

UC Berkeley

UC Berkeley Previously Published Works

Title

Candidate Phyla Radiation Roizmanbacteria From Hot Springs Have Novel and Unexpectedly Abundant CRISPR-Cas Systems

Permalink

<https://escholarship.org/uc/item/6k67r1p8>

Journal

Frontiers in Microbiology, 10(MAY)

ISSN

1664-302X

Authors

Chen, Lin-Xing
Al-Shayeb, Basem
Méheust, Raphaël
et al.

Publication Date

2019

DOI

10.3389/fmicb.2019.00928

Peer reviewed



Candidate Phyla Radiation Roizmanbacteria From Hot Springs Have Novel and Unexpectedly Abundant CRISPR-Cas Systems

Lin-Xing Chen¹, Basem Al-Shayeb², Raphaël Méheust^{1,3}, Wen-Jun Li⁴, Jennifer A. Doudna^{3,5} and Jillian F. Banfield^{1,3*}

¹ Department of Earth and Planetary Sciences, University of California, Berkeley, Berkeley, CA, United States, ² Department of Plant & Microbial Biology, University of California, Berkeley, Berkeley, CA, United States, ³ Innovative Genomics Institute, University of California, Berkeley, Berkeley, CA, United States, ⁴ State Key Laboratory of Biocontrol, Guangdong Provincial Key Laboratory of Plant Resources, School of Life Sciences, Sun Yat-sen University, Guangzhou, China, ⁵ The Department of Chemistry, University of California, Berkeley, Berkeley, CA, United States

OPEN ACCESS

Edited by:

Baolei Jia,
Chung-Ang University, South Korea

Reviewed by:

Gisle Vestergaard,
University of Copenhagen, Denmark
Olga Soutourina,
UMR9198 Institut de Biologie
Intégrative de la Cellule (I2BC), France

*Correspondence:

Jillian F. Banfield
jbanfield@berkeley.edu

Specialty section:

This article was submitted to
Evolutionary and Genomic
Microbiology,
a section of the journal
Frontiers in Microbiology

Received: 07 March 2019

Accepted: 12 April 2019

Published: 03 May 2019

Citation:

Chen L-X, Al-Shayeb B, Méheust R, Li W-J, Doudna JA and Banfield JF (2019) Candidate Phyla Radiation Roizmanbacteria From Hot Springs Have Novel and Unexpectedly Abundant CRISPR-Cas Systems. *Front. Microbiol.* 10:928. doi: 10.3389/fmicb.2019.00928

The Candidate Phyla Radiation (CPR) comprises a huge group of bacteria that have small genomes that rarely encode CRISPR-Cas systems for phage defense. Consequently, questions remain about their mechanisms of phage resistance and the nature of phage that infect them. The compact CRISPR-CasY system (Cas12d) with potential value in genome editing was first discovered in these organisms. Relatively few CasY sequences have been reported to date, and little is known about the function and activity of these systems in the natural environment. Here, we conducted a genome-resolved metagenomic investigation of hot spring microbiomes and recovered CRISPR systems mostly from Roizmanbacteria that involve CasY proteins that are divergent from published sequences. Within population diversity in the spacer set indicates current *in situ* diversification of most of the loci. In addition to CasY, some Roizmanbacteria genomes also encode large type I-B and/or III-A systems that, based on spacer targeting, are used in phage defense. CRISPR targeting identified three phage represented by complete genomes and a prophage, which are the first reported for bacteria of the Microgenomates superphylum. Interestingly, one phage encodes a Cas4-like protein, a scenario that has been suggested to drive acquisition of self-targeting spacers. Consistent with this, the Roizmanbacteria population that it infects has a CRISPR locus that includes self-targeting spacers and a fragmented CasY gene (fCasY). Despite gene fragmentation, the PAM sequence is the same as that of other CasY reported in this study. Fragmentation of CasY may avoid the lethality of self-targeting spacers. However, the spacers may still have some biological role, possibly in genome regulation. The findings expand our understanding of CasY diversity, and more broadly, CRISPR-Cas systems and phage of CPR bacteria.

Keywords: CPR, Roizmanbacteria, CRISPR-Cas, phage, hot spring

INTRODUCTION

The Candidate Phyla Radiation (CPR) comprises a huge fraction of Domain Bacteria. The scale of the radiation remains unclear, but it may include as much as 26–50% of all bacterial diversity (Hug et al., 2016; Parks et al., 2017; Schulz et al., 2017). The CPR bacteria uniformly have small genomes (often ~1 Mbp) and limited biosynthetic capacity (Brown et al., 2015; Anantharaman et al., 2016; Hug et al., 2016; Castelle and Banfield, 2018). Most are thought to be symbionts, in some cases cell surface attached (episymbionts), that depend on other bacteria for basic cellular building blocks (for review, see Castelle and Banfield, 2018).

A previous meta-analysis found that only 2.4% of organisms from the Parcubacteria (OD1) and Microgenomates (OP11) superphyla encode CRISPR-Cas systems in their genomes, as compared to 47.4% in archaea and 24.4% in non-CPR bacteria (Burstein et al., 2016). The authors noted that when CRISPR-Cas systems occur in CPR bacteria they tend to be different from those found in other bacteria. Four genomes from Dojkabacteria (WS6), Parcubacteria (OD1) and Roizmanbacteria were previously recognized to encode CRISPR-Cas12a (Cpf1) systems (Zetsche et al., 2015), and more recently, six genomes were reported encoding a newly recognized compact CasY effector enzyme that has genome editing potential (Burstein et al., 2017).

Several potential explanations for the low frequency of CRISPR-Cas systems in CPR bacteria have been suggested (Burstein et al., 2016). Small genome size may favor use of more compact restriction-modification systems for phage defense and low ribosome content may preclude sufficiently fast-acting CRISPR-Cas systems required for effective interference (Burstein et al., 2016). Symbiotic lifestyles, characterized by close association between multiple cells and a host cell, could lead to higher phage densities, which may cause selection of defense systems other than CRISPR-Cas (Westra et al., 2015). It has also been suggested that CPR bacteria may not have the RecBCD mechanism identified in non-CPR Bacteria to curtail self-targeting spacer acquisition (Levy et al., 2015; Castelle et al., 2018).

As few phage that infect CPR bacteria have been reported (Paez-Espino et al., 2016; Dudek et al., 2017), it is difficult to know how common phage that infect these bacteria might be. Phage particles in the process of infecting CPR bacterial cells have been observed via cryogenic electron microscopy (Luef et al., 2015). However, the sequences of phage associated with CPR bacteria are unusually difficult to identify in metagenomic datasets, in part due to the lack of CRISPR spacers that could be used to link them to host cells via CRISPR targeting (Andersson and Banfield, 2008). Further, like phage, CPR genomes encode a very high proportion of novel proteins (Castelle and Banfield, 2018), which obscures identification of potential prophage regions. Finally, phage structural proteins may be too divergent from those of well-studied phage to be identified. To date, phage have only been reported for bacteria from two CPR phyla, Absconditabacteria (previously SR1) and Saccharibacteria (previously TM7) (Paez-Espino et al., 2016; Dudek et al., 2017). Thus, there is a potentially huge knowledge gap related to the existence and diversity of

CPR phage. This motivates the search for new CPR genomes with CRISPR-Cas systems that could potentially provide links to additional examples of phage that replicate in these bacteria.

In the current study, we investigated the microbiomes of a series of hot springs in Tibet. CPR bacteria are relatively abundant in these thermal environments, and some of their genomes encode interesting and unusual CRISPR-Cas systems. Although uncommon overall, CRISPR-Cas systems are surprisingly frequently encoded in the genomes of members of the Roizmanbacteria, and multiple different systems coexist in some genomes. We identified many new examples of systems based on CasY and uncovered an intriguing example of a locus with self-targeting spacers and a fragmented CasY gene. We identified CPR phage for which complete, curated genomes were reconstructed, as well as prophage in other genomes. Thus, our analyses provide new insights into CPR biology, their phage and the diversity of the relatively unstudied CRISPR-CasY system.

MATERIALS AND METHODS

Study Site, Sampling and Physicochemical Determination

Hot spring (40.8–84.9°C) sediment samples were collected from Tibet Plateau (China) in August 2016 (**Supplementary Table S1**). As described previously (Song et al., 2013), sediment samples were collected from the hot spring pools using a sterile iron spoon into 50 ml sterile tubes, transported to the lab on dry ice, and stored at –80°C for DNA extraction. Temperature, dissolved oxygen (DO) and pH were determined *in situ* and the other physicochemical parameters were analyzed in the laboratory (**Supplementary Table S1**).

DNA Extraction, Sequencing, Quality Control and Metagenomic Assembly

Genomic DNA was extracted from sediment samples using the FastDNA SPIN kit (MP Biomedicals, Irvine, CA, United States) according to the manufacturer's instructions. The DNA samples were purified for library construction, and sequenced on an Illumina HiSeq2500 platform with PE (paired-end) 150 bp kits. The raw data of each metagenomic dataset were filtered to remove Illumina adapters, PhiX and other Illumina trace contaminants with BBTools, and low quality bases and reads using Sickle (version 1.33¹). The high-quality reads of each sample were assembled using metaSPADES (version 3.10.1) (Bankevich et al., 2012) with a kmer set of 21, 33, 55, 77, 99, and 127.

HMM-Based Search of CasY Proteins and Confirmation of CRISPR-CasY System

The six CasY proteins reported previously (Burstein et al., 2017) were aligned using Muscle (Edgar, 2004), and filtered to remove those columns comprising 95% or more gaps with

¹<https://github.com/najoshi/sickle>

TrimAL (Capella-Gutiérrez et al., 2009). A HMM model was built based on the filtered alignment using hmmbuild 2 (Eddy, 1998) with default parameters, hmmsearch was used to search all the proteins predicted by Prodigal from scaffolds. Those hits with an e -value $< 10^{-5}$ were manually checked, and the online tool CRISPRs finder (Grissa et al., 2008) was used to identify the Cas1 protein and CRISPR loci. Only those scaffolds detected with CasY, Cas1, and CRISPR array were retained for further analyses. Other CRISPR-Cas systems identified in these genomes based on the presence of Cas proteins and CRISPR arrays were also analyzed in this study.

Extension and Manual Curation of CasY Scaffolds

Those scaffolds with partial CasY representatives were manually extended as follows: (1) mapping the high quality reads to the corresponding scaffolds using bowtie2 with default parameters; (2) filtering the mapping files using mapped.py (part of the ra2 suite) to remove those PE reads with two or more mismatches to the assembled scaffold across both reads combined; (3) importing the filtered mapping files into Geneious and mapping using the “Map to Reference” function; (4) extending the scaffolds at the partial CasY protein ends; (5) performing the first 4 steps again (multiple times if necessary) until full length CasY proteins were obtained.

The extended scaffolds and other full-length CasY scaffolds were checked for any potential assembly errors using ra2.py², the general strategy was described previously (Brown et al., 2015). Errors reported as unresolved by ra2.py were fixed manually in Geneious using unplaced paired reads that were mapped to the scaffolding gaps.

Coverage Calculation, Genome Binning, Genome Curation and Completeness Assessment

The high quality reads were mapped to the corresponding assembled scaffolds using bowtie2 with default parameters and the coverage of each scaffold calculated as the total number of bases mapped to it divided by its length. For each sample, scaffolds over 2,500 bp were assigned to preliminary draft genome bins using MetaBAT with default parameters, considering both tetranucleotide frequencies (TNF) and scaffold coverage information. The clustering of scaffolds from the bins and the unbinned scaffolds was visualized using ESOM with a min length of 2,500 bp and max length of 5,000 bp as previously described (Dick et al., 2009). Misplaced scaffolds were removed from bins and unbinned scaffolds whose segments were placed within the bin areas of ESOMs were added to the bins. Scaffolds $\geq 1,000$ bp from each sample were uploaded to ggKbase³. The ESOM-curated bins with interesting CasY-bearing scaffolds were further evaluated based on consistency of GC content, coverage and taxonomic information, and scaffolds identified as contaminants were removed. The genome bins with CRISPR-CasY systems

were curated individually to fix local assembly errors using ra2.py, as described above. A total of 50 single copy genes (SCGs) that are commonly detected in CPR bacteria (Supplementary Table S2) were used to evaluate genome completeness.

Gene Prediction and Metabolic Prediction

The protein-coding genes of the curated genomes (see above) were predicted using Prodigal (-m single) (Hyatt et al., 2010), and searched against KEGG, UniRef100, and UniProt for annotation, and metabolic pathways were reconstructed. The 16S rRNA genes were predicted based on HMM models, as previously described (Brown et al., 2015). The ribosome binding site sequence was obtained via the Prodigal gene prediction results.

CRISPR Loci Reconstruction and Spacer Identification

For all the confirmed CRISPR-CasY and other CRISPR-Cas systems, the quality reads were aligned to the scaffolds from the corresponding sample using bowtie2 with default parameters (Brown et al., 2015; Langmead and Salzberg, 2012). Any unmapped reads of read pairs were mapped to the scaffolds in Geneious using the function of “Map to Reference,” then the CRISPR loci were manually reconstructed, allowing for spacer set diversification and loss of spacer-repeat units in some cells. Thus, it was possible to place most reads in an order that reflects the locus evolutionary history. For each CRISPR locus, all the reads that mapped were extracted, and spacers between two direct repeats were used for target searches (see below).

Spacers Target Search and Identification of (Pro)phage Scaffolds

All the spacer sequences from each CRISPR locus were dereplicated, then the sequences were searched against scaffolds from related samples using BLASTn with the following parameters: -task blastn-short, -dust no, -word_size 8. Those scaffolds with 0 mismatch and 100% alignment coverage to one or more spacers were manually checked for phage-specific proteins, including capsid, phage, virus, prophage, terminase, prohead, tape measure, tail, head, portal, DNA packaging, as described previously (Dudek et al., 2017).

In silico Determination of Protospacer Adjacent Motif (PAM)

To determine the PAM of the CRISPR-CasY systems in Roizmanbacteria genomes, for each CRISPR spacer with a target in two complete phage genomes from QZM (see section “Results”), the upstream 5 bp and downstream 5 bp of the targeted DNA strand were searched manually and the PAM was determined and visualized using Weblogo (Crooks et al., 2004). The PAM analyses for other CRISPR-Cas systems analyzed in this study were performed in the same way.

Phylogenetic Analyses

Phylogenetic analyses were performed using (1) 16 ribosomal proteins (16 RPs) and (2) 16S rRNA genes of genomes of interest

²https://github.com/christophertbrown/fix_assembly_errors/releases/tag/2.00

³<http://ggkbase.berkeley.edu/>

with CRISPR-CasY and/or other CRISPR systems (**Table 1**), (3) CasY proteins, (4) Cpf1 proteins, and (5) capsid proteins of CPR (pro)phage:

- (1) 16 RPs analyses: After preliminary classification based on the ribosomal protein S3 (rpS3) taxonomy, reference genomes were downloaded from NCBI (131 in total) and dereplicated using dRep (“-sa 0.95 -nc 0.5”) (Olm et al., 2017). A higher similarity threshold was used to perform dereplication of newly reconstructed genomes from hot spring sediment samples (“-sa 0.99 -nc 0.5”), to clarify the overall diversity. The 16 RPs (i.e., L2, L3, L4, L5, L6, L14, L15, L16, L18, L22, L24, S3, S8, S10, S17, and S19) were predicted from all the dereplicated genomes.
- (2) 16S rRNA genes sequences: The 16S rRNA genes were predicted from all the dereplicated genomes (see above) using HMM-based searches (Brown et al., 2015). All the insertion sequences with lengths > 10 bp were removed.
- (3) CasY proteins: all partial and full length CasY proteins from confirmed CRISPR-CasY systems in this study and the previously reported CasY proteins were included in a phylogenetic tree, with c2c3 proteins as the outgroup.
- (4) Cas12a (Cpf1) proteins: the Cas12a proteins in NCBI and our dataset were identified and used to construct a tree with Cas12c (C2c3) proteins as the outgroup.
- (5) CPR (pro)phage: the capsid protein was used as a marker to build phylogenetic trees for CPR (pro)phage. The capsid proteins identified in this study were searched against the NCBI RefSeq Phage Capsid proteins, the first 5 blast hits were used as reference proteins, along with those in previously reported in CPR phage genomes (Paez-Espino et al., 2016; Dudek et al., 2017).

For tree construction, protein sequences datasets were aligned using Muscle (Edgar, 2004). The 16S rRNA gene sequences were aligned using the SINA alignment algorithm (Edgar, 2004; Pruesse et al., 2012) through the SILVA web interface (Pruesse et al., 2007). All the alignments were filtered using TrimAL (Capella-Gutiérrez et al., 2009) to remove those columns comprising more than 95% gaps. For the 16 RP, ambiguously aligned C and N termini were removed and the amino acid sequences, which were concatenated in the order as stated above (alignment length, 2654 aa). The phylogenetic trees were reconstructed using RAXML version 8.0.26 with the following options: -m PROTGAMMALG (GTRGAMMAI for 16S rRNA phylogeny) -c 4 -e 0.001 -# 100 -f a (Capella-Gutiérrez et al., 2009; Stamatakis, 2014). All the trees were uploaded to iTOL v3 for visualization and formatting (Letunic and Bork, 2006).

Data Availability

The reconstructed CPR (representative genome of each group) and their infecting phage genomes reported in the current study were deposited at NCBI within BioProject PRJNA493250 (BioSample SUB4567369), under the accession numbers of SAMN10133524-SAMN10133631. The CPR and phage genomes, and also those four unbinned scaffolds detected with CasY systems can be explored and downloaded from

ggKbase⁴ following publication of this manuscript. Note that registration by provision of an email address is required prior to data download.

RESULTS

Newly Reconstructed Roizmanbacteria and Woesebacteria Genomes With CRISPR-Cas Systems

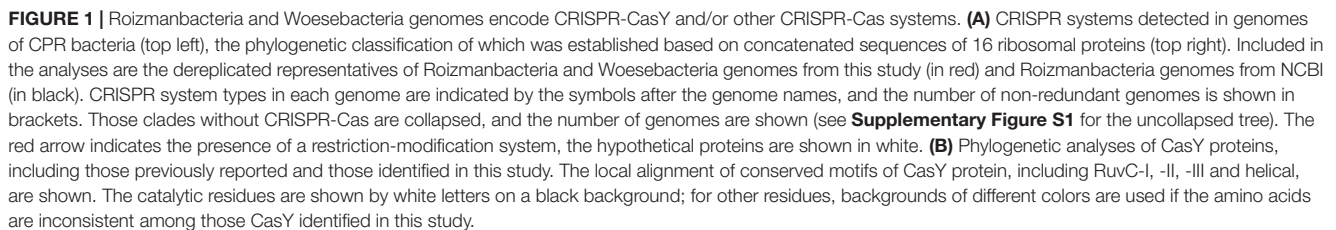
Candidate Phyla Radiation bacteria comprised up to 43.1% of the hot spring communities (9.5% on average) (**Supplementary Figure S1A** and **Supplementary Table S1**). Generally, we found that the samples from higher temperature hot springs had lower relative abundance of CPR (**Supplementary Figure S1B**). We selected 17 genomes that encode CRISPR-Cas systems for curation (**Figure 1A** and **Table 1**). Based on rpS3 protein taxonomic analysis, one Woesebacteria genome and 12 Roizmanbacteria genomes encode CRISPR-CasY systems, and four other Roizmanbacteria genomes encode only type III-A CRISPR-Cas systems. Both of these phylum-level groups place within the Microgenomates (OP11) (Brown et al., 2015; Hug et al., 2016). Phylogenetic analyses based on 16 RPs with published Roizmanbacteria genomes (43 dereplicated in total) indicated the divergence of the newly reconstructed Roizmanbacteria from previously published genomes (**Figure 1A**). The new Roizmanbacteria genomes were assigned to two distinct classes based on their 16S rRNA gene sequences (Yarza et al., 2014) and/or average nucleotide identity (ANI) (**Figure 1A** and **Supplementary Figure S2**). Five of the genomes represent two different strains, with an ANI of 98.39% (clade 1; **Figure 1A**), and the other 11 genomes belong to the same family (clade 2; **Figure 1A**). Genomes in clade 1 and 2 were assigned to groups (**Figure 1A** and **Table 1**).

CRISPR-CasY Detected in Roizmanbacteria and Woesebacteria Genomes

We identified 69 CasY candidates (see methods), 17 of which are on scaffolds with a Cas1 protein and CRISPR locus (**Supplementary Table S3**). Of these, 12 scaffolds could be assigned to Roizmanbacteria genomes and one to a Woesebacteria genome (**Table 1**). The other four scaffolds with CRISPR-CasY systems could not be binned, but were also included in our analyses (**Figure 1B**).

The CRISPR-CasY systems from Roizmanbacteria and Woesebacteria have a different architecture than those reported previously (Burstein et al., 2017), with CasY and Cas1 proteins on the same side of the CRISPR locus (**Figure 1A**). The Roizmanbacteria CasY proteins have similar lengths of 1,252–1,256 aa, whereas that found in the Woesebacteria is 1,304 aa (**Supplementary Table S3**), comparable to lengths of previously reported CasY [1,153–1,287 aa; (Burstein et al., 2017)]. Phylogenetic analyses of CasY proteins showed that the newly

⁴https://ggkbase.berkeley.edu/Tibet_CRISPR_CasY/organisms



CasY is an effector protein of Type V CRISPR-Cas systems. To date, all reported Type V CRISPR-Cas systems have RuvC-like nuclease domains (Burstein et al., 2017; Chen and Doudna, 2017). Comparative analyses of all CasY proteins reported in this study and CasY.1 with Cpf1, C2c1, and C2c3 references (Shmakov et al., 2015), identified all the catalytic residues within the three conserved motifs of RuvC-I, RuvC-II, and RuvC-III (**Figure 1B**), suggesting the RuvC domains in the new CasY proteins are active nucleases. On the other hand, we detected divergence in other regions of the CasY proteins from different sampling sites (**Figure 1B**).

A Type III-A system was detected in all 11 clade 2 Roizmanbacteria genomes, seven of which encode more than one type of system (**Figure 1A**, **Table 1**, and **Supplementary**

A Type I-B system was identified in two *Roizmanbacteria* genomes belonging to the same genus (C2-Gp5 and C2-Gp7), but not in the C2-Gp6 genomes, despite the fact that C2-Gp5 and C2-Gp7 are very closely related to C2-Gp6 (ANI = 99% and 16S similarity = 98.9%). Comparative genomic analyses showed that the Type I-B system is located between genes encoding a secreted

TABLE 1 | Summary of Roizmanbacteria and Woesebacteria genomes reconstructed in this study.

Clade	Group	Genome name	Genome size (kbp)	No. of scaffolds	CG%	Completeness		Bacterial 50 SCGs	No. of proteins	CRISPR-Cas systems	Prophage
						Completeness	Contamination				
Clade 1	1 (Cl-Gp1)	QZM_B4_Woesebacteria_36_36*	621.6	41	35.66	50/50	0		706	CasY	No
		GD2_1_Roizmanbacteria_31_27*	561.5	98	30.77	41/50	1/50		589	CasY	No
		QZM_A2_2_Roizmanbacteria_31_61*	895.7	74	31.00	43/50	0		775	CasY	No
		QZM_A1_Roizmanbacteria_31_22	753.1	62	30.72	39/50	1/50		678		
Clade2	3 (C2-Gp3)	QZM_A2_1_Roizmanbacteria_31_42	775.1	213	30.87	45/50	0		897		
		QZM_A2_3_Roizmanbacteria_31_19	500.3	126	30.74	27/50	2/50		588		
		QZM_B3_Roizmanbacteria_33_70*	994.5	17	33.41	41/50	0		985	Type III-A (degenerate)	Yes
		QZM_A1_Roizmanbacteria_33_14	822.0	125	33.30	40/50	1/50		978		
	4 (C2-Gp4)	QZM_A2_Roizmanbacteria_33_14	860.2	109	33.24	45/50	2/50		1020		
		QZM_A2_2_Roizmanbacteria_33_18	711.6	184	32.85	35/50	3/50		939		
		DGJ11_Roizmanbacteria_31_22*	655.0	139	31.23	40/50	3/50		779	CasY + Type III I-A	No
		GD2_I_Roizmanbacteria_32_73*	906.7	20	32.28	45/50	2/50		946	CasY + Type III-A + Type I-B	No
	6 (C2-Gp6)	QZM_BI_Roizmanbacteria_33_36*	816.6	133	33.15	43/50	2/50		948		No
		QZM_A1_Roizmanbacteria_33_28	647.1	82	33.16	42/50	3/50		740	CasY + Type III-A	
	7 (C2-Gp7)	QZM_A2_Roizmanbacteria_33_40	645.2	208	33.37	38/50	4/50		835		
		QZM_A2_2_Roizmanbacteria_33_54	792.5	150	32.83	36/50	3/50		947		
		QZM_B4_Roizmanbacteria_33_372*	891.6	13	33.22	46/50	0		919	CasY + Type III-A + Type I-B	No

*Representative genome of each group used in phylogenetic analyses. Figure 1 and Supplementary Figure S1 provide phylogeny and clade information.

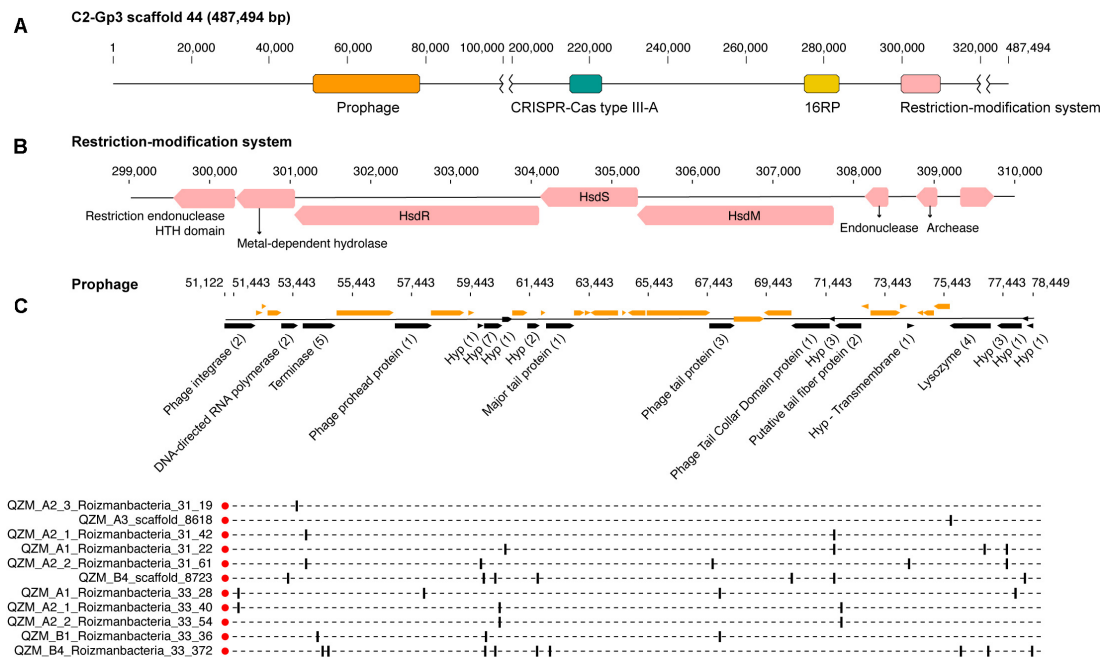


FIGURE 2 | Prophage and restriction modification systems are detected in the genomes of Roizmanbacteria C2-Gp3. **(A)** Scaffold 44 includes prophage, a restriction-modification system and an apparently degenerate Type III-A CRISPR-Cas system. **(B)** The proteins in the restriction-modification system shown in **(A)**. **(C)** Prophage genes targeted by spacers are shown in black and the number of spacers targeting each open reading frame (ORF) is listed in brackets following the annotation. Genes not targeted by CRISPR spacers are shown in orange (top panel). The genome affiliations of spacers targeting the prophage are indicated in the bottom panel.

cysteine-rich protein and a lamin tail domain protein that are present in both genomes (**Supplementary Figure S3B**). Two very short hypothetical proteins were detected between the cysteine-rich and lamin tail domain proteins in the C2-Gp6 genomes (**Supplementary Figure S3B**). However, NCBI BLAST and HMM searches indicate no homology of the hypothetical proteins to any known proteins or functional domains, respectively, and no significant similarity to the Cas proteins of Type I-B systems in the C2-Gp5 and C2-Gp7 genomes.

CRISPR-Cas12a Systems in Published Roizmanbacteria Genomes

We investigated 131 published Roizmanbacteria genomes available from NCBI to identify all CRISPR-Cas systems that occur in these bacteria (**Supplementary Table S4**). The CRISPR-Cas12a system (Cpf1), which was identified in one Roizmanbacteria genome (Zetsche et al., 2015), occurred in four Roizmanbacteria genomes from two classes (**Figure 1A** and **Supplementary Figures S2, S4**), one of them in the class containing Roizmanbacteria clade II with type I-B and III-A systems (see above).

Interestingly, the CRISPR-Cas12a systems reported in Zetsche et al. (2015) from *Candidatus Roizmanbacteria bacterium CG_4_9_14_0_2_um_filter_39_13* included two Cas12a proteins. We refer to the one near the CRISPR locus as Cas12a, and the other as Cas12a' (**Figure 1A**). Phylogenetic analyses of Cas12a and Cas12a' proteins (previously reported

and identified in this study) indicated those in CPR genomes could be assigned into at least three groups (**Supplementary Figure S4A**). Group 1 includes the Cas12a proteins from the two genomes with both Cas12a and Cas12a', and is highly divergent from other Cas12a proteins. Group 2 includes the Cas12a' of *Candidatus Roizmanbacteria bacterium CG_4_9_14_0_2_um_filter_39_13*, along with the Cas12a proteins from another two genomes. Group 3 includes Cas12a' of *Candidatus Roizmanbacteria bacterium GW2011_GWA2_37_7* (Zetsche et al., 2015) and clusters together with Cas12a from non-CPR Bacteria and Archaea.

The RuvC domains (-I, -II, and -III) of the CPR Cas12 and Cas12' group 2 and 3 proteins include all the conserved catalytic residues in **Supplementary Figure S4B**. However, in group 1 proteins, the conserved RuvC-II glutamic acid catalytic residue "E" was substituted by asparagine "N," and in RuvC-III asparagine "N" was substituted to valine "V." These substitutions suggest that the Cas12a in the systems with both Cas12a and Cas12a' may not perform cleavage as documented previously (Zetsche et al., 2015).

Roizmanbacteria-Infecting Phage From Podoviridae and Siphoviridae

A total of 1,118 spacers perfectly targeted (100% match and 100% alignment coverage; see section "Materials and Methods") 565 unique scaffolds. Of these, 156 were targeted by two or more spacers (153 from the QZM samples of the current

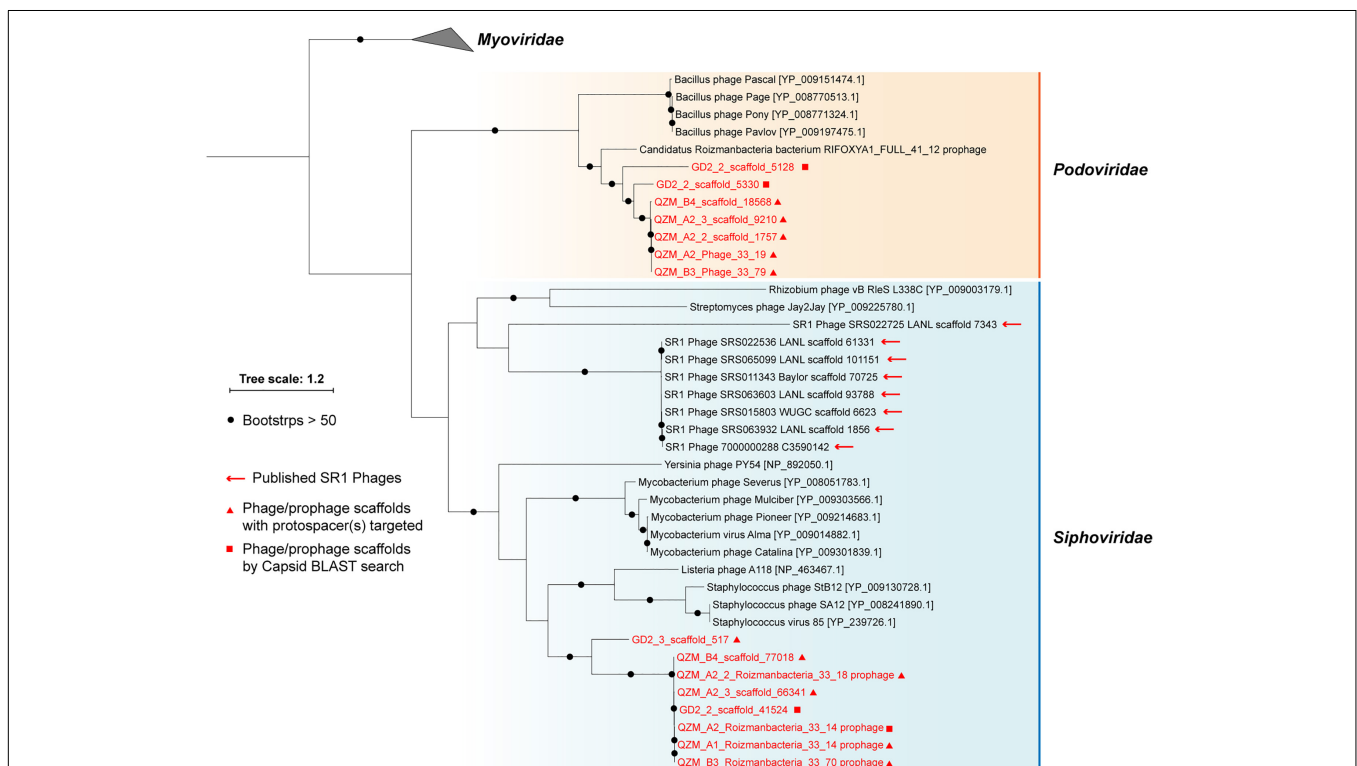
study) (**Supplementary Table S5**). Eleven of the CRISPR spacer-targeted scaffolds encode a phage capsid protein, which was used as a marker for phylogenetic analyses (**Figure 3**). Five additional scaffolds encoding a similar capsid protein were identified by a BLAST search. Capsid proteins were also predicted from the Absconditabacteria (SR1) phage (8 out of 17 with capsid genes identified) and included in the phylogenetic analyses. The (pro)phage identified in this study as well as the Absconditabacteria phage were assigned to either the *Podoviridae* clade or the *Siphoviridae* clade (**Figure 3**). The complete Saccharibacteria phage that lacks an identifiable capsid protein (Dudek et al., 2017) is most closely related to *Siphoviridae* phage based on comparison of its terminase with annotated sequences in the NCBI database.

One scaffold (QZM_B3_scaffold_44) from a C2-Gp3 Roizmanbacteria was targeted by multiple spacers. Detailed analyses indicate that this region is a prophage, with a length of approximately 27 kbp (**Figure 2C**), and is among the first prophage reported in CPR bacterial genomes. This prophage is predicted to encode 40 protein coding genes, including a phage integrase, terminase, prohead protein, major tail protein, tail tape measure protein, tail fiber protein and lysozyme. Nineteen of the ORFs were targeted by 41 CRISPR spacers from CasY-based systems, all of which were from Roizmanbacteria (**Figures 1, 4B**). BLAST comparison detected highly similar scaffolds in the other three genomes of the C2-Gp3 group (**Table 1** and **Figure 2**)

and also unbinned scaffolds in QZM_A2_1, QZM_A2_3, and QZM_A3, suggesting that this is a common Roizmanbacteria prophage. However, when reads of other QZM-related samples were mapped to QZM_B3_scaffold_44, the prophage region showed much higher coverage in QZM_B1 and QZM_B4 than the flanking region (**Supplementary Figure S5**). Further, a subset of reads that circularize the phage genome were detected. These observations indicate that the prophage existed as phage particles in these two samples.

One putative phage scaffold (**Supplementary Table S5**) could be circularized, and circularization of the genome was confirmed by paired-end read mapping. The length of complete phage genome QZM_A2_Phage_33_19 is 31,813 bp, with a GC content of 32.9% (**Figure 4A**). Another two scaffolds (**Supplementary Table S5**) were manually curated to generate another complete phage genome QZM_B3_Phage_33_79, with a length of 30,824 bp and GC content of 32.5% (**Figure 4B**). Phage QZM_A2_Phage_33_19 and QZM_B3_Phage_33_79 share high sequence similarity, and they are probably closely related strains.

A total of 53 and 52 open reading frames (ORFs) were predicted from QZM_A2_Phage_33_19 and QZM_B3_Phage_33_79, respectively (**Figures 4A,B** and **Supplementary Figure S6**). Of these, 46 shared an average amino acid identity of 98%. ORFs common to both phage encode capsid, terminase, lysozyme and tail proteins. Although these two genomes are highly similar, 7 and 6



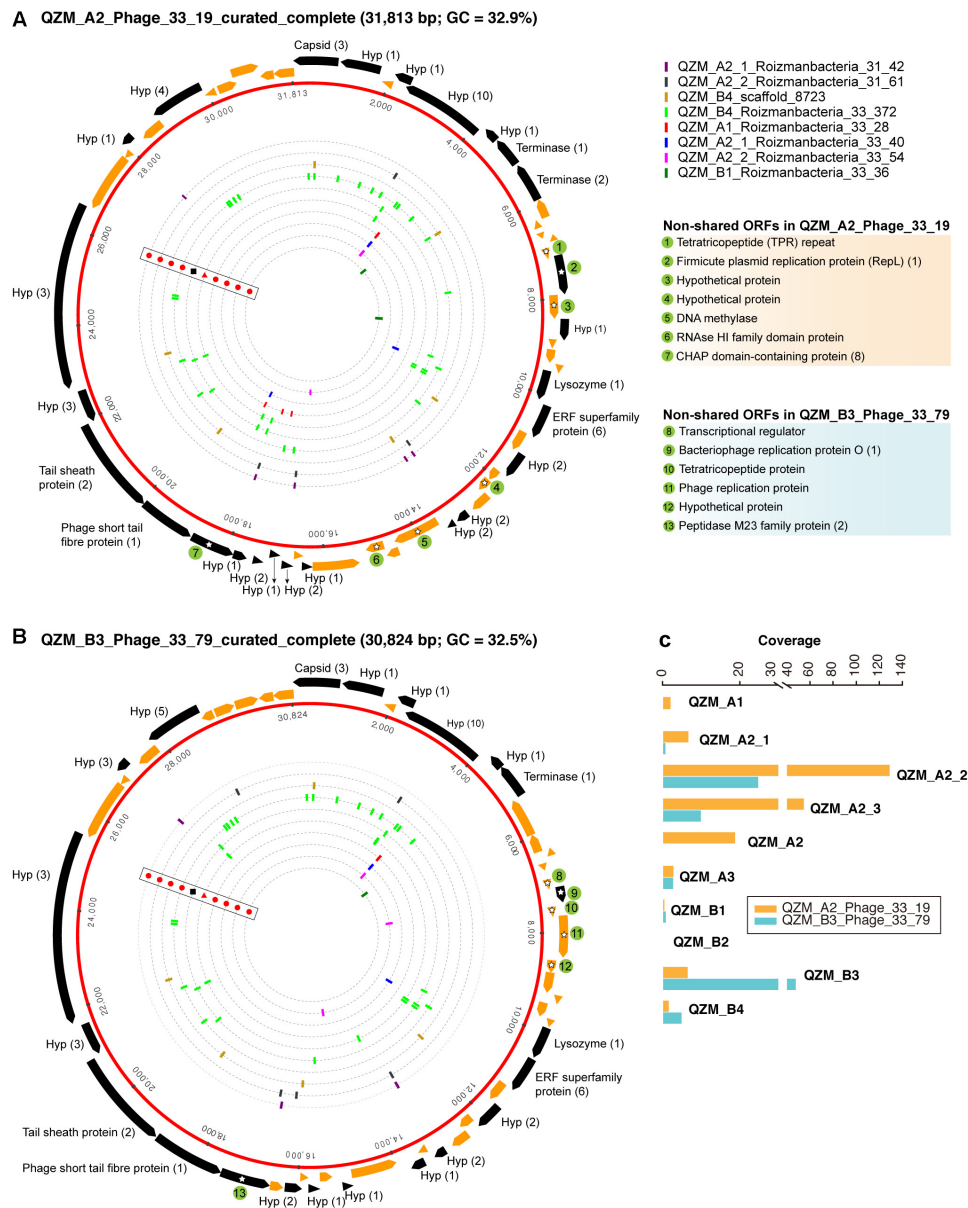


FIGURE 4 | Complete genomes of Roizmanbacteria-infecting phage. The red rings represent (A) QZM_A2_Phage_33_19 and (B) QZM_B3_Phage_33_79 phage genomes. The ORFs are shown outside the genomes, those targeted by at least one spacer are in black (genes not targeted are in orange). The total number of spacers that target each gene is listed in parentheses following the protein annotation. The spacers targeting the phage genome from a given CRISPR-Cas system are indicated by bars on the dotted inner rings (see Figure 1 for CRISPR-Cas system type). Bars are colored by genome of origin (see top right). The non-shared proteins between these two phage genomes are indicated by green circles and numbered, their annotations are shown at the right. Hyp, hypothetical protein. (C) The coverage information of these two phage genomes in QZM-related samples.

non-shared ORFs were detected in QZM_A2_Phage_33_19 and QZM_B3_Phage_33_79, respectively. Among those 13 non-shared ORFs, 4 are related to phage replication (Figures 4A,B), including one replication protein in QZM_A2_Phage_33_19, and one transcriptional regulator and two replication proteins in QZM_B3_Phage_33_79. We used the divergent region between the two genomes to calculate the coverage of the phage in all QZM-related samples and found that they co-occur in most samples (Figure 4C).

A total of 63 spacers targeted 26 ORFs in QZM_A2_Phage_33_19, and 52 spacers targeted 22 ORFs in QZM_B3_Phage_33_79 (Figures 4A,B), but no spacers targeted the intergenic regions of the two phage genomes. The majority of spacers with targets (49 and 39, respectively) were from the CRISPR-CasY systems in QZM Roizmanbacteria genomes, and all the other targeting spacers were from the Type I-B and III-A systems of C2-Gp7. No spacer from QZM_B4_Woesebacteria_36_36 (78 unique spacers)

and the other 10 type III-A systems targeted the two complete phage genomes.

For phylogenetic analyses, we searched the NCBI database for capsid proteins similar to those in the genomes reported here (Figure 3) and identified a scaffold containing a similar capsid ORF that was binned into a Roizmanbacteria genome [Candidatus Roizmanbacteria bacterium RIFOXYA1_FULL_41_12; (Anantharaman et al., 2016)] (Supplementary Figure S2). Comparative analyses showed a close relationship between the sequences of this prophage and the two complete phage mentioned above, including homologies for the capsid and two terminase proteins. In addition, these genes and several other hypothetical proteins share gene arrangements (Supplementary Figure S6). Thus, we conclude the two phage genomes reported are the full sequences for lysogenic (temperate) phage found in Roizmanbacteria genomes.

An Unusual CRISPR-CasY System With a Fragmented CasY Effector and Self-Targeting Spacers

Among the candidate CasY sequences from the Tibet hot springs predicted protein dataset were three adjacent partial proteins on a scaffold from sample GD2_1. In combination, the three ORFs appear to comprise a fragmented CasY protein (defined as “fCasY”). We identified Cas1 and a CRISPR locus adjacent to the fCasY (Figure 5A). Read mapping to the scaffold revealed that the CasY was fragmented by two mutations. One involves deletion of A (from “AAAAA” to “AAAA”) and introduces a TAA stop codon five amino acids downstream. This mutation occurred in all the mapped reads, indicating that all the cells have CasY fragmented at this position. The second mutation is a single nucleotide substitution from “C” to “T,” which introduces a TAA stop codon. This mutation was detected in 82% of the mapped reads. Interestingly, however, the three conserved motifs (RuvC-I, -II, and -III) are preserved in the largest protein fragment and all the catalytic residues are shared with functional CasY proteins (Figures 1B, 5A). We identified the ribosome binding site (RBS) for fragments 1 and 2 as TAA, the same RBS associated with 353 of 946 ORFs of this Roizmanbacteria genome. The longest fragment is predicted to have an RBS of AAT, which was only shared by 55 ORFs.

The fCasY locus includes 22 unique spacers, six of which were detected only once in the mapped reads (Figure 5B). We reconstructed the CRISPR locus (Figure 5B) and found that all of the single copy spacers are at the locus end that is closest to the Cas1 protein. As in prior studies, we infer that these were recently added to the diversifying end of the CRISPR locus in a subset of cells. Interestingly, 12 out of the 22 unique spacers target the scaffolds of the C2-Gp5 genome, which encodes the fCasY system (Figure 5C and Supplementary Table S6). In detail, 11 spacers targeted Roizmanbacteria genes, including those encoding a PINc domain ribonuclease, two permeases, a sigma-70 RNA polymerase and three hypothetical proteins with transmembrane domains. Only one spacer matched an intergenic region, which is next to two tRNAs (His and Thr). This spacer was recently

acquired, as it is encoded on three reads that also sampled part of the leader sequence (Figure 5B and Supplementary Table S6). Several of the self-targeting spacers are located in the old end of the locus (Figure 5B) and occurred in majority of the cells in the population. Thus, we infer that Roizmanbacteria with these self-targeting spacers have survived for a substantial period of time.

In addition to the fCasY locus, we identified type III-A and I-B CRISPR-Cas systems in the C2-Gp5 genome. Notably, one spacer from the type III-A and I-B systems and two fCasY spacers target a complete 34,706 bp phage genome GD2_3_Phage_34_19 (Supplementary Figure S7 and Supplementary Table S6) assigned to *Podoviridae*. A Cas4-like protein was detected in this phage genome (Supplementary Figure S7). As phage with Cas4-like proteins can induce their hosts to acquire self-targeting spacers (Hooton and Connerton, 2014), the presence of this protein may explain acquisition of self-targeting spacers by the C2-Gp5 genome.

Spacers from the loci of C2-Gp5 target other putative phage scaffolds (Supplementary Table S5). For example, one fCasY spacer targets GD2_3_scaffold_2486, which encodes a putative phage gene. Spacers from both fCasY and I-B systems target GD2_2_scaffold_18083, which encodes a phage tail tape measure protein. Two spacers from the type I-B system target GD2_3_scaffold_517, which encodes a capsid protein that is distantly related to that in the prophage of C2-Gp3 (Figure 3).

PAMs 5'-TA and 5'-TG Are Shared by CasY and fCasY Systems

The PAM is used for the acquisition of spacers into the CRISPR array and is important for target recognition and cleavage (Hille et al., 2018). We determined the probable PAM of the CasY systems reported here to target the two complete phage genomes (QZM_A2_Phage_33_19 and QZM_B3_Phage_33_79). Among all the 39 unique target locations on these two phage genomes (88 spacers in total), 20 had a potential 5' TA PAM and 14 had a potential 5' TG PAM (Supplementary Figure S8 and Supplementary Table S6). Moreover, the one spacer in the CRISPR-CasY system of the C1-Gp1 genome that targets GD2_3_Phage_34_19 also has a 5' TA PAM (Supplementary Figure S7). Previously, the PAM determined for the CasY.1 of *Candidatus* Katanobacteria using an *in vitro* approach was a 5' TA, and both 5' TA (dominant) and 5' TG PAMs occur, based on *in vivo* data (Burstein et al., 2017). For the fCasY, we checked to see if the self-targeting spacers have the same PAM as that of other CasY proteins. If this was not the case, the genomic region matching the spacer may not be recognized as a target by the fCasY CRISPR system. Among the 12 self-targeting spacers, 7 have 5' TA and 4 have 5' TG PAMs and one has a possible 5' AT PAM (Supplementary Table S6). Among the 5 fCasY spacers targets on phage scaffolds, two have 5' TA PAMs and two have 5' TG PAMs.

In combination the results indicate that both general CasY proteins and fCasY in this study use the 5' TA/TG PAM sequences for spacer acquisition and protospacer recognition. We identified a few targets with other PAM sequences (Supplementary Table S6), but it is possible that these targets

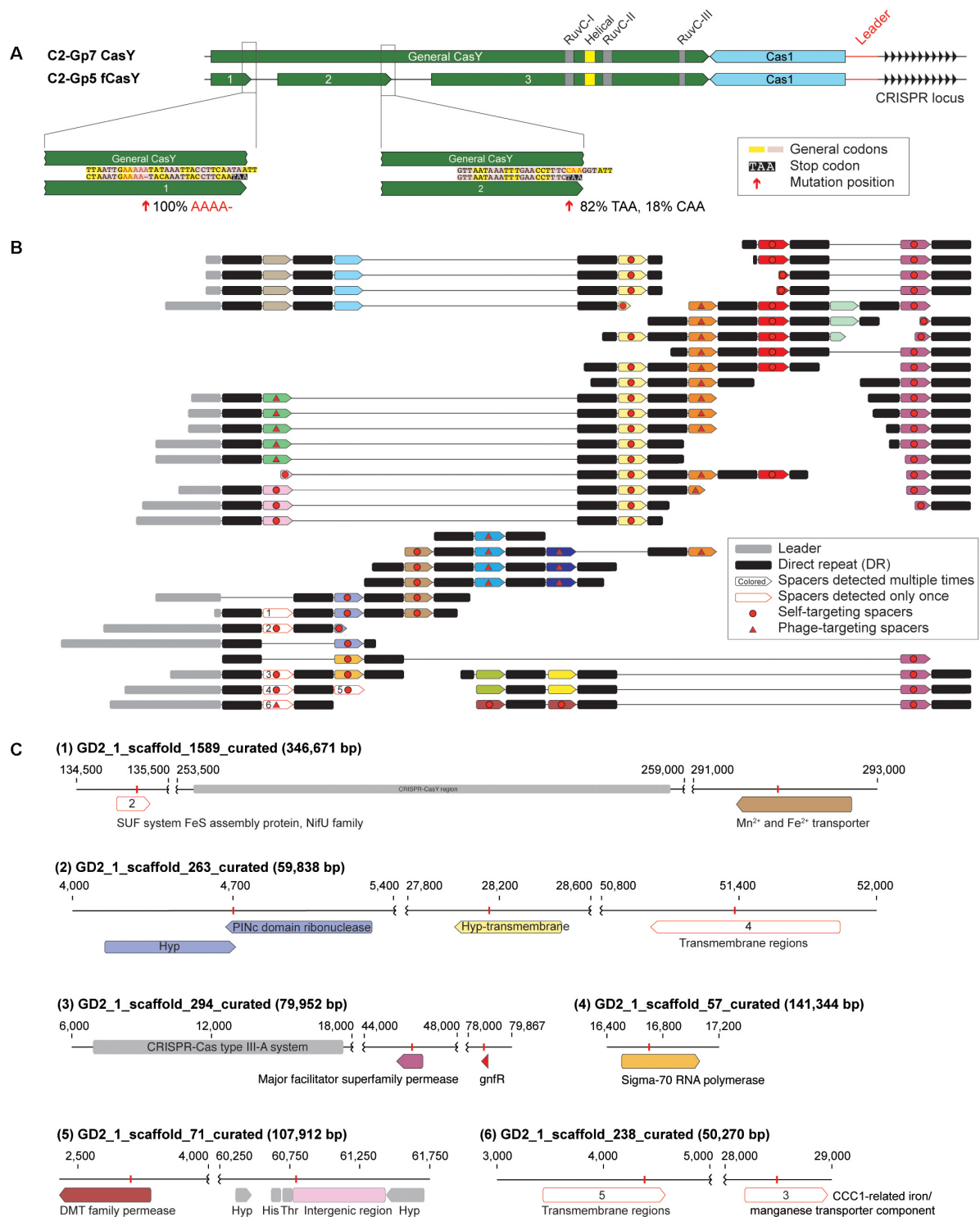


FIGURE 5 | One Roizmanbacteria genome encodes an unusual CRISPR system with a fragmented CasY (fCasY) protein and self-targeting spacers. **(A)** Mutations leading to fragmentation of CasY proteins into three pieces (red arrows) and their incidence in the population, and other features of the locus. **(B)** The reconstructed CRISPR locus showing the history of spacer acquisition and the distribution of self-targeted spacers (marked by red circles). **(C)** Scaffolds encoding genes and an intergenic region matching the self-targeting spacers. The targeted genes have the same color as the corresponding spacers in **(B)**, genes targeted by single copy spacers [white in **(B)**] are indicated by numbers, and CRISPR-Cas systems, tRNA and other genes on the scaffolds are shown in gray.

have mutated the PAM sites during their evolutionary history, as previously documented (Paez-Espino et al., 2015).

Potential Phage–Host Genetic Interactions

When examining the genomic context of CRISPR-CasY systems we noted four very short genes located next to the CRISPR array in the C2-Gp7 genome (**Supplementary Figure S9**). All four genes had at least one homolog in the three complete phage and one prophage (BLASTp *e*-value thresholds = $1e-5$) and when two or more homologs were identified in the same genome, they were together. However, homologs were not identified in the other newly reconstructed and previously reported Roizmanbacteria genomes (**Supplementary Table S4**). The four genes in the C2-Gp7 genome and phage and prophage shared >83% (up to 99%) nucleotide identity with >80% alignment coverage, but none had a NCBI blast hit with similarity >38% (>50 alignment coverage). Given this, and the deduction that QZM_A2_Phage_33_19 and QZM_B3_Phage_33_79 infect C2-Gp7 Roizmanbacteria (based on CRISPR spacer targeting), we conclude that there may have been lateral transfer of novel proteins related to phage–host interactions between Roizmanbacteria and their phage.

DISCUSSION

Candidate Phyla Radiation bacteria account for a huge amount of diversity within the Bacterial domain, but the mechanisms of their interactions with phage and the phage that infect them have remained largely undocumented. In part, this is due to scant information about their CRISPR-Cas systems, despite extensive genomic sampling from a wide variety of sites in nature (Burstein et al., 2016, 2017; Dudek et al., 2017; Castelle and Banfield, 2018). In this study, we report an unexpected diversity of CRISPR-Cas systems in the genomes of bacteria from the CPR phylum of Roizmanbacteria, both from newly reconstructed sequences from multiple hot spring sediments of Tibet, China (**Supplementary Table S1**) and some previously published genomes. Most of them are CasY-based systems (**Figure 1A** and **Table 1**). These new sequences constrain more and less highly conserved regions of CasY proteins, information that may be important in future efforts directed at tailoring the properties of genome-editing enzymes.

The finding that some of the Roizmanbacteria genomes encode multiple CRISPR-Cas systems, including the relatively large types I-B and III-A, is unexpected, given the overall paucity of systems in CPR bacteria, and their small genome sizes (**Figure 1B**). We infer that these systems are mostly active, given the identification of targets on potential phage scaffolds and evidence for locus diversification. Considering that majority of the spacers with targets on the three complete phage and one prophage were from CRISPR-CasY systems (**Figures 2C, 4A,B**), it seems that CasY is the primary CRISPR-Cas system used by these bacteria for phage defense. In the case of the Roizmanbacteria with only a degenerate Type III-A CRISPR-Cas system, defense may rely upon a restriction-modification

system, as suggested previously for CPR bacteria that lack any CRISPR-Cas system (Burstein et al., 2016) (**Figure 2**). In support of this correlation, restriction-modification systems were not detected in those Roizmanbacteria with seemingly functional CRISPR-Cas systems (**Figure 1A**). The discovery of two copies Cas12a proteins in a single system of two genomes is an additional case of unexpected investment in CRISPR-Cas-based phage defense by CPR bacteria (**Figure 1A** and **Supplementary Figure S2**). Overall, the genomes of Roizmanbacteria contained three of the six types of CRISPR-Cas systems reported so far (i.e., types I, III, and V), expanding our understanding of the investment of CPR bacteria in CRISPR-Cas-based defense.

The availability of a pool of CRISPR spacers enabled discovery of three Roizmanbacteria-infecting phage for which complete genomes were reconstructed, and one prophage (**Figures 2–4** and **Supplementary Figure S7**). These are the first reported phage infecting members of the Microgenomates superphylum of the CPR. All of these phage, along with the previously reported CPR phage, were assigned to *Podoviridae* and *Siphoviridae* of the *Caudovirales* order (**Figure 3**). The phylogenetic relatedness and genetic similarity among the *Podoviridae* phage obtained in this study and a Roizmanbacteria prophage deposited at NCBI (**Figure 3** and **Supplementary Figure S6**), and also the potential phage–host genetic interactions (**Supplementary Figure S9**), may indicate stable and similar host–phage relationships in a variety of habitats.

An interesting aspect of the CRISPR-CasY analyses was the fCasY system in one Roizmanbacteria that includes a locus with self-targeting spacers. It may be significant that a Cas4-like protein is encoded in the genome of a phage that replicates in this Roizmanbacteria, given that a Cas4-like protein in a *Campylobacter* sp. phage was suggested to facilitate acquisition of self-targeting spacers into the CRISPR-Cas system of its host (Hooton and Connerton, 2014). Roizmanbacteria lack the RecBCD mediated double-stranded DNA break repair complex, the only documented mechanism for avoidance of self-targeting spacer acquisition (Levy et al., 2015). Thus, it is plausible that the phage-encoded Cas4-like protein led to acquisition of the self-targeting spacers, which should result in autoimmunity (Stern et al., 2010).

Autoimmunity can be avoided via loss of cas genes, mutated repeats adjacent to self-targeting spacers, extended base-pairing with the upstream flanking repeat, and the absence of a PAM in the chromosomal region matched by the spacer (Stern et al., 2010), none of which were observed here. Autoimmunity also could be countered via loss of cas gene function. Interestingly, the fCasY harbored conserved RuvC domains and catalytic residues found in intact CasY proteins (**Figures 1B, 5**). However, given the relatively high abundance of Roizmanbacteria with fCasY in the community (1.37%), we infer that the fCasY protein fragmentation led to loss of cleavage function, preventing autoimmunity. It is possible that the region of the fCasY protein responsible for binding to the target sequence is encoded on a different gene fragment than that encoding the nuclease domain, so that the CRISPR RNA does not recruit the protein fragment with nuclease function.

The presence of old end CRISPR locus spacers that target the host chromosome suggests that the fCasY has been present in the genomes of the Roizmanbacteria C2-Gp5 population for some time. Why has this gene, or the entire locus, not been lost? It is possible that the spacers of the fCasY locus retain some function, for example in gene regulation (possibly involving binding of CRISPR RNAs to the DNA during transcription). Experiments will be required to determine whether fragments of fCasY can reassemble and bind to the genomic regions targeted by the self-targeting spacers (without cleavage) and to determine if the spacer-directed binding domain is on fragment 1 or 2 (Figure 5A).

CONCLUSION

CRISPR-Cas systems are unexpectedly common in a subset of CPR bacteria, and the number, variety and potential functional diversity of these systems is greater than expected. It is already established that CRISPR-CasY systems from these intriguing and enigmatic bacteria will have biotechnological value. Lessons from natural system studies such as reported here may provide information about CasY sequence variety and function that may be useful in enzyme engineering. Beyond this, the new information about CPR bacteria, their phage and the mechanisms of their interactions expands our understanding of the complex phenomena that shape the structure and functioning of natural microbial communities.

DATA AVAILABILITY

The datasets generated for this study can be found in NCBI, PRJNA493250.

REFERENCES

- Anantharaman, K., Brown, C. T., Hug, L. A., Sharon, I., Castelle, C. J., Probst, A. J., et al. (2016). Thousands of microbial genomes shed light on interconnected biogeochemical processes in an aquifer system. *Nat. Commun.* 7:13219. doi: 10.1038/ncomms13219
- Andersson, A. F., and Banfield, J. F. (2008). Virus population dynamics and acquired virus resistance in natural microbial communities. *Science* 320, 1047–1050. doi: 10.1126/science.1157358
- Bankevich, A., Nurk, S., Antipov, D., Gurevich, A. A., Dvorkin, M., Kulikov, A. S., et al. (2012). SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *J. Comput. Biol.* 19, 455–477. doi: 10.1089/cmb.2012.0021
- Brown, C. T., Hug, L. A., Thomas, B. C., Sharon, I., Castelle, C. J., Singh, A., et al. (2015). Unusual biology across a group comprising more than 15% of domain Bacteria. *Nature* 523, 208–211. doi: 10.1038/nature14486
- Burstein, D., Harrington, L. B., Strutt, S. C., Probst, A. J., Anantharaman, K., Thomas, B. C., et al. (2017). New CRISPR-Cas systems from uncultivated microbes. *Nature* 542, 237–241. doi: 10.1038/nature21059
- Burstein, D., Sun, C. L., Brown, C. T., Sharon, I., Anantharaman, K., Probst, A. J., et al. (2016). Major bacterial lineages are essentially devoid of CRISPR-Cas viral defence systems. *Nat. Commun.* 7:10613. doi: 10.1038/ncomms10613
- Capella-Gutiérrez, S., Silla-Martínez, J. M., and Gabaldón, T. (2009). trimAl: a tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* 25, 1972–1973. doi: 10.1093/bioinformatics/btp348

AUTHOR CONTRIBUTIONS

L-XC and JB designed the study. W-JL supported the metagenomic sequencing. L-XC performed the metagenomic assembly, HMM search, and scaffold extension and curation. L-XC and JB performed genome binning and curation. L-XC and JB conducted data analyses with input from BA-S, RM, and JD. L-XC and JB wrote the manuscript. All authors read and approved the final manuscript.

FUNDING

This research was supported by the Microbiology Program of the Innovative Genomics Institute. W-JL was financially supported by the Science and Technology Infrastructure work project (No. 2015FY110100), and the Natural Science Foundation of Guangdong Province, China (No. 2016A030312003). BA-S was supported by the National Science Foundation Graduate Research Fellowship (DGE 1752814).

ACKNOWLEDGMENTS

We thank the two reviewers for their helpful comments. This manuscript has been released as a preprint at bioRxiv (Chen et al., 2018).

SUPPLEMENTARY MATERIAL

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fmicb.2019.00928/full#supplementary-material>

- Castelle, C. J., and Banfield, J. F. (2018). Major new microbial groups expand diversity and alter our understanding of the tree of life. *Cell* 172, 1181–1197. doi: 10.1016/j.cell.2018.02.016
- Castelle, C. J., Brown, C. T., Anantharaman, K., Probst, A. J., Huang, R. H., and Banfield, J. F. (2018). Biosynthetic capacity, metabolic variety and unusual biology in the CPR and DPANN radiations. *Nat. Rev. Microbiol.* 16, 629–645. doi: 10.1038/s41579-018-0076-2
- Chen, J. S., and Doudna, J. A. (2017). The chemistry of Cas9 and its CRISPR colleagues. *Nat. Rev. Chem.* 1:0078. doi: 10.1038/s41570-017-0078
- Chen, L. X., Al-Shayeb, B., Meheust, R., Li, W. J., Doudna, J. A., and Banfield, J. F. (2018). Candidate phyla radiation roizmanbacteria from hot springs have novel, unexpectedly abundant, and potentially alternatively functioning CRISPR-Cas systems. *bioRxiv* [Preprint]. doi: 10.1101/448639
- Crooks, G. E., Hon, G., Chandonia, J. M., and Brenner, S. E. (2004). WebLogo: a sequence logo generator. *Genome Res.* 14, 1188–1190. doi: 10.1101/gr.849004
- Dick, G. J., Andersson, A. F., Baker, B. J., Simmons, S. L., Thomas, B. C., Yelton, A. P., et al. (2009). Community-wide analysis of microbial genome sequence signatures. *Genome Biol.* 10:R85. doi: 10.1186/gb-2009-10-8-r85
- Dudek, N. K., Sun, C. L., Burstein, D., Kantor, R. S., Aliaga Goltsman, D. S., Bik, E. M., et al. (2017). Novel microbial diversity and functional potential in the marine mammal oral microbiome. *Curr. Biol.* 27, 3752–3762.e6. doi: 10.1016/j.cub.2017.10.040
- Eddy, S. R. (1998). Profile hidden markov models. *Bioinformatics* 14, 755–763. doi: 10.1093/bioinformatics/14.9.755

- Edgar, R. C. (2004). MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* 32, 1792–1797. doi: 10.1093/nar/gkh340
- Grissa, I., Bouchon, P., Pourcel, C., and Vergnaud, G. (2008). On-line resources for bacterial micro-evolution studies using MLVA or CRISPR typing. *Biochimie* 90, 660–668. doi: 10.1016/j.biochi.2007.07.014
- Hille, F., Richter, H., Wong, S. P., Bratović, M., Ressel, S., and Charpentier, E. (2018). The biology of CRISPR-Cas: backward and forward. *Cell* 172, 1239–1259. doi: 10.1016/j.cell.2017.11.032
- Hooton, S. P. T., and Connerton, I. F. (2014). *Campylobacter jejuni* acquire new host-derived CRISPR spacers when in association with bacteriophages harboring a CRISPR-like Cas4 protein. *Front. Microbiol.* 5:744. doi: 10.3389/fmicb.2014.00744
- Hug, L. A., Baker, B. J., Anantharaman, K., Brown, C. T., Probst, A. J., Castelle, C. J., et al. (2016). A new view of the tree of life. *Nat. Microbiol.* 1:16048. doi: 10.1038/nmicrobiol.2016.48
- Hyatt, D., Chen, G.-L., Locascio, P. F., Land, M. L., Larimer, F. W., and Hauser, L. J. (2010). Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics* 11:119. doi: 10.1186/1471-2105-11-119
- Langmead, B., and Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. *Nat. Methods* 9, 357–359. doi: 10.1038/nmeth.1923
- Letunic, I., and Bork, P. (2006). Interactive Tree Of Life (iTOL): an online tool for phylogenetic tree display and annotation. *Bioinformatics* 23, 127–128. doi: 10.1093/bioinformatics/bt1529
- Levy, A., Goren, M. G., Yosef, I., Auster, O., Manor, M., Amitai, G., et al. (2015). CRISPR adaptation biases explain preference for acquisition of foreign DNA. *Nature* 520, 505–510. doi: 10.1038/nature14302
- Luef, B., Frischkorn, K. R., Wrighton, K. C., Holman, H.-Y. N., Birarda, G., Thomas, B. C., et al. (2015). Diverse uncultivated ultra-small bacterial cells in groundwater. *Nat. Commun.* 6:6372. doi: 10.1038/ncomms7372
- Núñez, J. K., Kranzusch, P. J., Noeske, J., Wright, A. V., Davies, C. W., and Doudna, J. A. (2014). Cas1–Cas2 complex formation mediates spacer acquisition during CRISPR–Cas adaptive immunity. *Nat. Struct. Mol. Biol.* 21, 528–534. doi: 10.1038/nsmb.2820
- Olm, M. R., Brown, C. T., Brooks, B., and Banfield, J. F. (2017). dRep: a tool for fast and accurate genomic comparisons that enables improved genome recovery from metagenomes through de-replication. *ISME J.* 11, 2864–2868. doi: 10.1038/ismej.2017.126
- Paez-Espino, D., Elie-Fadrosh, E. A., Pavlopoulos, G. A., Thomas, A. D., Huntemann, M., Mikhailova, N., et al. (2016). Uncovering Earth's virome. *Nature* 536, 425–430. doi: 10.1038/nature19094
- Paez-Espino, D., Sharon, I., Morovic, W., Stahl, B., Thomas, B. C., Barrangou, R., et al. (2015). CRISPR immunity drives rapid phage genome evolution in *Streptococcus thermophilus*. *mBio* 6:e262–15. doi: 10.1128/mBio.00262-15
- Parks, D. H., Rinke, C., Chuvochina, M., Chaumeil, P.-A., Woodcroft, B. J., Evans, P. N., et al. (2017). Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life. *Nat. Microbiol.* 2, 1533–1542. doi: 10.1038/s41564-017-0012-7
- Pruesse, E., Peplies, J., and Glöckner, F. O. (2012). SINA: accurate high-throughput multiple sequence alignment of ribosomal RNA genes. *Bioinformatics* 28, 1823–1829. doi: 10.1093/bioinformatics/bts252
- Pruesse, E., Quast, C., Knittel, K., Fuchs, B. M., Ludwig, W., Peplies, J., et al. (2007). SILVA: a comprehensive online resource for quality checked and aligned ribosomal RNA sequence data compatible with ARB. *Nucleic Acids Res.* 35, 7188–7196. doi: 10.1093/nar/gkm864
- Schulz, F., Elie-Fadrosh, E. A., Bowers, R. M., Jarett, J., Nielsen, T., Ivanova, N. N., et al. (2017). Towards a balanced view of the bacterial tree of life. *Microbiome* 5:140. doi: 10.1186/s40168-017-0360-9
- Shmakov, S., Abudayyeh, O. O., Makarova, K. S., Wolf, Y. I., Gootenberg, J. S., Semenova, E., et al. (2015). Discovery and functional characterization of diverse class 2 CRISPR-Cas systems. *Mol. Cell* 60, 385–397. doi: 10.1016/j.molcel.2015.10.008
- Song, Z.-Q., Wang, F.-P., Zhi, X.-Y., Chen, J.-Q., Zhou, E.-M., Liang, F., et al. (2013). Bacterial and archaeal diversities in yunnan and tibetan hot springs, China. *Environ. Microbiol.* 15, 1160–1175. doi: 10.1111/1462-2920.12025
- Stamatakis, A. (2014). RAXML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* 30, 1312–1313. doi: 10.1093/bioinformatics/btu033
- Stern, A., Keren, L., Wurtzel, O., Amitai, G., and Sorek, R. (2010). Self-targeting by CRISPR: gene regulation or autoimmunity? *Trends Genet.* 26, 335–340. doi: 10.1016/j.tig.2010.05.008
- Westra, E. R., van Houte, S., Oyesiku-Blakemore, S., Makin, B., Broniewski, J. M., Best, A., et al. (2015). Parasite exposure drives selective evolution of constitutive versus inducible defense. *Curr. Biol.* 25, 1043–1049. doi: 10.1016/j.cub.2015.01.065
- Yarza, P., Yilmaz, P., Pruesse, E., Glöckner, F. O., Ludwig, W., Schleifer, K.-H., et al. (2014). Uniting the classification of cultured and uncultured bacteria and archaea using 16S rRNA gene sequences. *Nat. Rev. Microbiol.* 12, 635–645. doi: 10.1038/nrmicro3330
- Zetsche, B., Gootenberg, J. S., Abudayyeh, O. O., Slaymaker, I. M., Makarova, K. S., Essletzbichler, P., et al. (2015). Cpf1 is a single RNA-guided endonuclease of a class 2 CRISPR-Cas system. *Cell* 163, 759–771. doi: 10.1016/j.cell.2015.09.038

Conflict of Interest Statement: JB is a founder of Metagenomi. JD is a co-founder of Caribou Biosciences, Inc., Editas Medicine, Intellia Therapeutics, Scribe Therapeutics, and Mammoth Biosciences. JD is a scientific advisory board member of Caribou Biosciences, Inc., Intellia Therapeutics, eFFECTOR Therapeutics, Scribe Therapeutics, Synthego, Mammoth Biosciences, and Inari. JD is a member of the board of directors at Driver and Johnson & Johnson and has sponsored research projects by Roche Biopharma and Biogen.

The remaining authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2019 Chen, Al-Shayeb, Méheust, Li, Doudna and Banfield. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.