

# Experimental Evolution of Pseudogenization and Gene Loss in a Plant RNA Virus

Mark P. Zwart,<sup>\*,1</sup> Anouk Willemsen,<sup>1</sup> José-Antonio Daròs,<sup>1</sup> and Santiago F. Elena<sup>1,2</sup>

<sup>1</sup>Instituto de Biología Molecular y Celular de Plantas, Consejo Superior de Investigaciones Científicas-UPV, València, Spain

<sup>2</sup>The Santa Fe Institute

\*Corresponding author: E-mail: marzwa@ibmcp.upv.es.

Associate editor: Howard Ochman

## Abstract

Viruses have evolved highly streamlined genomes and a variety of mechanisms to compress them, suggesting that genome size is under strong selection. Horizontal gene transfer has, on the other hand, played an important role in virus evolution. However, evolution cannot integrate initially nonfunctional sequences into the viral genome if they are rapidly purged by selection. Here we report on the experimental evolution of pseudogenization in virus genomes using a plant RNA virus expressing a heterologous gene. When long 9-week passages were performed, the added gene was lost in all lineages, whereas viruses with large genomic deletions were fixed in only two out of ten 3-week lineages and none in 1-week lineages. Illumina next-generation sequencing revealed considerable convergent evolution in the 9- and 3-week lineages with genomic deletions. Genome size was correlated to within-host competitive fitness, although there was no correlation with virus accumulation or virulence. Within-host competitive fitness of the 3-week virus lineages without genomic deletions was higher than for the 1-week lineages. Our results show that the strength of selection for a reduced genome size and the rate of pseudogenization depend on demographic conditions. Moreover, for the 3-week passage condition, we observed increases in within-host fitness, whereas selection was not strong enough to quickly remove the nonfunctional heterologous gene. These results suggest a demographically determined “sweet spot” might exist, where heterologous insertions are not immediately lost while evolution can act to integrate them into the viral genome.

**Key words:** genome evolution, plant virus, horizontal gene transfer, pseudogenization, fitness, next-generation sequencing.

## Introduction

Virus genomes are highly streamlined. Compared with more complex organisms, viruses tend to have small genomes with 1) a high percentage of coding sequences, 2) none or little intronic sequences, and 3) only short stretches of intergenic sequences (Lynch 2006; Belshaw et al. 2007). Moreover, viruses have evolved strategies to further compress their genomes, such as frameshifts and overlapping open reading frames (ORFs) (e.g., Belshaw et al. 2007; Chung et al. 2008). Field observations suggest that genome shrinkage sometimes occurs during epidemic spread and might be linked to increased within-host fitness and be adaptive for white spot syndrome virus (WSSV), a large DNA virus (Marks et al. 2005; Zwart, Dieu, et al. 2010). Moreover, it appears to be a very general observation that viruses expressing heterologous genes tend to be unstable (Chapman et al. 1992; Dolja et al. 1993; Guo et al. 1998; Pijlman et al. 2001; Chung et al. 2007; Paar et al. 2007). Furthermore, under conditions maximizing selection for competitive fitness—exemplified by undiluted serial passage in cultured cells—viruses tend to rapidly evolve defective interfering particles (DIPs): viruses with large genomic deletions are unable to replicate autonomously but with a replicative advantage at high multiplicities of infection (Huang 1973; Simon et al. 2004; Zwart et al. 2008; Pathak and Nagy 2009). All these observations suggest that genome size is under strong selection for viruses and that

having unnecessary genomic sequences has fitness costs. By contrast, striking cases of genome shrinkage have been found for obligate host-dependent species of bacteria, but this shrinkage appears to be the result of a mutational bias toward deletions and genetic drift (Ochman and Davalos 2006; Kuo and Ochman 2009).

Viruses play an important evolutionary role as vectors for horizontal gene transfer (HGT) in the genomes of their hosts (Canchaya et al. 2003; Belshaw et al. 2004; Routh et al. 2012). It is moreover becoming increasingly apparent that HGT is also widespread in most viruses and is a key mechanism in their evolution (Hughes and Friedman 2005; Dolja and Koonin 2011; Koonin and Dolja 2012; Liu et al. 2011, 2012; Yutin and Koonin 2012). Striking innovations such as the DNA-RNA virus hybrid (Diemer and Stedman 2012) and the cooption of an entire host immunity mechanism by a phage (Seed et al. 2013) exemplify how HGT can empower the evolutionary process. However, strong selection for genome size would in principle be an impediment to HGT. In order for a heterologous sequence to be beneficial to the recipient virus, it must be accommodated into the virus genome, transcriptome, and proteome. Mutation and selection must therefore act on newly transferred heterologous sequences, but to do so, these sequences must not be purged right upon acquisition because of selection for a streamlined genome. The mutational supply may favor the deletion of

© The Author 2013. Published by Oxford University Press on behalf of the Society for Molecular Biology and Evolution.

This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/3.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact [journals.permissions@oup.com](mailto:journals.permissions@oup.com)

Open Access

heterologous sequences; deletion of a sequence by recombination probably has a greater likelihood than the occurrence of beneficial mutations functionally integrating this element in the virus. If the effects of deletions of heterologous sequences are beneficial on the short term, how can HGT be common in viruses?

One possible answer to this question is that selection for genome size is not a very strong force. First, it is not at all clear that the metabolic cost of additional genetic material is sufficient to account for the expected fitness costs (Lynch 2007). Second, experimental results on the relationship between genome size and replicative fitness are ambiguous. When comparing phages with different genome organizations adapted for fast replication, the expected relationship was not found (Bull et al. 2004). When expressing different-size marker proteins in Sendai virus, an inverse relationship between insert size and replication was found in cultured cells (Sakai et al. 1999). However, this relationship was not observed in vivo (Sakai et al. 1999) and, furthermore, the use of sequences coding for different marker proteins makes the comparison troublesome (Majer et al. 2013). Moreover, instability of viruses expressing heterologous sequences (Chapman et al. 1992; Dolja et al. 1993; Guo et al. 1998; Chung et al. 2007) may depend on many environmental factors (Paar et al. 2007), and even subtleties of the heterologous sequence, such as guanine-uracil content (Lee et al. 2002). Finally, experiments corroborating the relationship between genome size and fitness for WSSV were performed with field isolates (Marks et al. 2005; Zwart, Dieu, et al. 2010), and hence, other genetic variation could be a confounding factor. In considering whether HGT is really implausible, we therefore need to ask whether increases in genome size really have appreciable fitness costs and what fitness components might be affected.

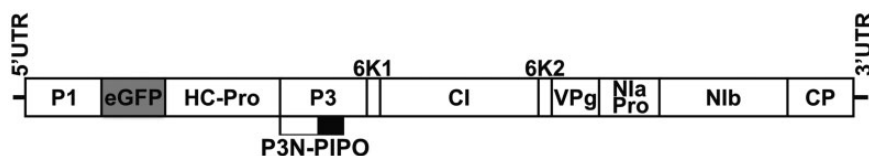
Here we explore the process of pseudogenization in virus genomes by means of experimental evolution. As a model system, we consider a plant RNA virus expressing a nontoxic heterologous gene, whose expression has been engineered to minimally disrupt the viral genome. We first looked for conditions under which the heterologous gene would be maintained in the genome. It has been shown that the time period between two consecutive transmission events, that is, the time a viral population has to expand between two consecutive bottlenecks, can be instrumental in determining genome stability in plant RNA viruses (Dolja et al. 1993). However, in this study, the heterologous gene was expressed as a fusion with one of the viral cistrons and, given the strong effects on viral accumulation (Dolja et al. 1993), must therefore be seen as a deleterious, rather than a

merely nonfunctional addition. We then consider whether the deletion of the heterologous gene was adaptive and what virus characteristics it modifies. Finally, we consider whether there are conditions that fulfill two requirements: 1) the heterologous gene has a high probability of being maintained in the evolved virus population and 2) there is evidence that the virus population is under positive selection and experiences increases in fitness. We think the combination of conditions is relevant to the context of HGT in viruses. If these two conditions are fulfilled, then in principle, a heterologous sequence can persist for long periods of time in the virus population, while increases in fitness imply that natural selection is acting on the population and could “tinker” with the heterologous gene, sequences regulating its expression, and other loci interacting with the heterologous gene, potentially and functionally integrating it into the viral genome. On the other hand, a heterologous gene may not be lost in a virus population subject to high levels of genetic drift, but it is then also unlikely that natural selection acts to functionally integrate it. Similarly, a virus population under strong positive selection in which the heterologous gene is lost is also a dead end for HGT. However, simultaneously having maintenance of the heterologous gene and increases in viral fitness suggests the occurrence of a “sweet spot” that could help explain how HGT occurs in viruses.

## Results and Discussion

### Results of Serial Passage Experiments

As a model system for virus expressing a heterologous gene without appreciable toxicity or disruption of viral replication, we used a variant of *Tobacco etch virus* (TEV; genus *Potyvirus*, family *Potyviridae*). TEV is a positive-sense single-stranded RNA virus that encodes a polyprotein autocatalytically cleaved into ten mature proteins (Riechmann et al. 1992) and a partially overlapping ORF with a +2 frameshift (Chung et al. 2008). The TEV variant we used expresses enhanced green fluorescent protein (eGFP) as a separate cistron between P1 and HC-Pro by introducing a second Nla-Pro proteolytic site downstream of the eGFP sequence while retaining the existing C-terminal site in P1 (fig. 1) (Zwart et al. 2011). Evolution experiments were performed in *Nicotiana tabacum* L. cv. Xanthi plants. In brief, 4-week-old plants were inoculated with high virus doses, and passages lasting either 1 week (for a total of 27 consecutive passages), 3 weeks (nine passages) or 9 weeks (three passages) were performed. Each “passage” is the infection of a single plant and harvesting of tissues at the designated time (i.e., the “passage duration”). Therefore, although the passage duration varied among



**Fig. 1.** Scheme of TEV-eGFP. Lines represent the viral 5′- and 3′-untranslated regions (5′-UTR and 3′-UTR), the gray box represents eGFP, open boxes represent the viral cistrons P1, HC-Pro, P3, 6K1, CI, 6K2, VPg, Nla-Pro, Nlb, and CP, whereas P3N-PIPO is indicated by the lower box.

treatments, each lineage evolved for the same total time (27 weeks) in *N. tabacum*. At the end of a passage, all the leaves above the inoculated leaf were collected, pooled, and used to obtain the inoculum for the next round of serial passaging. Ten independent lineages were generated and maintained for each passage duration (1, 3, or 9 weeks). An overview of the experimental setup used is given in figure 2, and further details are given in the Materials and Methods.

eGFP expression was readily apparent in infected plants (fig. 3A) and was used as a first indication of whether the heterologous gene was intact. Partial losses of fluorescence (fig. 3B) almost always preceded complete losses of fluorescence (fig. 3C). One out of ten 1-week lineages showed a partial loss of fluorescence, first observed at passage 7 and maintained until passage 27 (fig. 3E and F). Two out of ten 3-week lineages showed a loss of fluorescence. All but one 9-week passage showed decreased levels of fluorescence after a single passage, and all lineages showed a complete loss of fluorescence after two passages. Reverse transcription polymerase chain reaction (RT-PCR) with primers flanking eGFP (Materials and Methods) confirmed the occurrence of genomic deletions in all lineages with a partial or complete loss of fluorescence (fig. 3G). These results are congruent with previous work with TEV (Dolja et al. 1993), except that the deletion of the heterologous gene occurs much more slowly here, as anticipated.

We then inoculated *N. tabacum* with TEV-eGFP and, at 9 weeks of infection, harvested every fifth leaf up the stem. Because the infection progresses linearly as the plant grows, these leaves enable us to monitor a qualitative time course of evolution within the plant. All plants had by then reached the 45-leaf stage, except for one that had only 40 leaves. We performed RT-qPCR on individual leaves to ascertain at what leaf level deletions occurred, and if they were subsequently maintained in the population. This analysis can be performed since the virus moves mainly upward in the plant (Dolja et al. 1992). In most cases, once a deletion was detected in one leaf, it was maintained and fixed in the superior leaves (fig. 4). This observation suggests that selection for deletion variants is very strong in this experiment, being a stronger evolutionary force than genetic drift within the host. The first leaf in which a novel deletion was detected was not uniformly distributed over all tested leaves (leaves 5–45; one-sample Kolmogorov–Smirnov test:  $n = 13$ ;  $P = 0.019$ ); new deletion variants were usually first detected in higher leaves (mean  $\pm$  standard deviation [SD] =  $32.31 \pm 9.92$ ). This result suggests that passage duration is important to the dynamics of heterologous gene deletion because it regulates the amount of expansion that occurs between bottlenecks.

During virus infection of mechanically inoculated plants, genetic bottlenecks in the virus population can occur during primary infection of the inoculated leaf and the subsequent entry into systemically infected leaves (Hall et al. 2001; Sacristán et al. 2003; Gutiérrez et al. 2012). For TEV infection of tobacco plants, the number of primary infection foci in the inoculated leaf is a good estimator of the number of founders, whereas the subsequent bottlenecks during entry into systemically infected leaves do not appear to be severe

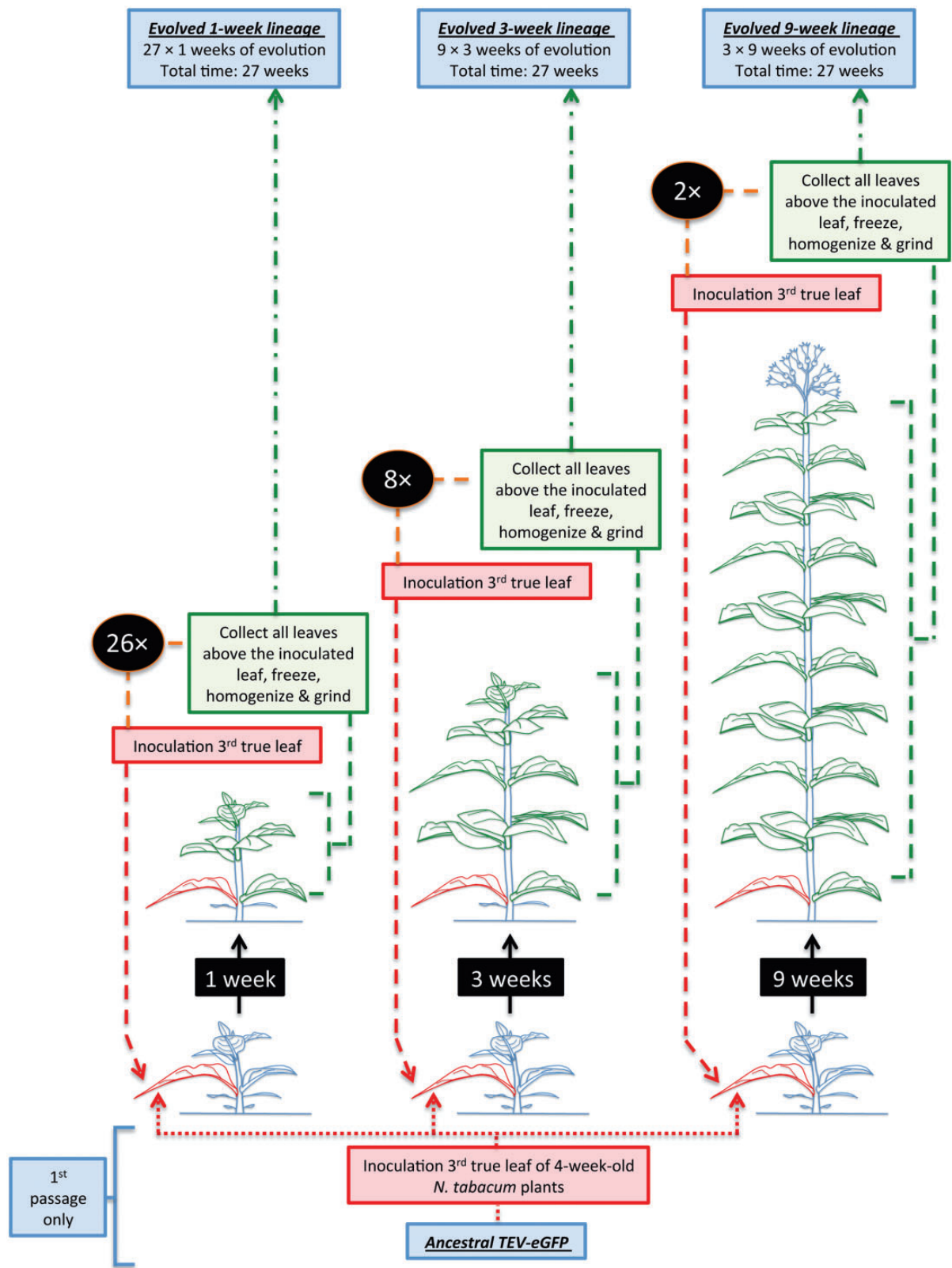
(Zwart et al. 2011). We therefore inoculated eight plants with TEV-eGFP using the same procedure and conditions, as during serial passaging (see Materials and Methods; homogenized tissue of plants infected with TEV-eGFP was used as an inoculum), and counted the number of primary infection foci (Zwart et al. 2011). The mean number of foci observed  $\pm$ SD was  $417 \pm 140$ , suggesting that although there is a bottleneck at the start of infection, it is not too severe. Nevertheless, this bottleneck could remove variation generated de novo during the previous passage and hereby limit the variation upon which selection can act. Longer passage duration would allow beneficial variation to increase in frequency, thus making it less likely to lose such variants due to genetic drift at the next transmission event. The occurrence of genetic bottlenecks therefore reinforces the idea that passage duration might be important to the evolutionary dynamics in this system.

### Genome Sequences of Evolved Lines with Deletions

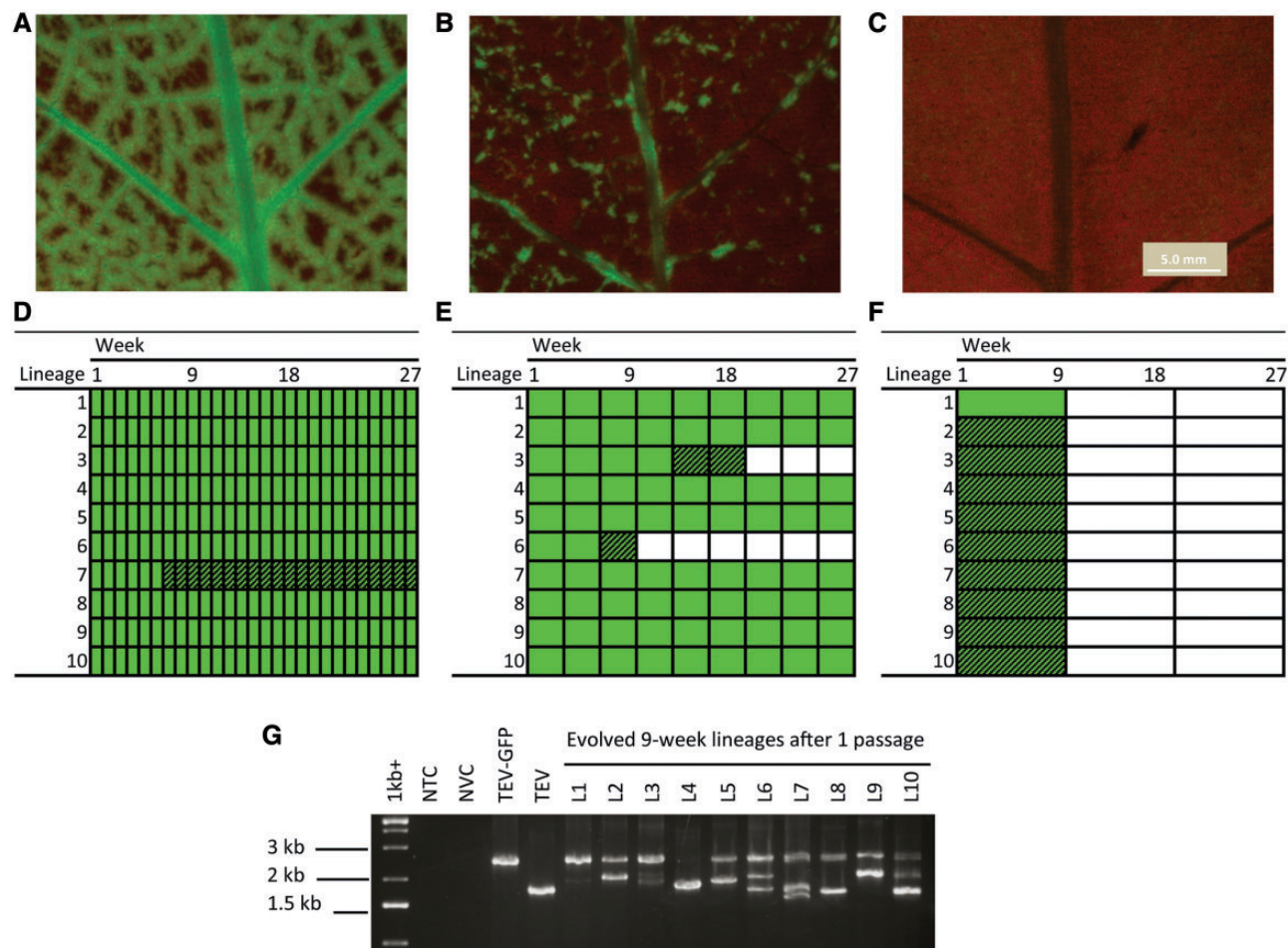
All evolved lineages in which deletions had been detected by RT-PCR were fully sequenced by Illumina next-generation sequencing (NGS). We developed an approach (Materials and Methods) for mapping large genomic deletions (i.e., deletions larger than the read size). We consistently saw pseudogenization or complete loss of eGFP in all these lineages (fig. 5A). None of these deletions included the C-terminus of P1, while for 7 out of 13 lineages, these deletions included N-terminal regions of HC-Pro, similar to the previous results (Dolja et al. 1993). The N-terminal region of HC-Pro is not essential for replication and movement (Dolja et al. 1993; Cronin et al. 1995) but has been implicated in vector-borne transmission (Thornbury et al. 1990; Atreya et al. 1992), which is not a selective force in our mechanical transmission passage experiments. In the seven lineages with deletions extending into HC-Pro, the remains of eGFP were fused with HC-Pro, while the proteolytic site between P1 and eGFP remained intact. The start of genomic deletion (5' end) was not uniformly distributed (one-sample Kolmogorov–Smirnov test:  $n = 13$ ;  $P = 0.006$ ), and there is clustering at the 5' of the eGFP cistron (fig. 5B and C), suggesting the existence of a hotspot for recombination or the unviability of any deletions in the P1-eGFP proteolytic site. On the other hand, the 3' end of the genomic deletion was uniformly distributed (one-sample Kolmogorov–Smirnov test:  $n = 13$ ;  $P = 0.130$ ).

We performed additional analyses to detect minority variants with different deletion sizes in sequenced lineages (Materials and Methods). Although minority variants were sometimes detected in the 3- and 9-week lineages, these were always present at low frequencies ( $<1.5\%$ ). Only in the case of the 1-week lineage with a partial eGFP loss was a minority variant present: 47.1% of the population was composed of a variant with intact eGFP. The deletion in the majority variant extends beyond the HC-Pro N-terminal regions nonessential for replication and movement (Cronin et al. 1995), suggesting this majority variant is not able to replicate without being complemented by the full-length variant. When plants were inoculated with low virus doses

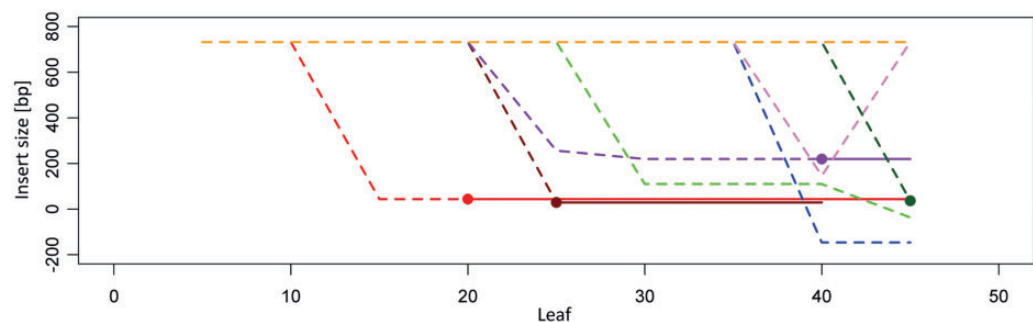




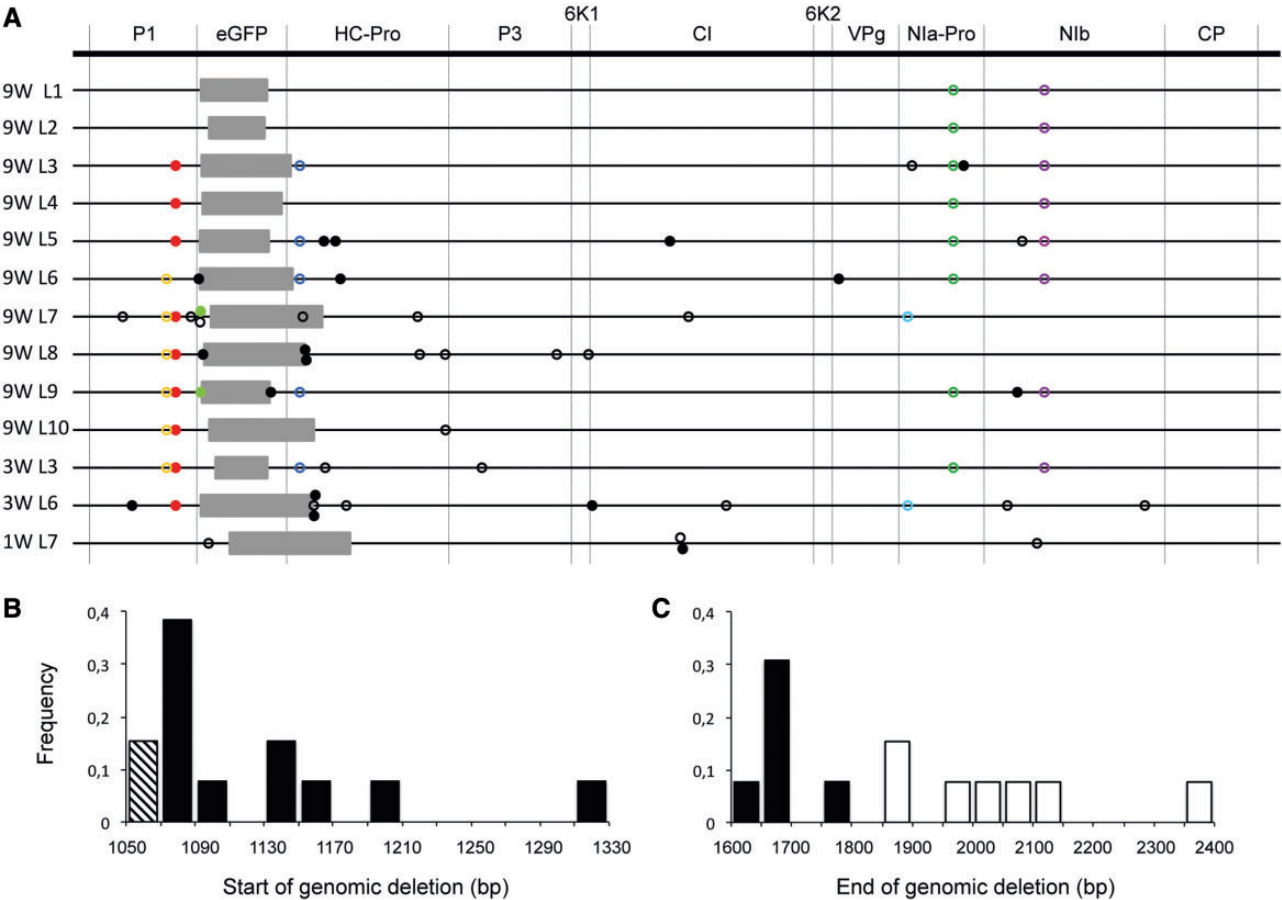
**FIG. 2.** Overview of the experimental setup employed in the study. At the start of the serial passage experiment, 4-week-old *N. tabacum* plants were mechanically inoculated with TEV-eGFP in the third true leaf (indicated in red above). At the end of the designated passage duration (1, 3, or 9 weeks), all leaves above the inoculated leaf, which are indicated in green above, were collected and stored at  $-80^{\circ}\text{C}$ . The frozen tissue was then homogenized, and a sample of the homogenized tissue was ground to a fine powder. For inoculation of subsequent passages, powder was resuspended in inoculation buffer and new *N. tabacum* plants were inoculated. Although the duration of the passages varied (1, 3, and 9 weeks), the number of passages was set so that the total time each lineage was present in plants was the same, being 27 weeks for all lineages. For each passage duration used, ten independent lineages were taken. Note that the figure is only schematic: the distance between leaf layers has been exaggerated, and after 9 weeks of infection, *N. tabacum* plants are in reality relatively taller than depicted here.



**Fig. 3.** Loss of eGFP fluorescence and sequence during serial passages: (A) An *N. tabacum* leaf that is completely symptomatic, indicating heavy virus infection, showed 1 week postinoculation with an evolved TEV-eGFP lineage with no loss of fluorescence. (B) A completely symptomatic leaf showed 1 week postinoculation with an evolved TEV-eGFP lineage with a partial loss of fluorescence (1-week passage lineage 7). (C) A completely symptomatic leaf showed 1 week postinoculation with an evolved TEV-eGFP lineage with a complete loss of fluorescence. (D) Observed fluorescence during serial passage of TEV-eGFP for 1-week passages. For panels D–E, green squares indicate no loss of eGFP fluorescence (as in panel A), hatched green squares indicate a partial loss of fluorescence (as in panel B) in parts or all of the plant, and white squares indicate no fluorescence was observed in the whole plant (as in panel C). (E) Observed fluorescence for 3-week passages. (F) Observed fluorescence for 9-week passages. (G) An agarose gel with RT-PCR products for deletions in the eGFP locus of TEV-eGFP is shown for the first passage of the 9-week lineages. The value 1 kb + indicates the lane with a 1 kb + DNA ladder, NTC is the nontemplate control, and NVC is the nonvirus control, a mock-inoculated healthy plant. TEV-eGFP, the ancestral virus for the evolution experiments, and TEV, the wild-type virus without the heterologous gene inserted, are included for comparison. Note that in each evolved lineage, deletions of eGFP are visible, although their frequency appears to vary. In lineage 1, the band corresponding to the ancestral virus is still very strong, whereas for lineage 4, it is not longer visible.



**Fig. 4.** A single 9-week passage of TEV-eGFP was performed in eight plants, and the smallest observed insert size at the eGFP locus (ordinate) was measured every fifth leaf (abscissa). Dotted lines indicate the full-length eGFP sequence was still detected, solid lines indicate it is no longer detected, and a circle indicates the point at which the full-length sequence is first no longer detected, and each replicate has a different color. Deletions were fixed in four out of eight replicates, and only in a single replicate were no deletions detected throughout infection (yellow). Only in one case (pink) was a deletion not maintained after it has first been observed (a deletion is observed in leaf 40 but not in leaf 45).



**Fig. 5.** Genome sequences of evolved lineages: (A) NGS data for the evolved lineages with deletions from the serial passage experiment (fig. 3D and F) are given. The names on the left identify lineages (e.g., 9W L1 is the final population of 9-week passage, lineage 1). Gray boxes indicate genomic deletions in the majority variant. Full circles and open circles are nonsynonymous and synonymous substitutions, respectively. Black substitutions occur in only one lineage, whereas color-coded substitutions are repeated in two or more lineages. For lineage 1W L7, only a single sequence is represented. However, note that in this lineage, another variant with the full-length TEV-eGFP genome is present at a frequency 47.1%. None of the single-nucleotide mutations in 1W L7 were fixed, suggesting they occurred after the genomic deletion and are present in only one of the two variants present. However, the sequence variation is minimal and the two variants are present at approximately same frequencies in this lineage. The single-nucleotide mutations can therefore not be assigned to the full-length or the deleted variant, although they are represented on the deletion variant in the figure. In all other lineages, the full-length virus was not detected and variants with other genomic deletions were present only at very low frequencies (<1.5%). (B) Histogram of the position of the start of genomic deletions in the evolved lines. For panels B and C, dark lines indicate regions in the eGFP cistron, white lines indicate regions in the viral genome, and hatched lines indicate the regions encompassing two cistrons (i.e., P1 and eGFP in panel B). (C) Histogram of the position of the end of the genomic deletion in evolved lines.

**Table 1.** In Vivo Cloning of 1-Week Passage Lineage 7 (1W L7).

Experiment	Mean Foci	Plants			
		Uninfected	eGFP Only	Mixture	No eGFP
1	1.022	9 (0.360)	12 (0.480)	4 (0.160)	0 (0.000)
2	1.514	18 (0.243)	36 (0.486)	20 (0.270)	0 (0.000)

NOTE.—Two replicate experiments were performed in which *Nicotiana tabacum* plants were infected with a 1:1,000 dilution of infectious sap of the final evolved population of 1W L7. This population had apparently harbored a virus variant that had lost fluorescence (fig. 3B), although this variant never went to fixation (fig. 3D). Even at low doses, all infected plants (as determined by symptoms 2 weeks after inoculation) contained only the eGFP-expressing virus (“eGFP only”; fluorescence patterns similar to fig. 3A) or both virus variants (“Mixture”; fig. 3B). A symptomatic plant without eGFP expression was never observed (“No eGFP”; fig. 3C), although the low dose resulted in a low mean number of primary infection foci of TEV-eGFP (“mean foci”) and many uninfected plants. RT-PCR was performed on five plants judged as being uninfected, infected only with the eGFP variant, or infected with a mixture of the two viruses, and the RT-PCR results were always congruent with microscopic observations. We therefore conclude the virus variant in the population that does not express eGFP, has therefore probably lost infectivity, or it has a very low infectivity. This variant is nevertheless surprisingly stable. The 1W L7 population was put through three 1-week passages or one 3-week passage, with ten replicates each. The presence of the virus variant not expressing eGFP could always be deduced from eGFP expression patterns (i.e., fig. 3B).



of this evolved lineage, we were indeed unable to in vivo clone the majority variant without intact eGFP (table 1).

When single-base substitutions were detected in sequenced lineages, there appeared to be convergent evolution in the 3- and 9-week lineages (fig. 5A). All these lineages contained at least one substitution present in another lineage, two substitutions were present in 8 out of 12 lineages, and one substitution was present in 9 out of 12 lineages. Overall, more than half of the substitutions found were present in other lineages, and some lineages contained only substitutions also present in other lineages (9-week lineages 1 and 2). In the single 1-week passage sequenced, none of the repeated substitutions were present, suggesting that convergent evolution did not occur under these conditions.

Of the substitutions found here, 28 were nonsynonymous and 53 were synonymous. Eleven nonsynonymous substitutions were convergent, while 29 synonymous substitutions were convergent. Synonymous substitutions were therefore more common than nonsynonymous substitutions, although both were equally likely among cases of convergence (Fisher's exact test  $P = 0.244$ ). Convergent, nonsynonymous substitutions were always found in the P1 cistron (A872G, N→S in 9 out of 13 lineages; all positions given relative to the original TEV-eGFP genome, GenBank KC918545), or the remaining 5' end of the eGFP cistron (A1085U, E→V in 2 lineages). Convergent synonymous substitutions were found in P1 (C795U in six lineages), HC-Pro (C1927A in five lineages), NIa-Pro (U7092C in two lineages and A7479C in eight lineages), and NIb cistrons (A8253C in eight lineages). The A7479C and A8253C substitutions always occurred together, suggesting synergistic epistasis or possibly even reciprocal sign epistasis, whereas the A7479C and A8253C never occurred together with U7092C, suggesting antagonistic epistasis. However, neither of these two effects was significant given the number of observations (table 2). Not a single mutation was detected in the coat protein (CP). The overall  $d_N/d_S$  ratio was significantly smaller than 1 (mean  $\pm$  SD =  $0.058 \pm 0.002$ ; z-test  $P < 0.001$ ), suggesting the polyprotein sequence is under purifying selection. Within-population single nucleotide polymorphisms (SNPs) were analyzed for every evolved lineage against their corresponding consensus sequence. Of the SNPs found, ten were synonymous and six were nonsynonymous. Five out of 13 evolved lineages contained the same synonymous SNP in the HC-Pro cistron (C1879A). None of the other synonymous and nonsynonymous SNPs were repeated in the evolved lineages.

Previous evolution experiments with TEV have also shown no genomic convergences for 1-week passages in *N. tabacum*:

after 15 weeks of evolution, no substitutions were repeated in different lineages (Bedhomme et al. 2012). These observations suggest that little adaptive evolution might occur when short 1-week serial passages are performed. Furthermore, the specific convergent mutations observed here were not observed in other 1-week passage experiments (Bedhomme et al. 2012; Tromas N, Zwart MP, Elena SF, unpublished manuscript). This suggests that these convergent mutations may be linked to the insertion of the eGFP gene. To test this possibility, we considered whether the most common mutations (C795U, A872G, A7479C, and A8253C) occurred in lineages of the wild-type TEV put through three 9-week passages in *N. tabacum* (Materials and Methods). In none of such lineages were any of these mutations found, strongly suggesting they are linked to the presence of eGFP. Although we sequenced only a small part of these evolved TEV genomes, we did find one mutation repeated in 3 out of 10 lineages (A6806G, K→E), suggesting there is at least some convergent evolution when the wild-type TEV is put through long passages.

Accumulation, Virulence, and Within-Host Competitive Fitness of Evolved Lineages

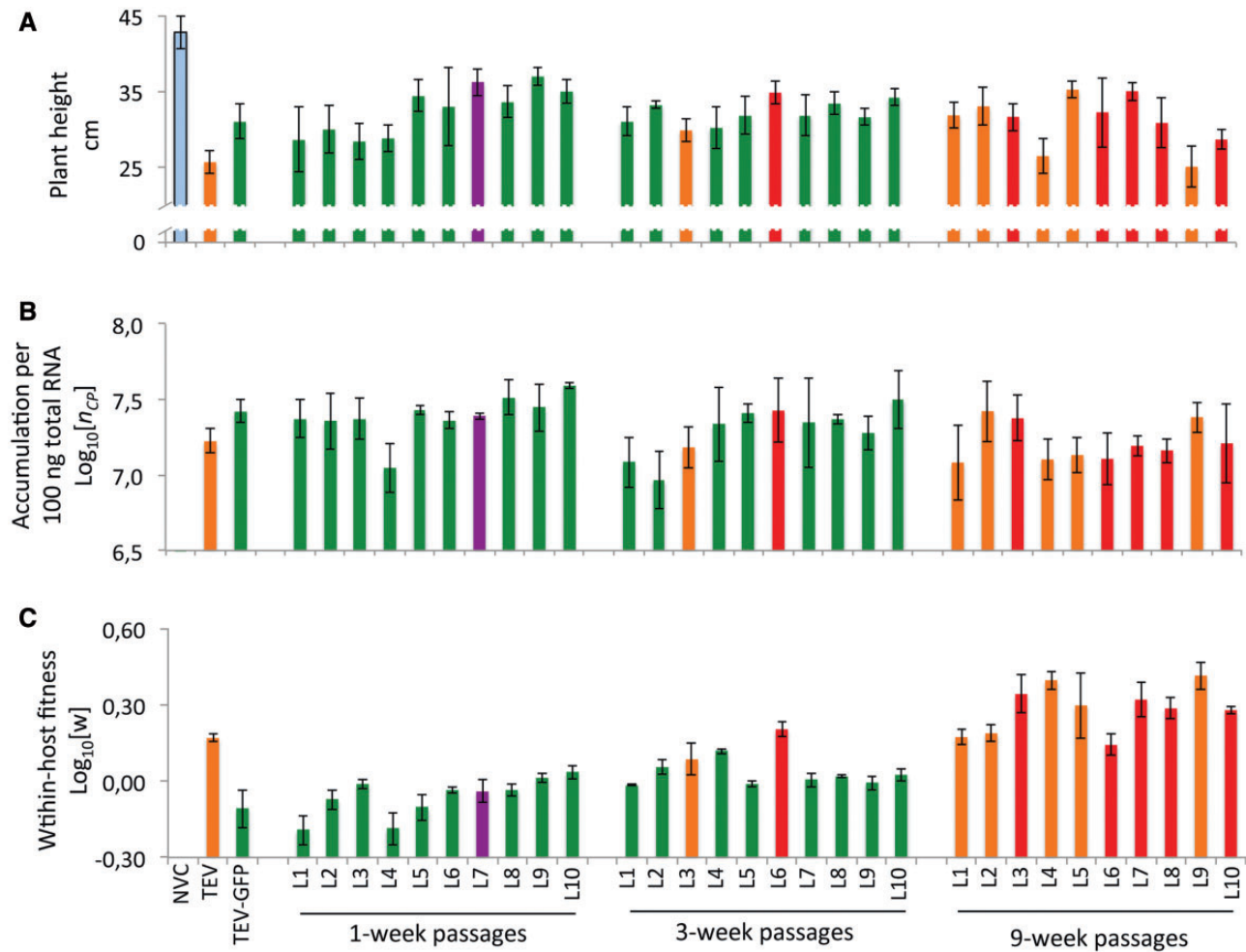
We biologically characterized all evolved lineages, in terms of their virulence and viral accumulation, and measured within-host competitive fitness (W; Materials and Methods). There was no effect of passage duration on either viral accumulation at 7 days post inoculation (dpi) or virulence (fig. 6A and B; table 3). On the other hand, there was a highly significant effect of passage duration on within-host competitive fitness, which increased significantly with passage duration (fig. 6C and table 3).

For the current experimental setup, we a priori expect that within-host competitive fitness will be under selection, because we are passaging a virus within a single host at high inoculation doses. Those virus variants that exist at the highest frequency at the end of infection are therefore most likely to be transferred, irrespective of accumulation levels of the entire population. We did not expect virulence to change, because we do not expect it to be under selection and it is probably not linked to within-host fitness in our model system (Carrasco, de la Iglesia, et al. 2007). Given that high inoculation doses are used, virus accumulation is not likely to be very important either, so long as it is enough to maintain infection in the next round of passaging. Note that a 1,000-fold dilution of an inoculum still causes moderate levels of infection, as shown by the in vivo cloning results (table 1). Published data (Carrasco, de la Iglesia, et al. 2007; Lalić et al.

Table 2. Cooccurrence of Single-Nucleotide Substitutions.

Substitution Combination	$\Pr(A) \Pr(B) = \Pr(A \cap B)$ .	$Obs(A \cap B)$ .	P Value
A7479C $\cap$ A8253C.	$(8/12)(8/12) = 0.444$	$8/12 = 0.667$	0.1501
(A7479C $\cap$ A8253C) $\cap$ U7092C.	$(8/12)(2/12) = 0.111$	$0/12 = 0$	0.3841

NOTE.—Here we test whether the cooccurrence, or lack thereof, is statistically significant for two groups of single-nucleotide substitution (i.e., substitution combination).  $\Pr(A) \Pr(B) = \Pr(A \cap B)$  is the product of the frequency of occurrence of the substitutions, the expected frequency at which we expect to see both substitutions in the absence of any epistatic interactions. We only considered the 3- and 9-week passages, because we do not think selection is acting on the 1-week lineage.  $Obs(A \cap B)$  gives the frequency at which the combination was observed, and P value is the significance as determined by comparison of predicted and observed values with an exact binomial test.



**Fig. 6.** Virulence, accumulation, and within-host fitness of evolved strains are given. (A) The height of control plants (NVC), and plants infected with TEV, TEV-eGFP (the ancestral virus for evolved lineages), and all lineages of evolved viruses are given. We consider the inverse of plant height as a proxy for virulence, though in the absence of significant differences in the data, we simply present the raw data. For all panels, green columns indicate no deletions in the heterologous gene (eGFP) were detected, orange indicates part or all of the eGFP was not present, and red indicates part of eGFP and the viral HC-Pro cistron are not present. The 1W L7 population is marked magenta because it contains a large deletion resulting in a virus apparently unable to infect on its own. (B) Virus accumulation, as measured by RT-qPCR. (C) The replicative advantage ( $W$ ) of the tested virus with respect to a common competitor, TEV-mCherry, is given, as determined by competition experiments and RT-qPCR (Materials and Methods). We consider  $W$  as an indicator of the within-host competitive fitness.

**Table 3.** Nested ANOVAs on Plant Height, Accumulation, and Within-Host Fitness of Evolved Lineages.

Trait	Source of Variation	SS	df	MS	F	P
Plant height	Treatment	66.040	2	33.020	0.791	0.463
	Lineage within treatment	1126.600	27	41.726	7.120	<0.001
	Error	703.200	120	5.826		
Accumulation	Treatment	41.440	2	20.720	0.391	0.680
	Lineage within treatment	1431.520	27	53.019	6.998	<0.001
	Error	909.200	120	7.577		
Fitness <sup>a</sup>	Treatment	1.881	2	0.940	48.534	<0.001
	Lineage within treatment	0.523	27	0.019	3.267	<0.001
	Error	0.356	60	0.006		
Fitness <sup>b</sup>	Treatment	0.100	1	0.100	7.868	0.013
	Lineage within treatment	0.203	16	0.013	4.096	<0.001
	Error	0.112	36	0.003		

<sup>a</sup>Fitness is a comparison of all lineages.

<sup>b</sup>Fitness compares only those lineages that have no fixed genomic deletions, ten of which are from the 1-week treatment and eight of which are from the 3-week treatment (see fig. 5). ANOVA was used for this comparison of two groups to allow lineage to be nested within treatment. Treatment is the passage duration (1, 3, or 9 weeks).



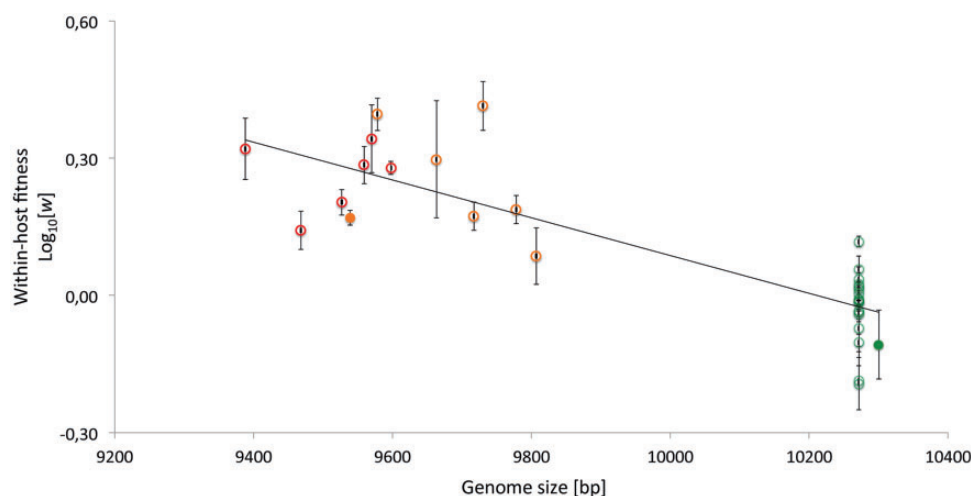
2011) show that the mutational effects on within-host fitness and virus accumulation at 7 dpi, expressed as the Malthusian growth rate per day (Lalić et al. 2011), are not correlated (Spearman correlation:  $\rho = 0.034$ , 19 df,  $P = 0.884$ ; see also [supplementary fig. S1, Supplementary Material](#) online). Therefore, accumulation was not expected to increase as a pleiotropic effect of increases in within-host fitness either. Accumulation will be an important parameter when each virus lineage is evolved in multiple host organisms, as higher accumulation can then lead to a higher frequency of a virus variant in the final population (e.g., Zwart, Van der Werf, et al. 2010).

For a plant virus, it is plausible that within-host fitness and accumulation are largely decoupled. Local infection by cell-to-cell movement can be achieved by a small number of virions transported to an adjacent cell (Miyashita and Kishino 2010). Therefore, those virus variants that spread rapidly need not necessarily accumulate a high number of virions per cell. Exclusion is moreover thought to play an important role in infection (Dietrich and Maiss 2003, Folimonova 2012), potentially allowing viruses that spread quickly to reach a high frequency and yet have relatively low accumulation. Furthermore, it should be noted that even if two virus variants have the same level of accumulation late in infection (e.g., 7 dpi), viruses with a high within-host competitive fitness may reach higher levels of accumulation early in infection (e.g., 3 dpi). When a virus rapidly exits the inoculated leaf, this can lead to higher levels of infection before infection levels saturate (Lafforgue et al. 2012; Zwart et al. 2012). Rapid replication and movement might therefore be the mechanism by which within-host fitness is increased in the evolved lineages, especially if genome size was directly linked to replication. Nevertheless, the key trait to measure from an evolutionary perspective—because it is expected to be under selection in this experimental setup—is competitive within-host fitness.

We then considered the relationship between genome size and within-host fitness for the evolved lineages ([fig. 7](#)) and also found a highly significant relationship (Spearman correlation:  $\rho = -0.877$ , 28 df,  $P < 0.001$ ). The mean fitness of evolved lineages with genomic deletions was higher than that of the ancestral virus without the heterologous gene (TEV) for 10 out of 12 lineages. However, when we performed pairwise comparisons between TEV and the evolved strains that fixed deletions, no significant differences were found ( $t$ -test with Holm–Bonferroni correction on the log-transformed  $W$  values). Statistical power is low when comparing individual lineages because individual-plant level variation is high, a limitation of our experimental system. We must therefore conclude that fitness of the wild-type TEV and evolved strains is similar. This result is, however, congruent with the observation that the convergent single-nucleotide mutations observed in the evolved TEV-eGFP lineages were not observed in wild-type TEV in this study and others (Bedhomme et al. 2012; Tromas N, Zwart MP, Elena SF, unpublished manuscript); it supports the suggestion that these mutations are specific for accommodating changes in the TEV-eGFP background and will probably not be beneficial in the wild-type virus background.

We then considered the within-host competitive fitness of those lineages without genomic deletions ([fig. 6C](#)). We found that 3-week lineages without genomic deletions had a significantly higher fitness than the 1-week lineages ([table 3](#)). The 3-week lineages appear to be at the sweet spot where the heterologous gene is maintained in many viral lineages, while there are concomitantly significant increases in viral fitness. This observation suggests that demography can play an important role in modulating the evolutionary outcome of HGT.

Although we have shown simultaneous maintenance of the heterologous gene and increases in fitness, the increases in



**Fig. 7.** The relationship between genome size (abscissa) and within-host competitive fitness (ordinate) is given. Green data points indicate no deletions in the heterologous gene (eGFP) were detected, orange indicates part or all of the eGFP was not present, and red indicates part of eGFP and of the viral HC-Pro cistron are not present. The data points for the ancestral TEV-eGFP and TEV have been filled, and TEV-eGFP has been shifted to the right so that it can be easily identified, even though its genome size is the same as other lineages without deletions. A linear regression line has been added only to emphasize the trend in the data. Note that most of the evolved viruses with deletions have a higher fitness than TEV, implicating the observed substitutions with increased fitness.

fitness are probably not related to retention of the eGFP sequence, because it is most unlikely to evolve any beneficial function to the virus. On the other hand, convergent single-nucleotide mutations occurring in TEV-eGFP did not occur in the wild-type virus, suggesting that these mutations are linked to the insertion of eGFP or the additional polyprotein cleavage site. The convergent evolution seen is therefore probably related to the heterologous sequence. Nevertheless, the next step is to perform evolution experiments with viruses carrying functional sequences. Although previous reports show heterologous sequences were generally unstable (e.g., Chapman et al. 1992; Dolja et al. 1993; Guo et al. 1998; Chung et al. 2007), here we have identified a demographic condition under which such sequences can be maintained and under which within-host fitness is still likely to be under selection. Indeed, preliminary results show that the cucumber mosaic virus 2b silencing suppressor is stably maintained in TEV genome when 3-week passages are made (Willemsen A, Zwart MP, Elena SF, unpublished data). On the other hand, the use of a heterologous sequence not expected to acquire any function has the advantage of focusing entirely on the process of pseudogenization and gene loss.

### Alternative Models That Explain the Data

Fixation of genomic deletions and increases in within-host fitness were not observed in the 1-week lineages. Moreover, we did not observe the convergent single-nucleotide mutations in the single 1-week lineage fully sequenced, and in another report, convergent evolution was not found after 1-week passages in *N. tabacum* (Bedhomme et al. 2012). Why is the outcome of short (1 week) and longer (3- and 9-week) passages so different? We discuss two conceptual models that may explain the data.

The first model focuses on the genetic bottlenecks during passaging and asserts that there may be too much drift for selection to operate effectively during the short-duration passages. Beneficial variation that arises de novo can be maintained in the virus population if its frequency increases to levels where it is likely to be sampled in the next round of infection. Therefore, the time between bottleneck events at the start of infection (i.e., passage duration) might be critical to the outcome of the evolutionary process. For the short 1-week passages, the passage duration is too short to allow any beneficial variation to increase to frequencies where it is likely to be sampled. For intermediate 3-week passages, this model necessarily postulates that the rate at which beneficial single-nucleotide mutations occur is higher than the rate at which recombination leading to the deletion of eGFP occurs, contrary to our expectations. Hence eGFP is lost in few intermediate-duration lineages, while beneficial mutations are fixed in many lineages. This leads to an overall increase in fitness, also for those lineages without deletions. Finally, for the long 9-week lineages, the passage duration is sufficiently long for both beneficial single-nucleotide mutations and recombinations to be selected to frequencies where they are maintained, and in most lineages, variants incorporating both predominate. This model predicts a sweet spot: functional

integration of a transgene will be most likely during intermediate-duration passages, where the heterologous sequence is maintained and passages are long enough for selection to act effectively on de novo variation. During longer passages, the heterologous sequence will simply be deleted. During shorter passages, the heterologous sequence will be maintained because selection cannot get rid of it. On the other hand, it will be doomed to remain nonfunctional because of the absence of effective selection.

The second model asserts that demographic conditions determine whether the heterologous sequence is detrimental or not. Early in infection, virus titers are low while the virus rapidly colonizes most plant tissues (Dolja et al. 1992; Zwart et al. 2012). For longer infections, virus titers will consistently be higher, and the virus continually expands into newly generated apical tissues. As a consequence, the cellular multiplicity of infection (MOI) is also likely to increase over the course of infection (González-Jara et al. 2009; Gutiérrez et al. 2010; Zwart et al. 2013). Different MOIs may lead to different selection pressures acting on the virus population, as exemplified by DIP viruses (e.g., Zwart et al. 2008). This second model predicts conditions under which the heterologous sequence can be maintained in the virus population. Integration of the heterologous sequence will depend on whether and how demography affects selection for improvement or new functions that can be recoded for by the heterologous sequence. However, in the absence of demographic effects on selection for transgene function, a sweet spot may still exist given that the strength of genetic drift will increase with shorter passages. Indeed, this may explain the lack of any convergent evolution in 1-week passages in tobacco (Bedhomme et al. 2012).

We think the first model is better supported by our data, because direct competition experiments suggest the fitness of wild-type TEV is higher than that of the TEV-GFP in 1-week passage (*t*-test on the log-transformed *W* values:  $t = 3.616$ , 4 df,  $P = 0.022$ ; see also fig. 6C). Hence, the occurrence of selection for smaller genome size appears to be independent of demographic conditions, although its intensity may not. Under both models, the proposed occurrence of a sweet spot for the evolutionary integration of heterologous genes will probably be pathosystem specific and, moreover, it may be sensitive to the exact virus genotype and environment used. For example, the stability of a heterologous gene and evolvability of the virus may both change with replication levels. Moreover, in some viruses in which inserts are highly unstable (Chung et al. 2007), there may always be a high number of lineages in which a heterologous gene is deleted, regardless of demography.

### Concluding Remarks

We found a clear relationship between genome size and within-host fitness for TEV, although other traits such as virus accumulation and virulence appear to be unaffected. This means that within-host fitness is sometimes under strong selection, apparently depending on the duration of infection between consecutive transmission events. This was exemplified by considering leaf-to-leaf evolution during

a 9-week infection period, where in some cases, deletions had already been fixed early in infection (e.g., leaf 15 of 45), and in most cases once a deletion occurred, it was maintained and fixed. We then asked whether there are conditions under which viral lineages evolved higher within-host fitness (suggesting that deterministic forces predominate over random forces in evolution) and the heterologous gene is concurrently maintained in most lineages. We found that this is the case for the 3-week passage condition. The within-host fitness of these lineages that had not lost the heterologous gene significantly increased, while 8 out of 10 lineages did not lose the heterologous gene after 27 weeks of evolution.

We therefore conclude that demographic conditions, in this case the duration of the infection period, can modulate the evolutionary process and possibly create conditions which one would expect to be more conducive for HGT. This suggested demographic sweet spot may help to resolve the paradox that viruses can have both streamlined genomes (Lynch 2006; Belshaw et al. 2007) and high levels of HGT (Dolja and Koonin 2011; Koonin and Dolja 2012). Moreover, our results suggest that demography may be critical for evolutionary innovation. However, observing the real-time integration of functional elements in the viral genome by experimental evolution would provide stronger support for these ideas.

These results also have implications for synthetic biology, in particular, the generation and employment of viral expression constructs. First, the significant difference in within-host competitive fitness between the 1- and 3-week lineages without genomic deletions suggests that an evolutionary approach could be used to optimize the fitness of expression vectors. All these lineages have intact eGFP expression (fig. 3), while the within-host fitness of the 3-week lineages is higher (table 3). In this respect, a key question that we have not addressed here is whether the evolved lineages with high fitness also have a higher stability of the eGFP insert. Second, the results suggest that intermediate-duration infections (i.e., approximately 3 weeks) can be suitable for expression of heterologous genes. This result contrasts with those of Dolja et al. (1993), who reported high instability of the  $\beta$ -glucuronidase marker in TEV. Both markers are approximately the same size. This difference might be explained by the expression strategies: eGFP was cleaved from both P1 and HC-Pro (Zwart et al. 2011) whereas  $\beta$ -glucuronidase was fused to HC-Pro (Dolja et al. 1992), possibly affecting HC-Pro functions and hereby lowering viral fitness. On the other hand, for many of the evolved TEV-eGFP variants, the N-terminal remains of eGFP were fused to HC-Pro (fig. 5A). Another possibility is therefore that expression of  $\beta$ -glucuronidase is harmful for viral replication.

Previous work has demonstrated that passage duration is of crucial importance to evolutionary outcomes for a plant virus (Dolja et al. 1993). Our results confirm the importance of passage duration for the evolution of genomic deletions, but we have also shown high levels of convergent evolution for single-nucleotide substitutions. In some lineages there were no unique mutations—although between-lineages variance

was high—making our results qualitatively comparable to other cases of convergent evolution in viruses (Bull et al. 1997). Our results therefore underscore that passage duration is crucial to the outcome of plant virus evolution and that adaptive and convergent evolution can both be observed in real time during long passages.

## Materials and Methods

### Plants, Virus Stocks, and Infections

*Nicotiana tabacum* was kept in a greenhouse at 24 °C with 16 h light. For the within-host competitive fitness assays and in vivo cloning experiment, plants were transferred to a growth chamber at 24 °C with 16 h light after inoculation. To generate a virus stock of the ancestral TEV-eGFP, we first transcribed RNA from the pMTEV-eGFP plasmid (Zwart et al. 2011) as described elsewhere (Carrasco, Daròs, et al. 2007). The third true leaf of 4-week-old *N. tabacum* plants was inoculated with 5  $\mu$ g of RNA. All systemically infected tissues were harvested 9 dpi, and virions were purified and stored as described elsewhere (Carrasco, Daròs, et al. 2007; Zwart et al. 2011).

Serial passage experiments were initiated by inoculating plants with 10  $\mu$ l of the purified virion stock by rubinoculating the third true leaf using carborundum. For passages 2 and onward, approximately 500 mg of homogenized infected tissue from the previous passage was diluted in 500  $\mu$ l phosphate buffer (50 mM potassium phosphate pH 7.0, 3% polyethylene glycol 6000). Fifty microliters were then rubinoculated to the third true leaf. For in vivo cloning, plants were inoculated with a 1:1,000 dilution in phosphate buffer of 1:1 mixture of infected tissue and phosphate buffer. eGFP fluorescence was observed with a Leica MZ16F stereomicroscope, using a 0.5 $\times$  objective and GFP2 filters (Leica).

### Reverse Transcription Polymerase Chain Reaction

To determine whether deletions had occurred at the eGFP locus, RNA was extracted from infected tissue using the RNeasy Plant kit (Qiagen). RT was performed using M-MuLV (Fermentas) and random hexamers. PCR was then performed with Taq DNA polymerase (Roche) and primers flanking the eGFP gene: forward primer 5'-CAATTG TTCGCAAGTGTGC-3', reverse primer 5'-ATGGTATGAAGAA TGCCTC-3'. One percent agarose gels were used to resolve PCR products. Note that this PCR assay was shown to have a high sensitivity for variants missing the eGFP gene in previous work (Zwart et al. 2011).

### Illumina NGS, SNP Calling, and Mapping of Genomic Deletions

For the sequencing of the experimentally evolved lineages containing deletions at the eGFP site, the virus genome was RT-PCR amplified using Accuscript RT (Agilent Technologies) and Phusion DNA polymerase (Thermo Scientific), with seven independent replicates that were pooled. The virus genome was amplified using three primer sets (set 1: 5'-GCAATCAAG CATTCTACTTC-3' and 5'-ATCCAACAGCACCTCTCAC-3'; set



2: 5'-TTGACGCTGAGCGGAGTGATGG-3' and 5'-AATGCTTC CAGAATATGCC-3'; set 3: 5'-TCATTACAAACAAGCACT TG-3' and 5'-CGCACTACATAGGAGAATTAG-3'), and equimolar mixtures of PCR products were made. Sequencing was performed at GenoScreen ([www.genoscreen.com](http://www.genoscreen.com)). Illumina HiSeq2000 2×100 bp paired-end libraries with multiplex adaptors were prepared along with an internal PhiX control. Sequencing quality control was performed by GenoScreen, based on PhiX error rate and Q30 values. Artifact filtering and read quality trimming (3' minimum Q20 and minimum read-length of 50 bp) was done using FASTX-Toolkit v0.0.13.2 ([hannonlab.cshl.edu/fastx\\_toolkit/index.html](http://hannonlab.cshl.edu/fastx_toolkit/index.html)). Dereplication of the reads and 5' quality trimming requiring a minimum of Q20 was done using PRINSEQ-lite v0.20.3 (Schmieder and Edwards 2011). Reads containing undefined nucleotides (N) were discarded. Mapping was done against the reference genome TEV-eGFP (GenBank KC918545) with Bowtie 2 v2.1.0 (Langmead and Salzberg 2012), which allows for gapped-read alignments. For every evolved lineage, SNPs were identified using SAMtools' mpileup (Li et al. 2009).

After the premapping step, the most common deletions observed were defined manually, and for every lineage, a new reference sequence was constructed masking each position of the defined deletion with the symbol N. These new reference genomes, together with the cleaned reads, were used as input for the program GapFiller v1.9 (Boetzer and Pirovano 2012), which reliably closes gaps within preassembled scaffolds using paired reads. GapFiller fills the gap from each edge in an iterative manner. In our case, it partially closed the gaps, base by base, until it could not extend any further given the difference between the a priori estimated deletion size and the actual size encountered. At both sides, overlapping sequences were manually identified and the ends were joined to reconstruct the new consensus sequences. These were error corrected using the software package Polisher v2.0.8 (available for academic use from the Joint Genome Institute). Accession numbers for the new consensus sequences are GenBank KC918546–KC918555 for 9-week lineages 1–10, GenBank KC918556 for 3-week lineage 3, GenBank KC918557 for 3-week lineage 6, and GenBank KC918558 for 1-week lineage 7.

The consensus sequences for the evolved TEV-eGFP lineages were reconstructed, and the cleaned reads were remapped against the corresponding consensus for every lineage (Materials and Methods). This remapping was efficient with about 93%–98% of the reads without a mate and about 96%–98% of the paired reads mapping exactly one time for the 9-week lineages. For the 3-week lineages, this was 91%–96% and 94%–97%, respectively. For the 1-week lineage, this was 85% and 87%, respectively, with 15% and 13% of the reads aligning zero times, for which most of these reads probably belong to the eGFP region of full-length TEV-eGFP variant present in the population. To detect other populations with a different deletion in the evolved lineages, we extracted all the reads that did not map exactly end-to-end, 100 bp around the deletion site, and mapped these reads against the original reference genome TEV-eGFP. In some regions, we found other deletions, but

these deletions occurred at very low frequencies ranging from 0.04% to 1.42%. The fact these frequencies are so low, and moreover that some of these deletions start at the same position, suggests that at least some of them could arise due to sequencing error or low-level contamination.

For each new consensus, SNPs were reidentified for every lineage using SAMtools' mpileup. Coverage by 35-bp windows distributions was generated for each new consensus. Statistically low-covered regions were searched for approximating the distribution to a normal and calculating a *P* value per window for a two-sided normal test. However, no regions with significantly lower coverage were identified, confirming that if other deletion variants exist, they are present at low frequencies.

### Evolution and Sequencing of Wild-Type TEV

To obtain wild-type TEV, plants were agroinoculated with pGTEVa (Bedoya et al. 2012). Three 9-week passages were performed as described for TEV-eGFP, with ten independent lineages. RNA was extracted from ground plant tissue of the evolved lineages and the 46–2,466 and 6,376–8,779 regions of the TEV genome (GenBank DQ986288) were RT-PCR amplified. The PCR products were sequenced with the 5'-GCAATCAAGCATTCTACTTC-3' and 5'-CCTGAT ATGTTTCCTGATAAC-3' primers, and the 5'-TCATTACAA CAAGCACTTG-3' and 5'-AGGCCCAACTCTCCGAAAG-3' primers, respectively.

### Virus Accumulation and Virulence Assay and Within-Host Competitive Fitness Assays

Four-week-old *N. tabacum* plants were infected with purified virions of the wild-type virus without the heterologous gene (TEV; derived from the pMTEV plasmid [Bedoya and Daròs 2010]), the ancestral virus for the evolution experiments (TEV-eGFP), and TEV marked with the mCherry reporter gene (TEV-mCherry, derived from the pMTEV-mCherry plasmid [Zwart et al. 2011]). Infected tissues were harvested after 1 week. Virus accumulation for these stocks of TEV, TEV-eGFP, and all evolved lineages were then determined. The Invitrap Spin Plant RNA Mini Kit (Stratag Molecular) was used to isolate total RNA. One-step RT-qPCR was then performed using the Primescript RT-PCR Kit II (Takara), in accordance with manufacturer instructions, and a PRISM Sequence Analyser 7500 (Applied Biosystems). Specific primers for the CP were used: 5'-TTGGTCTTGATGGC AACGTG-3' and reverse primer 5'-TGTGCCGTTTCAGTGTCTT CCT-3'. The 7500 Software version 2.0.4 (Applied Biosystems) was used to analyze the data. The concentration of genome equivalents could then be normalized to that of the sample with the lowest concentration, using phosphate buffer. Ten 4-week-old *N. tabacum* plants were then inoculated with 50 µl of these dilutions. Five randomly selected plants were harvested after 1 week, and virus accumulation was again determined using RT-qPCR, and these values were used as the viral accumulation of TEV, TEV-eGFP, and the evolved lineages. The height of the remaining five plants was measured 3 weeks postinoculation. To measure competitive

within-host fitness, TEV, TEV-eGFP, and all evolved lineages were again normalized to the sample of the lowest concentration, and 1:1 mixture of genome equivalents was made with TEV-mCherry. This virus has a similar insert size and within-host fitness compared with TEV-eGFP (Zwart et al. 2011) and was used as a common competitor for all virus strains. The 1:1 mixture of genome equivalents was rubinoculated in *N. tabacum* plants and infected tissues were harvested after 1 week. The mixture of TEV-eGFP and TEV-mCherry gave a 1:1 ratio of foci of primary infection (Zwart et al. 2011), validating the procedure to quantify genome equivalents and make mixtures. RT-qPCR for the CP was used to determine viral accumulation, while independent one-step RT-qPCR reactions were also performed for the mCherry sequence, using specific primers: 5'-CGGCGAGTTC ATCTACAAGG-3' and 5'-TGGTCTTCTTCTGCATTACGG-3'. The RT-qPCR method was otherwise identical to that used for the CP. The ratio of the evolved lineage to TEV-mCherry (R) is then:  $R = (n_{CP} - n_{mCherry}) / n_{mCherry}$ , where  $n_{CP}$  and  $n_{mCherry}$  are the RT-qPCR-measured copy numbers of the CP and mCherry, respectively. From the ratio at the start of the experiment ( $R_0 = 1$ ), we can then estimate the replicative advantage (W, see [Carrasco, Daròs, et al. 2007]) as:  $W = (R_t / R_0)^{1/t}$ , where  $t$  is the time in days and  $R_t$  the virus ratio at the end of the experiment. We consider the replicative advantage as a measure of within-host competitive fitness.

## Supplementary Material

Supplementary figure S1 is available at *Molecular Biology and Evolution* online (<http://www.mbe.oxfordjournals.org/>).

## Acknowledgments

The authors thank Alejandro Manzano Marín for his bioinformatics guidance with the Illumina analysis and Francisca de la Iglesia, Paula Agudo, and Àngels Pròsper for technical support. This project was made possible through the support of grant 22371 from the John Templeton Foundation to S.F.E. The opinions expressed in this publication are those of the authors and do not necessarily reflect the views of John Templeton Foundation. Additional support was received from the Spanish Dirección General de Investigación Científica y Técnica grants BFU2012-30805 to S.F.E, JCI2011-10379 to M.P.Z, and BIO2011-26741 to J.A.D., and by a Rubicon grant from the Netherlands Organization for Scientific Research ([www.nwo.nl](http://www.nwo.nl)) to M.P.Z.

## References

- Atreya CD, Atreya PL, Thornbury DW, Pirone TP. 1992. Site-directed mutations in the potyvirus HC-Pro gene affect helper component activity, virus accumulation, and symptom expression in infected tobacco plants. *Virology* 191:106–111.
- Bedhomme S, Lafforgue G, Elena SF. 2012. Multihost experimental evolution of a plant RNA virus reveals local adaptation and host-specific mutations. *Mol Biol Evol*. 29:1481–1492.
- Bedoya LC, Daròs JA. 2010. Stability of *Tobacco etch virus* infectious clones in plasmid vectors. *Virus Res*. 149:234–240.
- Bedoya LC, Martínez F, Orzáez D, Daròs JA. 2012. Visual tracking of plant virus infection and movement using a reporter MYB transcription factor that activates anthocyanin biosynthesis. *Plant Physiol*. 158:1130–1138.
- Belshaw R, Pereira V, Katzourakis A, Talbot G, Paces J, Burt A, Tristem M. 2004. Long-term reinfection of the human genome by endogenous retroviruses. *Proc Natl Acad Sci U S A*. 101:4894–4899.
- Belshaw R, Pybus OG, Rambaut A. 2007. The evolution of genome compression and genomic novelty in RNA viruses. *Genome Res*. 17:1496–1504.
- Boetzer M, Pirovano W. 2012. Toward almost closed genomes with GapFiller. *Genome Biol*. 13:R56.
- Bull JJ, Badgett MR, Springman R, Molineux IJ. 2004. Genome properties and the limits of adaptation in bacteriophages. *Evolution* 58:692–701.
- Bull JJ, Badgett MR, Wichman HA, Huelsenbeck JP, Hillis DM, Gulati A, Ho C, Molineux IJ. 1997. Exceptional convergent evolution in a virus. *Genetics* 147:1497–1507.
- Canchaya C, Proux C, Fournous G, Bruttin A, Brussow H. 2003. Prophage genomics. *Microbiol Mol Biol Rev*. 67:238–276.
- Carrasco P, Daròs JA, Agudelo-Romero P, Elena SF. 2007. A real-time RT-PCR assay for quantifying the fitness of *Tobacco etch virus* in competition experiments. *J Virol Methods*. 139:181–188.
- Carrasco P, de la Iglesia F, Elena SF. 2007. Distribution of fitness and virulence effects caused by single-nucleotide substitutions in *Tobacco etch virus*. *J Virol*. 81:12979–12984.
- Chapman S, Kavanagh T, Baulcombe D. 1992. *Potato virus X* as a vector for gene expression in plants. *Plant J*. 2:549–557.
- Chung BN, Canto T, Palukaitis P. 2007. Stability of recombinant plant viruses containing genes of unrelated plant viruses. *J Gen Virol*. 88:1347–1355.
- Chung BYW, Miller WA, Atkins JF, Firth AE. 2008. An overlapping essential gene in the *Potyviridae*. *Proc Natl Acad Sci U S A*. 105:5897–5902.
- Cronin S, Verchot J, Haldemanachill R, Schaad MC, Carrington JC. 1995. Long-distance movement factor—a transport function of the potyvirus helper component proteinase. *Plant Cell* 7:549–559.
- Diemer GS, Stedman KM. 2012. A novel virus genome discovered in an extreme environment suggests recombination between unrelated groups of RNA and DNA viruses. *Biol Direct*. 7:13.
- Dietrich C, Maiss E. 2003. Fluorescent labelling reveals spatial separation of potyvirus populations in mixed infected *Nicotiana benthamiana* plants. *J Gen Virol*. 84:2871–2876.
- Dolja VV, Herndon KL, Pirone TP, Carrington JC. 1993. Spontaneous mutagenesis of a plant potyvirus genome after insertion of a foreign gene. *J Virol*. 67:5968–5975.
- Dolja VV, Koonin EV. 2011. Common origins and host-dependent diversity of plant and animal viromes. *Curr Opin Virol*. 1:322–331.
- Dolja VV, McBride HJ, Carrington JC. 1992. Tagging of plant potyvirus replication and movement by insertion of  $\beta$ -glucuronidase into the viral polyprotein. *Proc Natl Acad Sci U S A*. 89:10208–10212.
- Folimonova SY. 2012. Superinfection exclusion is an active virus-controlled function that requires a specific viral protein. *J Virol*. 86:5554–5561.
- González-Jara P, Fraile A, Canto T, García-Arenal F. 2009. The multiplicity of infection of a plant virus varies during colonization of its eukaryotic host. *J Virol*. 83:7487–7494.
- Guo HS, López-Moya JJ, García JA. 1998. Susceptibility to recombination rearrangement of a chimeric *Plum pox potyvirus* genome after insertion of a foreign gene. *Virus Res*. 57:183–195.
- Gutiérrez S, Michalakakis Y, Blanc S. 2012. Virus population bottlenecks during within-host progression and host-to-host transmission. *Curr Opin Virol*. 2:546–555.
- Gutiérrez S, Yvon M, Thébaud G, Monsion B, Michalakakis Y, Blanc S. 2010. Dynamics of the multiplicity of cellular infection in a plant virus. *PLoS Pathog*. 6:e1001113.
- Hall JS, French R, Hein GL, Morris TJ, Stenger DC. 2001. Three distinct mechanisms facilitate genetic isolation of sympatric *Wheat streak mosaic virus* lineages. *Virology* 282:230–236.
- Huang AS. 1973. Defective interfering viruses. *Annu Rev Microbiol*. 27:101–117.

- Hughes AL, Friedman R. 2005. Poxvirus genome evolution by gene gain and loss. *Mol Phylogeny Evol.* 35:186–195.
- Koonin EV, Dolja VV. 2012. Expanding networks of RNA virus evolution. *BMC Biol.* 10:54.
- Kuo CH, Ochman H. 2009. Deletional bias across the three domains of life. *Genome Biol Evol.* 1:142–152.
- Lafforgue G, Tromas N, Elena SF, Zwart MP. 2012. Dynamics of the establishment of systemic potyvirus infection: independent yet cumulative action of primary infection sites. *J Virol.* 86:12912–12922.
- Lalić J, Cuevas JM, Elena SF. 2011. Effect of host species on the distribution of mutational fitness effects for an RNA virus. *PLoS Genet.* 7:e1002378.
- Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods.* 9:357–359.
- Lee SC, Kim DY, Hyun BH, Bae YS. 2002. Novel design architecture for genetic stability of recombinant poliovirus: the manipulation of G/C contents and their distribution patterns increases the genetic stability of inserts in a poliovirus-based RPS-Vax vector system. *J Virol.* 76:1649–1662.
- Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R, 1000 Genome Project Data Processing Subgroup. 2009. The sequence alignment/map format and SAMtools. *Bioinformatics* 25:2078–2079.
- Liu H, Fu Y, Li B, Yu X, Xie J, Cheng J, Ghabrial SA, Li G, Yi X, Jiang D. 2011. Widespread horizontal gene transfer from circular single-stranded DNA viruses to eukaryotic genomes. *BMC Evol Biol.* 11:91.
- Liu H, Fu Y, Xie J, Cheng J, Ghabrial SA, Li G, Peng Y, Yi X, Jiang D. 2012. Evolutionary genomics of mycovirus-related dsRNA viruses reveals cross-family horizontal gene transfer and evolution of diverse viral lineages. *BMC Evol Biol.* 12:276.
- Lynch M. 2006. Streamlining and simplification of microbial genome architecture. *Annu Rev Microbiol.* 60:327–349.
- Lynch M. 2007. The origins of genome architecture. Sunderland (MA): Sinauer Associates.
- Majer E, Daròs JA, Zwart MP. 2013. Stability and fitness impact of the visually discernible Rosea1 marker in the *Tobacco etch virus* genome. *Viruses* 5:2153–2168.
- Marks H, van Duijse JJA, Zuidema D, van Hulten MCW, Vlak JM. 2005. Fitness and virulence of an ancestral white spot syndrome virus isolate from shrimp. *Virus Res.* 110:9–20.
- Miyashita S, Kishino H. 2010. Estimation of the size of genetic bottlenecks in cell-to-cell movement of *Soil-borne wheat mosaic virus* and the possible role of the bottlenecks in speeding up selection of variations in trans-acting genes or elements. *J Virol.* 84:1828–1837.
- Ochman H, Davalos LM. 2006. The nature and dynamics of bacterial genomes. *Science* 311:1730–1733.
- Paar M, Schwab S, Rosenfellner D, Salmons B, Günzburg WH, Renner M, Portsmouth D. 2007. Effects of viral strain, transgene position, and target cell type on replication kinetics, genomic stability, and transgene expression of replication-competent *Murine leukemia virus*-based vectors. *J Virol.* 81:6973–6983.
- Pathak KB, Nagy PD. 2009. Defective interfering RNAs: Foes of viruses and friend of virologists. *Viruses* 1:895–919.
- Pijlman GP, van den Born E, Martens DE, Vlak JM. 2001. *Autographa californica* baculoviruses with large genomic deletions are rapidly generated in infected insect cells. *Virology* 283:132–138.
- Riechmann JL, Lain S, Garcia JA. 1992. Highlights and prospects of potyvirus molecular biology. *J Gen Virol.* 73:1–16.
- Routh A, Domitrovic T, Johnson JE. 2012. Host RNAs, including transposons, are encapsidated by eukaryotic single-stranded RNA virus. *Proc Natl Acad Sci U S A.* 109:1907–1912.
- Sacristán S, Malpica JM, Fraile A, García-Arenal F. 2003. Estimation of population bottlenecks during systemic movement of *Tobacco mosaic virus* in tobacco plants. *J Virol.* 77: 9906–9911.
- Sakai Y, Kiyotani K, Fukumura M, Asakawa M, Kato A, Shioda T, Yoshida T, Tanaka A, Hasegawa M, Nagai Y. 1999. Accommodation of foreign genes into the *Sendai virus* genome: sizes of inserted genes and viral replication. *FEBS Lett.* 456:221–226.
- Schmieder R, Edwards R. 2011. Quality control and preprocessing of metagenomic datasets. *Bioinformatics* 27:863–864.
- Seed KD, Lazinski DW, Calderwood SB, Camilli A. 2013. A bacteriophage encodes its own CRISPR/Cas adaptive response to evade host innate immunity. *Nature* 494:489–491.
- Simon AE, Roossinck MJ, Havelda Z. 2004. Plant virus satellite and defective interfering RNAs: new paradigms for a new century. *Annu Rev Phytopathol.* 42:415–437.
- Thornbury DW, Patterson CA, Dessens JT, Pirone TP. 1990. Comparative sequence of the helper component (HC) region of *Potato virus Y* and a HC-defective strain, *Potato virus C*. *Virology* 178:573–578.
- Yutin N, Koonin EV. 2012. Hidden evolutionary complexity of nucleocytoplasmic large DNA viruses of eukaryotes. *Virol J.* 9:161.
- Zwart MP, Daròs JA, Elena SF. 2011. One is enough: *in vivo* effective population size is dose-dependent for a plant RNA virus. *PLoS Pathog.* 7:e1002122.
- Zwart MP, Daròs JA, Elena SF. 2012. Effects of *Potyvirus* effective population size in inoculated leaves on viral accumulation and the onset of symptoms. *J Virol.* 86:9737–9747.
- Zwart MP, Dieu BTM, Hemerik L, Vlak JM. 2010. Evolutionary trajectory of *White spot syndrome virus* (WSSV) genome shrinkage during spread in Asia. *PLoS One* 5:e13400.
- Zwart MP, Erro E, Van Oers MM, De Visser JAGM, Vlak JM. 2008. Low multiplicity of infection *in vivo* results in purifying selection against baculovirus deletion mutants. *J Gen Virol.* 89:1220–1224.
- Zwart MP, Tromas N, Elena SF. 2013. Model-selection-based approach for calculating cellular multiplicity of infection during virus colonization of multi-cellular hosts. *PLoS One* 5:e64657.
- Zwart MP, Van der Werf W, Georgievskaya L, Van Oers MM, Vlak JM, Cory JS. 2010. Mixed-genotype infections of *Trichoplusia ni* larvae with *Autographa californica* multicapsid nucleopolyhedrovirus: speed of action and persistence of a recombinant in serial passage. *Biol Control.* 52:77–83.