npg

## ORIGINAL ARTICLE

# An antisense RNA in a lytic cyanophage links *psbA* to a gene encoding a homing endonuclease

Andrew D Millard[1], Gregor Gierga[2], Martha R J Clokie[3], David J Evans[1], Wolfgang R Hess[2] and David J Scanlan[1]

[1]Department of Biological Sciences, University of Warwick, Coventry, UK; [2]Faculty of Biology and Freiburg Initiative in Systems Biology, Institute of Biology III, University of Freiburg, Schänzlestraße 1, Freiburg, Germany and [3]Department of Infection, Immunity and Inflammation, Medical Sciences Building, University of Leicester, Leicester, UK

**Cyanophage genomes frequently possess the *psbA* gene, encoding the D1 polypeptide of photosystem II. This protein is believed to maintain host photosynthetic capacity during infection and enhance phage fitness under high-light conditions. Although the first documented cyanophage-encoded *psbA* gene contained a group I intron, this feature has not been widely reported since, despite a plethora of new sequences becoming available. In this study, we show that in cyanophage S-PM2, this intron is spliced during the entire infection cycle. Furthermore, we report the widespread occurrence of *psbA* introns in marine metagenomic libraries, and with *psbA* often adjacent to a homing endonuclease (HE). Bioinformatic analysis of the intergenic region between *psbA* and the adjacent HE gene F-CphI in S-PM2 showed the presence of an antisense RNA (asRNA) connecting these two separate genetic elements. The asRNA is co-regulated with *psbA* and F-CphI, suggesting its involvement with their expression. Analysis of scaffolds from global ocean survey datasets shows this asRNA to be commonly associated with the 3′ end of cyanophage *psbA* genes, implying that this potential mechanism of regulating marine 'viral' photosynthesis is evolutionarily conserved. Although antisense transcription is commonly found in eukaryotic and increasingly also in prokaryotic organisms, there has been no indication for asRNAs in lytic phages so far. We propose that this asRNA also provides a means of preventing the formation of mobile group I introns within cyanophage *psbA* genes.**

## Introduction

Viruses are the most abundant biological entities in the oceans, with numbers estimated to be more than $10^{30}$ (Suttle, 2005). As important agents of microbial mortality, they have critical roles in nutrient cycling and structuring microbial communities, while also contributing to horizontal gene transfer by mediating genetic exchange (Suttle, 2005, 2007). Bacteriophages infecting the picocyanobacterial genera *Synechococcus* (Waterbury and Valois, 1993; Suttle and Chan, 1994; Lu *et al.*, 2001; Marston and Sallee, 2003; Millard and Mann, 2006; Marston and Amrich, 2009) and *Prochlorococcus* (Sullivan *et al.*, 2003) are some of the most well-characterized marine viruses. Such cyanophages are widely distributed and abundant ($>10^5$ ml$^{-1}$; Suttle and Chan, 1994), with most isolates belonging to the family

myoviridae (Waterbury and Valois, 1993; Suttle and Chan, 1994; Lu *et al.*, 2001; Sullivan *et al.*, 2003; Millard and Mann, 2006; Millard *et al.*, 2009) and fewer representatives thus far known from the siphoviridae (Waterbury and Valois, 1993; Sullivan *et al.*, 2009) and podoviridae (Waterbury and Valois, 1993; Suttle and Chan, 1994; Sullivan *et al.*, 2003) families.

Although cyanophages, similar to several other viruses, can divert the flow of carbon through the microbial loop, they are unique in being considered to be able to directly contribute to the photosynthetic process through their possession of phage versions of the core photosystem II reaction centre polypeptides D1 and D2, encoded by *psbA* and *psbD*, respectively. A recent study using metagenomic data showed that cyanophages also carry genes encoding complete subunits of photosystem I (Sharon *et al.*, 2009). However, to date, most cyanophage research has focused on photosystem II. During photosynthesis, D1 is continually being turned over and replaced by newly synthesized D1. It is postulated that expression of the phage-encoded D1 protein provides a means to maintain

photosynthesis even after host protein synthesis in infected cells is diminished, thus ensuring a source of energy for virus production (Mann *et al.*, 2003; Lindell *et al.*, 2004, 2005, 2007; Millard *et al.*, 2004; Clokie *et al.*, 2006). This is supported both by evidence that cyanophage *psbA* transcripts can be detected throughout the infection cycle (Lindell *et al.*, 2005, 2007; Clokie *et al.*, 2006) and by the fact that the cyanophage D1 polypeptide increases in abundance during infection (Lindell *et al.*, 2007). Moreover, recent modelling studies have suggested that there is an increased fitness advantage for cyanophages possessing *psbA*, particularly under high-light conditions (Bragg and Chisholm, 2008; Hellweger, 2009).

Genome sequencing (Millard *et al.*, 2004, 2009; Mann *et al.*, 2005; Sullivan *et al.*, 2005; Weigele *et al.*, 2007) and PCR screening (Sullivan *et al.*, 2006; Wang and Chen, 2008; Marston and Amrich, 2009) efforts indicate that *psbA* is widely distributed in cyanophage isolates, whereas cyanophage-derived *psbA* transcripts can also be readily detected in the marine environment (Zeidner *et al.*, 2005; Sullivan *et al.*, 2006; Sharon *et al.*, 2007; Chenard and Suttle, 2008). Phylogenetic analysis of phage *psbA* suggests that it has been inherited from its cyanobacterial hosts on a number of occasions (Lindell *et al.*, 2004; Millard *et al.*, 2004; Zeidner *et al.*, 2005; Sullivan *et al.*, 2006), but with evidence of significant intragenic recombination between the phage and the host gene (Zeidner *et al.*, 2005; Sullivan *et al.*, 2006).

An unusual feature of the first viral *psbA* gene discovered, in cyanophage S-PM2, was that it contained a group I intron (Millard *et al.*, 2004). Although group I introns are common in other bacteriophage genomes, the sequencing of hundreds of other cyanophage (Sullivan *et al.*, 2006; Chenard and Suttle, 2008; Marston and Amrich, 2009) and host (Zeidner *et al.*, 2003; Sharon *et al.*, 2007) *psbA* genes has revealed only one more containing an intron (Millard *et al.*, 2004). This is likely due to the fact that the widely used reverse *psbA* PCR primer (Zeidner *et al.*, 2003) does not amplify *psbA* that contains an intron in the same position as found in S-PM2, thereby preventing detection of *psbA* genes with introns in the same position.

The origin of the *psbA* intron in S-PM2 is still unknown. Although introns are present in some *psbA* genes of chloroplasts (Maul *et al.*, 2002; Brouard *et al.*, 2008), they have thus far not been found in any of the cyanobacterial orthologues. The mobility of the *psbA* intron has previously been proposed to be mediated by an endonuclease in a process known as 'homing', which would transfer the intron and flanking DNA containing the endonuclease into an intron-less allele of *psbA* (Millard *et al.*, 2004). The recent characterization of a homing endonuclease (HE) situated immediately downstream of *psbA* in S-PM2, which is only able to cut intron-less copies of *psbA*, supports this idea (Zeng *et al.*, 2009). HEs are generally considered as selfish elements that target highly specific DNA target sites

of 14–40 bp in length (Jurica and Stoddard, 1999) and allow transfer of themselves, and the introns in which they often reside, to cognate sites within a population (Jurica and Stoddard, 1999). Paradoxically, they can tolerate sequence variation within their target site, allowing targeting of new hosts (Jurica and Stoddard, 1999). Although often found within introns, this is not always the case, with 'intron-less homing' observed between the bacteriophages T4 and T2 (Liu *et al.*, 2003).

Despite knowledge that S-PM2 *psbA* is expressed during the lytic cycle (Clokie *et al.*, 2006), it is not known whether the intron is spliced *in vivo*, how widespread these introns are in other cyanophage genomes or how they might have been acquired.

Another intriguing type of RNA molecule, which was discovered first in bacteriophages almost 40 years ago, is antisense RNA (asRNA). Such naturally occurring asRNAs were postulated first in bacteriophage-λ (Spiegelman *et al.*, 1972), and only afterwards observed in bacteria (Itoh and Tomizawa, 1980; Lacatena and Cesareni, 1981) and even later in eukaryotes. More recently, it was found that expression of the photosynthetic gene *isiA* in the cyanobacterium *Synechocystis* sp. PCC6803 is regulated by the 177-nucleotide (nt) long asRNA IsrR (Duhring *et al.*, 2006). However, asRNAs have not been described for any cyanophage gene thus far.

## Materials and methods

### Culturing

*Synechococcus* sp. WH7803 was cultured in the ASW medium (Wyman *et al.*, 1985) in 100 ml batch cultures in 250 ml conical flasks under constant illumination (5–36 µmol photons $m^{-2} s^{-1}$) at 25 °C. Larger volumes were grown in 0.5 l vessels under constant shaking (150 r.p.m.). Cyanophage S-PM2 stocks and phage titre were produced as reported previously (Wilson *et al.*, 1993).

### Host infection

The S-PM2 infection cycle has previously been well characterized with lysis of *Synechococcus* proceeding 9 h after infection (Wilson *et al.*, 1996; Clokie *et al.*, 2006). Therefore, 50 ml samples were collected before infection and then at 1, 3, 6 and 9 h after infection. Phage was added at a multiplicity of infection of ∼5, to ensure infection of all cells. Samples were immediately centrifuged at 8000 *g* to pellet the samples, which were then snap frozen in 0.5 ml of RNA extraction buffer (10 mM NaAc, pH 4.5; 200 mM sucrose, 5 mM EDTA) and stored at −80 °C until samples were further processed. Three biological replicates were collected.

### RNA extraction

Total RNA was extracted on the basis of a previously described method (Logemann *et al.*, 1987). Briefly,

frozen samples were gently thawed in 3 vol of Z buffer (8 M guanidinium hydrochloride; 50 mM β-mercaptoethanol; 20 mM EDTA) at room temperature for 30 min. Samples were extracted with the addition of $\frac{1}{2}$ vol of phenol (pH 4.5) at 65 °C for 30 min, followed by the addition of chloroform:isoamyl alcohol for 15 min. RNA was precipitated in 1 vol of isopropanol, followed by a wash in 70% (v/v) ethanol. RNA was treated with Turbo DNase I (Applied Biosystems, Warrington, UK) for 2 h at 37 °C, extracted with phenol:cholorform:isoamyl alcohol (25:24:1), re-precipitated with 3 M NaAc and tested for DNA contamination using PCR primers gp23F/R.

### In vivo *splicing*
RNA was extracted and cDNA synthesized. cDNA synthesis was carried out in 20 μl reaction volumes with 600 ng of total RNA. Each reaction contained 1 μl of 20 × dNTP mix (10 mM dGTP, dCTP, dATP and dTTP), 5 μg random hexamers (VHBio, Gateshead, UK) or 2 pM of gene-specific primer, 4 μl of 5 × buffer (250 mM Tris pH 8.3, 375 mM KCl, 15 mM $MgCl_2$), 2 μl of 0.1 M DTT, 200 Units Superscript III (Invitrogen, Paisley, UK) and water to a final volume of 20 μl. The RNA, water and random hexamers were mixed, heated to 65 °C for 10 min and thereafter cooled to 4 °C in a thermal cycler, before the addition of 5 × buffer, DTT and superscript, heated to 50 °C for 50 min, before finally heating to 70 °C for 10 min.

The primers psbA_F and psbA_R (see Supplementary Table S1) were used to amplify PCR products of 1080 and 1291 bp in length, respectively, depending on whether splicing had occurred. PCR was carried out with 0.03 Units ml$^{-1}$ Vent DNA polymerase in 1 × buffer (20 mM Tris-HCl, 10 mM $(NH_4)_2SO_4$), 2 mM $MgCl_2$, 0.2 mM dNTPs and 40 pM of each primer. PCR cycling conditions were 35 cycles of 94 °C for 10 s, 55 ± °C for 15 s and 72 °C for 20 s and a final incubation at 72 °C for 2 min. PCR products were sequenced in-house using an ABI 3730 automated sequencer (Applied Biosystems, Foster City, USA).

### Quantitative reverse transcriptase-PCR
PCR primers were designed using Primer Express (ABI, Warrington, UK). Sequences of PCR primers are reported in Supplementary Table S1. Various primer concentrations were tested and optimized to ensure that the amplification efficiency was within the required limits to implement relative quantification using $2^{\Delta\Delta CT}$ (Livak and Schmittgen, 2001). The amplification efficiencies for target and reference primers sets were tested by ensuring that the slope of the line was < 0.1, when log input DNA concentration was plotted against $^{\Delta}CT$.

A no reverse transcriptase control PCR reaction was used to assess DNA contamination of experimental samples. Any sample found to be contaminated was subjected to further DNAse treatment (see the 'RNA extraction' section above) and the process was repeated until the control PCR reaction proved negative. cDNA synthesis was then carried out as described above, with the gene-specific primer ncRNA_R used for the synthesis of cDNA from the ncRNA Cyanophage Functional RNA I (CfrI).

PCR reactions of 1 × power SYBR green mix (ABI), 150 μM forward primer, 150 μM reverse primer and 10 ng cDNA were used for the amplification of 16S rRNA, *psbA* and open reading frame (ORF)177 (F-CphI), whereas for psbA_ncRNA, 50 ng of cDNA was used per well. Thermal cycling was carried out in a 7500 sequence detector (Applied Biosystems) with an initial step of 95 °C for 10 min, followed by 40 cycles of 95 °C for 30 s, followed by 62 °C for 1 min and a final dissociation step. The fold change of each gene was determined using the $2^{\Delta\Delta CT}$ method (Livak and Schmittgen, 2001) using 16S rRNA as the calibrator. Absolute transcript abundance was calculated from a standard curve, while a dilution series of purified phage DNA was used to construct the curve.

### ncRNA prediction
Sequence-dependent RNA structure within the phage genome and scaffold sequences were identified by comparing the folding free energy of the native sequence with a large number of sequence order randomized controls. In practice, the scaffolds were divided into 200-nt fragments for both strands of DNA, overlapping by 190 bp, each of which was randomized 999 times using a method (designated NDR; implemented in the Simmonics suite of sequence analysis programs (Simmonds and Smith, 1999), available from http://www.picornavirus.org), which retained both the nucleotide and dinucleotide composition. Sequences were stored in a MySQL database, and the folding free energy for each was determined using hybrid-ss-min from the Unafold (http://dinamelt.bioinfo.rpi.edu/) suite of programs (Markham and Zuker, 2008), automated using Perl scripts. For each fragment, the mean folding energy difference, expressed as the percentage difference between native and the mean of the randomized sequences from the same fragment, was determined. In addition, the position of the native sequence in the distribution of energies of the randomized fragments—expressed as the Nth percentile—was calculated.

### 5′ *rapid amplification of cDNA ends*
5′ rapid amplification of cDNA ends (RACE) experiments were conducted on the basis of the protocol of Steglich *et al.* (2008). Briefly, RNA was treated with tobacco acid pyrophosphorylase (1 Unit/1 μg RNA; Epicentre, Madison, WI, USA) for 1 h at 37 °C, followed by phenol/chloroform extraction and ethanol precipitation. A synthetic RNA oligonucleotide (0.5 μl oligonucleotide (10 mM)/4 μg RNA; 5′-AUAUGCGCGAAUU CCUGUAGAACGAACACUAGAAGAAA-3′, Invitrogen,

Darmstadt, Germany) was ligated to RNA using T4 RNA ligase (3 Units/1 μg RNA; Fermentas, St Leon-Rot, Germany) for 1 h at 37 °C, followed by phenol/chloroform extraction and ethanol precipitation. Three control reactions were performed: (1) omitting tobacco acid pyrophosphorylase, (2) omitting tobacco acid pyrophosphorylase and RNA oligonucleotide and (3) dephosphorylating RNA before ligation with calf intestine alkaline phosphatase (0.1 Unit/1 μg RNA; Fermentas) at 37 °C for 1 h, followed by phenol/chloroform extraction and ethanol precipitation. For reverse transcription, 250 ng oligonucleotide-linked RNA per gene was incubated with 0.8 Units Omni-script reverse transcriptase (Qiagen, Hilden, Germany) in the provided reaction buffer, supplemented by 0.08 μM gene-specific primer and 1 mM dNTPs. Incubation was carried out at 42 °C for 2 h with a final inactivation step at 95 °C for 5 min. All reactions were performed in the presence of 40 Units Ribolock RNase Inhibitor (Fermentas). cDNA was amplified by PCR in GoTaq reaction buffer containing 1 Units GoTaq polymerase (Promega, Mannheim, Germany), 0.2 mM dNTPs, 3.5 mM MgCl$_2$, a gene-specific primer (0.2 μM) and an RNA oligonucleotide-specific primer (0.2 μM) with the following parameters: 93 °C/3 min, 35 cycles of 93 °C/30 s, 50 °C/30 s, 72 °C/45 s, followed by 72 °C/5 min. Amplified PCR fragments were gel-excised and purified on Nucleospin columns (Macherey & Nagel, Düren, Germany) and then cloned into plasmid pGEM-T (Promega). After transformation into *Escherichia coli* XL1-Blue, plasmid inserts were amplified by colony PCR, purified on Nucleospin columns and sequenced using an ABI 3130XL automatic DNA sequencer (Applied Biosystems).

### Northern blot analysis

RNA samples (20 μg) were denatured for 5 min at 65 °C in loading buffer (Fermentas), separated on 10% (w/v) urea polyacrylamide gels at 90 V for 16 h and transferred to Hybond-N$^+$ nylon membranes (GE Healthcare, Munich, Germany) by electroblotting for 1 h at 400 mA. The membranes were hybridized with single-stranded [α-32P]UTP-labelled transcripts. After pre-hybridization in 50% (v/v) deionized formamide, 7% (w/v) SDS, 250 mM NaCl and 120 mM Na(PO$_4$) pH 7.2 for 2 h, hybridization was performed at 62 °C overnight in the same buffer. The membranes were washed in 2 × SSC (3 M NaCl, 0.3 M sodium citrate, pH 7.0), 1% (w/v) SDS for 10 min; 1 × SSC, 0.5% (w/v) SDS for 10 min; and briefly in 0.1 × SSC, 0.1% (w/v) SDS. All wash steps were performed 5 °C below the hybridization temperature. Signals were detected and analysed on a Personal Molecular Imager FX system with Quantity One software (Bio-Rad, Munich, Germany).

### Bioinformatics analyses

Introns were initially identified in the global ocean survey (GOS) dataset using tblastx with the intron sequence of S-PM2 as the query sequence with an e-value cutoff of < 10$^{-3}$. Intron insertion sites (IISs) were determined manually by identifying the point at which the translated *psbA* sequences resulted in a premature stop codon or did not align with other highly conserved PsbA sequences. The intron sequence was then extracted. Using a custom Perl script, full-length scaffolds were blasted against a custom Basic Local Alignment Search Tool (BLAST) database containing uniprot 100 and all publically accessible cyanobacterial and viral genomes (as available in October 2008) to identify any other gene on the scaffolds. Again, a cutoff value of < 10$^{-3}$ was used to identify gene fragments. The same approach was used to identify homologues of F-CphI and genes adjacent to it.

## Results

### Confirmation of intron splicing *in vivo*

To determine whether the intron found within the S-PM2 *psbA* is spliced *in vivo* during infection, non-quantitative reverse transcriptase-PCR was used to amplify the *psbA* transcript from RNA extracted at 1, 3 and 9 h after infection. Two PCR products were observed, one of 1291 bp and the other 1080 bp in length (Figure 1). The 1291-bp product corresponds to the unspliced transcript (pre-mRNA) and the smaller product to the spliced transcript. Comparison of the sequence of the smaller product with the *psbA* gene sequence showed that splicing occurred between codons 334 and 335 as predicted previously (Millard *et al.*, 2004).

### Searching metagenomic data for intron sequences

After confirmation of splicing of the *psbA* intron, we searched metagenomic libraries for the presence of other introns similar to that found in S-PM2. The GOS data set (Rusch *et al.*, 2007) was searched using BLAST available through Community



**Figure 1** *In vivo* splicing of a group I intron within the cyanophage S-PM2 *psbA* gene. RNA isolated from S-PM2-infected *Synechococcus* sp. WH7803 was analysed by RT-PCR from samples collected at 1, 3 and 9 h after infection. No reverse transcriptase controls (nrtc) were used to test for contaminating DNA in purified RNA. Genomic DNA from S-PM2 was used a positive control (lane G). No template sample was used as a negative control (lane C-ve). The 1 kb and 1.5 kb size standards are marked.

Cyberinfrastructure for Advanced Marine Microbial Ecology Research and Analysis (CAMERA) (Seshadri *et al.*, 2007). A total of 16 scaffold sequences were identified as having introns that were similar to that of S-PM2 on the basis of sequence identity (Table 1). All intron sequences were localized within *psbA* genes and varied in length from 212 to 818 nt. With the exception of JCVI_S-CAF_10096627024160, the introns were not found to contain ORFs. Six different IISs were found, located throughout the length of the *psbA* gene. The most common IIS found in the 16 scaffolds (8/16) was located after codon 334 (Figure 2), which is the same position as reported in S-PM2 and S-RSM88 (Millard *et al.*, 2004). A single intron is positioned nearby, after codon 338 (Figure 2). In all, 8 of 16 of these introns were very similar in sequence to that of S-PM2 with percentage nucleotide identities ranging from 62% to 92%. In addition, there was conservation of the conserved paired helices (Supplementary Figure S1).

Intriguingly, three introns had IISs after codon 60, which is the same IIS as *psbA* intron 1 in the chloroplast genome of *Oedogonium cardiacum* (Brouard *et al.*, 2008) and *Chlamydomonas reinhardtii* (Maul *et al.*, 2002), with a further two introns located after codon 252, which is very close to the IIS of intron 4 in *C. reinhardtii*, which is inserted after codon 254 (Figure 2). The remaining introns were localized in IISs that have not been previously documented in *psbA* genes. The introns with an IIS close to or matching that of *C. reinhardtii* and *O. cardiacum* were substantially smaller than the introns found in these two chloroplast genomes. However, they did retain regions of sequence conservation at the 5′ end of the intron, with 52–59% nucleotide identity to the chloroplast introns (Supplementary Figure S2). Although these introns were detected on the basis of their sequence similarity to S-PM2, there is a small possibility that some chloroplast introns are present in the GOS dataset. However, most would be excluded by the <0.8-μm pore sized

**Table 1** Cyanophage genomes and global ocean survey scaffolds from which introns were identified

| Cyanophage/scaffold sequence | Presence of the 'G/KETTXXXSQ/H' motif in PsbA | Intron length | Phylogenetic classification | % ID to S-PM2 intron | Scaffold length | Area of isolation (Reference) |
|---|---|---|---|---|---|---|
| S-PM2 | ✔ | 212 | Cyanophage | 100 | N/A | Plymouth, United Kingdom (Wilson *et al.*, 1993) |
| S-RSM88 | ✔ | 212 | Cyanophage | 100 | N/A | Gulf of Aqaba, Red Sea (Millard *et al.*, 2006) |
| JCVI_SCAF_1101667034653 | ✔ | 212 | Unknown | 92 | 1018 | GS007—Northern Gulf of Maine (Rusch *et al.*, 2007) |
| JCVI_SCAF_1101669070555 | ✔ | 212 | Unknown | 92 | 959 | GS007—Northern Gulf of Maine (Rusch *et al.*, 2007) |
| JCVI_SCAF_1101668247417 | Fragment too short | 207 | Unknown | 73 | 1574 | GS020—Lake Gatun, Panama (Rusch *et al.*, 2007) |
| JCVI_SCAF_1101669425113 | Fragment too short | 207 | Unknown | 73 | 1574 | GS007—Northern Gulf of Maine (Rusch *et al.*, 2007) |
| JCVI_SCAF_1101669142352 | Fragment too short | 263 | Unknown | 70 | 566 | GS010—Cape May, NJ, USA (Rusch *et al.*, 2007) |
| JCVI_SCAF_1101667044432 | ✔ | 263 | Unknown | 70 | 566 | GS010—Cape May, NJ, USA (Rusch *et al.*, 2007) |
| JCVI_SCAF_1098315327957 | Fragment too short | 259 | Unknown | 70 | 1804 | MOVE858—Chesapeake Bay, USA (Rusch *et al.*, 2007) |
| JCVI_SCAF_1101669414852 | Fragment too short | 241 | Unknown | 70 | 1522 | GS020—Lake Gatun ,Panama (Rusch *et al.*, 2007) |
| JCVI_SCAF_1101668234973 | Fragment too short | 241 | Unknown | 70 | 1522 | GS020—Lake Gatun, Panama (Rusch *et al.*, 2007) |
| JCVI_SCAF_1097156666624 | ✔ | 204 | Cyanophage | 65 | 3436 | GS020—Lake Gatun ,Panama (Rusch *et al.*, 2007) |
| JCVI_SCAF_1096627283123 | ✔ | 204 | Cyanophage | 65 | 3437 | GS002—Gulf of Maine, Canada (Rusch *et al.*, 2007) |
| JCVI_SCAF_1097207205912 | Fragment too short | 204 | Unknown | 62 | 920 | GS020—Lake Gatun, Panama (Rusch *et al.*, 2007) |
| JCVI_SCAF_1096626190549 | Fragment too short | 207 | Unknown | 62 | 920 | GS020—Lake Gatun, Panama (Rusch *et al.*, 2007) |
| JCVI_SCAF_1096627024703 | Fragment too short | 236 | Unknown | 65 | 3221 | GS020—Lake Gatun, Panama (Rusch *et al.*, 2007) |
| JCVI_SCAF_1096627666661 | ✔ | 164 min | Unknown | | 1508 | GS020—Lake Gatun, Panama (Rusch *et al.*, 2007) |
| JCVI_SCAF_1096627024160 | ✔ | 818 | Cyanophage | 53 | 3755 | GS020—Lake Gatun, Panama (Rusch *et al.*, 2007) |

Abbreviation: N/A, not applicable.
'Fragment too short' denotes the *psbA* sequence was not long enough to cover the region where the cyanophage/cyanobacterial-specific PsbA motif R/KETTXXXSQ/H is found.

filter used in the collection of GOS samples, which would preclude collection of chloroplast containing eukaryotes.

To ascertain the origin of these intron-containing *psbA* genes, we searched for the cyanophage/cyano-bacterial-specific PsbA motif R/KETTXXXSQ/H (Sharon *et al.*, 2007). This motif was found in all *psbA* sequences (Table 1), whenever the sequence fragment was long enough to encompass this region, suggesting that the identified *psbA* genes are all of cyanobacterial or cyanophage origin and not from chloroplasts.

*Phylogenetic analysis of* psbA *sequences*
To further confirm the origin of the identified *psbA* sequences, phylogenetic analysis was carried out. Only sequences that were >920 bp in length were used, as this encompassed all regions where introns have been found to be inserted (the intron sequence itself was not included in the analysis). Sequences shorter than this were excluded from the analysis,

along with *psbA* fragments used in previous phylogenetic analyses (Sullivan *et al.*, 2006; Chenard and Suttle, 2008), which did not encompass the most common IISs due to the PCR primers used.

Phylogenetic analysis of *psbA* genes was essentially congruent with 16S rRNA phylogenies with eukaryotic algae clearly separate from cyanobacteria. Discrete clades of both high-light-adapted and low-light-adapted *Prochlorococcus* strains were discernible and were, in turn, distinct from *Synechococcus* (Figure 3). Phage isolates infecting *Prochlorococcus* formed a sister group to high-light-adapted *Prochlorococcus* strains as reported previously (Sullivan *et al.*, 2006). Phage isolates infecting *Synechococcus* formed a clade distinct from their *Synechococcus* hosts, whereas the intron-containing *psbA* sequences fell into two discrete clades that did not contain any cultured *Synechococcus* or *Prochlorococcus* strains, or cyanophage isolates (Figure 3). Clearly, the identified *psbA* genes that contain introns are cyanophage/cyanobacterial in origin as they do not group with eukaryotic algae.
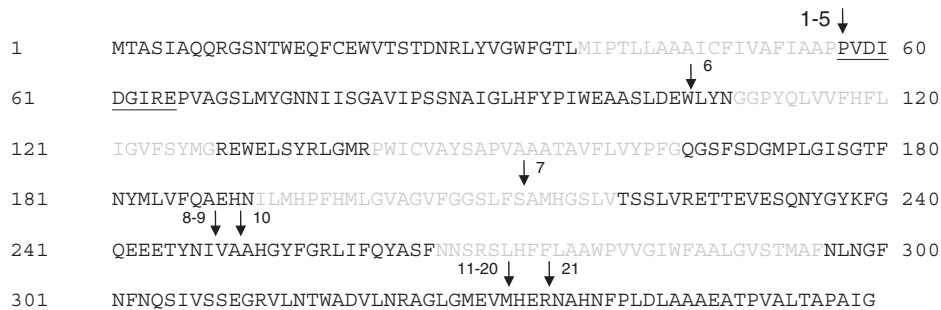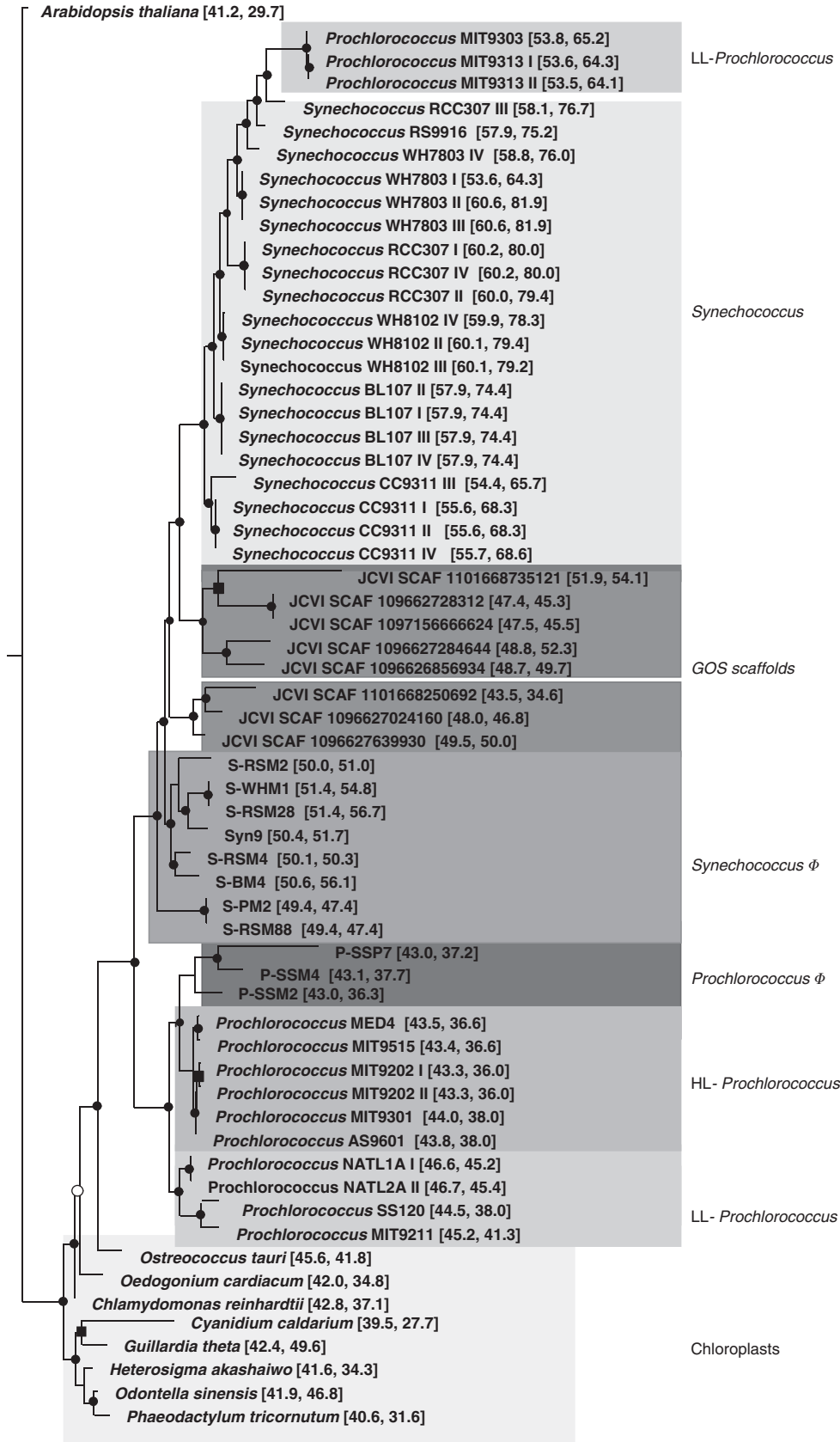


**Figure 2** Intron insertion sites (IISs) within *psbA*. Amino-acid sequences derived from *psbA* genes identified to have introns were aligned. Owing to the partial sequence of some of the *psbA* genes, the IIS is reported relative to the position of the full-length sequence of S-PM2. The trans-membrane domains of the D1 protein are marked by grey text. The amino-acid sequences targeted by the universal *psbA* primer set (Zeidner *et al.*, 2003) are underlined. IISs are marked by arrows, with numbers corresponding to the following sequences: 1: JCVI_SCAF1101669142352, 2: JCVI_SCAF_1101667044432, 3: JCVI_SCAF_1098315327957, 4: *Chlamydomonas reinhardtii* intron 4, 5: *Oedogonium cardiacum* intron 1, 6: JCVI_SCAF_1096627024160, 7: JCVI_SCAF_1096627666661, 8: JCVI_SCAF_1101668247417, 9: JCVI_SCAF_1101669425113, 10: *Chlamydomonas reinhardtii* intron 1, 11: S-PM2, :12 S-RSM88, 13: JCVI_SCAF_1101667034653, 14: JCVI_SCAF_1101669070555, 15: JCVI_SCAF_1097156666624, 16: JCVI_SCAF_1096627283123, 17: JCVI_SCAF_1097207205912, 18: JCVI_SCAF_1096626190594, 19: JCVI_SCAF_1101669414852, 20: JCVI_SCAF_1101668234973, 21: JCVI_SCAF_1096627024703.



**Figure 3** Phylogenetic relationships among *psbA* genes of cyanophages, cyanobacteria and plastids from cultures and environmental samples. Trees are based on an alignment of 925 nucleotides, clade support values are the result of 200 000 iterations and a burn-in of 25% using MRBAYES (Huelsenbeck and Ronquist, (2001). Clade support values >90 are marked by ●, >80 and < 90 by ■ and those >70 and <80 by a ○. GenBank accession numbers of *psbA* sequences used for phylogenetic analysis were as follows: *Synechococcus* (*Synechococcus* BL107: acc NA_AAT20000000, *Synechococcus* sp. WH8102: acc NC_005070, *Synechococcus* sp. WH7803: acc NC_00009481, *Synechococcus* sp. RCC307: acc NC_00009482, *Synechococcus* sp. RS9916: acc NZ_AA0A00000000, *Synechococcus* sp. CC9311: acc NC_008319); *Prochlorococcus* (*Prochlorococcus* sp. MIT9303: acc NC_008820, *Prochlorococcus* sp. MIT9319: acc NC_005071, *Prochlorococcus* sp. MED4: acc NC_005072, *Prochlorococcus* sp. MI9515: acc NC_008817, *Prochlorococcus* sp. MIT9202: acc NZ_ACDW00000000, *Prochlorococcus* sp. NAT2LA: acc NC_007335, *Prochlorococcus* sp. AS9601: acc NC_008816, *Prochlorococcus* sp. MIT9301: acc NC_009091, *Prochlorococcus* sp. MIT9211: acc NC_009976, *Prochlorococcus* sp. SS120: acc NC_00xxxx, ); *Synechococcus* phage (Syn9: acc NC_008296, S-RSM4: acc CAR63316.1, S-PM2 : acc NC_006820, S-RSM88: acc AJ629075, S-RSM2: acc AJ628768, S-WHM1:acc AJ628769, S-RSM28: acc AJ629221, S-BM4: acc AJ628858); *Prochlorococcus* phage (P-SSM4: acc NC_006884, P-SSM2: acc NC_006883, P-SSP7: acc NC_006882); plastids (*Ostreococcus tauri*: acc NC_008289, *Oedogonium cardiacum*: acc NC_011031, *Chlamydomonas reinhardtii*: acc NC_005353, *Cyanidium caldarium*: acc NC_001840, *Guillardia theta*: acc NC_000926, *Heterosigma akashaiwo*: acc NC_010772, *Odontella sinensis*: acc NC_001713, *Phaeodactylum tricornutum*: acc NC_008588). *Arabidopsis thaliana*: acc NC_009032 was used to root the tree. Roman numerals are used to denote the different copies of *psbA* found within the genomes of *Synechococcus* and *Prochlorococcus*. The numbers in square brackets are the average mol %GC content and the third base mol %GC content, respectively.

*Arabidopsis thaliana* [41.2, 29.7]

*Prochlorococcus* MIT9303 [53.8, 65.2]
*Prochlorococcus* MIT9313 I [53.6, 64.3]      LL-*Prochlorococcus*
*Prochlorococcus* MIT9313 II [53.5, 64.1]

*Synechococcus* RCC307 III [58.1, 76.7]
*Synechococcus* RS9916 [57.9, 75.2]
*Synechococcus* WH7803 IV [58.8, 76.0]
*Synechococcus* WH7803 I [53.6, 64.3]
*Synechococcus* WH7803 II [60.6, 81.9]
*Synechococcus* WH7803 III [60.6, 81.9]
*Synechococcus* RCC307 I [60.2, 80.0]
*Synechococcus* RCC307 IV [60.2, 80.0]
*Synechococcus* RCC307 II [60.0, 79.4]
*Synechococccus* WH8102 IV [59.9, 78.3]       *Synechococcus*
*Synechococcus* WH8102 II [60.1, 79.4]
*Synechococcus* WH8102 III [60.1, 79.2]
*Synechococcus* BL107 II [57.9, 74.4]
*Synechococcus* BL107 I [57.9, 74.4]
*Synechococcus* BL107 III [57.9, 74.4]
*Synechococcus* BL107 IV [57.9, 74.4]
*Synechococcus* CC9311 III [54.4, 65.7]
*Synechococcus* CC9311 I [55.6, 68.3]
*Synechococcus* CC9311 II [55.6, 68.3]
*Synechococcus* CC9311 IV [55.7, 68.6]

JCVI SCAF 1101668735121 [51.9, 54.1]
JCVI SCAF 109662728312 [47.4, 45.3]
JCVI SCAF 1097156666624 [47.5, 45.5]
JCVI SCAF 1096627284644 [48.8, 52.3]
JCVI SCAF 1096626856934 [48.7, 49.7]          *GOS scaffolds*

JCVI SCAF 1101668250692 [43.5, 34.6]
JCVI SCAF 1096627024160 [48.0, 46.8]
JCVI SCAF 1096627639930 [49.5, 50.0]

S-RSM2 [50.0, 51.0]
S-WHM1 [51.4, 54.8]
S-RSM28 [51.4, 56.7]
Syn9 [50.4, 51.7]
S-RSM4 [50.1, 50.3]                            *Synechococcus Φ*
S-BM4 [50.6, 56.1]
S-PM2 [49.4, 47.4]
S-RSM88 [49.4, 47.4]

P-SSP7 [43.0, 37.2]
P-SSM4 [43.1, 37.7]                            *Prochlorococcus Φ*
P-SSM2 [43.0, 36.3]

*Prochlorococcus* MED4 [43.5, 36.6]
*Prochlorococcus* MIT9515 [43.4, 36.6]
*Prochlorococcus* MIT9202 I [43.3, 36.0]
*Prochlorococcus* MIT9202 II [43.3, 36.0]      HL-*Prochlorococcus*
*Prochlorococcus* MIT9301 [44.0, 38.0]
*Prochlorococcus* AS9601 [43.8, 38.0]

*Prochlorococcus* NATL1A I [46.6, 45.2]
Prochlorococcus NATL2A II [46.7, 45.4]
*Prochlorococcus* SS120 [44.5, 38.0]           LL-*Prochlorococcus*
*Prochlorococcus* MIT9211 [45.2, 41.3]

*Ostreococcus tauri* [45.6, 41.8]
*Oedogonium cardiacum* [42.0, 34.8]
*Chlamydomonas reinhardtii* [42.8, 37.1]
*Cyanidium caldarium* [39.5, 27.7]
*Guillardia theta* [42.4, 49.6]                Chloroplasts
*Heterosigma akashaiwo* [41.6, 34.3]
*Odontella sinensis* [41.9, 46.8]
*Phaeodactylum tricornutum* [40.6, 31.6]

These newly identified *psbA* sequences fell into two clades, which are sister groups but clearly separated from the well-defined *Synechococcus* host clade. Their closer phylogenetic proximity to *psbA* genes from phage isolates infecting *Synechococcus* suggests that these sequences are of *Synechococcus* phage origin and not from their *Synechococcus* hosts (Figure 3). This is further supported by examination of both the average mol %guanine-cytosine (GC) content and third codon mol %GC content, with the newly identified *psbA* sequences possessing values that are markedly different from the *Synechococcus* host and much closer to that observed in known *Synechococcus* phages.

From phylogenetic analysis and/or detection of the cyanobacterial-specific *psbA* motif, it was possible to confirm that 7 of 16 introns were inserted into genes of cyanophage/cyanobacterial origin, with JCVI_SCAF_1101667044432 containing an intron inserted in the same IIS as found in the chloroplasts of algae and JCVI_SCAF_1096627024160 an intron at a unique site (Figure 2). However, for nine scaffolds, the origin of the *psbA* sequence could not be determined unequivocally as the *psbA* fragment was not long enough for phylogenetic analysis or its length did not extend to the region where the cyanobacterial-specific motif is located (Table 1). Where possible, all scaffolds were examined in further detail to identify the origin of those genes adjacent to *psbA* if any were present. JCVI_SCAF_1097207205912 and JCVI_SCAF_1096626190549 both have genes encoding homologues of the cyanophage protein F-CphI, thus suggesting that these are also phage-encoded copies of *psbA* (Supplementary Figure S3). JCVI_SCAF_1096627024703 contains *talC* and *gnd* genes that, although found in *Synechococcus* and *Prochlorococcus* host genomes, are also known to be widespread in cyanophage genomes (Millard *et al.*, 2009). Indeed, these genes have the highest sequence similarity to cyanophage-encoded versions of these genes, and phylogenetic analysis confirms that they are of cyanophage origin (Supplementary Figures S5 and S6). Again, this suggests that the associated scaffolds are also of cyanophage origin (Supplementary Figure S3). Unfortunately, for the remaining scaffolds, it was not possible to identify genes adjacent to *psbA* because of the limited size of the scaffold sequences. However, given that the *psbA* sequences on these scaffolds share higher sequence similarity with *psbA* from cyanophages, and IISs that are present on scaffolds of cyanophage origin, it is reasonable to assume that they are also likely to be of phage origin.

### Homing endonuclease

During the identification of intron sequences in the GOS dataset, it became apparent that genes with similarity to the HE (F-CphI) of S-PM2 were located adjacent to the intron-containing *psbA* genes (Supplementary Figure S3). Phylogenetic analysis showed that *psbA*-containing introns also grouped with *psbA* genes that are found adjacent to a HE (Figure 3). The arrangement of a HE adjacent to *psbA* has been observed in the genome of S-PM2 and S-RSM88 (Millard *et al.*, 2004). It has been suggested that this arrangement reflects the independent convergence of two separate genetic elements: (1) the intron within *psbA* and (2) the HE F-ChpI downstream of *psbA*, on a common DNA target in a process termed 'collaborative homing' and that this is the penultimate step in the proposed pathway for the formation of mobile group I introns (Bonocora and Shub, 2009).

In an effort to identify more intron sequences and to determine whether the arrangement of *psbA* adjacent to an HE-encoding gene is common, the GOS dataset was searched using the amino-acid sequence of F-CphI from S-PM2 as the query. A total of 89 scaffolds were identified using Basic Local Alignment Search Tool as having similarity to F-CphI, thus demonstrating that it is readily detected in the environment. The scaffold sequences in which HEs were found were extracted, and any genes adjacent to the HE were identified. This was possible for 23 scaffold sequences (Table 2). Of the genes adjacent to F-CphI homologues, 15 were identified as *psbA* and 4 as *psbD*. These *psbA* sequences were then searched for introns, but this did not reveal any intron sequences that had not been previously detected. The common occurrence of a HE located next to *psbA* suggests that this is not just a chance event, but that there is selective pressure to maintain this arrangement.

### Analysis of the psbA ORF177 (HE) intergenic region in S-PM2

In an effort to understand localization of the HE adjacent to *psbA*, we searched the corresponding region of the S-PM2 genome for elements that may maintain 'selective pressure' on the intergenic space between *psbA* and the HE-encoding gene. We used a bioinformatics approach that predicted the ability of the test sequence to form a stable secondary structure, a characteristic of non-coding RNAs (Backofen and Hess, 2010). We identified a possible transcript within the antisense strand of the S-PM2 genome starting at the 5′ end of ORF177 (F-ChpI) and ending within the 3′ half of *psbA*, close to the intron (Figure 4). This method also predicted a second putative transcript antisense to the 5′ end of *psbA* (Figure 4), although this prediction may merely be a reflection of the highly structured 5′ untranslated region of *psbA* on the sense strand.

### Experimental confirmation of the asRNA

As the bioinformatic analysis strongly suggested the presence of an asRNA in the intergenic region between *psbA* and F-ChpI, 5′ RACE was performed to test these predictions experimentally. 5′ RACE analysis generated two products that mapped to

**Table 2** GOS scaffolds on which homologues of the F-CphI homing endonuclease were detected

| Scaffold accession number | Scaffold size (nt) | *psbA* detected | Intron detected in *psbA* gene | Presence of the 'G/KETTXXXSQ/H' motif in PsbA | Phylogenetic classification of *psbA* | Site of DNA isolation |
|---|---|---|---|---|---|---|
| JCVI_SCAF_1096627879849 | 1071 | ✔ | Fragment too short | ✔ | Not determined | GS031—Upwelling, Fernandina Island |
| JCVI_SCAF_1101667164370 | 1026 | ✔ | x | ✔ | Not determined | GS020—Lake Gatun |
| JCVI_SCAF_1101667171453 | 755 | x | N/A | N/A | Not determined | GS020—Lake Gatun |
| JCVI_SCAF_1101667171626 | 732 | x | N/A | N/A | N/A | GS020—Lake Gatun |
| JCVI_SCAF_1101668541828 | 1657 | ✔ | Fragment too short | Fragment too short | Not determined | GS031—Upwelling, Fernandina Island |
| JCVI_SCAF_1096626190549 | 920 | ✔ | ✔ | Fragment too short | Not determined | GS020—Lake Gatun |
| JCVI_SCAF_1096626856934 | 3100 | ✔ | x | ✔ | Cyanophage | GS003—Browns Bank, Gulf of Maine |
| JCVI_SCAF_1096627283123 | 3436 | ✔ | ✔ | ✔ | Cyanophage | GS002—Gulf of Maine, Canada |
| JCVI_SCAF_1096627676525 | 1603 | x | N/A | N/A | N/A | GS020—Lake Gatun |
| JCVI_SCAF_1101668745121 | 1695 | ✔ | x | ✔ | Cyanophage | GS047—201 miles from F. Polynesia |
| JCVI_SCAF_1097156666624 | 3436 | ✔ | ✔ | ✔ | Cyanophage | GS020—Lake Gatun, Panama |
| JCVI_SCAF_1096627021912 | 1947 | ✔ | ✔ | ✔ | Not determined | GS020—Lake Gatun |
| JCVI_SCAF_1096627284644 | 2516 | ✔ | x | ✔ | Cyanophage | GS002—Gulf of Maine, Canada |
| JCVI_SCAF_1096627299009 | 3055 | x | N/A | N/A | N/A | GS012—Chesapeake Bay, MD |
| JCVI_SCAF_1096627313094 | 1528 | ✔ | Fragment too short | ✔ | Not determined | GS020—Lake Gatun |
| JCVI_SCAF_1096627639930 | 1706 | ✔ | Fragment too short | ✔ | Cyanophage | GS020—Lake Gatun |
| JCVI_SCAF_1096627674162 | 1845 | ✔ | N/A | N/A | N/A | GS020—Lake Gatun |
| JCVI_SCAF_1096627675073 | 1700 | ✔ | Fragment too short | ✔ | Not determined | GS020—Lake Gatun |
| JCVI_SCAF_1096627912725 | 1227 | x | N/A | N/A | N/A | GS031—Upwelling, Fernandina Island |
| JCVI_SCAF_1101668250692 | 1605 | ✔ | x | x | Cyanophage | GS020—Lake Gatun |
| JCVI_SCAF_1101668253187 | 1575 | x | N/A | N/A | N/A | GS020—Lake Gatun |
| JCVI_SCAF_1101668699797 | 1416 | ✔ | Fragment too short | Fragment too short | Not determined | GS035—Wolf Island |
| JCVI_SCAF_1096627024160 | 3075 | ✔ | ✔ | ✔ | Cyanophage | GS020—Lake Gatun |

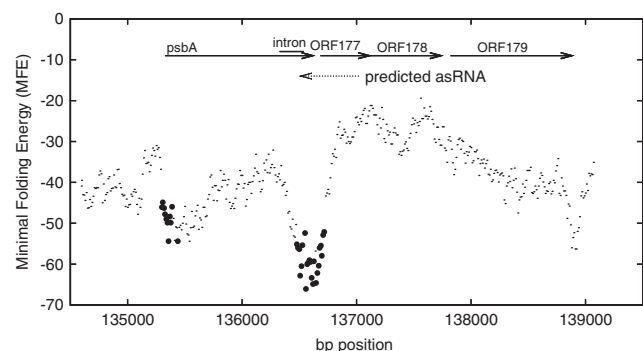Abbreviations: GOS, global ocean survey; N/A, not applicable; nt, nucleotide.



**Figure 4** Prediction of an ncRNA antisense to *psbA* and ORF177. The *psbA* region of S-PM2 analysed in 200-nt windows incrementing every 10 nt. The mean folding energy (MFE) for each was calculated and compared with 1000 scrambles of the same sequence, the MFE for each window is plotted (•) with those windows that had a MFE above the 99th percentile of the 1000 scrambles marked (•). The position of the genes *psbA*, ORF177 (encoding F-CphI), ORF178 and ORF179 (*psbD*) are marked by arrows. The position of the previously predicted ncRNA is marked with a dotted arrow.

positions 136 855 and 136 741 of the S-PM2 genome (Figure 5a). The reason for two RACE products is unclear. A possible explanation is that the transcript is processed to form a shorter product. Northern blotting with a probe specific to this putative asRNA confirmed its expression during the infection process (Figure 5b). A ca. 225-bp product was clearly detected. This fits with the 5′ RACE mapped position of 136 741 and the predicted 3′ terminator site (Figure 5a). Therefore, this experimental evidence confirmed the presence of the predicted asRNA, and we designated this unique element as Cyanophage Functional RNA I.

*Quantitative PCR analysis of CfrI expression*
The expression of CfrI, phage *psbA* and ORF177 (encoding F-ChpI) was monitored using quantitative PCR during the S-PM2 infection cycle. *psbA* expression peaked at 6 h after infection. This peak in expression was also common to CfrI and ORF177
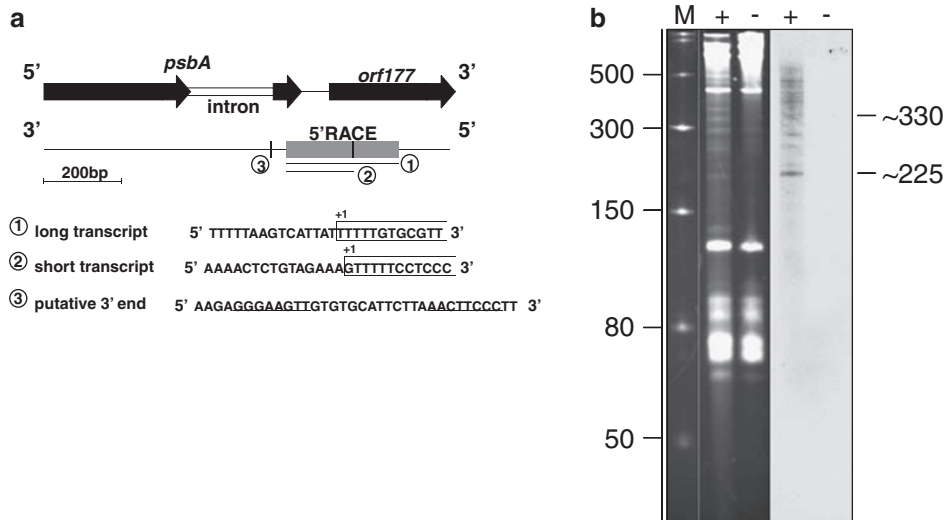
**Figure 5** Presence of an antisense RNA linking the S-PM2 endonuclease gene to the *psbA* second exon. (**a**) Experimentally verified 5′ ends of the antisense transcript overlapping the 3′ end of the intron, exon 2 of *psbA* and the 5′ end of the endonuclease gene ORF177 were mapped to positions 136 855 (long transcript) and 136 741 (short transcript) on the complementary strand. ORF177 was recently identified as a free-standing homing endonuclease gene (HE), targeting intron-less *psbA* genes of marine cyanobacteria. The sequence elements (136 526–136 560, complementary strand) that are predicted to form the terminator helix for the antisense transcript are underlined. (**b**) Separation of 20 μg of total RNA from phage-infected *Synechococcus* sp. WH7803 (+) and from non-infected cells (−) on a 10% (w/v) polyacrylamide gel. Northern hybridization (right) indicates a prominent band of ∼225 bp and some weaker bands of higher molecular weight in the RNA from phage-infected cells but not in the RNA from control cells. The band 225 bp in size corresponds to a transcript with the second mapped 5′ end (short transcript) and the predicted terminator. The blot was hybridized with a single-stranded RNA probe directed against the antisense transcript. An RNA molecular weight standard (M) is shown to the left.
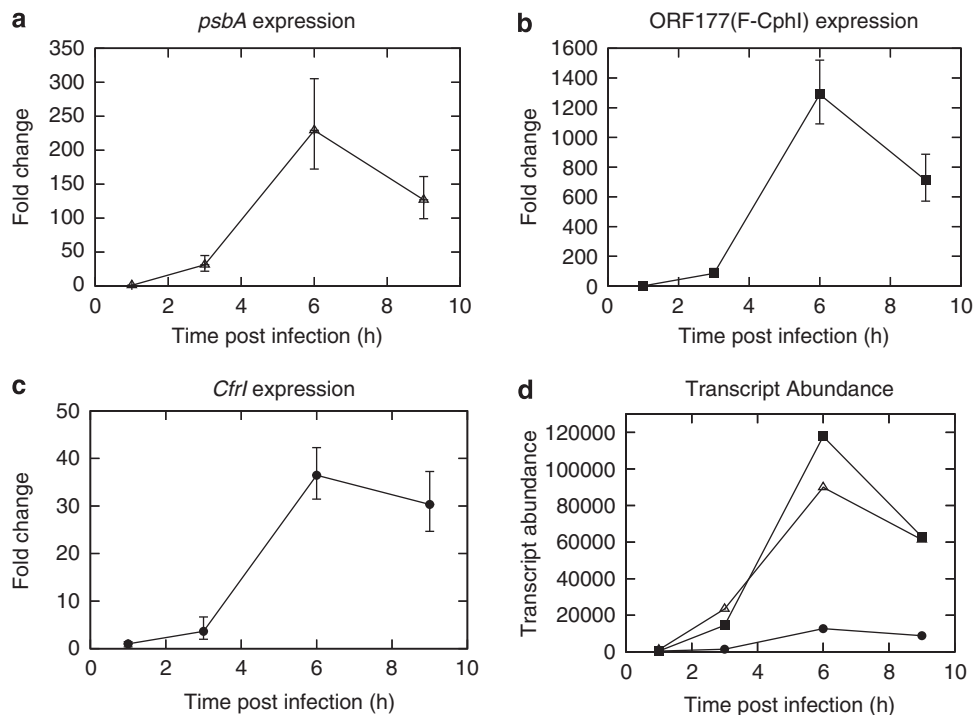


**Figure 6** Expression of *psbA* (▲), CfrI (●) and ORF177 (■). Plotted values are the mean of three independent biological replicates with error bars representing s.d. The relative expression of each transcript is plotted in **a**, **b** and **c** with absolute transcript abundance plotted in **d**. Cells begin to lyse after 9 h under the conditions used (see Clokie *et al.* (2006).

(Figure 6). ORF177 (F-ChpI) showed a large increase in expression between 3 and 6 h (Figure 6), with the absolute number of ORF177 transcripts exceeding those of *psbA* at this time point (Figure 6). CfrI has a temporal expression pattern similar to both ORF177 and *psbA* with a peak at 6 h. However, CfrI transcript

abundance was significantly lower than that of both *psbA* and ORF177 throughout the infection cycle (Figure 6).

*CfrI in other phage genomes*
By aligning the sequence of CfrI identified in S-PM2 with the GOS scaffolds, it was possible to identify CfrI on scaffolds JCVI_SCAF1101667164370, JCVI_SACF_1096627024160, JCVI_SCAF_1096627283123 and JCVI_SCAF_1097156666624 (Supplementary Figure S4). Subsequently, by applying the same bioinformatics approach that predicted CfrI in S-PM2 on a selection of GOS scaffolds, an asRNA was predicted in 9 of the 11 scaffolds tested (Supplementary Table S2). In addition, the *psbA* region of the four other currently sequenced cyano-myoviruses (namely Syn9, S-RMS4, P-SSM2 and P-SSM4) were also analysed for the presence of an asRNA. All four were predicted to encode an asRNA at the 3′ end of the *psbA* gene (Supplementary Table S1). However, unlike the situation in S-PM2, none of the other cyanophages encode a HE downstream of *psbA*, and the asRNA predicted do not appear to overlap the gene downstream of *psb*A.

## Discussion

The self-splicing group I intron in S-PM2 that interrupts *psbA* has been shown to be spliced throughout the infection cycle, presumably to maintain a supply of D1 throughout the infection cycle. In the absence of splicing and excision of the intron, it is assumed that a functional D1 poly-peptide would not be formed. This conclusion is supported by the situation in *C. reinhardtii* in which intron splicing was reduced by directed muta-genesis resulting in the loss of D1 production and consequent reduction in growth rate (Lee and Herrin, 2003). The detection of both spliced and unspliced transcripts during S-PM2 infection suggests that there may be a regulatory role for intron splicing. This would not be without precedent as light has been shown to regulate intron splicing in *C. reinhardtii* (Deshpande *et al.*, 1997).

We have shown that introns with sequence conservation can insert in multiple positions within *psbA* sequences. Why multiple IISs are found is unclear. The best strategy for group I introns to proliferate is to locate into highly conserved DNA sequences that have an essential biological role, which are often encountered in the gene pool, and that are conserved across the biological spectrum (Raghavan and Minnick, 2009). *psb*A fulfils these criteria and therefore provides an ideal 'home' for an intron. The multiple IISs may thus merely be a reflection of the highly conserved nature of this gene, and each site meets the requirements for introns to insert into. The origin of these introns remains unknown. The fact that some introns share IISs with introns found in the chloroplasts of

*C. reinhardtii* and *O. cardiacum* suggests that they may have had a common origin. However, the *psbA* genes they are now located within are all of phage origin. Phylogenetic analysis suggests that they are located within *psbA* genes from phage infecting *Synechococcus* rather than *Prochlorococcus*, with no evidence to suggest that these introns are present in their *Synechococcus* host. This is consistent with sequencing of numerous *Synechococcus* (Dufresne *et al.*, 2008; Scanlan *et al.*, 2009) and *Prochlorococcus* (Kettler *et al.*, 2007) genomes, in which no introns have been identified within *psbA* genes. This is surprising, given the intragenic recombination of *psbA* that has been proposed to occur between cyanophage and their hosts (Zeidner *et al.*, 2005; Sullivan *et al.*, 2006). However, this may be due to the numerical bias in the GOS data set that is dominated by sequences similar to those of *Prochlorococcus* and its infecting phage P-SSM4, with the consequent underrepresentation of *Synechococcus* (Williamson *et al.*, 2008). Previous studies that used only PCR to amplify *psbA* genes have not reported introns within cyanobacterial *psbA* genes, or their phages (Zeidner *et al.*, 2005; Sullivan *et al.*, 2006; Sharon *et al.*, 2007; Wang and Chen, 2008; Marston and Amrich, 2009). This may in part be due to the primers used; the widely used primers described in the study by Zeidner *et al.*, 2005 span the boundary between the two most common IISs and the *psbA*-coding sequence (Figure 2), thereby preventing amplification of any sequences that contain an intron at that particular IIS, and therefore lead to their underrepresentation in *psbA* gene datasets. The more recent primer set described by Wang and Chen (2008) would amplify the most common IISs. However, this primer set has been used to amplify <10 *psbA* genes.

In identifying introns, it became apparent that *psbA* is often localized next to a HE similar to that of F-CphI found in S-PM2. It has recently been suggested that localization of an intron in *psbA* and the presence of a HE adjacent to *psb*A is not accidental. Indeed, it has been proposed as the convergence of two genetic parasites on the same conserved region of DNA (Bonocora and Shub, 2009), with these two independent elements acting in a process of collaborative homing to proliferate within a population (Zeng *et al.*, 2009). In colla-borative homing, the HE targets the IIS as its cutting site, with the intron providing protection against HE nicking its own DNA, and HE providing mobility to transfer into intron-less alleles (Zeng *et al.*, 2009).

Bonocora and Shub (2009) have proposed that the HE will eventually integrate into the intron to form a mobile group I intron, which is the most stable entity as the HE can never be separated from the protective function of the intron. The proposed pathway for the formation of mobile group I introns suggests that both intron-less and intron-containing *psbA* genes adjacent to an F-CphI would have occurred over time. Both of these scenarios were

found in this dataset supporting the proposed model of Bonocora and Shub. However, the final step of integration of F-CphI into an intron was not observed. One intron was found to contain a HE. However, this was significantly different to F-CphI showing similarity to the HE found in the *psbA* intron of *O. cardiacum.* In addition, the gene immediately downstream of *psbA* was similar to the F-CphI found in S-PM2 (Supplementary Figure S3).

The failure to detect F-CphI within an intron-containing *psbA* gene may simply be due to the relatively small sample size of the GOS dataset compared with the total gene pool that is present in the oceans. Alternatively, there might be another selective pressure that has prevented the formation of a truly mobile group I intron within cyanophage *psbA* genes. We found that the expression of *psbA* in S-PM2 is consistent with that of previous reports (Clokie *et al.*, 2006) and fits with its proposed function of maintaining host photosynthetic function during infection (Mann *et al.*, 2003; Lindell *et al.*, 2004, 2005, 2007; Millard *et al.*, 2004; Bragg and Chisholm, 2008; Hellweger, 2009). We also measured expression of ORF177 encoding the HE F-CphI, which was found to be co-expressed with *psbA*. As the spread of the HE into intron-less alleles of *psbA* in other cyanophages is believed to occur during a mixed infection (Zeng *et al.*, 2009), it could be rationalized that the HE would only be expressed once DNA replication has begun, when copies of phage DNA are at their most abundant to provide a substrate for insertion. Thus, the protein would not be required until DNA has become abundant. Previous work has shown that genes involved in S-PM2 DNA replication are expressed 3 h into the infection cycle (Clokie *et al.*, 2006). DNA abundance is thus highest after this point, and before packaging into the protein head. Genes encoding head proteins are not expressed at maximal levels until 6 h and beyond (Clokie *et al.*, 2006). Therefore, DNA is likely to still be abundant and accessible at 6 h, which would explain the large increase in the expression of ORF177, and *psbA*, after 6 h. The identification of the cis-encoded asRNA, CFrI, is unprecedented in a lytic phage. *Cis*-encoded asRNAs have previously been reported in temperate phages, plasmids and bacterial chromosomes (Brantl, 2007), but not in lytic bacteriophages, or cyanophages. The target of an asRNA is often the mRNA that it is complementary to, with post-transcriptional regulation of gene expression being exerted by complementary base pairing (Brantl, 2007). asRNAs that overlap in substantial parts with other genes may be an elegant way to achieve a regulatory connection between neighbouring genes. Indeed, in the cyanobacterium *Anabaena* sp. PCC7120, gene alr1690 has a long 3′ overlap with *furA*, encoding a ferric uptake regulator, and controls the expression level in this manner (Hernandez *et al.*, 2006). It is worth mentioning that bacterial asRNAs not only trigger degradation of their target

mRNAs but also serve as terminators of transcription (Stork *et al.*, 2007) or as signals for RNA processing, triggering the discoordination of operons (Tramonti *et al.*, 2008). In S-PM2, CfrI joins the genetic elements of *psbA* and the gene encoding F-CphI. Presumably, for the gene encoding F-CphI to become integrated into the intron, it has to be removed from its current position. As HEs are normally found within intergenic regions, their inexact removal is unlikely to cause a detrimental effect (Raghavan and Minnick, 2009). However, in S-PM2, the 3′ end of *psbA* and the 5′ end of the gene encoding F-CphI are directly linked by the asRNA CfrI. Therefore, any rearrangement of the HE-encoding gene into the intron of *psbA* will cause disruption of the CfrI sequence, presumably leading to a lack of function. Thus, we propose the asRNA CfrI provides selective pressure to maintain the current *status quo* preventing the formation of a mobile group I intron, as removal of the endonuclease-encoding gene while maintaining the asRNA is likely to be a rare event. Thus, we propose that CfrI has prevented or slowed the evolution of the two genetic parasites of the intron and HE into a single mobile group I intron.

Although the function of CfrI is still unknown, the fact it is differentially expressed during the infection cycle suggests that it has a regulatory role. Given that asRNAs normally regulate the gene they are antisense to by complementary base pairing, it would be reasonable to assume it regulates *psbA* or ORF177 (F-CphI) gene expression. However, the predicted presence of an asRNA at the 3′ end of *psbA* genes in cyanophages that lack a homologue of F-CphI downstream, suggests that this asRNA specifically regulates expression of *psbA*. Given that early and late promoter motifs have already been identified upstream of cyanophage *psbA* (Mann *et al.*, 2005) such additional regulatory capacity may be important for phage infection under particular environmental conditions, for example, high-light intensities. This would be consistent with modelling studies that suggest that phage photosynthesis genes provide an increase in fitness in a manner that is correlated with irradiance (Bragg and Chisholm, 2008; Hellweger, 2009).

The archetypal example of a phage asRNA overlapping the 3′ end of a protein-gene is the 77-nt OOP asRNA of bacteriophage-λ. The OOP asRNA is complementary to the 3′ end of the λ cII-repressor mRNA. Overexpression of OOP asRNA from a plasmid vector results in RNAse III-dependent cleavage of cII mRNA (Krinke and Wulff, 1987). Regulation of the stress-inducible photosynthetic gene *isiA* by the asRNA IsrR in *Synechocystis* sp. PCC6803 is also consistent with this model, as accumulation of mRNA and asRNA follows inverse kinetics and is mutually exclusive (Duhring *et al.*, 2006).

In contrast, the function of CfrI is more likely to be protective as it appears to be co-ordinately expressed with *psbA*. It may act in a similar manner to

some asRNAs observed in *Prochlorococcus*. In *Prochlorococcus* sp. MED4, the asRNA Yfr15, accumulates during phage infection (Steglich *et al.*, 2008). Yfr15 overlaps the 3′ end of gene PMED4_07441 (PMM0686), the most highly upregulated host mRNA during phage infection. In contrast to this high level of expression, the vast majority of host-encoded mRNAs are rapidly degraded (Lindell *et al.*, 2007), implying that Yfr15 protects the PMED4_07441 mRNA, for example, by rendering RNase E recognition sites inaccessible. In this context, it is noteworthy that we detected accumulation of unspliced *psbA* precursor transcripts, indicating slow kinetics of intron splicing. This would imply that there would be a delay before exon 2 of *psbA* would—by physical occlusion by the translating ribosomes—become protected from endonuclease cleavage. This may also explain why a lower stoichiometric ratio of the asRNA relative to the mRNA may be sufficient, a hypothesis that can be tested in future experiments.

It is only recently that the role and importance of asRNAs in cyanobacterial regulation has become apparent (Steglich *et al.*, 2008; Georg *et al.*, 2009). In the cyanobacterium *Prochlorococcus*, which possesses a highly reduced genome, it is believed that *trans*-acting ncRNAs and *cis*-acting asRNAs have an important role in regulation (Steglich *et al.*, 2008). The co-evolution of virus and host, and transfer of genetic material between them, coupled with the relatively limited coding capacity of the phage genome, implies that similar genetically conservative ncRNAs and asRNAs remain to be identified in the genomes of lytic phages.

## Conclusions

The occurrence of introns inserted at multiple positions within cyanophage *psbA* genes appears to be a widespread phenomenon. These intron-containing *psbA* genes are often located adjacent to a gene encoding a HE, seemingly the result of the co-evolution of two genetic parasites on a single conserved sequence. Within cyanophage S-PM2, these two separate genetic elements are 'joined' by an asRNA, CfrI. CfrI is the first example of an asRNA in a lytic bacteriophage. Its co-expression with *psbA* points to a role in regulation. The discovery of sequences similar to CfrI in other cyanophage scaffolds suggests that asRNAs, and perhaps more generally other ncRNAs, are likely to be important in regulating cyanophage gene expression. However, CfrI also has the potentially unique property of preventing or slowing down the evolution of two genetic parasites, an intron and a HE into a single mobile group I intron.

## Acknowledgements

## References

Backofen R, Hess WR. (2010). Computational prediction of sRNAs and their targets in bacteria. *RNA Biol* **7**: 1–10.

Bonocora RP, Shub DA. (2009). A likely pathway for formation of mobile group I introns. *Curr Biol* **19**: 223–228.

Bragg JG, Chisholm SW. (2008). Modeling the fitness consequences of a cyanophage-encoded photosynthesis gene. *PLoS One* **3**: e3550.

Brantl S. (2007). Regulatory mechanisms employed by cis-encoded antisense RNAs. *Curr Opin Microbiol* **10**: 102–109.

Brouard JS, Otis C, Lemieux C, Turmel M. (2008). Chloroplast DNA sequence of the green alga *Oedogonium cardiacum* (*Chlorophyceae*): unique genome architecture, derived characters shared with the Chaetophorales and novel genes acquired through horizontal transfer. *BMC Genomics* **9**: 290.

Chenard C, Suttle CA. (2008). Phylogenetic diversity of sequences of cyanophage photosynthetic gene *psbA* in marine and freshwaters. *Appl Environ Microbiol* **74**: 5317–5324.

Clokie MR, Shan J, Bailey S, Jia Y, Krisch HM, West S *et al.* (2006). Transcription of a 'photosynthetic' T4-type phage during infection of a marine cyanobacterium. *Environ Microbiol* **8**: 827–835.

Deshpande NN, Bao Y, Herrin DL. (1997). Evidence for light/redox-regulated splicing of *psbA* pre-RNAs in *Chlamydomonas* chloroplasts. *RNA* **3**: 37–48.

Dufresne A, Ostrowski M, Scanlan DJ, Garczarek L, Mazard S, Palenik BP *et al.* (2008). Unraveling the genomic mosaic of a ubiquitous genus of marine cyanobacteria. *Genome Biol* **9**: R90.

Duhring U, Axmann IM, Hess WR, Wilde A. (2006). An internal antisense RNA regulates expression of the photosynthesis gene *isiA*. *Proc Natl Acad Sci USA* **103**: 7054–7058.

Georg J, Voss B, Scholz I, Mitschke J, Wilde A, Hess WR. (2009). Evidence for a major role of antisense RNAs in cyanobacterial gene regulation. *Mol Syst Biol* **5**: 305.

Hellweger FL. (2009). Carrying photosynthesis genes increases ecological fitness of cyanophage in silico. *Environ Microbiol* **11**: 1386–1394.

Hernandez JA, Muro-Pastor AM, Flores E, Bes MT, Peleato ML, Fillat MF. (2006). Identification of a *furA* cis antisense RNA in the cyanobacterium *Anabaena* sp. PCC 7120. *J Mol Biol* **355**: 325–334.

Huelsenbeck JP, Ronquist F. (2001). MRBAYES: Bayesian inference of phylogenetic trees. *Bioinformatics* **17**: 754–755.

Itoh T, Tomizawa J. (1980). Formation of an RNA primer for initiation of replication of ColE1 DNA by ribonuclease H. *Proc Natl Acad Sci USA* **77**: 2450–2454.

Jurica MS, Stoddard BL. (1999). Homing endonucleases: structure, function and evolution. *Cell Mol Life Sci* **55**: 1304–1326.

Kettler GC, Martiny AC, Huang K, Zucker J, Coleman ML, Rodrigue S *et al.* (2007). Patterns and implications of

1134

gene gain and loss in the evolution of *Prochlorococcus*. *PLoS Genet* **3**: e231.

Krinke L, Wulff DL. (1987). OOP RNA, produced from multicopy plasmids, inhibits lambda cII gene expression through an RNase III-dependent mechanism. *Genes Dev* **1**: 1005–1013.

Lacatena RM, Cesareni G. (1981). Base pairing of RNA I with its complementary sequence in the primer precursor inhibits ColE1 replication. *Nature* **294**: 623–626.

Lee J, Herrin DL. (2003). Mutagenesis of a light-regulated *psbA* intron reveals the importance of efficient splicing for photosynthetic growth. *Nucleic Acids Res* **31**: 4361–4372.

Lindell D, Jaffe JD, Coleman ML, Futschik ME, Axmann IM, Rector T *et al.* (2007). Genome-wide expression dynamics of a marine virus and host reveal features of co-evolution. *Nature* **449**: 83–86.

Lindell D, Jaffe JD, Johnson ZI, Church GM, Chisholm SW. (2005). Photosynthesis genes in marine viruses yield proteins during host infection. *Nature* **438**: 86–89.

Lindell D, Sullivan MB, Johnson ZI, Tolonen AC, Rohwer F, Chisholm SW. (2004). Transfer of photosynthesis genes to and from *Prochlorococcus* viruses. *Proc Natl Acad Sci USA* **101**: 11013–11018.

Liu Q, Belle A, Shub DA, Belfort M, Edgell DR. (2003). SegG endonuclease promotes marker exclusion and mediates co-conversion from a distant cleavage site. *J Mol Biol* **334**: 13–23.

Livak KJ, Schmittgen TD. (2001). Analysis of relative gene expression data using real-time quantitative PCR and the 2(-Delta Delta C(T)) method. *Methods* **25**: 402–408.

Logemann J, Schell J, Willmitzer L. (1987). Improved method for the isolation of RNA from plant tissues. *Anal Biochem* **163**: 16–20.

Lu J, Chen F, Hodson RE. (2001). Distribution, isolation, host specificity, and diversity of cyanophages infecting marine *Synechococcus* spp. in river estuaries. *Appl Environ Microbiol* **67**: 3285–3290.

Mann NH, Clokie MR, Millard A, Cook A, Wilson WH, Wheatley PJ *et al.* (2005). The genome of S-PM2, a 'photosynthetic' T4-type bacteriophage that infects marine *Synechococcus* strains. *J Bacteriol* **187**: 3188–3200.

Mann NH, Cook A, Millard A, Bailey S, Clokie M. (2003). Marine ecosystems: bacterial photosynthesis genes in a virus. *Nature* **424**: 741.

Markham NR, Zuker M. (2008). UNAFold: software for nucleic acid folding and hybridization. In: Keith JM (ed). *Bioinformatics Volume II. Structure, Functions and Applications*. Humana Press: Totowa, NJ, USA, pp 3–31.

Marston MF, Amrich CG. (2009). Recombination and microdiversity in coastal marine cyanophages. *Environ Microbiol* **11**: 2893–2903.

Marston MF, Sallee JL. (2003). Genetic diversity and temporal variation in the cyanophage community infecting marine *Synechococcus* species in Rhode Island's coastal waters. *Appl Environ Microbiol* **69**: 4639–4647.

Maul JE, Lilly JW, Cui L, dePamphilis CW, Miller W, Harris EH *et al.* (2002). The *Chlamydomonas reinhardtii* plastid chromosome: islands of genes in a sea of repeats. *Plant Cell* **14**: 2659–2679.

Millard A, Clokie MR, Shub DA, Mann NH. (2004). Genetic organization of the *psbAD* region in phages infecting marine *Synechococcus* strains. *Proc Natl Acad Sci USA* **101**: 11007–11012.

Millard A, Mann NH. (2006). A temporal and spatial investigation of cyanophage abundance in the Gulf of Aqaba, Red Sea. *J Mar Biol Assocn UK* **86**: 507–515.

Millard AD, Zwirglmaier K, Downey MJ, Mann NH, Scanlan DJ. (2009). Comparative genomics of marine cyanomyoviruses reveals the widespread occurrence of *Synechococcus* host genes localized to a hyperplastic region: implications for mechanisms of cyanophage evolution. *Environ Microbiol* **11**: 2370–2387.

Raghavan R, Minnick MF. (2009). Group I introns and inteins: disparate origins but convergent parasitic strategies. *J Bacteriol* **191**: 6193–6202.

Rusch DB, Halpern AL, Sutton G, Heidelberg KB, Williamson S, Yooseph S *et al.* (2007). The Sorcerer II Global Ocean Sampling expedition: Northwest Atlantic through eastern tropical Pacific. *PLoS Biol* **5**: e77.

Scanlan DJ, Ostrowski M, Mazard S, Dufresne A, Garczarek L, Hess WR *et al.* (2009). Ecological genomics of marine picocyanobacteria. *Microbiol Mol Biol Rev* **73**: 249–299.

Seshadri R, Kravitz SA, Smarr L, Gilna P, Frazier M. (2007). CAMERA: a community resource for metagenomics. *PLoS Biol* **5**: e75.

Sharon I, Alperovitch A, Rohwer F, Haynes M, Glaser F, Atamna-Ismaeel N *et al.* (2009). Photosystem I gene cassettes are present in marine virus genomes. *Nature* **461**: 258–262.

Sharon I, Tzahor S, Williamson S, Shmoish M, Man-Aharonovich D, Rusch DB *et al.* (2007). Viral photosynthetic reaction center genes and transcripts in the marine environment. *ISME J* **1**: 492–501.

Simmonds P, Smith DB. (1999). Structural constraints on RNA virus evolution. *J Virol* **73**: 5787–5794.

Spiegelman WG, Reichardt LF, Yaniv M, Heinemann SF, Kaiser AD, Eisen H. (1972). Bidirectional transcription and the regulation of Phage lambda repressor synthesis. *Proc Natl Acad Sci USA* **69**: 3156–3160.

Steglich C, Futschik ME, Lindell D, Voss B, Chisholm SW, Hess WR. (2008). The challenge of regulation in a minimal photoautotroph: non-coding RNAs in *Prochlorococcus*. *PLoS Genet* **4**: e1000173.

Stork M, Di Lorenzo M, Welch TJ, Crosa JH. (2007). Transcription termination within the iron transport-biosynthesis operon of *Vibrio anguillarum* requires an antisense RNA. *J Bacteriol* **189**: 3479–3488.

Sullivan MB, Coleman ML, Weigele P, Rohwer F, Chisholm SW. (2005). Three *Prochlorococcus* cyanophage genomes: signature features and ecological interpretations. *PLoS Biol* **3**: e144.

Sullivan MB, Krastins B, Hughes JL, Kelly L, Chase M, Sarracino D *et al.* (2009). The genome and structural proteome of an ocean siphovirus: a new window into the cyanobacterial 'mobilome'. *Environ Microbiol* **11**: 2935–2951.

Sullivan MB, Lindell D, Lee JA, Thompson LR, Bielawski JP, Chisholm SW. (2006). Prevalence and evolution of core photosystem II genes in marine cyanobacterial viruses and their hosts. *PLoS Biol* **4**: e234.

Sullivan MB, Waterbury JB, Chisholm SW. (2003). Cyanophages infecting the oceanic cyanobacterium *Prochlorococcus*. *Nature* **424**: 1047–1051.

Suttle CA. (2005). Viruses in the sea. *Nature* **437**: 356–361.

Suttle CA. (2007). Marine viruses—major players in the global ecosystem. *Nat Rev Microbiol* **5**: 801–812.

Suttle CA, Chan AM. (1994). Dynamics and distribution of cyanophages and their effect on marine *Synechococcus* spp. *Appl Environ Microbiol* **60**: 3167–3174.

Tramonti A, De Canio M, De Biase D. (2008). GadX/GadW-dependent regulation of the *Escherichia coli* acid fitness island: transcriptional control at the *gadY-gadW* divergent promoters and identification of four novel 42 bp GadX/GadW-specific binding sites. *Mol Microbiol* **70**: 965–982.

Wang K, Chen F. (2008). Prevalence of highly host-specific cyanophages in the estuarine environment. *Environ Microbiol* **10**: 300–312.

Waterbury JB, Valois FW. (1993). Resistance to co-occurring phages enables marine *Synechococcus* communities to coexist with cyanophages abundant in seawater. *Appl Environ Microbiol* **59**: 3393–3399.

Weigele PR, Pope WH, Pedulla ML, Houtz JM, Smith AL, Conway JF *et al.* (2007). Genomic and structural analysis of Syn9, a cyanophage infecting marine *Prochlorococcus* and *Synechococcus*. *Environ Microbiol* **9**: 1675–1695.

Williamson SJ, Rusch DB, Yooseph S, Halpern AL, Heidelberg KB, Glass JI *et al.* (2008). The Sorcerer II Global ocean sampling expedition: metagenomic characterization of viruses within aquatic microbial samples. *PLoS One* **3**: e1456.

Wilson WH, Carr NG, Mann NH. (1996). The effect of phosphate status on the kinetics of cyanophage infection in the oceanic cyanobacterium *Synechococcus* sp WH7803. *J Phycol* **32**: 506–516.

Wilson WH, Joint IR, Carr NG, Mann NH. (1993). Isolation and molecular characterization of five marine cyanophages propagated on *Synechococcus* sp. Strain WH7803. *Appl Environ Microbiol* **59**: 3736–3743.

Wyman M, Gregory RP, Carr NG. (1985). Novel role for phycoerythrin in a marine cyanobacterium, *Synechococcus* strain DC2. *Science* **230**: 818–820.

Zeidner G, Bielawski JP, Shmoish M, Scanlan DJ, Sabehi G, Beja O. (2005). Potential photosynthesis gene recombination between *Prochlorococcus* and *Synechococcus* via viral intermediates. *Environ Microbiol* **7**: 1505–1513.

Zeidner G, Preston CM, Delong EF, Massana R, Post AF, Scanlan DJ *et al.* (2003). Molecular diversity among marine picophytoplankton as revealed by *psbA* analyses. *Environ Microbiol* **5**: 212–216.

Zeng Q, Bonocora RP, Shub DA. (2009). A free-standing homing endonuclease targets an intron insertion site in the *psbA* gene of cyanophages. *Curr Biol* **19**: 218–222.

Supplementary Information accompanies the paper on The ISME Journal website (http://www.nature.com/ismej)