

## University of Groningen

### Jekyll or Hyde?

Ferguson, Kim B.; Visser, Sander; Dalíková, Martina; Provazníková, Irena; Urbaneja, Alberto; Perez-Hedo, Meritxell; Marec, František; Werren, John (Jack); Zwaan, B.J.; Pannebakker, Bart

*Published in:*  
 Insect Molecular Biology

*DOI:*  
[10.1111/imb.12688](https://doi.org/10.1111/imb.12688)

**IMPORTANT NOTE:** You are advised to consult the publisher's version (publisher's PDF) if you wish to cite from it. Please check the document version below.

*Document Version*  
 Publisher's PDF, also known as Version of record

*Publication date:*  
 2021

[Link to publication in University of Groningen/UMCG research database](#)

*Citation for published version (APA):*

Ferguson, K. B., Visser, S., Dalíková, M., Provazníková, I., Urbaneja, A., Perez-Hedo, M., Marec, F., Werren, J. J., Zwaan, B. J., Pannebakker, B., & Verhulst, E. (2021). Jekyll or Hyde? The genome (and more) of *Nesidiocoris tenuis*, a zoophytophagous predatory bug that is both a biological control agent and a pest. *Insect Molecular Biology*, 30, 188-209. <https://doi.org/10.1111/imb.12688>

#### Copyright

Other than for strictly personal use, it is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license (like Creative Commons).











The publication may also be distributed here under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license. More information can be found on the University of Groningen website: <https://www.rug.nl/library/open-access/self-archiving-pure/taverne-amendment>.

#### Take-down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Downloaded from the University of Groningen/UMCG research database (Pure): <http://www.rug.nl/research/portal>. For technical reasons the number of authors shown on this cover page is limited to 10 maximum.

# 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\*,  S. Vissert†‡, M. Dalíková†‡,   
I. Provazníková†‡\*\*,  A. Urbaneja§,   
M. Pérez-Hedo§,  F. Marec†,  J. H. Werren¶,   
B. J. Zwaan\*,  B. A. Pannebakker\*  and  
E. C. Verhulst|| 

\*Laboratory of Genetics, Wageningen University, Wageningen, The Netherlands; †Biology Centre CAS, Institute of Entomology, České Budějovice, Czech Republic; ‡Faculty of Science, University of South Bohemia, České Budějovice, Czech Republic; §Centro de Protección Vegetal y Biotecnología, Instituto Valenciano de Investigaciones Agrarias (IVIA), Moncada, Spain; ¶Department of Biology, University of Rochester, Rochester, New York, USA; ||Laboratory of Entomology, Wageningen University, Wageningen, The Netherlands; and \*\*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. From the family Miridae, it is an economically important insect yet very little is known in terms of genetic information and no genomic or transcriptomic studies have been published. Here, we use 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. In addition, 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*). We visualized the genome using various cytogenetic techniques, such as karyotyping, CGH and GISH, indicating a karyotype of  $2n = 32$ . Additional analyses include the localization of 18S rDNA and unique satellite probes as well as pooled sequencing to assess nucleotide diversity and neutrality of the commercial population. This is one of the first mirid genomes to be released and the first of a mirid biological control agent.

**Keywords:** biocontrol, Hemiptera, linked-read.

## Introduction

Hemiptera is the fifth largest insect order and the most speciose hemimetabolous order with over 82 000 described species (Panfilio and Angelini, 2018). While recent sequencing projects have presented a variety of information about hemipteran genomes, large families such as the plant bugs Miridae still lack genomic resources, with the exception of transcriptomic resources for some members (Tian *et al.*, 2015), and the more recent genome of *Apolygus lucorum* Meyer-Dür, a mirid pest that has a publicly available genome as of December 2019 (NCBI BioProject PRJNA526332). With the exception of *A. lucorum*, the lack of genomic resources for Miridae is in spite of the diverse life histories present, as it contains not only some of the most notorious agricultural pests but also predators that are often used in biological control (van Lenteren *et al.*, 2018; Pérez-Hedo *et al.*, 2020). In addition, Hemiptera is known for their intriguing karyotype evolution involving holocentric (holokinetic) chromosomes but there is a lack of cytogenetic information on Miridae. The absence of the ancestral TTAGG<sub>n</sub> telomeric repeat has been reported for mirids *Macrolophus* spp., *Deraeocoris* spp. and *Megaloceroea recticornis* (Geoffroy) (Jauset *et al.*, 2015; Grozeva *et al.*, 2011, 2019) but more knowledge of this trait is necessary for evolutionary studies of genomes and karyotypes. Furthermore, the taxonomic issues that lie within both Miridae and Hemiptera could

First published online 22 December 2020.

Correspondence: K. B. Ferguson, Wageningen University, Laboratory of Genetics, Wageningen, The Netherlands. email: kimbfergus@gmail.com

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

*Nesidiocoris tenuis* (Reuter) (Hemiptera: Miridae) is a zoophytophagous mirid used as a biological control agent worldwide, including in Spain, the Mediterranean Basin and China (Pérez-Hedo and Urbaneja, 2016; Xun *et al.*, 2016). Throughout the Mediterranean Basin, *N. tenuis* is used in tomato greenhouses and open fields against whiteflies (Hemiptera: Aleyrodidae), and the South American tomato pinworm, *Tuta absoluta* (Meyrick) (Lepidoptera: Gelechiidae) (Calvo *et al.*, 2009; Mollá *et al.*, 2014). In addition, due to its high degree of polyphagous behaviour, it is able to prey on other pest species such as thrips, leaf miners, leafhoppers, aphids, spider mites and lepidopteran pests (Pérez-Hedo and Urbaneja, 2016). While *N. tenuis* is an important biological control agent in Mediterranean countries (Urbaneja *et al.*, 2012; Pérez-Hedo *et al.*, 2020), it is often cited as a pest in other contexts and countries (Calvo *et al.*, 2009; Pérez-Hedo and Urbaneja, 2016; Moerkens *et al.*, 2020). When prey is scarce in tomatoes, due to its phytophagy, *N. tenuis* can cause plant lesions such as brown discoloration around tender stems, known as necrotic rings, in addition to leaf wilt, and flower abortion (Arnó *et al.*, 2010; Chinchilla-Ramírez *et al.*, 2021). This switch to phytophagy has been observed to be inversely proportional to the availability of prey (Sanchez, 2009). Therefore, much of the research thus far has focused on characterizing *N. tenuis* biology and ecology, classifying the induced damage and attempting to reduce it (Castañé *et al.*, 2011; Biondi *et al.*, 2016; Martínez-García *et al.*, 2016; Garantonakis *et al.*, 2018; Urbaneja-Bernat *et al.*, 2019; Pérez-Hedo *et al.*, 2020). Despite its associated plant damage, *N. tenuis* is widely used across south-eastern Spain as it is an efficient agent against the various pests it controls (Arnó *et al.*, 2010). Furthermore, the aforementioned phytophagy has been demonstrated to have benefits by triggering predator-induced defences, including attracting parasitoids, repulsing other herbivorous pests and restricting accumulation of viruses (Pérez-Hedo *et al.*, 2015; Pérez-Hedo *et al.*, 2018; Bouagga *et al.*, 2020). Recent

research with isofemale populations of wild-caught *N. tenuis* indicate a genetic basis for variation in feeding behaviour (both zoophagy and phytophagy) (Chinchilla-Ramírez *et al.*, 2020).

In recent years, the controversial success of *N. tenuis* has encouraged the scientific community to study this predatory mirid (Puentes *et al.*, 2018; Pérez-Hedo *et al.*, 2020). However, some issues remain to be addressed, such as the genetic variation in commercial stocks of similar biological control agents when compared to wild populations, with the former often diminished in comparison to the latter as seen in other biological control agents (Streito *et al.*, 2017; Rasmussen *et al.*, 2018; Paspati *et al.*, 2019; Pérez-Hedo *et al.*, 2020). In order to compare biological control stock to wild (or wild-caught) populations, determining the current diversity and genetic variation of the commercial stock is important. Finally, *N. tenuis* is known to host bacterial symbionts, including *Wolbachia* and *Rickettsia*, though the effect of these bacteria on their host is relatively unknown (Caspi-Fluger *et al.*, 2014). Sequence data can provide additional insight into potential symbionts as well as identify potential LGTs (lateral gene transfers) between host and symbiont.

With all of these fascinating avenues of research in mind, it may be surprising to learn that, aside from a mitogenome (Dai *et al.*, 2012), a regional population analysis (Xun *et al.*, 2016) and more recent work shedding light on evidence of LGT (Xu *et al.*, 2019b), little genomic information exists for *N. tenuis* and there is no published *N. tenuis* genome. Characteristics such as karyotype, sex chromosome system and presence or absence of telomeric repeats are currently unknown. A likely reason for this absence of genomic resources is that advances made in sequencing technology are often juxtaposed to the complexities of insect life cycles and difficulties in obtaining enough high quality genomic material due to size and exoskeleton (Richards and Murali, 2015; Leung *et al.*, 2020). In addition, current assembly tools have a hard time dealing with heterozygosity; therefore, a genome assembly is benefited by sequencing material of reduced genetic heterozygosity for a more contiguous assembly. Reduced heterozygosity is often difficult to achieve in diploid insects where the genetic variation within a population is unknown or the species cannot be inbred (Keeling *et al.*, 2013).

Generating the genomes of highly heterozygous, diploid and relatively small insects is tricky; researchers have to be prepared to balance their expectations and the available technology (Ellegren, 2014; Leung *et al.*, 2020). While a single diploid individual may yield enough material for an Illumina-only library, assembly may be difficult due to large repeat regions that extend beyond the insert size of the library. Conversely, enough material could be obtained for sequencing on a long-read platform, but may require pooling material from multiple individuals, potentially

complicating assembly due to the heterozygosity of the population. While possible solutions include estimating the heterozygosity or setting up inbred populations [which can be nearly impossible if deleterious effects of inbreeding need to be avoided or if the presence of a complementary sex determining system limits inbreeding (van Wilgenburg *et al.*, 2006; Szűcs *et al.*, 2019)], an alternative is to create a linked-read library. The 10x Genomics platform creates a microfluidic partitioned library that individually barcodes minute amounts of long strands of DNA for further amplification (10x Genomics Inc., Pleasanton, CA). This library is then sequenced on a short-read sequencing platform and then assembled using the barcodes to link reads together into the larger fragment (eg Chin *et al.*, 2016; Jones *et al.*, 2017). This method allows for a library to be constructed from a single individual that contains additional structural information to aid assembly (such as phasing), removing the need for pooling multiple individuals and avoiding assembly difficulties in repetitive regions. Additional information, such as karyotype, can further improve genomes in the assembly stage as well as inform further directions of research by providing chromosome-level context, encouraging further improvement of a genome beyond its initial release.

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

## Results

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

The presence of *Wolbachia* in the KBS (Koppert Biological Systems) biological control stock was confirmed (Electrophoresis gel image in supplementary material, S1.1). All sequence data generated, including raw reads, assembly and annotation, can be found in the EMBL-EBI European Nucleotide Archive (ENA) under BioProject PRJEB35378, further available within the INSDC initiative databases. An additional, complete annotation file (.gff) is also available (<https://doi.org/10.6084/m9.figshare.12073893.v1>).

### *Genome assembly and size*

The single adult female *N. tenuis* yielded 424 ng total DNA. The 10X Genomics Chromium reaction and subsequent

Illumina sequencing resulted in more than 212 million paired-end reads. The inferred heterozygosity, based on GenomeScope, was between 1.675% and 1.680% for a k-mer size of 21, and between 1.250% and 1.253% for a k-mer size of 48. Genome size estimates at this point were between 306 Mbp (k-mer = 21) and 320 Mbp (k-mer = 48). Following assembly with Supernova, assembly v1.0 was approximately 388 Mbp in size and comprised of 44 273 scaffolds (5.91% N's).

Assembly v1.0 was then assessed for contamination with a preliminary search against the NCBI for bacterial homology (see below for more details). Several scaffolds with high amounts of bacterial sequence contamination were identified, indicating that further decontamination of the assembly was required. A decontamination pipeline was used to identify and remove a total of 3043 scaffolds, while those identified as potential examples of LGT were kept. From the remainder, an additional 4717 were identified as being identical duplicates and were removed. At this point, the resulting assembly was finalized and designated v1.5. This assembly is 355 Mbp in size, consisting of 36 513 scaffolds (6.29% N's). Quality and completeness of v1.5 using BUSCO indicated a completeness of 87.5% (65.6% single copy orthologues, 21.9% duplicated orthologues), while 7.1% orthologues were fragmented and 5.4% were missing ( $n = 1658$ ). Further assessment for potential duplicate scaffolds (with a search for equal or greater than 95% identity and allowing for edits of up to 100 bp) identified 4942 scaffolds that would be considered extraneous in this definition. Of these scaffolds, 49 contain duplicate BUSCOs, while 85 contain single or fragmented BUSCOs. Of the duplicate BUSCOs, the majority (307) exist as two copies, while some are in triplicate (24) and one was found four separate times (albeit on two scaffolds). The list of potential scaffolds with greater than 95% identity to other scaffolds within the assembly, as well as the identified BUSCOs on those scaffolds can be found in supplementary materials (S1.2).

Initially, the genome size of *N. tenuis* was estimated by flow cytometry to be 232 Mbp, with a confidence interval of 20 Mbp (see supplementary material S1.2 for more details). Further estimates via k-mer analysis of sequence data in GenomeScope indicated an expected genome size of 306 Mbp (k-mer = 21) or 320 Mbp (k-mer = 48). Both the flow cytometry and sequence data estimates are smaller than the 355 Mbp of the final assembly (v1.5). In total, the *N. tenuis* genome has 36 513 scaffolds, with the largest scaffold being 1.39 Mbp, though the majority of scaffolds are under 50 000 bp in size. The number of gaps per 100 kbp is 6292.10 (6.29% of the genome). When continuous fragments of N's greater than 10 bp are removed from the assembly, the assembly size reduces slightly to 332 562 929 bp (6.4% reduction). Further details on the assemblies can be found in Table 1.

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

Assembly version	Size (bp)	No. of scaffolds	Contig N50 (bp)	Scaffold N50 (bp)	Largest scaffold (bp)	Scaffold length $\geq$ 50 000 (bp)	No. of N's per 100 kbp (% of genome)	BUSCO score, complete% (single%, duplicate%)
1.0	387 724 797	44 273	12 954	27 195	1 392 896	138 138 956	5912.60 (5.91)	81.3 (60.6, 20.7)
1.5 (final assembly)	355 120 802	36 513	13 374	28 732	1 392 896	130 346 782	6292.10 (6.29)	87.5 (65.6, 21.9)

### Assessment of potential symbionts and LGT candidates

**Potential symbionts.** The initial assembly (v1.0) was decontaminated using two bacterial decontamination pipelines: the first pipeline broadly utilized BLASTn to identify scaffolds with high amounts of bacterial sequences against the NCBI nr database, while the second pipeline is more specified and uses BLASTn against a list of known contaminants and symbionts (S1.4) and is adapted from previous work (Wheeler *et al.*, 2013). The first decontamination pipeline identified and removed 1 443 scaffolds with high bacterial content, and the second decontamination pipeline identified and removed an additional 1801 scaffolds alongside potential LGT events. All removed scaffolds are available in S1.4. The hits from the second pipeline (1801 scaffolds) were used to create a list of potential contaminants or symbionts of this particular *N. tenuis* individual used for whole genome sequencing according to probable family, base pair content and number of scaffolds affected (Table 2). The majority of these scaffolds (1470) are under 5 kbp in length, with an additional 61 scaffolds falling between 5 and 10 kbp. The 10 largest scaffolds contained matches to bacteria from the following genera: *Pantoea* and relatives (three of 5 617 472 bp, 205 621 bp and 131 905 bp), *Sodalis* (326 101 bp), *Erwinia* (254 660 bp; 220 307 bp; 154 581 bp) and *Citrobacter* (239 269 bp; 190 839 bp). We emphasize that these “calls” are very preliminary, as they are based on the most frequent hits in the

bacterial matches in each scaffold, rather than comprehensive gene annotations. Nevertheless, they do indicate a range of bacterial types associated with *N. tenuis*, and the scaffold assemblies are likely to contain some complete or near complete bacterial genomes of interest.

Sorting scaffolds across the range of bacterial families matches provides some substantial representation: *Erwiniaceae* (2 722 260 bp), *Enterobacteriaceae* (745 106 bp), *Pectobacteriaceae* (240 913 bp), *Rickettsiaceae* (204 934 bp), *Anaplasmataceae* (137 109 bp) and *Yersiniaceae* (136 728 bp) (Table 2). In addition to known symbiont *Wolbachia*, previously established via PCR, known symbiont *Rickettsia* is also present in the results (137 109 bp) (Table 2). Multiple genera of bacteria can be found on a single scaffold, likely due to misassembly. The full list of bacterial scaffolds and multiFASTA file is available with details in supplementary material (S1.4).

**LGT candidates.** We continued with our detection of potential LGT events by further assessing a handful of strong candidates. Two of these regions occur on scaffolds 22 012 and 22 013, which are of similar length (22 634 and 22 957 bp, respectively) and are highly similar on a nucleotide level. Scaffold 22 013 appears to have additional nucleotides on each flanking side, with some indels and SNPs between the two scaffolds. The putative LGT region in question is belonging to or gained from a *Sodalis* species, coding for phenazine biosynthesis protein PhzF (OIV46256.1). This region around the putative LGT also showed transcriptional activity, and is flanked by conserved insect genes, most immediately *Rab19* (NP\_523970.1) on one side and an uncharacterized protein, DmeI\_CG32112 (NP\_729820.2), on the other side. Two additional LGTs were found in the current assembly that match *Rickettsia* sequences (scaffolds 4712 and 27 281), which contain a segment of the rickettsial genes *elongation factor G* and *AAA family ATPase* genes, respectively. One corresponds to a gene model, while the other does not, and there is no evidence of expression for either in the current male, female and mixed sex juvenile RNA sequencing data. More information on these candidate regions can be found in S1.4.

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

To obtain a comprehensive set of transcripts for *N. tenuis*, three separate libraries of multiple individuals were

**Table 2.** Families of potential symbionts or contaminants as determined by second round of decontamination based on known contaminants and symbionts against *Nesidiocoris tenuis* assembly v1.0

Bacteria family	Total amount in genome (bp)	Number of scaffolds affected
Erwiniaceae	2 722 260	763
Enterobacteriaceae	745 106	223
Pectobacteriaceae	240 913	51
Rickettsiaceae	204 934	50
Anaplasmataceae	137 109	47
Yersiniaceae	136 728	80
Pseudoalteromonadaceae	43 206	120
Morganellaceae	29 759	231
Burkholderiaceae	11 195	24
Thermoproteaceae	3 420	3
Paenibacillaceae	3 070	8
Others	56 287	201
Total	4 333 987	1801

Note: 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 S1.3. Full list of hits available in supplementary material S1.4.

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

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

prepared—males, females and juveniles from different stages of mixed sex. More than 77 million 150 bp pair-end reads were generated. Filtering the reads for quality led to a slightly reduced total of 76 711 096 paired reads (male: 28413231 paired reads; female: 24075901 paired reads; juvenile: 24221964 paired reads) to be used for evidence-based gene finding. The mapping and assembling of reads of the three individual samples as well as the pooled reads resulted in four transcriptomes: male, female, juvenile and the combined transcriptome.

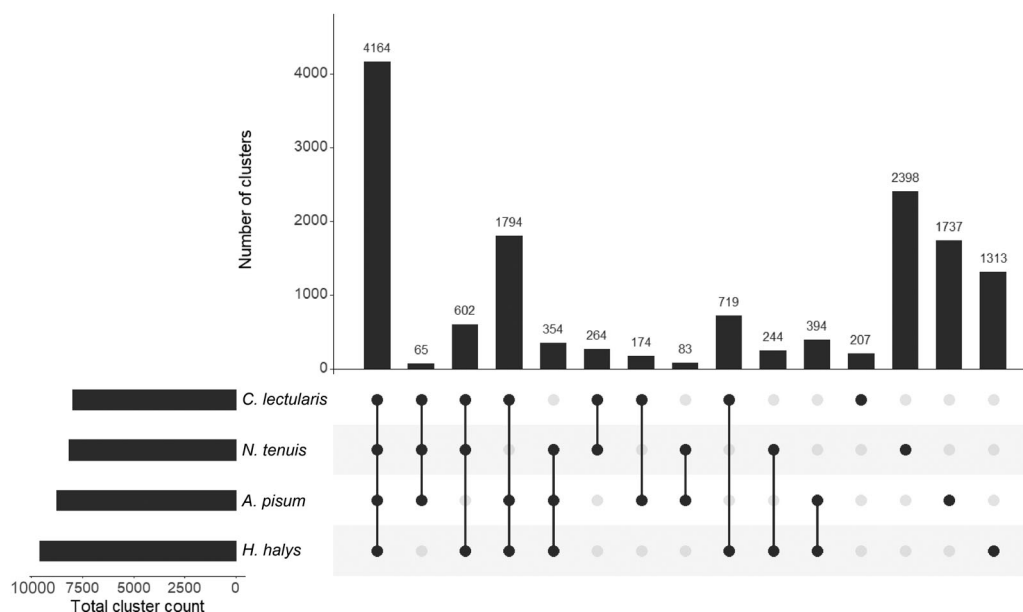
The male, female, juvenile and combined annotations from the evidence-based gene finding was used alongside homology-based findings and *ab initio* annotations in a weighted model, resulting in complete annotations for the assembly. When gene name assignment via the SwissProt database resulted in “no hit,” tracks are named “No\_blast\_hit.” This occurred in 1556 mRNA tracks and represents approximately 6% of the official gene set. The majority of tracks were annotated with reference to SwissProt or GenBank accession number of the top BLASTp hit.

CodingQuarry predicted 56 309 genes from the mapped transcript evidence, while *ab initio* gene finding using GlimmerHMM resulted in 39 888 genes and homology-based

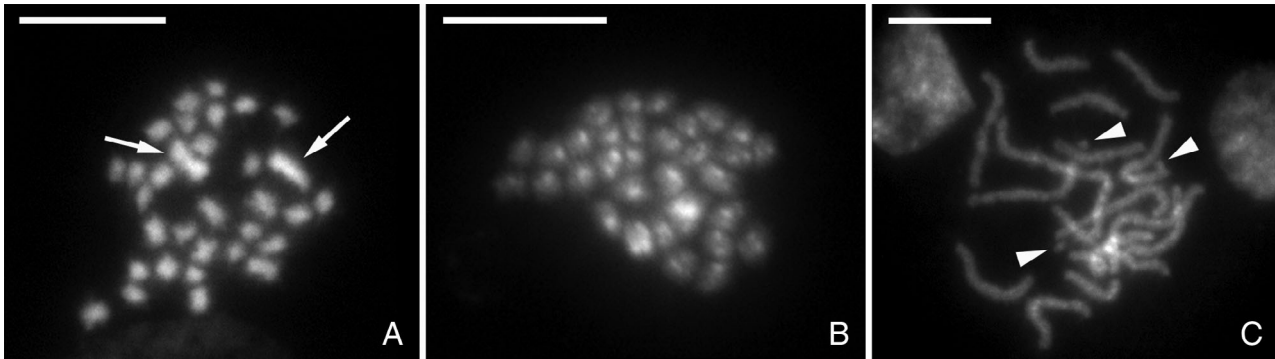
gene finding with GeMoMa resulted in 6028 genes. The complete gene set for *N. tenuis* was created using EVIDENCEModeler, where a weighted model using all three inputs resulted in a complete gene set of 24 688 genes.

#### Orthologue cluster analysis

The complete gene set of *N. tenuis* was compared to those of three additional species: the bed bug *Cimex lectularis* L. (Hemiptera: Cimicidae), the brown marmorated stinkbug *Halyomorpha halys* (Stål) (Hemiptera: Pentatomidae) and the pea aphid *Acyrtosiphon pisum* (Harris) (Hemiptera: Aphididae) using OrthoVenn2 (Xu et al., 2019a). The orthologue analysis summary is presented in Table 3 and visualized in Fig. 1. *N. tenuis* has a similar number of clusters (8174) as compared to *C. lectularis*, *H. halys* and *A. pisum* (7989; 9584 and 8765, respectively). In total 14 512 clusters are assigned, 12964 of which are orthologous clusters (contains at least two species), and the remaining 1548 are single-copy gene clusters (when there is a single protein for each species within a cluster). There are 9136 singletons (proteins that do not form clusters) in *N. tenuis*, 3573 in *C. lectularis*, 2170 in *H. halys* and 7298 in *A. pisum*. Just



**Figure 1.** Orthologue cluster analysis of *Nesidiocoris tenuis* with three other hemipterans (*Cimex lectularis*, *Halyomorpha halys* and *Acyrtosiphon pisum*). Numbers indicate the number of orthologue clusters in each grouping, with shared clusters indicated by connected nodes.



**Figure 2.** Cytogenetic analysis of *Nesidiocoris tenuis* karyotype. Chromosomes were counterstained with DAPI (grey). (A) Female mitotic metaphase consisting of 32 chromosomes ( $2n = 32$ ) with two large chromosomes indicated (arrows). (B) Male mitotic metaphase consisting of 32 chromosomes ( $2n = 32$ ). (C) Female pachytene nucleus with B chromosomes (arrowheads). Scale bar = 10  $\mu\text{m}$ .

over half of the orthologues cluster with *N. tenuis*, where 6338 clusters are outside of *N. tenuis* as compared to the 8174 clusters within *N. tenuis*. The final protein set from *N. tenuis* used in this analysis is available, see S1.5.

#### Karyotype analysis

Karyotype analysis revealed  $2n = 32$  chromosomes in both females and males (Fig. 2A, B). All chromosomes are relatively small with one larger pair of submetacentric chromosomes in females (Fig. 2A). In males (Fig. 2B), we were unable to obtain mitotic chromosomes of reasonable quality as in females, and therefore, we were unable to clearly identify these larger chromosomes. Screening of multiple nuclei showed sporadic deviations of the karyotype in some individuals. This was the result of supernumerary chromosomes (B chromosomes), which were clearly visible in (meiotic) pachytene stage as distinctly smaller chromosomes (Fig. 2C, three B chromosomes).

#### Analysis and localization of 18S rDNA

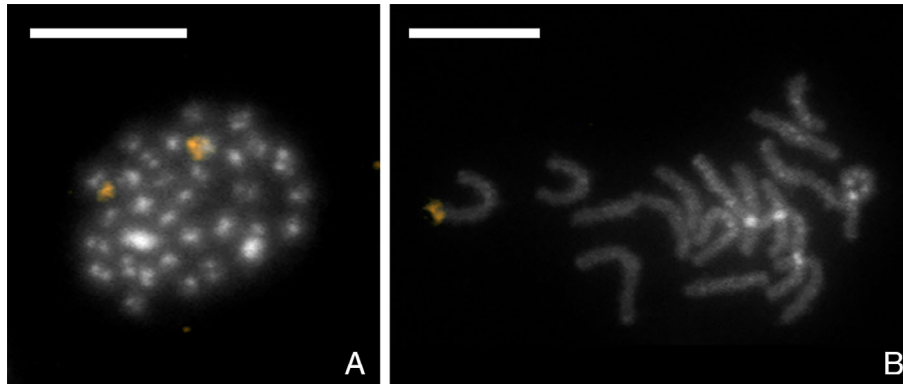
The 18S rDNA gene is often used as a cytogenetic marker in comparative evolutionary studies due to its ease of visualization on the chromosomes caused by high copy number and cluster organization in animal (Sochorová *et al.*, 2018) and plant (Gomez-Rodriguez *et al.*, 2013) genomes. The published partial 18S sequence of *N. tenuis* (GU194646.1) and the 18S sequence identified in this study were compared to each other revealing some differences between the sequences. The published sequence consists of two fragments of 869 and 739 bp, which are, respectively, 99.7% and 94.2% identical to our identified partial 18S sequence. Interestingly, the second half of our isolated 18S sequence has a higher identity with a *Macrolophus* sp. partial 18S sequence (EU683153.1), ie 97.8%, than to the previously published *N. tenuis* sequence. A BLAST search against the *N. tenuis* genome with either of the *N. tenuis* 18S sequences resulted in four

gene copies in both cases, each located on a different scaffold. However, RepeatExplorer analysis estimated 98 18S rDNA copies with the obtained genome size of 355 Mbp. Using fluorescence *in situ* hybridization (FISH) with the 18S rDNA probe we finally showed that the major rDNA forms a single cluster located terminally on a pair of homologous chromosomes (Fig. 3).

#### Identification of sex chromosomes

The common sex chromosome constitution in Miridae is the male-heterogametic XX/XY system. To identify the sex chromosome constitution and estimate sex chromosome differentiation in *N. tenuis* we employed genomic *in situ* hybridization (GISH) and comparative genomic hybridization (CGH) experiments. The GISH results clearly revealed a single chromosome densely labelled by the male-derived probe, caused by male-enriched repetitive DNA and/or male-specific sequences, which is typical for the Y chromosome (Fig. 4). In addition, the nucleolus organizer region (NOR; including 18S rDNA) was observed as well, as is often the case in GISH experiments due to the presence of highly repetitive sequences in the rDNA cluster. The NOR is clearly located terminally on a pair of autosomes, corroborating our 18S rDNA FISH results.

To further study the differentiation of the sex chromosomes, we carried out CGH experiments on chromosome preparations of both sexes (Fig. 5). All chromosomes were labelled evenly by the female and male probes with the exception of the largest chromosome pair. Both sex chromosomes were highlighted with DAPI (Fig. 5A, E), indicating that they are both A-T rich and largely composed of heterochromatin. In females, the largest chromosome pair was labelled more by the female probe than the male probe indicating that these chromosomes contain sequences with higher copy numbers in females, and are thus the X chromosomes, as seen in Fig. 5A–D. In male meiotic nuclei (Fig. 5E–H), two types of nuclei can be discerned, where



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

the largest chromosome was labelled more by either the female probe or the male probe corresponding to the X and Y chromosome, respectively, whereas the autosomes were labelled equally by both probes.

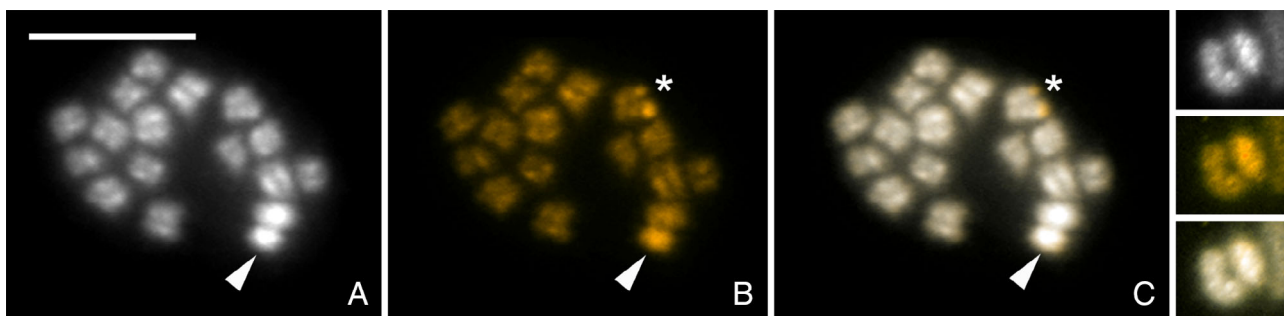
#### Identification and mapping of abundant repeats

RepeatExplorer software was used on reads with barcodes removed to identify the most abundant repeats in the genome of *N. tenuis* (results available in supplementary materials, see Table S1.6.1). The most abundant repeat, Nt\_rep1, makes up approximately 3% of the genome estimated by RepeatExplorer. Analysis on the assembled genome, using a coverage cut-off value of 70%, reveals that Nt\_rep1 is present on 3190 scaffolds (8.737% of the assembled scaffolds), with a maximum of 17 copies on a single scaffold. According to the assembled genome, Nt\_rep1 makes up approximately 0.8% of the entire genome (Fig. S1.6.2). We subsequently mapped Nt\_rep1 to the chromosomes of *N. tenuis* using FISH. The repeat

is located on most chromosomes and is accumulated in sub-telomeric regions (Fig. 6). Additional signals were identified on the X chromosome indicating a higher number of this repeat (Fig. 6A–C). This increase in frequency is specific to the X chromosome and is not found on the Y chromosome of *N. tenuis* (Fig. 6D–F).

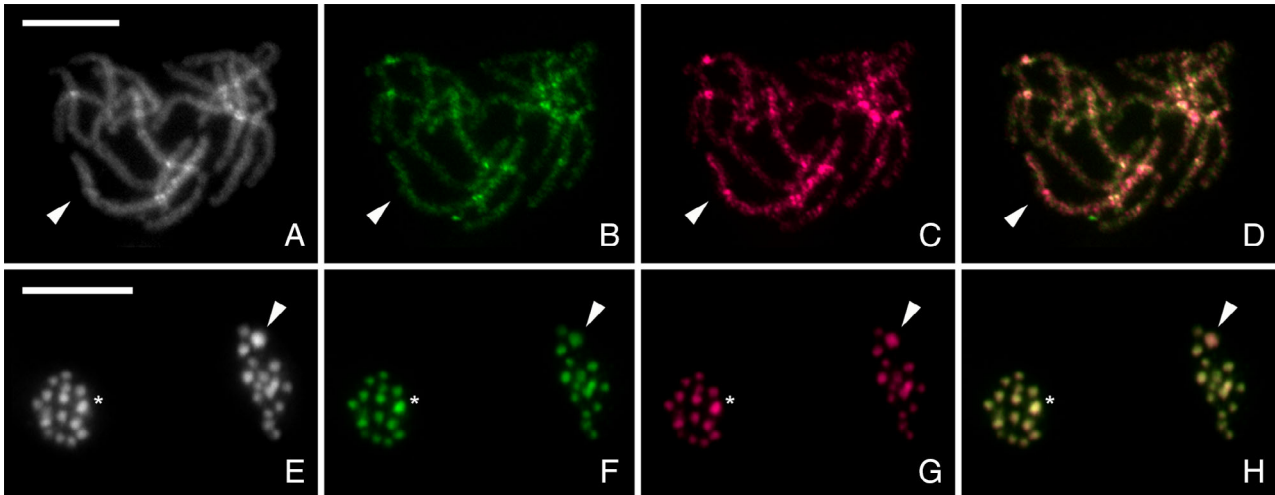
#### Testing of candidate telomere motifs

Analysis of the raw sequencing data and the assembled genome both revealed low numbers of the insect telomere motif (TTAGG)<sub>n</sub> (Frydrychová *et al.*, 2004) in *N. tenuis*, ie approximately 98 repeats per haploid genome. This translates into approximately three copies of the repeat per chromosome end, much lower than expected for a telomeric motif. These low copy numbers were additionally confirmed using Southern dot blot (Fig. S1.6.1). Other candidate telomere motifs previously identified by Tandem Repeat Finder (TRF) analysis, (TATGG)<sub>n</sub>, (TTGGG)<sub>n</sub> and (TCAGG)<sub>n</sub>, were examined by FISH for their distribution in the genome. They

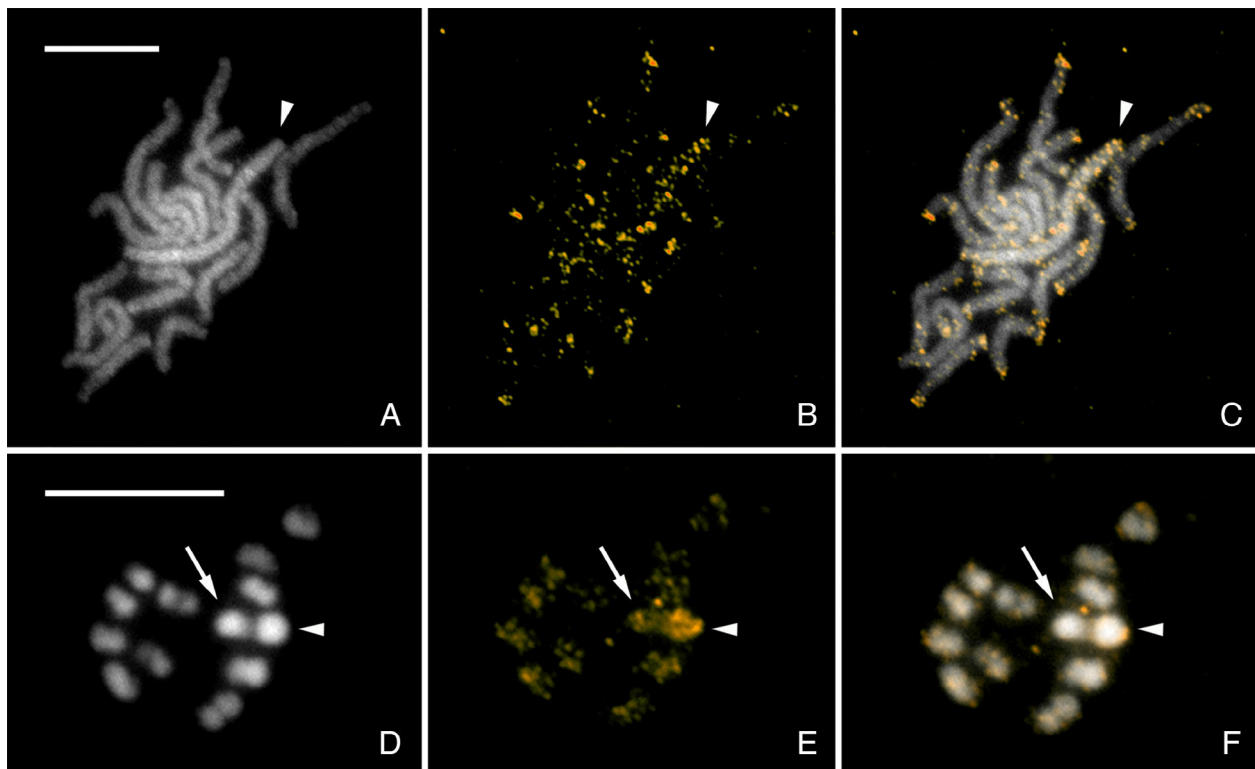


**Figure 4.** Genomic *in situ* hybridization (GISH) on male chromosomal preparation of *Nesidiocoris tenuis*. Panel (A) shows DAPI counterstaining (grey), panel (B) hybridisation signals of the male derived genomic probe labelled by Cy3 (gold) together with competitor generated from unlabelled female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

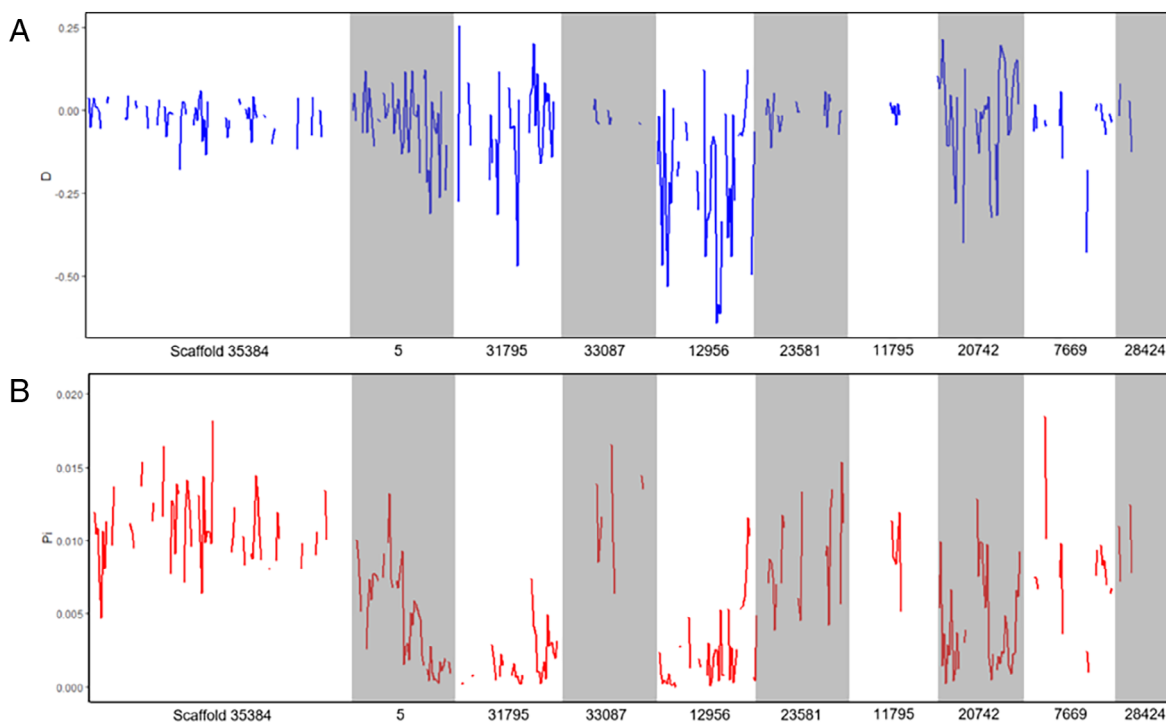




**Figure 5.** Comparative genomic hybridization (CGH) on female pachytene (A–D) and male meiotic metaphase II (E–H) chromosomes of *Nesidiocoris tenuis*. Panels (A, E) show chromosomes counterstained with DAPI (grey), panels (B, F) hybridization signals of the male derived genomic probe labelled by fluorescein (green), panels (C, G) hybridization signals of the female derived genomic probe labelled by Cy3 (magenta) and panels (D, H) merged images. (C, D) Note that the X chromosome bivalent (arrowhead) in female pachytene complement was highlighted more by female probe compared to the autosomal bivalents; (B, D) male probe labelled all chromosomes equally. (H) Two sister nuclei in meiotic metaphase II showed equal hybridization patterns of both probes on autosomes; in one of the forming nuclei, the X chromosome (arrowhead) was highlighted by female derived genomic probe (G, H) and in the second nucleus the Y chromosome (asterisk) was strongly highlighted by male derived genomic 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 most heterochromatic elements in the nucleus. Scale bar = 10  $\mu\text{m}$ . Note: Alternate coloration available in supplementary material (S1.6.3).



**Figure 6.** Fluorescence *in situ* hybridization with *Nt\_rep1* probe labelled by biotin (gold) on female (A–C) and male (D–F) chromosomes of *Nesidiocoris tenuis* counterstained with DAPI (grey). (A–C) Female pachytene chromosomes; *Nt\_rep1* probe highlighted the pair of X chromosomes (arrowhead). (B, C) Note that terminal regions of all bivalents were also labelled by probe, probably due to the presence of this sequence in sub-telomeric regions. (D–F) Incomplete male nucleus in meiotic metaphase I; probe highlighted the X chromosome (arrowhead) more compared to autosomes and Y chromosome (arrow). (B–F) Strong hybridization signals on X chromosomes in both sexes were caused by enrichment of *Nt\_rep1* sequence on the X chromosomes. Scale bar = 10  $\mu\text{m}$ .



**Figure 7.** Genetic diversity of Koppert Biological Systems commercial *Nesidiocoris tenuis* population according to Tajima's  $D$  (a), and nucleotide diversity,  $\pi$  (b). Scaffolds are ordered according to size, which each name beneath on the x-axis. [Colour figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)].

were found scattered throughout the genome but lacked a clear accumulation at the terminal regions of the chromosomes (not shown). Therefore, these sequences can also be excluded as telomeric motifs in *N. tenuis*.

#### Pooled sequencing analysis

We repurposed previously generated whole genome sequencing of 10 females from the KBS population, we were able to estimate genetic diversity of the commercial population via a pool-seq population analysis. Read coverage was randomly subsampled to 10X coverage (18 000 000 reads). In addition, we used a modified

v1.5 *N. tenuis* genome with scaffolds of less than 10 000 bp removed. This was to ensure that window sliding was not being inflated on scaffolds smaller than the window size. This reduced the genome from 36 513 scaffolds to 7076, however, the reduced genome still contained 72.23% of the genome in terms of size (256 487 768 bp).

Three runs of PoPoolation were performed with varied window size, step size and the masking of indel regions. The default setting, window size and step size of 10 000 bp, yielded similar results as the adjusted window size and step size of 5000 bp, while differences were apparent when indel regions were masked. As such, results of window size and step size 10 000 bp with indel

**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

Scaffold	Size (bp)	Windows (10 kbp)	Number of sufficiently covered windows	Average coverage of sufficiently covered windows	Average $\pi$ across scaffold	Average Tajima's $D$ across scaffold
35 384	1 392 896	140	70	0.68	0.010682	-0.02091
5	613 435	62	48	0.76	0.004932	-0.03869
31 795	577 751	58	38	0.74	0.001753	-0.04655
33 087	539 928	54	17	0.66	0.012571	-0.00458
12 956	519 254	52	39	0.80	0.002541	-0.21012
23 581	513 368	52	28	0.70	0.008417	-0.04658
11 795	508 155	51	14	0.68	0.011553	-0.01147
20 742	504 856	51	39	0.78	0.004165	-0.01963
7 669	488 533	49	24	0.67	0.007856	-0.0551
28 424	477 580	48	5	0.66	0.009563	-0.00998
Total	256 487 768	28 833	5 913	0.70	0.0080	-0.0355

regions masked are reported here (other results available in S1.7). The variance sliding programme created 28 833 windows of 10 000 bp with mapped reads, of which 5913 were sufficiently covered with reads to calculate values per window (coverage  $\geq 0.60$ ). Genome-wide, the nucleotide diversity ( $\pi$ ) is 0.0080 and *Tajima's D* is  $-0.0355$ . Fig. 7 shows the *Tajima's D* (7a) and  $\pi$  (7b) for the 10 largest scaffolds, all containing gene annotations, arranged in order of size. These 10 scaffolds represent approximately 1.7% of the genome (6 135 756 bp), and varied in terms of window coverage (from no coverage to full coverage) as well as both *Tajima's D* and  $\pi$ . These 10 scaffolds are a snapshot of the whole genome, summarized in Table 4, whereas genome-wide results can be found in S1.7.

## Discussion

### Assembly and annotation

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

While the final assembly offers a relatively high BUSCO score, with the final decontaminated assembly (v1.5) having a completeness of 87.5% of the insect\_odb9 orthologue dataset, the assembled genome size is approximately 120 Mbp larger than was expected based on flow cytometry data. As we do not know the B chromosome status of the individuals used for either gDNA or RNA sequencing or their genome size, part of the difference could be explained by the presence of these supernumerary chromosomes. However, most of the genome inflation is likely due to a combination of residual contamination along with faulty haploidization of the assembly. While the former is a relatively common feature of genome assemblies, the latter is a risk of using a linked-read method with a heterozygous sample. The final stage of the Supernova assembler pipeline takes artificially constructed phased genome regions and reduced it to a haploid structure. Here, we could infer faulty or incomplete haploidization of the assembly by the presence of these highly identical scaffolds (between 95% and 100% identity) that often have consecutive naming conventions. Returning to the BUSCO

results, this assembly returns a score of 21.9% duplicated BUSCOs (332 in total), with these duplicate BUSCOs often occurring on consecutive scaffolds. However, not all of the duplicate BUSCOs follow this convention or pattern of occurring on highly similar scaffolds, making it difficult to remove scaffolds as duplicate or potential duplicates. Therefore, we suggest that the assembly presented here can best be improved in terms of accuracy and contiguity with long reads, either from an inbred sample or using higher accuracy sequencing methods such as PacBio's High-Fidelity (HiFi) technology, in the future.

Annotation via evidence-based, homology-based and *ab initio* models resulted in 24 668 genes. Compared to other assemblies within the hemipteran order, such as *C. lectularis*, with a genome size of 650 Mbp and 12 699 genes (Thomas *et al.*, 2020) or *A. pisum*, with a draft genome size of 464 Mbp and 36 195 genes (Richards *et al.*, 2010), *N. tenuis* sits, in the middle in terms of genome size and number of genes.

Comparing the current gene set of *N. tenuis* to other hemipterans, the clustering identified considerable overlap as 71% of the clusters that are found in *N. tenuis* were shared between the other species in the comparative analysis. Despite being more closely related to *C. lectularis* in terms of phylogeny, in terms of lifestyle, *N. tenuis* is far more similar to *A. pisum* and *H. halys*, and this is likely reflected in absolute number of proteins and clusters shared between the four species. The 1548 single-copy gene clusters identified within the analysis would also be a place to continue further comparisons between these hemipterans. The remaining 29% of clusters, as well as the singleton proteins, are indications for proteins unique to either *N. tenuis* or Miridae in general. Through the OrthoVenn2 website (Xu *et al.*, 2019a), the analysis performed here can be easily replicated, altered with other species of interest and even improved upon if the complete gene sets are updated or with a newer software version. In our iteration, the 24 668 proteins of *N. tenuis* group into 8174 clusters. 2398 clusters are unique to *N. tenuis*, however, some of these genes have a relatively strong homology to genes of one of the other species used in the analysis and could be incorrectly flagged as being unique. Reasons could be poor gene annotation quality resulting in a poor *in silico* protein translation, or too stringent clustering settings. Regardless, these 2398 clusters may be of interest to researchers working on zoophytophagy, the negative effects of *N. tenuis* on tomato as compared to other mirids, or broader questions such as phylogeny of the Hemiptera.

### Characterizing the genome

Every sequence and assembly strategy has benefits and drawbacks, and the 10X Genomics linked-read strategy is no exception. The technique requires only few nanograms

of DNA for library preparation, which allowed us to use a single individual and removed the need for inbreeding to reduce variation in the sequencing population. However, using a single individual from a closed and proprietary rearing process presented other challenges. These challenges were threefold: we had to deal with bacterial contamination as antibiotic treatment is not possible, we had to ensure that the single individual-derived assembly reflects reality in terms of genes present and structure and we had to ensure that a single female-derived assembly is applicable for population-level analyses.

Contamination of genomes is a constant concern, and sequencing strategies should attempt to address the risks in the best way possible to deliver reliable genomes (Ekblom and Wolf, 2014). Equally so is the desire for inbred strains if multiple individuals are required to reach the micrograms of DNA necessary for NGS platforms. The inability to remove symbionts and microbiota using antibiotics administered to a few successive generations as well as the difficulty or inability to inbreed a strain is not restricted to *N. tenuis*. The sequencing strategy chosen for the mountain pine beetle, *Dendroctonus ponderosae* (Hopkins) (Coleoptera: Curculionidae), relied on assuming the relatedness of several individuals as well as isolating the gut during the extraction process, and still additional post-assembly decontamination was required (Keeling et al., 2013). A linked-read strategy with low input requirements, such as the 10X Genomics library preparation that was chosen here, negates the need for a pool of inbred samples or controlling for relatedness. However, another potential benefit of controlled rearing such as those used in inbreeding (as opposed to be limited to wild-caught specimens, eg) is the ability to treat with antibiotics for multiple generations. Without the ability to do so, sequencing and assembly strategies rely heavily on post-sequencing decontamination strategies (both pre- and post-assembly are possible). That such post-assembly filtering strategies as used here for *N. tenuis* can be successful as shown by a less contaminated assembly and by the identification of potential LGTs. However, as identified in the analysis of our assembly in relation to the BUSCO duplications and near duplications of scaffolds, as well as our repetitive sequence analysis for cytogenetics, the Supernova assembler may have difficulty in resolving the genome for haploidization when encountering heterozygosity in addition to losing biologically present repetitive regions (ie true repeats within the genome that are reduced in the assembly process). Furthermore, this may have implications on the final gene set, and requires further investigation by using longer reads for refining the assembly as well as manual annotation.

#### *Beyond the genome: potential symbionts and LGT events*

The list of potential symbionts or pathogens generated in Table 2 represent both insect and plant pathogens, as well

as potential environmental contaminants. In addition to the positive test for *Wolbachia* (family: Anaplasmataceae) in the KBS population used here, *N. tenuis* is known to potentially harbour *Rickettsia* (family: Rickettsiaceae) as an endosymbiont in addition to *Wolbachia* (Caspi-Fluger et al., 2014). *Rickettsia* genome sizes can range from 0.8 to 2.3 Mbp, also reflecting variation in levels of reductive evolution (Sachman-Ruiz and Quiroz-Castañeda, 2018). However, the total scaffold length identified here as potentially being from *Rickettsia* falls below this range, likely indicating incomplete recovery of the genome from the insect sequencing. The potential symbionts revealed included not only *Wolbachia* and *Rickettsia*, but also other known insect symbionts, in addition to the usual lab contamination suspects.

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

As for the family Yersiniaceae, *Serratia marcescens* is both a common Gram-negative human-borne pathogen and a causal agent of cucurbit yellow vine disease (CYVD) (Abreo and Altier, 2019; Bruton et al., 2007). It is worth noting that in cases of CYVD, the transmission of *Serratia marcescens* from its vector the squash bug, *Anasa tristis* (De Geer) (Hemiptera: Coreidae), to host crops is via the phloem. Other *Serratia* spp. have been identified as insect symbionts previously, as have other potential symbionts found in the contaminated scaffolds, such as *Cedecea* spp. (Jang and Nishijima, 1990). The putative presence of *Dickeya* (family: Pectobacteriaceae) is an interesting finding, as *Dickeya dadantii* has been established as a pathogen of *A. pisum*, while the pea aphid itself is a potential vector for the bacterium with regard to

plants (Costechareyre *et al.*, 2012). *Dickeya* spp. cause soft rot in various crops, including tomato. In a similar vein, the identification of members of the family Burkholderiaceae is noteworthy, as some members of this family, such as *Ralstonia solanacearum*, are soil-borne pathogens that causes wilt in several crop plants, including tomato (Lowe-Power *et al.*, 2018).

All of these described associations are preliminary, as follow-up analyses against the entire NCBI database and proper bacterial gene annotation are lacking. Nevertheless, these putative bacterial associations of *N. tenuis*, their distribution within the insect, and their possible biological significance, warrants further investigation. It is important to note that some of these bacterial “contaminants” may actually represent large LGTs, which can be confirmed by identifying flanking sequences (eg using long-read technologies) and/or *in situ* chromosome hybridization analyses, such as done for the large LGT in *Drosophila ananassae* (Doleschall) (Diptera: Drosophilidae) (Dunning Hotopp *et al.*, 2007). However, it is good to note that Miridae is not as well-represented in gene-finding and gene-annotating programmes, so it is difficult to parse LGT from contamination from yet-to-be annotated insect-derived genes that are specific to mirids in particular. Therefore, more research into the symbionts of *N. tenuis* via metagenomics would certainly shed some light on true symbionts (or pathogens) versus true contaminants, with potential implications for biological control and related research.

One of the LGT candidate genes that were detected following manual curation corresponds to phenazine biosynthesis protein PhzF (OIV46256.1), with the likely microbial source being a *Sodalis* species. Phenazines are heterocyclic metabolites with “antibiotic, antitumor, and antiparasitic activity,” but are also toxic when excreted by bacteria (Blankenfeldt *et al.*, 2004). This LGT region exhibits gene expression and is flanked by conserved insect genes, providing further support for it being a legitimate LGT, though further research into this region will be necessary to confirm this. The gene occurs on two different scaffolds, 22 012 and 22 013, which are highly similar to each other in some regions at the nucleotide level. These could represent homologous regions that differ sufficiently to assemble as different scaffolds, or alternatively a duplication in two different regions of the genome. Future work should focus on its expression patterns in different tissues (eg salivary glands, in interest of *PhzF*) and potential functional role in *N. tenuis*.

#### *Beyond the genome: cytogenetics*

We determined the karyotype of *N. tenuis* to be  $2n = 32$  (30 + XY in males) chromosomes, which is the second most common chromosome number in the family Miridae (Kuznetsova *et al.*, 2011). In addition, we have shown that

*N. tenuis* has an XX/XY sex chromosome constitution, with the sex chromosomes being the largest elements in the karyotype. This is different from the closely related *Macrolophus costalis* (Fieber) (Hemiptera: Miridae) ( $2n = 24 + X_1X_2Y$ ), and *M. pygmaeus* (Rambur) ( $2n = 26 + XY$ ) where two pairs of autosomes are larger than the sex chromosomes, yet similar to *M. melanotoma* (Costa), which only differs from *N. tenuis* in the number of autosomes,  $2n = 32 + XY$  (Jauset *et al.*, 2015). As we sequenced a single female, sequence information of the Y chromosome is missing from our genome assembly. While analysing the *N. tenuis* karyotype we discovered the sporadic presence of B chromosomes in the KBS population. B chromosomes are supernumerary chromosomes that are dispensable to the organism, and are often present in only a subset of individuals from a population (Banaei-Moghaddam *et al.*, 2013). Supernumerary chromosomes are common in Heteroptera, yet only a few species of Miridae have been identified to carry supernumerary chromosomes (Grozeva *et al.*, 2011). Presence of B chromosomes in high numbers within an individual is often found to be detrimental, though in lower numbers they are often considered neutral or, in some cases, beneficial (Jones and Rees, 1982; Camacho *et al.*, 2000). The abundance and origin of B chromosomes in *N. tenuis* biological control populations are currently unknown but determining their potential effects on fitness-relevant traits might reveal beneficial information for the optimization of mass-reared populations.

The hemizygous sex chromosomes of most organisms have a high content of repetitive DNA, consisting of multiple different repetitive sequences that are found less frequently on autosomes (Traut *et al.*, 1999; Charlesworth and Charlesworth, 2000). Therefore, the use of cytogenetic techniques, such as CGH and GISH, in the identification of hemizygous sex chromosomes is a powerful tool and is well established in different groups of organisms, eg *Lepidoptera* (Dalíková *et al.*, 2017a; Zrzavá *et al.*, 2018; Carabajal Paladino *et al.* 2019), Orthoptera (Jetybayev *et al.*, 2017), fish (Sember *et al.*, 2018) and frogs (Gatto *et al.*, 2018). However, to our knowledge, this is the first time these techniques have been used in the family Miridae. The X and Y chromosome of *N. tenuis* are similar in size, with the X chromosome being slightly bigger, and are difficult to distinguish from each other based solely on their appearance without special probing. Our CGH and GISH results showed relatively weak hybridization signals of genomic probes on the sex chromosomes compared to other species indicating little differentiation of sequence content between the X and Y chromosomes, and/or between the sex chromosomes and the autosomes. Though the hybridization signals are relatively weak, not only the Y chromosome but also the homogametic sex chromosome, the X chromosome, is distinguishable in the CGH results, which shows X-enriched or X-specific

repetitive DNA, similar to what was found on the Z chromosome in *Abraxas* spp. (Zrzavá *et al.*, 2018). Mapping of the most abundant repeat in the genome revealed that one such X-enriched repeat is Nt\_rep1, confirming the outcomes of our CGH results.

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

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

Screening of the genome and Southern blot assay suggest the absence of the ancestral insect telomere motif, (TTAGG)<sub>n</sub>, in *N. tenuis*, as the case in other species from the family Miridae (Kuznetsova *et al.*, 2011; Grozeva *et al.*, 2019). The telomeric motif was present in our Tandem Repeat Finder results, but in much lower numbers than expected for telomeric sequences. Additional attempts of identifying the telomeric repeat motif did not resolve this question. Three additional repeats we identified in the *N. tenuis* genome were tested via FISH, ie (TATGG)<sub>n</sub>, (TTGGG)<sub>n</sub>, and (TCAGG)<sub>n</sub>, but did not localize near the ends of the chromosomes. Notably though, mapping the most abundant repeat in the genome, Nt\_rep1,

did reveal accumulation in the sub-telomeric regions of chromosomes (Fig. 6). Therefore, our approach to identify potential telomere motifs, though presently unsuccessful, would presumably be effective if more repeats would be screened. In addition, a similar approach was used by Pita *et al.* (2016) in *T. infestans*, where the insect telomere motif, (TTAGG)<sub>n</sub>, was successfully identified from the raw sequencing data (Pita *et al.*, 2016). It must be noted, however, that the telomeres of *N. tenuis* might consist of different types of repeats other than short tandem repeats (as found in, eg, *Drosophila*; Traverse and Pardue, 1988), which would not be identified using Tandem Repeat Finder (Traverse and Pardue, 1988). Therefore, the identity (or even presence) of the telomeric repeat in *N. tenuis*, and by extension Miridae, remains unknown.

#### *Beyond the genome: population genomics*

Pooled sequence data of ten females from the KBS population were compared against the genome and provide interesting population-level effects. The overall negative Tajima's *D* would seem to indicate an abundance of rare alleles and is possible evidence of selective sweeps or population expansion, as seen in some populations of *Drosophila serrata* (Malloch) (Reddiex *et al.*, 2018), however, this generally results in more negative values (near  $-1$  or  $-2$ ). While overall negative, the absolute value of *D* in our results is small in comparison (total range from  $-0.89$  to  $0.56$ ). To best assess the state of the commercial population, monitoring the genetic variation over time would indicate if the population is undergoing an expansion after a bottleneck ( $D < 0$ ) or contracting ( $D > 0$ ), whereas when  $D = 0$ , we assume no selection. We can then assume that there is no selection currently at play in the commercial population. The few studies that have looked at genetic diversity within biological control populations have primarily been reduced representation analyses, such by genotyping with microsatellites (Paspati *et al.*, 2019). Here, a pool-seq approach offers a genome-wide look at the population and can give indications of the genetic diversity of the population; this could be a useful tool for monitoring population levels efficiently and determining, which regions of the genome are under selection in a biological control context.

Both genetic diversity values calculated here can also be used in population comparisons between the biological control population and wild populations. For instance, Xun *et al.* used mitochondrial and nuclear barcoding regions to haplotype 516 individuals across 37 populations into two regional groups, southwest China (SWC) and other regions in China (OC) (Xun *et al.*, 2016).  $\pi$  was 0.0048 (SWC) and 0.904 (OC), while *D* was  $-0.112$  (SWC) and  $-1.998$  (OC). It was concluded that the SWC population was stable while, similar to the KBS population here, the OC population was undergoing sudden population

explosion. Pooled sequencing could be a useful tool for comparing wild Mediterranean populations to the commercial population to determine disparities in genetic variation as well as to understand the dynamics of the wild populations.

There is a concern in using PoPoolation in this context: are 10 individual females sufficient for determining population variation? Here we used existing population sequence data (previously generated for another genome project, now available under accession ERX3938928) to better utilize resources, reduced to an appropriate coverage with masked indel regions. This enabled us to show population-level impacts at the very least, which can then pave the way for further studies, with better constructed sampling methods and sample sizes; the lack of perfect data should not preclude preliminary studies from being pursued. With this in mind, future population studies into wild *N. tenuis* populations can be interpolated with our population data as well as the genome for mapping and alignment.

## Conclusion

Reported here is the genome for *N. tenuis*, a mirid that is both used throughout the Mediterranean Basin as a biological control agent and reported as a greenhouse pest in other European countries. The assembled genome is 355 Mbp in length, composed of 36 513 scaffolds with an N50 of 28 732 bp. While the 10X linked-read strategy delivered a fairly complete genome, difficulties with contamination and potential faulty haploidization mean that there are some shortcomings of using linked-read sequencing on heterozygous samples, even if it is a single specimen. The goal of this project was to not only provide a genome, but also to highlight possible avenues of research now available with *N. tenuis*. A protein analysis has provided interesting prospects for mirid-specific proteins, while examples of potential LGT call for further inquiry. Putative symbionts were identified while filtering out contamination, creating a precursor for future metagenomic analysis. The cytogenetic analyses of *N. tenuis* here shed some light on mirid cytogenetics, such as the karyotype and sex determination system, but also solicits more questions. As for the commercial population, now that there is a baseline level of genetic variation documented through our pooled sequencing, what remains to be seen is how it compares to other populations, such as other commercial populations, wild or invasive populations. To this end, future exploration on these themes, among others, is now greatly facilitated with our release of this genome.

## Experimental procedures

### *Species origin and description*

Individuals of *N. tenuis* were received either from the commercial biological control stock at Koppert Biological Systems, S. L.

(Águilas, Murcia, Spain) (KBS) or from the population maintained for less than a year at Wageningen University and Research (WUR) Greenhouse Horticulture (Bleiswijk, The Netherlands), which in turn were originally sourced from the KBS commercial population. Material used for DNA sequencing, PCR testing, pooled sequencing and cytogenetics was from the KBS population, while material used for RNA sequencing was from the WUR Greenhouse Horticulture population. Additional species used for cytogenetic comparison purposes were sourced from two separate laboratory populations within the Biology Centre CAS in České Budějovice, Czech Republic: *Triatoma infestans* (Klug) (Hemiptera: Reduviidae) individuals were obtained from a laboratory colony at the Institute of Parasitology that was originally sourced from Bolivia (Schwarz *et al.*, 2014), while *Ephesia kuehniella* (Zeller) (Lepidoptera: Pyralidae) individuals were obtained from a wild-type laboratory colony at the Institute of Entomology (Marec and Shvedov, 1990). Species identification of the KBS population was confirmed via *COI* sequencing using a PCR amplification protocol (Itou *et al.*, 2013), in addition to testing for the presence of *Wolbachia* via PCR amplification protocol (Zhou *et al.*, 1998).

### *Flow cytometry*

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

### *gDNA extraction*

A single female *N. tenuis* was placed in a 1.5 ml safelock tube with 5–8 1 mm glass beads and frozen in liquid nitrogen and shaken for 30 s in a Silamat S6 shaker (Ivoclar Vivadent, Schaan, Liechtenstein). DNA was then extracted using the Qiagen MagAttract Kit (Qiagen, Hilden, Germany). Following an overnight lysis step with Buffer ATL and proteinase K at 56 °C, extraction was performed according to MagAttract Kit protocol. Elution was performed in two steps with 50 µl of Buffer AE (Tris-EDTA) each time, yielding 424 ng of genomic DNA (gDNA) in 100 µl as measured with an Invitrogen Qubit 2.0 fluorometer using the dsDNA HS Assay Kit (Thermo Fisher Scientific, Waltham).

### *Library preparation and sequencing*

Following extraction and size selection with magnetic beads, gDNA was further diluted to 1 ng/µl following the Chromium Genome Reagent Kits Version 1 User Guide (version CG-00022) (10x Genomics, Pleasanton). A library of Genome Gel Beads was combined with 1 ng of gDNA, Master Mix and partitioning oil to create Gel Bead-In-EMulsions (GEMs). The GEMs underwent an isothermal amplification step and barcoded DNA fragments

were recovered for Illumina library construction (Illumina, San Diego). The library was then sequenced on an Illumina HiSeq 2500 at the Bioscience Omics Facility at Wageningen University and Research (Wageningen, The Netherlands), yielding 212 910 509 paired-end reads with a read length of 150 bp. The first 23 bp of each forward read is a 10X GEM barcode used in the assembly process. Forward read quality was similar to that of the reverse reads, and no reads were flagged for poor quality in a FastQC assessment (Andrews *et al.*, 2015).

### Assembly

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

### Detecting contamination and LGT events

Lateral gene transfers from bacteria into metazoan genomes were once thought to be rare or non-existent, but are now known to be relatively common and can evolve into functional genes (Dunning Hotopp *et al.*, 2007; Husnik and McCutcheon, 2018). We, therefore, screened our insect genome for LGTs from bacteria. As insect genome assemblies often contain scaffolds from associated bacteria, we first screened for such “contaminating” scaffolds and moved them into a separate metagenomic multiFASTA (S1.4).

We used a DNA-based computational pipeline to both identify likely contaminating bacterial scaffolds in the assembly, and to detect potential LGT from bacteria into the insect genome. The LGT Pipeline was modified from an earlier version developed by David Wheeler and John Werren (Wheeler *et al.*, 2013), and has been used to screen for bacterial “contamination” and LGTs in a number of arthropod genomes before (eg bedbug *C. lectularis* (Benoit *et al.*, 2016), parasitoid wasp *Trichogramma pretiosum* Riley (Hymenoptera: Trichogrammatidae) (Lindsey *et al.*, 2018) and the milkweed bug *Oncopeltus fasciatus* (Dallas) (Hemiptera: Lygaeidae) (Panfilio *et al.*, 2019)). In some cases, entire or nearly

complete bacterial genomes have been retrieved from arthropod genome projects (eg (Benoit *et al.*, 2016; Lindsey *et al.*, 2016)).

**Detection of bacterial scaffolds in the assembly.** To detect bacterial contaminating scaffolds, the following method was used after the preliminary bacterial contamination assessment described above. First, each scaffold was broken into 1 kbp fragments and each fragment was subsequently searched with BLASTn (word size = 11 bp) against an in-house reference database that contains 2100 different bacterial species (complete list in S1.4), which was masked for low complexity regions using the NCBI Dustmasker function (Morgulis *et al.*, 2006). We recorded each bacterial match with bitscore >50, the number of bacterial matches, total bacterial coverage in the scaffold, proportion of the scaffold covered, total hit width of coverage (the distance between the leftmost and rightmost bacteria hit proportional to the scaffold size) and the bacterial species with the greatest number of matches within the scaffold from the in-house bacterial database. It should be noted that the latter method does not indicate the actual bacterial species from which the scaffold was derived, as it is based on similarity to a curated database—that determination would require follow-up analysis, which was not performed in this study.

Any criterion for deciding whether a scaffold comes from a bacterium is unavoidably arbitrary: too stringent and insect scaffolds are included; too lax and insect scaffolds are inappropriately removed. We applied a cut-off of  $\geq 0.40$  proportion bacterial hit width, which has performed well to remove contamination in a few test cases where we have manually examined scaffolds near the cut-off. All instances of contaminated scaffolds were removed from the assembly and are available in supplementary materials as a list (S1.4.2) and a multi-FASTA file (S1.4.3).

**Identifying LGT candidate regions.** We used the same DNA based computational pipeline to identify potential LGTs from bacteria into the insect genome. The basic method is as follows: as before, scaffolds from the genome assembly are broken into 1 kbp intervals, which are searched against a bacterial genome database. Any positive bacterial hit in a 1 kbp region (bitscore >50) was then searched against a database containing transcripts from the following eukaryotes: *Xenopus*, *Daphnia*, *Strongylocentrotus*, *Mus*, *Homo sapiens*, *Aplysia*, *Caenorhabditis*, *Hydra*, *Monosiga* and *Acanthamoeba* ([ftp://ftp.hgsc.bcm.edu/15K-pilot/LGT\\_analysis/All\\_species\\_genomes/lgt\\_finder\\_blastn\\_database\\_directories/](http://ftp.hgsc.bcm.edu/15K-pilot/LGT_analysis/All_species_genomes/lgt_finder_blastn_database_directories/)). The purpose of this eukaryotic screening is to identify highly conserved genes that are shared between eukaryotes and bacteria and exclude these for further analysis. To focus our attentions on the most likely LGT candidates, we selected hits with a bitscore = 0 in the corresponding reference eukaryote database and bitscore >75 from the bacterial database. We also screened the output for adjacent 1 kbp pieces that contain bacterial matches and reference eukaryote bitscore = 0 and fused these adjoining pieces for analysis.

LGT candidate regions were then manually curated as follows: each candidate region was searched with BLASTn to the NCBI nr/nt database. If this search indicated that the region’s nucleotide sequence was similar or identical to the nucleotide sequence of a known gene in related insects, it was discarded as a likely conserved insect gene. Regions were retained only when the matches to other insects were sporadic, as our experience has indicated



that these can be independent LGTs into different lineages. If no match was found, the region was additionally searched with BLASTx to the NCBI nr/nt database. If this second search also resulted in no hits to multiple insect proteins, it was called an LGT candidate. In this case, we additionally identified the best bacterial match using the NCBI nr and protein databases. Using the gene annotation information, we then evaluated the flanking genes within the scaffold to determine whether they were eukaryotic or bacterial, we determined whether the LGT region was associated with an annotated gene within the insect genome, and we observed with transcriptome data if RNA sequencing data showed evidence of transcriptional activity in the LGT region. This short list is available in the supplementary materials (S1.4).

#### RNA extraction, library construction and sequencing

Juveniles, adult males and adult females (approximately 4–5 of each) were prepared for RNAseq using the RNeasy Blood and Tissue Kit (Qiagen). Individuals were placed in a 1.5 ml safelock tube along with 5–8 1 mm glass beads placed in liquid nitrogen and then shaken for 30 s in a Silamat S6 shaker (Ivoclar Vivadent). RNeasy Blood and Tissue Kit (Qiagen) was used according to manufacturer's instructions. Samples were assessed for quality and RNA quantity using an Invitrogen Qubit 2.0 fluorometer and the RNA BR Assay Kit (Thermo Fisher Scientific). These three RNA samples were then processed by Novogene Bioinformatics Technology Co., Ltd., (Beijing, China) using poly(A) selection followed by cDNA synthesis with random hexamers and library construction with an insert size of 550–600 bp. Paired-end sequencing (100 bp) was performed on an Illumina HiSeq 4000 according to manufacturer's instruction.

#### Gene finding, transcriptome assembly and annotation

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

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

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

#### Orthologue cluster analysis and comparison

The complete gene set of *Nesidiocoris tenuis* was compared to those of three additional hemipteran species: the bed bug *C. lectularis*, the brown marmorated stinkbug *H. halys* and the pea aphid *A. pisum* using OrthoVenn2 (Xu *et al.*, 2019a). The gene set of *A. pisum* is the 2015 version from AphidBase (Legeai *et al.*, 2010; Richards *et al.*, 2010) as maintained on the OrthoVenn2 server. The *H. halys* 2.0 complete gene set was used (Lee *et al.*, 2009) along with the complete gene set of *C. lectularis* (Clec 2.1, OGSv1.3) (Benoit *et al.*, 2016; Thomas *et al.*, 2020), both of which were retrieved from the i5K Workspace (Poelchau *et al.*, 2015). An orthologue cluster analysis was performed on all four gene sets via OrthoVenn2 with the default settings of *E*-values of 1e–5 and an inflation value of 1.5. Results were visualized in an UpSet chart (Fig. 1) using UpSetR (Conway *et al.*, 2017).

#### Cytogenetic analysis

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

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

For preparation of the probe, we isolated gDNA from two *N. tenuis* females with the NucleoSpin DNA Insect Kit (Macherey-Nagel, Düren, Germany) according to the manufacturer's protocol. gDNA was used as template in PCR to amplify the 18S rDNA sequence using primers 18S-1 and 18S-4 as described in Jung and Lee (2012). Obtained products were purified using the Wizard SV Gel and PCR Clean-Up System (Promega, Madison, WI), and subsequently cloned using the pGEM-T Easy Vector System (Promega, Madison, WI) according to the manufacturer's protocol. Plasmids were extracted from positive clones with the NucleoSpin Plasmid kit (Macherey-Nagel) following the manufacturer's protocol, confirmed by sequencing (SEQme, Dobříš, Czech Republic) and used as template in PCR with the 18S-1 and 18S-4 primers. PCR-products were purified, and used as template for labelling by a modified nick translation protocol as described by Kato *et al.* (2006) with modifications described in Dalíková *et al.* (2017a), using biotin-16-dUTP (Jena Bioscience, Jena, Germany) and an incubation time of 35 min at 15 °C (Dalíková

*et al.*, 2017a; Kato *et al.*, 2006). Fluorescence *in situ* hybridization was performed as described in Sahara *et al.* (1999) with modifications described in Zrzavá *et al.* (2018) (Sahara *et al.*, 1999; Zrzavá *et al.*, 2018).

**Sex chromosome identification.** Determination of the sex chromosome constitution is important for the assembly of the *N. tenuis* genome to identify any potential missing information due to sequencing a single sex, as well as add to knowledge on sex chromosomes in Miridae. Comparative genomic hybridization, and GISH were, therefore, used to identify the sex chromosomes of *N. tenuis*. The reproductive organs of adult females were dissected out to avoid potential male gDNA contamination, as the mated status was unknown, after which remaining tissue was snap-frozen in liquid nitrogen and stored at  $-20^{\circ}\text{C}$  until further use. Adult males were not dissected but otherwise treated the same. Female and male gDNA was extracted from 10 to 20 pooled individuals using cetyltrimethylammonium bromide (CTAB) gDNA isolation with modifications (Doyle and Doyle, 1990). Samples were mechanically disrupted in extraction buffer (2% CTAB, 100 mM Tris-HCl pH 8.0, 40 mM EDTA, 1.4 M NaCl, 0.2%  $\beta$ -mercaptoethanol, 0.1 mg/ml proteinase K), and incubated overnight at  $60^{\circ}\text{C}$  with light agitation. An equal volume of chloroform was added, tubes were inverted for 2 min and samples were centrifuged 10 min at maximum speed. The aqueous phase was transferred to a new tube, RNase A (200 ng/ $\mu\text{l}$ ) was added and samples were incubated 30 min at  $37^{\circ}\text{C}$  to remove RNA. DNA was precipitated by adding 2/3 volume isopropanol, gently inverting the tubes and centrifugation for 15 min at maximum speed. Pellets were washed twice with 70% ethanol, air-dried briefly and dissolved overnight in sterile water. DNA was stored at  $-20^{\circ}\text{C}$  until further use. Probes were prepared with 1  $\mu\text{g}$  gDNA using Cy3-dUTP (for female gDNA), or fluorescein-dUTP (for male gDNA) (both Jena Bioscience, Jena, Germany) by nick translation mentioned above with an incubation time of 2–2.5 h at  $15^{\circ}\text{C}$ . CGH and GISH were performed according to Traut *et al.* (1999) with modifications described in Dalíková *et al.* (2017a) (Dalíková *et al.*, 2017b; Traut *et al.*, 1999).

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

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

**Repeat identification and visualization.** To assess the repetitive component of the *N. tenuis* genome, we used RepeatExplorer, version 2, on trimmed and quality-filtered reads with default parameters (Novák *et al.*, 2013). Repeats with high abundance in the genome were selected and amplified using PCR. These products, named Nt\_rep1, were additionally cloned and template for probe labelling was prepared from plasmid DNA as described above, see 18S rDNA probe preparation. Probes were labelled by PCR in a volume of 25  $\mu\text{l}$  consisting of 0.625 U Ex *Taq* polymerase (TaKaRa, Otsu, Japan), 1x Ex *Taq* buffer, 40  $\mu\text{M}$  dATP, dCTP and dGTP, 14.4  $\mu\text{M}$  dTTP, 25.6  $\mu\text{M}$  biotin-16-dUTP (Jena Bioscience), 400 nM of forward and reverse primer and 1 ng of purified PCR-product. The amplification programme consisted of an initial denaturing step at  $94^{\circ}\text{C}$  for 3 min, followed by 35 cycles of  $94^{\circ}\text{C}$  for 30 s,  $55^{\circ}\text{C}$  for 30 s,  $72^{\circ}\text{C}$  for 1 min and a final extension at  $72^{\circ}\text{C}$  for 2 min. The FISH procedure was performed as described for 18S rDNA. Abundance and distribution of Nt\_rep1 in the assembled genome was assessed using NCBI Genome Workbench version 2.13.0. The complete list of primers used in this study can be found in Table S1.6.2.

#### Pooled sequencing and population analysis

For this project, it was important to use existing data wherever possible to test the utility of the genome and possible research avenues. Therefore, we analysed whole genome sequence data originally generated for another *N. tenuis* genome assembly project that has not been published before (see accession ERX3938928). In the original set-up, 10 females were collected from the KBS population for pooled sequence analysis. DNA was isolated from this pooled cohort using the “salting-out” method as described in Sunnucks and Hales with a final volume of 20  $\mu\text{l}$  (Sunnucks and Hales, 1996) and then treated with 2  $\mu\text{l}$  RNase A. The paired-end library was sequenced on an Illumina HiSeq2500 platform by Macrogen Inc. (Seoul, Korea) with read sizes of 100 bp. Reads were assessed for quality using FastQC (Andrews *et al.*, 2015) and adapters were trimmed with Trimmomatic (Bolger *et al.*, 2014). Following quality filtering, reads with phred scores lower than 20 were discarded. Heterozygosity was calculated using jellyfish v2.3.0 and GenomeScope v1.0 with a k-

mer size of 21 and default parameters (Marçais and Kingsford, 2011; Vurture *et al.*, 2017).

Instead of genome assembly, these whole genome sequence reads were adjusted and subsequently used in a pooled sequencing (pool-seq) population analysis with our genome. First, the reads were randomly subsampled to a coverage of 10X (in a pool with 10 females, this results in approximately 1X coverage per female) using CLC Genomics Workbench 12 (Qiagen). Using the PoPoolation v1.2.2 pipeline (Kofler *et al.*, 2011), these reads were aligned to an adapted v1.5 genome, where scaffolds smaller than 10 000 bp were removed, and aligned reads were binned into windows using the bwa and samtools packages (Li *et al.*, 2009). Pileup files and the scripts from the PoPoolation pipeline were used to produce variance sliding windows analyses of neutrality, Tajima's *D* and nucleotide diversity,  $\pi$  ( $\pi$ ) with default settings and a pool size of 40. Window and step sizes of both 10 000 and 5000 were tested, as well as using the "basic-pipeline/mask-sam-indelregions.pl" pipeline to mask indel regions of the SAM file—this ensures that indel regions are not calculated. Of the 18 000 000 reads, 817 226 had regions of indels masked.

### Acknowledgements

We would like to express special thanks to Milena Chinchilla Ramírez for her advice and *N. tenuis* expertise. Thanks to Markus Knapp (Koppert BV), Javier Calvo (Koppert BS) and Gerben Messelink (WUR) for providing *N. tenuis* specimens. Thanks to José van de Belt, Frank Becker (WUR), and Carolina Gallego (IVIA) for their technical assistance in DNA and RNA extraction. We acknowledge the assistance of Magda Zrzavá, Anna Voleníková, Martina Flegrová and Diogo Cabral-de-Mello (Institute of Entomology BC CAS) for their cytogenetic expertise and assistance. JHW acknowledges Sammy Cheng for assistance with the LGT pipeline. KBF acknowledges Jetske de Boer for her assistance with the flow cytometry and Joost van den Heuvel for his assistance with PoPoolation. The authors sincerely thank the two anonymous reviewers whose comments helped improve this manuscript. Annotation was performed by GenomeScan B. V. (NL). Access to computing and storage facilities owned by parties and projects contributing to the National Grid Infrastructure MetaCentrum (CZ), provided under the programme "Projects of Large Research, Development, and Innovations Infrastructures" (CESNET LM2015042), is greatly appreciated. This project was funded by the European Union's Horizon 2020 research and innovation program under the Marie Skłodowska-Curie grant agreement no. 641456. JHW acknowledges the US-NSF-IOS-1456233, US-NSF-1950078, and Nathaniel and Helen Wisch Chair for funding support. The research leading to these results was partially funded by the Spanish Ministry of Economy and Competitiveness MINECO (RTA2017-00073-00-00). Cytogenetic experiments were financed by grants 17-13713S (SV and FM) and 17-17211S (MD and IP) of the Czech Science Foundation.

Open access funding enabled and organized by Projekt DEAL.

### Data availability

The data that support the findings of this study are openly available in the supplementary material of this article, in addition to the sequence data found in the European Nucleotide Archive under accession number PRJEB35378, as well as data available on figshare at <https://doi.org/10.6084/m9.figshare.12073893.v1>.

### References

- Abreo, E. and Altier, N. (2019) Pangenome of *Serratia marcescens* strains from nosocomial and environmental origins reveals different populations and the links between them. *Scientific Reports*, **9**(1), 46.
- Acland, A., Agarwala, R., Barrett, T., Beck, J., Benson, D.A., Bollin, C. *et al.* (2014) Database resources of the National Center for Biotechnology Information. *Nucleic Acids Research*, **42** (D1), 8–13.
- Andrews, S., Krueger, F., Semonds-Pichon, A., Biggins, F. and Wingett, S. (2015) In: *FastQC. A quality control tool for high throughput sequence data*. Babraham Bioinformatics. Cambridge: Babraham Institute.
- Arnó, J., Castañé, C., Riudavets, J. and Gabarra, R. (2010) Risk of damage to tomato crops by the generalist zoophytophagous predator *Nesidiocoris tenuis* (Reuter) (Hemiptera: Miridae). *Bulletin of Entomological Research*, **100**(1), 105–115.
- Banaei-Moghaddam, A.M., Meier, K., Karimi-Ashtiyani, R. and Houben, A. (2013) Formation and expression of pseudogenes on the B chromosome of rye. *Plant Cell*, **25**(7), 2536–2544.
- Benoit, J.B., Adelman, Z.N., Reinhardt, K., Dolan, A., Poelchau, M., Jennings, E.C. *et al.* (2016) Unique features of a global human ectoparasite identified through sequencing of the bed bug genome. *Nature Communications*, **7**(1), 10165.
- Benson, G. (1999) Tandem repeats finder: A program to analyze DNA sequences. *Nucleic Acids Research*, **27**(2), 573–580.
- Biondi, A., Zappalà, L., Di Mauro, A., Tropea Garzia, G., Russo, A., Desneux, N. *et al.* (2016) Can alternative host plant and prey affect phytophagy and biological control by the zoophytophagous mirid *Nesidiocoris tenuis*? *BioControl*, **61**(1), 79–90.
- Blankenföld, W., Kuzin, A.P., Skarina, T., Korniyenko, Y., Tong, L., Bayer, P. *et al.* (2004) Structure and function of the phenazine biosynthetic protein PhzF from *Pseudomonas fluorescens*. *Proceedings of the National Academy of Sciences of the United States of America*, **101**(47), 16431–16436.
- Bolger, A.M., Lohse, M. and Usadel, B. (2014) Trimmomatic: A flexible trimmer for Illumina sequence data. *Bioinformatics*, **30** (15), 2114–2120.
- Bouagga, S., Urbaneja, A., Depalo, L., Rubio, L. and Pérez-Hedo, M. (2020) Zoophytophagous predator-induced defences restrict accumulation of the tomato spotted wilt virus. *Pest Management Science*, **76**(2), 561–567.
- Boutet, E., Lieberherr, D., Tognolli, M., Schneider, M. and Bairoch, A. (2008) UniProtKB/Swiss-Prot: The manually annotated section of the UniProt KnowledgeBase. *Methods in Molecular Biology*, **406**(January 2014), 89–112.
- Boyd, B.M., Allen, J.M., Koga, R., Fukatsu, T., Sweet, A.D., Johnson, K.P. *et al.* (2016) Two Bacterial Genera, *Sodalis* and *Rickettsia*, Associated with the Seal Louse

- Proechinophthirus fluctus* (Phthiraptera: Anoplura). *Applied and Environmental Microbiology*, **82**(11), 3185–3197.
- Bruton, B.D., Mitchell, F., Fletcher, J., Pair, S.D., Wayadande, A., Melcher, U. et al. (2007) *Serratia marcescens*, a Phloem-Colonizing, Squash Bug-Transmitted Bacterium: Causal Agent of Cucurbit Yellow Vine Disease. *Plant Disease*, **87**(8), 937–944.
- Calvo, J.F., Bolckmans, K., Stansly, P.a. and Urbaneja, A. (2009) Predation by *Nesidiocoris tenuis* on Bemisia tabaci and injury to tomato. *BioControl*, **54**(2), 237–246.
- Camacho, J.P.M., Sharbel, T.F. and Beukeboom, L.W. (2000) B-chromosome evolution. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences*, **355** (1394), 163–178.
- Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K. et al. (2009) BLAST+: Architecture and applications. *BMC Bioinformatics*, **10**(1), 421.
- Carabajal Paladino, L.Z., Provazníková, I., Berger, M., Bass, C., Aratchige, N.S., López, S.N. et al. (2019) Sex chromosome turnover in moths of the diverse superfamily Gelechioidea. *Genome Biology and Evolution*, **11**(4), 1307–1319.
- Caspi-Fluger, A., Inbar, M., Steinberg, S., Friedmann, Y., Freund, M., Mozes-Daube, N. et al. (2014) Characterization of the symbiont *Rickettsia* in the mirid bug *Nesidiocoris tenuis* (Reuter) (Heteroptera: Miridae). *Bulletin of Entomological Research*, **104**(6), 681–688.
- Castañé, C., Arnó, J., Gabarra, R. and Alomar, O. (2011) Plant damage to vegetable crops by zoophytophagous mirid predators. *Biological Control*, **59**(1), 22–29.
- Charlesworth, B. and Charlesworth, D. (2000) The degeneration of Y chromosomes. *Philosophical Transactions of the Royal Society B: Biological Sciences*, **355**(1403), 1563–1572.
- Chin, C.-S., Peluso, P., Sedlazeck, F.J., Nattestad, M., Concepcion, G.T., Clum, A. et al. (2016) Phased diploid genome assembly with single-molecule real-time sequencing. *Nature Methods*, **13**(12), 1050–1054.
- Chinchilla-Ramírez, M., Pérez-Hedo, M., Pannebakker, B.A. and Urbaneja, A. (2020) Genetic Variation in the Feeding Behavior of Isofemale Lines of *Nesidiocoris tenuis*. *Insects*, **11**(8), 513.
- Chinchilla-Ramírez, M., Garzo, E., Fereres, A., Gavara-Vidal, J., ten Broeke, C.J., van Loon, J.J.A. et al. (2021) Plant feeding by *Nesidiocoris tenuis*: Quantifying its behavioral and mechanical components. *Biological Control*, **152**(March 2020), 104402.
- Conway, J.R., Lex, A. and Gehlenborg, N. (2017) UpSetR: an R package for the visualization of intersecting sets and their properties. *Bioinformatics*, **33**(18), 2938–2940.
- Costechareyre, D., Balmant, S., Condemine, G. and Rahbé, Y. (2012) *Dickeya dadantii*, a plant pathogenic bacterium producing cyt-like entomotoxins, causes septicemia in the pea aphid *Acyrtosiphon pisum*. *PLoS One*, **7**(1) available at, e30702. <https://doi.org/10.1371/journal.pone.0030702>.
- Dai, X., Xun, H., Chang, J., Zhang, J., Hu, B., Li, H. et al. (2012) The complete mitochondrial genome of the plant bug *Nesidiocoris tenuis* (Reuter) (Hemiptera: Miridae: Bryocorinae: Dicyphini). *Zootaxa*, **3554**(3554), 30–44.
- Dalíková, M., Zrzavá, M., Hladová, I., Nguyen, P., Šonský, I., Flegrová, M. et al. (2017a) New Insights into the Evolution of the W Chromosome in Lepidoptera. *Journal of Heredity*, **108** (7), 709–719.
- Dalíková, M., Zrzavá, M., Kubíčková, S. and Marec, F. (2017b) W-enriched satellite sequence in the Indian meal moth, *Plodia interpunctella* (Lepidoptera, Pyralidae). *Chromosome Research*, **25**(3–4), 241–252.
- De Boer, J.G., Ode, P.J., Vet, L.E.M., Whitfield, J.B. and Heimpel, G.E. (2007) Diploid males sire triploid daughters and sons in the parasitoid wasp *Cotesia vestalis*. *Heredity*, **99** (3), 288–294.
- Doyle, J. and Doyle, J. (1990) Isolation of plant DNA from fresh tissue. *Focus*, **12**(1), 13–15.
- Dunning Hotopp, J.C., Clark, M.E., Oliveira, D.C.S.G., Foster, J. M., Fischer, P., Muñoz Torres, M.C. et al. (2007) Widespread lateral gene transfer from intracellular bacteria to multicellular eukaryotes. *Science*, **317**(5845), 1753–1756.
- Ekblom, R. and Wolf, J.B.W. (2014) A field guide to whole-genome sequencing, assembly and annotation. *Evolutionary Applications*, **7**(7), 1026–1042.
- Ellegren, H. (2014) Genome sequencing and population genomics in non-model organisms. *Trends in Ecology & Evolution*, Elsevier Ltd, **29**(1), 51–63.
- Frydrychová, R., Grossmann, P., Trubac, P., Vítková, M. and Marec, F. (2004) Phylogenetic distribution of TTAGG telomeric repeats in insects. *Genome*, **47**(1), 163–178.
- Garantonakis, N., Pappas, M.L., Varikou, K., Skiada, V., Broufas, G.D., Kavroulakis, N. et al. (2018) Tomato inoculation with the endophytic strain *Fusarium solani* K results in reduced feeding damage by the zoophytophagous predator *Nesidiocoris tenuis*. *Frontiers in Ecology and Evolution*, **6**(AUG), 1–7.
- Gatto, K.P., Mattos, J.V., Seger, K.R. and Lourenço, L.B. (2018) Sex chromosome differentiation in the frog genus *Pseudis* involves satellite DNA and chromosome rearrangements. *Frontiers in Genetics*, **9**(AUG), 1–12.
- Gomez-Rodriguez, V.M., Rodriguez-Garay, B., Palomino, G., Martínez, J. and Barba-Gonzalez, R. (2013) Physical mapping of 5S and 18S ribosomal DNA in three species of Agave (Asparagales, Asparagaceae). *Comparative Cytogenetics*, **7** (3), 191–203.
- Grozeva, S., Kuznetsova, V.G. and Anokhin, B.A. (2011) Karyotypes, male meiosis and comparative FISH mapping of 18S ribosomal DNA and telomeric (TTAGG)<sub>n</sub> repeat in eight species of true bugs (Hemiptera, Heteroptera). *Comparative Cytogenetics*, **5**(4), 97–116.
- Grozeva, S., Anokhin, B.A., Simov, N. and Kuznetsova, V.G. (2019) New evidence for the presence of the telomere motif (TTAGG)<sub>n</sub> in the family Reduviidae and its absence in the families Nabidae and Miridae (Hemiptera, Cimicomorpha). *Comparative Cytogenetics*, **13**(3), 283–295.
- Gurevich, A., Saveliev, V., Vyahhi, N. and Tesler, G. (2013) QUASt: Quality assessment tool for genome assemblies. *Bioinformatics*, **29**(8), 1072–1075.
- Haas, B.J., Salzberg, S.L., Zhu, W., Pertea, M., Allen, J.E., Orvis, J. et al. (2008) Automated eukaryotic gene structure annotation using EvidenceModeler and the Program to Assemble Spliced Alignments. *Genome Biology*, **9**(1), 1–22.
- Hare, E.E. and Johnston, J.S. (2011) Molecular methods for evolutionary genetics. In: Orgogozo V., and Rockman M.V., (Eds.) *Molecular Methods for Evolutionary Genetics*. Totowa, NJ: Humana Press, pp. 3–12.

- Husnik, F. and McCutcheon, J.P. (2018) Functional horizontal gene transfer from bacteria to eukaryotes. *Nature Reviews Microbiology*, **16**(2), 67–79.
- Itou, M., Watanabe, M., Watanabe, E. and Miura, K. (2013) Gut content analysis to study predatory efficacy of *Nesidiocoris tenuis* (Reuter) (Hemiptera: Miridae) by molecular methods. *Entomological Science*, **16**(2), 145–150.
- Jang, E.B. and Nishijima, K.A. (1990) Identification and Attractancy of Bacteria Associated with *Dacus dorsalis* (Diptera: Tephritidae). *Environmental Entomology*, **19**(6), 1726–1731.
- Jauset, A.M., Edo-Tena, E., Castañé, C., Agustí, N., Alomar, O. and Grozeva, S. (2015) Comparative cytogenetic study of three *Macrolophus* species (Heteroptera, Miridae). *Comparative Cytogenetics*, **9**(4), 613–623.
- Jetybayev, I.Y., Bugrov, A.G., Ünal, M., Buleu, O.G. and Rubtsov, N.B. (2017) Molecular cytogenetic analysis reveals the existence of two independent neo-XY sex chromosome systems in Anatolian Pamphagidae grasshoppers. *BMC Evolutionary Biology*, **17**(Suppl 1), 1–12.
- Jones, R.N. and Rees, H. (1982) In: *B Chromosomes*. New York: Academic Press.
- Jones, S., Taylor, G., Chan, S., Warren, R., Hammond, S., Bilobram, S. et al. (2017) The Genome of the Beluga Whale (*Delphinapterus leucas*). *Genes*, **8**(12), 378.
- Jung, S. and Lee, S. (2012) Molecular phylogeny of the plant bugs (Heteroptera: Miridae) and the evolution of feeding habits. *Cladistics*, **28**(1), 50–79.
- Kamber, T., Smits, T.H.M., Rezzonico, F. and Duffy, B. (2012) Genomics and current genetic understanding of *Erwinia amylovora* and the fire blight antagonist *Pantoea vagans*. *Trees - Structure and Function*, **26**(1), 227–238.
- Kato, A., Albert, P.S., Vega, J.M. and Birchler, J.A. (2006) Sensitive fluorescence *in situ* hybridization signal detection in maize using directly labeled probes produced by high concentration DNA polymerase nick translation. *Biotechnic and Histochemistry*, **81**(2–3), 71–78.
- Keeling, C.I., Yuen, M.M.S., Liao, N.Y., Roderick Docking, T., Chan, S.K., Taylor, G.A. et al. (2013) Draft genome of the mountain pine beetle, *Dendroctonus ponderosae* Hopkins, a major forest pest. *Genome Biology*, **14**(3), R27.
- Keilwagen, J., Wenk, M., Erickson, J.L., Schattat, M.H., Grau, J. and Hartung, F. (2016) Using intron position conservation for homology-based gene prediction. *Nucleic Acids Research*, **44**(9), e89. <https://doi.org/10.1093/nar/gkw092>.
- Kofler, R., Orozco-terWengel, P., de Maio, N., Pandey, R.V., Nolte, V., Futschik, A. et al. (2011) Popoolation: A toolbox for population genetic analysis of next generation sequencing data from pooled individuals. *PLoS One*, **6**(1), e15925. <https://doi.org/10.1371/journal.pone.0015925>.
- Kuznetsova, V., Grozeva, S., Nokkala, S. and Nokkala, C. (2011) Cytogenetics of the true bug infraorder Cimicomorpha (Hemiptera, Heteroptera): a review. *ZooKeys*, **154**(December), 31–70.
- Lee, W., Kang, J., Jung, C., Hoelmer, K., Lee, S.H. and Lee, S. (2009) Complete mitochondrial genome of brown marmorated stink bug *Halyomorpha halys* (Hemiptera: Pentatomidae), and phylogenetic relationships of hemipteran suborders. *Molecules and Cells*, **28**(3), 155–165.
- Legeai, F., Shigenobu, S., Gauthier, J.P., Colbourne, J., Rispe, C., Collin, O. et al. (2010) AphidBase: A centralized bioinformatic resource for annotation of the pea aphid genome. *Insect Molecular Biology*, **19**(SUPPL. 2), 5–12.
- van Lenteren, J.C., Bolckmans, K., Köhl, J., Ravensberg, W.J. and Urbaneja, A. (2018) Biological control using invertebrates and microorganisms: plenty of new opportunities. *BioControl*, **63**(1), 39–59.
- Leung, K., Ras, E., Ferguson, K.B., Ariëns, S., Babendreier, D., Bijma, P. et al. (2020) Next-generation biological control: the need for integrating genetics and genomics. *Biological Reviews*, **95**(6), 1838–1854. <https://doi.org/10.1111/brv.12641>.
- Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N. et al. (2009) The sequence alignment/map format and SAMtools. *Bioinformatics*, **25**(16), 2078–2079.
- Lindsey, A.R.I., Werren, J.H., Richards, S. and Stouthamer, R. (2016) Comparative genomics of a parthenogenesis-inducing *Wolbachia* symbiont. *G3: Genes, Genomes, Genetics*, **6**(7), 2113–2123.
- Lindsey, A.R.I., Kelkar, Y.D., Wu, X., Sun, D., Martinson, E.O., Yan, Z. et al. (2018) Comparative genomics of the miniature wasp and pest control agent *Trichogramma pretiosum*. *BMC Biology*, **16**(1), 1–20.
- Lowe-Power, T.M., Khokhani, D. and Allen, C. (2018) How *Ralstonia solanacearum* exploits and thrives in the flowing plant Xylem environment. *Trends in Microbiology*, **26**(11), 929–942.
- Majoros, W.H., Pertea, M. and Salzberg, S.L. (2004) TigrScan and GlimmerHMM: Two open source ab initio eukaryotic gene-finders. *Bioinformatics*, **20**(16), 2878–2879.
- Marçais, G. and Kingsford, C. (2011) A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics*, **27**(6), 764–770.
- Marec, F. and Shvedov, A.N. (1990) Yellow eye, a new pigment mutation in *Ephesia kuehniella* Zeller (Lepidoptera: Pyralidae). *Hereditas*, **113**(2), 97–100.
- Martínez-García, H., Román-Fernández, L.R., Sáenz-Romo, M. G., Pérez-Moreno, I. and Marco-Mancebón, V.S. (2016) Optimizing *Nesidiocoris tenuis* (Hemiptera: Miridae) as a biological control agent: mathematical models for predicting its development as a function of temperature. *Bulletin of Entomological Research*, **106**(02), 215–224.
- May, C.M., Heuvel, J., Doroszuk, A., Hoedjes, K.M., Flatt, T. and Zwaan, B.J. (2019) Adaptation to developmental diet influences the response to selection on age at reproduction in the fruit fly. *Journal of Evolutionary Biology*, **32**(5), 425–437.
- Mediouni, J., Fuková, I., Frydrychová, R., Marec, F., Fuková, I., Marec, F. et al. (2004) Karyotype, sex chromatin and sex chromosome differentiation in the carob moth, *Ectomyelois ceratoniae* (Lepidoptera: Pyralidae). *Caryologia*, **57**(2), 184–194.
- Moerkens, R., Pekas, A., Bellinkx, S., Hanssen, I., Huysmans, M., Bosmans, L. et al. (2020) *Nesidiocoris tenuis* as a pest in Northwest Europe: Intervention threshold and influence of Pepino mosaic virus. *Journal of Applied Entomology*, **144**(7), 566–577.
- Mollá, O., Biondi, A., Alonso-Valiente, M. and Urbaneja, A. (2014) A comparative life history study of two mirid bugs preying on *Tuta absoluta* and *Ephesia kuehniella* eggs on tomato crops: Implications for biological control. *BioControl*, **59**(2), 175–183.

- Morgulis, A., Gertz, E.M., Sch affer, A.A. and Agarwala, R. (2006) A fast and symmetric DUST implementation to mask low-complexity DNA sequences. *Journal of Computational Biology*, **13**(5), 1028–1040.
- Nov ak, P., Neumann, P., Pech, J., Steinhaisl, J. and MacAs, J. (2013) RepeatExplorer: A Galaxy-based web server for genome-wide characterization of eukaryotic repetitive elements from next-generation sequence reads. *Bioinformatics*, **29**(6), 792–793.
- Panfilio, K.A. and Angelini, D.R. (2018) By land, air, and sea: hemipteran diversity through the genomic lens. *Current Opinion in Insect Science*, **25**, 106–115.
- Panfilio, K.A., Vargas Jentzsch, I.M., Benoit, J.B., Erezylmaz, D., Suzuki, Y., Colella, S. et al. (2019) Molecular evolutionary trends and feeding ecology diversification in the Hemiptera, anchored by the milkweed bug genome. *Genome Biology*, **20**(1), 64.
- Paspati, A., Ferguson, K.B., Verhulst, E.C., Urbaneja, A., Gonz alez-Cabrera, J. and Pannebakker, B.A. (2019) Effect of mass rearing on the genetic diversity of the predatory mite *Amblyseius swirskii* Athias-Henriot (Acari: Phytoseiidae). *Entomologia Experimentalis et Applicata*, **167**, 670–681.
- P rez-Hedo, M. and Urbaneja, A. (2016) The zoophytophagous predator *Nesidiocoris tenuis*: A successful but controversial biocontrol agent in tomato crops. Horowitz, A & Ishaaya, I., In: *Advances in Insect Control and Resistance Management*. Cham: Springer International Publishing, pp. 121–138.
- P rez-Hedo, M., Urbaneja-Bernat, P., Jaques, J.A., Flors, V. and Urbaneja, A. (2015) Defensive plant responses induced by *Nesidiocoris tenuis* (Hemiptera: Miridae) on tomato plants. *Journal of Pest Science*, **88**(3), 543–554.
- P rez-Hedo, M., Arias-Sanguino,  .M. and Urbaneja, A. (2018) Induced tomato plant resistance against *Tetranychus urticae* Triggered by the Phytophagy of *Nesidiocoris tenuis*. *Frontiers in Plant Science*, **9**(October), 1–8.
- P rez-Hedo, M., Riahi, C. and Urbaneja, A. (2020) Use of zoophytophagous mirid bugs in horticultural crops: current challenges and future perspectives. *Pest Management Science*, ps.6043. **77**(1), 33–42.
- Pita, S., Panzera, F., Mora, P., Vela, J., Palomeque, T. and Lorite, P. (2016) The presence of the ancestral insect telomeric motif in kissing bugs (Triatominae) rules out the hypothesis of its loss in evolutionarily advanced Heteroptera (Cimicomorpha). *Comparative Cytogenetics*, **10**(3), 427–437.
- Poelchau, M., Childers, C., Moore, G., Tsavatapalli, V., Evans, J., Lee, C.Y. et al. (2015) The i5k Workspace@NAL-enabling genomic data access, visualization and curation of arthropod genomes. *Nucleic Acids Research*, **43**(D1), D714–D719.
- Prado, S.S. and Almeida, R.P.P. (2009) Phylogenetic placement of pentatomid stink bug gut symbionts. *Current Microbiology*, **58**(1), 64–69.
- Puentes, A., Stephan, J.G. and Bj rkman, C. (2018) A systematic review on the effects of plant-feeding by omnivorous arthropods: Time to catch-up with the mirid-tomato bias? *Frontiers in Ecology and Evolution*, **6**(218). <https://doi.org/10.3389/fevo.2018.00218>.
- Rasmussen, L.B., Jensen, K., S rensen, J.G., Sverrisd ttir, E., Nielsen, K.L., Overgaard, J. et al. (2018) Are commercial stocks of biological control agents genetically depauperate?—A case study on the pirate bug *Orius majusculus* Reuter. *Biological Control*, **127**(June), 31–38.
- Reddiex, A.J., Allen, S.L. and Chenoweth, S.F. (2018) A Genomic Reference Panel for *Drosophila serrata*. *G3: Genes, Genomes Genetics*, **8**(4), 1335–1346.
- Richards, S. and Murali, S.C. (2015) Best practices in insect genome sequencing: what works and what doesn't. *Current Opinion in Insect Science*, **7**, 1–7.
- Richards, S., Gibbs, R.A., Gerardo, N.M., Moran, N., Nakabachi, A., Stern, D. et al. (2010) Genome Sequence of the Pea Aphid *Acyrtosiphon pisum*. *PLoS Biology*, **8**(2), e1000313.
- Sachman-Ruiz, B. and Quiroz-Casta eda, R.E. (2018) Genomics of Rickettsiaceae: An update. Quiroz-Casta eda, R.E., In: *Farm Animals Diseases, Recent Omic Trends and New Strategies of Treatment*, Vol. I London: InTechOpen, pp. 13.
- Sahara, K., Marec, F. and Traut, W. (1999) TTAGG telomeric repeats in chromosomes of some insects and other arthropods. *Chromosome Research*, **7**(6), 449–460.
- Sanchez, J.A. (2009) Density thresholds for *Nesidiocoris tenuis* (Heteroptera: Miridae) in tomato crops. *Biological Control*, **51**(3), 493–498.
- Santos-Garcia, D., Silva, F.J., Morin, S., Dettner, K. and Kuechler, S.M. (2017) The all-rounder *Sodalis*: A new bacteriome-associated endosymbiont of the lygaeoid bug *Henestaris halophilus* (Heteroptera: Henestarinae) and a critical examination of its evolution. *Genome Biology and Evolution*, **9**(10), 2893–2910.
- Schwarz, A., Medrano-Mercado, N., Schaub, G.A., Struchiner, C. J., Barges, M.D., Levy, M.Z. et al. (2014) An updated insight into the sialotranscriptome of *Triatoma infestans*: developmental stage and geographic variations. *PLoS Neglected Tropical Diseases*, **8**(12), e3372.
- Sember, A., Bertollo, L.A.C., R b, P., Yano, C.F., Hatanaka, T., de Oliveira, E.A. et al. (2018) Sex chromosome evolution and genomic divergence in the fish *Hoplias malabaricus* (Characiformes, Erythrinidae). *Frontiers in Genetics*, **9**(MAR), 1–12.
- Sim o, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E.V. and Zdobnov, E.M. (2015) BUSCO: Assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*, **31**(19), 3210–3212.
- Sobreira, T.J.P., Durham, A.M. and Gruber, A. (2006) TRAP: Automated classification, quantification and annotation of tandemly repeated sequences. *Bioinformatics*, **22**(3), 361–362.
- Sochorov a, J., Garcia, S., G lvez, F., Symonov a, R. and Kovařík, A. (2018) Evolutionary trends in animal ribosomal DNA loci: introduction to a new online database. *Chromosoma*, **127**(1), 141–150.
- Streito, J.C., Clouet, C., Hamdi, F. and Gauthier, N. (2017) Population genetic structure of the biological control agent *Macrolophus pygmaeus* in Mediterranean agroecosystems. *Insect Science*, **24**(5), 859–876.
- Sunnucks, P. and Hales, D.F. (1996) Numerous transposed sequences of mitochondrial cytochrome oxidase I-II in aphids of the genus *Sitobion* (Hemiptera: Aphididae). *Molecular Biology and Evolution*, **13**(3), 510–524.
- Sz cs, M., Vercken, E., Bitume, E.V. and Hufbauer, R.A. (2019) The implications of rapid eco-evolutionary processes for

- biological control—a review. *Entomologia Experimentalis et Applicata*, eea.12807. **167**, 598–615.
- Testa, A.C., Hane, J.K., Ellwood, S.R. and Oliver, R.P. (2015) CodingQuarry: highly accurate hidden Markov model gene prediction in fungal genomes using RNA-seq transcripts. *BMC Genomics*, **16**(1), 1–12.
- Thomas, G.W.C., Dohmen, E., Hughes, D.S.T., Murali, S.C., Poelchau, M., Glastad, K. *et al.* (2020) Gene content evolution in the arthropods. *Genome Biology*, **21**(1), 15.
- Tian, C., Tek Tay, W., Feng, H., Wang, Y., Hu, Y. and Li, G. (2015) Characterization of *Adelphocoris suturalis* (Hemiptera: Miridae) transcriptome from different developmental stages. *Scientific Reports*, **5**(1), 11042.
- Trapnell, C., Pachter, L. and Salzberg, S.L. (2009) TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics*, **25**(9), 1105–1111.
- Trapnell, C., Williams, B.A., Pertea, G., Mortazavi, A., Kwan, G., van Baren, M.J. *et al.* (2010) Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nature Biotechnology*, Nature Publishing Group, **28**(5), 511–515.
- Traut, W. (1976) Pachytene mapping in the female silkworm, *Bombyx mori* L. (Lepidoptera). *Chromosoma*, **58**(3), 275–284.
- Traut, W., Sahara, K., Otto, T.D. and Marec, F. (1999) Molecular differentiation of sex chromosomes probed by comparative genomic hybridization. *Chromosoma*, **108**(3), 173–180.
- Traverse, K.L. and Pardue, M.L. (1988) A spontaneously opened ring chromosome of *Drosophila melanogaster* has acquired He-T DNA sequences at both new telomeres. *Proceedings of the National Academy of Sciences of the United States of America*, **85**(21), 8116–8120.
- Urbaneja, A., González-Cabrera, J., Arnó, J. and Gabarra, R. (2012) Prospects for the biological control of *Tuta absoluta* in tomatoes of the Mediterranean basin. *Pest Management Science*, **68**(9), 1215–1222.
- Urbaneja-Bernat, P., Bru, P., González-Cabrera, J., Urbaneja, A. and Tena, A. (2019) Reduced phytophagy in sugar-provisioned mirids. *Journal of Pest Science*, **92**(3), 1139–1148.
- Vurture, G.W., Sedlazeck, F.J., Nattestad, M., Schatz, M.C., Gurtowski, J., Underwood, C.J. *et al.* (2017) GenomeScope: fast reference-free genome profiling from short reads. *Bioinformatics*, **33**(14), 2202–2204.
- Weisenfeld, N.I., Kumar, V., Shah, P., Church, D.M. and Jaffe, D. B. (2017) Direct determination of diploid genome sequences. *Genome Research*, **27**, 757–767.
- Wheeler, D., Redding, A.J. and Werren, J.H. (2013) Characterization of an Ancient Lepidopteran Lateral Gene Transfer. *PLoS One*, **8**(3), e59262.
- van Wilgenburg, E., Driessen, G. and Beukeboom, L.W. (2006) Single locus complementary sex determination in Hymenoptera: an 'unintelligent' design? *Frontiers in Zoology*, **3**, 1.
- Xu, L., Dong, Z., Fang, L., Luo, Y., Wei, Z., Guo, H. *et al.* (2019a) OrthoVenn2: A web server for genome wide comparison and annotation of orthologous clusters across multiple species. *Nucleic Acids Research*, Oxford University Press, **47**(W1), W52–W58.
- Xu, P., Lu, B., Liu, J., Chao, J., Donkersley, P., Holdbrook, R. *et al.* (2019b) Duplication and expression of horizontally transferred polygalacturonase genes is associated with host range expansion of mirid bugs. *BMC Evolutionary Biology*, BMC Evolutionary Biology, **19**(1), 12.
- Xun, H., Li, H., Li, S., Wei, S., Zhang, L., Song, F. *et al.* (2016) Population genetic structure and post-LGM expansion of the plant bug *Nesidiocoris tenuis* (Hemiptera: Miridae) in China. *Scientific Reports*, **6**(26755). <https://doi.org/10.1038/srep26755>.
- Zhang, Y. and Qiu, S. (2015) Examining phylogenetic relationships of *Erwinia* and *Pantoea* species using whole genome sequence data. *Antonie van Leeuwenhoek*, **108**(5), 1037–1046.
- Zhou, W., Rousset, F. and O'Neill, S. (1998) Phylogeny and PCR-based classification of *Wolbachia* strains using *wsp* gene sequences. *Proceedings of the Royal Society B: Biological Sciences*, **265**(1395), 509–515.
- Zrzavá, M., Hladová, I., Dalíková, M., Šichová, J., Ůunap, E., Kubíčková, S. *et al.* (2018) Sex Chromosomes of the Iconic Moth *Abraxas grossulariata* (Lepidoptera, Geometridae) and Its Congener *A. sylvata*. *Genes*, **9**(279). <https://doi.org/10.3390/genes9060279>.

### Supporting Information

Additional supporting information may be found online in the Supporting Information section at the end of the article.

**Appendix S1:** Supporting Information.