

## **Jekyll or Hyde? The genome (and more) of *Nesidiocoris tenuis*, a zoophytophagous predatory bug that is both a biological control agent and a pest**

K. B. Ferguson<sup>a</sup>, S. Visser<sup>b,c</sup>, M. Dalíková<sup>b,c</sup>, I. Provazníková<sup>b,c,†</sup>, A. Urbaneja<sup>d</sup>, M. Pérez-Hedo<sup>d</sup>, F. Marec<sup>b</sup>, J. H. Werren<sup>e</sup>, B. J. Zwaan<sup>a</sup>, B. A. Pannebakker<sup>a</sup>, and E. C. Verhulst<sup>f</sup>

<sup>a</sup> Wageningen University, Laboratory of Genetics, Wageningen, The Netherlands

<sup>b</sup> Biology Centre CAS, Institute of Entomology, České Budějovice, Czech Republic

<sup>c</sup> University of South Bohemia, Faculty of Science, České Budějovice, Czech Republic

<sup>d</sup> Instituto Valenciano de Investigaciones Agrarias (IVIA), Centro de Protección Vegetal y Biotecnología, Moncada, Valencia, Spain

<sup>e</sup> University of Rochester, Department of Biology, Rochester, NY, USA

<sup>f</sup> Wageningen University, Laboratory of Entomology, Wageningen, The Netherlands

<sup>†</sup> Current address: European Molecular Biology Laboratory, Heidelberg, Germany

### **Abstract**

*Nesidiocoris tenuis* (Reuter) is an efficient predatory biological control agent used throughout the Mediterranean Basin in tomato crops but regarded as a pest in northern European countries. Belonging to the family Miridae, it is an economically important insect yet very little is known in terms of genetic information – no published genome, population studies, or RNA transcripts. It is a relatively small and long-lived diploid insect, characteristics that complicate genome sequencing. Here, we circumvent these issues by using a linked-read sequencing strategy on a single female *N. tenuis*. From this, we assembled the 355 Mbp genome and delivered an *ab initio*, homology-based, and evidence-based annotation. Along the way, the bacterial “contamination” was removed from the assembly, which also revealed potential symbionts. Additionally, bacterial lateral gene transfer (LGT) candidates were detected in the *N. tenuis* genome. The complete gene set is composed of 24,688 genes; the associated proteins were compared to other hemipterans (*Cimex lectularis*, *Halyomorpha halys*, and *Acyrtosiphon pisum*), resulting in an initial assessment of unique and shared protein clusters. We visualised the genome using various cytogenetic techniques, such as karyotyping, CGH and GISH, indicating a

18 karyotype of  $2n=32$  with a male-heterogametic XX/XY system. Additional analyses  
19 include the localization of 18S rDNA and unique satellite probes via FISH techniques.  
20 Finally, population genomics via pooled sequencing further showed the utility of this  
21 genome. This is one of the first mirid genomes to be released and the first of a mirid  
22 biological control agent, representing a step forward in integrating genome  
23 sequencing strategies with biological control research.

## 24 **Introduction**

25 Hemiptera is the fifth largest insect order and the most speciose hemimetabolous order  
26 with over 82,000 described species (Panfilio and Angelini, 2018). While recent  
27 sequencing projects have presented a variety of information about hemipteran  
28 genomes, large families such as the plant bugs Miridae still lack genomic resources,  
29 with the exception of transcriptomic resources for some members (Tian et al., 2015),  
30 and the more recent genome of *Apolygus lucorum*, a mirid pest that has a publicly  
31 available genome as of December 2019 (NCBI BioProject PRJNA526332). With the  
32 exception of *A. lucorum*, the lack of genomic resources for Miridae is in spite of the  
33 diverse life histories present, as it contains not only some of the most notorious  
34 agricultural pests but also predators that are often used in biological control (van  
35 Lenteren et al., 2018). In addition, Hemiptera are known for their intriguing karyotype  
36 evolution involving holocentric (holokinetic) chromosomes but there is a lack of  
37 cytogenetic information on Miridae. The absence of the ancestral TTAGG<sub>n</sub> telomeric  
38 repeat have been reported for mirids *Macrolophus* spp., *Deraeocoris* spp., and  
39 *Megaloceroea relicticornis* (Geoffroy) (Grozeva et al., 2019, 2011; Jauset et al., 2015)  
40 but more knowledge of this trait is necessary for evolutionary studies of genomes and  
41 karyotypes. Furthermore, the taxonomic issues that lie within both Miridae and  
42 Hemiptera could better be resolved using protein and transcriptome-based analysis,

43 but there is a noted lack of data in this regard as well (Panfilio and Angelini, 2018). While  
44 there is a relatively large amount of research into mirids and their use in biological  
45 control compared to other predators (Puentes et al., 2018), sequencing projects, if  
46 any, often focus on pest species and not on biological control agents (Panfilio and  
47 Angelini, 2018). For more advanced molecular methods such as RNAi and CRISPR-  
48 based genome editing strategies, it is necessary to have access to genomic and  
49 transcriptomic resources of the target species, and so these methods are currently out  
50 of reach for *N. tenuis* researchers. This lack in resources on both agricultural pests and  
51 biological control agents in the Miridae prompted us to generate genomic and  
52 cytogenetic resources of a mirid species that is both.

53 *Nesidiocoris tenuis* (Reuter) (Hemiptera: Miridae) is a zoophytophagous mirid used as  
54 a biological control agent worldwide, including in Spain, the Mediterranean Basin, and  
55 China (Pérez-Hedo and Urbaneja, 2016; Xun et al., 2016). Throughout the  
56 Mediterranean Basin, *N. tenuis* is used in tomato greenhouses and open fields against  
57 whiteflies (Hemiptera: Aleyrodidae), and the South American tomato pinworm, *Tuta*  
58 *absoluta* (Meyrick) (Lepidoptera: Gelechiidae) (Calvo et al., 2009; Mollá et al., 2014).  
59 In addition, due to its high degree of polyphagous behaviour, it is able to prey on other  
60 pest species such as thrips, leaf miners, leafhoppers, aphids, spider mites, and  
61 lepidopteran pests (Pérez-Hedo and Urbaneja, 2016). While *N. tenuis* is an important  
62 biological control agent in Mediterranean countries (Urbaneja et al., 2012), it is often  
63 cited as a pest in other contexts and countries (Calvo et al., 2009; Pérez-Hedo and  
64 Urbaneja, 2016). When prey is scarce in tomatoes, due to its phytophagy, *N. tenuis*  
65 can cause plant lesions such as brown discolouration around tender stems, known as  
66 necrotic rings, in addition to leaf wilt, and flower abortion (Arnó et al., 2010). This switch  
67 to phytophagy has been observed to be inversely proportional to the availability of  
68 prey (Sanchez, 2009). Therefore, much of the research thus far has focused on

69 characterizing *N. tenuis* biology and ecology, classifying the induced damage, and  
70 attempting to reduce it (Biondi et al., 2015; Castañé et al., 2011; Garantonakis et al.,  
71 2018; Martínez-García et al., 2016; Urbaneja-Bernat et al., 2019). Despite its associated  
72 plant damage, *N. tenuis* is widely used across South-eastern Spain as it is an efficient  
73 agent against the various pests it controls (Arnó et al., 2010). Furthermore, the  
74 aforementioned phytophagy has been demonstrated to have benefits by triggering  
75 predator-induced defences, including attracting parasitoids, repulsing other  
76 herbivorous pests, and restricting accumulation of viruses (Bouagga et al., 2019; Pérez-  
77 Hedó et al., 2018, 2015).

78 In recent years, the controversial success of *N. tenuis* has encouraged the scientific  
79 community to study this predatory mirid (Puentes et al., 2018). However, some issues  
80 remain to be addressed, such as the genetic variation in commercial stocks of similar  
81 biological control agents when compared to wild populations, with the former often  
82 diminished in comparison to the latter as seen in other biological control agents  
83 (Paspati et al., 2019; Rasmussen et al., 2018; Streito et al., 2017). In order to compare  
84 biological control stock to wild (or wild-caught) populations, determining the current  
85 diversity and genetic variation of the commercial stock is important. Finally, *N. tenuis*  
86 is known to host bacterial symbionts, including *Wolbachia* and *Rickettsia*, though the  
87 effect of these bacteria on their host is relatively unknown (Caspi-Fluger et al., 2014).  
88 Sequence data can provide additional insight into potential symbionts as well as  
89 identify potential LGTs (lateral gene transfers) between host and symbiont.

90 With all of these fascinating avenues of research in mind, it may be surprising to learn  
91 that, aside from a mitogenome (Dai et al., 2012), a regional population analysis (Xun  
92 et al., 2016), and more recent work shedding light on evidence of LGT (P. Xu et al.,  
93 2019), little genomic information exists for *N. tenuis* and there is no published *N. tenuis*

94 genome. Characteristics such as karyotype, sex chromosome system, and presence  
95 or absence of telomeric repeats are currently unknown. A likely reason for this absence  
96 of genomic resources is that advances made in sequencing technology are often  
97 juxtaposed to the complexities of insect life cycles and difficulties in obtaining enough  
98 high quality genomic material due to size and exoskeleton (Leung et al., 2019; Richards  
99 and Murali, 2015). Additionally, current assembly tools have a hard time dealing with  
100 heterozygosity; therefore, a genome assembly is benefited by sequencing material of  
101 reduced genetic heterozygosity for a more contiguous assembly. Reduced  
102 heterozygosity is often difficult to achieve in diploid insects where the genetic variation  
103 within a population is unknown or the species cannot be inbred (Keeling et al., 2013).

104 Generating the genomes of highly heterozygous, diploid, and relatively small insects is  
105 tricky; researchers have to be prepared to balance their expectations and the  
106 available technology (Ellegren, 2014; Leung et al., 2019). While a single diploid  
107 individual may yield enough material for an Illumina-only library, assembly may be  
108 difficult due to large repeat regions that extend beyond the insert size of the library.  
109 Conversely, enough material could be obtained for sequencing on a long-read  
110 platform, but may require pooling material from multiple individuals, potentially  
111 complicating assembly due to the heterozygosity of the population. While possible  
112 solutions include estimating the heterozygosity or setting up inbred populations (which  
113 can be nearly impossible if deleterious effects of inbreeding need to be avoided or if  
114 the presence of a complementary sex determining system limits inbreeding (Szűcs et  
115 al., 2019; van Wilgenburg et al., 2006)), an alternative is to create a linked-read library.  
116 The 10x Genomics platform creates a microfluidic partitioned library that individually  
117 barcodes minute amounts of long strands of DNA for further amplification (10x  
118 Genomics Inc., Pleasanton, CA, USA). This library is then sequenced on a short-read  
119 sequencing platform and then assembled using the barcodes to link reads together

120 into the larger fragment (i.e. Chin et al. 2016; Jones et al. 2017). This method allows for  
121 a library to be constructed from a single individual that contains additional structural  
122 information to aid assembly (such as phasing), removing the need for pooling multiple  
123 individuals and avoiding assembly difficulties in repetitive regions. Additional  
124 information, such as karyotype, can further improve genomes in the assembly stage  
125 as well as inform further directions of research by providing chromosome-level context,  
126 encouraging further improvement of a genome beyond its initial release.

127 Here we present the genome of *Nesidiocoris tenuis* achieved by sequencing a linked-  
128 read library of a single adult female bug, along with an annotation based on  
129 transcriptome, homology-based, and *ab initio* predictions. In addition to the genome,  
130 various avenues for future research are initiated to raise the profile of *N. tenuis* as a  
131 research organism, including cytogenetic analyses, protein cluster analysis, and a  
132 genome-wide pooled sequencing population genetics analysis. These resources  
133 benefit biological control research, as more knowledge becomes available to use in  
134 research as well as knowledge of the species for taxonomic and phylogenetic  
135 purposes.

## 136 **Methods**

### 137 **Species origin and description**

138 Individuals of *N. tenuis* were received either from the commercial biological control  
139 stock at Koppert Biological Systems, S. L. (Águilas, Murcia, Spain) (KBS) or from the  
140 population maintained for less than a year at Wageningen University and Research  
141 (WUR) Greenhouse Horticulture (Bleiswijk, The Netherlands), which in turn were  
142 originally sourced from the KBS commercial population. Material used for DNA  
143 sequencing, PCR testing, pooled sequencing, and cytogenetics was from the KBS

144 population, while material used for RNA sequencing was from the WUR Greenhouse  
145 Horticulture population. Additional species used for cytogenetic comparison purposes  
146 were sourced from two separate laboratory populations within the Biology Centre CAS  
147 in České Budějovice, Czech Republic: *Triatoma infestans* (Klug) (Hemiptera:  
148 Reduviidae) individuals were obtained from a laboratory colony at the Institute of  
149 Parasitology that was originally sourced from Bolivia (Schwarz et al., 2014), while  
150 *Ephestia kuehniella* (Zeller) (Lepidoptera: Pyralidae) individuals were obtained from a  
151 wild-type laboratory colony at the Institute of Entomology (Marec and Shvedov, 1990).  
152 Species identification of the KBS population was confirmed via COI sequencing using  
153 a PCR amplification protocol (Itou et al., 2013), in addition to testing for the presence  
154 of *Wolbachia* via PCR amplification protocol (Zhou et al., 1998).

#### 155 **Flow cytometry**

156 Genome size was estimated with flow cytometry on propidium-iodide stained nuclei.  
157 Individuals from a mixed *Drosophila melanogaster* (Meigen) (Diptera: Drosophilidae)  
158 laboratory population (May et al., 2019) were used as the standard for genome size  
159 comparison. Following established preparation protocols (De Boer et al., 2007), three  
160 samples of single *D. melanogaster* heads, two samples of single *N. tenuis* heads, and  
161 one sample of a single *N. tenuis* head pooled with a single *D. melanogaster* head  
162 were analysed in a FACS flow cytometer (BD FACSAria™ III Fusion Cell Sorter, BD  
163 Biosciences, San Jose, USA). With the known genome size of *D. melanogaster* of 175  
164 Mbp, we could calculate an approximate genome size relative to the amount of  
165 fluorescence (Hare and Johnston, 2011).

#### 166 **gDNA Extraction**

167 A single female *N. tenuis* was placed in a 1.5 mL safelock tube with 5-8 one mm glass  
168 beads and frozen in liquid nitrogen and shaken for 30 s in a Silamat S6 shaker (Ivoclar  
169 Vivadent, Schaan, Liechtenstein). DNA was then extracted using the Qiagen  
170 MagAttract Kit (Qiagen, Hilden, Germany). Following an overnight lysis step with Buffer  
171 ATL and proteinase K at 56°C, extraction was performed according to MagAttract Kit  
172 protocol. Elution was performed in two steps with 50 µL of Buffer AE (Tris-EDTA) each  
173 time, yielding 424 ng of genomic DNA (gDNA) in 100 µL as measured with an Invitrogen  
174 Qubit 2.0 fluorometer using the dsDNA HS Assay Kit (Thermo Fisher Scientific, Waltham,  
175 USA).

### 176 **Library Preparation and Sequencing**

177 Following extraction, gDNA was further diluted to 1 ng/µl following the Chromium  
178 Genome Reagent Kits Version 1 User Guide (version CG-00022) (10x Genomics,  
179 Pleasanton, USA). A library of Genome Gel Beads was combined with 1 ng of gDNA,  
180 Master Mix, and partitioning oil to create Gel Bead-In-EMulsions (GEMs). The GEMs  
181 underwent an isothermal amplification step and barcoded DNA fragments were  
182 recovered for Illumina library construction (Illumina, San Diego, USA). The library was  
183 then sequenced on an Illumina HiSeq 2500 at the Bioscience Omics Facility at  
184 Wageningen University and Research (Wageningen, The Netherlands), yielding  
185 212,910,509 paired-end reads with a read length of 150 bp. The first 23 bp of each  
186 forward read is a 10X GEM barcode used in the assembly process. Forward read  
187 quality was similar to that of the reverse reads, and no reads were flagged for poor  
188 quality in a FastQC assessment (Andrews et al., 2015).

### 189 **Assembly**



190 Using the reads, a k-mer count analysis was performed using GenomeScope on k-mer  
191 sizes of 21 and 48, which was used to infer heterozygosity (Vurture et al., 2017).  
192 Assembly was performed using all available reads with the GEM barcodes  
193 incorporated during the Chromium library preparation in Supernova v2.1.1 (10X  
194 Genomics, Pleasanton, USA), with default settings (Weisenfeld et al., 2017). This  
195 assembly, v1.0, underwent a preliminary decontamination using NCBI BLASTn v2.2.31+  
196 against the NCBI nucleotide collection (nt) focusing on scaffolds with over 95%  
197 homology to bacteria (Camacho et al., 2009), followed by the more elaborate  
198 method described below (Detecting contamination and LGT events). Finally, 100%  
199 duplicate scaffolds were identified using the *dedupe* tool within BBTools  
200 ([sourceforge.net/projects/bbmap/](https://sourceforge.net/projects/bbmap/)), and removed alongside the contaminated  
201 scaffolds, resulting in assembly v1.5. Attempts at further deduplication by adjusting the  
202 threshold (such as 95% duplication) resulted in further deletions, but at larger scaffold  
203 size, percentage is a rather blunt tool and any percentage is an arbitrary cut-off, so  
204 we decided to only remove true duplicates. Assembly completeness for both  
205 assemblies were determined using BUSCO v3.0.2 and the insect\_odb9 ortholog set  
206 (Simão et al., 2015), while assembly statistics were determined using QUAST (Gurevich  
207 et al., 2013).

## 208 **Detecting contamination and LGT events**

209 Lateral gene transfers (LGTs) from bacteria into metazoan genomes were once  
210 thought to be rare or non-existent, but are now known to be relatively common and  
211 can evolve into functional genes (Dunning Hotopp et al., 2007; Husnik and  
212 McCutcheon, 2018). We therefore screened our insect genome for LGTs from  
213 bacteria. As insect genome assemblies often contain scaffolds from associated

214 bacteria, we first screened for such “contaminating” scaffolds and moved them into  
215 a separate metagenomic multiFASTA (S1.3.4).

216 We used a DNA-based computational pipeline to both identify likely contaminating  
217 bacterial scaffolds in the assembly, and to detect potential LGT from bacteria into the  
218 insect genome. The LGT Pipeline was modified from an earlier version developed by  
219 David Wheeler and John Werren (Wheeler et al., 2013), and has been used to screen  
220 for bacterial “contamination” and LGTs in a number of arthropod genomes before  
221 (e.g. bedbug *Cimex lectularis* L. (Hemiptera: Cimicidae) (Benoit et al., 2016), parasitoid  
222 wasp *Trichogramma pretiosum* Riley (Hymenoptera: Trichogrammatidae) (Lindsey et  
223 al., 2018), and the milkweed bug *Oncopeltus fasciatus* (Dallas) (Hemiptera:  
224 Lygaeidae) (Panfilio et al., 2019)). In some cases, entire or nearly complete bacterial  
225 genomes have been retrieved from arthropod genome projects (e.g. (Benoit et al.,  
226 2016; Lindsey et al., 2016)).

#### 227 *Detection of Bacterial Scaffolds in the Assembly*

228 To detect bacterial contaminating scaffolds, the following method was used after the  
229 preliminary bacterial contamination assessment described above. First, each scaffold  
230 was broken into 1 kbp fragments and each fragment was subsequently searched with  
231 BLASTn against an in-house reference database that contains 2,100 different bacterial  
232 species (complete list in S1.3.1) which was masked for low complexity regions using the  
233 NCBI Dustmasker function (Morgulis et al., 2006). We recorded each bacterial match  
234 with bitscore > 50, the number of bacterial matches, total bacterial coverage in the  
235 scaffold, proportion of the scaffold covered, total hit width of coverage (the distance  
236 between the leftmost and rightmost bacteria hit proportional to the scaffold size) and  
237 the bacterial species with the greatest number of matches within the scaffold from the  
238 in-house bacterial data base. It should be noted that the latter method does not

239 indicate the actual bacterial species from which the scaffold was derived, as it is  
240 based on similarity to a curated database – that determination would require follow-  
241 up analysis, which was not performed in this study.

242 Any criterion for deciding whether a scaffold comes from a bacterium is unavoidably  
243 arbitrary: Too stringent and insect scaffolds are included; too lax and insect scaffolds  
244 are inappropriately removed. We applied a cut-off of  $\geq 0.40$  proportion bacterial hit  
245 width, which has performed well to remove contamination in a few test cases where  
246 we have manually examined scaffolds near the cut-off. All instances of contaminated  
247 scaffolds were removed from the assembly and are available in supplementary  
248 materials as a list (S1.3.2) and a multi-FASTA file (S1.3.3).

#### 249 *Identifying LGT Candidate Regions*

250 We used the same DNA based computational pipeline to identify potential LGTs from  
251 bacteria into the insect genome. The basic method is as follows: as before, scaffolds  
252 from the genome assembly are broken into 1 kbp intervals, which are searched  
253 against a bacterial genome database. Any positive bacterial hit in a 1 kbp region  
254 (bitscore > 50) was then searched against a database containing transcripts from the  
255 following eukaryotes: *Xenopus*, *Daphnia*, *Strongylocentrotus*, *Mus*, *Homo sapiens*,  
256 *Aplysia*, *Caenorhabditis*, *Hydra*, *Monosiga*, and *Acanthamoeba*  
257 ([ftp://ftp.hgsc.bcm.edu/I5K-](ftp://ftp.hgsc.bcm.edu/I5K-pilot/LGT_analysis/All_species_genomes/lgt_finder_blastn_database_directories/)  
258 [pilot/LGT\\_analysis/All\\_species\\_genomes/lgt\\_finder\\_blastn\\_database\\_directories/](ftp://ftp.hgsc.bcm.edu/I5K-pilot/LGT_analysis/All_species_genomes/lgt_finder_blastn_database_directories/)).

259 The purpose of this eukaryotic screening is to identify highly conserved genes that are  
260 shared between eukaryotes and bacteria and exclude these from further analysis. To  
261 focus our attentions on the most likely LGT candidates, we selected hits with a bitscore  
262 = 0 in the corresponding reference eukaryote database and bitscore > 75 from the  
263 bacterial database. We also screened the output for adjacent 1 kbp pieces that

264 contain bacterial matches and reference eukaryote bitscore = 0 and fused these  
265 adjoining pieces for analysis.

266 LGT candidate regions were then manually curated as follows: each candidate region  
267 was searched with BLASTn to the NCBI nr/nt database. If this search indicated that the  
268 region's nucleotide sequence was similar or identical to the nucleotide sequence of a  
269 known gene in related insects, it was discarded as a likely conserved insect gene.  
270 Regions were retained only when the matches to other insects were sporadic, as our  
271 experience has indicated that these can be independent LGTs into different lineages.  
272 If no match was found, the region was additionally searched with BLASTx to the NCBI  
273 nr/nt database. If this second search also resulted in no hits to multiple insect proteins,  
274 it was called an LGT candidate. In this case, we additionally identified the best  
275 bacterial match using the NCBI nr and protein databases. Using the gene annotation  
276 information, we then evaluated the flanking genes within the scaffold to determine  
277 whether they were eukaryotic or bacterial, we determined whether the LGT region  
278 was associated with an annotated gene within the insect genome, and we observed  
279 with transcriptome data if RNA sequencing data showed evidence of transcriptional  
280 activity in the LGT region. This short list is available in the supplementary materials  
281 (S1.3.4.)

## 282 **RNA extraction, library construction, and sequencing**

283 Juveniles, adult males, and adult females (approximately 4-5 of each) were prepared  
284 for RNAseq using the RNeasy Blood and Tissue Kit (Qiagen). Individuals were placed in  
285 a 1.5 mL safelock tube along with 5-8 one mm glass beads placed in liquid nitrogen  
286 and then shaken for 30 s in a Silamat S6 shaker (Ivoclar Vivadent). RNeasy Blood and  
287 Tissue Kit (Qiagen) was used according to manufacturer's instructions. Samples were  
288 assessed for quality and RNA quantity using an Invitrogen Qubit 2.0 fluorometer and

289 the RNA BR Assay Kit (Thermo Fisher Scientific). These three RNA samples were then  
290 processed by Novogene Bioinformatics Technology Co., Ltd., (Beijing, China) using  
291 poly(A) selection followed by cDNA synthesis with random hexamers and library  
292 construction with an insert size of 550-600 bp. Paired-end sequencing was performed  
293 on an Illumina HiSeq 4000 according to manufacturer's instruction.

#### 294 **Gene finding, transcriptome assembly, and annotation**

295 For the *ab initio* gene finding, a training set was established using the reference  
296 genome of *D. melanogaster* (Genbank: GCA\_000001215.4; Release 6 plus ISO1 MT)  
297 and the associated annotation. The training parameters were used by GlimmerHMM  
298 v3.0.1 for gene finding in the *N. tenuis* genome assembly v1.5 (Majoros et al., 2004). For  
299 homology-based gene prediction, GeMoMa v1.6 was used with the *D. melanogaster*  
300 reference genome alongside our RNAseq data as evidence for splice site prediction  
301 (Keilwagen et al., 2016). For evidence-based gene finding, each set of RNAseq data  
302 (male, female, and juvenile) was mapped to the *N. tenuis* genome separately with  
303 TopHat v2.0.14 with default settings (Trapnell et al., 2009). After mapping, Cufflinks  
304 v2.2.1 was used to assemble transcripts (Trapnell et al., 2010). CodingQuarry v1.2 was  
305 used for gene finding in the genome using the assembled transcripts, with the  
306 strandness setting set to 'unstranded' (Testa et al., 2015).

307 The tool EvidenceModeler (EVM) v1.1.1 was used to combine the *ab initio*, homology-  
308 based, and evidence-based information, with evidence-based weighted 1, *ab initio*  
309 weighted 2, and homology-based weighted 3 (Haas et al., 2008). The resulting amino  
310 acid sequences were searched with BLASTp v2.2.31+ on a custom database  
311 containing all SwissProt and Refseq genes of *D. melanogaster* (Acland et al., 2014;  
312 Boutet et al., 2008; Camacho et al., 2009). The top hit for each amino acid  
313 sequence/gene was retained and its Genbank accession number and name are

314 found within the annotation. If no hit was found, an additional search in the NCBI non-  
315 redundant protein database (nr) was performed to obtain additional homology data.

### 316 **Functional annotation and GO term analysis**

317 Gene attributes from the annotation were used to construct a list of genes to be used  
318 in Gene Ontology (GO) term classification. Duplicate accession numbers were  
319 removed alongside cases where no BLAST hit was found. The remaining accession IDs  
320 were converted into UniProtKB accession IDs using the UniProt ID mapping feature  
321 (Huang et al., 2011). These UniProtKB accession IDs were in turn used with the DAVID  
322 6.8 Functional Annotation Tool to assign GO terms to each accession ID with the *D.*  
323 *melanogaster* background and generate initial functional analyses (Huang et al.,  
324 2009a, 2009b).

### 325 **Ortholog cluster analysis and comparison**

326 The complete gene set of *Nesidiocoris tenuis* was compared to those of three  
327 additional hemipteran species: the bed bug *Cimex lectularis*, the brown marmorated  
328 stinkbug *Halyomorpha halys* (Stål) (Hemiptera: Pentatomidae), and the pea aphid  
329 *Acyrtosiphon pisum* (Harris) (Hemiptera: Aphididae) using OrthoVenn2 (L. Xu et al.,  
330 2019). The gene set of *A. pisum* is the 2015 version from AphidBase (Legeai et al., 2010;  
331 Richards et al., 2010) as maintained on the OrthoVenn2 server. The *H. halys* 2.0  
332 complete gene set was used (Lee et al., 2009) along with the complete gene set of *C.*  
333 *lectularis* (Clec 2.1, OGSv1.3) (Benoit et al., 2016; Thomas et al., 2020), both of which  
334 were retrieved from the i5K Workspace (Poelchau et al., 2015). An ortholog cluster  
335 analysis was performed on all four gene sets via OrthoVenn2 with the default settings  
336 of E-values of 1e-5 and an inflation value of 1.5.

### 337 **Cytogenetic analysis**

338 *Slide preparations*

339 To determine karyotype, *N. tenuis* individuals were obtained from the KBS population  
340 and prepared for cytogenetic experiments. Chromosomal preparations were  
341 prepared from the female and male reproductive organs of adults and juveniles by  
342 spreading technique according to Traut (1976) with modifications from Mediouni *et al.*  
343 2004 (Mediouni *et al.*, 2004; Traut, 1976). After inspection via stereomicroscope to  
344 confirm presence of chromosomes, slides were dehydrated in ethanol series (70, 80,  
345 and 100%, 30 s each) and stored at -20°C for future use.

346 *18S rDNA probe preparation and Fluorescence In Situ Hybridisation (FISH)*

347 To confirm the presence of 18S rDNA sequences in the assembled genome, the  
348 previously published partial 18S rDNA sequence of *N. tenuis* (GU194646, Jung and Lee,  
349 2012) was used as a BLAST query against the *N. tenuis* v1.5 genome. To verify sequence  
350 homology, the obtained 18S rDNA sequences were subsequently compared to the  
351 previously published sequence.

352 For preparation of the probe, we isolated gDNA from two *N. tenuis* females with the  
353 NucleoSpin DNA Insect Kit (Macherey-Nagel, Düren, Germany) according the  
354 manufacturer's protocol. gDNA was used as template in PCR to amplify the 18S rDNA  
355 sequence using primers 18S-1 and 18S-4 as described in Jung and Lee (2012).  
356 Obtained products were purified using the Wizard SV Gel and PCR Clean-Up System  
357 (Promega, Madison, WI, USA), and subsequently cloned using the pGEM-T Easy Vector  
358 System (Promega, Madison, WI) according to the manufacturer's protocol. Plasmids  
359 were extracted from positive clones with the NucleoSpin Plasmid kit (Macherey-Nagel)  
360 following the manufacturer's protocol, confirmed by sequencing (SEQme, Dobříš,  
361 Czech Republic), and used as template in PCR with the 18S-1 and 18S-4 primers. PCR-

362 products were purified, and used as template for labelling by a modified nick  
363 translation protocol as described by Kato *et al.* (2006) with modifications described in  
364 Dalíková *et al.* (2017), using biotin-16-dUTP (Jena Bioscience, Jena, Germany) and an  
365 incubation time of 35 minutes at 15°C (Dalíková *et al.*, 2017a; Kato *et al.*, 2006).  
366 Fluorescence *in situ* hybridization (FISH) was performed as described in Sahara *et al.*  
367 (1999) with modifications described in Zrzavá *et al.* (2018) (Sahara *et al.*, 1999; Zrzavá  
368 *et al.*, 2018).

### 369 *Sex chromosome identification*

370 Determination of the sex chromosome constitution is important for the assembly of the  
371 *N. tenuis* genome to identify any potential missing information due to sequencing a  
372 single sex, as well as add to knowledge on sex chromosomes in Miridae. Comparative  
373 Genomic Hybridization (CGH), and Genomic *In Situ* Hybridization (GISH) were,  
374 therefore, used to identify the sex chromosomes of *N. tenuis*. The reproductive organs  
375 of adult females were dissected out to avoid potential male gDNA contamination, as  
376 the mated status was unknown, after which remaining tissue was snap-frozen in liquid  
377 nitrogen and stored at -20°C until further use. Adult males were not dissected but  
378 otherwise treated the same. Female and male gDNA was extracted from 10-20 pooled  
379 individuals using cetyltrimethylammonium bromide (CTAB) gDNA isolation with  
380 modifications (Doyle and Doyle, 1990). Samples were mechanically disrupted in  
381 extraction buffer (2% CTAB, 100 mM Tris-HCl pH 8.0, 40 mM EDTA, 1.4 M NaCl, 0.2%  $\beta$ -  
382 mercaptoethanol, 0.1 mg/mL proteinase K), and incubated overnight at 60°C with  
383 light agitation. An equal volume chloroform was added, tubes were inverted for 2 min,  
384 and samples were centrifuged 10 min at maximum speed. The aqueous phase was  
385 transferred to a new tube, RNase A (200 ng/ $\mu$ L) was added and samples were  
386 incubated 30 min at 37°C to remove RNA. DNA was precipitated by adding 2/3



387 volume isopropanol, gently inverting the tubes, and centrifugation for 15 min at  
388 maximum speed. Pellets were washed twice with 70% ethanol, air-dried briefly, and  
389 dissolved overnight in sterile water. DNA was stored at -20°C until further use. Probes  
390 were prepared with 1 µg gDNA using Cy3-dUTP (for female gDNA), or fluorescein-dUTP  
391 (for male gDNA) (both Jena Bioscience, Jena, Germany) by nick translation  
392 mentioned above with an incubation time of 2-2.5 hours at 15°C. CGH and GISH were  
393 performed according to Traut et al. (1999) with modifications described in Dalíková et  
394 al. (2017) (Dalíková et al., 2017b; Traut et al., 1999).

### 395 *Detecting a telomeric motif*

396 Initially, we searched both the raw sequencing data and the assembled genome for  
397 presence of the ancestral insect telomere motif (TTAGG)<sub>n</sub>, which is known to be absent  
398 in several Miridae (Grozeva et al., 2019; Kuznetsova et al., 2011), and tested for its  
399 presence using Southern dot blot in *N. tenuis*. gDNA was isolated from *N. tenuis*, and  
400 positive controls *E. kuehniella* and *T. infestans*, using CTAB DNA isolation described  
401 above. DNA concentrations were measured by Qubit 2.0 (Broad Spectrum DNA Kit)  
402 (Invitrogen) and diluted to equalize concentrations, after which 500 ng and 150 ng of  
403 each specimen was spotted on a membrane and hybridized as described in (Dalíková  
404 et al., 2017b). As a negative control, an equal amount of sonicated DNA from the  
405 chum salmon *Oncorhynchus keta* (Walbaum) (Salmoniformes: Salmonidae) (Sigma-  
406 Aldrich, St. Louis, MO, USA), was spotted on the same membrane. Probe template was  
407 prepared using non-template PCR according to Sahara et al. (1999) and labelling with  
408 digoxigenin-11-dUTP (Jena Bioscience) was performed using nick translation, with an  
409 incubation time of 50 min according to Dalíková et al. (2017) (Dalíková et al., 2017a;  
410 Sahara et al., 1999). Absence of the insect telomere motif (TTAGG)<sub>n</sub> was confirmed by  
411 dot blot, and three sequence motifs, (TATGG)<sub>n</sub>, (TTGGG)<sub>n</sub>, and (TCAGG)<sub>n</sub>, were

412 selected as potential telomeric motifs in *N. tenuis* based on high copy numbers in the  
413 genome and sequence similarity to the ancestral insect telomere motif. Copy numbers  
414 were determined by Tandem Repeat Finder (TRF, version 4), on collapsed quality  
415 filtered reads corresponding to 0.5x coverage with default numeric parameters  
416 except maximal period size, which was set to 25 bp (Benson, 1999). TRF output was  
417 further analysed using Tandem Repeat Analysis Program (Sobreira et al., 2006). Probe  
418 template, and subsequent labelling of the probes, was done as described above with  
419 slight alterations. To obtain optimal length of fragments for labelling, non-template  
420 PCR was performed with reduced primer concentrations (50 nM for each primer). In  
421 addition, probes were labelled by biotin-16-dUTP (Jena Bioscience) using nick  
422 translation as described above, with an incubation time of 50 min. FISH was performed  
423 as described above for 18S rDNA.

#### 424 *Repeat identification and visualization*

425 To assess the repetitive component of the *N. tenuis* genome, we used RepeatExplorer,  
426 version 2, on trimmed and quality-filtered reads with default parameters (Novák et al.,  
427 2013). Repeats with high abundance in the genome were selected and amplified  
428 using PCR. These products, named Nt\_rep1, were additionally cloned and template  
429 for probe labelling was prepared from plasmid DNA as described above, see 18S rDNA  
430 probe preparation. Probes were labelled by PCR in a volume of 25 µL consisting of  
431 0.625 U Ex Taq polymerase (TaKaRa, Otsu, Japan), 1x Ex Taq buffer, 40 µM dATP, dCTP,  
432 and dGTP, 14.4 µM dTTP, 25.6 µM biotin-16-dUTP (Jena Bioscience), 400 nM of forward  
433 and reverse primer, and 1 ng of purified PCR-product. The amplification program  
434 consisted of an initial denaturing step at 94°C for 3 min, followed by 35 cycles of 94°C  
435 for 30 s, 55°C for 30 s, 72°C for 1 min, and a final extension at 72°C for 2 min. The FISH  
436 procedure was performed as described for 18S rDNA. Abundance and distribution of

437 Nt\_rep1 in the assembled genome was assessed using NCBI Genome Workbench  
438 version 2.13.0. The complete list of primers used in this study can be found in Table  
439 S1.6.2.

#### 440 **Pooled sequencing and population analysis**

441 For this project, it was important to use existing data wherever possible to test the utility  
442 of the genome and possible research avenues. Therefore, we analysed whole  
443 genome sequence data originally generated for another *N. tenuis* genome assembly  
444 project that has not been published before. In the original set-up, ten females were  
445 collected from the KBS population for pooled sequence analysis. DNA was isolated  
446 from this pooled cohort using the “salting-out” method as described in Sunnucks and  
447 Hales with a final volume of 20  $\mu$ L (Sunnucks and Hales, 1996) and then treated with 2  
448  $\mu$ L RNase. The paired-end library was sequenced on an Illumina HiSeq2500 platform by  
449 MacroGen Inc. (Seoul, Korea) with read sizes of 100 bp. Reads were assessed for quality  
450 using FastQC (Andrews et al., 2015) and adapters were trimmed with Trimmomatic  
451 (Bolger et al., 2014). Following quality filtering, reads with phred scores lower than 20  
452 were discarded. Heterozygosity was calculated using jellyfish v2.3.0 and  
453 GenomeScope v1.0 with a k-mer size of 21 and default parameters (Marçais and  
454 Kingsford, 2011; Vurture et al., 2017).

455 Instead of genome assembly, these whole genome sequence reads were adjusted  
456 and subsequently used in a pooled sequencing (pool-seq) population analysis with  
457 our genome. First, the reads were randomly subsampled to a coverage of 10X (in a  
458 pool with 10 females, this results in approximately 1X coverage per female) using CLC  
459 Genomics Workbench 12 (Qiagen). Using the PoPoolation v1.2.2 pipeline (Kofler et al.,  
460 2011), these reads were aligned to an adapted v1.5 genome, where scaffolds smaller  
461 than 10,000 bp were removed, and aligned reads were binned into windows using the

462 bwa and samtools packages (Li et al., 2009). Pileup files and the scripts from the  
463 PoPoolation pipeline were used to produce variance sliding windows analyses of  
464 neutrality, *Tajima's D*, and nucleotide diversity, *Tajima's pi* ( $\pi$ ) with default settings and  
465 a pool size of 40. Window and step sizes of both 10,000 and 5,000 were tested, as well  
466 as using the "basic-pipeline/mask-sam-indelregions.pl" pipeline to mask indel regions  
467 of the SAM file – this ensures that indel regions are not calculated. Of the 18,000,000  
468 reads, 817,226 had regions of indels masked.

469

## 470 **Results**

### 471 **Species origin, description, and data availability**

472 The presence of *Wolbachia* in the KBS biological control stock was confirmed  
473 (Electrophoresis gel image in supplementary material, S1.1). All sequence data  
474 generated, including raw reads, assembly, and annotation, can be found in the EMBL-  
475 EBI European Nucleotide Archive (ENA) under BioProject PRJEB35378. An additional,  
476 complete annotation file (.gff) is also available (Ferguson, 2020).

### 477 **Genome assembly and size**

478 The single adult female *N. tenuis* yielded 424 ng total DNA. The 10X Genomics  
479 Chromium reaction and subsequent Illumina sequencing resulted in more than 212  
480 million paired-end reads. The inferred heterozygosity, based on GenomeScope, was  
481 between 1.675% and 1.680% for a k-mer size of 21, and between 1.250% and 1.253%  
482 for a k-mer size of 48. Genome size estimates at this point were between 306 Mbp (k-  
483 mer=21) and 320 Mbp (k-mer=48). Following assembly with Supernova, assembly v1.0  
484 was approximately 388 Mbp in size and comprised of 44,273 scaffolds (5.91%  
485 ambiguous nucleotides).

486 Assembly v1.0 was then assessed for contamination with a preliminary search against  
487 the NCBI for bacterial homology (see below for more details). Several scaffolds with  
488 high amounts of bacterial sequence contamination were identified, indicating that  
489 further decontamination of the assembly was required. A decontamination pipeline  
490 was used to identify and remove a total of 3,043 scaffolds, while those identified as  
491 potential examples of LGT were kept. From the remainder, an additional 4,717 were  
492 identified as being identical duplicates and were removed. At this point, the resulting  
493 assembly was finalised and designated v1.5. This assembly is 355 Mbp in size, consisting  
494 of 36,513 scaffolds (6.29% ambiguous nucleotides). Quality and completeness of v1.5  
495 using BUSCO indicated a completeness of 87.5% (65.6% single copy orthologs, 21.9%  
496 duplicated orthologs), while 7.1% orthologs were fragmented and 5.4% were missing  
497 (n=1658).

498 Initially, the genome size of *N. tenuis* was estimated by flow cytometry to be 232 Mbp,  
499 with a confidence interval of 20 Mbp (See supplementary material S1.2 for more  
500 details). Further estimates via k-mer analysis of sequence data in GenomeScope  
501 indicated an expected genome size of 306 Mbp (k-mer=21) or 320 Mbp (k-mer=48).  
502 Both the flow cytometry and sequence data estimates are smaller than the 355 Mbp  
503 of the final assembly (v1.5). In total, the *N. tenuis* genome has 36,513 scaffolds, with  
504 the largest scaffold being 1.39 Mbp, though the majority of scaffolds are under 50,000  
505 bp in size. The number of gaps per 100 kbp is 6292.10 (6.29% of the genome). Details  
506 on the assemblies can be found in Table 1.

## 507 **Assessment of potential symbionts and LGT candidates**

### 508 *Potential symbionts*

509 The initial assembly (v1.0) was decontaminated using two bacterial decontamination  
510 pipelines: the first pipeline broadly utilised BLASTn to identify scaffolds with high  
511 amounts of bacterial sequences against the NCBI nr database, while the second  
512 pipeline is more specified and uses BLASTn against a list of known contaminants and  
513 symbionts (S1.3) and is adapted from previous work (Wheeler et al., 2013). The first  
514 decontamination pipeline identified and removed 1,443 scaffolds with high bacterial  
515 content, and the second decontamination pipeline identified and removed an  
516 additional 1,600 scaffolds alongside potential LGT events. All removed scaffolds are  
517 available in S1.3. The hits from the second pipeline were used to create a list of  
518 potential contaminants or symbionts of this particular *N. tenuis* individual used for  
519 whole genome sequencing according to genus, base pair content, and number of  
520 scaffolds affected (Table 2). The majority of these scaffolds (1470) are under 5 kbp in  
521 length, with an additional 61 scaffolds falling between 5-10 kbp. The ten largest  
522 scaffolds are putatively associated with *Pantoea* and relatives (three of 561,7472 bp,  
523 205,621 bp, and 131,905 bp), *Sodalis* (326,101 bp), *Erwinia* (254,660 bp; 220,307 bp;  
524 154,581 bp), and *Citrobacter* (239,269 bp; 190,839 bp). We emphasize that these  
525 “calls” are very preliminary, as they are based on the most frequent hits in the bacterial  
526 matches in each scaffold, rather than comprehensive gene annotations.  
527 Nevertheless, they do indicate a range of bacterial types associated with *N. tenuis*,  
528 and the scaffold assemblies are likely to contain some complete or near complete  
529 bacterial genomes of interest.

530 Sorting scaffolds across the range of bacterial genera matches gives 131 genera with  
531 some substantial representation: *Erwinia* (2,078,531 bp), *Pantoea* (2,226,778 bp),  
532 *Citrobacter* (594,902 bp), *Sodalis* (355,847 bp), *Cronobacter* (314,511 bp), and  
533 *Rickettsia* (483,217 bp) (Table 2). In addition to known symbiont *Rickettsia*, previously  
534 established via PCR and known symbiont *Wolbachia* is also present in the results

535 (137,109 bp) (Table 2). Multiple genera of bacteria can be found on a single scaffold,  
536 likely due to misassembly. The full list of bacterial scaffolds and multiFASTA file is  
537 available with details in supplementary materials S1.3.2 and S1.3.3.

### 538 *LGT Candidates*

539 We continued with our detection of potential LGT events by further assessing a handful  
540 of strong candidates. Two of these regions occur on scaffolds 22012 and 22013, which  
541 are of similar length (22,634 bp and 22,957 bp, respectively) and are highly similar on  
542 a nucleotide level. Scaffold 22013 appears to have additional nucleotides on each  
543 flanking side, with some indels and SNPs between the two scaffolds. The putative LGT  
544 region in question is belonging to or gained from a *Sodalis* species, coding for  
545 phenazine biosynthesis protein *PhzF* (OIV46256.1). This region also showed  
546 transcriptional support, and is flanked by conserved insect genes, most immediately  
547 *Rab19* (NP\_523970.1) on one side and an uncharacterized protein, Dmel\_CG32112  
548 (NP\_729820.2), on the other side. Two additional LGTs were found in the current  
549 assembly that match *Rickettsia* sequences (scaffolds 4712 and 27281), which contain  
550 a segment of the rickettsial genes *elongation factor G* and *AAA family ATPase* genes,  
551 respectively. One corresponds to a gene model, while the other does not, and there  
552 is no evidence of expression for either in the current male, female, and mixed sex  
553 juvenile RNA sequencing data. More information on these candidate regions can be  
554 found in S1.3.4.

### 555 ***Ab initio* gene finding, transcriptome assembly, and annotation**

556 To obtain a comprehensive set of transcripts for *N. tenuis*, three separate libraries of  
557 multiple individuals were prepared – males, females, and juveniles from different  
558 stages of mixed sex. More than 77 million 150 bp pair-end reads were generated.

559 Filtering the reads for quality led to a slightly reduced total of 76,711,096 paired reads  
560 (male: 28,413,231 paired reads; female: 24,075,901 paired reads; juvenile: 24,221,964  
561 paired reads) to be used for evidence-based gene finding. The mapping and  
562 assembling of reads of the three individual samples as well as the pooled reads  
563 resulted in four transcriptomes: male, female, juvenile, and the combined  
564 transcriptome.

565 The male, female, juvenile, and combined annotations from the evidence-based  
566 gene finding was used alongside homology-based findings and *ab initio* annotations  
567 in a weighted model, resulting in complete annotations for the assembly. When gene  
568 name assignment via the SwissProt database resulted in “no hit,” tracks are named  
569 “No\_blast\_hit.” This occurred in 1,556 mRNA tracks and represents approximately 6%  
570 of the official gene set. The majority of tracks were annotated with reference to  
571 SwissProt or GenBank accession number of the top BLASTp hit.

572 CodingQuarry predicted 56,309 genes from the mapped transcript evidence, while  
573 *ab initio* gene finding using GlimmerHMM resulted in 39,888 genes and homology-  
574 based gene finding with GeMoMa resulted in 6,028 genes. The complete gene set for  
575 *N. tenuis* was created using EVIDENCEModeler, where a weighted model using all  
576 three inputs resulted in a complete gene set of 24,688 genes.

### 577 **Functional annotation and GO term analysis**

578 The complete gene set of 24,668 genes was deduplicated and genes with no  
579 correlating BLASTp hit were removed. The remaining 11,724 genes were mapped to  
580 UniProtKB IDs, resulting in 11,261 genes with a matching ID after another round of  
581 deduplication (80 duplicates found). The remaining 383 genes either did not match to  
582 a UniProt KB ID or were considered obsolete proteins within the UniParc database.



583 DAVID used 8,920 genes for the functional annotation analysis, of which 78.4% (6503)  
584 contribute to 19 biological processes, 75.8% (6826) contribute to 100 different cellular  
585 components, and 72.8% (6032) contribute to 91 categories of molecular functions  
586 (genes can code to multiple GO terms). The remaining genes were uncategorized.  
587 Data linking the genes to the GO terms, the DAVID Gene List Report, and the DAVID  
588 Gene Report are available in S1.4.

### 589 **Ortholog cluster analysis**

590 The complete gene set of *N. tenuis* was compared to those of three additional species:  
591 the bed bug *C. lectularis*, the brown marmorated stinkbug *H. halys*, and the pea aphid  
592 *A. pisum* using OrthoVenn2. The ortholog analysis summary is presented in Table 3 and  
593 visualized in Figure 1. *N. tenuis* has a similar number of clusters (8,174) as compared to  
594 *C. lectularis*, *H. halys* and *A. pisum* (7,989; 9,584; and 8,765, respectively). In total 14,512  
595 clusters are assigned, 12,964 of which are orthologous clusters (contains at least two  
596 species), and the remaining 1,548 are single-copy gene clusters. There are 9,136  
597 singleton clusters in *N. tenuis*, 3,573 in *C. lectularis*, 2,170 in *H. halys*, and 7,298 in *A.*  
598 *pisum*. The amount of singleton clusters, i.e. proteins that do not cluster, indicate that  
599 *N. tenuis* differs the most from the other species, as 37.04% of the proteins are  
600 singletons. Just over half of the orthologs cluster with *N. tenuis*, where 6,338 clusters are  
601 outside of *N. tenuis* as compared to the 8,174 clusters within *N. tenuis*. The final protein  
602 set from *N. tenuis* used in this analysis is available, see S1.5.

### 603 **Karyotype analysis**

604 Karyotype analysis revealed  $2n=32$  chromosomes in both females and males (Figure  
605 2a and b). All chromosomes are relatively small with one larger pair of submetacentric  
606 chromosomes in females (Figure 2a). In males (Figure 2b), we were unable to obtain

607 mitotic chromosomes of reasonable quality as in females and therefore we were  
608 unable to clearly identify these larger chromosomes. Screening of multiple nuclei  
609 showed sporadic deviations of the karyotype in some individuals. This was the result of  
610 supernumerary chromosomes (B chromosomes) which were clearly visible in (meiotic)  
611 pachytene stage as distinctly smaller chromosomes (Figure 2c, three B chromosomes).

## 612 **Analysis and localization of 18S rDNA**

613 The 18S rDNA gene is often used as a cytogenetic marker in comparative evolutionary  
614 studies due to its ease of visualization on the chromosomes caused by high copy  
615 number and cluster organisation in animal (Sochorová et al., 2018) and plant (Gomez-  
616 Rodriguez et al., 2013) genomes. The published partial 18S sequence of *N. tenuis*  
617 (GU194646.1) and the 18S sequence identified in this study were compared to each  
618 other revealing some differences between the sequences. The published sequence  
619 consists of two fragments of 869 bp and 739 bp, which are, respectively, 99.7% and  
620 94.2% homologous to our identified partial 18S sequence. Interestingly, the second half  
621 of our isolated 18S sequence is more homologous to a *Macrolophus sp.* partial 18S  
622 sequence (EU683153.1), i.e. 97.8%, than to the previously published *N. tenuis* sequence.  
623 A BLAST search against the *N. tenuis* genome with either of the *N. tenuis* 18S sequences  
624 resulted in four gene copies in both cases, each located on a different scaffold.  
625 However, RepeatExplorer analysis estimated 98 18S rDNA copies with the obtained  
626 genome size of 355 Mbp. Using FISH with the 18S rDNA probe we finally showed that  
627 the major rDNA forms a single cluster located terminally on a pair of homologous  
628 chromosomes (Figure 3).

## 629 **Identification of sex chromosomes**

630 The common sex chromosome constitution in Miridae is the male-heterogametic  
631 XX/XY system. To identify the sex chromosome constitution and estimate sex  
632 chromosome differentiation in *N. tenuis* we employed GISH and CGH experiments. The  
633 GISH results clearly revealed a single chromosome densely labelled by the male-  
634 derived probe, caused by male-enriched repetitive DNA and/or male-specific  
635 sequences which is typical for the Y chromosome (Figure 4). In addition, the Nucleolus  
636 Organizer Region (NOR; including 18S rDNA) was observed as well, as is often the case  
637 in GISH experiments due to the presence of highly repetitive sequences in the rDNA  
638 cluster. The NOR is clearly located terminally on a pair of autosomes, corroborating our  
639 18S rDNA FISH results.

640 To further study the differentiation of the sex chromosomes we carried out CGH  
641 experiments on chromosome preparations of both sexes (Figure 5). All chromosomes  
642 were labeled evenly by the female and male probes with the exception of the largest  
643 chromosome pair. Both sex chromosomes were highlighted with DAPI (Figure 5a, e),  
644 indicating that they are both A-T rich and largely composed of heterochromatin. In  
645 females, the largest chromosome pair was labelled more by the female probe than  
646 the male probe indicating that these chromosomes contain sequences with higher  
647 copy numbers in females, and are thus the X chromosomes, as seen in Figure 5a-d. In  
648 male meiotic nuclei (Figure 5e-h), two types of nuclei can be discerned, where the  
649 largest chromosome was labelled more by either the female probe or the male probe  
650 corresponding to the X, and Y chromosome, respectively, whereas the autosomes  
651 were labelled equally by both probes.

## 652 **Identification and mapping of abundant repeats**

653 RepeatExplorer software was used on reads with GEM barcodes removed to identify  
654 the most abundant repeats in the genome of *N. tenuis* (results available in

655 supplementary materials, see Table S1.6.1). The most abundant repeat, Nt\_rep1,  
656 makes up approximately 3% of the genome estimated by RepeatExplorer. Analysis on  
657 the assembled genome, using a coverage cut-off value of 70%, reveals that Nt\_rep1  
658 is present on 3190 scaffolds (8.737% of the assembled scaffolds), with a maximum of  
659 17 copies on a single scaffold. According to the assembled genome, Nt\_rep1 makes  
660 up approximately 0.8% of the entire genome (Figure S1.6.2). We subsequently mapped  
661 Nt\_rep1 to the chromosomes of *N. tenuis* using FISH. The repeat is located on most  
662 chromosomes and is accumulated in sub-telomeric regions (Figure 6). Additional  
663 signals were identified on the X chromosome indicating a higher number of this repeat  
664 (Figure 6a-c). This increase in frequency is specific to the X chromosome and is not  
665 found on the Y chromosome of *N. tenuis* (Figure 6d-f).

#### 666 **Testing of candidate telomere motifs**

667 Analysis of the raw sequencing data and the assembled genome both revealed low  
668 numbers of the insect telomere motif (TTAGG)<sub>n</sub> (Frydrychová et al., 2004) in *N. tenuis*,  
669 i.e. approximately 98 repeats per haploid genome. This translates into approximately  
670 three copies of the repeat per chromosome end, much lower than expected for a  
671 telomeric motif. These low copy numbers were additionally confirmed using Southern  
672 dot blot (Figure S1.6.1). Other candidate telomere motifs previously identified by TRF  
673 analysis, (TATGG)<sub>n</sub>, (TTGGG)<sub>n</sub>, and (TCAGG)<sub>n</sub>, were examined by FISH for their  
674 distribution in the genome. They were found scattered throughout the genome but  
675 lacked a clear accumulation at the terminal regions of the chromosomes (not shown).  
676 Therefore, these sequences can also be excluded as telomeric motifs in *N. tenuis*.

#### 677 **Pooled sequencing analysis**

678 Using previously generated whole genome sequencing of ten females from the KBS  
679 population, we were able to estimate genetic diversity of the commercial population  
680 via a pool-seq population analysis. Read coverage was randomly subsampled to 10X  
681 coverage (18,000,000 reads). Additionally, we used a modified v1.5 *N. tenuis* genome  
682 with scaffolds of less than 10,000 bp removed. This was to ensure that window sliding  
683 was not being inflated on scaffolds smaller than the window size. This reduced the  
684 genome from 36,513 scaffolds to 7,076, however, the reduced genome still contained  
685 72.23% of the genome in terms of size (256,487,768 bp).

686 Three runs of PoPoolation were performed with varied window size, step size, and the  
687 masking of indel regions. The default setting, window size and step size of 10,000 bp,  
688 yielded similar results as the adjusted window size and step size of 5,000 bp, while  
689 differences were apparent when indel regions were masked. As such, results of  
690 window size and step size 10,000 bp with indel regions mapped are reported here  
691 (other results available in S1.7). The variance sliding program created 28,833 windows  
692 of 10,000 bp with mapped reads, of which 5,913 were sufficiently covered with reads  
693 to calculate values per window (coverage  $\geq 0.60$ ). Genome-wide, the nucleotide  
694 diversity ( $\pi$ ) is 0.0080 and *Tajima's D* is -0.0355. Figure 7 shows the *Tajima's D* (7a) and  
695  $\pi$  (7b) for the ten largest scaffolds, all containing gene annotations, arranged in order  
696 of size. These ten scaffolds represent approximately 1.7% of the genome (6,135,756  
697 bp), and varied in terms of window coverage (from no coverage to full coverage) as  
698 well as both *Tajima's D* and  $\pi$ . These ten scaffolds are a snapshot of the whole  
699 genome, summarised in Table 4, whereas genome-wide results can be found in S1.7.

700

701 **Discussion**

702 **Assembly and Annotation**

703 Presented here is the genome of *N. tenuis*, a biological control agent used throughout  
704 the Mediterranean in tomato crops. We chose to use 10X Genomics linked-read  
705 sequencing strategy as it best suited the challenges that come with working with a  
706 relatively small and long-lived mirid such as *N. tenuis*. Assembling a genome is easiest  
707 with reduced heterozygosity in the input sample, often through single individual  
708 sampling or inbreeding (Ekblom and Wolf, 2014; Richards and Murali, 2015). This proved  
709 an initial challenge for the sequencing strategy of *N. tenuis*, as they are too small for a  
710 single individual to yield the minimum amount of DNA required for a traditional NGS  
711 library, and an inbred population was not readily available for sequencing. Therefore,  
712 10X Genomics linked-read sequencing was the immediate solution for which a small  
713 amount of input DNA from a single individual would yield a highly contiguous genome.

714 Assemblies v1.0 and v1.5 contain 5.91% and 6.29% ambiguous nucleotides, while still  
715 offering a relatively high BUSCO score, with the final decontaminated assembly (v1.5)  
716 having a completeness of 87.5% of the insect\_0db9 ortholog dataset. However, the  
717 final assembled genome size is approximately 150 Mbp larger than was expected  
718 based on flow cytometry data, and we suggest the assembly presented here can best  
719 be improved in terms of accuracy and contiguity with long reads from an inbred  
720 sample. This discrepancy between estimated genome size and assembled genome  
721 size may also be due to the ambiguous nucleotides inserted into the genome during  
722 the assembly process. Making up just over 6% of the final assembled genome, that is  
723 approximately 2.2 Mbp of ambiguous nucleotides. However, most of the genome  
724 inflation is likely due to residual contamination along with duplicate scaffolds that  
725 remain after removing 100% identical ones.

726 Annotation via evidence-based, homology-based, and *ab initio* models resulted in  
727 24,668 genes. Compared to other assemblies within the hemipteran order, such as *C.*  
728 *lectularis*, with a genome size of 650 Mbp and 12,699 genes (Thomas et al., 2020) or *A.*  
729 *pisum*, with a draft genome size of 464 Mbp and 36,195 genes (Richards et al., 2010),  
730 *N. tenuis* sits, in the middle in terms of genome size and number of genes. It is worth  
731 noting that of the 24,668 genes within the complete gene set, only 11,261 (45.7%)  
732 remained after UniProtKB mapping, of which 8,920 (36.2% of total) were used by DAVID  
733 for functional analysis. This is relatively low compared to similar genome projects, such  
734 as *Aphys gossypii* (Glover) (Hemiptera: Aphididae), where 49.2% of the gene set could  
735 be used for GO term analysis (Quan et al., 2019). However, we used different methods  
736 of gene prediction and annotation which may explain the difference. The next step  
737 for the *N. tenuis* genome is manual annotation and curation, which would likely  
738 improve the GO term analysis, but this requires time and expertise. Still, we hope that  
739 other researchers will use and add to the annotation.

740 Comparing the current gene set of *N. tenuis* to other Hemipterans, the clustering  
741 identified considerable overlap, as 71% of the clusters that are found in *N. tenuis* were  
742 shared between the other species in the comparative analysis. Despite being more  
743 closely related to *C. lectularis* in terms of phylogeny, in terms of lifestyle, *N. tenuis* is far  
744 more similar to *A. pisum* and *H. halys*, and this is likely reflected in absolute number of  
745 proteins and clusters shared between the four species. The remaining 29% of clusters,  
746 as well as the singleton proteins, are indications for proteins unique to either *N. tenuis*  
747 or Miridae in general. Through the OrthoVenn2 website, the analysis performed here  
748 can be easily replicated, altered with other species of interest, and even improved  
749 upon if the complete gene sets are updated or with a newer software version. In our  
750 iteration, the 24,668 proteins of *N. tenuis* group into 8,174 clusters. 2,398 clusters are  
751 unique to *N. tenuis*, however, some of these genes have a relatively strong homology

752 to genes of one of the other species used in the analysis and could be incorrectly  
753 flagged as being unique. Reasons could be poor gene annotation quality resulting in  
754 a poor *in silico* protein translation, or too stringent clustering settings. Regardless, these  
755 2,398 clusters may be of interest to researchers working on zoophytophagy, the  
756 negative effects of *N. tenuis* on tomato as compared to other mirids, or broader  
757 questions such as phylogeny of the Hemiptera.

## 758 **Characterizing the Genome**

759 Every sequence and assembly strategy has benefits and drawbacks, and the 10X  
760 Genomics linked-read strategy is no exception. The technique requires only few  
761 nanograms of DNA for library preparation which allowed us to use a single individual  
762 and removed the need for inbreeding to reduce variation in the sequencing  
763 population. However, using a single individual from a closed and proprietary rearing  
764 process presented other challenges. These challenges were threefold: we had to deal  
765 with bacterial contamination as antibiotic treatment is not possible, we had to ensure  
766 that the single individual-derived assembly reflects reality in terms of genes present  
767 and structure, and we had to ensure that a single female-derived assembly is  
768 applicable for population-level analyses.

769 Contamination of genomes is a constant concern, and sequencing strategies should  
770 attempt to address the risks in the best way possible to deliver reliable genomes  
771 (Ekblom and Wolf, 2014). Equally so is the desire for inbred strains if multiple individuals  
772 are required to reach the micrograms of DNA necessary for NGS platforms. The inability  
773 to remove symbionts and microbiota using antibiotics administered to a few  
774 successive generations as well as the difficulty or inability to inbreed a strain is not  
775 restricted to *N. tenuis*. The sequencing strategy chosen for the mountain pine beetle,  
776 *Dendroctonus ponderosae* (Hopkins) (Coleoptera: Curculionidae), relied on assuming



777 the relatedness of several individuals as well as isolating the gut during the extraction  
778 process, and still additional post-assembly decontamination was required (Keeling et  
779 al., 2013). A linked-read strategy with low input requirements, such as the 10X  
780 Genomics library preparation that was chosen here, negates the need for a pool of  
781 inbred samples or controlling for relatedness. However, another potential benefit of  
782 controlled rearing such as those used in inbreeding (as opposed to be limited to wild-  
783 caught specimens, for example) is the ability to treat with antibiotics for multiple  
784 generations. Without the ability to do so, sequencing and assembly strategies rely  
785 heavily on post-sequencing decontamination strategies (both pre- and post-assembly  
786 are possible). That such post-assembly filtering strategies as used here for *N. tenuis* can  
787 be successful was shown by a less contaminated assembly and by the identification  
788 of potential LGTs.

### 789 **Beyond the Genome: Potential symbionts and LGT events**

790 The list of potential symbionts or pathogens generated in Table 2 represent both insect  
791 and plant pathogens, as well as potential environmental contaminants. In addition to  
792 the positive test for *Wolbachia* in the KBS population used here, *N. tenuis* is known to  
793 potentially harbour *Rickettsia* as an endosymbiont in addition to *Wolbachia* (Caspi-  
794 Fluger et al., 2014). *Rickettsia* genome sizes can range from 0.8 to 2.3 Mbp, also  
795 reflecting variation in levels of reductive evolution (Sachman-Ruiz and Quiroz-  
796 Castañeda, 2018). However, the total scaffold length identified here as from a  
797 *Rickettsia* falls below this range, likely indicating incomplete recovery of the genome  
798 from the insect sequencing. The potential symbionts revealed included not only  
799 *Wolbachia* and *Rickettsia*, but also other known insect symbionts, in addition to the  
800 usual lab contamination suspects.

801 *Sodalis* is a genus of bacterium symbiotic with various insects, including the tsetse fly  
802 and louse fly, louse and hemipteran species (Boyd et al., 2016). Genome sizes of  
803 *Sodalis* and close relatives range from 0.35 to 4.57 Mbp (Santos-Garcia et al., 2017).  
804 The relatively small total scaffold size found in our results (0.36 Mbp) likely reflects  
805 incomplete genome recovery in the assembly, but could also be due to genome size  
806 reduction, and is worthy of further investigation. *Erwinia* and *Pantoea* are closely  
807 related bacteria that are associated with plant pathology (Kamber et al., 2012; Zhang  
808 and Qiu, 2015) and both have been found in the midgut of stink bugs as vertically  
809 transferred plant-associated bacteria that become temporary endosymbionts of stink  
810 bugs until later replacement with another endosymbiont (Prado and Almeida, 2009).  
811 The genome sizes of *Erwinia* and *Pantoea* species typically range from 3.8-5.1 Mb. Our  
812 total scaffold size for the *Erwinia* and *Pantoea* are substantially smaller (2.078 Mbp and  
813 2.22 Mbp), but it is possible that these scaffolds belong to the same bacterium. In any  
814 case, the association of a zoophytophagous mirid bug with potential plant pathogens  
815 is noteworthy, especially in a biological control context.

816 As for *Serratia*, *Serratia marcescens* is both a common Gram-negative human-borne  
817 pathogen and a causal agent of cucurbit yellow vine disease (CYVD) (Abreo and  
818 Altier, 2019; Bruton et al., 2007). It is worth noting that in cases of CYVD, the transmission  
819 of *Serratia marcescens* from its vector the squash bug, *Anasa tristis* (De Geer)  
820 (Hemiptera: Coreidae), to host crops is via the phloem. Other *Serratia* spp. have been  
821 identified as insect symbionts previously, as have other potential symbionts found in  
822 the contaminated scaffolds, such as *Cedecea* spp. (Jang and Nishijima, 1990). The  
823 presence of *Dickeya* is an interesting find, as *Dickeya dadantii* has been established  
824 as a pathogen of *A. pisum*, while the pea aphid itself is a potential vector for the  
825 bacterium with regards to plants (Costechareyre et al., 2012). *Dickeya* spp. cause soft  
826 rot in various crops, including tomato. In a similar vein, the identification of *Ralstonia*,

827 as some members of this order, such as *Ralstonia solanacearum*, are soil-borne  
828 pathogens that causes wilt in several crop plants, including tomato (Lowe-Power et  
829 al., 2018). However, both *Dickeya* spp. and *R. solanacearum* infect the xylem while *N.*  
830 *tenuis* is a phloem-feeder. The scaffold lengths for *Citrobacter* and *Cronobacter* are  
831 also considerably below their typical genome sizes, and likely represent incomplete  
832 sequence recovery in the metagenomic sample.

833 All of these described associations are preliminary, as follow-up analyses against the  
834 entire NCBI database and proper bacterial gene annotation are lacking.  
835 Nevertheless, these putative bacterial associations of *N. tenuis*, their distribution within  
836 the insect, and their possible biological significance, warrants further investigation. It is  
837 important to note that some of these bacterial “contaminants” may actually represent  
838 large LGTs, which can be confirmed by identifying flanking sequences (e.g. using long-  
839 read technologies) and/or *in situ* chromosome hybridization analyses, such as done  
840 for the large LGT in *Drosophila ananassae* (Doleschall) (Diptera: Drosophilidae)  
841 (Dunning Hotopp et al., 2007). More research into the symbionts of *N. tenuis* via  
842 metagenomics would certainly shed some light on true symbionts (or pathogens)  
843 versus true contaminants, with potential implications for biological control and related  
844 research.

845 One of the LGT candidate genes that were detected following manual curation  
846 corresponds to phenazine biosynthesis protein *PhzF* (OIV46256.1), with the likely  
847 microbial source being a *Sodalis* species. Phenazines are heterocyclic metabolites  
848 with “antibiotic, antitumor, and antiparasitic activity,” but are also toxic when  
849 excreted by bacteria (Blankenfeldt et al., 2004). This LGT region exhibits gene  
850 expression and is flanked by conserved insect genes, providing further support for it  
851 being a legitimate LGT, though further research into this region will be necessary to

852 confirm this. The gene occurs on two different scaffolds, 22012 and 22013, which are  
853 highly similar to each other in some regions at the nucleotide level. These could  
854 represent homologous regions that differ sufficiently to assemble as different scaffolds,  
855 or alternatively a duplication in two different regions of the genome. Future work  
856 should focus on its expression patterns in different tissues (e.g. salivary glands, in interest  
857 of *PhzF*) and potential functional role in *N. tenuis*.

## 858 **Beyond the Genome: Cytogenetics**

859 We determined the karyotype of *N. tenuis* to be  $2n=32$  ( $30+XY$  in males) chromosomes  
860 which is the second most common chromosome number in the family Miridae  
861 (Kuznetsova et al., 2011). In addition, we have shown that *N. tenuis* has an  $XX/XY$  sex  
862 chromosome constitution, with the sex chromosomes being the largest elements in the  
863 karyotype. This is different from the closely related *Macrolophus costalis* (Fieber)  
864 (Hemiptera: Miridae) ( $2n=24+X1X2Y$ ), and *M. pygmaeus* (Rambur) ( $2n=26+XY$ ) where  
865 two pairs of autosomes are larger than the sex chromosomes, yet similar to *M.*  
866 *melanotoma* (Costa) which only differs from *N. tenuis* in the number of autosomes,  
867  $2n=32+XY$  (Jauset et al., 2015). As we sequenced a single female, sequence  
868 information of the Y chromosome is missing from our genome assembly. While  
869 analyzing the *N. tenuis* karyotype we discovered the sporadic presence of B  
870 chromosomes in the KBS population. B chromosomes are supernumerary  
871 chromosomes that are dispensable to the organism, and are often present in only a  
872 subset of individuals from a population (Banaei-Moghaddam et al., 2013).  
873 Supernumerary chromosomes are common in Heteroptera, yet only a few species of  
874 Miridae have been identified to carry supernumerary chromosomes (Grozeva et al.,  
875 2011). Presence of B chromosomes in high numbers within an individual is often found  
876 to be detrimental, though in lower numbers they are often considered neutral or, in

877 some cases, beneficial (Camacho et al., 2000; Jones and Rees, 1982). The abundance  
878 of B chromosomes in *N. tenuis* biological control populations is currently unknown but  
879 determining their potential effects on fitness-relevant traits might reveal beneficial  
880 information for the optimization of mass-reared populations.

881 The hemizygous sex chromosomes of most organisms have a high content of repetitive  
882 DNA, consisting of multiple different repetitive sequences that are less frequent found  
883 on autosomes (Charlesworth and Charlesworth, 2000; Traut et al., 1999). Therefore, the  
884 use of cytogenetic techniques, such as CGH and GISH, in the identification of  
885 hemizygous sex chromosomes is a powerful tool and is well established in different  
886 groups of organisms, e.g. Lepidoptera (Carabajal Paladino et al., 2019; Dalíková et al.,  
887 2017a; Zrzavá et al., 2018), Orthoptera (Jetybayev et al., 2017), fish (Sember et al.,  
888 2018), and frogs (Gatto et al., 2018). However, to our knowledge, this is the first time  
889 these techniques have been used in the family Miridae. The X and Y chromosome of  
890 *N. tenuis* are similar in size, with the X chromosome being slightly bigger, and are  
891 difficult to distinguish from each other based solely on their appearance without  
892 special probing. Our CGH and GISH results showed relatively weak hybridization signals  
893 of genomic probes on the sex chromosomes compared to other species indicating  
894 little differentiation of sequence content between the X and Y chromosomes, and/or  
895 between the sex chromosomes and the autosomes. Though the hybridization signals  
896 are relatively weak, not only the Y chromosome but also the homogametic sex  
897 chromosome, the X chromosome, is distinguishable in the CGH results, which shows X-  
898 enriched or X-specific repetitive DNA, similar to what was found on the Z chromosome  
899 in *Abraxas* spp. (Zrzavá et al., 2018). Mapping of the most abundant repeat in the  
900 genome revealed that one such X- enriched repeats is Nt\_rep1, confirming the  
901 outcomes of our CGH results.

902 The low copy numbers of 18S rDNA identified in the assembled genome were  
903 surprising. The NOR is usually composed of tens to hundreds of copies, and is therefore  
904 used in heteropteran cytogenetic studies due to its easy visualization (Kuznetsova et  
905 al., 2011). Analysis of the raw data estimates 98 copies of 18S rDNA are present in the  
906 genome, yet the majority of these copies are missing from the final assembly. The FISH  
907 results show that 18S rDNA is present as a single cluster in the genome, indicating that  
908 there is a limit to the genome assembler Supernova, and 10X Genomics by extension,  
909 and its ability to assemble highly repetitive regions of the genome.

910 Similarly, the FISH results of Nt\_rep1 and the analysis of the copy numbers and  
911 distribution of the repeat in the genome assembly do not corroborate. Though many  
912 copies of the repeat are present in the assembled genome, most scaffolds contain  
913 one or few copies of the repeat. The FISH results, however, show multiple clusters  
914 scattered across most chromosomes each containing high copy numbers, revealing  
915 a lack of scaffolds containing high copy numbers of Nt\_rep1 in the assembled  
916 genome. Therefore, analyses on repetitive DNA content are currently more reliable  
917 using the short sequence reads rather than the assembled genome as it  
918 underestimates repeat content. Long read sequencing methods would be able to  
919 overcome such problems with repetitive DNA, not only in *N. tenuis* but in any species,  
920 and would be better suited to analyse repetitive regions of genomes. As mentioned  
921 before, a hybrid assembly strategy combining our 10X sequencing data with long  
922 reads, obtained by e.g. Oxford Nanopore or PacBio sequencing, would presumably  
923 improve the assembly, though in this aspect for particular segments of the genome  
924 that are high in repetitive DNA. This should be kept in mind for other 10X  
925 Genomics/Supernova-derived genomes: the true number of repeats may be  
926 underestimated.

927 Screening of the genome and Southern blot assay suggests the absence of the  
928 ancestral insect telomere motif, (TTAGG)<sub>n</sub>, in *N. tenuis*, as the case in other species  
929 from the family Miridae (Grozeva et al., 2019; Kuznetsova et al., 2011). The telomeric  
930 motif was present in our Tandem Repeat Finder results, but in much lower numbers  
931 than expected for telomeric sequences. Additional attempts of identifying the  
932 telomeric repeat motif did not resolve this question. Three additional repeats we  
933 identified in the *N. tenuis* genome were tested via FISH, i.e. (TATGG)<sub>n</sub>, (TTGGG)<sub>n</sub>, and  
934 (TCAGG)<sub>n</sub>, but did not localise near the ends of the chromosomes. Notably though,  
935 mapping the most abundant repeat in the genome, Nt\_rep1, did reveal accumulation  
936 in the sub-telomeric regions of chromosomes (Figure 6). Therefore, our approach to  
937 identify potential telomere motifs, though presently unsuccessful, would presumably  
938 be effective if more repeats would be screened. In addition, a similar approach was  
939 used by Pita et al. (2016) in *T. infestans*, where the insect telomere motif, (TTAGG)<sub>n</sub>, was  
940 successfully identified from the raw sequencing data (Pita et al., 2016). It must be  
941 noted, however, that the telomeres of *N. tenuis* might consist of different types of  
942 repeats other than short tandem repeats (as found in, for example, *Drosophila*;  
943 Traverse & Pardue, 1988) which would not be identified using Tandem Repeat Finder  
944 (Traverse and Pardue, 1988). Therefore, the identity (or even presence) of the  
945 telomeric repeat in *N. tenuis*, and by extension Miridae, remains unknown.

#### 946 **Beyond the Genome: Population Genomics**

947 Pooled sequence data of ten females from the KBS population were compared  
948 against the genome and provide interesting population-level effects. The overall  
949 negative *Tajima's D* would seem to indicate an abundance of rare alleles and is  
950 possible evidence of selective sweeps or population expansion, as seen in some  
951 populations of *Drosophila serrata* (Malloch) (Reddiex et al., 2018), however, this

952 generally results in more negative values (near -1 or -2). While overall negative, the  
953 absolute value of  $D$  in our results is small in comparison (total range from -0.89 to 0.56).  
954 To best assess the state of the commercial population, monitoring the genetic  
955 variation over time would indicate if the population is undergoing an expansion after  
956 a bottleneck ( $D < 0$ ) or contracting ( $D > 0$ ), whereas when  $D = 0$ , we assume no  
957 selection. We can then assume that there is no selection currently at play in the  
958 commercial population. The few studies that have looked at genetic diversity within  
959 biological control populations have primarily been reduced representation analyses,  
960 such by genotyping with microsatellites (Paspati et al., 2019). Here, a pool-seq  
961 approach offers a genome-wide look at the population and can give indications of  
962 the genetic diversity of the population; this could be a useful tool for monitoring  
963 population levels efficiently and determining which regions of the genome are under  
964 selection in a biological control context.

965 Both genetic diversity values calculated here can also be used in population  
966 comparisons between the biological control population and wild populations. For  
967 instance, Xun *et al.* used mitochondrial and nuclear barcoding regions to haplotype  
968 516 individuals across 37 populations into two regional groups, southwest China (SWC)  
969 and other regions in China (OC) (Xun et al., 2016).  $\pi$  was 0.0048 (SWC) and 0.904 (OC),  
970 while  $D$  was -0.112 (SWC) and -1.998 (OC). It was concluded that the SWC population  
971 was stable while, similar to the KBS population here, the OC population was  
972 undergoing sudden population explosion. Pooled sequencing could be a useful tool  
973 for comparing wild Mediterranean populations to the commercial population to  
974 determine disparities in genetic variation as well as to understand the dynamics of the  
975 wild populations.



976 There is a concern in using PoPoolation in this context: are ten individual females  
977 sufficient for determining population variation? Here we used existing population  
978 sequence data to better utilize resources, reduced to an appropriate coverage with  
979 masked indel regions. This enabled us to show population-level impacts at the very  
980 least, which can then pave the way for further studies, with better constructed  
981 sampling methods and sample sizes; the lack of perfect data should not preclude  
982 preliminary studies from being pursued.

983

## 984 **Conclusion**

985 Reported here is the genome for *N. tenuis*, a mirid that is both used throughout the  
986 Mediterranean Basin as a biological control agent and reported as a greenhouse pest  
987 in other European countries. The assembled genome is 355 Mbp in length, composed  
988 of 42,830 scaffolds with an N50 of 27,074 bp. The goal of this project was to not only  
989 provide a genome, but also to highlight possible avenues of research now available  
990 with *N. tenuis*. A protein analysis has provided interesting prospects for mirid-specific  
991 proteins, while examples of potential LGT call for further inquiry. Putative symbionts  
992 were identified while filtering out contamination, creating a precursor for future  
993 metagenomic analysis. The cytogenetic analyses of *N. tenuis* here shed some light on  
994 Mirid cytogenetics, such as the karyotype and sex determination system, but also  
995 solicits more questions. As for the commercial population, now that there is a baseline  
996 level of genetic variation documented through our pooled sequencing, what remains  
997 to be seen is how it compares to other populations, such as other commercial  
998 populations, wild, or invasive populations. To this end, future exploration on these  
999 themes, among others, are now greatly facilitated with our release of this genome.

1000

1001 **Acknowledgements**

1002  
1003 We would like to express special thanks to Milena Chinchilla Ramírez for her advice and *N.*  
1004 *tenuis* expertise. Thanks to Markus Knapp (Koppert BV), Javier Calvo (Koppert BS) and Gerben  
1005 Messelink (WUR) for providing *N. tenuis* specimens. Thanks to both José van de Belt, Frank  
1006 Becker (WUR), and Carolina Gallego (IVIA) for their technical assistance in DNA and RNA  
1007 extraction. We acknowledge the assistance of Magda Zrzavá, Anna Voleníková, Martina  
1008 Flegrová, and Diogo Cabral-de-Mello (Institute of Entomology BC CAS) for their cytogenetic  
1009 expertise and assistance. JHW acknowledges Sammy Cheng for assistance with the LGT  
1010 pipeline. KBF acknowledges Jetske de Boer for her assistance with the flow cytometry and Joost  
1011 van den Heuvel for his assistance with PoPoolation. Annotation was performed by  
1012 GenomeScan B. V (NL). Access to computing and storage facilities owned by parties and  
1013 projects contributing to the National Grid Infrastructure MetaCentrum (CZ), provided under the  
1014 programme "Projects of Large Research, Development, and Innovations Infrastructures"  
1015 (CESNET LM2015042), is greatly appreciated. This project was funded by the European Union's  
1016 Horizon 2020 research and innovation program under the Marie Skłodowska-Curie grant  
1017 agreement no. 641456. JHW acknowledges the US-NSF-IOS-1456233 and Nathaniel & Helen  
1018 Wisch Chair for funding support. The research leading to these results was partially funded by  
1019 the Spanish Ministry of Economy and Competitiveness MINECO (RTA2017-00073-00-00).  
1020 Cytogenetic experiments were financed by grants 17-13713S (SV and FM) and 17-17211S (MD  
1021 and IP) of the Czech Science Foundation.

1022  
1023 **Conflict of Interest**

1024  
1025 The authors declare no conflict of interest.

1026 **Tables**

Table 1. Assembly statistics for both versions of the *Nesidiocoris tenuis* assembly, pre- and post-decontamination

<b>Assembly Version</b>	<b>Size (bp)</b>	<b>No. of Scaffolds</b>	<b>N50 (bp)</b>	<b>Largest scaffold (bp)</b>	<b>No. of N's per 100 kbp (% of genome)</b>	<b>BUSCO score, Complete% (Single%, Duplicate%)</b>
1.0	387,724,797	44,273	27,195	1,392,896	5912.60 (5.91)	81.3 (60.6, 20.7)
1.5 (final assembly)	355,120,802	36,513	28,732	1,392,896	6292.10 (6.29)	87.5 (65.6, 21.9)

1027

1028

1029

Table 2. Genera of potential symbionts or contaminants as determined by decontamination pipeline based on known contaminants and symbionts against *Nesidiocoris tenuis* assembly v1.0. Identification is according to largest hit percentage, multiple bacterial sections possible in each scaffold. Affected scaffolds were removed leading to assembly v1.5 and are available in S1.2. Full list of hits available in supplementary material S1.4.

<b>Bacteria Genus</b>	<b>Total amount in genome (bp)</b>	<b>Number of scaffolds affected</b>
<i>Pantoea</i>	1,379,962	307
<i>Erwinia</i>	1,342,298	456
<i>Citrobacter</i>	349,562	40
<i>Rickettsia</i>	204,934	50
<i>Sodalis</i>	189,471	17
<i>Cronobacter</i>	180,427	33
<i>Wolbachia</i>	137,109	47
<i>Enterobacter</i>	110,104	70
<i>Serratia</i>	100,197	49
<i>Klebsiella</i>	44,248	30
<i>Pseudoalteromonas</i>	39,779	118
<i>Pectobacterium</i>	28,130	16
<i>Shigella</i>	25,221	24
<i>Yersinia</i>	24,052	21
<i>Dickeya</i>	23,312	18
<i>Salmonella</i>	19,243	15
<i>Photorhabdus</i>	18,283	14
<i>Escherichia</i>	16,301	11
<i>Rahnella</i>	12,479	10
<i>Burkholderia</i>	7,825	17
<i>Xenorhabdus</i>	6,406	8
<i>Arsenophonus</i>	5,070	8
<i>Pseudomonas</i>	3,427	2
<i>Vulcanisaeta</i>	3,420	3
<i>Ralstonia</i>	3,370	7
<i>Paenibacillus</i>	3,070	8
Others (105)	56,287	201
<b>Total</b>	<b>4,333,987</b>	<b>1,600</b>

1030

1031

1032

1033

Table 3. Output of OrthoVenn2 ortholog cluster analysis of *Nesidiocoris tenuis*, *Cimex lectularis*, *Halyomorpha halys*, and *Acyrtosiphon pisum*.

<b>Species</b>	<b>Proteins</b>	<b>Clusters</b>	<b>Singletons</b>	<b>Source of complete gene set</b>
<i>N. tenuis</i>	24,668	8,174	9,136	This work
<i>C. lectularis</i>	12,699	7,989	7,989	Poelchau <i>et al.</i> , 2015; Benoit <i>et al.</i> , 2016; Thomas <i>et al.</i> , 2020
<i>H. halys</i>	25,026	9,584	2,170	Lee <i>et al.</i> , 2009; Poelchau <i>et al.</i> , 2015
<i>A. pisum</i>	36,195	8,765	7,298	Legeai <i>et al.</i> , 2010; L. Xu <i>et al.</i> , 2019

1034

1035

1036

1037

Table 4. PoPoolation analysis on commercial Koppert Biological Systems *Nesidiocoris tenuis* population (n=10 females), with 10 largest scaffolds according to size. Coverage is  $\geq 0.60$  and indel regions are masked.

<b>Scaffold</b>	<b>Size (bp)</b>	<b>Windows (10 kbp)</b>	<b>Number of sufficiently covered windows</b>	<b>Average coverage of sufficiently covered windows</b>	<b>Average <math>\pi</math> across scaffold</b>	<b>Average Tajima's D across scaffold</b>
35384	1,392,896	140	70	0.68	0.010682	-0.02091
5	613,435	62	48	0.76	0.004932	-0.03869
31795	577,751	58	38	0.74	0.001753	-0.04655
33087	539,928	54	17	0.66	0.012571	-0.00458
12956	519,254	52	39	0.80	0.002541	-0.21012
23581	513,368	52	28	0.70	0.008417	-0.04658
11795	508,155	51	14	0.68	0.011553	-0.01147
20742	504,856	51	39	0.78	0.004165	-0.01963
7669	488,533	49	24	0.67	0.007856	-0.0551
28424	477,580	48	5	0.66	0.009563	-0.00998
<b>Total</b>	<b>256,487,768</b>	<b>28,833</b>	<b>5,913</b>	<b>0.70</b>	<b>0.0080</b>	<b>-0.0355</b>

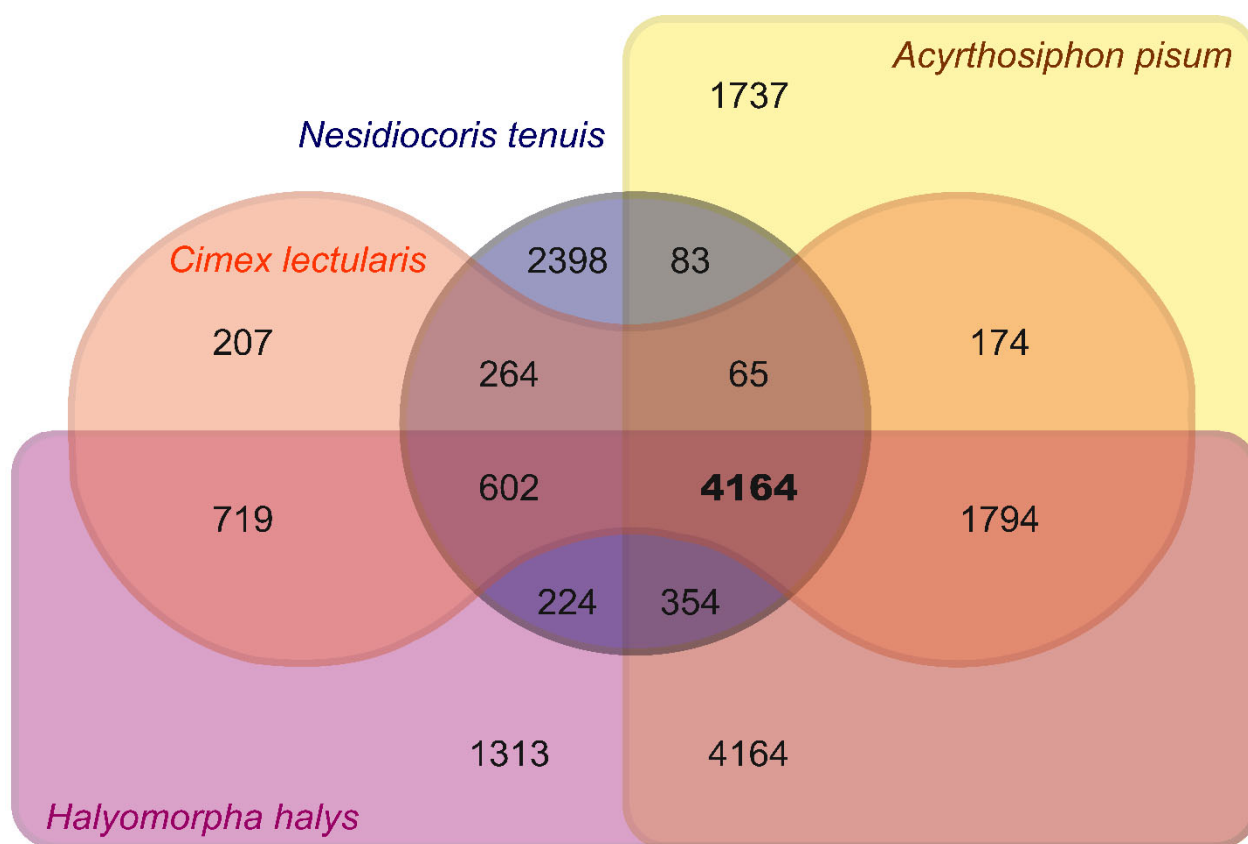
1038

1039

1040

1041 **Figures**

1042

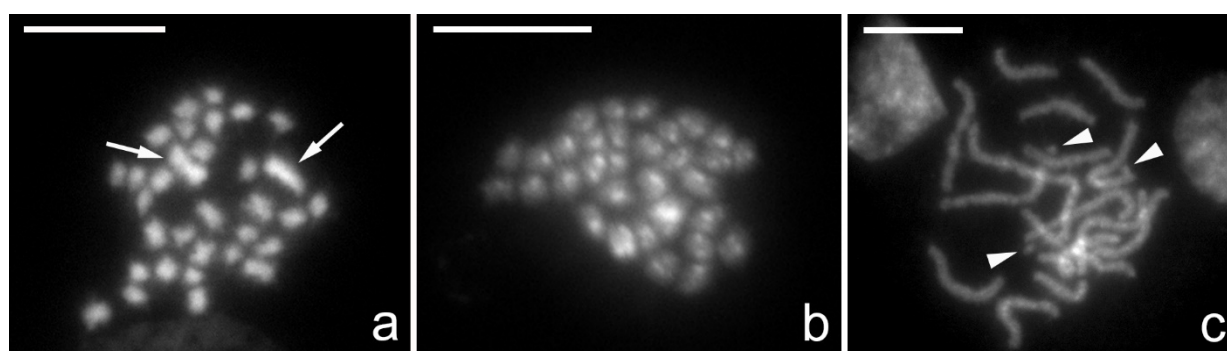


1043

1044 **Figure 1.** Ortholog cluster analysis of *Nesidiocoris tenuis* with three other hemipterans (*Cimex*  
 1045 *lectularis*, *Halyomorpha halys*, and *Acyrthosiphon pisum*). Numbers indicate the number of  
 1046 ortholog clusters in each grouping, with the clusters shared by all four species in bold.

1047

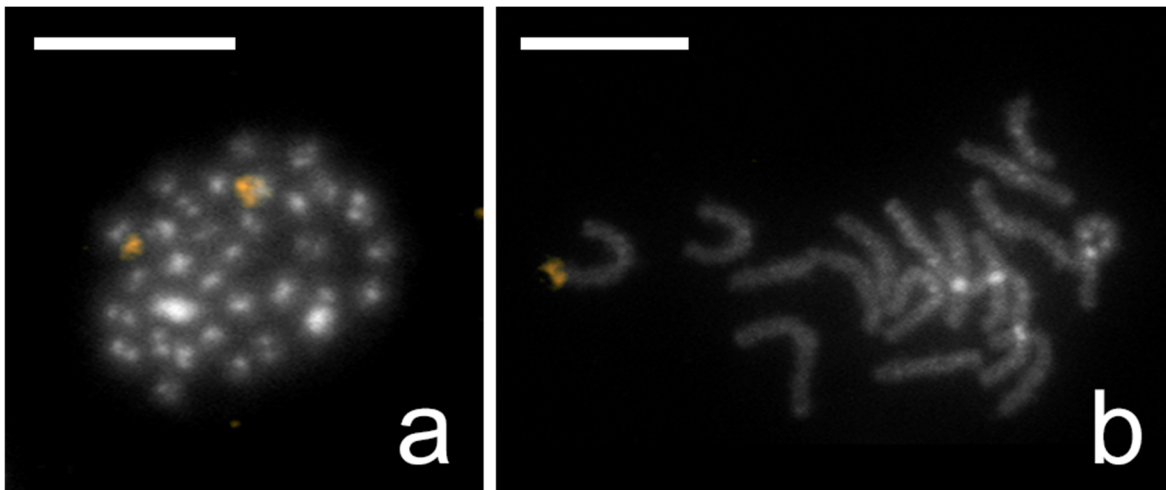
1048



1049

1050 **Figure 2:** Cytogenetic analysis of *Nesidiocoris tenuis* karyotype. Chromosomes were  
 1051 counterstained by DAPI (grey). **(a)** Female mitotic metaphase consisting of 32 chromosomes  
 1052 ( $2n=32$ ) with two large chromosomes indicated (arrows). **(b)** Male mitotic metaphase  
 1053 consisting of 32 chromosomes ( $2n=32$ ). **(c)** Female pachytene nucleus with B chromosomes  
 1054 (arrowheads). Scale bar = 10  $\mu\text{m}$ .

1055

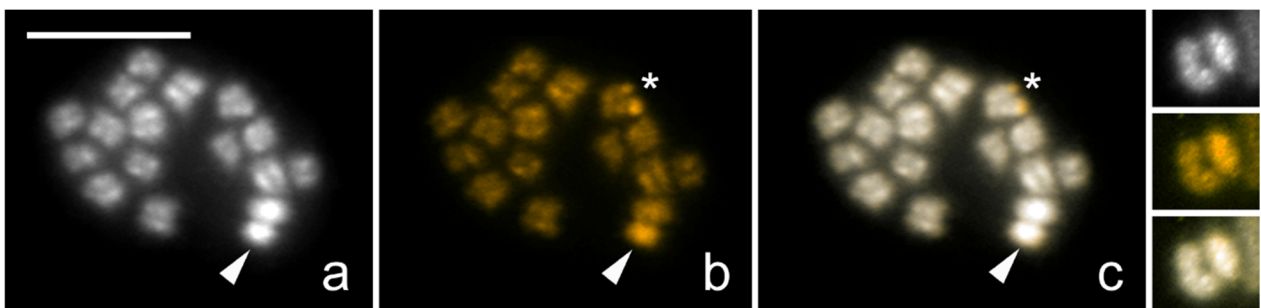


1056

1057 **Figure 3:** Results of *Nesidiocoris tenuis* fluorescence *in situ* hybridization with 18S rDNA probe  
1058 labelled by biotin and visualised by detection with Cy3-conjugated streptavidin (gold).  
1059 Chromosomes were counterstained by DAPI (grey). **(a)** Male mitotic metaphase; probe  
1060 identified a cluster of 18S rDNA on two homologues chromosomes. **(b)** Female pachytene  
1061 complement with one terminal cluster of 18S rDNA genes on a bivalent. Scale bar = 10  $\mu\text{m}$ .

1062

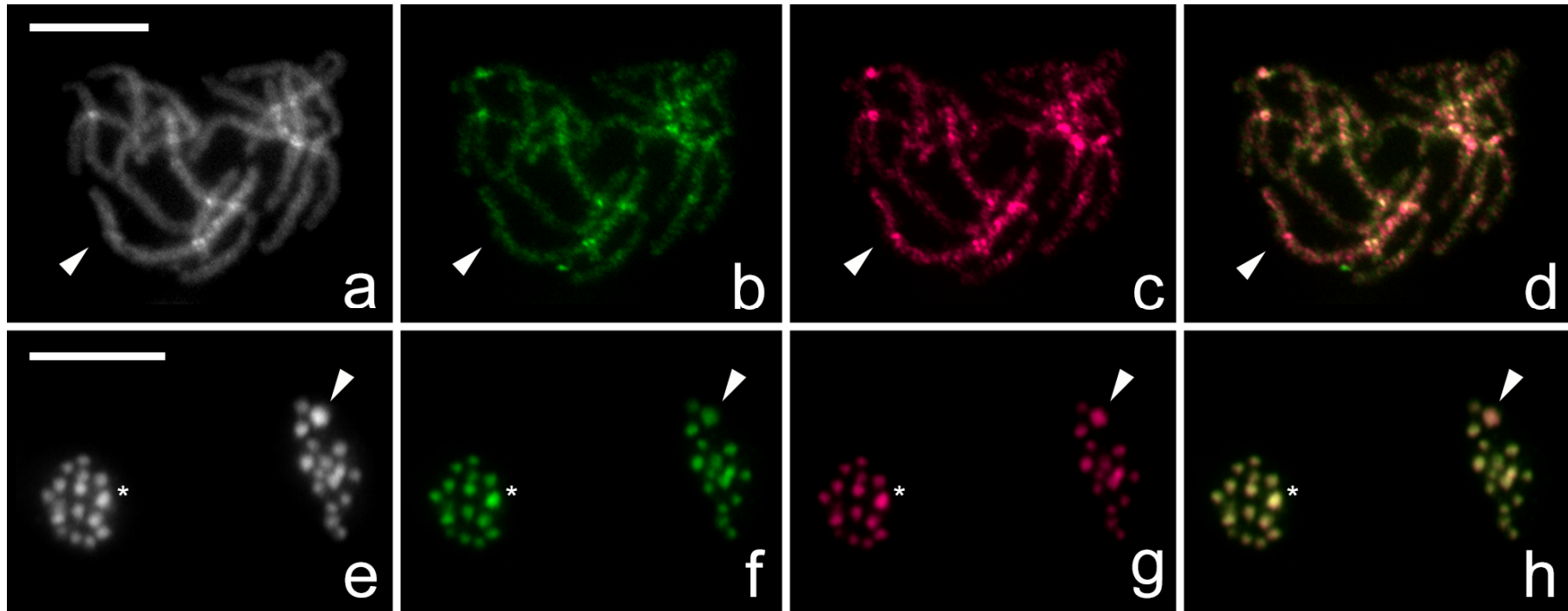
1063



1064

1065 **Figure 4:** Genomic *in situ* hybridization (GISH) on male chromosomal preparation of  
1066 *Nesidiocoris tenuis*. Panel **(a)** shows DAPI counterstaining (grey), panel **(b)** hybridisation signals  
1067 of the male derived genomic probe labelled by Cy3 (gold) together with competitor  
1068 generated from unlabelled female genomic DNA, and panel **(c)** a merged image. **(a, b, c,**  
1069 **detail)** Meiotic metaphase I, male derived probe highlighted the Y chromosome (arrowhead)  
1070 more **(b, c)** compared to autosomes and the X chromosome. Note highlighted terminal  
1071 regions of one of the bivalents caused by presence of major rDNA genes (asterisk). **(detail)**  
1072 Detail picture of XY bivalent; Y chromosome labelled by male derived probe. Note that the Y  
1073 chromosome is smaller in size and showing more heterochromatin compared to the X  
1074 chromosome (and autosomes). Scale bar = 10  $\mu\text{m}$ .

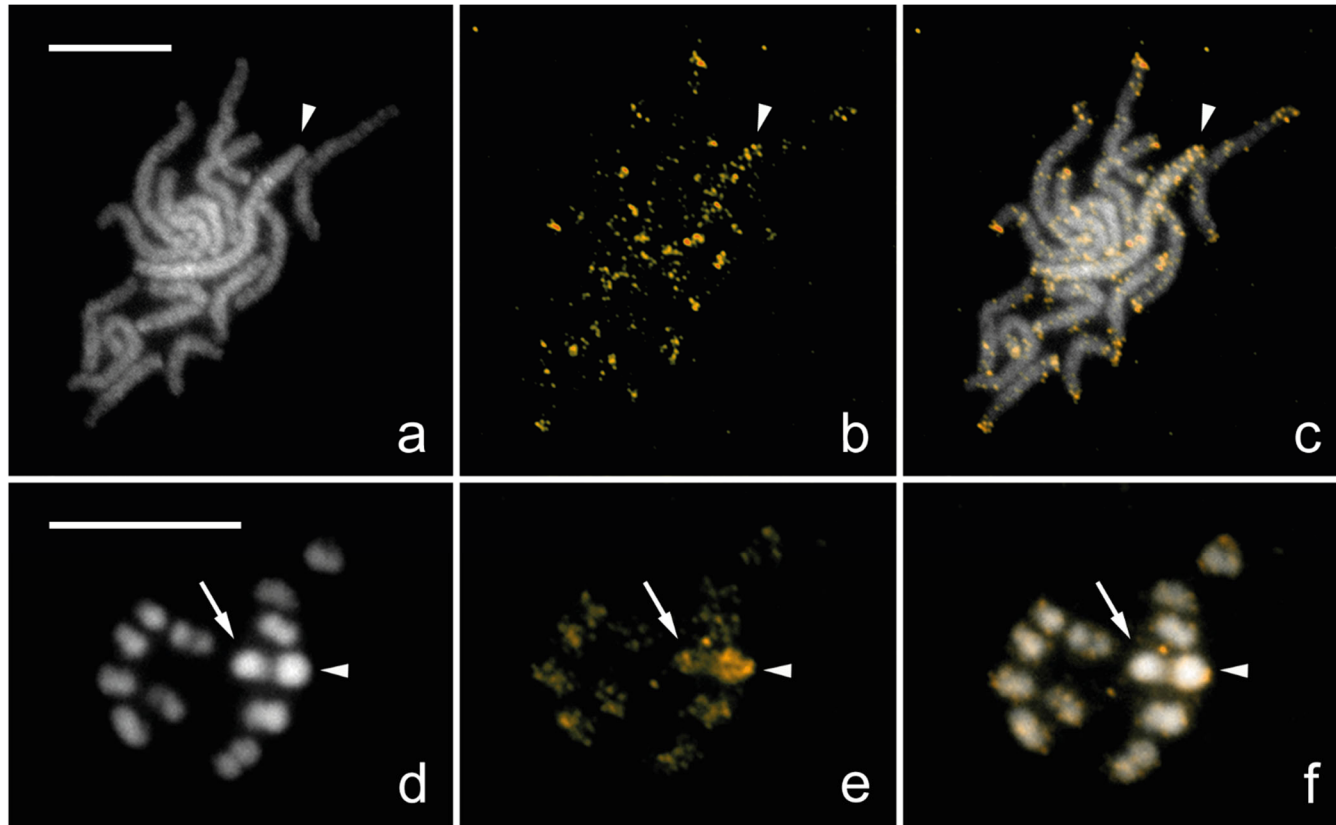




1076

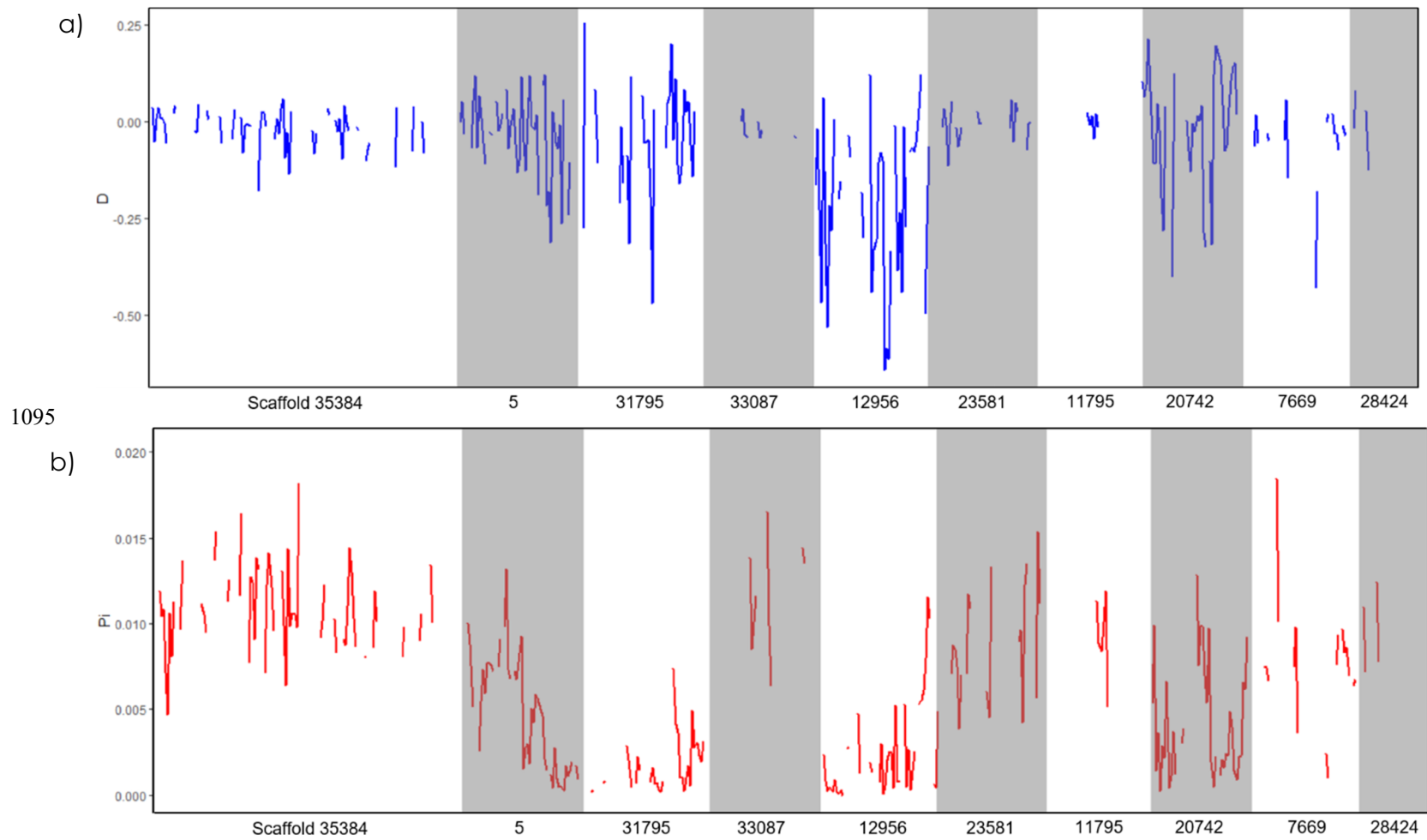
1077 **Figure 5:** Comparative genomic hybridization (CGH) on female (**a, b, c, d**) and male meiotic metaphase II (**e, f, g, h**) chromosomes of *Nesidiocoris*  
 1078 *tenuis*. Panels (**a, e**) show chromosomes counterstained by DAPI (grey), panels (**b, f**) hybridization signals of the male derived genomic probe  
 1079 labelled by fluorescein (blue), panels (**c, g**) hybridization signals of the female derived genomic probe labelled by Cy3 (gold), and panels (**d, h**)  
 1080 merged images. (**c, d**) Note that the X chromosome bivalent (arrowhead) in female pachytene complement was highlighted more by female  
 1081 probe compared to the autosomal bivalents; (**b, d**) male probe labelled all chromosomes equally. (**h**) Two sister nuclei in meiotic metaphase II  
 1082 showed equal hybridization patterns of both probes on autosomes; in one of the forming nuclei, the X chromosome (arrowhead) was highlighted  
 1083 by female derived genomic probe (**g, h**) and in the second nucleus the Y chromosome (asterisk) was strongly highlighted by male derived genomic  
 1084 probe compared to autosomes (**f, h**) and less highlighted by female derived probe (**g, h**). (**e**) Note that the sex chromosomes are the biggest and  
 1085 most heterochromatic elements in the nucleus. Scale bar = 10  $\mu$ m. **Note: Alternate colouration available in supplementary material (S1.6.3).**

1086



1088

1089 **Figure 6:** Fluorescence *in situ* hybridization with Nt\_rep1 probe labelled by biotin (gold) on female (**a, b, c**) and male (**d, e, f**) chromosomes of  
 1090 *Nesidiocoris tenuis* counterstained with DAPI (grey). (**a, b, c**) Female pachytene chromosomes; Nt\_rep1 probe highlighted the pair of X  
 1091 chromosomes (arrowhead). (**b, c**) Note that terminal regions of all bivalents were also labelled by probe, probably due to presence of this  
 1092 sequence in sub-telomeric regions. (**d, e, f**) Incomplete male nucleus in meiotic metaphase I; probe highlighted the X chromosome (arrowhead)  
 1093 more compared to autosomes and Y chromosome (arrow). (**b, c, e, f**) Strong hybridization signals on X chromosomes in both sexes were caused  
 1094 by enrichment of Nt\_rep1 sequence on the X chromosomes. Scale bar = 10



1095

1096

1097

1098 **Figure 7.** Genetic diversity of Koppert Biological Systems commercial *Nesidiocoris tenuis* population according to Tajima's  $D$  **(a)**, and nucleotide  
 1099 diversity,  $\pi$  **(b)**. Scaffolds are ordered according to size, which each name beneath on the x-axis.

## 1100 References

- 1101 Abreo, E., Altier, N., 2019. Pangenome of *Serratia marcescens* strains from nosocomial and  
1102 environmental origins reveals different populations and the links between them. *Sci. Rep.*  
1103 9, 1–8. <https://doi.org/10.1038/s41598-018-37118-0>
- 1104 Acland, A., Agarwala, R., Barrett, T., Beck, J., Benson, D.A., Bollin, C., Bolton, E., Bryant, S.H.,  
1105 Canese, K., Church, D.M., Clark, K., Dicuccio, M., Dondoshansky, I., Federhen, S., Feolo,  
1106 M., Geer, L.Y., Gorelenkov, V., Hoepfner, M., Johnson, M., Kelly, C., Khotomlianski, V.,  
1107 Kimchi, A., Kimelman, M., Kitts, P., Krasnov, S., Kuznetsov, A., Landsman, D., Lipman, D.J.,  
1108 Lu, Z., Madden, T.L., Madej, T., Maglott, D.R., Marchler-Bauer, A., Karsch-Mizrachi, I.,  
1109 Murphy, T., Ostell, J., O'Sullivan, C., Panchenko, A., Phan, L., Pruitt, D.P.K.D., Rubinstein,  
1110 W., Sayers, E.W., Schneider, V., Schuler, G.D., Sequeira, E., Sherry, S.T., Shumway, M.,  
1111 Sirotkin, K., Siyan, K., Slotta, D., Soboleva, A., Soussov, V., Starchenko, G., Tatusova, T.A.,  
1112 Trawick, B.W., Vakatos, D., Wang, Y., Ward, M., John Wilbur, W., Yaschenko, E., Zbicz, K.,  
1113 2014. Database resources of the National Center for Biotechnology Information. *Nucleic*  
1114 *Acids Res.* 42, 8–13. <https://doi.org/10.1093/nar/gkt1146>
- 1115 Andrews, S., Krueger, F., Seccombe, P., Biggins, F., Wingett, S., 2015. FastQC. A quality  
1116 control tool for high throughput sequence data. Babraham Bioinformatics. Babraham  
1117 Inst.
- 1118 Arnó, J., Castañé, C., Riudavets, J., Gabarra, R., 2010. Risk of damage to tomato crops by the  
1119 generalist zoophytophagous predator *Nesidiocoris tenuis* (Reuter) (Hemiptera: Miridae).  
1120 *Bull. Entomol. Res.* 100, 105–115. <https://doi.org/10.1017/S0007485309006841>
- 1121 Banaei-Moghaddam, A.M., Meier, K., Karimi-Ashtiyani, R., Houben, A., 2013. Formation and  
1122 expression of pseudogenes on the B chromosome of rye. *Plant Cell* 25, 2536–2544.  
1123 <https://doi.org/10.1105/tpc.113.111856>
- 1124 Benoit, J.B., Adelman, Z.N., Reinhardt, K., Dolan, A., Poelchau, M., Jennings, E.C., Szuter, E.M.,  
1125 Hagan, R.W., Gujar, H., Shukla, J.N., Zhu, F., Mohan, M., Nelson, D.R., Rosendale, A.J.,  
1126 Derst, C., Resnik, V., Wernig, S., Menegazzi, P., Wegener, C., Peschel, N., Hendershot,  
1127 J.M., Blenau, W., Predel, R., Johnston, P.R., Ioannidis, P., Waterhouse, R.M., Nauen, R.,  
1128 Schorn, C., Ott, M.-C., Maiwald, F., Johnston, J.S., Gondhalekar, A.D., Scharf, M.E.,  
1129 Peterson, B.F., Raje, K.R., Hottel, B.A., Armisen, D., Crumière, A.J.J., Refki, P.N., Santos,  
1130 M.E., Sghaier, E., Viala, S., Khila, A., Ahn, S.-J., Childers, C., Lee, C.-Y., Lin, H., Hughes,  
1131 D.S.T., Duncan, E.J., Murali, S.C., Qu, J., Dugan, S., Lee, S.L., Chao, H., Dinh, H., Han, Y.,  
1132 Doddapaneni, H., Worley, K.C., Muzny, D.M., Wheeler, D., Panfilio, K.A., Vargas Jentzsch,  
1133 I.M., Vargo, E.L., Booth, W., Friedrich, M., Weirauch, M.T., Anderson, M.A.E., Jones, J.W.,  
1134 Mittapalli, O., Zhao, C., Zhou, J.-J., Evans, J.D., Attardo, G.M., Robertson, H.M., Zdobnov,  
1135 E.M., Ribeiro, J.M.C., Gibbs, R.A., Werren, J.H., Palli, S.R., Schal, C., Richards, S., 2016.  
1136 Unique features of a global human ectoparasite identified through sequencing of the  
1137 bed bug genome. *Nat. Commun.* 7, 10165. <https://doi.org/10.1038/ncomms10165>
- 1138 Benson, G., 1999. Tandem repeats finder: A program to analyze DNA sequences. *Nucleic*  
1139 *Acids Res.* 27, 573–580. <https://doi.org/10.1093/nar/27.2.573>
- 1140 Biondi, A., Zappalà, L., Di Mauro, A., Tropea Garzia, G., Russo, A., Desneux, N., Siscaro, G.,  
1141 2015. Can alternative host plant and prey affect phytophagy and biological control by  
1142 the zoophytophagous mirid *Nesidiocoris tenuis*? *BioControl* 79–90.  
1143 <https://doi.org/10.1007/s10526-015-9700-5>
- 1144 Blankenfeldt, W., Kuzin, A.P., Skarina, T., Korniyenko, Y., Tong, L., Bayer, P., Janning, P.,  
1145 Thomashow, L.S., Mavrodi, D. V., 2004. Structure and function of the phenazine  
1146 biosynthetic protein PhzF from *Pseudomonas fluorescens*. *Proc. Natl. Acad. Sci. U. S. A.*  
1147 101, 16431–16436. <https://doi.org/10.1073/pnas.0407371101>

- 1148 Bolger, A.M., Lohse, M., Usadel, B., 2014. Trimmomatic: A flexible trimmer for Illumina sequence  
1149 data. *Bioinformatics* 30, 2114–2120. <https://doi.org/10.1093/bioinformatics/btu170>
- 1150 Bouagga, S., Urbaneja, A., Depalo, L., Rubio, L., Pérez-Hedo, M., 2019. Zoophytophagous  
1151 predator-induced defences restrict accumulation of the tomato spotted wilt virus. *Pest*  
1152 *Manag. Sci.* ps.5547. <https://doi.org/10.1002/ps.5547>
- 1153 Boutet, E., Lieberherr, D., Tognolli, M., Schneider, M., Bairoch, A., 2008. UniProtKB/Swiss-Prot:  
1154 The manually annotated section of the UniProt KnowledgeBase. *Methods Mol. Biol.* 406,  
1155 89–112. <https://doi.org/10.1007/978-1-59745-535-0>
- 1156 Boyd, B.M., Allen, J.M., Koga, R., Fukatsu, T., Sweet, A.D., Johnson, K.P., Reed, D.L., 2016. Two  
1157 Bacterial Genera, *Sodalis* and *Rickettsia*, Associated with the Seal Louse  
1158 *Proechinophthirus fluctus* (Phthiraptera: Anoplura). *Appl. Environ. Microbiol.* 82, 3185–  
1159 3197. <https://doi.org/10.1128/AEM.00282-16>
- 1160 Bruton, B.D., Mitchell, F., Fletcher, J., Pair, S.D., Wayadande, A., Melcher, U., Brady, J., Bextine,  
1161 B., Popham, T.W., 2007. *Serratia marcescens*, a Phloem-Colonizing, Squash Bug -  
1162 Transmitted Bacterium: Causal Agent of Cucurbit Yellow Vine Disease. *Plant Dis.* 87, 937–  
1163 944. <https://doi.org/10.1094/pdis.2003.87.8.937>
- 1164 Calvo, J.F., Bolckmans, K., Stansly, P. a., Urbaneja, A., 2009. Predation by *Nesidiocoris tenuis* on  
1165 *Bemisia tabaci* and injury to tomato. *BioControl* 54, 237–246.  
1166 <https://doi.org/10.1007/s10526-008-9164-y>
- 1167 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., Madden, T.L.,  
1168 2009. BLAST+: Architecture and applications. *BMC Bioinformatics* 10, 421.  
1169 <https://doi.org/10.1186/1471-2105-10-421>
- 1170 Camacho, J.P.M., Sharbel, T.F., Beukeboom, L.W., 2000. B-chromosome evolution. *Philos.*  
1171 *Trans. R. Soc. London. Ser. B Biol. Sci.* 355, 163–178. <https://doi.org/10.1098/rstb.2000.0556>
- 1172 Carabajal Paladino, L.Z., Provazníková, I., Berger, M., Bass, C., Aratchige, N.S., López, S.N.,  
1173 Marec, F., Nguyen, P., 2019. Sex Chromosome Turnover in Moths of the Diverse  
1174 Superfamily Gelechioidea. *Genome Biol. Evol.* 11, 1307–1319.  
1175 <https://doi.org/10.1093/gbe/evz075>
- 1176 Caspi-Fluger, A., Inbar, M., Steinberg, S., Friedmann, Y., Freund, M., Mozes-Daube, N., Zchori-  
1177 Fein, E., 2014. Characterization of the symbiont *Rickettsia* in the mirid bug *Nesidiocoris*  
1178 *tenuis* (Reuter) (Heteroptera: Miridae). *Bull. Entomol. Res.* 104, 681–688.  
1179 <https://doi.org/10.1017/S0007485314000492>
- 1180 Castañé, C., Arnó, J., Gabarra, R., Alomar, O., 2011. Plant damage to vegetable crops by  
1181 zoophytophagous mirid predators. *Biol. Control* 59, 22–29.  
1182 <https://doi.org/10.1016/j.biocontrol.2011.03.007>
- 1183 Charlesworth, B., Charlesworth, D., 2000. The degeneration of Y chromosomes. *Philos. Trans. R.*  
1184 *Soc. B Biol. Sci.* 355, 1563–1572. <https://doi.org/10.1098/rstb.2000.0717>
- 1185 Chin, C.-S., Peluso, P., Sedlazeck, F.J., Nattestad, M., Concepcion, G.T., Clum, A., Dunn, C.,  
1186 O'Malley, R., Figueroa-Balderas, R., Morales-Cruz, A., Cramer, G.R., Delledonne, M., Luo,  
1187 C., Ecker, J.R., Cantu, D., Rank, D.R., Schatz, M.C., 2016. Phased diploid genome  
1188 assembly with single-molecule real-time sequencing. *Nat. Methods* 13, 1050–1054.  
1189 <https://doi.org/10.1038/nmeth.4035>
- 1190 Costechareyre, D., Balmand, S., Condemine, G., Rahbé, Y., 2012. *Dickeya dadantii*, a plant  
1191 pathogenic bacterium producing cyt-like entomotoxins, causes septicemia in the pea  
1192 aphid *Acyrtosiphon pisum*. *PLoS One* 7. <https://doi.org/10.1371/journal.pone.0030702>

- 1193 Dai, X., Xun, H., Chang, J., Zhang, J., Hu, B., Li, H., Yuan, X., Cai, W., 2012. The complete  
1194 mitochondrial genome of the plant bug *Nesidiocoris tenuis* (Reuter) (Hemiptera: Miridae:  
1195 Bryocorinae: Dicyphini). *Zootaxa* 3554, 30–44. <https://doi.org/10.11646/zootaxa.3554.1.2>
- 1196 Dalíková, M., Zrzavá, M., Hladová, I., Nguyen, P., Šonský, I., Flegrová, M., Kubičková, S.,  
1197 Voleníková, A., Kawahara, A.Y., Peters, R.S., Marec, F., Sayres, M.W., 2017a. New Insights  
1198 into the Evolution of the W Chromosome in Lepidoptera. *J. Hered.* 108, 709–719.  
1199 <https://doi.org/10.1093/jhered/esx063>
- 1200 Dalíková, M., Zrzavá, M., Kubičková, S., Marec, F., 2017b. W-enriched satellite sequence in the  
1201 Indian meal moth, *Plodia interpunctella* (Lepidoptera, Pyralidae). *Chromosom. Res.* 25,  
1202 241–252. <https://doi.org/10.1007/s10577-017-9558-8>
- 1203 De Boer, J.G., Ode, P.J., Vet, L.E.M., Whitfield, J.B., Heimpel, G.E., 2007. Diploid males sire  
1204 triploid daughters and sons in the parasitoid wasp *Cotesia vestalis*. *Heredity (Edinb.)*. 99,  
1205 288–294. <https://doi.org/10.1038/sj.hdy.6800995>
- 1206 Doyle, JJ, Doyle, JL, 1990. Isolation of plant DNA from fresh tissue. *Focus (Madison)*. 12, 13–15.
- 1207 Dunning Hotopp, J.C., Clark, M.E., Oliveira, D.C.S.G., Foster, J.M., Fischer, P., Muñoz Torres,  
1208 M.C., Giebel, J.D., Kumar, N., Ishmael, N., Wang, S., Ingram, J., Nene, R. V., Shepard, J.,  
1209 Tomkins, J., Richards, S., Spiro, D.J., Ghedin, E., Slatko, B.E., Tettelin, H., Werren, J.H., 2007.  
1210 Widespread lateral gene transfer from intracellular bacteria to multicellular eukaryotes.  
1211 *Science (80- )*. 317, 1753–1756. <https://doi.org/10.1126/science.1142490>
- 1212 Ekblom, R., Wolf, J.B.W., 2014. A field guide to whole-genome sequencing, assembly and  
1213 annotation. *Evol. Appl.* 1026–1042. <https://doi.org/10.1111/eva.12178>
- 1214 Ellegren, H., 2014. Genome sequencing and population genomics in non-model organisms.  
1215 *Trends Ecol. Evol.* 29, 51–63. <https://doi.org/10.1016/j.tree.2013.09.008>
- 1216 Ferguson, K., 2020. *Nesidiocoris tenuis* PRJEB35378 linked-read genome annotation..  
1217 <https://doi.org/10.6084/m9.figshare.12073893.v1>
- 1218 Frydrychová, R., Grossmann, P., Trubac, P., Vítková, M., Marec, F., 2004. Phylogenetic  
1219 distribution of TTAGG telomeric repeats in insects. *Genome* 47, 163–178.  
1220 <https://doi.org/10.1139/g03-100>
- 1221 Garantonakis, N., Pappas, M.L., Varikou, K., Skiada, V., Broufas, G.D., Kavroulakis, N.,  
1222 Papadopoulou, K.K., 2018. Tomato Inoculation With the Endophytic Strain *Fusarium solani*  
1223 K Results in Reduced Feeding Damage by the Zoophytophagous Predator *Nesidiocoris*  
1224 *tenuis*. *Front. Ecol. Evol.* 6, 1–7. <https://doi.org/10.3389/fevo.2018.00126>
- 1225 Gatto, K.P., Mattos, J. V., Seger, K.R., Lourenço, L.B., 2018. Sex chromosome differentiation in  
1226 the frog genus *Pseudis* involves satellite DNA and chromosome rearrangements. *Front.*  
1227 *Genet.* 9, 1–12. <https://doi.org/10.3389/fgene.2018.00301>
- 1228 Gomez-Rodriguez, V.M., Rodriguez-Garay, B., Palomino, G., Martínez, J., Barba-Gonzalez, R.,  
1229 2013. Physical mapping of 5S and 18S ribosomal DNA in three species of *Agave*  
1230 (*Asparagales*, *Asparagaceae*). *Comp. Cytogenet.* 7, 191–203.  
1231 <https://doi.org/10.3897/compcytogen.v7i3.5337>
- 1232 Grozeva, S., Anokhin, B.A., Simov, N., Kuznetsova, V.G., 2019. New evidence for the presence  
1233 of the telomere motif (TTAGG) n in the family Reduviidae and its absence in the families  
1234 Nabidae and Miridae (Hemiptera, Cimicomorpha). *Comp. Cytogenet.* 13, 283–295.  
1235 <https://doi.org/10.3897/compcytogen.v13i3.36676>
- 1236 Grozeva, S., Kuznetsova, V.G., Anokhin, B.A., 2011. Karyotypes, male meiosis and  
1237 comparative FISH mapping of 18S ribosomal DNA and telomeric (TTAGG) n repeat in

- 1238 eight species of true bugs (Hemiptera, Heteroptera). *Comp. Cytogenet.* 5, 97–116.  
1239 <https://doi.org/10.3897/CompCytogen.v5i4.2307>
- 1240 Gurevich, A., Saveliev, V., Vyahhi, N., Tesler, G., 2013. QUILT: Quality assessment tool for  
1241 genome assemblies. *Bioinformatics* 29, 1072–1075.  
1242 <https://doi.org/10.1093/bioinformatics/btt086>
- 1243 Haas, B.J., Salzberg, S.L., Zhu, W., Pertea, M., Allen, J.E., Orvis, J., White, O., Robin, C.R.,  
1244 Wortman, J.R., 2008. Automated eukaryotic gene structure annotation using  
1245 EVIDENCEModeler and the Program to Assemble Spliced Alignments. *Genome Biol.* 9, 1–  
1246 22. <https://doi.org/10.1186/gb-2008-9-1-r7>
- 1247 Hare, E.E., Johnston, J.S., 2011. Genome Size Determination Using Flow Cytometry of  
1248 Propidium Iodide-Stained Nuclei, *Molecular Methods for Evolutionary Genetics*.  
1249 <https://doi.org/10.1007/978-1-61779-228-1>
- 1250 Huang, D.W., Sherman, B.T., Lempicki, R.A., 2009a. Bioinformatics enrichment tools: Paths  
1251 toward the comprehensive functional analysis of large gene lists. *Nucleic Acids Res.* 37,  
1252 1–13. <https://doi.org/10.1093/nar/gkn923>
- 1253 Huang, D.W., Sherman, B.T., Lempicki, R.A., 2009b. Systematic and integrative analysis of large  
1254 gene lists using DAVID bioinformatics resources. *Nat. Protoc.* 4, 44–57.  
1255 <https://doi.org/10.1038/nprot.2008.211>
- 1256 Huang, H., McGarvey, P.B., Suzek, B.E., Mazumder, R., Zhang, J., Chen, Y., Wu, C.H., 2011. A  
1257 comprehensive protein-centric ID mapping service for molecular data integration.  
1258 *Bioinformatics* 27, 1190–1191. <https://doi.org/10.1093/bioinformatics/btr101>
- 1259 Husnik, F., McCutcheon, J.P., 2018. Functional horizontal gene transfer from bacteria to  
1260 eukaryotes. *Nat. Rev. Microbiol.* 16, 67–79. <https://doi.org/10.1038/nrmicro.2017.137>
- 1261 Itou, M., Watanabe, M., Watanabe, E., Miura, K., 2013. Gut content analysis to study  
1262 predatory efficacy of *Nesidiocoris tenuis* (Reuter) (Hemiptera: Miridae) by molecular  
1263 methods. *Entomol. Sci.* 16, 145–150. <https://doi.org/10.1111/j.1479-8298.2012.00552.x>
- 1264 Jang, E.B., Nishijima, K.A., 1990. Identification and Attractancy of Bacteria Associated with  
1265 *Dacus dorsalis* (Diptera: Tephritidae). *Environ. Entomol.* 19, 1726–1731.  
1266 <https://doi.org/10.1093/ee/19.6.1726>
- 1267 Jauset, A.M., Edo-Tena, E., Castañé, C., Agustí, N., Alomar, O., Grozeva, S., 2015.  
1268 Comparative cytogenetic study of three *Macrolophus* species (Heteroptera, Miridae).  
1269 *Comp. Cytogenet.* 9, 613–623. <https://doi.org/10.3897/CompCytogen.v9i4.5530>
- 1270 Jetybayev, I.Y., Bugrov, A.G., Ünal, M., Buleu, O.G., Rubtsov, N.B., 2017. Molecular  
1271 cytogenetic analysis reveals the existence of two independent neo-XY sex chromosome  
1272 systems in Anatolian Pamphagidae grasshoppers. *BMC Evol. Biol.* 17, 1–12.  
1273 <https://doi.org/10.1186/s12862-016-0868-9>
- 1274 Jones, R.N., Rees, H., 1982. B Chromosomes. Academic Press, New York.
- 1275 Jones, Steven, Taylor, G., Chan, S., Warren, R., Hammond, S., Bilobram, S., Mordecai, G.,  
1276 Suttle, C., Miller, K., Schulze, A., Chan, A., Jones, Samantha, Tse, K., Li, I., Cheung, D.,  
1277 Mungall, K., Choo, C., Ally, A., Dhalla, N., Tam, A., Troussard, A., Kirk, H., Pandoh, P.,  
1278 Paulino, D., Coope, R., Mungall, A., Moore, R., Zhao, Y., Birol, I., Ma, Y., Marra, M.,  
1279 Haulena, M., 2017. The Genome of the Beluga Whale (*Delphinapterus leucas*). *Genes*  
1280 (Basel). 8, 378. <https://doi.org/10.3390/genes8120378>
- 1281 Jung, S., Lee, S., 2012. Molecular phylogeny of the plant bugs (Heteroptera: Miridae) and the  
1282 evolution of feeding habits. *Cladistics* 28, 50–79. <https://doi.org/10.1111/j.1096->

- 1283 0031.2011.00365.x
- 1284 Kamber, T., Smits, T.H.M., Rezzonico, F., Duffy, B., 2012. Genomics and current genetic  
1285 understanding of *Erwinia amylovora* and the fire blight antagonist *Pantoea vagans*.  
1286 *Trees - Struct. Funct.* 26, 227–238. <https://doi.org/10.1007/s00468-011-0619-x>
- 1287 Kato, A., Albert, P.S., Vega, J.M., Birchler, J.A., 2006. Sensitive fluorescence in situ hybridization  
1288 signal detection in maize using directly labeled probes produced by high concentration  
1289 DNA polymerase nick translation. *Biotech. Histochem.* 81, 71–78.  
1290 <https://doi.org/10.1080/10520290600643677>
- 1291 Keeling, C.I., Yuen, M.M.S., Liao, N.Y., Roderick Docking, T., Chan, S.K., Taylor, G.A., Palmquist,  
1292 D.L., Jackman, S.D., Nguyen, A., Li, M., Henderson, H., Janes, J.K., Zhao, Y., Pandoh, P.,  
1293 Moore, R., Sperling, F.A.H., W Huber, D.P., Birol, I., Jones, S.J.M., Bohlmann, J., 2013. Draft  
1294 genome of the mountain pine beetle, *Dendroctonus ponderosae* Hopkins, a major  
1295 forest pest. *Genome Biol.* 14, R27. <https://doi.org/10.1186/gb-2013-14-3-r27>
- 1296 Keilwagen, J., Wenk, M., Erickson, J.L., Schattat, M.H., Grau, J., Hartung, F., 2016. Using intron  
1297 position conservation for homology-based gene prediction. *Nucleic Acids Res.* 44.  
1298 <https://doi.org/10.1093/nar/gkw092>
- 1299 Kofler, R., Orozco-terWengel, P., de Maio, N., Pandey, R.V., Nolte, V., Futschik, A., Kosiol, C.,  
1300 Schlötterer, C., 2011. Popoolation: A toolbox for population genetic analysis of next  
1301 generation sequencing data from pooled individuals. *PLoS One* 6.  
1302 <https://doi.org/10.1371/journal.pone.0015925>
- 1303 Kuznetsova, V.G., Grozeva, S.M., Nokkala, S., Nokkala, C., 2011. Cytogenetics of the true bug  
1304 infraorder cimicomorpha (hemiptera, heteroptera): A review. *Zookeys* 154, 31–70.  
1305 <https://doi.org/10.3897/zookeys.154.1953>
- 1306 Lee, W., Kang, J., Jung, C., Hoelmer, K., Lee, S.H., Lee, S., 2009. Complete mitochondrial  
1307 genome of brown marmorated stink bug *Halyomorpha halys* (Hemiptera:  
1308 Pentatomidae), and phylogenetic relationships of hemipteran suborders. *Mol. Cells* 28,  
1309 155–165. <https://doi.org/10.1007/s10059-009-0125-9>
- 1310 Legeai, F., Shigenobu, S., Gauthier, J.P., Colbourne, J., Rispe, C., Collin, O., Richards, S., Wilson,  
1311 A.C.C., Murphy, T., Tagu, D., 2010. AphidBase: A centralized bioinformatic resource for  
1312 annotation of the pea aphid genome. *Insect Mol. Biol.* 19, 5–12.  
1313 <https://doi.org/10.1111/j.1365-2583.2009.00930.x>
- 1314 Leung, K., Ras, E., Ferguson, K.B., Ariëns, S., Babendreier, D.B., Bijma, P., Bourtzis, K., Brodeur, J.,  
1315 Bruins, M., Centurión, A., Chattington, S., Chinchilla-Ramírez, M., Dicke, M., Fatouros, N.,  
1316 González Cabrera, J., Groot, T., Haye, T., Knapp, M., Koskinioti, P., Le Hesran, S., Lirakis,  
1317 M., Paspati, A., Pérez-Hedo, M., Plouvier, W., Schlötterer, C., Stahl, J., Thiel, A., Urbaneja,  
1318 A., van de Zande, L., Verhulst, E., Vet, L., Visser, S., Werren, J., Xia, S., Zwaan, B.,  
1319 Magalhães, S., Beukeboom, L., Pannebakker, B., 2019. Next Generation Biological  
1320 Control: the Need for Integrating Genetics and Evolution. *Preprints* 1–34.  
1321 <https://doi.org/10.20944/preprints201911.0300.v1>
- 1322 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G.,  
1323 Durbin, R., 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25,  
1324 2078–2079. <https://doi.org/10.1093/bioinformatics/btp352>
- 1325 Lindsey, A.R.I., Kelkar, Y.D., Wu, X., Sun, D., Martinson, E.O., Yan, Z., Rugman-Jones, P.F.,  
1326 Hughes, D.S.T., Murali, S.C., Qu, J., Dugan, S., Lee, S.L., Chao, H., Dinh, H., Han, Y.,  
1327 Doddapaneni, H.V., Worley, K.C., Muzny, D.M., Ye, G., Gibbs, R.A., Richards, S., Yi, S. V.,  
1328 Stouthamer, R., Werren, J.H., 2018. Comparative genomics of the miniature wasp and  
1329 pest control agent *Trichogramma pretiosum*. *BMC Biol.* 16, 1–20.  
1330 <https://doi.org/10.1186/s12915-018-0520-9>



- 1331 Lindsey, A.R.I., Werren, J.H., Richards, S., Stouthamer, R., 2016. Comparative Genomics of a  
1332 Parthenogenesis-Inducing *Wolbachia* Symbiont. *G3 Genes, Genomes, Genet.* 6, 2113–  
1333 2123. <https://doi.org/10.1534/g3.116.028449>
- 1334 Lowe-Power, T.M., Khokhani, D., Allen, C., 2018. How *Ralstonia solanacearum* Exploits and  
1335 Thrives in the Flowing Plant Xylem Environment. *Trends Microbiol.* 26, 929–942.  
1336 <https://doi.org/10.1016/j.tim.2018.06.002>
- 1337 Majoros, W.H., Pertea, M., Salzberg, S.L., 2004. TigrScan and GlimmerHMM: Two open source  
1338 ab initio eukaryotic gene-finders. *Bioinformatics* 20, 2878–2879.  
1339 <https://doi.org/10.1093/bioinformatics/bth315>
- 1340 Marçais, G., Kingsford, C., 2011. A fast, lock-free approach for efficient parallel counting of  
1341 occurrences of k-mers. *Bioinformatics* 27, 764–770.  
1342 <https://doi.org/10.1093/bioinformatics/btr011>
- 1343 Marec, F., Shvedov, A.N., 1990. Yellow eye, a new pigment mutation in *Ephestia kuehniella*  
1344 Zeller (Lepidoptera: Pyralidae). *Hereditas* 113, 97–100. <https://doi.org/10.1111/j.1601-5223.1990.tb00072.x>
- 1346 Martínez-García, H., Román-Fernández, L.R., Sáenz-Romo, M.G., Pérez-Moreno, I., Marco-  
1347 Mancebón, V.S., 2016. Optimizing *Nesidiocoris tenuis* (Hemiptera: Miridae) as a  
1348 biological control agent: mathematical models for predicting its development as a  
1349 function of temperature. *Bull. Entomol. Res.* 106, 215–224.  
1350 <https://doi.org/10.1017/S0007485315000978>
- 1351 May, C.M., Heuvel, J., Doroszuk, A., Hoedjes, K.M., Flatt, T., Zwaan, B.J., 2019. Adaptation to  
1352 developmental diet influences the response to selection on age at reproduction in the  
1353 fruit fly. *J. Evol. Biol.* 32, 425–437. <https://doi.org/10.1111/jeb.13425>
- 1354 Mediouni, J., Fuková, I., Frydrychová, R., Marec, F., Fuková, I., Marec, F., Dhouibi, M.H., 2004.  
1355 Karyotype, sex chromatin and sex chromosome differentiation in the carob moth,  
1356 *Ectomyelois ceratoniae* (Lepidoptera: Pyralidae). *Caryologia* 57, 184–194.  
1357 <https://doi.org/10.1080/00087114.2004.10589391>
- 1358 Mollá, O., Biondi, A., Alonso-Valiente, M., Urbaneja, A., 2014. A comparative life history study  
1359 of two mirid bugs preying on *Tuta absoluta* and *Ephestia kuehniella* eggs on tomato  
1360 crops: Implications for biological control. *BioControl* 59, 175–183.  
1361 <https://doi.org/10.1007/s10526-013-9553-8>
- 1362 Morgulis, A., Gertz, E.M., Schäffer, A.A., Agarwala, R., 2006. A fast and symmetric DUST  
1363 implementation to mask low-complexity DNA sequences. *J. Comput. Biol.* 13, 1028–1040.  
1364 <https://doi.org/10.1089/cmb.2006.13.1028>
- 1365 Novák, P., Neumann, P., Pech, J., Steinhaisl, J., MacAs, J., 2013. RepeatExplorer: A Galaxy-  
1366 based web server for genome-wide characterization of eukaryotic repetitive elements  
1367 from next-generation sequence reads. *Bioinformatics* 29, 792–793.  
1368 <https://doi.org/10.1093/bioinformatics/btt054>
- 1369 Panfilio, K.A., Angelini, D.R., 2018. By land, air, and sea: hemipteran diversity through the  
1370 genomic lens. *Curr. Opin. Insect Sci.* 25, 106–115.  
1371 <https://doi.org/10.1016/j.cois.2017.12.005>
- 1372 Panfilio, K.A., Vargas Jentsch, I.M., Benoit, J.B., Erezyilmaz, D., Suzuki, Y., Colella, S., Robertson,  
1373 H.M., Poelchau, M.F., Waterhouse, R.M., Ioannidis, P., Weirauch, M.T., Hughes, D.S.T.,  
1374 Murali, S.C., Werren, J.H., Jacobs, C.G.C., Duncan, E.J., Armisén, D., Vreede, B.M.I., Baa-  
1375 Puyoulet, P., Berger, C.S., Chang, C., Chao, H., Chen, M.-J.M., Chen, Y.-T., Childers, C.P.,  
1376 Chipman, A.D., Cridge, A.G., Crumière, A.J.J., Dearden, P.K., Didion, E.M., Dinh, H.,  
1377 Doddapaneni, H.V., Dolan, A., Dugan, S., Extavour, C.G., Febvay, G., Friedrich, M.,

- 1378 Ginzburg, N., Han, Y., Heger, P., Holmes, C.J., Horn, T., Hsiao, Y., Jennings, E.C., Johnston,  
1379 J.S., Jones, T.E., Jones, J.W., Khila, A., Koelzer, S., Kovacova, V., Leask, M., Lee, S.L., Lee,  
1380 C.-Y., Lovegrove, M.R., Lu, H., Lu, Y., Moore, P.J., Munoz-Torres, M.C., Muzny, D.M., Palli,  
1381 S.R., Parisot, N., Pick, L., Porter, M.L., Qu, J., Refki, P.N., Richter, R., Rivera-Pomar, R.,  
1382 Rosendale, A.J., Roth, S., Sachs, L., Santos, M.E., Seibert, J., Sghaier, E., Shukla, J.N.,  
1383 Stancliffe, R.J., Tidswell, O., Traverso, L., van der Zee, M., Viala, S., Worley, K.C., Zdobnov,  
1384 E.M., Gibbs, R.A., Richards, S., 2019. Molecular evolutionary trends and feeding ecology  
1385 diversification in the Hemiptera, anchored by the milkweed bug genome. *Genome Biol.*  
1386 20, 64. <https://doi.org/10.1186/s13059-019-1660-0>
- 1387 Paspati, A., Ferguson, K.B., Verhulst, E.C., Urbaneja, A., González-Cabrera, J., Pannebakker,  
1388 B.A., 2019. Effect of mass rearing on the genetic diversity of the predatory mite  
1389 *Amblyseius swirskii* Athias-Henriot (Acari: Phytoseiidae). *Entomol. Exp. Appl.* 167, 670–681.  
1390 <https://doi.org/10.1111/eea.12811>
- 1391 Pérez-Hedo, M., Arias-Sanguino, Á.M., Urbaneja, A., 2018. Induced tomato plant resistance  
1392 against *tetranychus urticae* triggered by the phytophagy of *nesidiocoris tenuis*. *Front.*  
1393 *Plant Sci.* 9, 1–8. <https://doi.org/10.3389/fpls.2018.01419>
- 1394 Pérez-Hedo, M., Urbaneja-Bernat, P., Jaques, J.A., Flors, V., Urbaneja, A., 2015. Defensive plant  
1395 responses induced by *Nesidiocoris tenuis* (Hemiptera: Miridae) on tomato plants. *J. Pest*  
1396 *Sci.* (2004). 88, 543–554. <https://doi.org/10.1007/s10340-014-0640-0>
- 1397 Pérez-Hedo, M., Urbaneja, A., 2016. The Zoophytophagous Predator *Nesidiocoris tenuis*: A  
1398 Successful But Controversial Biocontrol Agent in Tomato Crops, in: *Advances in Insect*  
1399 *Control and Resistance Management*. Springer International Publishing, Cham, pp. 121–  
1400 138. [https://doi.org/10.1007/978-3-319-31800-4\\_7](https://doi.org/10.1007/978-3-319-31800-4_7)
- 1401 Pita, S., Panzera, F., Mora, P., Vela, J., Palomeque, T., Lorite, P., 2016. The presence of the  
1402 ancestral insect telomeric motif in kissing bugs (Triatominae) rules out the hypothesis of its  
1403 loss in evolutionarily advanced Heteroptera (Cimicomorpha). *Comp. Cytogenet.* 10,  
1404 427–437. <https://doi.org/10.3897/compcytogen.v10i3.9960>
- 1405 Poelchau, M., Childers, C., Moore, G., Tsavatapalli, V., Evans, J., Lee, C.Y., Lin, H., Lin, J.W.,  
1406 Hackett, K., 2015. The i5k Workspace@NAL-enabling genomic data access, visualization  
1407 and curation of arthropod genomes. *Nucleic Acids Res.* 43, D714–D719.  
1408 <https://doi.org/10.1093/nar/gku983>
- 1409 Prado, S.S., Almeida, R.P.P., 2009. Phylogenetic placement of pentatomid stink bug gut  
1410 symbionts. *Curr. Microbiol.* 58, 64–69. <https://doi.org/10.1007/s00284-008-9267-9>
- 1411 Puentes, A., Stephan, J.G., Björkman, C., 2018. A Systematic Review on the Effects of Plant-  
1412 Feeding by Omnivorous Arthropods: Time to Catch-Up With the Mirid-Tomato Bias? *Front.*  
1413 *Ecol. Evol.* 6. <https://doi.org/10.3389/fevo.2018.00218>
- 1414 Quan, Q., Hu, X., Pan, B., Zeng, B., Wu, N., Fang, G., Cao, Y., Chen, X., Li, X., Huang, Y., Zhan,  
1415 S., 2019. Draft genome of the cotton aphid *Aphis gossypii*. *Insect Biochem. Mol. Biol.* 105,  
1416 25–32. <https://doi.org/10.1016/j.ibmb.2018.12.007>
- 1417 Rasmussen, L.B., Jensen, K., Sørensen, J.G., Sverrisdóttir, E., Nielsen, K.L., Overgaard, J.,  
1418 Holmstrup, M., Kristensen, T.N., 2018. Are commercial stocks of biological control agents  
1419 genetically depauperate? – A case study on the pirate bug *Orius majusculus* Reuter.  
1420 *Biol. Control* 127, 31–38. <https://doi.org/10.1016/j.biocontrol.2018.08.016>
- 1421 Reddiex, A.J., Allen, S.L., Chenoweth, S.F., 2018. A Genomic Reference Panel for *Drosophila*  
1422 *serrata*. *G3 Genes, Genomes, Genet.* 8, 1335–1346.  
1423 <https://doi.org/10.1534/g3.117.300487>
- 1424 Richards, S., Gibbs, R.A., Gerardo, N.M., Moran, N., Nakabachi, A., Stern, D., Tagu, D., Wilson,

- 1425 A.C.C., Muzny, D., Kovar, C., Cree, A., Chacko, J., Chandrabose, M.N., Dao, M.D., Dinh,  
1426 H.H., Gabisi, R.A., Hines, S., Hume, J., Jhangian, S.N., Joshi, V., Lewis, L.R., Liu, Y.S., Lopez,  
1427 J., Morgan, M.B., Nguyen, N.B., Okwuonu, G.O., Ruiz, S.J., Santibanez, J., Wright, R.A.,  
1428 Fowler, G.R., Hitchens, M.E., Lozado, R.J., Moen, C., Steffen, D., Warren, J.T., Zhang, J.,  
1429 Nazareth, L. V., Chavez, D., Davis, C., Lee, S.L., Patel, B.M., Pu, L.L., Bell, S.N., Johnson,  
1430 A.J., Vattathil, S., Williams, R.L., Shigenobu, S., Dang, P.M., Morioka, M., Fukatsu, T., Kudo,  
1431 T., Miyagishima, S.Y., Jiang, H., Worley, K.C., Legeai, F., Gauthier, J.P., Collin, O., Zhang, L.,  
1432 Chen, H.C., Ermolaeva, O., Hlavina, W., Kapustin, Y., Kiryutin, B., Kitts, P., Maglott, D.,  
1433 Murphy, T., Pruitt, K., Sapojnikov, V., Souvorov, A., Thibaud-Nissen, F., Câmara, F., Guigó,  
1434 R., Stanke, M., Solovyev, V., Kosarev, P., Gilbert, D., Gabaldón, T., Huerta-Cepas, J.,  
1435 Marcet-Houben, M., Pignatelli, M., Moya, A., Rispe, C., Ollivier, M., Quesneville, H.,  
1436 Permal, E., Llorens, C., Futami, R., Hedges, D., Robertson, H.M., Alioto, T., Mariotti, M.,  
1437 Nikoh, N., McCutcheon, J.P., Burke, G., Kamins, A., Latorre, A., Ashton, P., Calevro, F.,  
1438 Charles, H., Colella, S., Douglas, A.E., Jander, G., Jones, D.H., Febvay, G., Kamphuis, L.G.,  
1439 Kushlan, P.F., Macdonald, S., Ramsey, J., Schwartz, J., Seah, S., Thomas, G., Vellozo, A.,  
1440 Cass, B., Degnan, P., Hurwitz, B., Leonardo, T., Koga, R., Altincicek, B., Anselme, C.,  
1441 Atamian, H., Barribeau, S.M., De Vos, M., Duncan, E.J., Evans, J., Ghanim, M., Heddi, A.,  
1442 Kaloshian, I., Vincent-Monegat, C., Parker, B.J., Pérez-Brocal, V., Rahbé, Y., Spragg, C.J.,  
1443 Tamames, J., Tamarit, D., Tamborineguy, C., Vilcinskis, A., Bickel, R.D., Brisson, J.A.,  
1444 Butts, T., Chang, C.C., Christiaens, O., Davis, G.K., Duncan, E., Ferrier, D., Iga, M., Janssen,  
1445 R., Lu, H.L., McGregor, A., Miura, T., Smagghe, G., Smith, J., Van Der Zee, M., Velarde, R.,  
1446 Wilson, M., Dearden, P., Edwards, O.R., Gordon, K., Hilgarth, R.S., Rider, S.D., Srinivasan,  
1447 D., Walsh, T.K., Ishikawa, A., Jaubert-Possamai, S., Fenton, B., Huang, W., Rizk, G.,  
1448 Lavenier, D., Nicolas, J., Smadja, C., Zhou, J.J., Vieira, F.G., He, X.L., Liu, R., Rozas, J., Field,  
1449 L.M., Campbell, P., Carolan, J.C., Fitzroy, C.I.J., Reardon, K.T., Reeck, G.R., Singh, K.,  
1450 Wilkinson, T.L., Huybrechts, J., Abdel-Latif, M., Robichon, A., Veenstra, J.A., Hauser, F.,  
1451 Cazzamali, G., Schneider, M., Williamson, M., Stafflinger, E., Hansen, K.K.,  
1452 Grimmelikhuijzen, C.J.P., Price, D.R.G., Caillaud, M., Van Fleet, E., Ren, Q., Gatehouse,  
1453 J.A., Brault, V., Monsion, B., Diaz, J., Hunnicutt, L., Ju, H.J., Pechuan, X., Aguilar, J., Cortés,  
1454 T., Ortiz-Rivas, B., Martínez-Torres, D., Dombrovsky, A., Dale, R.P., Davies, T.G.E., Williamson,  
1455 M.S., Jones, A., Sattelle, D., Williamson, S., Wolstenholme, A., Cottret, L., Sagot, M.F.,  
1456 Heckel, D.G., Hunter, W., 2010. Genome sequence of the pea aphid *Acyrtosiphon*  
1457 *pisum*. *PLoS Biol.* 8. <https://doi.org/10.1371/journal.pbio.1000313>
- 1458 Richards, S., Murali, S.C., 2015. Best practices in insect genome sequencing: what works and  
1459 what doesn't. *Curr. Opin. Insect Sci.* 7, 1–7. <https://doi.org/10.1016/j.cois.2015.02.013>
- 1460 Sachman-Ruiz, B., Quiroz-Castañeda, R.E., 2018. Genomics of Rickettsiaceae: An Update, in:  
1461 Farm Animals Diseases, Recent Omic Trends and New Strategies of Treatment. InTech, p.  
1462 13. <https://doi.org/10.5772/intechopen.74563>
- 1463 Sahara, K., Marec, F., Traut, W., 1999. TTAGG telomeric repeats in chromosomes of some  
1464 insects and other arthropods. *Chromosom. Res.* 7, 449–460.  
1465 <https://doi.org/10.1023/A:1009297729547>
- 1466 Sanchez, J.A., 2009. Density thresholds for *Nesidiocoris tenuis* (Heteroptera: Miridae) in tomato  
1467 crops. *Biol. Control* 51, 493–498. <https://doi.org/10.1016/j.biocontrol.2009.09.006>
- 1468 Santos-Garcia, D., Silva, F.J., Morin, S., Dettner, K., Kuechler, S.M., 2017. The all-rounder *sodalis*:  
1469 A new bacteriome-associated endosymbiont of the lygaeoid bug *henestaris halophilus*  
1470 (heteroptera: Henestarinae) and a critical examination of its evolution. *Genome Biol.*  
1471 *Evol.* 9, 2893–2910. <https://doi.org/10.1093/gbe/evx202>
- 1472 Schwarz, A., Medrano-Mercado, N., Schaub, G.A., Struchiner, C.J., Bargues, M.D., Levy, M.Z.,  
1473 Ribeiro, J.M.C., 2014. An Updated Insight into the Sialotranscriptome of *Triatoma*  
1474 *infestans*: Developmental Stage and Geographic Variations. *PLoS Negl. Trop. Dis.* 8.  
1475 <https://doi.org/10.1371/journal.pntd.0003372>

- 1476 Sember, A., Bertollo, L.A.C., Ráb, P., Yano, C.F., Hatanaka, T., de Oliveira, E.A., Cioffi, M. de B.,  
1477 2018. Sex chromosome evolution and genomic divergence in the fish *Hoplias*  
1478 *malabaricus* (Characiformes, Erythrinidae). *Front. Genet.* 9, 1–12.  
1479 <https://doi.org/10.3389/fgene.2018.00071>
- 1480 Simão, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E. V., Zdobnov, E.M., 2015. BUSCO:  
1481 Assessing genome assembly and annotation completeness with single-copy orthologs.  
1482 *Bioinformatics* 31, 3210–3212. <https://doi.org/10.1093/bioinformatics/btv351>
- 1483 Sobreira, T.J.P., Durham, A.M., Gruber, A., 2006. TRAP: Automated classification, quantification  
1484 and annotation of tandemly repeated sequences. *Bioinformatics* 22, 361–362.  
1485 <https://doi.org/10.1093/bioinformatics/bti809>
- 1486 Sochorová, J., Garcia, S., Gálvez, F., Symonová, R., Kovařík, A., 2018. Evolutionary trends in  
1487 animal ribosomal DNA loci: introduction to a new online database. *Chromosoma* 127,  
1488 141–150. <https://doi.org/10.1007/s00412-017-0651-8>
- 1489 Streito, J.C., Clouet, C., Hamdi, F., Gauthier, N., 2017. Population genetic structure of the  
1490 biological control agent *Macrolophus pygmaeus* in Mediterranean agroecosystems.  
1491 *Insect Sci.* 24, 859–876. <https://doi.org/10.1111/1744-7917.12370>
- 1492 Sunnucks, P., Hales, D.F., 1996. Numerous transposed sequences of mitochondrial cytochrome  
1493 oxidase I-II in aphids of the genus *Sitobion* (Hemiptera: Aphididae). *Mol. Biol. Evol.* 13,  
1494 510–524. <https://doi.org/10.1093/oxfordjournals.molbev.a025612>
- 1495 Szűcs, M., Vercken, E., Bitume, E. V., Hufbauer, R.A., 2019. The implications of rapid eco-  
1496 evolutionary processes for biological control - a review. *Entomol. Exp. Appl.* eea.12807.  
1497 <https://doi.org/10.1111/eea.12807>
- 1498 Testa, A.C., Hane, J.K., Ellwood, S.R., Oliver, R.P., 2015. CodingQuarry: Highly accurate hidden  
1499 Markov model gene prediction in fungal genomes using RNA-seq transcripts. *BMC*  
1500 *Genomics* 16, 1–12. <https://doi.org/10.1186/s12864-015-1344-4>
- 1501 Thomas, G.W.C., Dohmen, E., Hughes, D.S.T., Murali, S.C., Poelchau, M., Glastad, K., Anstead,  
1502 C.A., Ayoub, N.A., Batterham, P., Bellair, M., Binford, G.J., Chao, H., Chen, Y.H., Childers,  
1503 C., Dinh, H., Doddapaneni, H.V., Duan, J.J., Dugan, S., Esposito, L.A., Friedrich, M., Garb,  
1504 J., Gasser, R.B., Goodisman, M.A.D., Gundersen-Rindal, D.E., Han, Y., Handler, A.M.,  
1505 Hatakeyama, M., Hering, L., Hunter, W.B., Ioannidis, P., Jayaseelan, J.C., Kalra, D., Khila,  
1506 A., Korhonen, P.K., Lee, C.E., Lee, S.L., Li, Y., Lindsey, A.R.I., Mayer, G., McGregor, A.P.,  
1507 McKenna, D.D., Misof, B., Munidasa, M., Munoz-Torres, M., Muzny, D.M., Niehuis, O., Osuji-  
1508 Lacy, N., Palli, S.R., Panfilio, K.A., Pechmann, M., Perry, T., Peters, R.S., Poynton, H.C., Prpic,  
1509 N.-M., Qu, J., Rotenberg, D., Schal, C., Schoville, S.D., Scully, E.D., Skinner, E., Sloan, D.B.,  
1510 Stouthamer, R., Strand, M.R., Szucsich, N.U., Wijeratne, A., Young, N.D., Zattara, E.E.,  
1511 Benoit, J.B., Zdobnov, E.M., Pfrender, M.E., Hackett, K.J., Werren, J.H., Worley, K.C., Gibbs,  
1512 R.A., Chipman, A.D., Waterhouse, R.M., Bornberg-Bauer, E., Hahn, M.W., Richards, S.,  
1513 2020. Gene content evolution in the arthropods. *Genome Biol.* 21, 15.  
1514 <https://doi.org/10.1186/s13059-019-1925-7>
- 1515 Tian, C., Tek Tay, W., Feng, H., Wang, Y., Hu, Y., Li, G., 2015. Characterization of *Adelphocoris*  
1516 *suturalis* (Hemiptera: Miridae) Transcriptome from Different Developmental Stages. *Sci.*  
1517 *Rep.* 5, 11042. <https://doi.org/10.1038/srep11042>
- 1518 Trapnell, C., Pachter, L., Salzberg, S.L., 2009. TopHat: Discovering splice junctions with RNA-  
1519 Seq. *Bioinformatics* 25, 1105–1111. <https://doi.org/10.1093/bioinformatics/btp120>
- 1520 Trapnell, C., Williams, B.A., Pertea, G., Mortazavi, A., Kwan, G., van Baren, M.J., Salzberg, S.L.,  
1521 Wold, B.J., Pachter, L., 2010. Transcript assembly and quantification by RNA-Seq reveals  
1522 unannotated transcripts and isoform switching during cell differentiation. *Nat.*  
1523 *Biotechnol.* 28, 511–5. <https://doi.org/10.1038/nbt.1621>

- 1524 Traut, W., 1976. Pachytene mapping in the female silkworm, *Bombyx mori* L. (Lepidoptera).  
1525 *Chromosoma* 58, 275–284. <https://doi.org/10.1007/BF00292094>
- 1526 Traut, W., Sahara, K., Otto, T.D., Marec, F., 1999. Molecular differentiation of sex chromosomes  
1527 probed by comparative genomic hybridization. *Chromosoma* 108, 173–180.  
1528 <https://doi.org/10.1007/s004120050366>
- 1529 Traverse, K.L., Pardue, M.L., 1988. A spontaneously opened ring chromosome of *Drosophila*  
1530 *melanogaster* has acquired He-T DNA sequences at both new telomeres. *Proc. Natl.*  
1531 *Acad. Sci. U. S. A.* 85, 8116–8120. <https://doi.org/10.1073/pnas.85.21.8116>
- 1532 Urbaneja-Bernat, P., Bru, P., González-Cabrera, J., Urbaneja, A., Tena, A., 2019. Reduced  
1533 phytophagy in sugar-provisioned mirids. *J. Pest Sci.* (2004). 92, 1139–1148.  
1534 <https://doi.org/10.1007/s10340-019-01105-9>
- 1535 Urbaneja, A., González-Cabrera, J., Arnó, J., Gabarra, R., 2012. Prospects for the biological  
1536 control of *Tuta absoluta* in tomatoes of the Mediterranean basin. *Pest Manag. Sci.* 68,  
1537 1215–1222. <https://doi.org/10.1002/ps.3344>
- 1538 van Lenteren, J.C., Bolckmans, K., Köhl, J., Ravensberg, W.J., Urbaneja, A., 2018. Biological  
1539 control using invertebrates and microorganisms: plenty of new opportunities. *BioControl*  
1540 63, 39–59. <https://doi.org/10.1007/s10526-017-9801-4>
- 1541 van Wilgenburg, E., Driessen, G., Beukeboom, L.W., 2006. Single locus complementary sex  
1542 determination in Hymenoptera: an “unintelligent” design? *Front. Zool.* 3, 1.  
1543 <https://doi.org/https://doi.org/10.1186/1742-9994-3-1>
- 1544 Vurture, G.W., Sedlazeck, F.J., Nattestad, M., Schatz, M.C., Gurtowski, J., Underwood, C.J.,  
1545 Vurture, G.W., Fang, H., 2017. GenomeScope: fast reference-free genome profiling from  
1546 short reads. *Bioinformatics* 33, 2202–2204. <https://doi.org/10.1093/bioinformatics/btx153>
- 1547 Weisenfeld, N.I., Kumar, V., Shah, P., Church, D.M., Jaffe, D.B., 2017. Direct determination of  
1548 diploid genome sequences. *Genome Res.* 27, 757–767.  
1549 <https://doi.org/10.1101/gr.214874.116.Freely>
- 1550 Wheeler, D., Redding, A.J., Werren, J.H., 2013. Characterization of an Ancient Lepidopteran  
1551 Lateral Gene Transfer. *PLoS One* 8. <https://doi.org/10.1371/journal.pone.0059262>
- 1552 Xu, L., Dong, Z., Fang, L., Luo, Y., Wei, Z., Guo, H., Zhang, G., Gu, Y.Q., Coleman-Derr, D., Xia,  
1553 Q., Wang, Y., 2019. OrthoVenn2: A web server for genome wide comparison and  
1554 annotation of orthologous clusters across multiple species. *Nucleic Acids Res.* 47, W52–  
1555 W58. <https://doi.org/10.1093/nar/gkz333>
- 1556 Xu, P., Lu, B., Liu, J., Chao, J., Donkersley, P., Holdbrook, R., Lu, Y., 2019. Duplication and  
1557 expression of horizontally transferred polygalacturonase genes is associated with host  
1558 range expansion of mirid bugs. *BMC Evol. Biol.* 19, 12. <https://doi.org/10.1186/s12862-019-1351-1>
- 1560 Xun, H., Li, H., Li, S., Wei, S., Zhang, L., Song, F., Jiang, P., Yang, H., Han, F., Cai, W., 2016.  
1561 Population genetic structure and post-LGM expansion of the plant bug *Nesidiocoris*  
1562 *tenuis* (Hemiptera: Miridae) in China. *Sci. Rep.* 6, 26755.  
1563 <https://doi.org/10.1038/srep26755>
- 1564 Zhang, Y., Qiu, S., 2015. Examining phylogenetic relationships of *Erwinia* and *Pantoea* species  
1565 using whole genome sequence data. *Antonie van Leeuwenhoek, Int. J. Gen. Mol.*  
1566 *Microbiol.* 108, 1037–1046. <https://doi.org/10.1007/s10482-015-0556-6>
- 1567 Zhou, W., Rousset, F., O'Neill, S., 1998. Phylogeny and PCR-based classification of *Wolbachia*  
1568 strains using *wsp* gene sequences. *Proc. R. Soc. B Biol. Sci.* 265, 509–515.

1569 <https://doi.org/10.1098/rspb.1998.0324>

1570 Zrzavá, M., Hladová, I., Dalíková, M., Šíchová, J., Ůunap, E., Kubičková, S., Marec, F., 2018. Sex  
1571 chromosomes of the iconic moth *Abraxas grossulariata* (Lepidoptera, Geometridae)  
1572 and its congener *A. sylvata*. *Genes* (Basel). 9, 1–16.

1573 <https://doi.org/10.3390/genes9060279>

1574