and the curation of their gene  
109 models (referred hereinafter as b-curated (bitacora-curated) gene models or proteins).

110 Specifically, the pipeline launches BLASTP and HMMER searches (Altschul, 1997;  
111 Eddy, 2011) against the proteins predicted from the features in the input GFF using the  
112 FPDB protein sequences and HMM profiles as queries; the resulting alignments are  
113 filtered for quality (i.e. BLASTP hits covering at least two-thirds of the length of query  
114 sequences or including at least 80% of the complete protein used as a subject are  
115 retained). The results from both searches are combined into a single integrated result for  
116 every single protein (gene model). Then, BITACORA trims the original models based  
117 on these combined results (retaining only the aligned sequence) and reports new gene  
118 coordinates (b-curated models) in a new updated GFF (uGFF), fixing for example all  
119 chimeric annotations. In addition, the proteins encoded by these b-curated models are  
120 incorporated into the FPDB (updated FPDB or uFPDB), to be used in an additional  
121 search round.

#### 122 *Identifying new genomic regions encoding gene family members*

123 In the second step, BITACORA uses TBLASTN to search the genome sequences for  
124 regions encoding homologs of the proteins included in the uFPDB but not annotated in  
125 the uGFF. BITACORA implements two different approaches for generating novel gene  
126 models from TBLASTN results (set with the “gemoma” parameter). In the first  
127 approach, BITACORA executes the GeMoMa tool, a homology-based gene prediction  
128 program that uses amino acid sequence and intron position conservation to reconstruct  
129 gene models from BLAST hits (Keilwagen, Hartung, & Grau, 2019; Keilwagen,  
130 Hartung, Paulini, Twardziok, & Grau, 2018; Keilwagen et al., 2016). The second  
131 approach is based on a “close proximity” strategy. Under this strategy, all independent  
132 TBLASTN hits (i.e., after merging all alignments that overlap in TBLASTN results)  
133 located in the same scaffold and separated by less than a predetermined distance (set  
134 with the “intron distance” parameter), are connected to form a unique gene model. This

135 step intends to join all coding exons of the same gene based on the average intron length  
136 in the focal genome. We provide some scripts to estimate this average length from the  
137 input GFF (see Supplementary Material).

138 Finally, to avoid reporting inaccurate gene models due to artifactual gene fusions in  
139 dense gene clusters or any other possible errors (regardless of which algorithm of the  
140 abovementioned has been applied), BITACORA will check for the presence of the gene  
141 family-specific protein domain (using the HMM profile in FPDB), and only reports in  
142 the curated dataset those gene models containing the domain. In addition, all proteins  
143 are tagged with a label that indicates the number of different domains in the sequence  
144 (Ndom). This final filtering step can be relaxed using the BITACORA "genomicblastp"  
145 option, which evaluates the presence of positive hits in either HMMER, or BLASTP  
146 searches against the proteins in FPDB (see Supplementary Material for details).

#### 147 *Optional search round and final output*

148 Finally, BITACORA can also be used to perform a second search round using as the  
149 input data all proteins obtained in steps 1 and 2 (sFPDB database). This additional step  
150 (step 3 in Fig 1) is especially useful for searching remote homologs undetected in the  
151 first round. The final BITACORA outcome will include 1) an updated GFF file with  
152 both b-curated and b-novel gene models. 2) All non-redundant proteins predicted from  
153 these feature annotations (in a FASTA file). 3) Two BED files, one with the coordinates  
154 of all independent TBLASTN hits found in the genome sequence, and the other with  
155 only those hits that would encode novel putative exons and, 4) all protein sequences  
156 found in all steps.

#### 157 *Additional features*

158 BITACORA could be also used in the absence of either a reference genome for the  
159 target species (e.g. for transcriptomic studies; Protein mode) or a precompiled GFF (e.g.  
160 for non-annotated genomes; Genome mode); in these cases, the input should be a  
161 FASTA file with the set of predicted proteins or the genome sequences, respectively  
162 (see Supplementary Material for alternative usage modes). With BITACORA, we also  
163 distribute a series of scripts to perform some useful tasks, such as estimating intron  
164 length statistics from a GFF, converting GFF to GTF format, and retrieving all protein  
165 sequences encoded by the features of a GFF file. Furthermore, to better adjust to the  
166 particularities of each genome, BITACORA allows the user to specify the values of the  
167 most important parameters, such as the *E*-value for BLAST and HMMER searches, the  
168 number of threads in BLAST runs, and the algorithm to build novel gene models from  
169 TBLASN hits.

170

### 171 **BITACORA application example**

172 To demonstrate the performance of BITACORA in annotating gene family members in  
173 a group of genomes of different assembly quality, we present an extended report of the  
174 results in Vizuela et al., (2018). Specifically, we selected two of the arthropod  
175 chemosensory gene families, insect gustatory receptors (GR) and Niemann-Pick type  
176 C2 (NPC2) proteins (Pelosi, Iovinella, Felicioli, & Dani, 2014; Robertson, 2015) in a  
177 subset of seven of the eleven chelicerate genomes surveyed in this study (Table 1; Fig.  
178 2). We selected these gene families since they widely differ in the number of members  
179 and protein length. Whereas the GR is a large gene family that encode seven-  
180 transmembrane receptors of about 400 amino acids long, the NPC2 have few members  
181 and encode shorter proteins (an average of about 150 amino acids); despite the different  
182 length, both gene families have a similar average number of exons per gene in the



183 surveyed species. Furthermore, to validate the accuracy of our software in gold standard  
184 annotated genomes, we checked the performance of BITACORA in the annotation of  
185 GR and NPC2 members in the genome of *Drosophila melanogaster* (Adams et al.,  
186 2000) and of the C2H2 zinc finger domain (PF00096) in human and mouse genomes  
187 (Lander et al., 2001; Waterston et al., 2002). The last corresponds to a very short  
188 domain (about 20-30 amino acids) that is present in multiple adjacent copies (usually 2-  
189 3, but up to 16) in C2H2 zinc finger proteins, an important family of higher eukaryotic  
190 transcription factors that represent about 3% of the human genes (Cassandri et al.,  
191 2017).

192 For the analysis, we retrieved genome sequences, annotations and predicted peptides of  
193 *D. melanogaster* (r6.31, FlyBase; Adams et al., 2000); the scorpions *Centruroides*  
194 *sculpturatus* (bark scorpion, genome assembly version v1.0, annotation version v0.5.3;  
195 Human Genome Sequencing Center (HGSC)) and *Mesobuthus martensii* (v1.0,  
196 Scientific Data Sharing Platform Bioinformatics (SDSPB)) (Cao et al., 2013); the  
197 spiders *Acanthoscurria geniculata* (tarantula, v1, NCBI Assembly, BGI) (Sanggaard et  
198 al., 2014), *Stegodyphus mimosarum* (African social velvet spider, v1, NCBI Assembly,  
199 BGI) (Sanggaard et al., 2014), *Latrodectus hesperus* (western black widow, v1.0,  
200 HGSC), *Parasteatoda tepidariorum* (common house spider, v1.0 Augustus 3,  
201 SpiderWeb and HGSC) (Schwager et al., 2017) and *Loxosceles reclusa* (brown recluse,  
202 v1.0, HGSC); human (GRCh38; Lander et al., 2001) and mouse (GRCm38; Waterston  
203 et al., 2002).

204 In addition, and with a benchmarking purpose, we compared the performance of  
205 BITACORA with Augustus PPX, a method that also uses protein profiles to improve  
206 automatic annotations of gene family members (--proteinprofile; Keller et al., 2011;  
207 Mario Stanke, Schöffmann, Morgenstern, & Waack, 2006), in annotating GR and NPC2

208 copies in the same seven chelicerate genomes. Strikingly, BITACORA uncovered the  
209 identification of thousands of new gene models previously undetected in chelicerates,  
210 even after applying Augustus-PPX (Table 1; see also supplementary data in Vizueta et  
211 al. 2018 to find the BITACORA curated sequences). For instance, in the bark scorpion  
212 *Centruroides sculpturatus*, the automatic annotation pipelines show 24 GR encoding  
213 sequences, while BITACORA was able to identify and annotate 1,234 genes or gene  
214 fragments (1,210 in addition to the 24 previously annotated genes), for the only 307  
215 recovered with Augustus-PPX (Table 1; Supplementary table S1). Globally,  
216 BITACORA identified, annotated and curated 3,570 sequences encoding GR proteins  
217 across the seven chelicerate genomes (3,466 of which were absent in the available GFF  
218 for this species), while Augustus-PPX only predicted 1,638 gene models for this family  
219 (Table 1; Supplementary table S1). It is largely known that this gene family evolves  
220 rapidly in arthropods, both in terms of sequence change and repertoire size, encoding in  
221 the same genome very recent and distantly related receptors as well as pseudogenes.  
222 Since some of these receptors show a very restricted gene expression pattern (expressed  
223 in specialized cells and tissues involved in chemoreception), their transcripts are often  
224 missing in RNA-seq data sets, which are one of the evidences used for the automatic  
225 annotation of genomes (Joseph & Carlson, 2015; Robertson, 2015; Vizueta et al., 2017;  
226 Zhang, Zheng, Li, & Fan, 2014). This fact, together with the huge divergence that  
227 exhibit many copies (old duplication events and/or rapid evolution), are probably the  
228 causes of the low accuracy of both automatic annotation and Augustus-PPX.

229 The members of the NPC2 family, on the contrary, are much more conserved at the  
230 sequence level and show higher levels of gene expression in arthropods (Pelosi et al.,  
231 2014). As expected, the number of newly identified copies is much lower than in the  
232 case of GRs. Even then, BITACORA was able to detect 44 novel NPC2 encoding

233 sequences, raising the total annotated repertoire in these species from 75 to 119 (Table  
234 1). In this case, Augustus-PPX was able to recover 97 gene models for this gene family,  
235 which improves the performance of previous automatic annotations, but still is  
236 outperformed by BITACORA. Importantly, Augustus-PPX predicted thousands of gene  
237 models that are not real members of the focal gene family (Supplementary table S1),  
238 requiring further actions to separate gene family copies from false allocations. Finally,  
239 both methods correctly annotated all members of the GR and NPC2 families in the *D.*  
240 *melanogaster* genome. It is worth noting, however, that a non-negligible number of  
241 these novel identified genes in chelicerate genomes are incomplete (about 40% and 63%  
242 of the GR and NPC2 members, respectively). This feature can be partially explained by  
243 the poor genome assembly quality (indicated by the N50 and number of scaffolds), or  
244 by the low number of annotated proteins in the input GFF. Although BITACORA can  
245 be useful under such low-quality data, it will compromise its performance in terms of  
246 complete gene models.

247 We identified 4,510 and 3,068 annotated proteins containing C2H2 zinc finger domains  
248 in the human and mouse genomes, respectively (Supplementary table S2). These  
249 proteins correspond to 709 human and 708 mouse genes, of which 645 in human and  
250 278 in mice are curated proteins of the Uniprot Swiss-Prot database. In addition to the  
251 pseudogenes annotated in these species, BITACORA detected 44 and 133 putative  
252 novel genes encoding C2H2 domain sequences in the human and mouse genomes,  
253 respectively (i.e. absent in the last version of the GFFs for these species); these genes  
254 could be false positives caused by the very short length and repetitive structure of the  
255 query domain or deprecated models resulting from curated genome annotations. In any  
256 case, BITACORA was able to correctly identify all human and mouse genes reported to  
257 contain the C2H2 domain in NCBI, in addition to 90 mouse members of the C2H2 zinc

258 finger family initially annotated as just zinc finger proteins. We compared the results of  
259 BITACORA with those of Augustus *ab initio* (through BRAKER1 pipeline) and  
260 Augustus-PPX using the optimized parameters for vertebrates. Like BITACORA, these  
261 annotation tools identified all C2H2 zinc finger genes reported in NCBI (in both human  
262 and mice) plus some additional gene models, which, as in our pipeline, could represent  
263 false positives or deprecated models (Supplementary table S2). As expected,  
264 BRAKER1, but specially Augustus-PPX, gene models are more fragmented than those  
265 found in BITACORA since these pipelines perform a completely *de novo* prediction,  
266 resulting in a higher number of shorter genes. Altogether, these results clearly  
267 demonstrate the utility of BITACORA in low quality genome drafts for which  
268 annotation pipelines are not optimized. First, BITACORA demonstrates a similar  
269 performance than these pipelines in high quality annotated genomes as different as those  
270 of insects and vertebrates, and for families with very different characteristics and  
271 repertory sizes. Second, our software is able to fetch information from the automatic  
272 annotations generated by these pipelines to validate and curate existing gene family  
273 members and detect new family copies in poor-quality genome drafts.

## 274 **Discussion**

275 Gene families are one of the most abundant and dynamic components of eukaryotic  
276 genomes. Therefore, having curated genomic data is fundamental not only to carry out  
277 comprehensive comparative or functional genomics studies on gene families, but also to  
278 understand global genome architecture and biology. During the last decades, the rapid  
279 development of sequencing technologies has enabled the large accumulation of genome  
280 sequences of non-model organisms. These projects, which often address very specific  
281 molecular ecology studies or are in the context of large comparative genomics analyses,  
282 typically rely on automatic annotation pipelines and very little efforts are devoted to

283 curate these annotations. The proteins predicted by automatic annotation tools often  
284 contain systematic errors, such as incomplete or chimeric gene models, which are  
285 especially notable in gene families given the repetitive nature of their members.  
286 Besides, since new copies commonly arise by unequal crossing-over, they are  
287 frequently found in physically close tandem arrays of similar sequences, further  
288 complicating annotations (Clifton et al., 2020; Vieira et al., 2007).

289 With this in mind, we have developed a bioinformatics tool that helps researchers to  
290 access these automatic annotations, extract the information of focal gene families,  
291 curate and update gene models and identify new copies from DNA sequences. Using  
292 BITACORA, gene family annotations can be substantially improved using both HMM  
293 profiles and iterative searches that incorporate the new variability found in previous  
294 searches. Indeed, we validated our tool by comparing its performance with a method  
295 developed to improve the annotation of gene family members matching a protein  
296 profile, Augustus-PPX (Keller et al., 2011b; Mario Stanke et al., 2006). BITACORA  
297 not only outperforms the annotations of Augustus-PPX in the examples shown here, but  
298 it also demonstrated to be more accurate in its predictions.

299 The estimation of gene gains and losses, and the associated birth and death rates  
300 analyses, are very sensitive to the quality of genome annotations. The example of the  
301 GR family in chelicerates demonstrates the importance of refining annotations using  
302 BITACORA. Indeed, using unsupervised annotations in low quality genome drafts of  
303 non-model organisms directly to estimate turnover rates might produce very erroneous  
304 results, not only in terms of gene counts but also in calculations biased to highly  
305 expressed and/or very recent copies. BITACORA can be used to considerably reduce  
306 these errors and make more accurate and robust inferences about the age/origin of the  
307 family and of its mode of evolution.

308 On the other hand, the curation of both existing and new identified members of a family  
309 with BITACORA might be also crucial for further analysis on their sequence evolution.  
310 The quality of multiple sequence alignments, which are used to determine orthology  
311 groups, to obtain divergence estimates or to detect the footprint of natural selection in  
312 gene family members, is strongly compromised by the presence of badly annotated  
313 copies, including chimeras and incorrectly annotated fragments. Using BITACORA we  
314 can detect these artifacts and either fix or discard them from further analyses.

315 Despite its proven utility, we are aware that BITACORA does not provide perfect  
316 annotations for a gene family. The use of GeMoMa algorithm is more sensitive than the  
317 close-proximity method generating more accurate gene models, although, in the  
318 presence of assembly errors or highly fragmented genomes, this approach might fail to  
319 identify genes, and especially putative pseudogenes. In these cases, the close-proximity  
320 method could help to detect these cases and report them in final output. Consequently,  
321 the combination of different genome annotation tools, such as general automatic  
322 pipelines (e.g. BRAKER and MAKER2), with software specifically designed to  
323 annotate gene families, such as Augustus-PPX or BITACORA, would be highly  
324 recommended for most of the poor-quality genome drafts of non-model organisms. In  
325 this sense, the advantage of BITACORA is that it is able to process and curate already  
326 existing gene models in addition to identifying totally novel family members in genome  
327 sequences.

328 Furthermore, to overcome putative gene model errors, BITACORA implements some  
329 filtering steps to determine if the predicted coding sequences are correct. The program  
330 carries out a HMMER search to identify the protein family domain in all new annotated  
331 sequences. In addition, if the HMMER search is negative, BITACORA can relax this  
332 step by checking if the novel genes show significant BLASTP hits in a search against

333 FPDB proteins. In this case, the sensitivity of the annotations will increase at the  
334 expense of specificity (i.e. it could generate false allocations to the focal family in the  
335 presence of repetitive regions or FPDB contaminations, for instance). It is important to  
336 note that BITACORA generates homology-based predictions that could require  
337 different levels of experimental validation depending on the nature of further  
338 downstream analyses.

339 Notwithstanding such filtering steps, BITACORA offers an output directly readable in  
340 genome editor tools, such as Apollo, which facilitate researchers to improve gene  
341 models. Fig. 3 shows an example of the annotation tracks generated by BITACORA  
342 (GFF3 and BED files) for a cluster of three members of the NPC2 family in the genome  
343 of the spider *P. tepidariorum*. The automatic annotation of this region using MAKER2  
344 (track Ptep\_v0.5.3-Models), generated a chimeric gene model (two different genes are  
345 fused) which could be easily curated using BITACORA. Additionally, despite  
346 TBLASTN searches having detected a putative novel exon in the gene encoding  
347 NPC2\_5, GeMoMa did not include this sequence in the final gene model due to the  
348 presence of an in-frame stop codon. In order to decide if this stop codon is an  
349 annotation, assembly or sequencing artifact, it would be necessary, for instance, to  
350 verify if the exon exists in other species, if that region is transcribed, or if the gene is  
351 under selective constraints.

352

### 353 **Conclusion**

354 Genome annotation, especially in medium to low quality drafts of non-model  
355 organisms, is still a drawback for the increasingly large number of evolutionary and  
356 functional genomic analyses in the context of molecular ecology studies. To assist this

357 task, we developed a comprehensive pipeline that facilitates the identification and  
358 curation of existing models and the annotation of new gene family copies in novel  
359 genome assemblies. The improved annotations generated with our pipeline can be used  
360 either directly to perform downstream analyses or as a baseline for further manual  
361 curation in genome annotation editors. Future directions should focus on including  
362 novel sources of evidence, such as RNA-seq data, in BITACORA searches or  
363 integrating the pipeline as a part of genome annotation editors, which will greatly  
364 facilitate the annotation of large gene families in collaborative genome projects.

365

### 366 **Acknowledgements**

367 We would like to thank Paula Escuer and Vadim Pisarenco for helpful discussions, and  
368 three anonymous reviewers whose comments have greatly improved this manuscript.  
369 This work was supported by the Ministerio de Economía y Competitividad of Spain  
370 (CGL2013-45211, CGL2016-75255) and the Comissió Interdepartamental de Recerca I  
371 Innovació Tecnològica of Catalonia, Spain (2017SGR1287). J.V. was supported by a  
372 FPI grant (Ministerio de Economía y Competitividad of Spain, BES-2014-068437).

373

### 374 **Author contributions**

375 J.V., A.S.-G and J.R. conceived the work. J.V. wrote the scripts, did the analyses and  
376 wrote the first version of the manuscript. All authors checked and confirmed the final  
377 version of the manuscript.

378

### 379 **Data accessibility**



380 BITACORA is available from <http://www.ub.edu/softevol/bitacora>, and  
381 <https://github.com/molevol-ub/bitacora>

382

### 383 **References**

384 Adams, M. D., Celniker, S. E., Holt, R. A., Evans, C. A., Gocayne, J. D., Amanatides,  
385 P. G., ... Venter, J. C. (2000). The genome sequence of *Drosophila melanogaster*.  
386 *Science*, 287(5461), 2185–95. Retrieved from  
387 <http://www.ncbi.nlm.nih.gov/pubmed/10731132>

388 Altschul, S. (1997). Gapped BLAST and PSI-BLAST: a new generation of protein  
389 database search programs. *Nucleic Acids Research*, 25(17), 3389–3402.  
390 doi:10.1093/nar/25.17.3389

391 Cassandri, M., Smirnov, A., Novelli, F., Pitolli, C., Agostini, M., Malewicz, M., ...  
392 Raschellà, G. (2017). Zinc-finger proteins in health and disease. *Cell Death*  
393 *Discovery*. Springer Nature. doi:10.1038/cddiscovery.2017.71

394 Clifton, B. D., Jimenez, J., Kimura, A., Chahine, Z., Librado, P., Sánchez-Gracia, ...  
395 Ranz, J. M. (2020). Understanding the early evolutionary stages of a tandem D.  
396 melanogaster-specific gene family: a structural and functional population study.  
397 *Molecular Biology and Evolution*, XX, (in press). doi: 10.1093/molbev/msaa109

398 Dominguez Del Angel, V., Hjerde, E., Sterck, L., Capella-Gutierrez, S., Notredame, C.,  
399 Vinnere Pettersson, O., ... Lantz, H. (2018). Ten steps to get started in Genome  
400 Assembly and Annotation. *F1000Research*, 7, ELIXIR-148.  
401 doi:10.12688/f1000research.13598.1

402 Eddy, S. R. (2011). Accelerated Profile HMM Searches. *PLoS Computational Biology*,

403 7(10), e1002195. doi:10.1371/journal.pcbi.1002195

404 Finn, R. D., Bateman, A., Clements, J., Coghill, P., Eberhardt, R. Y., Eddy, S. R., ...  
405 Punta, M. (2014). Pfam: the protein families database. *Nucleic Acids Research*,  
406 42(Database issue), D222–D230. doi:10.1093/nar/gkt1223

407 Gremme, G., Brendel, V., Sparks, M. E., & Kurtz, S. (2005). Engineering a software  
408 tool for gene structure prediction in higher organisms. *Information and Software  
409 Technology*, 47(15), 965–978. doi:10.1016/J.INFSOF.2005.09.005

410 Hoff, K. J., Lange, S., Lomsadze, A., Borodovsky, M., & Stanke, M. (2016).  
411 BRAKER1: Unsupervised RNA-Seq-Based Genome Annotation with GeneMark-  
412 ET and AUGUSTUS. *Bioinformatics*, 32(5), 767–769.  
413 doi:10.1093/bioinformatics/btv661

414 Holt, C., & Yandell, M. (2011). MAKER2: an annotation pipeline and genome-database  
415 management tool for second-generation genome projects. *BMC Bioinformatics*,  
416 12(1), 491. doi:10.1186/1471-2105-12-491

417 Joseph, R. M., & Carlson, J. R. (2015). *Drosophila* Chemoreceptors: A Molecular  
418 Interface Between the Chemical World and the Brain. *Trends in Genetics : TIG*,  
419 31(12), 683–695. doi:10.1016/j.tig.2015.09.005

420 Keilwagen, J., Hartung, F., & Grau, J. (2019). GeMoMa: Homology-based gene  
421 prediction utilizing intron position conservation and RNA-seq data. In *Methods in  
422 Molecular Biology* (Vol. 1962, pp. 161–177). Humana Press Inc. doi:10.1007/978-  
423 1-4939-9173-0\_9

424 Keilwagen, J., Hartung, F., Paulini, M., Twardziok, S. O., & Grau, J. (2018).  
425 Combining RNA-seq data and homology-based gene prediction for plants, animals

426 and fungi. *BMC Bioinformatics*, 19(1), 189. doi:10.1186/s12859-018-2203-5

427 Keilwagen, J., Wenk, M., Erickson, J. L., Schattat, M. H., Grau, J., & Hartung, F.  
428 (2016). Using intron position conservation for homology-based gene prediction.  
429 *Nucleic Acids Research*, 44(9), 89. doi:10.1093/nar/gkw092

430 Keller, O., Kollmar, M., Stanke, M., & Waack, S. (2011a). A novel hybrid gene  
431 prediction method employing protein multiple sequence alignments.  
432 *Bioinformatics*, 27(6), 757–763. doi:10.1093/bioinformatics/btr010

433 Keller, O., Kollmar, M., Stanke, M., & Waack, S. (2011b). A novel hybrid gene  
434 prediction method employing protein multiple sequence alignments.  
435 *Bioinformatics*, 27(6), 757–763. doi:10.1093/bioinformatics/btr010

436 Korf, I. (2004). Gene finding in novel genomes. *BMC Bioinformatics*, 5, 59.  
437 doi:10.1186/1471-2105-5-59

438 Lander, E. S., Linton, L. M., Birren, B., Nusbaum, C., Zody, M. C., Baldwin, J., ...  
439 Morgan, M. J. (2001). Initial sequencing and analysis of the human genome.  
440 *Nature*, 409(6822), 860–921. doi:10.1038/35057062

441 Lee, E., Helt, G. A., Reese, J. T., Munoz-Torres, M. C., Childers, C. P., Buels, R. M.,  
442 ... Lewis, S. E. (2013). Web Apollo: a web-based genomic annotation editing  
443 platform. *Genome Biology*, 14(8), R93. doi:10.1186/gb-2013-14-8-r93

444 Lomsadze, A., Burns, P. D., & Borodovsky, M. (2014). Integration of mapped RNA-  
445 Seq reads into automatic training of eukaryotic gene finding algorithm. *Nucleic*  
446 *Acids Research*, 42(15), e119–e119. doi:10.1093/nar/gku557

447 Pelosi, P., Iovinella, I., Felicioli, A., & Dani, F. R. (2014). Soluble proteins of chemical  
448 communication: an overview across arthropods. *Frontiers in Physiology*,

449 5(August), 320. doi:10.3389/fphys.2014.00320

450 Robertson, H. M. (2015). The Insect Chemoreceptor Superfamily Is Ancient in  
451 Animals. *Chemical Senses*, 40(9), 609–614. doi:10.1093/chemse/bjv046

452 Slater, G. S. C., & Birney, E. (2005). Automated generation of heuristics for biological  
453 sequence comparison. *BMC Bioinformatics*, 6, 31. doi:10.1186/1471-2105-6-31

454 Stanke, M., & Waack, S. (2003). Gene prediction with a hidden Markov model and a  
455 new intron submodel. *Bioinformatics*, 19(Suppl 2), ii215–ii225.  
456 doi:10.1093/bioinformatics/btg1080

457 Stanke, Mario, Diekhans, M., Baertsch, R., & Haussler, D. (2008). Using native and  
458 syntenically mapped cDNA alignments to improve de novo gene finding.  
459 *Bioinformatics*, 24(5), 637–644. doi:10.1093/bioinformatics/btn013

460 Stanke, Mario, Schöffmann, O., Morgenstern, B., & Waack, S. (2006). Gene prediction  
461 in eukaryotes with a generalized hidden Markov model that uses hints from  
462 external sources. *BMC Bioinformatics*, 7(1), 62. doi:10.1186/1471-2105-7-62

463 Vieira, F. G., Sánchez-Gracia, A., & Rozas, J. (2007). Comparative genomic analysis of  
464 the odorant-binding protein family in 12 *Drosophila* genomes: purifying selection  
465 and birth-and-death evolution. *Genome Biology*, 8(11), R235. doi:10.1186/gb-  
466 2007-8-11-r235

467 Vizuela, J., Frías-López, C., Macías-Hernández, N., Arnedo, M. A., Sánchez-Gracia, A.,  
468 & Rozas, J. (2017). Evolution of chemosensory gene families in arthropods:  
469 Insight from the first inclusive comparative transcriptome analysis across spider  
470 appendages. *Genome Biology and Evolution*, 9(1), 178–196.  
471 doi:10.1093/gbe/evw296

472 Vizuela, J., Rozas, J., & Sánchez-Gracia, A. (2018). Comparative Genomics Reveals  
473 Thousands of Novel Chemosensory Genes and Massive Changes in  
474 Chemoreceptor Repertoires across Chelicerates. *Genome Biology and Evolution*,  
475 *10*(5), 1221–1236. doi:10.1093/gbe/evy081

476 Waterston, R. H., Lindblad-Toh, K., Birney, E., Rogers, J., Abril, J. F., Agarwal, P., ...  
477 Lander, E. S. (2002). Initial sequencing and comparative analysis of the mouse  
478 genome. *Nature*, *420*(6915), 520–562. doi:10.1038/nature01262

479 Yandell, M., & Ence, D. (2012). A beginner’s guide to eukaryotic genome annotation.  
480 *Nature Reviews Genetics*, *13*(5), 329–342. doi:10.1038/nrg3174

481 Yohe, L. R., Davies, K. T. J., Simmons, N. B., Sears, K. E., Dumont, E. R., Rossiter, S.  
482 J., & Dávalos, L. M. (2019). Evaluating the performance of targeted sequence  
483 capture, RNA-Seq, and degenerate-primer PCR cloning for sequencing the largest  
484 mammalian multigene family. *Molecular Ecology Resources*. doi:10.1111/1755-  
485 0998.13093

486 Zhang, Y., Zheng, Y., Li, D., & Fan, Y. (2014). Transcriptomics and identification of  
487 the chemoreceptor superfamily of the pupal parasitoid of the oriental fruit fly,  
488 *Spalangia endius* Walker (Hymenoptera: Pteromalidae). *PloS One*, *9*(2), e87800.  
489 doi:10.1371/journal.pone.0087800

490

491

492 **Tables**

493 **Table 1.** Summary of the number of GRs and NPC2 genes identified by BITACORA  
494 and Augustus-PPX in genome assemblies.

495

496 **Figures**

497 **Fig. 1.** Schematic representation of the BITACORA workflow.

498

499 **Fig. 2.** Phylogenetic relationships among the seven chelicerate species surveyed for the  
500 GR and the NPC2 families.

501

502 **Fig. 3.** Example of the visualization in the Apollo genome editor of the BITACORA  
503 output. The example includes the annotation features of three genes encoding NPC2  
504 proteins that are arranged in tandem in the spider *P. tepidariorum*. Current automatic  
505 annotation of this genomic region obtained with MAKER2 (track PTEP\_v0.5.3-  
506 Models), produced a chimeric gene model (PtepTmpM024154-RA; an artifactual two  
507 gene fusion), which is effectively curated by BITACORA (NPC2\_5 and NPC2\_6 gene  
508 models). The next three tracks are generated by BITACORA. The  
509 GFF3\_NPC2\_BITACORA track, which includes the final gene models, both curated or  
510 newly identified by the program, and the BED\_NPC2\_All and BED\_NPC2\_Novel  
511 tracks showing the position of all independent TBLASTN hits found in sequence  
512 similarity-based searches, or only those involving novel putative exons, respectively.  
513 Note that a novel coding sequence (not predicted in automatic annotations) is predicted  
514 by the program.

515

516 **Supplementary Material**

517

518 **Table S1.** Summary of the genome information and the number of GRs and NPC2  
519 genes identified by BITACORA and Augustus-PPX in the genome assemblies of the  
520 seven surveyed chelicerates, and in *D. melanogaster*.

521

522 **Table S2.** Summary of the number of C2H2 zinc finger genes identified by  
523 BITACORA, BRAKER1 and Augustus-PPX in human and mouse genome assemblies.

524

525

526 **Supplementary documentation**

527 BITACORA Documentation

528