# Hymenoptera associated eukaryotic virome lacks host

# specificity

Ward Deboutte*[1], Leen Beller[1], Claude Kwe Yinda[1,2], Chenyan Shi[1], Lena Smets[1],

Bert Vanmechelen[1], Nadia Conceição-Neto[1], Kai Dallmeier[3], Piet Maes[1], Dirk C de

Graaf[4], Jelle Matthijnssens*[1]

**Author Affiliations**

[1] *KU Leuven - University of Leuven, Department of Microbiology, Immunology and Transplantation,*

*Rega Institute for Medical Research, Division of Clinical and Epidemiological Virology, 3000, Leuven,*

*Belgium*

[2] *NIAID/NIH, Rocky Mountain Laboratories, Laboratory of Virology, Virus Ecology Unit,59840, Montana,*

*USA*

[3] *KU Leuven - University of Leuven, Department of Microbiology, Immunology and Transplantation,*

*Rega Institute for Medical Research, Laboratory of Virology and Chemotherapy, 3000 Leuven, Belgium*

[4] *UGent – Ghent University, Department of Biochemistry and Microbiology, Laboratory of Molecular*

*Entomology and Bee Pathology (L-MEB),9000, Ghent, Belgium*

*Correspondence: ward.deboutte@kuleuven.be; jelle.matthijnssens@kuleuven.be

## Abstract

Recent advancements in sequencing technologies and metagenomic studies have increased the knowledge of the virosphere associated with honey bees tremendously. In this study, viral-like particle enrichment and deep sequencing was deployed to detect viral communities in managed Belgian honey bees. A substantial number of previously undescribed divergent virus genomes was detected, including a rhabdovirus and a recombinant virus possessing a divergent *Lake Sinai Virus* capsid and a Hepe-like polymerase. Furthermore, screening > 5,000 public sequencing datasets for the retrieved set of viral genomes revealed an additional plethora of undetected, divergent viruses present in a wide range of Hymenoptera species. The unexpected high number of shared viral genomes within the Apidae family and across different families within the order Hymenoptera suggests that many of these viruses are highly promiscuous, that virus sharing within and between Hymenoptera families occurs frequently, and that the concept of species-specific viral taxa inside the Hymenoptera should be revisited. In particular, this estimation implies that sharing of several viral species, thought to be specific for bees, across other eukaryotic taxa is rampant. This study provides important insights on the host taxonomical breadth of some of the known "bee viruses" and might have important implications on strategies to combat viruses that are relevant to pollinators.

## Introduction

The European honey bee (*Apis Mellifera*) forms a central hub in ecosystem maintenance, resilience and diversity. Aside from the economically valuable products, such as honey and nectar (1,2), managed bee colonies together with other insects

42   contribute tremendously to pollination (3) and play a key role in global agricultural

43   production (4). In the past decades, pressures on both managed and wild bees have

44   increased vastly and there is evidence for declining trends in pollinator populations

45   globally (5,6). These pressures encompass ecological factors such as habitat loss (7),

46   pollution (8), pesticide use (9,10) and adverse agricultural practices (11), but biological

47   factors including bacterial, parasitic, and viral infections (12–15), also play a pivotal

48   role. Recently, more attention is being given to the microbiota and their influence on

49   bee health, development and homeostasis (16–18), and it has been shown that the

50   microbiota can be exploited to protect bees from other pathogens (19). The influence

51   these factors have can be cumulative or even synergistic. For example, it has been

52   shown that pesticide use can perturb the expression of essential immunocompetence

53   genes, increasing the probability of microbial infections (20). Perhaps the best

54   example for mutual synergistic factors detrimental for bee health confine parasitic and

55   viral infections. The worldwide spread of the *Varroa destructor* parasite facilitated

56   *Deformed wing* virus (DWV) infections by acting as an active vector (where the virus

57   can replicate in both the vector and the host) (21). Parallel to its role as viral vector, it

58   has been shown that the *V. destructor* parasite can also influence the immune status

59   of its host (22). Globalization of *V. destructor* and concomitant DWV infections raised

60   the question what influence DWV plays in colony health. Recent studies have revealed

61   an association between DWV infections and colony health status (23–25). Despite the

62   worldwide dominance of DWV, other RNA viruses have been shown to be highly

63   virulent, resulting in a strong phenotype in infected bees. Acute bee paralysis virus

64   (ABPV), Black queen cell virus (BQCV) and Sacbrood virus (SBV) are all members of

65   the order *Picornavirales* that have a detrimental effect on colony health once they

66   infect a hive (26). Scattered information suggests that some of these viruses are not

67  to be restricted to honey bees, but also infect and replicate in other members of the

68  Apidae family. Spill-over events from managed honey bees into bumblebee species

69  have been described for DWV, BQCV, ABPV, SBV and Lake Sinai viruses (LSV) (27–

70  30), whereas honey bee viruses have also been described in ants (Formicidae) (31)

71  and wasps (Vespidae) (32). Recent advancements in sequencing technologies and

72  metagenomics have accelerated virus discovery in bees and a number of studies have

73  attempted to describe the viral diversity associated with bees. These studies were able

74  to expand the range of known honey bee viruses significantly and aside from

75  numerous viruses belonging to the order *Picornavirales*, numerous other RNA viruses

76  have been discovered belonging to the orders *Bunyavirales*, *Mononegavirales*

77  (containing the family *Rhabdoviridae*) and *Articulavirales* (containing the family

78  *Orthomyxoviridae*), and several unclassified RNA viruses such as LSV (33–37). DNA

79  viruses have also been described, such as Apis mellifera Filamentous virus (AmFV)

80  (38), and numerous single-stranded DNA viruses (39). While these sequencing efforts

81  have vastly increased the number of known honey bee related viruses, the relevance

82  of most of these viruses remains enigmatic. In this study, we first describe the

83  eukaryotic viruses present in > 300 Belgian bee colonies collected in the framework

84  of the EpiloBEE study (40) in 2012 and 2013. We place these results in the context of

85  other known insect viruses. Finally, by screening more than 5,000 public RNA

86  sequencing datasets, we shed light on the sharing of (bee) viruses between different

87  members of the order Hymenoptera and within the Apidae lineage.

# Results

**Eukaryotic virus identification yields previously known and unknown honey bee viruses**

Viral-like particle enrichment (41) and Illumina sequencing was performed on pooled samples derived from 300 weak and healthy (as defined by the EpiloBEE study (40)) managed honey bee colonies in Flanders, Belgium as described before (42). After sequencing and *de novo* assembly of the individual libraries, redundancy of the retrieved contigs was removed by collapsing sequences with 97% nucleotide identity over 80% of their length. Subsequently, the non-redundant contig set was annotated using DIAMOND (43) against NCBI's NR database. Viruses were taxonomically classified using the lowest-common ancestor algorithm implemented in Kronatools (44). Sequences showing similarity to bacteriophages were omitted from this analysis. Genome coverage values were obtained by mapping the sequencing reads per sample back to the non-redundant contig set. Clustering analysis on the viral coverage matrix revealed a distinct clustering pattern between samples derived from weak and healthy colonies, although with a very small biological relevance (adonis test, $R^2$ = 0.035, p-value = 0.0042) (Fig. 1A). The log-transformed coverage matrix showed that the vast majority of viral reads could be attributed to the family *Iflaviridae*, of which DWV is a member (Fig. 1B). The second most prevalent viral family was the family *Orthomyxoviridae*. Several families containing plant and fungal viruses, such as *Partitiviridae*, *Chrysoviridae*, and *Tymoviridae*, were also recovered. The clustering pattern of the coverage matrix reflected the adonis test results, showing most of the samples being dispersed by health status and although one healthy cluster containing mainly unclassified reads exists, the lack of monophyly implies no clear differences in composition with respect to the health status. In terms of absolute contig count, the

113    most prevalent orders were (apart from unclassified sequences) *Picornavirales*,

114    *Tymovirales* and *Mononegavirales* (supplemental fig. S1A) and the most prevalent

115    families were (next to unclassified sequences) *Partitiviridae*, *Comoviridae,* and

116    *Parvoviridae* (supplemental fig. S1B). There was no significant difference between the

117    number of non-redundant contigs present in healthy and weak samples (Mann-

118    Whitney U test, p-value = 0.32) (supplemental fig. S2). Only 30% of the non-redundant

119    contigs had an amino acid similarity percentage with the best hit in the NR database

120    higher than 90%, reflecting the divergent nature of the retrieved sequences

121    (supplemental fig. S3). Species accumulation curves revealed a near horizontal

122    asymptote, implying that viral sequence space was probed sufficiently (supplemental

123    fig. S4). The relatively short length of the majority of retrieved viral sequences

124    hampered a complete phylogenetic analysis (supplemental fig. S3). Therefore, an all-

125    by-all TBlastX search was conducted using the retrieved non-redundant contig set

126    complemented with a filtered viral Refseq set (see methods) as both query and bait.

127    The resulting blast output was converted into a network using sequences as vertices,

128    and hits as edges. A minimized-nested block network was constructed and visualized

129    using the taxonomical information of the reference sequences (Fig. 1C). The vast

130    majority of retrieved sequences clustered together in blocks with the order

131    *Picornavirales*, although the orders *Bunyavirales*, *Mononegavirales* and *Tymovirales*

132    were also represented substantially. Several contigs could not be assigned to any

133    known order and represented unclassified (ds)RNA viruses or unclassified (circular)
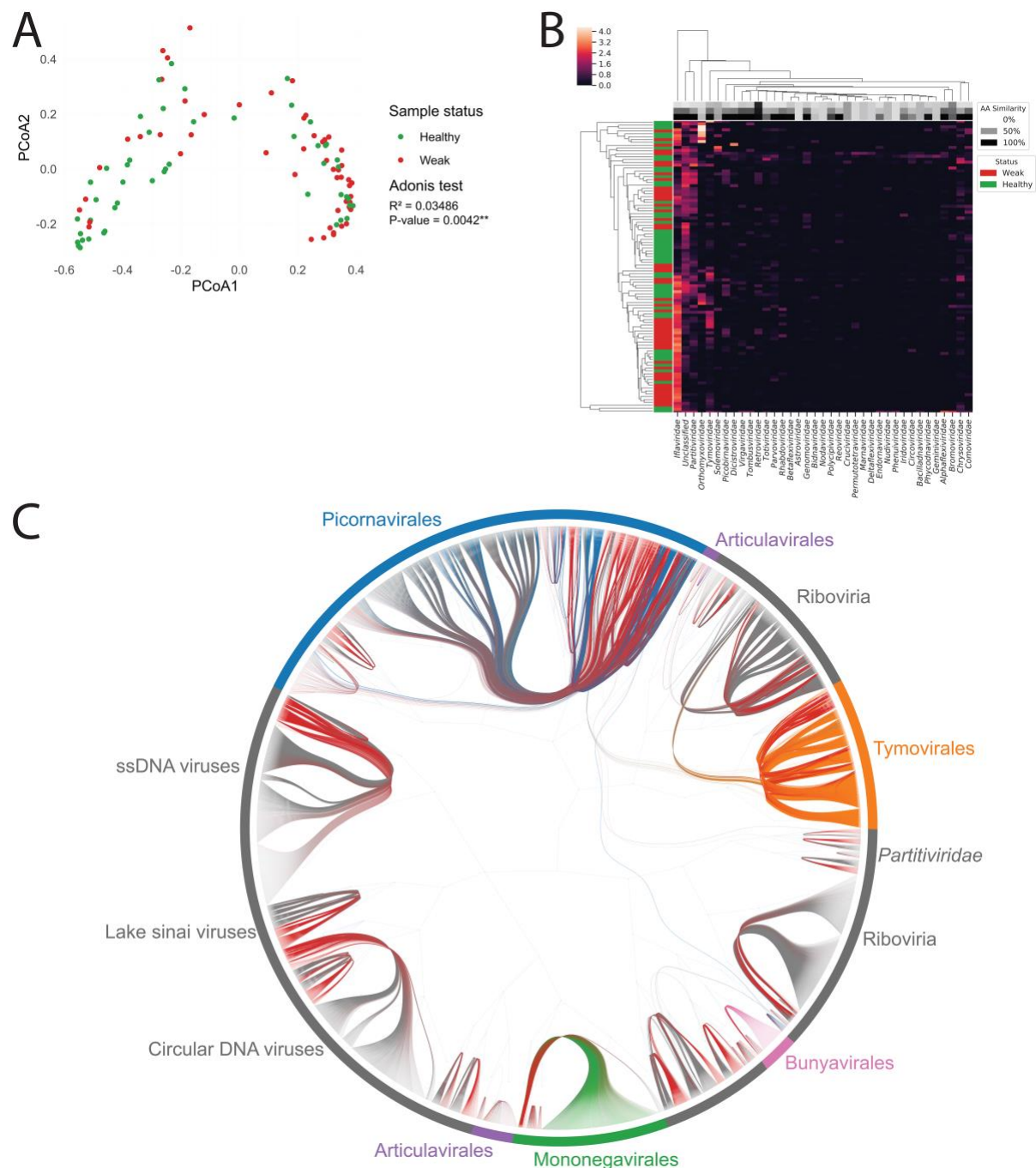
134    DNA viruses.

**Fig. 1. Belgian honey bees harbor a diverse range of known and novel viruses.**
(A) PCoA clustering using Bray-Curtis distances calculated on the viral coverage matrix derived from the Belgian samples (n = 102). Green dots reflect samples derived from healthy colonies; red dots reflect samples derived from weak colonies. The R squared and p-value obtained from the Adonis test are indicated on the right. (B) Average values per viral family of the log-transformed viral coverage matrix are depicted in a heatmap, clustered using Euclidian distances. The left column depicts samples derived from healthy (green) and weak (red) samples. The first three rows indicate the minimum (top), the average (middle) and the maximum (bottom) percentage of amino acid similarity of the contigs per viral family. (C) Minimized nested block network using retrieved sequences in this study (red) and known Refseq viruses (all other colors). Known orders are indicated in colors and unclassified reference sequences are indicated in gray.

**Phylogenetic analysis confirms the presence of known and divergent**

**eukaryotic viruses in Belgium**

To investigate the phylogenetic placement of a subset of the retrieved near-complete

viral genomes, maximum clade credibility trees (MCC) were created using BEAST (45)

(Fig. 2). Retrieved genomes from this study (orange tip labels, and listed in

supplemental table 1) and from the short-read sequencing archive (SRA, NCBI)

search (blue tip labels, see below), as well as reference sequences (green tips), were

included based on sequence length and based on a BlastP search (see methods).

Phylogenies were created for Rhabdo-like, Picorna-like, Bunya-like, Orthomyxo-like,

Sinai-like, Partiti-like, Toti-like and Tymo-like viruses. One of the retrieved rhabdo-like

viruses (Apis rhabdovirus1-Belgium) was nearly identical to the recently identified Apis

rhabdovirus 1 (34), while the other Rhabdo-like virus (Apis rhabdovirus3-Belgium) has

Diachasmimorpha longicaudata rhabdovirus as closest relative (but only had 38%

amino acid identity for the L protein). The retrieved Picorna-like viruses reflect known

bee pathogens clading in the families *Iflaviridae* and *Dicistroviridae*, such as DWV,

SBV and ABPV, but also include more divergent sequences related to Nora-like

viruses. A number of sequences clading together with plant infecting picornaviruses,

such as several comoviruses were also retrieved. The retrieved Orthomyxo-like

viruses are three closely related viruses (*Apis orthomyxovirus 1, 2* and *3-Belgium*),

clustering together with other known thogotoviruses. These three viruses are nearly

identical to the recently discovered Varroa Orthomyxovirus, with the exception of the

nucleoprotein (35). Furthermore, five LSV-like viruses were retrieved, out of which four

were very similar to other known Lake Sinai viruses (between 94% and 97% nucleotide

similarity). Interestingly, the fifth identified *Lake Sinai virus* was initially identified as an

Astro-like virus, but was shown to be a divergent recombinant virus with a 'Hepe-like'

172  polymerase region (31% amino acid similarity with the non-structural protein of Culex

173  Bastrovirus-like virus), and a Lake Sinai virus-like capsid (Lake Sinai virus, 35% amino

174  acid similarity) (supplemental fig. S5). Sequence depth profiling indicated that this

175  sequence was a true recombinant rather than an assembly artefact. The other

176  retrieved (near-) complete viral genomes were most likely plant derived eukaryotic

177  viruses, including Partiti-like viruses (24 sequences), Toti-like viruses (two sequences)

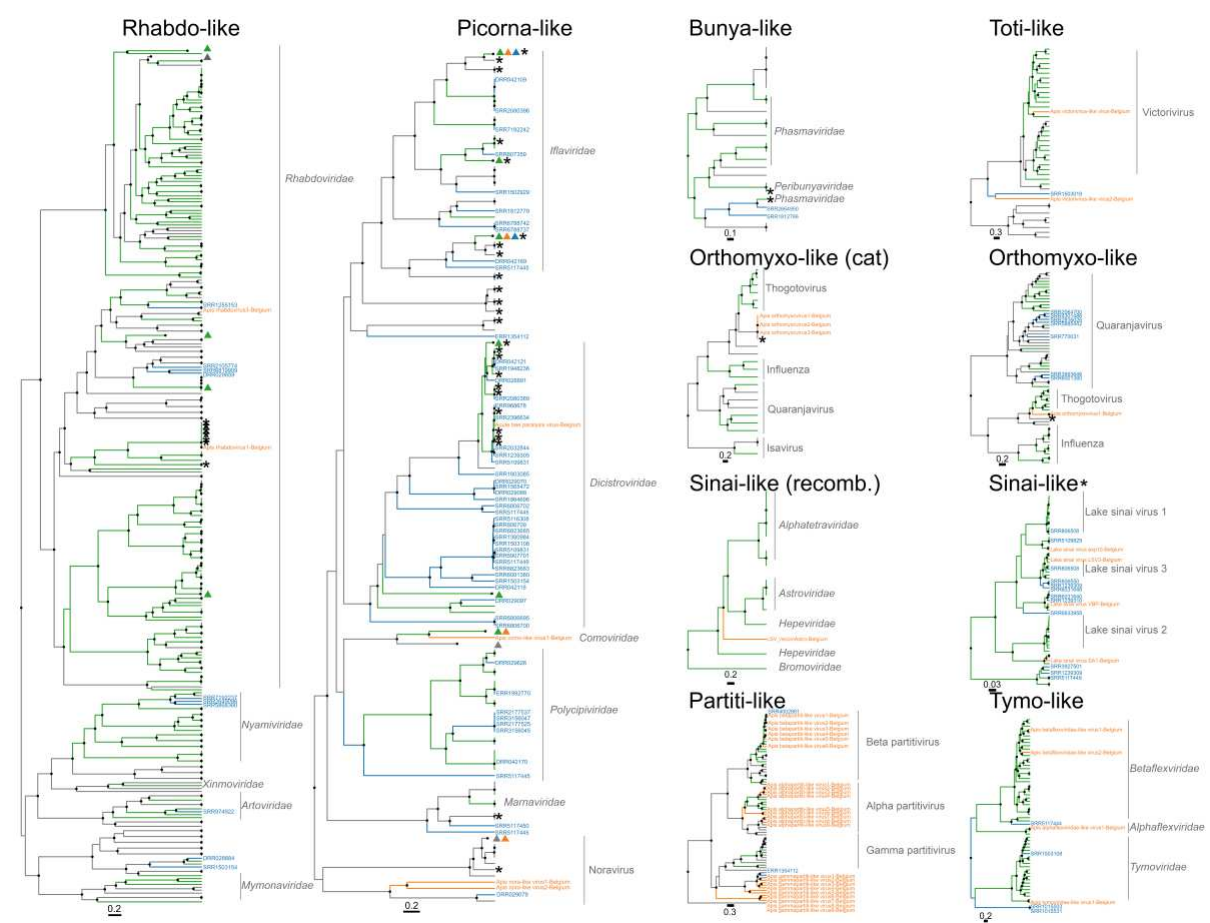178  and Tymo-like viruses (four sequences).



179

**Fig. 2. Phylogenetic analysis highlights the vast diversity of viruses identified in the Belgian samples.**
Maximum clade credibility trees for the best-represented groups of viruses retrieved in this study. Black circles on the nodes indicate posterior support values > 0.9. Viruses identified in the Belgian samples are indicated with orange tip labels, those identified through SRA searches are indicated with blue tip labels. Reference sequences that belong to a classified viral family or genus are indicated with green branches. Known honey bee viruses are indicated with an asterisk. One or more triangles indicate collapsed clades, and the colors are equivalent to the tip and branch colors. Clades that belong to the same family or genus are indicated with a gray line and in text. The 'Rhabdo-like' tree is built using the putative L protein. The 'Picorna-like' tree is built using the putative polyprotein (monocistronic viruses), the putative ORF1 (dicistronic viruses) or the putative replication polyprotein (Nora-like viruses). The 'Bunya-like' tree is built using the putative L protein. The 'Toti-like' tree and the 'Partiti-like' tree are built

192  using the putative RdRP gene. The 'Orthomyxo-like (cat)' tree is built using a concatenated protein
193  alignment of the putative PB2 – PB1 – PA – NP genes, while the 'Orthomyxo-like' tree is built using
194  only the PB2 segment. The 'Sinai-like (recomb.)' tree is built using the putative polymerase gene of the
195  astrovirus-LSV recombinant virus, while the 'Sinai-like' tree is built using the putative polymerase region
196  of all the known LSV viruses (not including the recombinant). The 'Tymo-like' tree is built using the
197  putative polyprotein gene.
198

199  **Re-screening of existing RNA sequencing datasets reveals untapped viral**

200  **diversity within the Hymenoptera lineage**

201  Since the recovered viral sequences included most of the known honey bee viral

202  sequence space (Fig. 1C), the assumption was made that the non-redundant viral

203  dataset we recovered was a good reflection of all known honey bee viruses. This

204  dataset was used as bait to map a total of 5,246 RNA sequencing datasets found in

205  the SRA database when using the query 'Hymenoptera + RNA'. A dataset was

206  considered to be 'virus enriched' when at least 100,000 reads mapped to the bait set.

207  All datasets that met this criterium (1,331) were individually *de novo* assembled using

208  SKESA (46) and viral sequences were identified and clustered as was described for

209  the Belgian samples. An additional clustering step was performed, collapsing the non-

210  redundant SRA-derived sequences together with the non-redundant Belgian

211  sequence dataset. This resulted in the recovery of nearly 10,000 non-redundant

212  putative viral contigs, out of which only 42.8% had an amino acid similarity with

213  proteins in Genbank higher than 90% (supplemental fig. S6). Forward model selection

214  analysis revealed that together, putative host taxonomy and location of the dataset

215  could explain 33% of the variability observed within the coverage matrix (Fig. 3A). This

216  result was further validated by the observation that hierarchical clustering on Euclidian

217  distances revealed clusters of both eukaryotic host families and location within the

218  coverage matrix (Fig. 3B). Viral taxonomy analysis revealed that the majority of the

219  recovered viruses could be assigned to the orders *Picornavirales* and

220  *Mononegavirales* (Fig. 3C). The retrieved viral contigs that fell below the

221   abovementioned threshold of 90% amino acid similarity were included in the

222   phylogenetic analysis and revealed ten previously undescribed Rhabdo-like viruses,

223   and more than 50 previously undescribed Picorna-like viruses (Fig. 2, blue tip labels).

224   Both these groups span multiple viral families. Another striking finding was the fact

225   that seven previously undescribed PB2 segments of *Orthomyxoviridae*-like sequences

226   were recovered (most closely related to the *Quaranjavirus* genus), indicating that this

227   viral family is more strongly represented within the Hymenoptera lineage than was

228   previously known. Furthermore, also Bunya-like, Toti-like, Sinai-like, Partiti-like and

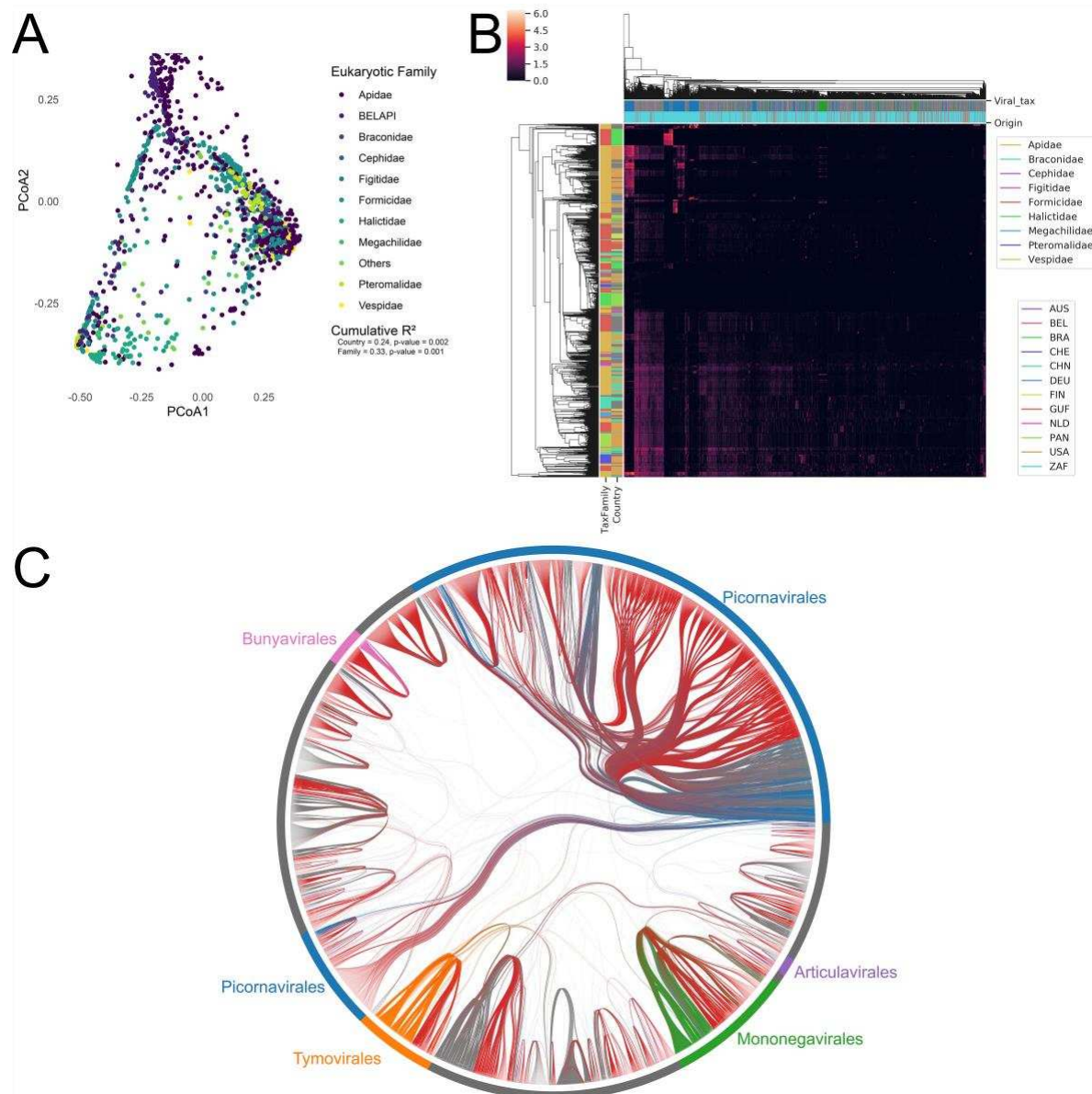229   Tymo-like viruses were recovered (Fig. 2, blue tip labels).

230

**Fig. 3. SRA searches shed light on the hymenoptera virosphere and reveal the wealth of undescribed viral sequences present in public datasets.**
(A) PCoA clustering using Bray-Curtis distances calculated on the viral coverage matrix derived from the Belgian samples clustered with the SRA screening results. Dots are colored per hymenoptera family. The cumulative $R^2$ values reported are calculated by forward model selection using the OrdiR2step function after distance-based redundancy analysis. (B) Heatmap depiction of the log transformed viral coverage matrix, clustered using euclidian distances. Leftmost columns indicate the hymenoptera families and the geographical location of the samples. The top two rows indicate the viral taxonomical classification (with the same color per viral order as Fig. 3C) and the origin of the viral sequence (light blue indicates an SRA sample as origin, red indicates viral sequences found in the Belgian samples). (C) Minimized nested block network using the non-redundant sequences retrieved from the SRA searches (red) and known Refseq viruses (all other colors).

243

**Virus-sharing networks show a large number of virus sharing within the order of Hymenoptera**

Although the discrimination between samples based on their location and eukaryotic taxonomy was significant (p-value 0.002 and 0.001, respectively) (Fig. 3A). The cumulative $R_2$ value (0.33) indicates that a large majority of the variances within the datasets cannot be explained by aforementioned variables. This observation could imply that a large number of viruses are shared across hymenoptera families, that the variance within eukaryotic families in a specific country is large relative to the variance between these parameters, or a combination of both. To investigate the first possibility, the assumption was made that the host of a specific virus sequence was that of the sample of which the sequence cluster representative was derived. Hymenoptera families of which less than ten samples were obtained were grouped together into an 'Others' group and virus sharing was calculated in a pair-wise manner for all the possible combinations within the eukaryotic host families, and within the Apidae lineage. A substantial number of viruses were found to be present not only within the Apidae lineage but also shared over multiple eukaryotic host families (Fig. 4 A, B). Within the family Apidae, most viral sequences were shared between *Apis Mellifera* and *Lepidotrigona* species (1,050 viral sequences shared), between *Apis Mellifera* and *Apis florea* (938 sequences shared), and between *Apis Mellifera* and *Ceratina* species (729 sequences shared) (Fig. 4A). The majority of these shared sequences could be traced back to the order *Picornavirales*, with a total of 224 (21.3%), 497 (53.0%) and 130 Picorna-like sequences (17.8%) shared between these groups, respectively (Fig. 4C, blue edges). Beyond the family Apidae, substantial virus sharing was detected between the families Apidae and Pteromalidae (1,066 viral sequenced shared), the families Apidae and Cephidae (742 sequences shared), and the families Apidae and

269    Braconidae (737 sequences shared). Concomitant with the situation between different

270    Apidae species, the majority of shared viral sequences could be assigned to the order

271    *Picornavirales*, with 201 (18.8%), 137 (18.4%) and 111 Picorna-like sequences

272    (15.0%) shared between these groups, respectively (Fig. 4D, blue edges). Aside from

273    Picorna-like sequences, evidence could also be found for sharing of viruses predicted

274    to belong to the orders *Mononegavirales* (Fig. 4C,D, green edges) and *Tymovirales*

275    (Fig. 4C,D, orange lines), although the number of shared viral sequences was on

276    general an order of magnitude lower than those of the *Picornavirales* (39 Mononega-

277    like viral sequences shared between Formicidae and Pteromalidae, and 27 Tymo-like

278    viral sequences shared between *Apis Mellifera* and *Lepidotrogona*). Since a fraction

279    of the recovered viruses are most likely infecting plants or reflect viruses not relevant

280    for bees (Fig. 2), an additional analysis was ran with a number of the retrieved, near-

281    complete, known bee viruses (AMFV, ABPV, BQCV, Kashmir Bee virus (KBV), DWV,

282    LSV, *Apis Rhabdovirus* and *Apis Orthomyxovirus*), as well as the retrieved Nora-like

283    viruses and other Orthomyxo-like viruses. Calculation of the fraction of positive

284    samples revealed that most of the previously thought bee-specific viruses occur in

285    multiple Apidae species but are also found within other Hymenopteran families

286    (supplemental fig. S7). An attempt was made to quantify the host specificity of these

287    viruses by calculating an Apidae specificity index (ASI), and an *Apis Mellifera*

288    specificity index (AMSI) (Table 1). These indices revealed that some of the established

289    bee viruses (ABPV, AMFV, BQCV and Quaranja-like orthomyxoviruses) show a low

290    specificity for *Apis Mellifera* within the Apidae family (characterized by a low AMSI),

291    and (with the exception of BQCV) were not restricted within the family Apidae

292    (characterized by a low ASI). Other "established honey bee species" were shown to

293    be highly specific for *Apis Mellifera*, and revealed a high AMSI (KBV, DWV and LSV).

294　　The recently discovered Apis rhabdoviruses and Nora-like viruses are found

295　　exclusively in *Apis Mellifera* within the family Apidae. The Apis rhabdoviruses are

296　　restricted within the family Apidae, but the retrieved Nora-like viruses are also highly

297　　prevalent in other Hymenoptera families (ASI 0.01, Table 1). Finally, the retrieved

298　　Quaranja orthomyxo-like viruses were highly prevalent in other Hymenoptera families

299　　and only to a limited extent in *Apis Mellifera* and the family Apidae (ASI of 0.04 and

300　　AMSI of 0.05, respectively). On the other hand, *Apis Orthomyxovirus 1* was slightly

301　　more honey bee specific, with an ASI and AMSI of 0.31 and 0.21, respectively.
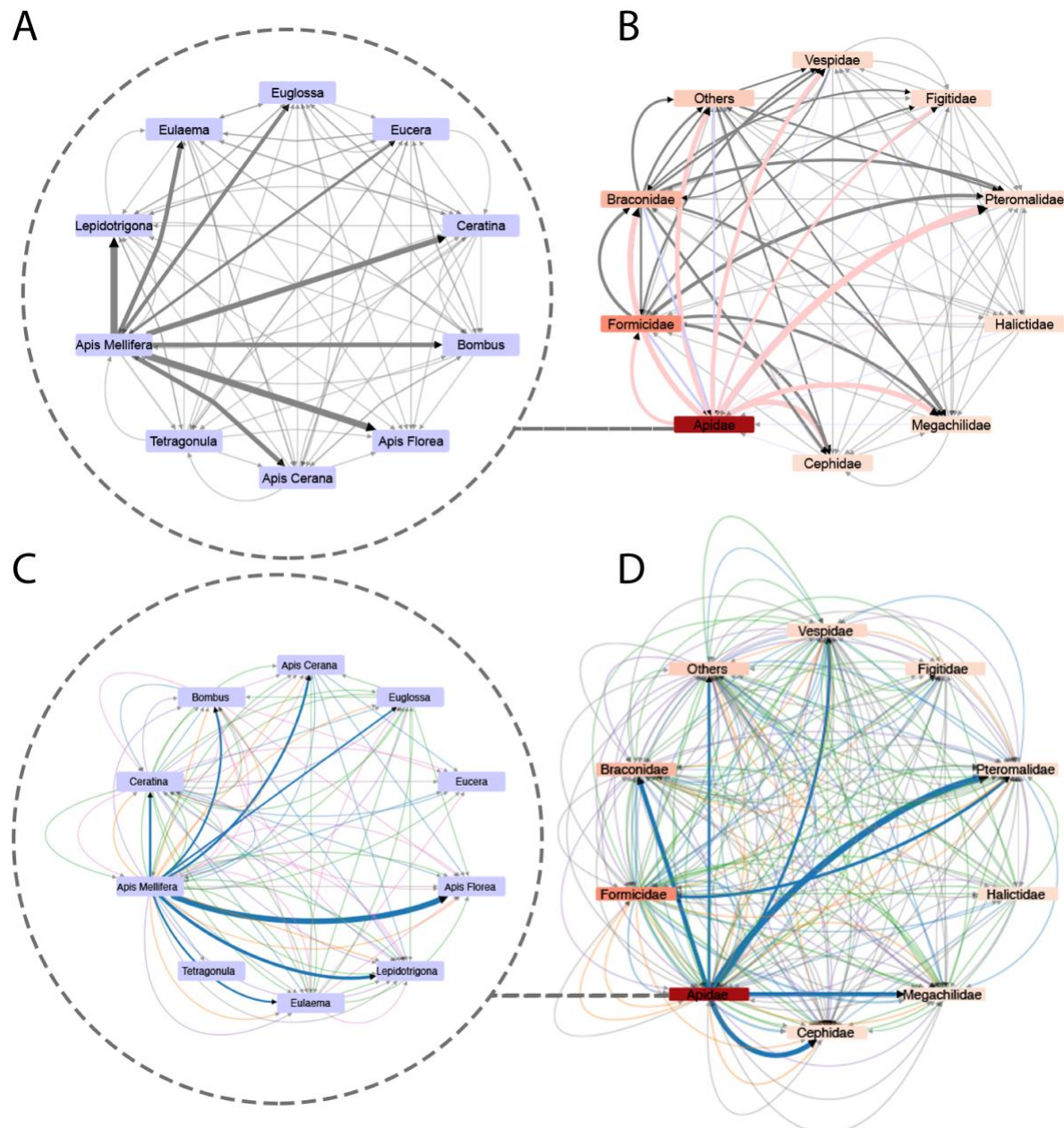
**Fig. 4. Cumulative viral sequence sharing network reflect the aspecificity of hymenoptera associated viruses.**

Networks reflecting the cumulative sharing of viral contigs between eukaryotic lineages. The networks inside the dashed circle (A,C) reflect viral sequence sharing within the family Apidae. The networks on the right (B,D) reflect sharing over different families within the order hymenoptera. Panels A and C and panels B and D both reflect the same networks, but both panels C and D reflect cumulative shared viral sequences broken up per assigned viral order (using the same color code as fig. 3C). Nodes in panels B and D are colored by number of representative virus contigs per eukaryotic lineage (ranging from 18 contigs (Cephidae) to 5,662 contigs (Apidae)). Edge thickness reflects the total shared contig count, ranging from 1 to 1,050 contigs (panel A), from 1 to 1,066 contigs (panel B), from 1 to 497 contigs (panel C), and from 1 to 201 contigs (panel D). Edge arrows indicate directionality, of which the root is the predicted host (the taxonomical group of which the virus sequence representative was derived from).

317 **Table 1. Host (a)specificity of a selection of known bee viruses.**
318 The values reflecting how specific a known bee virus is for *Apis mellifera* (AMSI) and for Apidae (ASI).
319 A value of 1 reflects complete lineage restriction. The number of viral contigs included per viral species
320 is indicated with Contig number.
321

| Virus | Virus Abbreviation | Contig number | ASI | AMSI |
|---|---|---|---|---|
| **Acute Bee Paralysis virus** | ABPV | 7 | 0.13 | 0.04 |
| **Apis Mellifera Filamentous virus** | AMFV | 8 | 0.06 | 0.05 |
| **Kashmir Bee virus** | KBV | 2 | 0.01 | 1.00 |
| **Deformed Wing virus** | DWV | 11 | 0.84 | 0.71 |
| **Black Queen Cell virus** | BQCV | 7 | 0.88 | 0.09 |
| **Lake Sinai viruses** | Sinaiviruses | 20 | 0.07 | 1.00 |
| **Apis Rhabdovirus (1 and 2)** | Rhabdo | 4 | 1.00 | 1.00 |
| **Quaranja-like Orthomyxoviruses** | Quaranja | 15 | 0.05 | 0.04 |
| **Apis Orthomyxovirus 1** | Thogoto | 3 | 0.31 | 0.21 |
| **Nora-like viruses** | Nora | 3 | 0.01 | 1.00 |

322

# Discussion

324 This study, combined with other recent sequencing efforts, provides new insights into

325 known and previously undescribed viruses associated with *Apis Mellifera*. Variance

326 analysis revealed a significant, but biologically limited difference in the viral

327 composition between weak and healthy colonies, and no significant difference in the

328 total number of viral sequences derived from healthy and weak colonies could be

329 detected. Genomes from a large number of viral families could be retrieved, of which

330 a substantial part most likely includes plant viruses. While it cannot be excluded that

331 some of the recovered divergent plant viruses constitute viruses actually infecting the

332 bee, it is likely that the majority of these viruses reflect environmental contaminants.

333 The host of the most closely related viral sequence can give an indication if these

334 sequences are environmental contamination. The fact that numerous viral sequences

335  belonging to families solely infecting plants were recovered in a large scale viral

336  discovery study in insects indicates that this assumption does not necessarily hold

337  true (47). The recent detection of viruses belonging to plant-specific viral families in

338  mosquitoes reinforces this observation (48). Some of the retrieved viral sequences

339  were very similar to recently discovered viruses (*Apis Rhabdovirus 1*, *Apis*

340  *Orthomyxovirus 1*), increasing the likelihood that these are true honey bee viruses,

341  and further confirming their presence in Belgium. Interestingly, a divergent

342  recombinant Lake Sinai virus was found, comprised of a Hepe-like polymerase region,

343  and a divergent Lake Sinai virus capsid. A novel divergent rhabdovirus (*Apis*

344  *Rhabdovirus 3*) was also described. Additionally, full genomes for an orthomyxovirus

345  (*Apis Orthomyxovirus 1*), very similar to a virus from a previous study from Levin *et al.*

346  (35), was found in multiple individual libraries, and evidence for the presence of this

347  virus was found in other Hymenoptera families (Table 1). Multiple sequence

348  alignment-free network analysis implied that, despite the species accumulation curves

349  reaching a plateau, many of the putative viral sequences retrieved were too

350  fragmented to be included in a phylogenetic analysis (the number of sequences that

351  made the threshold to be included in phylogenies was 188, while the network reflected

352  5,224 retrieved sequences). Furthermore, this analysis implies that the actual viral

353  diversity exceeds what can be captured by regular phylogenetic analysis. Larger

354  sample sizes and especially deeper sequencing efforts could help to fully elucidate

355  the viral diversity associated with honey bees. Since the retrieved non-redundant viral

356  sequence set encapsulates nearly all of the known and even more recently described

357  viruses, this set was used to probe pre-existing Hymenopteran sequencing datasets

358  for any bee-related viral signal. A total of 1,331 virus-rich RNA sequencing datasets

359  were *de novo* assembled and screened for viruses. This approach revealed that these

360    datasets harbor a substantial number of viruses that have been previously described

361    (roughly 40%), but also that the amount of undescribed, divergent viruses is rampant.

362    In concordance with the previous results, the viral sequences retrieved from the SRA

363    search also suffer from fragmentation and incomplete sequencing. This observation is

364    most likely the result of the fact that most RNA sequencing datasets included in the

365    SRA search are transcriptome studies rather than metagenomic analyses, and that for

366    most of them no wet-lab procedures for microbial or viral enrichment were performed.

367    Despite this setback, multiple-sequence alignment free network analysis implied a

368    massive hidden viral diversity within the Hymenoptera lineage (roughly 60% of the

369    retrieved contigs were less than 90% similar to any other known virus in Genbank).

370    Constrained ordination analysis showed that both the geographical origin and the

371    taxonomical lineage of the host organism sequenced could explain a biologically

372    relevant proportion (cumulative $R_2 = 0.33$) of the variance within the viral coverage

373    matrix. Since the included samples constitute a wide range of taxonomical host

374    lineages, this proportion was below expectations and implies a substantial amount of

375    viral sequences to be shared over eukaryotic Apidae species and Hymenoptera

376    families. This hypothesis was confirmed by cumulative counting of the viral sequences

377    over the different lineages included, based on a rather rigorous coverage threshold for

378    presence/absence. This analysis revealed a non-trivial number of viral sequences,

379    spanning all of the viral orders previously associated with honey bees, being shared

380    across different lineages within the Apidae, but also over other families belonging to

381    the Hymenoptera. Of all the sequences present in the total non-redundant viral

382    dataset, 53% were shared with another taxonomical lineage (5139 shared sequences,

383    9655 in total). Since the included SRA dataset suffers strongly from sampling bias,

384    this percentage is most likely an underestimation. Given this strikingly high number of

385    virus sharing, the dataset was revisited with a subset of previously described honey

386    bee specific viruses. Surprisingly, none of the tested viruses were lineage restricted to

387    *Apis Mellifera*, with the exception of Apis rhabdovirus. Other viruses, such as Nora-

388    like viruses, KBV and LSV were restricted to *Apis Mellifera* within the Apidae lineage

389    (AMSI = 1.00) but were underrepresented relative to non-Apidae Hymenopteran

390    families. The only viruses that were bee specific, *i.e.* having both a high AMSI and

391    ASI, were DWV and Apis rhabdoviruses. These results imply that despite the recent

392    sequencing efforts, many unknowns remain on viral diversity within the Hymenoptera

393    lineage. Finally, the concept of "honey bee specific viruses" should be revisited, since

394    most of the previously described viruses are not bee specific, neither are they

395    restricted to the Apidae lineage.

396

## Methods

397

### Data and code availability

398

399    All relevant (intermediate) output files, metadata tables, fasta sequences, R code,

400    Python code and jupyter notebooks are available on Github through the URL

401    https://github.com/Matthijnssenslab/Bee_euvir. Intermediary output files too large to

402    be hosted on Github are available through Zenodo (10.5281/zenodo.3979324). The

403    raw sequencing data is available through the SRA database under project accession

404    PRJNA579886. Accession numbers of the viral sequences included in the phylogenies

405    will be made available in supplemental table S1. Accession information for the public

406    datasets screened in this study are available in supplemental table S2.

407

408 **Sample preparation, pooling, VLP-sequencing and read processing**

409 Samples were pooled and prepared for Illumina sequencing as described before, and

410 the prokaryotic viruses in these pools were described previously (42). Briefly, samples

411 were taken from the Flanders EpiloBEE study (40), from both sampling years (2012

412 and 2013), and 102 pools were constructed based on health status (defined

413 retrospectively within the EpiloBEE study, with "strong" hives surviving winter and

414 "weak" hives not surviving winter), subspecies and geographical location. Pooling

415 information and SRA accession numbers were described before (42). After

416 sequencing, reads were quality controlled using Trimmomatic (49), version 0.38.

417 Subsequently, *de novo* assemblies were made for the individual libraries using

418 SPAdes (50), version 3.12.0, with kmer sizes 21, 33, 55 and 77 in the metagenomic

419 mode. To remove redundancy, the resulting contigs larger than 500 bp were collapsed

420 if they showed 97% nucleotide identity over at least 80% of the contig lengths, using

421 ClusterGenomes (https://bitbucket.org/MAVERICLab/docker-clustergenomes).

422 Putative eukaryotic viruses were identified using the BlastX method implemented in

423 DIAMOND (43) version 0.9.22, using the 'c 1' and 'sensitive' flags, against the NR

424 database (NCBI), downloaded on 30 september 2018. Taxonomical paths were

425 parsed with the KtClassifyBLAST algorithm implemented in Kronatools (44). All

426 contigs that fell under taxID '10239' (Viruses) were included in the analysis. Contigs

427 that could be annotated as bacteriophages (as described before (42)) were excluded.

428 Coverage values per sample were obtained by mapping the reads per sample back to

429 the viral dataset, using BWA-mem version 0.7.16a (51), filtering the obtained

430 alignments for an identity of 97% over a coverage of 70% using BAMM

431 (https://github.com/Ecogenomics/BamM). Coverage values were calculated by

432 dividing the readcounts per contig by the contig length.

433

**SRA searches**

The SRA database was searched by using the query 'Hymenoptera + RNA', and the resulting 5,246 fastQ files were retrieved by using the prefetch and fastq-dump tools implemented in the SRA toolkit (NCBI). The previously obtained viral dataset was used as an index and retrieved fastQ files were mapped back using BWA-mem (51), version 0.7.16a. Only samples that had a cumulative read count of at least 100,000 reads (1,331 samples) were included downstream. Samples were then *de novo* assembled using SKESA (46) and annotated and clustered as described above. Information on the included samples is provided in supplemental table 2.

**Phylogenetic analysis**

Viral sequences were included based on an *ad hoc* determined length cut-off depending on the expected genome length of each virus (supplemental table S3). Reference sequences were included by using the retrieved viral sequences as query and performing a TBlastX search (52) with an e-value cutoff of 1E-10 against the nt database (NCBI), downloaded on 1 october 2019. For the Partiti-like, Tymo-like and Toti-like trees only Refseq sequences were included. Significant hits were also filtered on the abovementioned alignment length cut-off specific for a viral group (supplemental table 3). Next, proteins were predicted from both the queries and the significant hits, using prodigal (53), version 2.6.3. Predicted proteins were submitted to an all-to-all BlastP search, with an e-value cut-off of 1E-10. The output was then transformed into a network and the largest connected component was extracted using the networkx library (54) implemented in Python. Proteins within the largest connected component were subsequently aligned with MAFFT (55), version 7.313, using the L-INS-I setting and trimmed using trimAL (56), version 1.4.1, using the gappyout setting.

458    Model selection was performed using Prottest (57), version 3.4.2. Bayesian

459    phylogenetic analysis was performed using BEAST (45), version 1.10.4, using the

460    predicted protein models (supplemental table 3) under a strict clock and constant

461    population size prior. The respective analysis was ran until all the effective samples

462    sizes were above 200, and maximum clade credibility trees were calculated using

463    TreeAnnotator, implemented in the BEAST package. Final trees were plotted in R

464    using the ggtree package (58).

465    **Network and contig sharing analysis**

466    Networks were created from the retrieved viral sequence data by using TBlastX

467    against the Refseq nt database, downloaded on 1 october 2019. The Refseq database

468    was filtered by removing entries containing the keyword 'phage' (for bacteriophages)

469    or 'herpes' in the header, and by removing sequences longer than 15000 nt and

470    shorter than 500 nt. These cut-offs were implemented to reduce 'noisy' hits, where for

471    example herpes polymerases have significant hits to other viral polymerases. The

472    remaining sequences were clustered on 80% nucleotide identity over 80% of the

473    length, by using CDhit, version 4.8.1 (59). The tBlastX search was performed with an

474    E-value cutoff of 1E-10 and an alignment length cut-off of 300 positions, and was ran

475    in two iterations to include reference sequences that only made the cut-off when

476    aligning to other reference sequences. The resulting blast output was then converted

477    into a minimized nested block network, using the graph-tool package (60),

478    implemented in Python. Virus sharing over the eukaryotic families belonging to the

479    Hymenoptera and within the Apidae was determined by using the coverage matrix. A

480    viral sequence was assumed to originate from the taxonomical lineage of the sample

481    of which the cluster representative (the longest contig inside a cluster) was derived, in

482    order to determine directionality. A viral representative sequence was assumed to be

483    present in a sample when the coverage was above 0.1. For the cumulative virus

484    sharing, an additional threshold was imposed were at least 10% of the included

485    samples of a specific taxonomical host lineage had to be positive before the viral

486    sequence was assumed to be present within that lineage. Resulting networks were

487    visualized in Cytoscape (61), version 3.7.1. Percentages of positive samples were

488    calculated using the same relative count cutoff as mentioned before and the ASI and

489    AMSI were calculated by taking the ratio of the fraction of positive samples for a

490    specific bee virus within Apidae or *Apis Mellifera* samples, divided by the fraction of

491    positive samples in other eukaryotic families or other *Apidae* species, respectively.

492    **QUANTIFICATION AND STATISTICAL ANALYSIS**

493    PCoA analysis was performed in R (62) version 3.5.3, with the pcoa function

494    implemented in the 'ape' library (63). Variance analysis and distance-based

495    redundancy analysis was performed on the coverage matrix using Bray-Curtis

496    distances, using the adonis test and capscale function implemented in vegan (64).

497    Cumulative explanation power of the location (country of origin) and eukaryotic

498    taxonomy (on family level) covariates was calculated using the ordiR2 function

499    (vegan). The difference in absolute numbers of contigs was calculated using the

500    Mann-Whitney U test implemented in scipy (65), in Python.

501
502    **References**

503  1.    vanEngelsdorp D, Meixner MD. A historical review of managed honey bee populations
504        in Europe and the United States and the factors that may affect them. J Invertebr
505        Pathol. 2010 Jan;103:S80–95.

506  2.    Gallai N, Salles J-M, Settele J, Vaissière BE. Economic valuation of the vulnerability of
507        world agriculture confronted with pollinator decline. Ecol Econ. 2009 Jan;68(3):810–21.

508  3.    Garibaldi LA, Steffan-Dewenter I, Winfree R, Aizen MA, Bommarco R, Cunningham SA,
509        et al. Wild Pollinators Enhance Fruit Set of Crops Regardless of Honey Bee Abundance.
510        Science. 2013 Mar 29;339(6127):1608–11.

511    4.    Aizen MA, Garibaldi LA, Cunningham SA, Klein AM. How much does agriculture depend
512           on pollinators? Lessons from long-term trends in crop production. Ann Bot. 2009
513           Jun;103(9):1579–88.

514    5.    Hallmann CA, Sorg M, Jongejans E, Siepel H, Hofland N, Schwan H, et al. More than 75
515           percent decline over 27 years in total flying insect biomass in protected areas. Lamb
516           EG, editor. PLOS ONE. 2017 Oct 18;12(10):e0185809.

517    6.    Potts SG, Biesmeijer JC, Kremen C, Neumann P, Schweiger O, Kunin WE. Global
518           pollinator declines: trends, impacts and drivers. Trends Ecol Evol. 2010 Jun;25(6):345–
519           53.

520    7.    Foley JA. Global Consequences of Land Use. Science. 2005 Jul 22;309(5734):570–4.

521    8.    Lusebrink I, Girling RD, Farthing E, Newman TA, Jackson CW, Poppy GM. The Effects of
522           Diesel Exhaust Pollution on Floral Volatiles and the Consequences for Honey Bee
523           Olfaction. J Chem Ecol. 2015 Oct;41(10):904–12.

524    9.    Henry M, Beguin M, Requier F, Rollin O, Odoux J-F, Aupinel P, et al. A Common
525           Pesticide Decreases Foraging Success and Survival in Honey Bees. Science. 2012 Apr
526           20;336(6079):348–50.

527    10.    Decourtye A, Devillers J, Genecque E, Le Menach K, Budzinski H, Cluzeau S, et al.
528           Comparative sublethal toxicity of nine pesticides on olfactory learning performances of
529           the honeybee Apis mellifera. Arch Environ Contam Toxicol. 2005 Feb;48(2):242–50.

530    11.    Nicholls CI, Altieri MA. Plant biodiversity enhances bees and other insect pollinators in
531           agroecosystems. A review. Agron Sustain Dev. 2013 Apr;33(2):257–74.

532    12.    Fünfhaus A, Ebeling J, Genersch E. Bacterial pathogens of bees. Curr Opin Insect Sci.
533           2018 Apr;26:89–96.

534    13.    Forfert N, Natsopoulou ME, Frey E, Rosenkranz P, Paxton RJ, Moritz RFA. Parasites and
535           Pathogens of the Honeybee (Apis mellifera) and Their Influence on Inter-Colonial
536           Transmission. Rueppell O, editor. PLOS ONE. 2015 Oct 9;10(10):e0140337.

537    14.    Grozinger CM, Flenniken ML. Bee Viruses: Ecology, Pathogenicity, and Impacts. Annu
538           Rev Entomol. 2019 Jan 7;64(1):205–26.

539    15.    McMenamin AJ, Flenniken ML. Recently identified bee viruses and their impact on bee
540           pollinators. Curr Opin Insect Sci. 2018 Apr;26:120–9.

541    16.    Engel P, Kwong WK, McFrederick Q, Anderson KE, Barribeau SM, Chandler JA, et al. The
542           Bee Microbiome: Impact on Bee Health and Model for Evolution and Ecology of Host-
543           Microbe Interactions. mBio [Internet]. 2016 May 4 [cited 2019 Mar 15];7(2). Available
544           from: http://mbio.asm.org/lookup/doi/10.1128/mBio.02164-15

545    17.    Kešnerová L, Emery O, Troilo M, Liberti J, Erkosar B, Engel P. Gut microbiota structure
546           differs between honeybees in winter and summer. ISME J. 2020 Mar;14(3):801–14.

547  18.  Raymann K, Moran NA. The role of the gut microbiome in health and disease of adult
548      honey bee workers. Curr Opin Insect Sci. 2018 Apr;26:97–104.

549  19.  Leonard SP, Powell JE, Perutka J, Geng P, Heckmann LC, Horak RD, et al. Engineered
550      symbionts activate honey bee immunity and limit pathogens. Science. 2020 Jan
551      31;367(6477):573–6.

552  20.  Grassl J, Holt S, Cremen N, Peso M, Hahne D, Baer B. Synergistic effects of pathogen
553      and pesticide exposure on honey bee (Apis mellifera) survival and immunity. J
554      Invertebr Pathol. 2018 Nov;159:78–86.

555  21.  Villalobos EM. The mite that jumped, the bee that traveled, the disease that followed.
556      Science. 2016 Feb 5;351(6273):554–6.

557  22.  Annoscia D, Brown SP, Di Prisco G, De Paoli E, Del Fabbro S, Frizzera D, et al.
558      Haemolymph removal by *Varroa* mite destabilizes the dynamical interaction between
559      immune effectors and virus in bees, as predicted by Volterra's model. Proc R Soc B Biol
560      Sci. 2019 Apr 24;286(1901):20190331.

561  23.  Benaets K, Van Geystelen A, Cardoen D, De Smet L, de Graaf DC, Schoofs L, et al. Covert
562      deformed wing virus infections have long-term deleterious effects on honeybee
563      foraging and survival. Proc R Soc B Biol Sci. 2017 Feb 8;284(1848):20162149.

564  24.  Natsopoulou ME, McMahon DP, Doublet V, Frey E, Rosenkranz P, Paxton RJ. The
565      virulent, emerging genotype B of Deformed wing virus is closely linked to overwinter
566      honeybee worker loss. Sci Rep. 2017 Dec;7(1):5242.

567  25.  Tehel A, Vu Q, Bigot D, Gogol-Döring A, Koch P, Jenkins C, et al. The Two Prevalent
568      Genotypes of an Emerging Infectious Disease, Deformed Wing Virus, Cause Equally Low
569      Pupal Mortality and Equally High Wing Deformities in Host Honey Bees. Viruses. 2019
570      Jan 29;11(2):114.

571  26.  McMenamin AJ, Genersch E. Honey bee colony losses and associated viruses. Curr Opin
572      Insect Sci. 2015 Apr;8:121–9.

573  27.  Alger SA, Burnham PA, Boncristiani HF, Brody AK. RNA virus spillover from managed
574      honeybees (Apis mellifera) to wild bumblebees (Bombus spp.). Rueppell O, editor.
575      PLOS ONE. 2019 Jun 26;14(6):e0217822.

576  28.  Fürst MA, McMahon DP, Osborne JL, Paxton RJ, Brown MJF. Disease associations
577      between honeybees and bumblebees as a threat to wild pollinators. Nature. 2014
578      Feb;506(7488):364–6.

579  29.  Genersch E, Yue C, Fries I, de Miranda JR. Detection of Deformed wing virus, a honey
580      bee viral pathogen, in bumble bees (Bombus terrestris and Bombus pascuorum) with
581      wing deformities. J Invertebr Pathol. 2006 Jan;91(1):61–3.

582  30.  Dolezal AG, Hendrix SD, Scavo NA, Carrillo-Tripp J, Harris MA, Wheelock MJ, et al.
583       Honey Bee Viruses in Wild Bees: Viral Prevalence, Loads, and Experimental Inoculation.
584       Rueppell O, editor. PLOS ONE. 2016 Nov 10;11(11):e0166190.

585  31.  Sébastien A, Lester PJ, Hall RJ, Wang J, Moore NE, Gruber MAM. Invasive ants carry
586       novel viruses in their new range and form reservoirs for a honeybee pathogen. Biol
587       Lett. 2015 Sep 30;11(9):20150610.

588  32.  Mordecai GJ, Brettell LE, Pachori P, Villalobos EM, Martin SJ, Jones IM, et al. Moku
589       virus; a new Iflavirus found in wasps, honey bees and Varroa. Sci Rep. 2016
590       Dec;6(1):34983.

591  33.  Galbraith DA, Fuller ZL, Ray AM, Brockmann A, Frazier M, Gikungu MW, et al.
592       Investigating the viral ecology of global bee communities with high-throughput
593       metagenomics. Sci Rep [Internet]. 2018 Dec [cited 2019 Mar 15];8(1). Available from:
594       http://www.nature.com/articles/s41598-018-27164-z

595  34.  Remnant EJ, Shi M, Buchmann G, Blacquière T, Holmes EC, Beekman M, et al. A Diverse
596       Range of Novel RNA Viruses in Geographically Distinct Honey Bee Populations. Ross SR,
597       editor. J Virol [Internet]. 2017 Aug 15 [cited 2019 Mar 15];91(16). Available from:
598       http://jvi.asm.org/lookup/doi/10.1128/JVI.00158-17

599  35.  Levin S, Sela N, Erez T, Nestel D, Pettis J, Neumann P, et al. New Viruses from the
600       Ectoparasite Mite Varroa destructor Infesting Apis mellifera and Apis cerana. Viruses.
601       2019 Jan 24;11(2):94.

602  36.  Daughenbaugh K, Martin M, Brutscher L, Cavigli I, Garcia E, Lavin M, et al. Honey Bee
603       Infecting Lake Sinai Viruses. Viruses. 2015 Jun 23;7(6):3285–309.

604  37.  Schoonvaere K, Smagghe G, Francis F, de Graaf DC. Study of the Metatranscriptome of
605       Eight Social and Solitary Wild Bee Species Reveals Novel Viruses and Bee Parasites.
606       Front Microbiol. 2018 Feb 14;9:177.

607  38.  Gauthier L, Cornman S, Hartmann U, Cousserans F, Evans J, de Miranda J, et al. The
608       Apis mellifera Filamentous Virus Genome. Viruses. 2015 Jul 9;7(7):3798–815.

609  39.  Kraberger S, Cook CN, Schmidlin K, Fontenele RS, Bautista J, Smith B, et al. Diverse
610       single-stranded DNA viruses associated with honey bees (Apis mellifera). Infect Genet
611       Evol. 2019 Jul;71:179–88.

612  40.  Jacques A, Laurent M, EPILOBEE Consortium, Ribière-Chabert M, Saussac M, Bougeard
613       S, et al. A pan-European epidemiological study reveals honey bee colony survival
614       depends on beekeeper education and disease control. Chaline N, editor. PLOS ONE.
615       2017 Mar 9;12(3):e0172591.

616  41.  Conceição-Neto N, Zeller M, Lefrère H, De Bruyn P, Beller L, Deboutte W, et al. Modular
617       approach to customise sample preparation procedures for viral metagenomics: a
618       reproducible protocol for virome analysis. Sci Rep [Internet]. 2015 Dec [cited 2018 Oct
619       4];5(1). Available from: http://www.nature.com/articles/srep16532

620   42.  Deboutte W, Beller L, Yinda CK, Maes P, de Graaf DC, Matthijnssens J. Honey-bee–
621        associated prokaryotic viral communities reveal wide viral diversity and a profound
622        metabolic coding potential. Proc Natl Acad Sci. 2020 Apr 27;201921859.

623   43.  Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND. Nat
624        Methods. 2015 Jan;12(1):59–60.

625   44.  Ondov BD, Bergman NH, Phillippy AM. Interactive metagenomic visualization in a Web
626        browser. BMC Bioinformatics. 2011 Dec;12(1):385.

627   45.  Suchard MA, Lemey P, Baele G, Ayres DL, Drummond AJ, Rambaut A. Bayesian
628        phylogenetic and phylodynamic data integration using BEAST 1.10. Virus Evol
629        [Internet]. 2018 Jan 1 [cited 2020 Feb 21];4(1). Available from:
630        https://academic.oup.com/ve/article/doi/10.1093/ve/vey016/5035211

631   46.  Souvorov A, Agarwala R, Lipman DJ. SKESA: strategic k-mer extension for scrupulous
632        assemblies. Genome Biol. 2018 Dec;19(1):153.

633   47.  Shi M, Lin X-D, Tian J-H, Chen L-J, Chen X, Li C-X, et al. Redefining the invertebrate RNA
634        virosphere. Nature. 2016 Dec;540(7634):539–43.

635   48.  Shi C, Beller L, Deboutte W, Yinda KC, Delang L, Vega-Rúa A, et al. Stable distinct core
636        eukaryotic viromes in different mosquito species from Guadeloupe, using single
637        mosquito viral metagenomics. Microbiome. 2019 Dec;7(1):121.

638   49.  Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence
639        data. Bioinformatics. 2014 Aug 1;30(15):2114–20.

640   50.  Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes: A
641        New Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing. J
642        Comput Biol. 2012 May;19(5):455–77.

643   51.  Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler
644        transform. Bioinformatics. 2009 Jul 15;25(14):1754–60.

645   52.  Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+:
646        architecture and applications. BMC Bioinformatics. 2009;10(1):421.

647   53.  Hyatt D, Chen G-L, LoCascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal: prokaryotic
648        gene recognition and translation initiation site identification. BMC Bioinformatics. 2010
649        Dec;11(1):119.

650   54.  Hagberg AA, Schult DA, Swart PJ. Exploring Network Structure, Dynamics, and Function
651        using NetworkX. In: Varoquaux G, Vaught T, Millman J, editors. Proceedings of the 7th
652        Python in Science Conference. Pasadena, CA USA; 2008. p. 11–5.

653   55.  Katoh K. MAFFT: a novel method for rapid multiple sequence alignment based on fast
654        Fourier transform. Nucleic Acids Res. 2002 Jul 15;30(14):3059–66.

56. Capella-Gutierrez S, Silla-Martinez JM, Gabaldon T. trimAl: a tool for automated alignment trimming in large-scale phylogenetic analyses. Bioinformatics. 2009 Aug 1;25(15):1972–3.

57. Darriba D, Taboada GL, Doallo R, Posada D. ProtTest 3: fast selection of best-fit models of protein evolution. Bioinformatics. 2011 Apr 15;27(8):1164–5.

58. Yu G, Lam TT-Y, Zhu H, Guan Y. Two Methods for Mapping and Visualizing Associated Data on Phylogeny Using *Ggtree*. Battistuzzi FU, editor. Mol Biol Evol. 2018 Dec 1;35(12):3041–3.

59. Fu L, Niu B, Zhu Z, Wu S, Li W. CD-HIT: accelerated for clustering the next-generation sequencing data. Bioinformatics. 2012 Dec;28(23):3150–2.

60. Peixoto TP. The graph-tool python library [Internet]. Figshare; 2017 [cited 2019 Oct 8]. Available from: https://figshare.com/articles/graph_tool/1164194

61. Shannon P. Cytoscape: A Software Environment for Integrated Models of Biomolecular Interaction Networks. Genome Res. 2003 Nov 1;13(11):2498–504.

62. R Core Team. R: A Language and Environment for Statistical Computing [Internet]. Vienna, Austria: R Foundation for Statistical Computing; 2013. Available from: http://www.R-project.org/

63. Paradis E, Schliep K. ape 5.0: an environment for modern phylogenetics and evolutionary analyses in R. Schwartz R, editor. Bioinformatics. 2019 Feb 1;35(3):526–8.

64. Oksanen J, Blanchet FG, Friendly M, Kindt R, Legendre P, McGlinn D, et al. vegan: Community Ecology Package [Internet]. 2019. Available from: https://CRAN.R-project.org/package=vegan

65. Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, et al. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. Nat Methods. 2020;