

# Gene copy number variations at the within-host population level modulate gene expression in a multipartite virus

Romain Gallet,<sup>1,2</sup> J  r  my Di Mattia,<sup>1</sup> S  bastien Ravel,<sup>1</sup> Jean-Louis Zeddam,<sup>1</sup> Renaud Vitalis,<sup>2</sup> Yannis Michalakakis,<sup>3,  ,  </sup> and St  phane Blanc<sup>1,  ,  </sup>

<sup>1</sup>PHIM, Univ Montpellier, INRAE, CIRAD, IRD, Institut Agro, Montpellier, France, <sup>2</sup>CBGP, Univ Montpellier, INRAE, CIRAD, IRD, Institut Agro, Montpellier, France and

<sup>3</sup>MIVEGEC, Univ Montpellier, CNRS, IRD, Montpellier, France

  Equally contributed to this work.

  <https://orcid.org/0000-0003-1929-0848>

  <https://orcid.org/0000-0002-3412-0989>

\*Corresponding author: E-mail: [stephane.blanc@inrae.fr](mailto:stephane.blanc@inrae.fr)

## Abstract

Multipartite viruses have a segmented genome, with each segment encapsidated separately. In all multipartite virus species for which the question has been addressed, the distinct segments reproducibly accumulate at a specific and host-dependent relative frequency, defined as the ‘genome formula’. Here, we test the hypothesis that the multipartite genome organization facilitates the regulation of gene expression via changes of the genome formula and thus via gene copy number variations. In a first experiment, the faba bean necrotic stunt virus (FBNSV), whose genome is composed of eight DNA segments each encoding a single gene, was inoculated into faba bean or alfalfa host plants, and the relative concentrations of the DNA segments and their corresponding messenger RNAs (mRNAs) were monitored. In each of the two host species, our analysis consistently showed that the genome formula variations modulate gene expression, the concentration of each genome segment linearly and positively correlating to that of its cognate mRNA but not of the others. In a second experiment, twenty parallel FBNSV lines were transferred from faba bean to alfalfa plants. Upon host switching, the transcription rate of some genome segments changes, but the genome formula is modified in a way that compensates for these changes and maintains a similar ratio between the various viral mRNAs. Interestingly, a deep-sequencing analysis of these twenty FBNSV lineages demonstrated that the host-related genome formula shift operates independently of DNA-segment sequence mutation. Together, our results indicate that nanoviruses are plastic genetic systems, able to transiently adjust gene expression at the population level in changing environments, by modulating the copy number but not the sequence of each of their genes.

**Key words:** virus; multipartite; gene expression; gene copy number; copy number variation; genome formula.

## Introduction

Multipartite viruses are intriguing genetic systems whose biology is difficult to explain within the current conceptual framework of virology (Iranzo and Manrubia 2012; Luc  a-Sanz and Manrubia 2017; Michalakakis and Blanc 2020). Their genomes consist of two or more nucleic acid segments, the most striking feature being their individual encapsidation into distinct virus particles. Such a split genome packaging has an obvious cost, which is the increased risk of segment loss at each transmission event, and thus of failed inoculation. Even though the means by which the myriads of multipartite viral species efficiently manage this cost are poorly uncovered (Gilmer, Ratti, and Michel 2018; Sicard et al. 2019), it is unanimously acknowledged that they all face the same issue of maintenance of the genome integrity (Iranzo and Manrubia 2012; Luc  a-Sanz, Aguirre, and Manrubia 2018; Michalakakis and Blanc 2020). Counterpart benefits of multipartite genome architecture, in contrast, are highly debated and as yet

no proposition reached a consensus (Sicard et al. 2016; Luc  a-Sanz, Aguirre, and Manrubia 2018; Zwart et al. 2021). Most hypotheses recognize the smaller size of the genome segments as conferring an advantage to the system, through faster replication (Nee 1987), mutation escape (Pressing and Reanney 1984), genetic exchange via segment reassortments (Chao 1991), or particle stability (Ojosnegros et al. 2011). These proposals have several drawbacks: (1) most do not explain the separate encapsidation of the distinct segments; (2) they are not specific to multipartite genome architecture and similarly apply to viruses encapsidating all segments together; and (3) none constrains the relative frequency of the segments, which should thus evolve toward the situation of minimum cost where all segments accumulate at equal copy number. Available studies estimating the relative amount of distinct genome segments in hosts infected by multipartite viruses together indicate that this situation of minimum cost is never reached (French and Ahlquist 1988; Sicard et al. 2013;

Hu et al. 2016; Wu et al. 2017; Mansourpour et al. 2022). Although other explanations are imaginable, one possibility is that the actual benefits in these viral systems are related to the differential accumulation of the distinct segments. We and others accordingly proposed that multipartite viruses can tune gene expression in fluctuating environments by modifying gene (or segment) copy number (Sicard et al. 2013; Michalakakis and Blanc 2020) and that this capacity could be adaptive (Zwart and Elena 2020).

Gene copy number (GCN) and copy number variations (CNV), defined as variations of the number of copies of one or several genes across individuals, have a strong impact on gene expression and phenotypes in all organisms (Todd and Selmecki 2020). That CNV-induced changes in gene expression contribute to adaptation in fluctuating environments, particularly upon colonization of new niches and host switching in host-pathogen interactions, have been consistently demonstrated through experimental evolution of fungi (Todd and Selmecki 2020), bacteria (Elliott, Cuff, and Neidle 2013), and viruses (Bayer, Brennan, and Geballe 2018). Some general features related to CNV are highlighted by the corresponding literature. First, there is a rampant generation of copy number polymorphism, sometimes occurring at a rate higher than the mutation rate, that selection can act on within a population (Kondrashov 2012; Belikova et al. 2020). In all cases, CNV polymorphism is generated by recombination and DNA repair machineries, preferentially acting on specific features of sequences flanking the amplified regions. Accordingly, depending on these flanking sequences, some genome regions are more prone to amplification than others (Elliott, Cuff, and Neidle 2013; Bayer, Brennan, and Geballe 2018; Sasani et al. 2018; Belikova et al. 2020). Second, depending on the regulatory network within a genome, even small-scale amplification can have large effects. The amplification of one region can either increase or decrease the expression of specific genes, located within or outside this region (Mileyko, Joh, and Weitz 2008; Kondrashov 2012), with drastic non-linear changes and even bifurcation in the behavior of the network when amplification thresholds are reached (Mileyko, Joh, and Weitz 2008). Experimental observations of such non-linear effects of gene amplification have been reported and are discussed in Mileyko, Joh, and Weitz (2008). Third, CNV repeatedly proved immediately adaptive due to a simple gene dosage effect (Kondrashov 2012; Elliott, Cuff, and Neidle 2013; Bayer, Brennan, and Geballe 2018); i.e. following environmental changes, the expression of a gene may be deregulated (maladapted) and a simple adjustment of its copy number (gene dosing) alleviates the defect and the corresponding GCN variant is selected for (Elde et al. 2012; Belikova et al. 2020; Todd and Selmecki 2020). Fourth, such a coarse mechanism of amplification-mediated gene expression tuning (AMGET) (Tomanek et al. 2020) is based on gene expression heterogeneity within the population and can evolve rapidly at a pace where transcriptional regulation has no chance to emerge/adapt *de novo* (Elde et al. 2012; Kondrashov 2012; Brennan et al. 2014; Sasani et al. 2018; Todd and Selmecki 2020; Tomanek et al. 2020). Fifth, gene amplification is costly (Kondrashov 2012; Elliott, Cuff, and Neidle 2013; Bayer, Brennan, and Geballe 2018). Consequently, whenever selection pressure is removed, when the organism is back in the benign nonrestrictive environment or when a regulatory mutation occurs, the extra copies of the gene are rapidly deleted. This phenomenon of genome expansion and contraction, named 'genomic accordion', has been empirically observed in fungi, bacteria, and viruses (Elde et al. 2012; Belikova et al. 2020; Todd and Selmecki 2020). Sixth, and finally, the fleeting

nature of genomic accordion often leaves no sequence signature, and therefore its importance in the adaptation and evolution of pathogenic fungi, bacteria, and viruses has likely been underestimated (Elliott, Cuff, and Neidle 2013; Bayer, Brennan, and Geballe 2018; Tomanek et al. 2020).

Related studies on viruses have focused on monopartite large double-stranded DNA (dsDNA) viruses such as phage T4 (Wu and Black 1995), baculoviruses (Ardissou-Araújo et al. 2018), herpesviruses (Arias et al. 2014), poxviruses (Elde et al. 2012), or even giant viruses (Shukla, Chatterjee, and Kondabagil 2018) because they can accommodate relatively large genome size variations. It is commonly assumed that physical packaging constraints do not allow such genome size variations for other viruses, as for example RNA or single-stranded DNA (ssDNA) viruses, and thus preclude any possibility of genomic accordion-like adaptation processes (Bayer, Brennan, and Geballe 2018). In this context, it is astounding that multipartite viruses have not been envisaged as potentially specialized genetic systems for AMGET, and this is the hypothesis we experimentally addressed in this study. Because each segment is separately packaged in its own virus particle, the genome-length constraint on segment copy number is totally absent in a multipartite architecture, opening the way to GCN-regulated gene expression and GCN-driven adaptation for ssDNA/RNA viruses where it is usually deemed impossible. Similarly, because genome segments are by definition physically separated, each could be amplified independently with no requirement for recombination/nucleic acid repair machineries. Precisely because no sequence rearrangement appears necessary, the system may not require any sequence modification at all and therefore be extremely conducive to CNV. Would this be so, the importance of CNV in the way of life of multipartite viruses may have long been overlooked because of the total absence of genomic (sequence) signatures, even the transient ones involved in CNV and genomic accordion in other organisms.

Using the faba bean necrotic stunt virus (FBNSV, family *Nanoviridae*), where each of the eight genome segments encodes a single gene and where each segment is encapsidated individually, we provide support for the use of segment AMGET as an everyday lifestyle. We earlier reported that the FBNSV segments each accumulate in specific amounts, reproducibly yielding a host-dependent frequency pattern designated as the genome formula (Sicard et al. 2013). We then speculated that the copy number of each of the segments could contribute to the control of gene expression, and others theoretically supported the idea that a genome formula producing a gene expression pattern better adapted to a given environment can be selected extremely rapidly (Zwart and Elena 2020). Here, we empirically show that the copy number of each DNA segment correlates positively and linearly to the concentration of its encoded messenger RNA (mRNA), but rarely (if at all) to that of the other viral mRNAs, indicating that GCN variations drive gene expression in FBNSV. This effect of the genome formula on gene expression is consistently verified in two distinct host species. Further and unanticipatedly, our results reveal that the genome formula modifications observed upon host switching compensate for distinct rates of mRNA production and maintain a relatively constant stoichiometry in the viral transcriptome. Finally, high-throughput sequencing of twenty parallel viral lines demonstrates that the host-dependent FBNSV genome formula shift is not associated with positive/negative selection of sequence variants but rather illustrates a *bona fide* mutation-free CNV.

## Materials and methods

### Viral strain and plant infection procedures

In all experiments, we used the FBNSV isolate JKI-2000 provided by the Gronenborn lab and described in Grigoros et al. (2009). Faba bean (*Vicia faba*, cv 'Seville') plants were agroinoculated with cultures of *Agrobacterium tumefaciens* COR308 strain, each carrying a pBin19 plasmid containing a tandem repeat of one of the eight FBNSV segments. All eight *A. tumefaciens* cultures were mixed together at equal proportions and inoculated into plants as described in Sicard et al. (2013). For practical reasons, alfalfa (*Medicago truncatula*) plants were infected via aphid transmission as described earlier (Sicard et al. 2013). We have previously shown that the same genome formula is reached whether plants are infected through agroinoculation or aphid transmission and whatever the initial frequency of inoculated segments (Sicard et al. 2013).

### Experiment 1

#### Overview: concomitant quantification of viral DNAs and mRNAs

The genome formula was characterized as the median relative frequency of each segment across several plant replicates. The idea of this experiment was to embrace the across-replicate variation in the relative frequency of each segment in order to see whether it is translated into across-replicate variation in the relative frequency of the corresponding mRNAs. We thus estimated the relative concentrations of both viral DNA segments and viral mRNAs in each plant sample analyzed. For plant viruses, the viral gene expression is stopped at some point of the infection in fully infected tissues (Gutiérrez et al. 2015). Consequently, to ensure capturing the transient expression of mRNAs, we repeated this experiment at two different time points. The first replicate (Trial A) was performed on sixteen faba beans and twenty-eight alfalfa plants. Samples were collected at different dates, on the first day where the individual infected plants showed symptoms, i.e. 10–15 days post-infection (dpi) for faba bean and 13–18 dpi for alfalfa. The second replicate (Trial B) was performed on twenty-one faba beans and twenty alfalfa plants, and the samples were collected at one single later date for each plant species, once all plants of the species expressed symptoms, i.e. 21 days post-infection (dpi) for faba bean and 20 dpi for alfalfa. In both trials, the infection of each of these plants with FBNSV was independent.

#### Extraction of ssDNA and mRNA from each plant sample

On each infected plant, an apical leaflet was sampled. Approximately 100 mg of leaf tissue was placed in a microtube containing two sterile glass beads and frozen immediately in liquid nitrogen. Samples were homogenized mechanically using a mixer mill MM 301 (four cycles of 20 s at 30 Hz).

To extract nucleic acids, 900 µl of GHLC buffer (6.5 M guanidinium hydrochloride, 100 mM Tris-HCl pH 8.0, 100 mM sodium acetate pH 5.5, 0.1 M β-mercaptoethanol) were added to the homogenized samples. Tubes were vortexed and then centrifuged at 10,000 *g* for 10 min at 4°C in a 5415R Eppendorf (Hamburg, Germany) centrifuge. Nine hundred microliters of TRI Reagent (Sigma-Aldrich) warmed at 65°C were added to supernatants. Tubes were vortexed gently over three cycles of 30 s, and 200 µl of chloroform was added. After vortexing, samples were incubated for 10 min at room temperature and centrifuged at 12,000 *g* for 15 min at 4°C. For each sample, 1,200 µl of aqueous upper phase containing nucleic acids was retrieved and divided into two tubes each containing 600 µl. These 600 µl were mixed with 560 µl of cold

(–20°C) isopropanol and centrifuged at 12,000 *g* for 20 min at 4°C. Supernatants were discarded and pellets were washed with 70 per cent ethanol at 4°C. Finally, nucleic acids were resuspended in 50 µl RNase-free water, and the two tubes for each sample pooled back together to obtain 100 µl of nucleic acid solution. This nucleic acid extraction step allowed retrieving both ssDNA and mRNA from the plant tissues.

#### Estimation of the genome and transcriptome formulas by quantitative PCR

Quantitative PCRs (qPCRs) were directly performed on these nucleic acid samples in order to infer the FBNSV genome (DNA) formula. Inferring the viral transcriptome (RNA) formula by qPCR was more tedious as it first required complete removal of the viral DNA and then reverse transcription of the mRNAs into cDNA.

Total elimination of the viral DNA could be achieved by using two treatments, a DNase I digestion followed by the purification of mRNAs. The DNase digestion was conducted by mixing 16 µl of the nucleic acid sample with 2 µl of 10× DNase buffer (400 mM Tris-HCl pH 8.0, 100 mM MgSO<sub>4</sub>, and 10 mM CaCl<sub>2</sub>) and 2 µl of DNase I (Promega). This mix was incubated at 37°C for 30 min. A subsequent 10 min incubation at 65°C inactivated the DNase. After digestion, the Dynabeads mRNA purification kit (Ambion—ThermoFisher) was used on the nucleic acid samples following the manufacturer's recommendations. Control qPCRs were performed on these samples and confirmed the complete degradation/elimination of viral DNA (Supplementary Fig. S1).

For production of the cDNAs, 10 µl of mRNA samples was mixed with 1 µl dNTP (10 mM) incubated for 5 min at 65°C and later placed on ice. A mix composed of 4 µl of 5× buffer (250 mM Tris-HCl pH 8.3, 375 mM KCl, and 15 mM MgCl<sub>2</sub>), 2 µl of DTT (100 mM), and forty units of RNasin ribonuclease inhibitor (Promega) was added to the mRNA sample and incubated for 2 min at 42°C. Two hundred units of SuperScript™ II Reverse Transcriptase (RT) (Invitrogen) were added to the mix, followed by an additional 50 min incubation at 42°C. The RT was inactivated with a final 15 min incubation at 70°C. The newly formed cDNAs were diluted 10 times so that the buffer does not affect the following qPCR reactions.

All qPCR reactions (40 cycles of 95°C for 10 s, 60°C for 10 s, and 72°C for 10 s) were carried out using a LightCycler 480 thermocycler (Roche) and the LightCycler FastStart DNA Master Plus SYBR green I kit (Roche), following the manufacturer's instructions. The nucleic acid sample (1.2 µl of a 10-fold dilution of either total nucleic acid extracts or cDNA preparation) was added to the qPCR mix (5 µl of Roche 2× qPCR mastermix, 3.5 µl of H<sub>2</sub>O, 0.3 µl of primer mix, and 8.8 µl total) after distribution in 384-well microtiter plates. Primers (Sicard et al. 2013) were used at a final concentration of 0.3 µM for amplifications of the C, M, and S segments and 0.5 µM for amplifications of the N, R, U1, U2, and U4 segments.

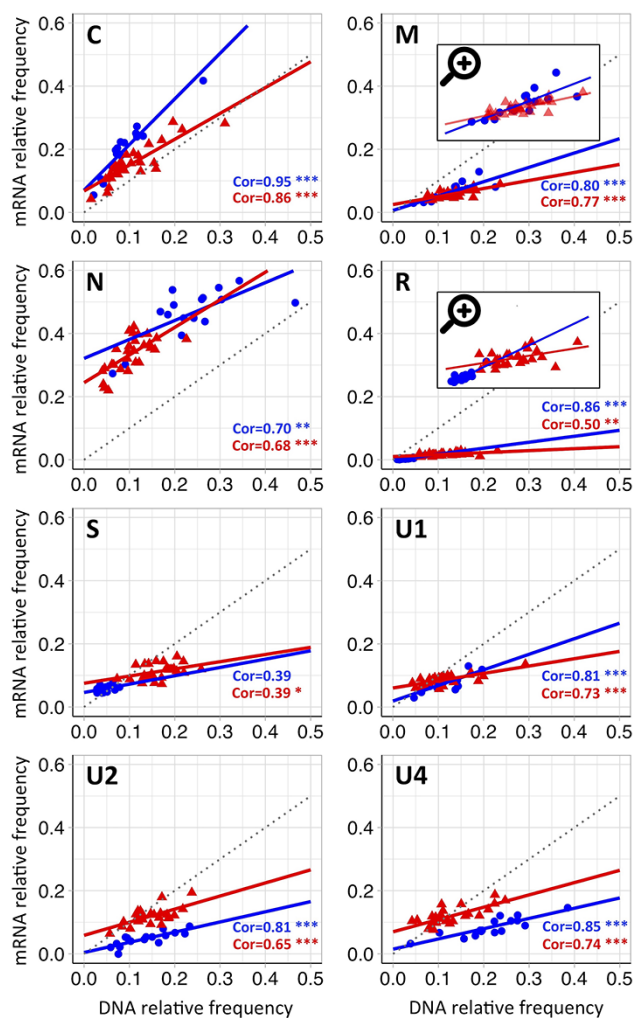
Serial dilutions of plitmus28 plasmids each carrying one of the eight FBNSV segments (Grigoros et al. 2009) were placed on each qPCR plate (eight serial dilutions per PCR plate in total, one for each FBNSV segment). These were used as an internal control in order to draw a standard curve for each segment and for each qPCR plate, alleviating any potential bias related to between-qPCR plate variations. Fluorescence data were first analyzed with the LinRegPCR program (Ruijter et al. 2009) and later converted into ng of DNA by using the standard curves. Both DNA and RNA formulas could then be inferred by computing the relative proportions



of each segment or mRNA as described in [Sicard et al. \(2013\)](#). All qPCR reactions were duplicated (two wells on the same PCR plate).

### Statistical analyses

To investigate the relationship between gene expression and the concentration of DNA, we first calculated the Pearson correlations between the frequency of each segment and the frequency of its corresponding mRNA in each host plant and trial. Because of the large number of correlations and tests, we applied the Benjamini–Hochberg False Discovery Rate (FDR) correction to the correlations across all segments for each host plant species and trial. These results are reported in [Fig. 1](#) and [Supplementary Table S1](#). We then calculated the Pearson correlations between the frequency of each segment and that of the seven non-cognate mRNA, but in this case we did not apply FDR corrections, as further commented in the Results section. We used the R software version 3.1.3 ([R Core Team 2021](#)) to calculate all these correlations.



**Figure 1.** Correlations between the relative frequency of FBNSV genome segments and that of their respective mRNAs. The data used here are those from Trial A, Experiment 1 (see Materials and Methods). Each panel shows the correlation between the relative frequency of an FBNSV segment and the relative frequency of the corresponding mRNA. Data points, linear regressions, correlation coefficients, and P-values are shown in blue and red for FBNSV infecting faba bean and alfalfa, respectively. ‘\*\*\*\*’, ‘\*\*\*’, and ‘\*’ correspond to P-value  $\leq 0.001$ , 0.01, and 0.05, respectively. The dotted line illustrates a slope of 1.

For further characterizing the relationship between the frequency of each genome segment and that of its cognate mRNA, we compared linear and quadratic model fits to the data and applied model selection using the Akaike Information Criterion (AIC) as described in [Burnham and Anderson \(2002\)](#) and [Bolker \(2008\)](#).

To study how the different factors (segment, host plant, and DNA formula) interact, we modeled the concentration of mRNA of each segment as a function of the segment, the host plant, and the concentration of DNA of the segment. Performing such analyses on relative frequencies of mRNA and DNA would provide intuitively interpretable results on intuitively normalized quantities: the frequencies. Unfortunately, such an analysis would be flawed by the fact that the frequencies of the DNA and mRNA of different segments within each replicate are not independent since they sum to one; because of this, the regression coefficients linking them would also be correlated since their mean should also equal one. We thus opted for the following approach: (1) to investigate the interaction between the DNA formula and the host plant species we performed separate analyses on each segment, modeling the logit frequency of the segment’s mRNA as a function of the host plant species, the logit of the frequency of the segment’s DNA and their interaction; these analyses are reported in [Supplementary Table S2](#); (2) we run a full model on the concentrations, and not the frequencies, of the DNA and mRNA of each segment because the concentrations are not parametrically constrained. To comply with the analysis of variance assumptions, these concentrations were first transformed using the Johnson Sb transformation. The transformed values were analyzed in a mixed linear model whose dependent variable was the concentration of mRNA, and the explanatory variables were ‘replicate’, declared as a random factor, and ‘segment’, ‘host plant’, and ‘concentration of DNA’ declared as fixed variables (and all the multiple interactions among the latter three declared as fixed variables). This analysis is reported in [Supplementary Table S3](#). The analyses mentioned in this paragraph were performed using JMP 13.2.1 (SAS Institute 2016).

A distance between DNA formulas and between RNA formulas was calculated to compare the situation in faba bean and in alfalfa host plants. This distance between formulas was calculated as follows:

$$d = \sum |f_i^{faba} - f_i^{alfalfa}|$$

where  $f_i$  is the relative frequency of the  $i$ th segment (or mRNA) in the formula.

All distances between DNA formulas in faba bean and in alfalfa and between RNA formulas in faba bean and in alfalfa were calculated ( $16 \times 28 = 448$  distances between DNA formulas and 448 distances between RNA formulas in Trial A;  $21 \times 20 = 420$  distances between DNA formulas and 420 distances between RNA formulas in Trial B). As the formula of each plant was used several times to calculate all possible distances (e.g. the genomic formula of the faba bean plant 1 was used twenty-eight times to calculate all distances between this formula and all alfalfa formulas), not all distances in the dataset are independent. In order to take this pseudo-replication into account, we analyzed these distances with a mixed model with the factor ‘faba bean plant identity’ and ‘alfalfa plant identity’ as random factors and the ‘nucleic acid’ (DNA vs. mRNA) and time (Trial A vs. Trial B) as fixed effect factors. This statistical analysis was performed with JMP 13.2.1 (SAS Institute 2016).



## Experiment 2

### Overview: monitoring viral polymorphism in populations passing from faba bean to alfalfa

This experiment has been described in a previous technical paper estimating the various possible quantitative biases during amplification steps and ultra-deep-sequencing of these viral populations (Gallet et al. 2017). In the present study, the same experiment and thus the same deep-sequencing dataset are used to monitor polymorphism in twenty independent viral populations passing from faba bean to alfalfa host plants.

Briefly, fifteen aphids were placed on each of twenty FBNSV-infected faba bean plants, 3 weeks post-infection. Three days later, ten of these aphids were used from each plant to transmit the FBNSV to a set of twenty alfalfa plants, thus creating twenty independent viral populations. During this experiment, total DNA extraction was performed on systemically infected faba bean (21 days post-infection, just before aphids were placed on the plants) and alfalfa plants (26 days post-inoculation by the aphids). qPCR were first performed on all forty DNA extracts in order to measure the FBNSV genome formula in the two host species. Then, a rolling circle amplification (RCA, amplifying single-stranded circular FBNSV DNA segments) was performed in order to enrich the samples with viral DNA sequences, and the forty RCA products were sent for deep-sequencing (for full details see Gallet et al. 2017). The full sequence data set is available upon request.

### Candidate mutations for genome formula variation

To be considered a mutation impacting the FBNSV genome formula when the virus is passed from faba bean to alfalfa, the mutation should (1) show a significant increase in frequency between faba bean and alfalfa samples, beyond that expected under drift alone; (2) this change in frequency should be consistent across replicates; and (3) this increase should correlate with the variation in the genome formula. We describe below how mutations under selection have been searched for. The other two requirements, repeated occurrence in parallel viral populations and correlation with genome formula changes, are reported in the Results section.

In order to identify mutations whose frequency changed between faba bean and alfalfa samples beyond what is expected under drift alone, we tested for the homogeneity of temporal differentiation across nucleotide sites for each viral population passed from faba bean to alfalfa, using a procedure inspired by Goldringer and Bataillon (2004). The rationale of this analysis is that if all sites are selectively neutral, they should provide identically distributed estimates of temporal differentiation. However, if some sites are targeted by selection (or if they are in linkage disequilibrium with selected variants), then some heterogeneity in site-specific measures of temporal differentiation should be observed. To identify those sites that show outstanding differentiation compared to neutral expectation, we simulated the dynamics of nucleotide frequency change between the faba bean and the alfalfa samples, conditionally on the initial nucleotide counts in the faba bean sample and on the strength of genetic drift during the experiment.

To that end, we first estimated the haploid effective size of the viral population ( $N_e$ ) using approximate Bayesian computation (ABC) (see, e.g. Beaumont 2010). Because each segment is transmitted independently and since the genome formula may reflect different rates of genetic drift during transmission (Gallet et al. 2018), ABC analyses were performed (and therefore  $N_e$  estimates were computed) independently for each segment.

The data consisted of the observed number of A, T, C, and G counts obtained by deep-sequencing in all 40 FBNSV populations (twenty in faba bean and twenty in alfalfa). Yet, to lessen the impact of sequencing errors in deep-sequencing data, we discarded all the variants with an observed frequency of the most frequent allele (MAF, computed as the overall frequency across the faba bean and the alfalfa samples) falling above 0.97, thereby assuming a variant calling threshold of 0.03 (see, e.g. Sobel Leonard et al. 2019). We ended up with 269 polymorphic sites (out of 7,907 sites  $\times$  20 replicates = 158,140 sites), corresponding to 173 unique sites. For each segment-specific analysis, all polymorphic sites (with  $MAF \leq 0.97$ ) were pooled.

We then simulated individual nucleotide frequency trajectories as follows: suppose that we observe a vector  $\mathbf{y} \equiv (y_A, y_C, y_G, y_T)$  of nucleotide counts, out of the total coverage  $n_{fb} \equiv y_A + y_C + y_G + y_T$  in the faba bean sample. We assume that these observed counts correspond to a (biallelic) Single Nucleotide Polymorphism (SNP) with sequencing errors, and we denote by  $y_{fb}$  the counts for the major (most frequent) allele. We further assume (following Frachon et al. 2017) that  $y_{fb}$  is drawn from a binomial distribution  $B(n_{fb}, \pi_{fb})$ , where  $\pi_{fb}$  is the (unknown) allele frequency of the major allele in the faba bean population. Assuming a (uniform) Beta(1,1) prior distribution for  $\pi_{fb}$  and using the Bayes inversion formula, the posterior distribution of  $\pi_{fb}$  is distributed as  $Beta(y_{fb} + 1, n_{fb} - y_{fb} + 1)$ . For each nucleotide site and for each ABC simulation, we therefore draw the initial allele frequencies in the faba bean sample  $\tilde{\pi}_{fb}$ , at random from a  $Beta(y_{fb} + 1, n_{fb} - y_{fb} + 1)$  distribution. We then draw 'pseudo-observed' allele counts using a random binomial draw from  $B(n_{fb}, \tilde{\pi}_{fb})$ . This procedure allows accounting for the sampling variance in initial allele frequencies. Then, we simulate  $\tau$  generations of drift, using successive binomial draws with parameters  $N_e$  (the segment-specific effective population size) and the nucleotide frequencies in the previous generation. In the last generation, a sample of nucleotide counts is drawn from a binomial distribution with parameters  $n_M$  (the total observed coverage in the alfalfa sample) and  $\tilde{\pi}_M$  (the simulated nucleotide frequencies in the last generation). In what follows, we considered a single generation of drift (i.e.  $\tau = 1$ ). Finally, sequencing errors were modeled (for both the faba bean and the alfalfa samples) by means of multinomial draws, with probabilities  $(1 - \epsilon)$  not to mutate, and  $\epsilon/3$  to mutate to any other state. For each segment, a total of 1,000,000 ABC simulations were performed assuming a log-uniform prior for  $N_e$  in the [1; 1,000] range and a log-uniform prior for the error rate  $\epsilon$  in the [0.001; 0.1] interval. To avoid any bias, all simulations with a major allele frequency larger than or equal to 0.97 were discarded. The summary statistics considered to compare observed and simulated data were (1) the mean, variance, skewness, and kurtosis of single-locus estimates of  $F_{ST}$  (Weir and Cockerham 1984) computed between the faba bean and the alfalfa samples at each SNP (with a major allele frequency  $\leq 0.97$ ) within a segment and (2) the allele frequency difference of the MAF between the faba bean and the alfalfa samples at each SNP within a segment. Posterior distributions of  $N_e$  and  $\epsilon$  were computed using the abc package for R (Csilléry, François, and Blum 2012) with the local linear regression model (MA, Zhang, and Balding 2002) and a tolerance threshold of 0.001.

In a second step, for each segment and for each variant, we tested the null hypothesis that the locus-specific differentiation measured at this focal marker was only due to genetic drift. For this purpose, we computed the expected distribution of  $F_{ST}$  at each site, conditional upon the estimated effective population size for the segment, the inferred error rate, and the allele frequencies

at the focal site in the faba bean sample. To do so, we simulated individual nucleotide frequency trajectories following the same rationale as for the ABC simulations, drawing  $N_e$  and  $\epsilon$  estimates from their ABC posterior distributions. For each simulated trajectory, we computed site-specific estimates of temporal  $F_{ST}$  from the simulated nucleotide counts at the initial and last generation. The whole procedure was repeated at least 1,000,000 times for each of the 269 polymorphic sites. Finally, we assigned a  $P$ -value to each site, computed as the proportion of simulations giving a site-specific estimate of  $F_{ST}$  larger than or equal to the observed value at the focal nucleotide site. As above, all simulations with a major allele frequency larger than or equal to 0.97 were discarded.

All codes and R scripts, as well as the SNP counts data, specifically developed and used for these analyses are publicly accessible at (10.57745/ILFCP4).

## Results

### GCN drives gene expression in FBNSV

To investigate whether the FBNSV gene expression is affected by GCN, we assessed whether variation of the relative concentration of the viral mRNA produced by each segment across different individual plants of a given host species could be explained by variation of the genome formula across these same individual plants (Experiment 1 as described in Materials and Methods). In each plant sample analyzed, we thus estimated the relative concentrations of both viral DNA segments and their cognate mRNAs, which we hereafter, respectively, designate genome formula and transcriptome formula. It has been shown in various biological systems that the viral gene expression is stopped at some point of the infection (Gutiérrez et al. 2015). To maximize our chances to capture the transient expression of viral mRNAs, we thus repeated this experiment at two different time points: early in Trial A, as soon as infection symptoms were visible on each individual plant, and later in Trial B, at the same time post-infection for all individual plants once all had exhibited symptoms. Because an mRNA half-life can be short, we were aware that the two trials could differ in their capacity to potentially reveal a correlation between the genome and transcriptome formulas.

Figure 1 (Supplementary Table S1) shows that the relative frequency of each of the eight mRNAs of the FBNSV is positively correlated to that of its encoding segment in Trial A, both in faba bean and in alfalfa host plants (except for the S segment in faba bean for which the correlation is not significant). Trial B provided consistent observations, with six and four segments, respectively, on faba bean and alfalfa, showing significant positive relationships (Supplementary Fig. S2 and Table S1). The segment-by-segment analyses identified statistically significant effects of either the DNA formula or its interaction with the host plant species for all segments in Trial A and for six segments in Trial B (Supplementary Table S2), further indicating that a change in a segment frequency and thus of the genome formula induces a change of the gene expression.

The slopes of the linear regressions between mRNA and DNA relative frequencies vary with both the segments and the host species (Fig. 1). To assess the statistical significance of this slope variation across hosts, we analyzed the plant species effect on the DNA/mRNA correlation for each segment separately. A statistically significant effect was observed for segments C, R, and U1 in both trials and for segment M in Trial A (Supplementary Table S2), indicating that these segments are differentially expressed in the two host plant species. The slope variation across segments is further supported by the statistically significant segment-by-plant

interaction in the full model using mRNA and DNA concentrations (Supplementary Table S3).

### The relationship between GCN and gene expression is remarkably simple

Two observations indicate a simple relationship between genome formula and gene expression in this viral system. First, the relative abundance of any specific genome segment does not strikingly depart from a simple positive and linear relationship with that of its cognate mRNA. Second, most correlation tests between the frequency of any specific segment and that of each of the seven non-cognate viral mRNAs proved nonsignificant.

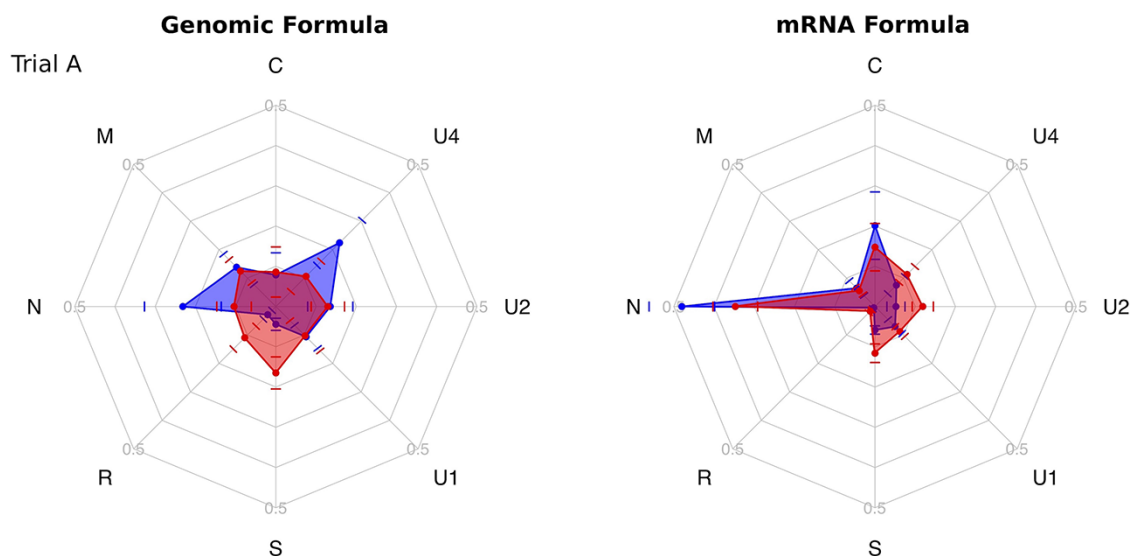
To substantiate the first observation, we verified whether incorporating quadratic terms in the regressions better explains the data than the regressions reported in Fig. 1 and Supplementary Fig. S2, which only contain terms linear in DNA concentration. Across all trials in faba bean and alfalfa, this proved very rarely true, i.e. for five regressions out of thirty-two (for full detail see Supplementary Table S4). Adding a quadratic term explained the data better solely for C, N, and R in faba bean Trial A (in the case of segments N and R, after removing the point with the highest DNA concentration—rightmost in Fig. 1—this was no longer true), for N in alfalfa Trial A, for none of the segments in faba bean Trial B, and for U4 in alfalfa Trial B (here also, after removing the rightmost point, the quadratic term is no longer statistically significant).

For the second observation, we calculated all possible Pearson's correlations between viral DNAs and mRNAs in faba bean and alfalfa and in Trials A and B (256 correlation tests; see Supplementary Table S5). Concentrations of genome segments nearly systematically correlated positively with those of their cognate mRNAs, as already presented in the previous section, but rarely with non-cognate mRNAs. More specifically, in faba bean Trial A, 87.5 per cent (7/8) cognate correlations were statistically significant vs. 14 per cent (8/56) non-cognate (one-tailed Fisher exact test  $P < 0.0001$ ). In faba bean Trial B, the corresponding numbers were 75 per cent (6/8) for cognate vs. 23 per cent (13/56) for non-cognate (one-tailed Fisher exact test  $P = 0.0055$ ), in alfalfa Trial A, 100 per cent (8/8) for cognate vs. 16 per cent (9/56) for non-cognate (one-tailed Fisher exact test  $P < 0.0001$ ), and in alfalfa Trial B, 50 per cent (4/8) for cognate vs. 5 per cent (3/56) for non-cognate (one-tailed Fisher exact test  $P = 0.0033$ ). In order to make our conclusions as conservative as possible, no corrections for multiple tests and related false discovery rates were performed in this analysis.

All together, these results suggest that changes in the frequency of a given FBNSV genome segment positively and linearly affect the expression of the corresponding gene, while poorly affecting the others.

### Different genome formulas in faba bean and alfalfa produce similar transcriptome formulas

We plotted and compared genome and transcriptome formulas when estimated from faba bean and from alfalfa plants (Fig. 2 for Trial A, Supplementary Fig. S3 for Trial B). As already observed in a previous study (Sicard et al. 2013), the FBNSV genome formulas in faba bean and alfalfa are clearly distinct. However, the transcriptome formulas observed in the two host species appear more similar. To confirm this observation, we compared the distance between genome formulas and between transcriptome formulas in these two hosts (see Materials and Methods). Our statistical analysis formally established that the distance between faba bean and alfalfa transcriptome formulas was significantly smaller than



**Figure 2.** Radar plot of FBNSV genome and transcriptome formulas in Trial A (Experiment 1). The median relative frequencies of each FBNSV segment (left) or of their corresponding transcripts (right) are represented on one of the eight axes composing the radar plot (formulas calculated from the sixteen faba bean and twenty-eight alfalfa plants in Trial A). The FBNSV formulas observed in faba bean and alfalfa are represented in blue and red, respectively. SDs are represented by colored bars. The distances between the transcriptome formulas observed in faba bean and alfalfa are significantly smaller than those between the corresponding genome formulas (Table 1).

**Table 1.** Statistical analysis of the distance between genome and transcriptome formulas in faba bean and alfalfa (Experiment 1, Trial A). The nature of the nucleic acid (NatNA), DNA or mRNA, was a fixed factor. We used the individual plant, faba bean or alfalfa, and its interaction with the nucleic acid as random factors to account for pseudo-replication. See Materials and Methods for more explanations.

| Model adjusted R <sup>2</sup> = 0.8351 |                    |                 |                    |           |              |            |
|--|--------------------|-----------------|--------------------|-----------|--------------|------------|
| Fixed effect tests                     |                    |                 |                    |           |              |            |
| Source                                 | Nparm <sup>a</sup> | DF <sup>b</sup> | DFDen <sup>c</sup> | F ratio   | P > F        |            |
| NatNA                                  | 1                  | 1               | 38.59              | 48.1650   | <0.0001      |            |
| REML variance component estimates      |                    |                 |                    |           |              |            |
| Random effect                          | Variance component | Std. Error      | 95% lower          | 95% upper | Wald P-value | % of total |
| Faba                                   | 0.0017             | 0.0014          | −0.0011            | 0.0044    | 0.2334       | 9.416      |
| Alfalfa                                | 0.0004             | 0.0015          | −0.0026            | 0.0034    | 0.7895       | 2.291      |
| NatNA * faba                           | 0.0033             | 0.0013          | 0.0008             | 0.0058    | 0.0095       | 18.657     |
| NatNA * alfalfa                        | 0.0072             | 0.0020          | 0.0032             | 0.0112    | 0.0004       | 40.424     |
| Residual                               | 0.0052             | 0.0006          | 0.0048             | 0.0057    |              | 29.211     |
| Total                                  | 0.0177             | 0.0021          | 0.0143             | 0.0226    |              | 100.000    |

<sup>a</sup>Number of parameter.

<sup>b</sup>Degrees of freedom.

<sup>c</sup>Denominator degrees of freedom.

the distance between faba bean and alfalfa genome formulas in both trials (Table 1 for Trial A and Supplementary Table S6 for Trial B). These results demonstrate that while the relative copy number of the genome segments changes drastically when FBNSV switches from faba bean to alfalfa, the relative proportions (or stoichiometry) of the eight mRNAs tend to be conserved. This interesting observation is further discussed later.

### Looking for adaptive mutations in the FBNSV sequence

To investigate whether the change in genome formula when FBNSV switches hosts is due, or not, to the selection of mutations in the sequence of one or several segments, we reanalyzed deep-sequencing data from twenty independent FBNSV populations passed from faba bean to alfalfa (Experiment 2 described

in Materials and Methods). Supplementary Fig. S4 shows that, just like in Experiment 1 and in earlier reports (Sicard et al. 2013; Gallet et al. 2017), the FBNSV genome formula was clearly different in faba bean and alfalfa, confirming that the expected host-dependent genome formula shift has occurred.

The modification of a phenotype during viral infection could either have a genetically determined basis or be due to a plastic response. To distinguish between these two possibilities, we aimed at identifying mutations showing outstanding differentiation between faba bean and alfalfa samples (as compared to what is expected under genetic drift alone) that could be interpreted as evidence of selection and whose frequency variation across host species could explain the genome formula variation.

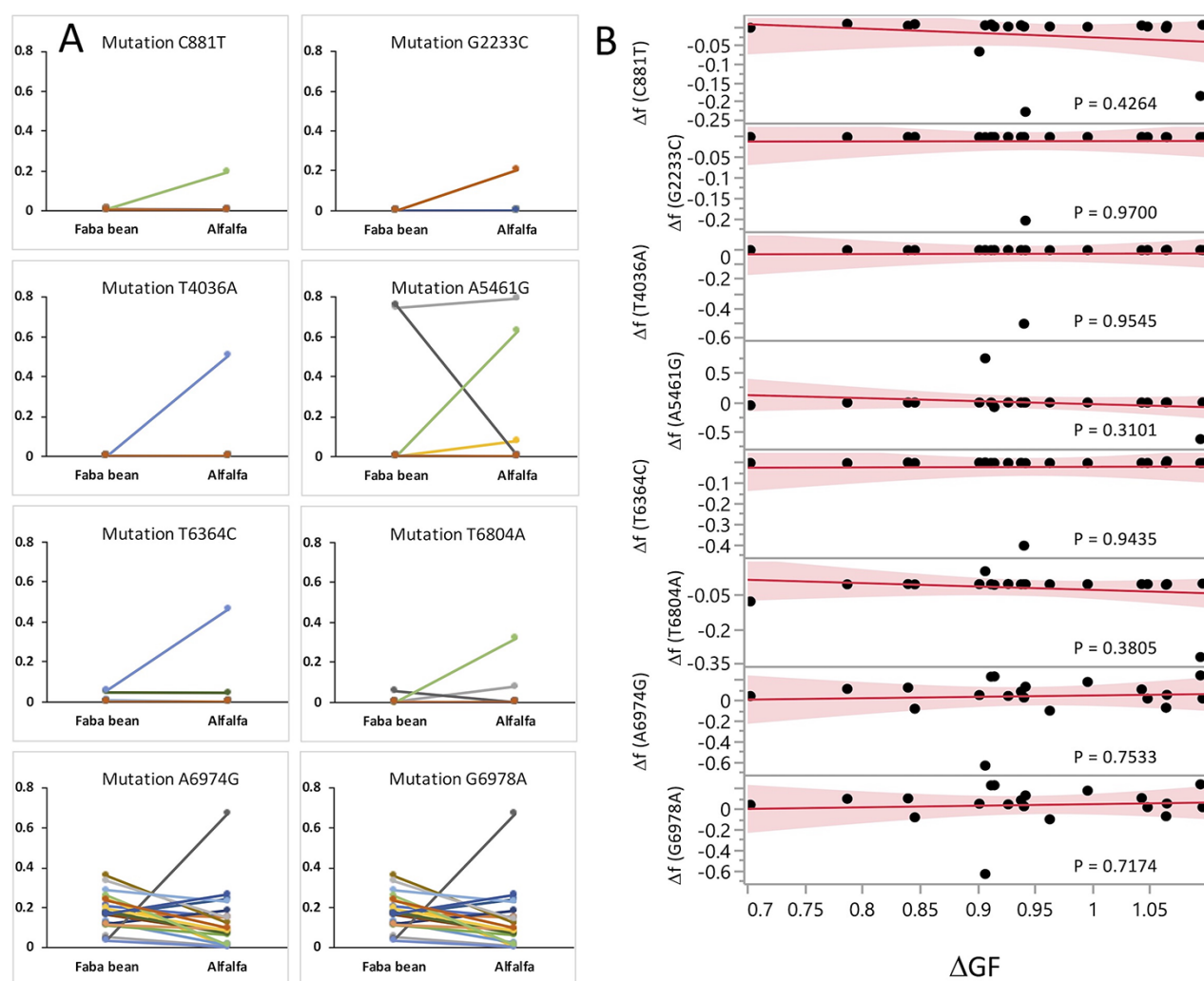
Over the 7,907 nucleotide positions in the concatenated FBNSV genome and the twenty replicated viral populations monitored in faba bean and alfalfa samples, we detected 269 variants



(i.e. with a major allele frequency  $\leq 0.97$ ), corresponding to 173 distinct sites. From the ABC analysis, we then inferred the effective population size for each segment, as well as the error rate (Supplementary Table S7), and used these estimates to simulate allele frequency dynamics in order to test whether the extent of differentiation observed in our viral lines passing from faba bean to alfalfa could be explained by drift only. It is noticeable that these estimates of effective population size for each segment, although using a totally distinct approach, are very similar to those reported earlier (Gallet et al. 2018). Interestingly, we found only eight sites at which the observed differentiation departed from the expected distribution under neutrality ( $P \leq 0.01$ ). Among these eight sites, two were revealed in two out of the twenty parallel viral lines and six were revealed only once. The position of these sites on the FBNSV genome segments, whether they are in coding regions, synonymous or non, is indicated in Supplementary Table S8. Figure 3A illustrates that the frequency of each of the

corresponding mutations can follow very diverse trajectories in the twenty parallel FBNSV lines, increasing, decreasing, or not changing at all, pleading against a deterministic process.

We finally tested whether the frequency variations (observed between faba bean and alfalfa samples) of each of these eight candidate mutations were correlated to changes of the genome formula. For this, we calculated the distance between the genome formula in faba bean and that in alfalfa for each of the twenty FBNSV lines and plotted these distances against the variation of the mutation frequency in each corresponding line (Fig. 3B). All regressions proved nonsignificant, further confirming that even the extremely rare sites identified as eventually showing higher differentiation than expected under drift alone cannot account for the genome formula shift of FBNSV. We thus conclude that this is a mutation-independent process, and whether it is to be considered a plastic or genetically-driven phenomenon is not trivial and is further discussed in the next section.



**Figure 3.** No adaptive mutations can explain the host-dependent genome formula shift of FBNSV. (A) Changes in the frequency of mutations at the eight sites detected as possibly under selection in at least one of the twenty parallel FBNSV lines passing from faba bean to alfalfa. Each viral population is represented with a specific color. All twenty populations are represented in all graphs; when less than twenty populations are visible, it is because several are superimposed. For all eight mutations, the frequency sometimes increases, decreases, or does not change, depending on the population considered. (B) Distances between the genome formula in faba bean and that in alfalfa ( $\Delta GF$ ) in the twenty parallel FBNSV lines plotted against the frequency of mutations ( $\Delta f$ ) at the eight sites detected as possibly under selection. The P-values of the regressions are indicated in each case.

## Discussion

After the discovery of the genome formula of nanoviruses (Sicard et al. 2013; Mansourpour et al. 2022), additional studies performed on other multipartite viruses showed that their genomic segments also accumulate at different frequencies (Hu et al. 2016), in a host-dependent manner (Wu et al. 2017). While this phenomenon appears general in multipartite (Gutiérrez and Zwart 2018) and perhaps even in segmented viruses (Diefenbacher, Sun, and Brooke 2018; Moreau et al. 2020), the mechanisms leading to the establishment of the genome formula as well as its actual function remain a mystery. We first hypothesized (Sicard et al. 2013, 2016) that the multipartite genome architecture would allow the adjustment of gene expression through the modulation of the GCN, and this proposition was further discussed (Wu et al. 2017; Diefenbacher, Sun, and Brooke 2018; Gutiérrez and Zwart 2018; Moreau et al. 2020) and even theoretically supported (Zwart and Elena 2020) by others. The proposed process (Sicard et al. 2013, 2016; Gutiérrez and Zwart 2018) is that infection sites could randomly differ in the proportions of the different segments and that within-host selection would act on this variation to favor the replication/dissemination of those sites with a genome formula producing the gene expression pattern better adapted to the specific host. According to this process, viral populations of a given virus genotype could rapidly converge to the genome formula that is best adapted to a given environment (Zwart and Elena 2020). What is appealing with this hypothesis on the mechanism that can generate the set-point genome formula is that it provides an astonishingly versatile means to regulate gene expression that perfectly matches, or even magnifies, the general conclusion enounced from studies on CNV in other organisms: gene amplification is based on rampant generation of copy number polymorphism and allows rapid and graded response for populations in heterogeneous and changing environments, which can tune gene expression when promoters are not adequately regulated at a pace where regulatory sequences have no time to evolve (Elliott, Cuff, and Neidle 2013; Bayer, Brennan, and Geballe 2018; Todd and Selmecki 2020; Tomanek et al. 2020). Baseline experimental support for such a function of FBNSV genome formula lies in three key points: (1) the demonstration of a correlation between GCN and gene expression—i.e. a correlation between the relative segment frequencies and those of their respective mRNA, (2) the ability to adjust the GCN when the environment changes, and (3) the demonstration that this is a mutation-free process, confirming that no *de novo* regulatory sequences have evolved. Thus, our results demonstrate the functional role of the genome formula and its variation, although its potential role in adaptation to host switches, or or to changes of the physiological state of its host, remains to be empirically demonstrated.

The first point is consistently verified in both faba bean and alfalfa. In the 'early' Trial A, the only segment that did not show a statistically significant correlation was segment S in the faba bean background (i.e. one nonsignificant correlation out of sixteen—Fig. 1 and Supplementary Table S1). This is probably due to the relative scarcity of this segment in the faba bean DNA formula, which leads to overall low relative frequencies of both S mRNA and DNA, and consequently low between-plant variation. In Trial B, statistically significant correlations could be observed in ten instances out of sixteen. As already commented in both the Materials and Methods and the Results, we anticipated a possible distinct turnover for DNA segments and for their cognate mRNAs, which might bias the assessment of their correlated accumulation at some stages of the infection. Despite this potential

drawback apparently affecting more Trial B, the DNA-segment frequencies proved to significantly impact the cognate mRNA production in most cases. Since the set-point genome formula has been reported (Sicard et al. 2013) or modeled (Zwart and Elena 2020) to be reached early during the onset of the infection and then remain constant, we assume that there is more experimental noise in Trial B, related to shorter-lived mRNA when compared to viral DNA.

The fact that different FBNSV segments had different levels of mRNA production (different slopes for different DNA/mRNA regressions; Fig. 1) may simply reflect a different efficiency of the segments' respective promoter, as earlier reported for the related banana bunchy top virus (BBTV, genus *Babuvirus*, family *Nanoviridae*) (Yu et al. 2019), or different mRNA half-lives. In FBNSV and in nanoviruses in general, although totally uncharacterized, each gene likely has its own promoter strength because sequences flanking the transcription start are not highly conserved across segments. We note, however, that whatever the regulatory sequences on distinct FBNSV segments, the effect of GCN variation reported here impacts gene expression patterns. More, interestingly, the mRNA production could also vary for a given segment between the faba bean and alfalfa backgrounds (Fig. 1), indicating that its promoter may not be equally compatible with the host plant species' respective transcriptional machinery or that the stability of mRNA may vary across hosts. Considering that the FBNSV genome formula is different and modulates the expression of the viral mRNAs in the two host plant species and that a given segment does not produce the same amount of mRNA in these two environments, one intuitively expects the relative proportion of the distinct viral mRNAs (transcriptome formula) to also greatly vary, at least reflecting differences observed at the DNA level and perhaps even more. Surprisingly, however, our results reveal that the transcriptome formula tends to be more conserved in the two hosts. We believe this observation supports the second key point listed in the first paragraph of the Discussion. The potential importance of the stability of gene expression patterns and how CNVs can maintain a dosage balance between interacting genes has earlier been discussed (Kondrashov 2012). Here, the interactions between FBNSV and the mRNA metabolism machineries in faba bean and in alfalfa should modify the viral mRNA frequency pattern (different slopes of segments DNA/mRNA regressions in the two hosts). Our results suggest that the modification of the genome formula in the two hosts allows the maintenance of a dosage balance between FBNSV genes, resulting in a similar transcriptome formula.

The third point lies in the demonstration that the host-related genome formula shift is a mutation-free phenomenon, which both greatly advances our understanding of this genetic system and adds one additional enigma. The advance is the discovery that the genome formula changes of a given viral isolate upon host switching are plastic and not traceable on a sequence basis. In diverse organisms, it is a classical view that gene amplification is often transient and followed by gene contraction (genomic accordion) ultimately leaving no genomic signature (Elde et al. 2012; Belikova et al. 2020; Todd and Selmecki 2020; Tomanek et al. 2020). Remarkably, for multipartite viruses (at least for FBNSV), not only gene expansion/contraction leaves no genomic signature, but it does not even require transient sequence modification. The other face of the coin is that this discovery adds one dimension to the puzzle: In monopartite viruses, bacteria, and eucaryotes, it is a genome that is modified with a portion amplified. The corresponding genotype then represents a genetic

innovation that is selected for or against. In multipartite viruses, when the genome formula changes, what exactly is the genetic innovation? A group of interacting segments may represent the genetic innovation. Indeed, as discussed above, we earlier proposed that such a group of interacting segments could be the unit of selection (Sicard et al. 2013, 2016). This possibility has been theoretically formalized and supported (Gutiérrez and Zwart 2018; Zwart and Elena 2020), but experimental evidence is still missing.

One additional observation that we found particularly intriguing is the type of relationship between FBNSV GCN and gene expression. As already documented in the Introduction, the impact of gene duplication/amplification on gene expression has been empirically reported in eucaryotes, procaryotes, and monopartite large dsDNA viruses where, to the best of our knowledge, a correlated increase of the expression of the corresponding gene has been evidenced but not characterized in detail. In their seminal theoretical paper, Mileyko, Joh, and Weitz (2008) considered only two to three interacting genes, all located collinearly in the same amplified region, thus all similarly amplified. Despite these simple virtual gene networks and as already explained in the Introduction, a remarkable diversity of possible gene expression changes was revealed. In our experimental system, with all eight FBNSV genes physically separated on distinct genome segments, and all differentially amplified, considering the complete lack of data on the interaction network between these genes, we had no ground for educated guesses, but we did not expect something as simple as what we observed. The copy number of each of the DNA segments has a positive and linear relationship with the production of its cognate mRNA with little impact on the expression of other viral genes. Again inspired by and in line with the same theoretical study (Mileyko, Joh, and Weitz 2008), we propose that the FBNSV may have evolved away from gene-amplification thresholds leading to drastic changes/bifurcation in the behavior of the expression network. Such a simple behavior of the gene expression network might be a condition for this virus to operate AMGET as an everyday lifestyle without jeopardizing the system at every possible change of the genome formula.

All the above considerations highlight the conceptual problem of what exactly constitutes the genome of multipartite viruses: should it be just the set of the sequences of their segments or should a definition of their genome also include the number of copies of each segment? The genome formula shift accompanying a host switch could be viewed as an indication of viral plasticity (as assumed in the previous paragraph). The situation is, however, unclear: the genome formula shift is in essence a modification of the copy number of specific genes. In this respect, it is not conceptually different from the 'genomic accordion' process. In the case of the monopartite viruses, bacteria, and eucaryotes, it is obvious to everyone that the determinism of the phenotypic adaptation to the environmental challenge via 'genomic accordions' is genetically based. The fact that the genes of multipartite viruses are carried by separate segments offers them great flexibility in the number of copies. This flexibility muddles our conceptual characterizations: because we have difficulty in defining what their genome really is, we have difficulty in deciding on the nature of their responses to environmental changes. Is it genetic, plastic, or epigenetic? At this point, we cannot but leave this question open, only the unraveling of the mechanisms underlying genome formula changes, the identification of the unit of selection in these systems, will provide the answer.

## Supplementary data

Supplementary data are available at Virus Evolution online.

## Acknowledgements

We are grateful to Sophie Leblaye for plant production. R.G., R.V., and S.B. acknowledge support from INRAE Plant Health and Environment Division, S.R. from CIRAD, J.L.Z. from IRD, and Y.M. from CNRS and IRD.

## Funding

This work was funded by French national institutions INRAE, CNRS, IRD, the French National Research Agency (ANR), and Montpellier University of Excellence (MUSE). This work was supported by the ANR grants 'Nano' (ANR-14-CE02-0014) and 'Nanovirus' (ANR-18-CE92-0028-01) and by MUSE (project BLANC-MUSE2020-Multivir).

**Conflict of interest:** The authors declare no competing interest.

## References

- Ardissou-Araújo, D.M.P., et al. (2018) 'A Novel Betabaculovirus Isolated from the Monocot Pest *Mocis latipes* (Lepidoptera: Noctuidae) and the Evolution of Multiple-Copy Genes', *Viruses*, 10: 134.
- Arias, C. et al. (2014) 'KSHV 2.0: A Comprehensive Annotation of the Kaposi's Sarcoma-Associated Herpesvirus Genome Using Next-Generation Sequencing Reveals Novel Genomic and Functional Features', *PLoS Pathogens*, 10: e1003847.
- Bayer, A., Brennan, G., and Geballe, A. P. (2018) 'Adaptation by Copy Number Variation in Monopartite Viruses', *Current Opinion in Virology*, 33: 7–12.
- Beaumont, M. A. (2010) 'Approximate Bayesian Computation in Evolution and Ecology', *Annual Review of Ecology, Evolution, and Systematics*, 41: 379–406.
- Beaumont, M. A., Zhang, W., and Balding, D. J. (2002) 'Approximate Bayesian Computation in Population Genetics', *Genetics*, 162: 2025–35.
- Belikova, D. et al. (2020) "'Gene Accordions" Cause Genotypic and Phenotypic Heterogeneity in Clonal Populations of *Staphylococcus Aureus*', *Nature Communications*, 11: 3526.
- Bolker, B. M. (2008) *Ecological Models and Data in R*. Princeton: Princeton University Press.
- Brennan, G. et al. (2014) 'Adaptive Gene Amplification as an Intermediate Step in the Expansion of Virus Host Range', *PLoS Pathogens*, 10: e1004002.
- Weir, B. S., and Cockerham, C. C. (1984) 'Estimating F-Statistics for the Analysis of Population Structure', *Evolution: International Journal of Organic Evolution*, 38: 1358–70.
- Burnham, K. P., and Anderson, D. R. (2002) *Model Selection and Multi-model Inference: A Practical Information-Theoretic Approach*, 2nd edn. Springer: New York.
- Chao, L. (1991) 'Levels of Selection, Evolution of Sex in RNA Viruses, and the Origin of Life', *Journal of Theoretical Biology*, 153: 229–46.
- Csilléry, K., François, O., and Blum, M. G. B. (2012) 'Abc: An R Package for Approximate Bayesian Computation (ABC)', *Methods in Ecology and Evolution*, 3: 475–9.
- Diefenbacher, M., Sun, J., and Brooke, C. B. (2018) 'The Parts Are Greater Than the Whole: The Role of Semi-infectious Particles in Influenza A Virus Biology', *Current Opinion in Virology*, 33: 42–6.



- Elde, N. C. et al. (2012) 'Poxviruses Deploy Genomic Accordions to Adapt Rapidly Against Host Antiviral Defenses', *Cell*, 150: 831–41.
- Elliott, K. T., Cuff, L. E., and Neidle, E. L. (2013) 'Copy Number Change: Evolving Views on Gene Amplification', *Future Microbiology*, 8: 887–99.
- Frachon, L. et al. (2017) 'Intermediate Degrees of Synergistic Pleiotropy Drive Adaptive Evolution in Ecological Time', *Nature Ecology & Evolution*, 1: 1551–61.
- French, R., and Ahlquist, P. (1988) 'Characterization and Engineering of Sequences Controlling in Vivo Synthesis of Brome Mosaic Virus Subgenomic RNA', *Journal of Virology*, 62: 2411–20.
- Gallet, R. et al. (2017) 'The Number of Target Molecules of the Amplification Step Limits Accuracy and Sensitivity in Ultradeep-Sequencing Viral Population Studies', *Journal of Virology*, 91.
- et al. (2018) 'Small Bottleneck Size in a Highly Multipartite Virus during a Complete Infection Cycle', *Journal of Virology*, 92.
- Gilmer, D., Ratti, C., and Michel, F. (2018) 'Long-Distance Movement of Helical Multipartite Phytoviruses: Keep Connected or Die?' *Current Opinion in Virology*, 33: 120–8.
- Goldringer, I., and Bataillon, T. (2004) 'On the Distribution of Temporal Variations in Allele Frequency: Consequences for the Estimation of Effective Population Size and the Detection of Loci Undergoing Selection', *Genetics*, 168: 563–8.
- Grigoras, I. et al. (2009) 'Reconstitution of Authentic Nanovirus from Multiple Cloned DNAs', *Journal of Virology*, 83: 10778–87.
- Gutiérrez, S. et al. (2015) 'The Multiplicity of Cellular Infection Changes Depending on the Route of Cell Infection in a Plant Virus', *Journal of Virology*, 89: 9665–75.
- Gutiérrez, S., and Zwart, M. P. (2018) 'Population Bottlenecks in Multicomponent Viruses: First Forays into the Uncharted Territory of Genome-formula Drift', *Current Opinion in Virology*, 33: 184–90.
- Hu, Z. et al. (2016) 'Genome Segments Accumulate with Different Frequencies in Bombyx Mori Bidesovirus', *Journal of Basic Microbiology*, 56: 1338–43.
- Iranzo, J., and Manrubia, S. C. (2012) 'Evolutionary Dynamics of Genome Segmentation in Multipartite Viruses', *Proceedings Biological Sciences*, 279: 3812–9.
- Kondrashov, F. A. (2012) 'Gene Duplication as a Mechanism of Genomic Adaptation to a Changing Environment', *Proceedings Biological Sciences*, 279: 5048–57.
- Lucía-Sanz, A., Aguirre, J., and Manrubia, S. (2018) 'Theoretical Approaches to Disclosing the Emergence and Adaptive Advantages of Multipartite Viruses', *Current Opinion in Virology*, 33: 89–95.
- Lucía-Sanz, A., and Manrubia, S. (2017) 'Multipartite Viruses: Adaptive Trick or Evolutionary Treat?' *NPJ Systems Biology and Applications*, 3: 34.
- Mansourpour, M. et al. (2022) 'Effects of an Alphasatellite on Life Cycle of the Nanovirus Faba Bean Necrotic Yellow Virus', *Journal of Virology*, 96: e0138821.
- Michalakakis, Y., and Blanc, S. (2020) 'The Curious Strategy of Multipartite Viruses', *Annual Review of Virology*, 7: 203–18.
- Milevko, Y., Joh, R. I., and Weitz, J. S. (2008) 'Small-Scale Copy Number Variation and Large-Scale Changes in Gene Expression', *Proceedings of the National Academy of Sciences*, 105: 16659–64.
- Moreau, Y. et al. (2020) 'The Genome Segments of Bluetongue Virus Differ in Copy Number in a Host-Specific Manner', *Journal of Virology*, 95: e01834–20.
- Nee, S. (1987) 'The Evolution of Multicompartmental Genomes in Viruses', *Journal of Molecular Evolution*, 25: 277–81.
- Ojosnegros, S. et al. (2011) 'Viral Genome Segmentation Can Result from a Trade-Off Between Genetic Content and Particle Stability', *PLoS Genetics*, 7: e1001344.
- Pressing, J., and Reanney, D. C. (1984) 'Divided Genomes and Intrinsic Noise', *Journal of Molecular Evolution*, 20: 135–46.
- R Core Team. (2021) *R: A Language and Environment for Statistical Computing*. Vienna: R Foundation for Statistical Computing.
- Ruijter, J. M. et al. (2009) 'Amplification Efficiency: Linking Baseline and Bias in the Analysis of Quantitative PCR Data', *Nucleic Acids Research*, 37: e45.
- Sasani, T. A. et al. (2018) 'Long Read Sequencing Reveals Poxvirus Evolution through Rapid Homogenization of Gene Arrays', *eLife*, 7: e35453.
- Shukla, A., Chatterjee, A., and Kondabagil, K. (2018) 'The Number of Genes Encoding Repeat Domain-Containing Proteins Positively Correlates with Genome Size in Amoebal Giant Viruses', *Virus Evolution*, 4: vex039.
- Sicard, A. et al. (2016) 'The Strange Lifestyle of Multipartite Viruses', *PLOS Pathogens*, 12: e1005819.
- et al. (2019) 'A Multicellular Way of Life for a Multipartite Virus', *eLife*, 8.
- et al. (2013) 'Gene Copy Number Is Differentially Regulated in a Multipartite Virus', *Nature Communications*, 4: 2248.
- Sobel Leonard, A. et al. (2019) 'Transmission Bottleneck Size Estimation from Pathogen Deep-Sequencing Data, with an Application to Human Influenza A Virus', *Journal of Virology*, 93: e00936–19.
- Todd, R. T., and Selmecki, A. (2020) 'Expandable and Reversible Copy Number Amplification Drives Rapid Adaptation to Antifungal Drugs', *eLife*, 9: e58349.
- Tomanek, I. et al. (2020) 'Gene Amplification as a Form of Population-Level Gene Expression Regulation', *Nature Ecology & Evolution*, 4: 612–25.
- Wu, B. et al. (2017) 'Within-Host Evolution of Segments Ratio for the Tripartite Genome of Alfalfa Mosaic Virus', *Scientific Reports*, 7: 5004.
- Wu, C. H., and Black, L. W. (1995) 'Mutational Analysis of the Sequence-Specific Recombination Box for Amplification of Gene 17 of Bacteriophage T4', *Journal of Molecular Biology*, 247: 604–17.
- Yu, N. T. et al. (2019) 'Independent Modulation of Individual Genomic Component Transcription and a cis-Acting Element Related to High Transcriptional Activity in a Multipartite DNA Virus', *BMC Genomics*, 20: 573.
- Zwart, M. P. et al. (2021) 'Unresolved in Spatially Structured Environments' Advantages of Multipartitism, *Virus Evolution*, 7: veab004.
- Zwart, M. P., and Elena, S. F. (2020) 'Modeling Multipartite Virus Evolution: The Genome Formula Facilitates Rapid Adaptation to Heterogeneous Environments', *Virus Evolution*, 6: veaa022.