1 **Genome dynamics across the evolutionary transition to endosymbiosis**

2

3 Stefanos Siozios (1+), Pol Nadal Jimenez (1), Tal Azagi (2), Hein Sprong (2), Crystal L Frost (1), Steven

4 R Parratt (1), Graeme Taylor (3), Laura Brettell (4), Kwee Chin Liew (5), Larry Croft (6), Kayla C King

5 (1,7), Michael A Brockhurst (1,8), Václav Hypša (9), Eva Novakova (9), Alistair C Darby (1), and

6 Gregory DD Hurst (1+).

7

8 [1] Institute of Infection, Veterinary and Ecological Sciences, University of Liverpool, Liverpool L59 7ZB

9 UK.

10 [2] Centre for Infectious Diseases Research, National Institute for Public Health and the Environment,

11 3720 BA Bilthoven, the Netherlands.

12 [3] Department of Biology, University of Victoria, Victoria, BC, Canada

13 [4] Liverpool School of Tropical Medicine

14 [5] NSW Health Pathology Infectious Diseases Department, Wollongong Hospital, NSW, Australia

15 [6] School of Medicine, Deakin University, 75 Pigdons Road, Waurn Ponds, Victoria 3216, Australia.

16 [7] Department of Biology, University of Oxford

17 [8] Division of Evolution, Infection and Genomics | Faculty of Biology, Medicine and Health | University

18 of Manchester

19 [9] Department of Parasitology, Faculty of Science, University of South Bohemia, České Budějovice,

20 Czech Republic

21

22 + For correspondence

23 siozioss@liverpool.ac.uk; g.hurst@liverpool.ac.uk

24

25

26 Further information and requests for resources and reagents should be directed to and will be fulfilled

27 by the lead contact,  Greg Hurst (g.hurst@liverpool.ac.uk)

28

29

36

37

38

39

40

41

42

43    *Summary*

44    Endosymbiosis – where a microbe lives and replicates within a host – is an important contributor to
45    organismal function that has accelerated evolutionary innovations and catalysed the evolution of
46    complex life. The evolutionary processes associated with transitions to endosymbiosis, however, are
47    poorly understood. Here, we use comparative genomics of the genus *Arsenophonus* to reveal the
48    complex processes that occur on evolution of an endosymbiotic lifestyle. We compared the genomes
49    of 38 strains spanning diverse lifestyles from environmentally acquired infections to obligate inter-
50    dependent endosymbionts. We observed recent endosymbionts had larger genome sizes than
51    closely related environmentally acquired strains, consistent with evolutionary innovation and rapid
52    gain of new function. Increased genome size was a consequence of prophage and plasmid
53    acquisition including a cargo of type III effectors, and concomitant loss of CRISPR-Cas genome
54    defence systems enabling mobile genetic element expansion. Persistent endosymbiosis was also
55    associated with loss of type VI secretion, likely reflecting reduced microbe-microbe competition.
56    Thereafter, the transition to stable endosymbiosis and vertical inheritance was associated with the
57    expected relaxation of purifying selection, pseudogenisation of genes and reduction of metabolism,
58    leading to genome reduction. However, reduced %GC that is typically considered a progressive linear
59    process was observed only in obligate interdependent endosymbionts. We argue that a combination
60    of the need for rapid horizontal gene transfer-mediated evolutionary innovation together with
61    reduced phage predation in endosymbiotic niches drives loss of genome defence systems and rapid
62    genome expansion upon adoption of endosymbiosis. These remodelling processes precede the
63    reductive evolution traditionally associated with adaptation to endosymbiosis.

64

65    **Results and Discussion**

66

67    Animals live in a microbial world. Their interactions with microbes range from the antagonistic with
68    pathogenic symbionts, through to the mutualistic with beneficial symbionts, with evolutionary
69    transitions occurring commonly between these states [1, 2]. Transitions in symbiotic interactions can
70    further select for evolution of key symbiotic traits, such as vertical transmission and eventual
71    integration of symbionts into host anatomy and physiology. These transitions correspondingly alter
72    selection pressures. For example, vertical transmission relaxes selection for traits necessary for
73    external survival whilst also correlating microbe transmission with host fitness and thus favouring
74    beneficial function(s). Concurrently, population bottlenecks associated with vertical transmission
75    limit within-host symbiont diversity, selecting for lower virulence [3]. Further, these bottlenecks
76    intensify genetic drift [4], reducing the efficiency of purifying selection for function. These processes
77    collectively drive genome degradation through pseudogenization, genome reduction and lowered
78    %GC [5].

79

80    Endosymbiosis is the state in which the symbiont lives within the body or cells of the host organism.
81    Evolutionary transitions from environmental to endosymbiotic lifestyles are common across the tree
82    of microbial diversity, and have occurred with microeukaryotic, fungal, plant and animal hosts [2].
83    However, our understanding of the evolutionary processes associated with transitions to
84    endosymbiosis is limited. In particular, there have been few studies of the initial transition to host
85    association and vertical transmission inherent to becoming an endosymbiont. Fewer chart the entire
86    transition from environmentally acquired associations through to obligate vertically inherited
87    endosymbiont via a facultative endosymbiotic stage.

88

89    To gain a precise view of the tempo and mode of evolution during the transition to persistent
90    endosymbiosis, we leveraged the wide diversity of host-associated lifestyles in the
91    gammaproteobacterial genus *Arsenophonus*. Within this clade, there are environmentally acquired
92    extracellular pathogens, extracellular endosymbionts with mixed modes of transmission,
93    intracellular facultative endosymbionts with vertical transmission, and obligate vertically transmitted

2

94    endosymbionts where the partners are co-dependent [6-10]. Importantly, extracellular pathogens
95    and endosymbionts with mixed modes of transmission are very closely related. This genus thus
96    provides a unique opportunity to investigate the genomic adaptations enabling the emergence of an
97    endosymbiotic host association, building from work in other clades with more coarse-grained
98    comparisons [11, 12].
99
100   To gain a high-resolution view of the evolutionary transition to endosymbiosis in an insect-
101   associated bacterial endosymbiont, we first completed closed genomes for seven new strains: three
102   *Arsenophonus nasoniae* from *Nasonia vitripennis*, one from each of the parasitic wasps
103   *Pachycrepoideus vindemmiae* and *Ixodiphagus hookeri*, one from the blue butterfly *Polyommatus*
104   *bellargus*, and a strain of *Arsenophonus apicola* isolated from Australian *Apis mellifera*. We also
105   completed a draft genome for *Ca.* A. triatominarum. In addition, novel draft genomes for a further
106   17 *Arsenophonus* strains were assembled from Sequence Read Archive (SRA) deposits from a variety
107   of insect genome sequencing projects. Genome assembly data, alongside sample details and current
108   understanding of transmission mode and nature of symbiosis, are given in Supplementary Table 1.
109
110   We then estimated the phylogenetic affiliation of the strains through core genome analysis (Figure
111   1), maximizing the common gene set for phylogenetic inference by excluding the strains with highly
112   reduced genomes. This analysis revealed three main clades according to lifestyle; the *nasoniae* clade
113   where characterized members have mixed modes of transmission or have recently become
114   facultative vertically transmitted symbionts; the *apicola* clade where characterized strains are
115   infectiously transmitted, and the *triatominarum* clade where characterized members are vertically
116   transmitted. When the obligate co-dependent symbionts are included in the phylogeny
117   (Supplementary Figure S1), they do not form a monophyletic clade as was previously shown [13],
118   indicating independent evolutionary transitions to obligate host dependence. Culturable
119   representatives are common in the *apicola/nasoniae* clades, likely reflecting a current or recent
120   requirement to live outside the endosymbiotic environment. Members of the bee genus *Bombus*
121   have members from both *apicola* and *nasoniae* clades, indicating recurrent colonization by the
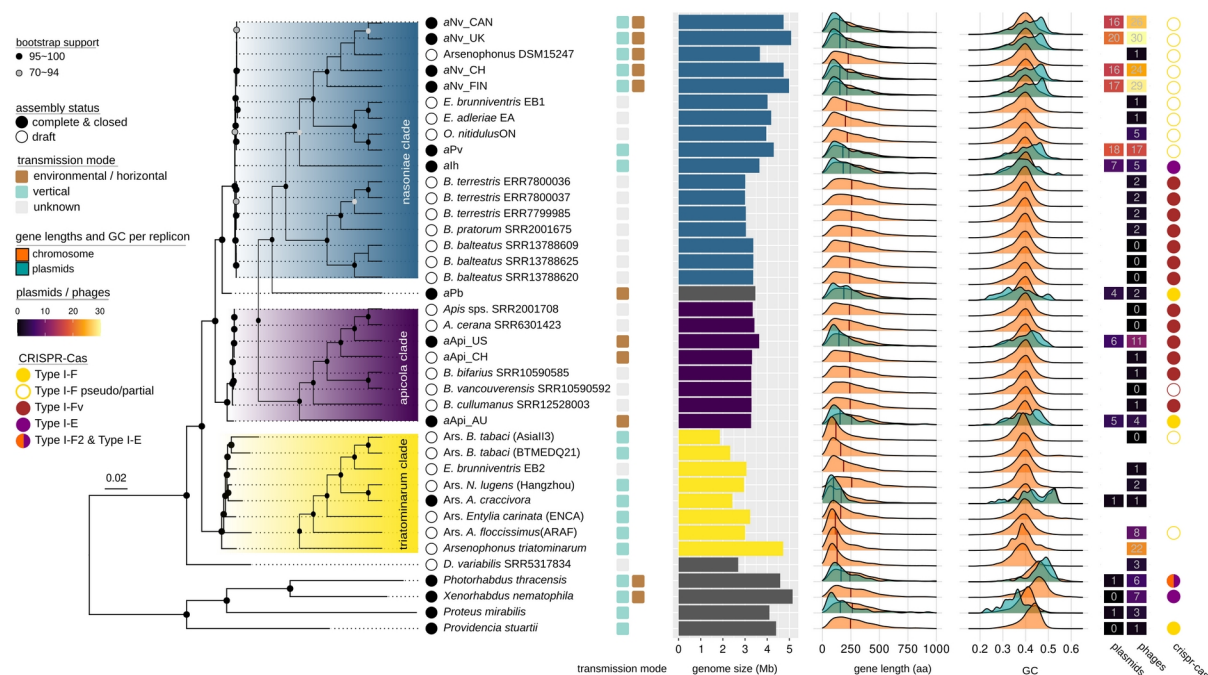122   endosymbiont of this host group.
123

**Figure 1** *Core genome phylogeny and genome features of the Arsenophonus clade. The phylogenetic relationships between the Arsenophonus strains were inferred using maximum likelihood on the concatenated set of 230 single copy core protein sequences in IQ-TREE v2.1.4 under the JTTDCMut+F+R3 model. Only bootstrap support values >= 70 are shown. Inset cladograms are used to improve tree readability. The nasoniae, apicola and triatominarum clades are highlighted in blue, magenta and yellow respectively, while complete genomes are indicated by black circles in front of the tip labels. The transmission mode of each strain is indicated with the coloured squares (green: vertical, brown: horizontal and grey: currently unknown). The horizontal bar-plot shows the genome size in MB. The ridgeline plots show the distribution of CDS length (in amino acids) and GC content fraction for genes of chromosomal (orange) or extra-chromosomal (green) origin. The vertical lines in the gene length ridgeplot represent median values. The heatmap shows the number of plasmids and intact phages, while the coloured circles indicate the type and intactness of the CRISPR-Cas system (filled circle: intact, non-filled circle: pseudogenised, no circle: not present). A Bayesian phylogenetic analysis including the four obligate Arsenophonus genomes (Arsenophonus of Lipoptena fortisetosa, Arsenophonus of Aleurodicus dispersus, Arsenophonus of Ceratovacuna japonica and Arsenophonus of Melophagus ovinus) and Ca. Riesia is shown in the Supplementary Figure S1.*

125   *Arsenophonus* strains varied in genome size from 663,125 to 5,080,918 bp. It is generally considered
126   that symbiont genome size is an inverse function of time evolving as endosymbiont under vertical
127   transmission, with strains with long history of vertical transmission having smaller genomes. As
128   expected, the smallest genomes in the clade were indeed from interdependent obligate
129   endosymbionts and members of the triatominarum clade where all members are vertically
130   transmitted and require a host for replication (Figure 2A). Unexpectedly, the five strains with the
131   largest closed genomes were not those with environmental transmission, but related strains with
132   either mixed modes of transmission or vertical transmission. This pattern reflected larger genome
133   size in the nasoniae clade (where characterized members have protracted host association and
134   vertical transmission and genome size ranges from 3.65-5.1MB) compared to the apicola clade
135   (where characterized members are environmentally acquired pathogens and genome size is 3.27-
136   3.63 Mb) (Figure 2A). The increase in genome size in early-stage endosymbionts is driven largely by
137   an increase in mobile genetic elements, notably prophage and plasmid content, in strains which
138   have recently entered endosymbiosis and lost environmental transmission (Figure 1).
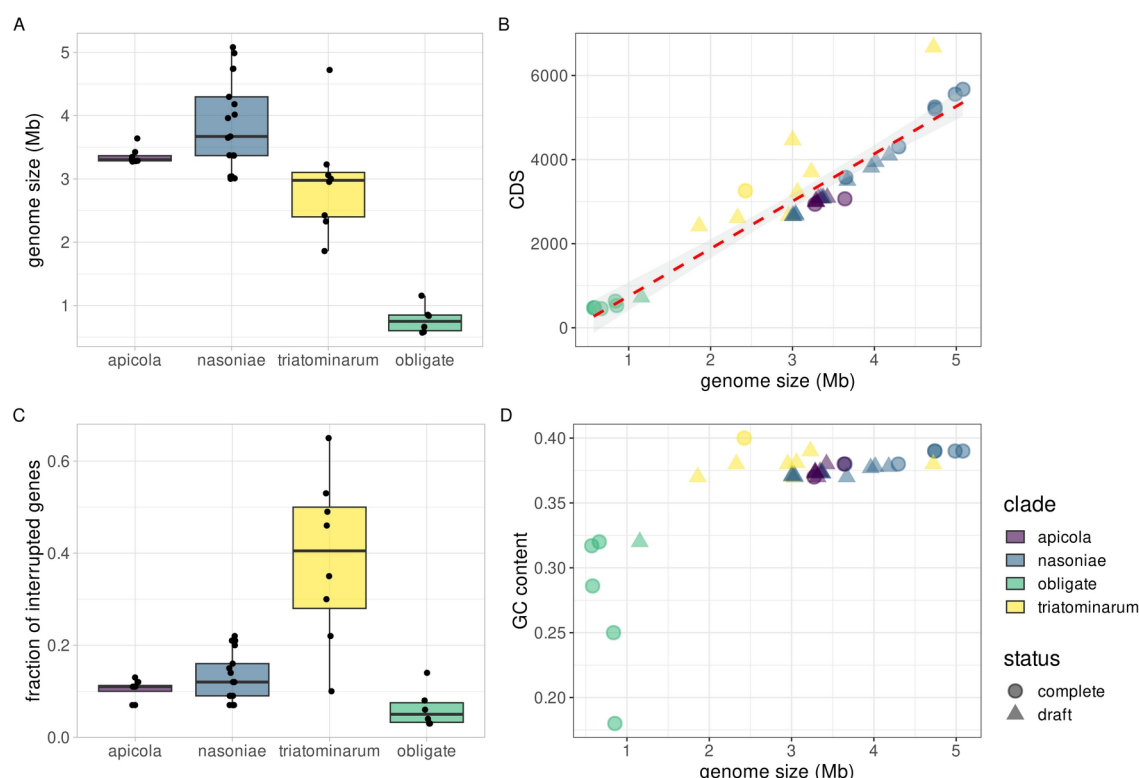
***Figure 2*** *Genomic characteristics of the Arsenophonus/Riesia clades. **A**) Genome size distribution across Arsenophonus clades, **B**) Association between genome size and the number of coding sequences (CDS). The red dashed line in panel B represents a fitted linear trend line with confidence intervals shown as grey shading. **C**) the fraction of interrupted genes across Arsenophonus clade. This was estimated by calculating the fraction of proteins with length <80% of the length of their top hit in the Swiss-Prot database. **D**) Association between genome size and the fraction of GC content. Although the number of predicted protein-coding genes shows, as expected, a linear relationship with the genome size across the Arsenophonus/Riesia clades, this association does not hold for GC content contrary to the classical observations between free-living and symbiotic microbes. Box plots: center line, median; box limits, 25th and 75th percentiles; whiskers, ±1.5x interquartile range; data points are shown with the black dots.*

139

140    Acquisition of mobile genetic elements, like prophages and plasmids in bacteria, is often
141    accompanied with horizontal gene transfer of accessory genes that are important in microbial
142    virulence and adaptation [14, 15]. Consistent with this pattern, we observed a gain in Type III
143    secretion (T3SS) associated effectors in nasoniae group strains that have recently become
144    endosymbionts, compared to environmentally acquired strains (Figure 3). These data support
145    previous work in *Sodalis* arguing symbionts repurpose the T3SS on transition from pathogenesis to
146    enable persistent intracellular host association [16]. Further to this hypothesis, our data indicate that
147    the process may in some cases involve acquisition of new effectors. Contrastingly, T3SS are either
148    absent or heavily pseudogenized in the triatominarum clade where vertical transmission is well
149    established. The loss of T3SS systems in highly derived vertically inherited endosymbiont genomes
150    suggests that these traits become redundant or costly once the host and endosymbiont have
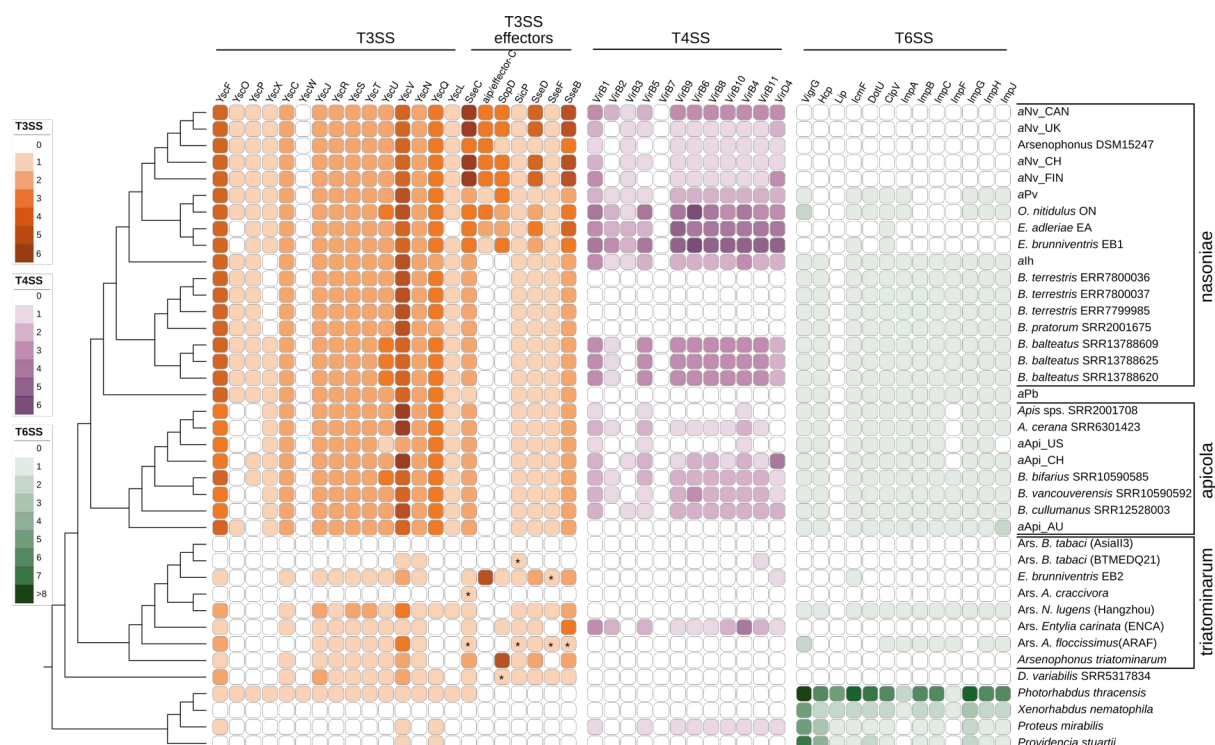151    become tightly coevolved.

**Figure 3** *Comparative analysis of secretion systems across the Arsenophonus clades. Core components of the Type III (orange), type IV (magenta) and type VI (green) systems as predicted by BlastKOALA are shown. The absence of genes is indicated by empty squares. Identified type III effectors are also shown. The relationship of the Arsenophonus strains is shown with the cladogram based on the core genome phylogeny (Figure 1). The asterisks indicate potential pseudogenes. Not shown in the Figure: Secretion systems are absent in the obligate Arsenophonus lineages (Arsenophonus of Lipoptena fortisetosa, Arsenophonus of Aleurodicus dispersus, Arsenophonus of Ceratovacuna japonica and Arsenophonus of Melophagus ovinus) and Ca. Riesia.*

153

154  We next investigated potential drivers of prophage and plasmid accumulation. Past analysis has
155  indicated that genome defence systems, such as CRISPR-Cas systems which protect the genome
156  from mobile DNA, are less commonly found in symbiotic microbes than in free-living ones [17]. We
157  observed that the identity, distribution and completeness of CRISPR-Cas genome defence systems
158  varies extensively among *Arsenophonus* genomes (Figure 4). We identified three types of CRISPR-Cas
159  systems (types I-F, I-Fv and I-E) with variable gene content across all strains, suggesting dynamic
160  turnover of these important genome defence systems (Figure 4A). This turnover is likely mediated by
161  horizontal gene transfer and recombination (Supplementary Figure S2), consistent with variable
162  mobile genetic element-mediated selection for genome defence across niches. Variable mobile
163  element exposure is further supported by the distinct lack of spacer matches between different
164  *Arsenophonus* clades. This result suggests that different clades have been exposed to distinct mobile
165  genetic element communities (Figure 4B).

166

167  Intact CRISPR-Cas systems were predominantly found in environmentally acquired strains. Notably,
168  closely related vertically inherited strains contained recently pseudogenized CRISPR-Cas systems
169  (evidenced by intact and shorter CRISPR arrays but insertions within the Cas systems), while obligate
170  intracellular symbionts carried either highly degraded fragmentary CRISPR-Cas systems or none
171  (Figure 1 and Figure 4A). CRISPR-Cas is known to have metabolic and autoreactive costs [18].
172  Theoretical and experimental data suggest that rapid loss of CRISPR-Cas function is under selection
173  as a means to maintain horizontally transferred genetic elements that are beneficial for microbial
174  adaptation [19, 20]. Thus, the extensive pseudogenization and loss of CRISPR-Cas systems outside of
175  environmentally acquired strains likely reflects purifying selection against maintaining this genome

176  defence systems within host environments. Selection against CRISPR-Cas within host-associated
177  environments could reflect the selection to retain beneficial mobile genetic elements whilst avoiding
178  autoimmunity or may be due to lower rates of phage attack within hosts. Aside CRISPR-based
179  defences, other predicted anti-phage systems are also more diverse in the environmentally acquired
180  strains compared to strains with mixed modes of transmission. These are fragmentary in strains
181  exhibiting vertical transmission (Supplementary Figure S3 and Supplementary Table S2).
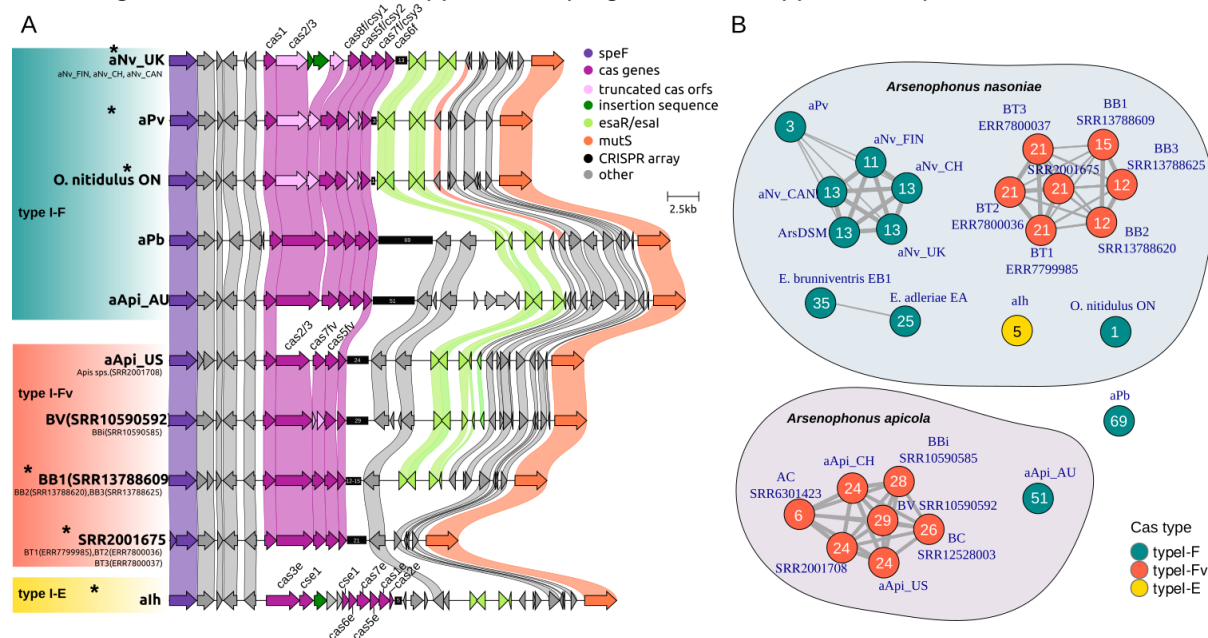


**Figure 4** *The CRISPR-Cas phage defence system in Arsenophonus and early signs of pseudogenisation within the nasoniae clade. **A**) Gene order and genomic context of the Arsenophonus CRISPR-Cas systems plotted with clinker software. Groups of homologous genes are connected by coloured ribbons. Arsenophonus strains belonging to the nasoniae clade are highlighted with an asterisk. **B**) Networks showing the relatedness of Arsenophonus genomes based on shared CRISPR spacers. Nodes correspond to genomes and the edges are scaled based on the spacer repertoire relatedness (see details in Methods section). Absence of edge correspond to no shared spacers between the genomes. Nodes are coloured according to the Cas type identified in each genome. The numbers in the nodes represent the number of spacers identified by CRISPRCasFinder tool.*

182
183  Within intracellular and other endosymbiotic host-associated niches, bacteria are likely to encounter
184  fewer competitors. To test if this weakened selection to retain anticompetitor weapons, we
185  examined the presence and integrity of type VI secretion systems (T6SS), a contact-dependent
186  system for killing competitor microbes [21]. Whereas T6SS were present in the environmentally
187  acquired strains, these were incomplete in nasoniae clade strains with either mixed modes or
188  vertical transmission and absent or fragmentary in all but one member of the triatominarum clade. It
189  is in this last clade where all members show vertical transmission (Figure 3). T6SS are multiprotein
190  complexes and thus likely to be costly to produce and use ([22] but see [23]). Their loss in
191  endosymbionts supports a model where reduced competition within the endosymbiotic niche
192  reducing the benefits of maintaining T6SS.
193
194  Classically, the transition to endosymbiosis has been associated with reduced metabolic capability of
195  endosymbionts. The stable nutritional environment within host cells reduces the need for metabolic
196  plasticity. Contrastingly, a recent study of Chlamydiae concluded that the transition to intracellular
197  endosymbiosis was associated with increased predicted metabolic capability of the symbionts [24]).
198  In the *Arsenophonus* clade, predicted metabolic functions are reduced compared to free living
199  ancestors. However, metabolic capabilities are not markedly distinct between vertically transmitted
200  endosymbionts and environmentally-acquired strains. The completeness of metabolic pathways

201   broadly reflects the abilities of strains to grow in *in vitro* cell-free culture [9]. One notable exception
202   is beta oxidation of fatty acids, which was only intact in environmentally acquired strains; this result
203   was sustained when draft genomes for the triatominarum group were additionally examined
204   (Supplementary Figure S4). Aspects of cofactor synthesis, amino acid synthesis and carbohydrate
205   metabolism were degraded in strains that live intracellularly, and loss of function in these pathways
206   was most pronounced for obligate co-dependent endosymbiont strains. Together, these patterns
207   suggest that the transition to intracellular endosymbiosis, rather than endosymbiosis *per se*, is
208   associated with evolved losses of metabolic functions. Whether this degradation is a function of
209   longer evolutionary time under relaxed selection or is actively selected for by the intracellular niche
210   is unclear.
211
212   We further examined gene loss processes over the life course of symbiosis. As expected, the number
213   of predicted CDS was an approximately linear function of genome size (Figure 2B). The fraction of
214   pseudogenized genes was low in both environmentally acquired strains and in interdependent
215   obligate symbionts, at intermediate levels in the nasoniae group containing vertically transmitted
216   strains that retain infectious transfer and culturability, and at highest levels in unculturable
217   facultative endosymbiont strains (Figure 2C). This pattern was only apparent in non-core genes
218   (genes not present in all strains) (Supplementary figure S5). These data suggest a process of
219   pseudogenization that accelerates on entering into endosymbiosis, as systems for environmental
220   survival become redundant in the host context. The number of pseudogenes then increases over
221   time until purifying selection drives their loss and genomes become streamlined to core genes.
222
223   A key hypothesis in symbiont evolution is that vertical transmission leads to bottlenecks in symbiont
224   population size. This process is expected to increase the importance of drift over selection and thus
225   weaken the capacity of purifying selection to maintain function [25]. We analysed the pattern of
226   molecular evolution of highly conserved single copy genes that are critical for microbial function (see
227   methods for details). All three clades (nasoniae, apicola and triatominarum) showed evidence of
228   relaxed selection compared to the autonomous free-living outgroup. Relaxed selection was most
229   pronounced in the triatominarum clade comprising vertically transmitted intracellular symbionts
230   with a long history of endosymbiosis. Relaxation of selection was also observed in the other two
231   clades. Notably, it was somewhat more pronounced in the nasoniae clade where the strains have
232   mixed modes of transmission or vertical transmission than the apicola clade where environmental
233   transmission is common (Figure 5). These data corroborate our current thinking of evolution
234   patterns in vertically transmitted symbiosis, which combines adaptive gene loss through redundancy
235   with gene degradation through fixation of mildly deleterious alleles, the latter permitted by
236   increased primacy of drift processes. Notably, this signature can be detected shortly following the
237   evolution of vertical transmission.
238
239   Finally, we examined how overall genomic features vary between strains. Unlike previous work [26],
240   our data do not support the linear association between genome size and %GC content during the
241   transition to endosymbiosis, at least in the genus *Arsenophonus* (Figure 2D). Reduced %GC is only
242   markedly observed in the obligate co-dependent endosymbionts with highly reduced genomes. For
243   the other strains, %GC content is relatively consistent at 37-40%. Analysis of coding sequences of
244   non−obligate strains did not support a relationship between %GC and genome size (Null hypothesis
245   of no association: $F_{1,33}$ = 1.141, p=0.29; see Supplementary Figure S6). Our data suggest therefore
246   that reduced %GC may be restricted to obligate interdependent symbioses in this group. This may
247   also reflect changes in DNA repair systems which are ablated in obligate strains (Supplementary
248   figure S7), increased primacy of genetic drift associated with the pronounced bottlenecks that
249   accompany obligate interdependent symbioses, or a combination of these factors.
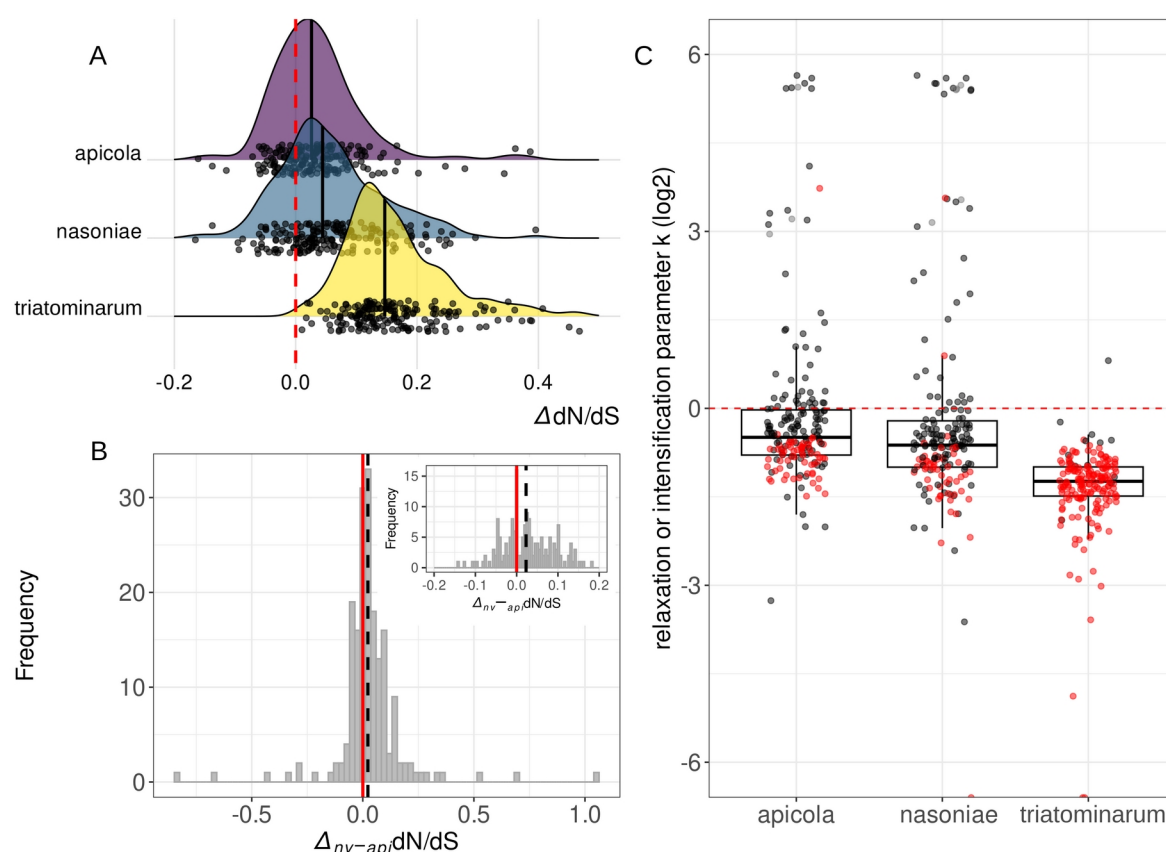
***Figure 5*** *Evidence of relaxation of selection between Arsenophonus clades with contrasting modes of transmission.* ***A****) Distribution of gene-wise dN/dS ratios in the three main Arsenophonus clades as compared to the outgroup clade comprised of free-living species (Providencia stuartii, Proteus mirabilis, Morganella morganii and Moellerella wisconsensis) (dN/dS$_{test\ clade}$ – dN/dS$_{outgroup}$) for 188 highly conserved single-copy BUSCO marker genes. Individual values are shown as jitter points. Black solid lines represent the median of the distribution. In all clades the median is shifted to the right with the triatominarum clade (vertical transmission) showing the largest shift followed by the nasoniae clade (mixed mode of transmission) and last the apicola clade (environmental transmission) suggesting a gradual increase in the dN/dS ratios as we progress towards protract symbiosis.* ***B****) Distribution of gene-wise differences in dN/dS ratios between nasoniae and apicola clades for the same 188 highly conserved BUSCO marker genes. The black dotted vertical line represents the median of the distribution which is shifted to the right (median = 0.0228146, Wilcoxon signed rank test V = 11373, P = 4.018e-06) indicating that nasoniae clade has higher dN/dS ratios compared to apicola clade which is mostly characterised by environmental mode of transmission. A narrower range of the same data between values -0.2 and 0.2 is shown in the inset plot on the top right corner.* ***C****) Distribution of relaxation or intensification parameter k (log2) per gene as calculated by the RELAX method in HyPhy package (v2.3) compared to the outgroup clade. Values below zero indicate that selection strength has been relaxed while values above zero indicate an intensification of the selection strength. Genes with statistically significant k values (FDR; q<0.1) are shown as red jitter dots.*

251

252    In summary, our examination of genome evolution across symbiotic lifestyles in the genus

253    *Arsenophonus* reveals a new model for the evolution of endosymbiosis. Becoming a persistent

254    endosymbiont requires rapid evolutionary innovation fuelled by horizontal gene transfer, notably

255    the gain of new functions for host manipulation. This rapid genome expansion is achieved through

256    the acquisition of prophage and plasmid mobile genetic elements, which is enabled by the loss of

257    genome defence systems, including CRISPR-Cas. An initial establishment of endosymbiosis phase is

258    associated with enrichment for T3SS effector toxins but this is followed by their loss in strains that

259    become vertically transmitted and more highly adapted to the host intracellular environment. Whilst

260    our model is based on data from a single genus, the centrality of CRISPR defence loss reflects recent
261    studies of *Mycoplasma* evolution following a host switch event [27], and the presence of intact or
262    recently pseudogenized CRISPR in culturable, but not unculturable, aphid-associated *Serratia* [12]. In
263    *Mycoplasma*, the authors argued the loss of CRISPR systems enabled adaptation to the novel host
264    species. We argue that the evolutionary transition from free-living to endosymbiosis may commonly
265    be associated with a more complex genome dynamics than previously reported, with genome
266    expansion and remodelling preceding the processes of reductive evolution traditionally associated
267    with endosymbiosis.

268
269
270

271    **Method details**
272    **Isolate collection, culture, sequencing and assembly**
273    Seven *Arsenophonus* isolates were isolated to pure culture and sequenced within this project: *A.*
274    *nasoniae* isolates aNv_UK, aNv_CH and aNv_CAN derived from *Nasonia vitripennis* from the UK,
275    Switzerland and Canada respectively, *A. nasoniae aPv* from *Pachycrepoideus vindemmiae* and the
276    *Arsenophonus* strain *a*Pb from the butterfly *Polyommatus bellargus* [9], *Arsenophonus nasoniae*
277    strain *a*Ih previously identified in the parasitoid wasp *Ixodiphagus hookeri* from a questing *Ixodes*
278    *ricinus* tick in The Netherlands [28] and *Arsenophonus apicola* strain aApi_AU from Australian honey
279    bees described in [29]. A further strain, *Ca*. A. triatominarum, was isolated into cell culture. Details of
280    isolation, culture, sequencing and assembly are given in Supplementary methods.

281

282    ***Arsenophonus* genomes assembled from publicly available SRA deposits.** We screened publicly
283    available SRA datasets (Source: DNA, Platform: Illumina, Strategy: genome) originated mainly from
284    the Apoidea superfamily (containing bees and bumblebees) as well as Parasitoida infraorder and
285    Ixodida order for the presence of *Arsenophonus* reads. We did that by performing a "*Mash screen*"
286    using Mash v2.3 [30, 31] to measure the containment of a local database of reference Arsenophonus
287    genomes within the unassembled SRA read sets. SRA datasets with at least 80% containment were
288    taken for downstream processing. An initial metagenomic assembly of the short reads from the
289    identified SRA datasets were performed using MEGAHIT v1.2.8 [32] and the assembled contigs >=
290    1.5kb were binned based on their deferential tetranucleotide frequencies using MetaBAT2 v2.12.1
291    under default parameters [33]. The *Arsenophonus* bins were identified and completeness was
292    assessed using CheckM v1.0.18 [34]. To identify *Arsenophonus* contigs potentially missed from the
293    initial binning process the original contigs from the metagenomic assembly were screened using
294    blastn (-task megablast) against a local database consisted of all available and complete
295    *Arsenophonus* genomes using BLAST 2.12.0+ [35]. Contigs >1kb with significant matches (e-value <
296    1e-25) to *Arsenophonus* genomes were extracted and included in the metabat bins. The augmented
297    bins were quality inspected and further refined in anvio v7 [36] by identifying and removing
298    potential contaminant contigs based on atypical coverage and gene-level taxonomic classification.
299    The original BioProject and SRA accessions from which the draft Arsenophonus was obtained are
300    shown in Supplementary Table S3.

301

302    **Comparative analysis of the metabolic potential across the *Arsenophonus* clades.**
303    To avoid inconsistencies stemming from draft and incomplete genomes, only the metabolic potential
304    of complete *Arsenophonus* genomes was estimated. To these we included for comparison the
305    genomes of the closely related and obligate symbionts *Ca*. Riesia pediculicola and *Ca*. Riesia
306    pediculischaeffi as well as the genomes of the four outgroup species used in the phylogenetic
307    analysis. All genomes were annotated for functions and metabolic pathways on the basis of KEGG
308    database using the anvi-run-kegg-kofams command in anvio v7. KEGG MODULE metabolism was
309    finally estimated using the anvi-estimate-metabolism pipeline as described in [37]. A KEGG module

310    was considered to be complete in a given genome when at least 75% of the steps involved were
311    present.
312
313    **Phylogenomic analysis.**
314    For the maximum likelihood phylogenomic analysis the highly divergent genomes from the obligate
315    Arsenophonus strains (Ca. *Arsenophonus lipoptenae*, *Arsenophonus* of *Aleurodicus dispersus*, *Ca.*
316    Arsenophonus melophagi and *Arsenophonus* of *Ceratovacuna japonica*) including *Ca*. Riesia
317    pediculola were excluded as their placement can be affected by strong compositional heterogeneity
318    and long branch attraction. The phylogenetic relationship of the remaining 35 *Arsenophonus*
319    genomes was estimated on the concatenated set of 230 single-copy core protein sequences
320    representing highly conserved gammaproteobacterial BUSCO v4.1.4 markers [38]. These were
321    identified through Orthofinder v2.3.11 [39]. Alignments of individual protein orthologs were
322    performed using mafft program [40] as implemented in Orthofinder and screened for recombination
323    based on the Pairwise Homoplasy Index (PHI) test using the Phipack package [41] revealing no
324    significant evidence. The alignments were quality trimmed using ClipKIT alignment trimming tool
325    v1.3.0 [42] under the *smart-gap* mode before concatenated into a super matrix using seqkit [43].
326    Best protein model (JTTDCMut+F+R3) was identified using ModelFinder [44] and a ML phylogenetic
327    tree was reconstructed in IQ-TREE v1.6.12 [45]. The genomes from the related species *Proteus*
328    *mirabilis* strain HI4320 (NC_010554), *Providencia stuartii* strain MRSN 2154 (NC_017731),
329    *Photorhabdus thracensis* strain DSM 15199 (NZ_CP011104) and *Xenorhabdus nematophila* ATCC
330    19061 (NC_014228) were used as an outgroup. All genomes were pre-annotated using prokka v1.13
331    [46] for consistency.
332
333    To more precisely estimate the relationships between the *Arsenophonus* strains including the
334    obligate and highly diverse lineages a separate Bayesian phylogenetic analysis was conducted based
335    on the concatenated set of 77 manually curated single-copy core protein clusters using PhyloBayes-
336    MPI v1.9 [47] and the CAT-Poisson model. Briefly, the alignments of all single copy protein clusters
337    (101) were manually inspected to minimize alignments with high gap content before concatenation.
338    Two independent chains were run in parallel for at least 30,000 cycles until convergence was
339    observed (rel diff < 0.1 and minimum effective size > 300 for all trace file metrics) assessed by
340    running bpcomp and tracecomp commands in PhyloBayes.
341
342    **Annotation of genomic features.**
343    The proportion of pseudogenised genes were estimated by calculating the fraction of interrupted
344    proteins using the Ideel method against the UniProt/Swiss-Prot database
345    (https://github.com/mw55309/ideel). Prophage regions were annotated using the PHAge Search
346    Tool Enhanced Release (PHASTER) web server [48]. Annotation of CRISPR arrays and cas genes was
347    performed with the CRISPRCasFinder program using the default parameters [49]. CRISPR spacer
348    relatedness was calculated as the number of shared spacers by two genomes divided by the number
349    of spacers in the smallest array [50]. Shared spacers were identified based on an all-vs-all blastn
350    search allowing for at least 98% similarity and 97% coverage between individual spacers. Apart from
351    CRISPR-Cas systems we screened for additional phage defense systems using DefenseFinder [51].
352
353    **Analysis of relaxation of selection strength between Arsenophonus clades.**
354    We search for signatures of relaxation of selection on a set of 188 highly conserved single-copy
355    BUSCO orthologs across 39 *Arsenophonus* strains using the HyPhy (Hypothesis Testing using
356    Phylogenies) software package version v2.5.8 [52]. The obligate and highly diverged *Arsenophonus*
357    strains including *Ca*. Riesia were excluded from the analyses. Four, mostly environmental,
358    Morganellaceae (Enterobacterales) species (*Proteus mirabilis*, *Providencia stuartii*, *Morganella*
359    *morganii* and *Moellerella wisconsensis*) were used as reference/outgroup since *Photorhabdus*
360    *thracensis* and *Xenorhabdus nematophila* that were used in the phylogenetic analyses above are

361  having complex modes of transmission involving both vertical and horizontal transmission. To this
362  end, single-copy protein clusters were identified as before using Orthofinder v2.3.11 and 188
363  clusters representing highly conserved gammaproteobacterial BUSCO v4.1.4 markers were selected
364  for downstream analyses. Alignment of protein sequences was performed with mafft using the "-
365  auto" option and back translated to nucleotide alignments using pal2nal [53]. A phylogenetic tree
366  was estimated using the JTTDCMut+F+I+G4 model in IQ-TREE v1.6.12 on the concatenated data of
367  the same set of 188 orthologues protein clusters. This tree was then used to identify genes with
368  significant evidences of relaxation or intensification of selection using the RELAX hypothesis testing
369  framework in HyPhy package [54]. Three sets of branches were selected as the "test branches"
370  (nasoniae clade, apicola clade and the triatominarum clade, see Supplementary Figure S8) and
371  compared to the reference/outgroup set of branches to estimate the selection intensity parameter k
372  for each gene. A value of k > 1 indicates that selection on the test branches is intensified compared
373  to the reference branches while a value < 1 indicates a relaxation of selection. Statistically significant
374  values were assessed through a likelihood ratio test (LRT) followed by a false discovery rate (fdr)
375  correction to account for multiple comparisons. Branch specific dN/dS ratios were estimated for
376  each individual gene using the partitioned MG94xREV model which fits a single dN/dS value to each
377  branch partition as implemented in RELAX method in HyPhy. The differences of dN/dS ratios
378  between branches (ΔdN/dS) were statistically assessed using a Wilcoxon signed-rank test in R v4.2.2
379  [55].
380

381  **Data visualization**
382  Tools used for data visualization: R v4.2.2 [55], ggplot2 [56], patchwork [57], igraph v1.4.0 [58]. A
383  presence/absence heatmap for metabolic pathways was generated using the pheatmap v1.0.12
384  (https://rdrr.io/cran/pheatmap/) package in R v4.2.2. Phylogenetic trees were drawn and annotated
385  using the ggtree package v4.2 [59]. The gene order comparison of CRISPR-Cas systems was visualized
386  in clinker v0.0.24 [60].
387

388  **Data and code availability**
389  All genomes and the raw reads generated in this study are deposited in GenBank database under the
390  BioProject accession number PRJNA956975. The genome for *Ca.* Arsenophonus triatominarum can
391  be found in GenBank under the BioProject accession number PRJNA311587. The code and source
392  data for the various analyses in this study can be found on GitHub
393  (https://github.com/SioStef/Arsenophonus-comparative-genomics).
394
395

396    References
397

398    1.    Sachs, J.L., Skophammer, R.G., and Regus, J.U. (2011). Evolutionary transitions in bacterial
399          symbiosis. Proceedings of the National Academy of Sciences of the United States of America
400          *108*, 10800-10807.
401    2.    Drew, G.C., Stevens, E.J., and King, K.C. (2021). Microbial evolution and transitions along the
402          parasite–mutualist continuum. Nature Reviews Microbiology *19*, 623-638.
403    3.    Leeksl, A., dos Santos, M., and West, S.A. (2019). Transmission, relatedness, and the
404          evolution of cooperative symbionts. Journal of Evolutionary Biology *32*, 1036-1045.
405    4.    Bennett, G.M., and Moran, N.A. (2015). Heritable symbiosis: The advantages and perils of an
406          evolutionary rabbit hole. Proceedings of the National Academy of Sciences *112*, 10169-
407          10176.
408    5.    Toft, C., and Andersson, S.G. (2010). Evolutionary microbial genomics: insights into bacterial
409          host adaptation. Nat Rev Genet *11*, 465-475.
410    6.    Gherna, R.L., Werren, J.H., Weisburg, W., Cote, R., Woese, C.R., Mandelco, L., and Brenner,
411          D.J. (1991). Arsenophonus-Nasoniae Gen-Nov, Sp-Nov, the Causative Agent of the Son-Killer
412          Trait in the Parasitic Wasp Nasonia-Vitripennis. International Journal Of Systematic
413          Bacteriology *41*, 563-565.
414    7.    Novakova, E., Hypsa, V., and Moran, N.A. (2009). Arsenophonus, an emerging clade of
415          intracellular symbionts with a broad host distribution. BMC Microbiol. *9*, 143.
416    8.    Parratt, S.R., Frost, C.L., Schenkel, M.A., Rice, A., Hurst, G.D.D., and King, K.C. (2016).
417          Superparasitism Drives Heritable Symbiont Epidemiology and Host Sex Ratio in a Wasp. Plos
418          Pathogens *12*, e1005629.
419    9.    Nadal-Jimenez, P., Parratt, S.R., Siozios, S., and Hurst, G.D.D. (2023). Isolation, culture and
420          characterization of Arsenophonus symbionts from two insect species reveal loss of infectious
421          transmission and extended host range. Frontiers in Microbiology *14*, 1089143.
422    10.   Nadal-Jimenez, P., Siozios, S., Frost, C.L., Court, R., Chrostek, E., Drew, G.C., Evans, J.D.,
423          Hawthorne, D.J., Burritt, J.B., and Hurst, G.D.D. (2022). Arsenophonus apicola sp. nov.,
424          isolated from the honeybee Apis mellifera. International Journal of Systematic and
425          Evolutionary Microbiology *72*, *e*005469.
426    11.   Santos-Garcia, D., Silva, F.J., Morin, S., Dettner, K., and Kuechler, S.M. (2017). The All-
427          Rounder Sodalis: A New Bacteriome-Associated Endosymbiont of the Lygaeoid Bug
428          Henestaris halophilus (Heteroptera: Henestarinae) and a Critical Examination of Its
429          Evolution. Genome Biol Evol *9*, 2893-2910.
430    12.   Renoz, F., Foray, V., Ambroise, J., Baa-Puyoulet, P., Bearzatto, B., Mendez, G.L., Grigorescu,
431          A.S., Mahillon, J., Mardulyn, P., Gala, J.-L., et al. (2021). At the Gate of Mutualism:
432          Identification of Genomic Traits Predisposing to Insect-Bacterial Symbiosis in Pathogenic
433          Strains of the Aphid Symbiont Serratia symbiotica. Frontiers in Cellular and Infection
434          Microbiology *11*, e660007.
435    13.   Yorimoto, S., Hattori, M., Kondo, M., and Shigenobu, S. (2022). Complex host/symbiont
436          integration of a multi-partner symbiotic system in the eusocial aphid Ceratovacuna japonica.
437          iScience *25*, e105478.
438    14.   Fortier, L.C., and Sekulovic, O. (2013). Importance of prophages to evolution and virulence of
439          bacterial pathogens. Virulence *4*, 354-365.
440    15.   Dragoš, A., Andersen, A.J.C., Lozano-Andrade, C.N., Kempen, P.J., Kovács, Á.T., and Strube,
441          M.L. (2021). Phages carry interbacterial weapons encoded by biosynthetic gene clusters.
442          Current Biology *31*, 3479-3489.e3475.
443    16.   Dale, C., Young, S.A., Haydon, D.T., and Welburn, S.C. (2001). The insect endosymbiont
444          Sodalis glossinidius utilizes a type III secretion system for cell invasion. Proceedings of the
445          National Academy of Sciences of the United States of America *98*, 1883-1888.

446  17.  Burstein, D., Sun, C.L., Brown, C.T., Sharon, I., Anantharaman, K., Probst, A.J., Thomas, B.C.,
447       and Banfield, J.F. (2016). Major bacterial lineages are essentially devoid of CRISPR-Cas viral
448       defence systems. Nature Communications *7*, 10613.
449  18.  Zaayman, M., and Wheatley, R.M. (2022). Fitness costs of CRISPR-Cas systems in bacteria.
450       Microbiology *168*, 001209.
451  19.  Gandon, S., and Vale, P.F. (2014). The evolution of resistance against good and bad
452       infections. J Evol Biol *27*, 303-312.
453  20.  Jiang, W., Maniv, I., Arain, F., Wang, Y., Levin, B.R., and Marraffini, L.A. (2013). Dealing with
454       the Evolutionary Downside of CRISPR Immunity: Bacteria and Beneficial Plasmids. PLOS
455       Genetics *9*, e1003844.
456  21.  Unni, R., Pintor, K.L., Diepold, A., and Unterweger, D. (2022). Presence and absence of type
457       VI secretion systems in bacteria. Microbiology *168*, e001151.
458  22.  Septer, A.N., Sharpe, G., and Shook, E.A. (2023). The Vibrio fischeri type VI secretion system
459       incurs a fitness cost under host-like conditions. bioRxiv, 2023.2003.2007.529561.
460  23.  Zhang, C., Ratcliff, W.C., and Hammer, B.K. (2023). Constitutive expression of the Type VI
461       secretion system carries no measurable fitness cost in Vibrio cholerae. bioRxiv,
462       2023.2003.2024.534098.
463  24.  Dharamshi, J.E., Köstlbacher, S., Schön, M.E., Collingro, A., Ettema, T.J.G., and Horn, M.
464       (2023). Gene gain facilitated endosymbiotic evolution of Chlamydiae. Nature Microbiology *8*,
465       40-54.
466  25.  Moran, N.A. (1996). Accelerated evolution and Muller's rachet in endosymbiotic bacteria.
467       Proceedings of the National Academy of Sciences *93*, 2873-2878.
468  26.  Lo, W.-S., Huang, Y.-Y., and Kuo, C.-H. (2016). Winding paths to simplicity: genome evolution
469       in facultative insect symbionts. FEMS Microbiol Rev *40*, 855-874.
470  27.  Ipoutcha, T., Tsarmpopoulos, I., Gourgues, G., Baby, V., Dubos, P., Hill, G.E., Dowling, A., Arfi,
471       Y., Lartigue, C., Thebault, P., et al. (2023). Evolution of the CRISPR-Cas9 defence system
472       following a bacterial host shift. bioRxiv, 2023.2003.2014.532377.
473  28.  Krawczyk, A.I., Bakker, J.W., Koenraadt, C.J.M., Fonville, M., Takumi, K., Sprong, H., and
474       Demir, S. (2020). Tripartite Interactions among Ixodiphagus hookeri, Ixodes ricinus and Deer:
475       Differential Interference with Transmission Cycles of Tick-Borne Pathogens. Pathogens *9*,
476       9050339.
477  29.  Liew, K.C., Graves, S., Croft, L., Brettell, L.E., Cook, J., Botes, J., and Newton, P. (2022). First
478       human case of infection with Arsenophonus nasoniae, the male killer insect pathogen.
479       Pathology *54*, 664-666.
480  30.  Ondov, B.D., Treangen, T.J., Melsted, P., Mallonee, A.B., Bergman, N.H., Koren, S., and
481       Phillippy, A.M. (2016). Mash: fast genome and metagenome distance estimation using
482       MinHash. Genome Biology *17*, 132.
483  31.  Ondov, B.D., Starrett, G.J., Sappington, A., Kostic, A., Koren, S., Buck, C.B., and Phillippy, A.M.
484       (2019). Mash Screen: high-throughput sequence containment estimation for genome
485       discovery. Genome Biology *20*, 232.
486  32.  Li, D., Liu, C.-M., Luo, R., Sadakane, K., and Lam, T.-W. (2015). MEGAHIT: an ultra-fast single-
487       node solution for large and complex metagenomics assembly via succinct de Bruijn graph.
488       Bioinformatics *31*, 1674-1676.
489  33.  Kang, D.D., Li, F., Kirton, E., Thomas, A., Egan, R., An, H., and Wang, Z. (2019). MetaBAT 2: an
490       adaptive binning algorithm for robust and efficient genome reconstruction from
491       metagenome assemblies. PeerJ *7*, e7359.
492  34.  Parks, D.H., Imelfort, M., Skennerton, C.T., Hugenholtz, P., and Tyson, G.W. (2015). CheckM:
493       assessing the quality of microbial genomes recovered from isolates, single cells, and
494       metagenomes. Genome Res *25*, 1043-1055.
495  35.  Altschul, S.F., Gish, W., Miller, W., Myers, E.W., and Lipman, D.J. (1990). Basic local
496       alignment search tool. J. mol. biol. *215*, 403-410.

497   36.   Eren, A.M., Kiefl, E., Shaiber, A., Veseli, I., Miller, S.E., Schechter, M.S., Fink, I., Pan, J.N.,
498         Yousef, M., Fogarty, E.C., et al. (2021). Community-led, integrated, reproducible multi-omics
499         with anvi'o. Nature Microbiology *6*, 3-6.
500   37.   Muto, A., Kotera, M., Tokimatsu, T., Nakagawa, Z., Goto, S., and Kanehisa, M. (2013).
501         Modular architecture of metabolic pathways revealed by conserved sequences of reactions.
502         J Chem Inf Model *53*, 613-622.
503   38.   Manni, M., Berkeley, M.R., Seppey, M., and Zdobnov, E.M. (2021). BUSCO: Assessing
504         Genomic Data Quality and Beyond. Curr Protoc *1*, e323.
505   39.   Emms, D.M., and Kelly, S. (2019). OrthoFinder: phylogenetic orthology inference for
506         comparative genomics. Genome Biology *20*, 238.
507   40.   Katoh, K., and Standley, D.M. (2013). MAFFT Multiple Sequence Alignment Software Version
508         7: Improvements in Performance and Usability. Molecular Biology and Evolution *30*, 772-
509         780.
510   41.   Bruen, T.C., Philippe, H., and Bryant, D. (2006). A simple and robust statistical test for
511         detecting the presence of recombination. Genetics *172*, 2665-2681.
512   42.   Steenwyk, J.L., Buida, T.J., III, Li, Y., Shen, X.-X., and Rokas, A. (2020). ClipKIT: A multiple
513         sequence alignment trimming software for accurate phylogenomic inference. PLOS Biology
514         *18*, e3001007.
515   43.   Shen, W., Le, S., Li, Y., and Hu, F. (2016). SeqKit: A Cross-Platform and Ultrafast Toolkit for
516         FASTA/Q File Manipulation. PLOS ONE *11*, e0163962.
517   44.   Kalyaanamoorthy, S., Minh, B.Q., Wong, T.K.F., von Haeseler, A., and Jermiin, L.S. (2017).
518         ModelFinder: fast model selection for accurate phylogenetic estimates. Nature Methods *14*,
519         587-589.
520   45.   Nguyen, L.-T., Schmidt, H.A., von Haeseler, A., and Minh, B.Q. (2015). IQ-TREE: A Fast and
521         Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. Molecular
522         Biology and Evolution *32*, 268-274.
523   46.   Seemann, T. (2014). Prokka: rapid prokaryotic genome annotation. Bioinformatics *30*, 2068-
524         2069.
525   47.   Lartillot, N., Rodrigue, N., Stubbs, D., and Richer, J. (2013). PhyloBayes MPI: Phylogenetic
526         Reconstruction with Infinite Mixtures of Profiles in a Parallel Environment. Systematic
527         Biology *62*, 611-615.
528   48.   Arndt, D., Grant, J.R., Marcu, A., Sajed, T., Pon, A., Liang, Y., and Wishart, D.S. (2016).
529         PHASTER: a better, faster version of the PHAST phage search tool. Nucleic Acids Research *44*,
530         W16-W21.
531   49.   Couvin, D., Bernheim, A., Toffano-Nioche, C., Touchon, M., Michalik, J., Néron, B., Rocha,
532         E.P.C., Vergnaud, G., Gautheret, D., and Pourcel, C. (2018). CRISPRCasFinder, an update of
533         CRISRFinder, includes a portable version, enhanced performance and integrates search for
534         Cas proteins. Nucleic Acids Res *46*, W246-w251.
535   50.   Touchon, M., Cury, J., Yoon, E.J., Krizova, L., Cerqueira, G.C., Murphy, C., Feldgarden, M.,
536         Wortman, J., Clermont, D., Lambert, T., et al. (2014). The genomic diversification of the
537         whole Acinetobacter genus: origins, mechanisms, and consequences. Genome Biol Evol *6*,
538         2866-2882.
539   51.   Tesson, F., Hervé, A., Mordret, E., Touchon, M., d'Humières, C., Cury, J., and Bernheim, A.
540         (2022). Systematic and quantitative view of the antiviral arsenal of prokaryotes. Nature
541         Communications *13*, 2561.
542   52.   Kosakovsky Pond, S.L., Poon, A.F.Y., Velazquez, R., Weaver, S., Hepler, N.L., Murrell, B.,
543         Shank, S.D., Magalis, B.R., Bouvier, D., Nekrutenko, A., et al. (2020). HyPhy 2.5-A
544         Customizable Platform for Evolutionary Hypothesis Testing Using Phylogenies. Mol Biol Evol
545         *37*, 295-299.

546    53.    Suyama, M., Torrents, D., and Bork, P. (2006). PAL2NAL: robust conversion of protein
547           sequence alignments into the corresponding codon alignments. Nucleic Acids Research *34*,
548           W609-W612.
549    54.    Wertheim, J.O., Murrell, B., Smith, M.D., Kosakovsky Pond, S.L., and Scheffler, K. (2015).
550           RELAX: detecting relaxed selection in a phylogenetic framework. Mol Biol Evol *32*, 820-832.
551    55.    R core development team (2013). R: A language and environment for statistical computing. .
552           (Vienna, Austria: R Foundation for Statistical Computing).
553    56.    Wickham, H. (2016). ggplot2: Elegant Graphics for Data Analysis. (New York: Springer-
554           Verlag).
555    57.    Pedersen, T. (2022). patchwork: The Composer of Plots. https://patchwork.data-
556           imaginist.com
557    58.    Csardi, G., and Nepusz, T. (2006). The igraph software package for complex network
558           research. InterJournal, Complex Systems *1695*, 1-9.
559    59.    Yu, G., Smith, D.K., Zhu, H., Guan, Y., and Lam, T.T.-Y. (2017). ggtree: an r package for
560           visualization and annotation of phylogenetic trees with their covariates and other associated
561           data. Methods in Ecology and Evolution *8*, 28-36.
562    60.    Gilchrist, C.L.M., and Chooi, Y.H. (2021). clinker & clustermap.js: automatic generation of
563           gene cluster comparison figures. Bioinformatics *37*, 2473-2475.
564