

PAPER

# Polymorphism of Genetic Ambigrams

Gytis Dudas,<sup>1,†</sup> Greg Huber,<sup>2,†</sup> Michael Wilkinson<sup>2,3,\*†</sup> and David Yllanes<sup>2,†</sup>

<sup>1</sup>Gothenburg Global Biodiversity Centre, Carl Skottsbergs gata 22B, 413 19, Gothenburg, Sweden, <sup>2</sup>Chan Zuckerberg Biohub, 499 Illinois Street, San Francisco, CA 94158, USA and <sup>3</sup>School of Mathematics and Statistics, The Open University, Walton Hall, Milton Keynes, MK7 6AA, UK

\*Corresponding author. [michael.wilkinson@czbiohub.org](mailto:michael.wilkinson@czbiohub.org) †Authors in alphabetical order

FOR PUBLISHER ONLY Received on Date Month Year; revised on Date Month Year; accepted on Date Month Year

## Abstract

Double synonyms in the genetic code can be used as a tool to test competing hypotheses regarding ambigrammatic narnavirus genomes. Applying the analysis to recent observations of *Culex narnavirus 1* and *Zhejiang mosquito virus 3* ambigrammatic viruses indicates that the open reading frame on the complementary strand of the segment coding for RNA-dependent RNA polymerase does *not* code for a functional protein. *Culex narnavirus 1* has been shown to possess a second segment, also ambigrammatic, termed ‘Robin’. We find a comparable segment for *Zhejiang mosquito virus 3*, a moderately diverged relative of *Culex narnavirus 1*. Our analysis of Robin polymorphisms suggests that its reverse open reading frame also does not code for a protein. We make a hypothesis about its role.

## 1 Introduction

2 Of all the various types of viruses catalogued, narnaviruses  
3 rank among the simplest and most surprising (Cobián Güemes  
4 et al., 2016). Narnaviruses (a contraction of ‘naked RNA virus’)  
5 are examples of a minimal blueprint for a virus: no capsid,  
6 no envelope, no apparent assembly of any kind. The known  
7 narnavirus blueprint appeared for all intents and purposes to  
8 be a single gene, that which codes for an RNA-dependent RNA  
9 polymerase, abbreviated as RdRp, (Hillman and Cai, 2013).  
10 However, some narnaviruses have been found to have a genome  
11 with an open reading frame (i.e., a reading frame without stop  
12 codons) on the strand complementary to that coding for the  
13 RdRp gene, calling into question the general hypothesis of a  
14 one-gene blueprint (DeRisi et al., 2019; Dinan et al., 2020;  
15 Cepelewicz, 2020). This reverse open reading frame (rORF)  
16 has codon boundaries aligned with the forward reading frame.  
17 Because the genome can be translated in either direction,  
18 we say that these narnaviruses are *ambigrammatic*. The  
19 significance of an ambigrammatic genome is an open problem.  
20 In this paper we discuss how polymorphisms of sampled  
21 sequences can distinguish between competing hypotheses on  
22 the function and nature of ambigrammatic viral genomes. Our  
23 methods are applied to known ambigrammatic narnavirus genes  
24 and to the newly discovered ambigrammatic second segment of  
25 some narnaviruses, termed *Robin* (Batson et al., 2020).

Our discussion is based upon two rules about the genetic  
code and its relation to ambigrammatic sequences. Both of  
these *ambigram rules* are concerned with the availability of  
synonyms within the genetic code, which allow coding of the  
same amino acid with a different codon. The first rule states  
that for any sequence of amino acids coded by the forward  
strand, it is possible to use individual synonymous substitutions  
to remove all stop codons on the complementary strand (this  
result was discussed already in DeRisi et al., 2019). The second  
ambigram rule, described below, states that the genetic code  
contains double synonyms that allow polymorphisms, accessible  
by single-base mutations, even when the amino acids coded by  
both the forward and the complementary strands are fixed.

The first of these rules addresses the ‘how’ of ambigrammatic  
genomes, by showing that stop codons on the complementary  
strand can be removed by single-point mutations, without  
altering the protein (in narnaviruses, the RdRp) coded in the  
forward direction. Here we argue that the second rule can help  
to resolve the ‘why’ of ambigrammatic genomes: the origin of  
ambigrammaticity itself. There are two distinct reasons why  
there might be an evolutionary advantage for a virus to evolve  
an ambigrammatic sequence. The first possibility is that the  
complementary strand might code for a functionally significant  
protein, for example, one that might interfere with host defence  
mechanisms. The second possibility is that the lack of stop  
codons on the complementary strand is significant, even if the

52 amino acid sequence that is coded is irrelevant. In particular,  
 53 the lack of stop codons may promote the association between  
 54 ribosomes and the complementary strand viral RNA (produced  
 55 as part of its replication cycle). It is possible that a ‘polysome’  
 56 formed by a covering of ribosomes helps to shield the virus from  
 57 degradation or from detection by cellular defence mechanisms  
 58 (Cepelewicz, 2020; Retallack et al., 2020; Wilkinson et al.,  
 59 2021). The second ambigram rule combined with data on the  
 60 polymorphism of the virus genome can help distinguish whether  
 61 the complementary strand codes for a functional protein. We  
 62 shall argue that in the case of *Culex narnavirus 1* and *Zhejiang*  
 63 *mosquito virus 3*, the evidence is in favour of this second  
 64 hypothesis, namely that the open reading frame (ORF) on the  
 65 complementary strand does not code for a functional protein.

66 After describing the genetic ambigram rules, we discuss how  
 67 the existence of double synonyms can be used to assess whether  
 68 the open reading frame on the complementary chain codes for  
 69 functional protein. It is well known that, because RdRp is a  
 70 highly-conserved gene, non-synonymous mutations are likely to  
 71 be detrimental, so that most of the observed diversity consists  
 72 of synonymous changes. Some of these synonymous mutations  
 73 have the potential to be synonymous in the complementary  
 74 strand. If the complementary strand also codes for a functional  
 75 protein, we expect that doubly synonymous mutations will  
 76 be favoured. In fact, there would be mutational ‘hotspots’  
 77 corresponding to the potential doubly-synonymous loci. We  
 78 introduce two tests for whether the complementary strand is  
 79 coding, based respectively on looking for mutational ‘hotspots’,  
 80 and upon the mutational frequencies at loci which have double  
 81 synonyms. We used these tests to analyse sequences for two  
 82 different ambigrammatic narnaviruses: 46 RdRp segments of  
 83 *Culex narnavirus 1* and 12 RdRp segments of *Zhejiang*  
 84 *mosquito virus 3*, abbreviated to CNV and ZMV respectively.  
 85 We find that neither of our tests supports the hypothesis that  
 86 the translated sequence of the complementary strand of RdRp  
 87 is under purifying selection. We also applied these tests to the  
 88 second segment, termed *Robin*, which is found to be closely  
 89 associated with this ambigrammatic narnavirus infection in  
 90 mosquitos (Batson et al., 2020; Retallack et al., 2020). We also  
 91 found that the complementary open reading frame of Robin  
 92 does not appear to be under purifying selection. The discovery  
 93 of Robin suggested that ambigrammatic companions may exist  
 94 for other ambigrammatic viruses. Accordingly, we searched the  
 95 assembled contigs of studies reporting the detection of ZMV,  
 96 the only other ambigrammatic narnavirus observed multiple  
 97 times in numerous locations, and discovered an ambigrammatic  
 98 segment with similar properties to CNV Robin. Thus we  
 99 consider four viral segments, denoted CNV-RdRp, CNV-Robin,  
 100 ZMV-RdRp, ZMV-Robin. We shall report evidence that Robin  
 101 does code for a protein in its forward direction, but that its  
 102 complementary strand is non-coding. We find evidence that  
 103 Robin segments are under detectable purifying selection. Figure  
 104 1 illustrates the phylogenetic relationship of CNV and ZMV,  
 105 and ORF-wide dN/dS values of all their segments and coding  
 106 directions (discussed in detail below).

107 Some careful consideration is required to reconcile our  
 108 observations with results recently reported in Retallack et al.  
 109 (2020), where it was shown that introducing mutations which  
 110 are non-synonymous on the reverse open reading frame of  
 111 *Culex narnavirus 1* can reduce the fitness of this virus. In  
 112 the concluding section, we consider the interpretation of these  
 113 observations, and discuss whether there may be implications  
 114 for other viral families.

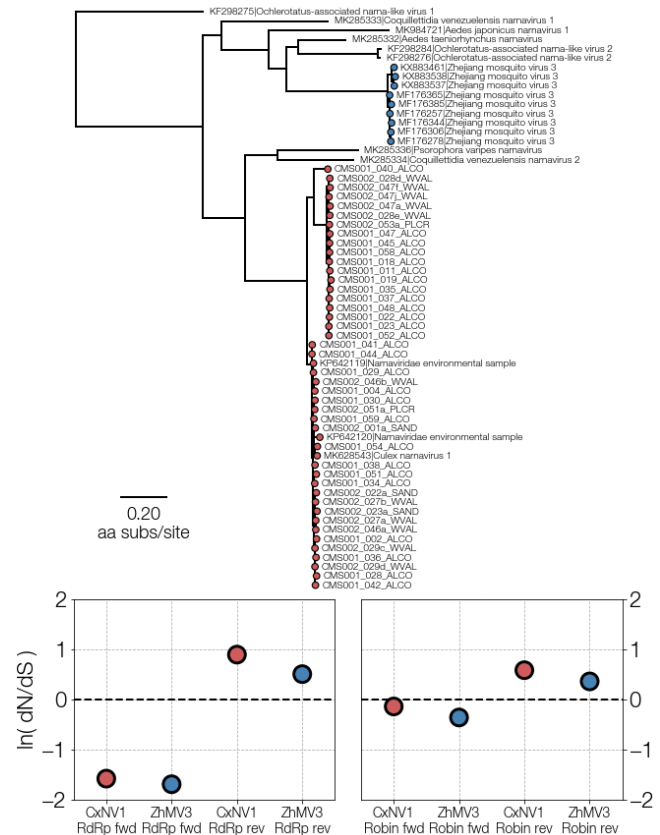


Fig. 1. a A maximum-likelihood tree illustrating the relationship between CNV (*Culex narnavirus 1*) (red) and ZMV (*Zhejiang mosquito virus 3*) (blue). b ORF-wide dN/dS values for forward and reverse directions of RdRp and Robin segments for both viruses.

115 There are many examples of overlapping viral genes with  
 116 staggered reading frames: this was first clearly described in  
 117 Barrell et al. (1976), and has been reviewed in Chirico et al.  
 118 (2010). Recent work by Nelson, Ardern and Wei (Nelson et al.,  
 119 2020) discusses how these can be identified. Our investigations  
 120 indicate that the ambigrammatic ORFs discussed in this work  
 121 are a different phenomenon, because they are non-coding.  
 122 Our approach to analysing the ambigrammatic sequences is  
 123 quite distinct from the rather complex machinery proposed in  
 124 Nelson et al. (2020), because it emphasises the role of double  
 125 synonyms as an unambiguous discriminant of the role of the  
 126 ambigrammatic sequences.

## Ambigram rules and their significance

We start by describing the two genetic ambigram rules.

**Rule 1** All complementary-strand stops are removable

130 Consider the reading frame on the complementary strand  
 131 that has its codons aligned with those on the forward  
 132 strand. Every codon on the forward strand corresponds to a  
 133 complementary-strand codon read in the reverse direction. The  
 134 rule states that any stop codon on the complementary strand  
 135 can be removed by a single-point mutation which leaves the  
 136 amino acid specified by the forward-read codon unchanged.

137 This result is demonstrated by the following argument, as  
 138 discussed in DeRisi et al. (2019). Reversing the read direction  
 139 and taking the pairing complement, the stop codons UAA,  
 140 UAG, UGA in the standard genetic code become, respectively,  
 141 UUA, CUA, UCA, for which the amino acids are Leu, Leu,  
 142 Ser. It is only instances of leucine and serine in the forward  
 143 sequence that can result in stop codons in the reverse read.  
 144 The synonyms of Leu are CUN, UUA, UUG (where N means  
 145 any base). The synonyms of Ser are UCN, AGU, AGC. The  
 146 undesirable Leu codon UUA can be transformed to UUG by  
 147 a single substitution. Similarly, the Leu codon CUA can be  
 148 transformed to CUU, CUG or CUC by single substitutions. And  
 149 the Ser codon UCA is transformed to UCU, UCG or UCC by  
 150 single substitutions. We conclude that every stop codon on the  
 151 reverse reading frame can be removed by a synonymous, single  
 152 site nucleotide mutation.

153 Furthermore, it is found that complementary-strand stops  
 154 cannot always be removed by synonymous substitutions in the  
 155 other two read frames for the complementary strand (each  
 156 case requires a separate and somewhat involved argument, also  
 157 given in DeRisi et al., 2019). As a consequence of these two  
 158 arguments, we need discuss only the complementary read frame  
 159 with aligned codons.

160 **Rule 2** *There exist double synonyms*

161 Most synonymous mutations of the forward strand produce  
 162 a non-synonymous change in the complementary strand, but  
 163 the genetic code does include a number of double synonyms,  
 164 where the reverse complement of a synonymous mutation is  
 165 also a synonym. For example codon AGG (Arg) can become  
 166 CGG (Arg) via a synonymous mutation, while the reverse  
 167 complement of AGG, which is CCU (Pro) transforms to CCG  
 168 (Pro) under the same mutation.

169 The full set of double synonyms in the standard genetic code  
 170 are as follows:

- 171 • Two of the six synonyms of Ser are double synonyms, with  
 172 reverse complements coding Arg. Conversely, two of the  
 173 six synonyms of Arg are double synonyms, with reverse  
 174 complement coding Ser.
- 175 • Two more of the six synonyms of Arg are double synonyms,  
 176 with reverse complement Pro. Conversely, two of the four  
 177 synonyms of Pro are double synonyms coding for Arg.
- 178 • Two of the six synonyms of Leu are double synonyms, with  
 179 reverse complement Gln. Conversely, both synonyms of Gln  
 180 are double synonyms, with reverse complement coding Leu.

181 Table 1 lists the sets of single and double synonyms for those  
 182 amino acids that can have double synonyms. (We exclude the  
 183 two synonyms of Ser and the one synonym of Leu for which  
 184 the reverse complement is Stop, because these do not occur in  
 185 ambigrammatic genes.)

186 **Implications**

187 Our first rule shows that an ambigrammatic version of any  
 188 gene can evolve, without making any changes to the amino acid  
 189 sequence. This establishes how ambigrammatic sequences can  
 190 arise, but it does not illuminate why they are favoured.

191 Combined with observed polymorphisms of narnaviruses,  
 192 the second ambigram rule can give an indication of the utility  
 193 of ambigrammatic sequences. In studies on the (usual) non-  
 194 ambigrammatic genomes, the ratio of synonymous to non-  
 195 synonymous mutations is used as an indicator of whether

**Table 1.** For each amino acid (AA) that can have double-synonym mutations, we list all of the possible codons which do not code for Stop on the complementary strand, indicating their reverse complement (Comp. AA). The codons that have a double synonym are marked with an asterisk. For each of these codons, we list the number of mutations which are synonymous, and the number of double synonym mutations. In each case the numbers of single (double) mutations are written  $S^{(n)} + S^{(v)}$  ( $D^{(n)} + D^{(v)}$ ), where the superscript n denotes transitions, and superscript v transversions. Also, double synonyms are counted in the list of single synonyms.

AA	Codon	$S^{(n)} + S^{(v)}$	$D^{(n)} + D^{(v)}$	Comp. AA
Leu	UUG*	1 + 0	1 + 0	Gln
	CUU	1 + 1	0 + 0	Lys
	CUC	1 + 1	0 + 0	Glu
	CUG*	1 + 2	1 + 0	Gln
Pro	CCU*	1 + 2	0 + 1	Arg
	CCC	1 + 2	0 + 0	Gly
	CCA	1 + 2	0 + 0	Trp
	CCG*	1 + 2	0 + 1	Arg
Gln	CAA*	1 + 0	1 + 0	Leu
	CAG*	1 + 0	1 + 0	Leu
Arg	CGU	1 + 2	0 + 0	Thr
	CGC	1 + 2	0 + 0	Ala
	CGA*	1 + 3	0 + 1	Ser
	CGG*	1 + 3	0 + 1	Pro
	AGA*	1 + 1	0 + 1	Ser
	AGG*	1 + 1	0 + 1	Pro
Ser	UCU*	1 + 1	0 + 1	Arg
	UCC	1 + 1	0 + 0	Gly
	UCG*	1 + 2	0 + 1	Arg
	AGU	1 + 0	0 + 0	Thr
	AGC	1 + 0	0 + 0	Ala

196 the nucleotide sequence codes for a protein: non-synonymous  
 197 mutations are likely to be deleterious if the sequence codes  
 198 for a functional protein. We shall adapt this approach to our  
 199 study of ambigrammatic narnavirus genes. We assume that the  
 200 forward direction is a coding sequence (usually for RdRp), and  
 201 confine attention to those mutations which are synonymous in  
 202 the forward direction. If the complementary strand codes for  
 203 a functional protein, most of these synonymous mutations will  
 204 inevitably result in changes of the complementary amino acid  
 205 sequence. However, at many loci the evolutionarily favoured  
 206 amino acid will be one that allows double synonyms. In these  
 207 cases, there can be non-deleterious mutations between a pair  
 208 of codons that preserve the amino acid sequence of both the  
 209 forward and the complementary strands.

210 If the complementary strand codes for a functional protein,  
 211 we expect studies of the polymorphism of the gene would show  
 212 that these double-synonym loci will be mutational ‘hotspots’,  
 213 where mutations occur more frequently. In addition, the double-  
 214 synonym pairs would be represented far more frequently than  
 215 other mutations at these loci. These observations lead to two  
 216 distinct tests for whether there is evolutionary pressure on the  
 217 translated sequence of the complementary strand.

## 218 Ambigrammatic narnavirus genes

219 We analysed data from samples of two ambigrammatic  
220 narnaviruses, *Culex narnavirus 1* (CNV, with 46 genomes)  
221 and *Zhejiang mosquito virus 3* (ZMV, with 10 genomes).  
222 Both narnaviruses have an ambigrammatic RdRp coding gene,  
223 denoted CNV-RdRp and ZMV-RdRp respectively. The reverse  
224 open reading frame has its codons aligned with the forward  
225 frame. In both forward and reverse reading frames any stop  
226 codons are close to the 3' end of the respective frame.

227 The ambigrammatic feature is certainly a puzzle. There  
228 appear to be two classes of plausible explanations:

- 229 1. **The reverse open reading frame codes a protein.**  
230 This is logically possible, but if the RdRp gene is strongly  
231 conserved, there is very little flexibility in the rORF.  
232 However, in the absence of any additional evidence it  
233 is the explanation which requires the fewest additional  
234 hypotheses.
- 235 2. **The reverse open reading frame facilitates association**  
236 **of ribosomes with RNA.** This could conceivably convey  
237 advantages by providing a mechanism to protect viral  
238 RNA from degradation, but without further evidence this  
239 requires additional hypotheses.

240 Recently, additional evidence has emerged which may  
241 provide support for the second of these explanations.  
242 Specifically, the CNV infection has recently been shown  
243 to be associated with another ambigrammatic viral RNA  
244 segment, termed *Robin* (Batson et al., 2020; Retallack et al.,  
245 2020). It was reported that this segment, CNV-Robin, is  
246 ambigrammatic, with forward and reverse codons aligned, over  
247 very nearly the entire length (about 850 nt), where direction  
248 designation is determined by which amino acid sequence  
249 appears more conserved. Again, any stop codons occur close  
250 to the 3' end. Neither forward nor reverse directions of Robin  
251 are homologous with known sequences.

252 Because ambigrammatic genes are rare, finding two of them  
253 in the same system is a strong indication that their occurrence  
254 has a common explanation. This observation makes it appear  
255 unlikely that the reverse open reading frame is a device to 'pack  
256 in' an additional protein coding gene, and more likely that the  
257 ambigrammatic feature is associated with allowing ribosomes  
258 to associate with both strands of the viral RNA.

259 This reasoning suggests that the Robin gene may play a  
260 role in selecting for the ambigrammatic property (for example,  
261 it may facilitate protection by ribosomes of the viral RNA). If  
262 this surmise is correct, we should expect to see a version of the  
263 Robin gene associated with other ambigrammatic narnaviruses.  
264 It is possible that this might be detected by a search of archived  
265 sequence data. Only *Zhejiang mosquito virus 3* appeared to  
266 be observed multiple times to make detection of an additional  
267 Robin segment possible, so we concentrated on that system.

268 We were able to find evidence of an ambigrammatic RNA, of  
269 length approximately 900 nt, that co-occurs with ZMV RdRp  
270 segment across multiple samples recovered by at least two  
271 studies that, like CNV Robin, bears no recognisable homology  
272 to publicly available sequences or CNV Robin itself. Given the  
273 conjunction of these unusual features we strongly believe this  
274 ambigrammatic RNA to be the equivalent of a Robin segment  
275 in ZMV.

## 276 Methods

277 Tests for whether the complementary strand is coding

278 We have argued that doubly-synonymous mutations will give  
279 a signature of the reverse strand coding for a functional  
280 protein. If the reverse-direction code is functional, then the  
281 only assuredly non-deleterious mutations would be the double-  
282 synonym ones, where one codon is transformed by a single-  
283 nucleotide substitution to another codon which preserves the  
284 amino acid coded in both the forward and the reverse directions.

285 Assume that we have  $M$  sequences of an ambigrammatic  
286 gene, fully sequenced and maximally aligned with each other,  
287 and that one strand, referred to as the 'forward' strand, codes  
288 for a functional protein. We identify a 'consensus' codon at each  
289 of the  $N$  loci, and then enumerate the set of variant codons at  
290 each amino acid locus. If the consensus codon at a locus is  
291 one of the twelve double-synonym codons listed in table 1, we  
292 term this a *doubly-synonymous locus*. The number of doubly-  
293 synonymous loci is  $N_{ds}$ .

294 There are two different approaches to testing whether double  
295 synonyms indicate that the complementary strand is coding:

296 Look for the existence of mutational 'hotspots'

297 We can look for evidence that the doubly-synonymous loci  
298 experience more substitutions than other loci.

299 For each codon locus  $k$ , we can determine the number of  
300 elements of the variant set,  $n(k)$ , and also the fraction of codons  
301  $f(k)$  which differ from the consensus codon. We then determine  
302 the averages of these quantities,  $\langle n(k) \rangle$  and  $\langle f(k) \rangle$ , for the  
303 doubly-synonymous loci and for the other loci. If the ratios

$$R_n = \frac{\langle n(k) \rangle |_{\text{double syn. loci}}}{\langle n(k) \rangle |_{\text{other loci}}}, \quad R_f = \frac{\langle f(k) \rangle |_{\text{double syn. loci}}}{\langle f(k) \rangle |_{\text{other loci}}} \quad (1)$$

304 are large, this is evidence that the complementary strand is  
305 coding.

306 The null hypothesis, indicating that the reverse open  
307 reading frame is non-coding, is that the ratios  $R_n$  and  $R_f$  are  
308 sufficiently close to unity that the difference may be explained  
309 by statistical fluctuations.  
310

311 Mutation frequencies test

312 We can also look at codon frequencies for different mutations at  
313 doubly-synonymous loci. If the complementary strand is coding,  
314 we expect to find that the frequency of mutations observed  
315 at doubly-synonymous loci will heavily favour double-synonym  
316 codons over single synonyms. We consider the subset of double-  
317 synonym loci where mutations are observed (that is, where  
318  $n(k) > 1$ ). For each of these  $N_a$  *variable doubly-synonymous*  
319 *loci*, we can determine two numbers:  $n_s(k)$  is the numbers of  
320 singly-synonymous variants at locus  $k$ , and  $n_d(k)$  is the number  
321 of these variants which are also doubly-synonymous. (Clearly  
322  $n(k) \geq n_s(k) \geq n_d(k)$ ). If  $n_d(k) = n_s(k)$ , that means that  
323 the mutations preserve the complementary-strand amino acid,  
324 which is an indication that the reverse strand is coding. If  $\{k^*\}$   
325 is the set of variable doubly-synonymous loci, we then calculate

$$N_s = \sum_{k \in \{k^*\}} n_s(k), \quad N_d = \sum_{k \in \{k^*\}} n_d(k). \quad (2)$$

327 If the complementary strand is coding, we expect

$$R \equiv \frac{N_s}{N_d} \quad (3)$$

329 to be close to unity.

330 However, there will also be beneficial or neutral mutations  
 331 which do change the amino acids, so that not all mutations  
 332 will be between sets of doubly-synonymous codons. We need  
 333 to be able to quantify the extent to which finding other than  
 334 double-synonym mutations is an indication that the reverse  
 335 strand is non-coding. We must do this by comparison with a  
 336 null hypothesis, in which the reverse strand is non-coding.

337 Null hypothesis for mutation frequencies

338 Let  $R_0$  be the value of the ratio  $R$  that is derived from  
 339 this null hypothesis that the complementary strand is non-  
 340 coding. In order to compute the expected  $N_s/N_d$  ratio,  $R_0$ ,  
 341 we adopt the following approach. We assume that the  $M$   
 342 sequences are sufficiently similar that only a small fraction of  
 343 loci have undergone mutations. We adopt the Kimura model  
 344 (Kimura, 1980), which assumes that the mutation rate  $r_n$  for  
 345 transitions ( $A \leftrightarrow G$  or  $C \leftrightarrow U$ ) is different from the rate  
 346  $r_v$  for transversions (other single-nucleotide mutations), and  
 347 negligible for other types of mutation. The ratio of these rates  
 348 is

$$349 \quad \alpha \equiv \frac{r_n}{r_v}. \quad (4)$$

350 If the numbers of single (double) synonyms of the consensus  
 351 nucleotide at locus  $k$  leading to transitions or transversions are  
 352 respectively  $S_k^{(n)}$  and  $S_k^{(v)}$  ( $D_k^{(n)}$ ,  $D_k^{(v)}$ ), then we estimate

$$353 \quad R_0 = \frac{\sum_{k \in \{k^*\}} \alpha S_k^{(n)} + S_k^{(v)}}{\sum_{k \in \{k^*\}} \alpha D_k^{(n)} + D_k^{(v)}} \quad (5)$$

354 The numbers  $S_k^{(n)}$ ,  $S_k^{(v)}$ ,  $D_k^{(n)}$ ,  $D_k^{(v)}$  are given in table 1 for all  
 355 of the double-synonym codons.

356 Finding the Robin segment of *Zhejiang mosquito virus 3*

357 We looked through assembled contig datasets from two  
 358 metagenomic mosquito studies (three from China and six from  
 359 Australia) (Shi et al., 2016, 2017), kindly provided to us by  
 360 Mang Shi and Edward C Holmes. We clustered contigs from  
 361 the nine datasets by similarity using CD-HIT (Fu et al., 2012)  
 362 with a threshold of 90% and looked for clusters that contained  
 363 contigs from at least 6 samples, that did not have standard  
 364 deviation in contig length greater than 1200, and had fewer  
 365 than 200 contigs. Of the hundreds of clusters filtered this way  
 366 only a handful also possessed sequences ambigrammatic across  
 367 at least 90% of their length and only two clusters were mostly  
 368 comprised of ambigrammatic sequences, while the rest were  
 369 clearly recognisable as mosquito contigs. Of the two clusters  
 370 one was identifiable as the RdRp of *Zhejiang mosquito virus*  
 371 *3*, while we presume the other to be an unrecognisably distant  
 372 orthologue of *Culex narnavirus 1* Robin, on account of its  
 373 co-occurrence with ZMV RdRp, ambigrammaticity, and length.

## 374 Results

375 Next we report the results of our studies of polymorphism of the  
 376 four ambigrammatic narnavirus genes. We discuss what can be  
 377 learned from applying standard techniques, before discussing  
 378 the results of our tests for whether the reverse open reading  
 379 frame codes for a protein.

380 Forward reading frame

381 Each sequence was trimmed to a length of  $3N$  nucleotides. We  
 382 identified a consensus nucleotide at each locus, and determined

the set of variant nucleotides at each locus. We determined the  
 total number of transition and transversion mutations which are  
 observed,  $N_n$  and  $N_v$  respectively. We also determined the total  
 number of mutations at each position in the codon,  $(n_1, n_2, n_3)$ .  
 We estimated the average number of variable sites  $r$  as the total  
 number of nucleotide variants, divided by the product of the  
 number of sequences and alignment length. We also estimated  
 the ratio  $\alpha$  of the rate of selected transition mutations to the  
 rate of transversions:

$$392 \quad r \equiv \frac{n_1 + n_2 + n_3}{3NM}, \quad \alpha \equiv \frac{r_n}{r_v} = \frac{2N_n}{N_v} \quad (6)$$

(recall that there are twice as many transversions as  
 transitions). We also determined a ‘normalised’ triplet of  
 variable sites for each position within the codon:  $(z_1 : z_2 : z_3) = 3(n_1 : n_2 : n_3)/(n_1 + n_2 + n_3)$ . Our results on the  
 nucleotide-level investigation of polymorphism are summarised  
 in table 2.

We then assigned a consensus codon at each codon locus,  
 selecting the frame by the criterion of minimising the number  
 of stop codons. For each of the  $N$  codons, we determined  
 the variant set of codons which were observed in each of  
 the  $M$  sequences. The total number of synonymous and non-  
 synonymous single-nucleotide changes in the variant sets was  
 $N_{sy}$  and  $N_{ns}$  respectively. The total number of mutations  
 encountered in the variant sets where two or three nucleotides  
 were changed was  $N_{mult}$ . For each codon there are numbers  
 of possible non-synonymous mutations which are transitions  
 and transversions,  $n_k^{(n)}$  and  $n_k^{(v)}$ , and numbers of synonymous  
 mutations which are transitions and transversions,  $s_k^{(n)}$  and  $s_k^{(v)}$   
 (with  $s_k^{(n)} + n_k^{(n)} + s_k^{(v)} + n_k^{(v)} = 9$ ). Under the null hypothesis  
 that the sequence is non-coding, the expected value of the ratio

$$399 \quad R = \frac{N_{ns}}{N_{sy}} \quad (7)$$

is

$$403 \quad R_{exp} = \frac{\sum_{k=1}^N \alpha n_k^{(n)} + n_k^{(v)}}{\sum_{k=1}^N \alpha s_k^{(n)} + s_k^{(v)}}. \quad (8)$$

We also determined the fraction of codons where multi-  
 nucleotide mutations are observed,  $f_{mult} = N_{mult}/N$ . We  
 present our results for the codon-level mutations in table  
 3, which includes information for both the forward and the  
 complementary read directions (with codon boundaries aligned  
 for the complementary direction).

The alignments are *ambigrammatic*, in the sense that there  
 are no stop codons in the interior of the sequence. None of  
 the individual sequences had stop codons in the body of the  
 sequence in either direction.

We also computed ORF-wide  $dN/dS$  values (plotted in  
 figure 1(b)), by assuming that every mutation in the alignment  
 has occurred only once to be conservative. This was motivated  
 by the presence of pairs of sites with four haplotypes between  
 them (4G sites), an indication that recombination may be  
 a potential issue with narnavirus sequences. Normalising  
 the number of observed non-synonymous and synonymous  
 mutations was done by assuming a transition/transversion  
 ratio of 2, consistent with equation (6). These values  $dN/dS$   
 values are slightly different from the  $R/R_{exp}$  ratios in table  
 3 because the latter excludes mutations where more than one  
 base differs from the consensus codon. In all but one of the cases  
 $dN/dS$  is higher than  $R/R_{exp}$ , because the multiple nucleotide  
 mutations which are included in  $dN/dS$  are predominantly  
 non-synonymous.

441 Based upon these tables, we can make the following  
442 observations and deductions:

- 443 1. **Diversity.** We observe that both RdRp and Robin  
444 segments are comparable in their diversity, for both  
445 CNV and ZMV. As expected, RdRp sequences are highly  
446 conserved at the amino acid level. Robin, on the other  
447 hand, appears far more relaxed at the amino acid level and,  
448 consistent with this, diverged beyond recognition between  
449 CNV and ZMV.
- 450 2. **Relative mutation rate by codon position.** For RdRp  
451 sequences, more mutations are observed at the third  
452 nucleotide in each codon, as expected for a sequence  
453 that preserves the amino acid sequence (because most  
454 synonymous mutations involve the third nucleotide of a  
455 codon). In the case of Robin sequences, the frequencies  
456 of mutation are much closer to being equal, to the extent  
457 that for CNV-Robin the null hypothesis that the rates are  
458 equal is not definitively rejected. However, mutations at  
459 different codon sites are sufficiently weighted towards the  
460 third position that we shall assume that Robin does code  
461 for a functional protein.  
462 While the values of  $(z_1 : z_2 : z_3)$  are very different for  
463 RdRp and Robin, their values are comparable for CNV and  
464 ZMV, which is an indication that the selective pressures on  
465 both viruses are the same.
- 466 3. **Rate of multiple-nucleotide mutations.** The fraction of  
467 multiple-nucleotide mutations is higher for Robin sequences  
468 than it is for RdRp sequences. This may be an indication  
469 that the Robin sequence is under strong selective pressure,  
470 because some amino acid substitutions can only be achieved  
471 through multiple nucleotide mutations.
- 472 4. **Transition to transversion ratio.** Three of the values  
473 of  $\alpha$  were similar to each other, while the value for ZMV-  
474 RdRp was higher than the others. Because transitions occur  
475 at a higher intrinsic rate, a lower value of  $\alpha$  indicates  
476 that observed mutations are biased in favour of the rarer  
477 transversions, which is an indication of unusual selective  
478 pressures. The fact that the values of  $\alpha$  for the Robin  
479 segments are comparable to, or lower than, the values for  
480 RdRp are a further indication that Robin is under similar  
481 selective pressure too.
- 482 5. **Ratio of non-synonyms to synonyms.** For the RdRp  
483 segments the values of  $R = N_{ns}/N_s$  are much smaller  
484 than the values  $R_0$  predicted (equation (8)) by the null  
485 hypothesis that mutations are random. This indicates that  
486 the selective pressure on RdRp acts to preserve the amino  
487 acid sequence. For Robin segments, the values of  $R$  are  
488 much larger, but still smaller than the prediction from  
489 the null hypothesis. This indicates that while points 1-4  
490 above indicate that Robin is under some selective pressure,  
491 the amino acid sequence is not strongly conserved. This is  
492 consistent with the hypothesis that the selection acting on  
493 Robin is relaxed.

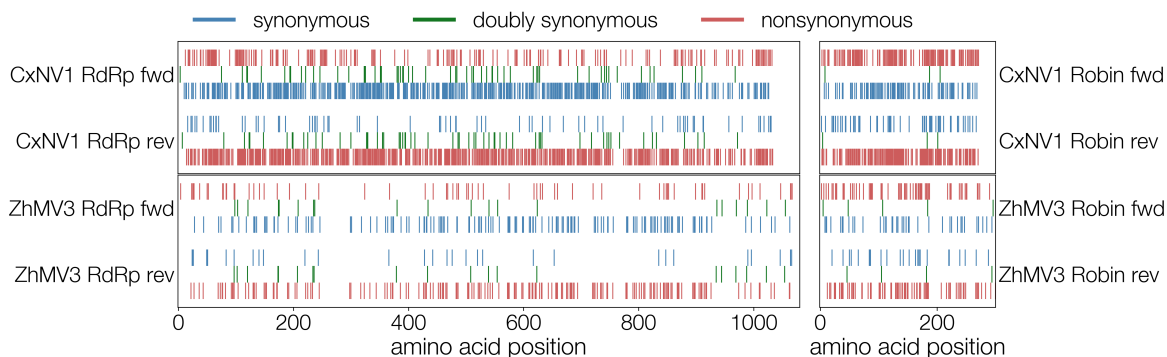
494 Figure 2 illustrates the distribution of mutations across the  
495 forward and reverse reading frames of all four ORFs for both  
496 CNV and ZMV. As expected, there is evidence that some  
497 regions accumulate mutations more readily than others. The  
498 pattern is consistent with what would be expected from the  
499 statistical reductions in the tables.

Complementary reading frame

500 We determined the set of  $N_{ds}$  doubly-synonymous codons in  
501 the consensus sequence, and the subset of  $N_a$  of these which  
502 have variant codons.  
503

- 504 1. **Mutational hotspots test.** We applied the mutational  
505 hotspots test to all four sequences, as described by  
506 equations (1) above. The results (tables 4) show no evidence  
507 that the doubly synonymous sites are undergoing more  
508 frequent mutations, or that their mutations are more widely  
509 spread across the dataset.
- 510 2. **Mutation rate test.** We examined the number of  
511 mutations in the set of  $N_a$  doubly-synonymous sites  
512 which were variable. We found (table 5) that many  
513 more of the observed mutations at these sites are only  
514 singly synonymous, when a doubly-synonymous mutation is  
515 possible, which is further evidence that the complementary  
516 strand is non-coding. The numbers of doubly-synonymous  
517 mutations were quite low, and so it was not possible to  
518 make a reliable comparison of the ratio  $N_s/N_d$  with the  
519 null hypothesis.
- 520 3. **Ratio of non-synonyms to synonyms.**

521 The ratios of non-synonymous to synonymous mutations,  
522 presented in table 3 and figure 1(b), were lower than the  
523 null hypothesis for the forward direction. This is readily  
524 explained as an indication that the forward ORF codes for  
525 a functional protein. However the  $N_{ns}/N_s$  ratios for the  
526 reverse direction were all higher than the null hypothesis.  
527 This observation is explained, qualitatively, as follows. If  
528 the forward direction strictly conserves the amino acid  
529 sequence, then all of the mutations which are synonymous  
530 on the reverse strand are doubly-synonymous. Because  
531 only 12 of the 64 codons allow for doubly-synonymous  
532 mutations, the  $N_{ns}/N_s$  ratio would be very high for the  
533 complementary strand if the forward sequence were to be  
534 exactly conserved. We computed this ratio, and found  
535 11.2 for CNV-RdRp, and similar values for the other  
536 sequences. This theoretical ratio is considerably higher than  
537 the measured value of 4.97, because the forward sequence  
538 is not exactly conserved. For Robin segments, the value of  
539  $R$  for the reverse ORF is only slightly higher than the null  
540 hypothesis, because the amino acid sequence is only weakly  
541 conserved.



**Fig. 2.** Distribution of synonymous (blue), non-synonymous (red) substitutions, and doubly synonymous sites (green) in CNV (upper plots) and ZMV (lower plots) RdRp (left) and Robin (right) segments in both directions (forward towards top, reverse towards bottom). Translated reverse ORFs are shown backwards (segment coordinate space). Double synonyms don't overlap perfectly because forward and reverse ORFs differ in length and begin and end at different positions along the segment.

**Table 2.** Nucleotide-level statistics of mutations. The consensus sequence has  $N$  codons. Among the mutations observed in  $M$  polymorphs, there are  $N_n$  transitions,  $N_v$  transversions, with overall rate  $r$  and transition/transversion rate ratio  $\alpha$ . The numbers total mutations at each base position is  $(n_1 : n_2 : n_3)$ , and normalising these to ratios via equation (6) yields  $(z_1 : z_2 : z_3)$ .

Strand	$N$	$M$	$N_n$	$N_v$	$r$	$\alpha$	$(n_1, n_2, n_3)$	$(z_1 : z_2 : z_3)$
CNV-RdRp	1033	46	606	362	0.0068	3.35	(181, 140, 645)	(0.56 : 0.44 : 2.00)
ZMV-RdRp	1075	12	210	39	0.0064	10.80	(47, 29, 173)	(0.57 : 0.35 : 2.08)
CNV-Robin	272	46	213	146	0.0096	2.92	(107, 100, 152)	(0.89 : 0.84 : 1.27)
ZMV-Robin	304	10	84	48	0.0145	3.50	(35, 31, 66)	(0.80 : 0.70 : 1.50)

## 542 Discussion

543 We have argued that doubly synonymous codons provide a key  
 544 to understanding whether ambigrammatic viral RNA segments  
 545 code for two functional proteins. If there were two coding  
 546 genes, doubly synonymous mutations would be mutational  
 547 hotspots, because they are unambiguously non-deleterious. We  
 548 applied our analysis to recent observations of polymorphisms  
 549 in two ambigrammatic narnaviruses: *Culex narnavirus 1* and  
 550 *Zhejiang mosquito virus 3*. There was no evidence that  
 551 doubly synonymous sites are mutational hotspots, or that  
 552 there is a prevalence of mutations to other doubly-synonymous  
 553 codons at these sites. Other, circumstantial, evidence favours  
 554 the interpretation that the complementary strand is non-  
 555 coding. Ambigrammatic sequences have been observed in other  
 556 narnaviruses, but they are undoubtedly a rare phenomenon.  
 557 If the rORF (reverse open reading frame) of both RdRp and  
 558 Robin segments had evolved to code for a functional protein,

559 each RNA segment would code for two genes. Given that  
 560 ambigrammatic sequences are rare (DeRisi et al., 2019), finding  
 561 a system where two had evolved independently would be highly  
 562 improbable. Moreover, because the ambigrams are full length,  
 563 each of the ambigrammatically coded sequences would code for  
 564 two genes which have the same length as each other.

565 An observation of the simultaneous detection of two or more  
 566 ambigrammatic genes would strongly favour models where there  
 567 is an advantage in evolving an ambigrammatic sequence which  
 568 is independent of whether the reverse open reading frames are  
 569 translated into functional proteins. This argument led us to  
 570 discover the Robin segment of ZMV, and suggests that more  
 571 ambigrammatic narnaviruses with at least two segments will  
 572 be discovered by metagenomic surveys, when suitable data sets  
 573 become available. Similarly, the elusive Robin segment should  
 574 already be hiding in datasets of narnaviruses descended from  
 575 the common ancestor of CNV and ZMV.

**Table 3.** Summary of results for codon-level mutations. The numbers of single-nucleotide synonymous and non-synonymous mutations are  $N_{sy}$  and  $N_{ns}$  respectively,  $N_{mult}$  is the number of mutations with more than one base changed,  $R_{exp}$  is the null value of  $R = N_{ns}/N_{sy}$ , and  $f_{mult}$  if the fraction of mutations which have multiple-nucleotide changes.

Strand	$N_{sy}$	$N_{ns}$	$N_{mult}$	$R = N_{ns}/N_{sy}$	$R_{exp}$	$R/R_{exp}$	$f_{mult}$
CNV-RdRp-fwd	623	189	123	0.303	2.37	0.128	0.12
ZMV-RdRp-fwd	170	59	13	0.347	2.14	0.162	0.012
CNV-Robin-fwd	112	141	89	1.26	2.34	0.538	0.45
ZMV-Robin-fwd	49	61	14	1.24	2.35	0.528	0.046
CNV-RdRp-comp	136	676	123	4.97	2.43	2.04	0.12
ZMV-RdRp-comp	50	179	13	3.58	2.14	1.67	0.012
CNV-Robin-comp	66	187	89	2.83	2.39	1.23	0.45
ZMV-Robin-comp	32	78	14	2.43	2.28	1.07	0.046

**Table 4.** Summary of results of the mutational hotspots test. Left panel: values of the average number of elements of the variant set,  $\langle n(k) \rangle$  and of the average fraction of non-consensus codons,  $\langle f(k) \rangle$ , for double-synonym sites, and for the other sites. Right panel:  $N$  is the number of loci in the alignment,  $N_{ds}$  is the number of double-synonym loci, and  $R_n$ ,  $R_f$  are the ratios of  $\langle n(k) \rangle$  and  $\langle f(k) \rangle$  at double-synonym sites to their values at other sites. The differences of these ratios from unity do not appear significant.

Sample	$\langle n(k) \rangle$	$\langle f(k) \rangle$
Double syns., CNV-RdRp	0.954	0.161
Other codons, CNV-RdRp	0.968	0.155
Double syns., ZMV-RdRp	1.20	0.042
Other codons, ZMV-RdRp	1.23	0.050
Double syns, CNV-Robin	1.76	0.195
Other codons, CNVRobin	1.48	0.169
Double syns, ZMV-Robin	0.889	0.096
Other codons, ZMV-Robin	0.960	0.097

Gene	$N$	$N_{ds}$	$R_n$	$R_f$
CNV-RdRp	1033	220	0.986	1.044
ZMV-RdRp	1075	219	0.975	0.840
CNV-Robin	272	54	1.19	1.16
ZMV-Robin	304	81	0.926	0.978

**Table 5.** Results for the mutational codon frequency test:  $N$  is the number of loci in the alignment,  $N_a$  is the number of mutationally active double-synonym loci, and  $N_s$ ,  $N_d$  are, respectively, the numbers of single and double synonym mutations.

Sample	$N$	$N_a$	$N_s$	$N_d$	$R$	$R_0$	$R/R_0$
CNV-RdRp	1033	136	151	60	2.51	3.02	0.83
ZMV-RdRp	1075	219	33	20	1.65	3.21	0.51
CNV-Robin	272	40	24	3	8.00	3.21	2.49
ZMV-Robin	304	59	20	4	4.00	4.04	0.99

576 Our studies of polymorphisms in the forward direction  
 577 indicate that both RdRp and Robin are under purifying  
 578 selection. In the case of RdRp the amino acid sequence is  
 579 strongly conserved, but the Robin sequence is not.

580 The role of the RdRp coding fragment is already understood.  
 581 This makes it plausible that the other fragment plays a role  
 582 which facilitates the evolution of ambigrams. If the lack of  
 583 stop codons on the complementary strand is not required to  
 584 allow protein synthesis, we can surmise that its role is to allow  
 585 ribosomes to associate with the complementary strand. Having  
 586 RNA segments able to be covered by ribosomes may provide  
 587 some protection for the viral RNA against degradation.

588 Recent experiments indicate that ambigrammatic narnavirus  
 589 genes display unusual ribosome profiles, with a ‘plateau’  
 590 structure (Retallack et al., 2020). It has been argued (Wilkinson  
 591 et al., 2021) that the plateaus indicate that the ribosomes  
 592 attached to the viral RNA become stalled, creating a cover  
 593 (see also Cepelewicz (2020)). The ambigram property allows  
 594 binding of ribosomes to both strands, hiding the viral RNA  
 595 from host defence and degradation mechanisms. We can surmise  
 596 that there exists a molecule which binds to the 3’ end of the  
 597 viral RNA, preventing release of ribosomes (Wilkinson et al.,  
 598 2021). It is possible that Robin plays a role in this process, by  
 599 creating a protein which blocks ribosome detachment at 3’ end.  
 600 Alternatively, it might be proposed that the ribosome ‘traffic  
 601 jam’ is a consequence of the structure of the RdRp itself, due to  
 602 formation of RNA hairpins. However, these would have to trade  
 603 off against RdRp function. The proposed mechanism involving  
 604 Robin making a blocking protein has the advantage that the  
 605 RdRp works efficiently when the viral RNA concentration is  
 606 small. Later, after it has duplicated many copies of itself and of  
 607 Robin, the Robin protein attaches to the viral RNA and creates  
 608 stalled polysomes, protecting the viral RNA from degradation.

609 There may, however, be additional viral genes involved  
 610 in ambigrammatic narnavirus infections, and there are many

possible roles for the Robin gene. It could code a protein which  
 611 inhibits the mechanism of ‘no-go-decay’, which releases stalled  
 612 ribosomes, play a role in the viral suppression of RNAi (Mierlo  
 613 et al., 2014) or in formation of syncytia or viral particles.  
 614 Without a better understanding of the narnavirus lifecycle in  
 615 arthropods it is not certain whether Robin does code for a  
 616 protein which blocks detachment of ribosomes.

We did search the CNV dataset for further fragments of  
 618 ambigrammatic viral RNA, which might be candidates for  
 619 coding additional genes. A search for additional ambigrammatic  
 620 sequences greater than 200nt in length did not produce any  
 621 candidates.

A recent preprint (Retallack et al., 2020) presents evidence  
 623 that inserting mutations in the RdRp sequence which are  
 624 synonymous in the forward reading frame but introduce stop  
 625 codons in the reverse frame reduces the fitness of the virus. The  
 626 mutations were clustered close to the 3’ end of the RdRp gene.  
 627 These observations could be interpreted as indicating that the  
 628 reverse reading frame codes for a functional protein or that all  
 629 ORFs in the cell may be translated in a ‘leaky’ way. However,  
 630 changing the RNA sequence may also interfere with the action  
 631 of molecules which bind to the RdRp strand.  
 632

## 633 Competing interests

There is NO Competing Interest. 634

## 635 Author contributions

GD devised and directed the search for an analog of Robin  
 636 in the ZMV sequence archive. MW produced a draft of the  
 637 manuscript following discussions with the other authors about  
 638 the recent discovery of a narnavirus system which has two  
 639 ambigrammatic genes. All authors contributed to writing the  
 640 manuscript, and reviewed the manuscript before submission.  
 641

## 642 Acknowledgments

We thank Hanna Retallack and Joe DeRisi for discussions of  
 643 their experimental studies of narnaviruses and Amy Kistler  
 644 for assistance with narnaviral genomes and for comments on  
 645 a draft. We would like to thank Mang Shi and Edward C  
 646 Holmes for sharing assembled contigs from Australian and  
 647 Chinese mosquito metagenomic datasets. G.H. and D.Y. were  
 648 supported by the Chan Zuckerberg Biohub; MW thanks the  
 649 Chan Zuckerberg Biohub for its hospitality.  
 650



## 651 Data availability

## 652 References

- 653 B. G. Barrell, G. M. Air, and C. A. Hutchison. Overlapping  
654 genes in bacteriophage phiX174. *Nature*, 264:34–41, 1976.  
655 doi: 10.1038/264034a0.
- 656 J. Batson, G. Dudas, E. Haas-Stapleton, A. L. Kistler, L. M. Li,  
657 P. Logan, K. Ratnasiri, and H. Retallack. Single mosquito  
658 metatranscriptomics recovers mosquito species, blood meal  
659 sources, and microbial cargo, including viral dark matter.  
660 bioRxiv: <https://doi.org/10.1101/2020.02.10.942854>, 2020.
- 661 J. Cepelewicz. New clues about ‘ambigram’ viruses  
662 with strange reversible genes. *Quanta Magazine*,  
663 2020. URL [https://www.quantamagazine.org/  
664 new-clues-about-ambigram-viruses-with-strange-reversible-genes-20200212/](https://www.quantamagazine.org/new-clues-about-ambigram-viruses-with-strange-reversible-genes-20200212/).
- 665 N. Chirico, A. Vianelli, and R. Belshaw. Why genes overlap in  
666 viruses. *Proc Biol Sci.*, 277:1701, 2010. doi: 10.1098/rspb.  
667 2010.1052.
- 668 A. G. Cobián Güemes, M. Youle, V. A. Cantú, B. Felts,  
669 J. Nulton, and F. Rohwer. Viruses as winners in the game  
670 of life. *Annual Review of Virology*, 3(1):197–214, 2016.  
671 doi: 10.1146/annurev-virology-100114-054952. URL <https://doi.org/10.1146/annurev-virology-100114-054952>.
- 672 J. DeRisi, G. Huber, A. Kistler, H. Retallack, M. Wilkinson,  
673 and D. Yllanes. An exploration of ambigrammatic sequences  
674 in narnaviruses. *Sci. Rep.*, 9:17982, 2019. doi: 10.  
675 1038/s41598-019-54181-3. URL [https://doi.org/10.1038/  
676 s41598-019-54181-3](https://doi.org/10.1038/s41598-019-54181-3).
- 677 A. M. Dinan, N. I. Lukhovitskaya, I. Olenraite, and A. E.  
678 Firth. A case for a negative-strand coding sequence in a group  
679 of positive-sense rna viruses. *Virus Evolution*, 6:veaa007,  
680 2020. doi: <https://doi.org/10.1093/ve/veaa007>.
- 681 L. Fu, B. Niu, Z. Zhu, S. Wu, and W. Li. CD-HIT:  
682 accelerated for clustering the next-generation sequencing  
683 data. *Bioinformatics*, 28(23):3150–3152, Dec. 2012. ISSN  
684 1367-4811. doi: 10.1093/bioinformatics/bts565.
- 685 B. I. Hillman and G. Cai. The family *narnaviridae*: simplest  
686 of RNA viruses. In S. A. Ghabrial, editor, *Mycoviruses*,  
687 volume 86 of *Advances in Virus Research*, pages 149–176.  
688 2013. doi: 10.1016/B978-0-12-394315-6.00006-4.
- 689 M. Kimura. A simple method for estimating evolutionary  
690 rates of base substitutions through comparative studies of  
691 nucleotide sequences. *J. Molecular Evolution*, 16:111–20,  
692 1980.
- 693 J. T. v. Mierlo, G. J. Overheul, B. Obadia, K. W. R. v.  
694 Cleef, C. L. Webster, M.-C. Saleh, D. J. Obbard, and  
695 R. P. v. Rij. Novel Drosophila Viruses Encode Host-Specific  
696 Suppressors of RNAi. *PLOS Pathogens*, 10(7):e1004256,  
697 July 2014. ISSN 1553-7374. doi: 10.1371/journal.ppat.  
698 1004256. URL [https://journals.plos.org/plospathogens/  
699 article?id=10.1371/journal.ppat.1004256](https://journals.plos.org/plospathogens/article?id=10.1371/journal.ppat.1004256).
- 700 C. W. Nelson, Z. Ardern, and X. Wei. Olgene: Estimating  
701 natural selection to predict functional overlapping genes.  
702 *Molecular Biology and Evolution*, 37:2440–2449, 2020. doi:  
703 <https://doi.org/10.1093/molbev/msaa087>.
- 704 H. Retallack, K. D. Popova, M. T. Laurie, S. Sunshine, and J. L.  
705 DeRisi. Persistence of ambigrammatic narnaviruses requires  
706 translation of the reverse open reading frame. bioRxiv  
707 preprint, doi: <http://10.1101/2020.12.18.423567> 2020.
- 708 M. Shi, X.-D. Lin, J.-H. Tian, L.-J. Chen, X. Chen, C.-X.  
709 Li, X.-C. Qin, J. Li, J.-P. Cao, J.-S. Eden, J. Buchmann,  
710 W. Wang, J. Xu, E. C. Holmes, and Y.-Z. Zhang. Redefining  
711 the invertebrate RNA virosphere. *Nature*, Nov. 2016. ISSN  
712 1476-4687. doi: 10.1038/nature20167.
- M. Shi, P. Neville, J. Nicholson, J.-S. Eden, A. Imrie, and E. C.  
Holmes. High-Resolution Metatranscriptomics Reveals the  
Ecological Dynamics of Mosquito-Associated RNA Viruses in  
Western Australia. *Journal of Virology*, 91(17), Sept. 2017.  
ISSN 1098-5514. doi: 10.1128/JVI.00680-17.
- M. Wilkinson, D. Yllanes, and G. Huber. Polysomally  
protected viruses. 2021. URL [https://arxiv.org/abs/2102.  
00316](https://arxiv.org/abs/2102.00316). arXiv:2102.00316.