

# Type IV CRISPR–Cas systems are highly diverse and involved in competition between plasmids

Rafael Pinilla-Redondo<sup>1,2,\*</sup>, David Mayo-Muñoz<sup>1,†</sup>, Jakob Russel<sup>1,†</sup>, Roger A. Garrett<sup>3</sup>, Lennart Randau<sup>4</sup>, Søren J. Sørensen<sup>1,\*</sup> and Shiraz A. Shah<sup>5,\*</sup>

<sup>1</sup>Section of Microbiology, University of Copenhagen, Universitetsparken 15, 2100 Copenhagen, Denmark,

<sup>2</sup>Department of Technological Educations, University College Copenhagen, Sigurdsgade 26, 2200 Copenhagen,

Denmark, <sup>3</sup>Danish Archaea Centre, Department of Biology, University of Copenhagen, University of Copenhagen,

Ole Maaløes Vej 5, 2200 Copenhagen, Denmark, <sup>4</sup>Philipps-Universität Marburg, Faculty of Biology,

Hans-Meerwein-Straße 6, 35032 Marburg, Germany and <sup>5</sup>Copenhagen Prospective Studies on Asthma in Childhood (COPSAC), Herlev and Gentofte Hospital, University of Copenhagen, Ledreborg Alle 34, 2820 Gentofte, Denmark

Received November 05, 2019; Revised December 02, 2019; Editorial Decision December 07, 2019; Accepted December 13, 2019

## ABSTRACT

**CRISPR–Cas systems provide prokaryotes with adaptive immune functions against viruses and other genetic parasites. In contrast to all other types of CRISPR–Cas systems, type IV has remained largely overlooked. Here, we describe a previously uncharted diversity of type IV gene cassettes, primarily encoded by plasmid-like elements from diverse prokaryotic taxa. Remarkably, via a comprehensive analysis of their CRISPR spacer content, these systems were found to exhibit a strong bias towards the targeting of other plasmids. Our data indicate that the functions of type IV systems have diverged from those of other host-related CRISPR–Cas immune systems to adopt a role in mediating conflicts between plasmids. Furthermore, we find evidence for cross-talk between certain type IV and type I CRISPR–Cas systems that co-exist intracellularly, thus providing a simple answer to the enigmatic absence of type IV adaptation modules. Collectively, our results lead to the expansion and reclassification of type IV systems and provide novel insights into the biological function and evolution of these elusive systems.**

## INTRODUCTION

Clustered regularly interspaced short palindromic repeats (CRISPR), together with their CRISPR-associated (Cas) genes, constitute a diverse family of nucleic acid-based adaptive immune systems that protect archaea and bacteria

against invading mobile genetic elements (MGEs). These defence systems are classified by virtue of their modular composition and structure, into two major groups, Class 1 and Class 2, that are respectively subdivided into types I, III, IV and types II, V, VI (1).

Over the last decade, our knowledge regarding CRISPR–Cas systems has expanded at an exceptional rate, mainly driven by a strong effort to harness their biotechnological potential (2–4). To date, the functions and mechanisms of action of all known CRISPR–Cas types have been characterized in detail, except for type IV for which the biological function(s) remain enigmatic. Importantly, type IV CRISPR–Cas modules have recently been reported to be primarily encoded by plasmids or, occasionally, by prophage genomes, evidencing the recurrent transfer of the CRISPR–Cas machinery to and from MGEs (5). Furthermore, although type IV cas operons are frequently associated with CRISPR arrays, they lack certain hallmark components of other CRISPR–Cas systems, including the highly conserved adaptation module and an effector nuclease (1). Consequently, these reduced systems have been proposed to exhibit altered CRISPR–Cas functions or to be functionally defective (6).

To date, type IV CRISPR–Cas loci are classified into two distinct subtypes, IV-A and IV-B, both of which share a common set of effector module proteins, including a highly diverged Cas7 (Csf2), Cas5 (Csf3) and a smaller version of Cas8 (Csf1) (1). Moreover, subtype IV-A loci encode a DinG family helicase (Csf4), a type IV-specific Cas6-like protein (Csf5), and they typically co-locate with a CRISPR array. In contrast, subtype IV-B loci lack *dinG*, *csf5* and an associated CRISPR array but they encode a putative ‘small subunit’ (Cas11) and they often neighbour a *cysH* gene (7,

\*To whom correspondence should be addressed. Tel: +45 3867 7360; Email: shiraz.shah@dbac.dk  
Correspondence may also be addressed to Rafael Pinilla-Redondo. Email: rafael.pinilla@bio.ku.dk  
Correspondence may also be addressed to Søren J. Sørensen. Email: sjs@bio.ku.dk

†The authors wish it to be known that, in their opinion, the second and third authors should be regarded as Joint Second Authors.

8). A recent structural and biochemical analysis of a subtype IV-A CRISPR–Cas system demonstrated the essential role of the Cas6-like enzyme in both the maturation of crRNAs and in the subsequent formation of a Cascade-like crRNA-guided effector complex, composed of Csf1, Csf3, Csf5 and multiple copies of Csf2 (9). These data suggest that the subtype IV-A effector complexes, as in other CRISPR–Cas systems, survey the cellular environment searching for matching nucleic acid targets. However, the study concluded that the spacers of the associated CRISPR arrays yielded no clear spacer–protospacer matches (9), but an earlier larger-scale analysis reported putative sequence matches to MGEs of which 72% were reported to be of viral origin (10).

In summary, it is plausible that subtype IV-A systems perform a defensive role, although the apparent absence of an effector nuclease suggests that the mechanism of interference differs significantly from those of other CRISPR–Cas systems. Consistent with this view, alternative functions have been suggested for type IV systems, including their involvement in plasmid propagation mechanisms, and in the enhancement of recombination events with other nucleic acids (7,9). In particular, the absence of CRISPR arrays linked to the minimal subtype IV-B system provides support for the effector module machinery participating in alternative cellular functions (7). In the present study, we have undertaken a comparative genomics approach to survey all publicly available bacterial and archaeal genomes for type IV CRISPR–Cas systems. The collected type IV systems were then subjected to an in-depth bioinformatic characterization to obtain insights into their biology and evolution.

## MATERIALS AND METHODS

### Detection, clustering and classification of type IV modules

Bacterial and archaeal complete and draft genomes were obtained from GenBank and scanned with the TIGR03115 Csf2 model (11) using HMMER3 (12). Protein sequences from two genes upstream and downstream of the detected *csf2* gene along with the Csf2 sequence itself were pooled and subjected to an all-against-all sequence comparison using FASTA (13). A neighbour-joining tree was constructed using distances derived from the aggregate similarities between each module pair using a previously described method (14). The tree was used to pick diverse representative type IV systems, which were then annotated manually using PSI-BLAST (15) searches. Hits on type III CRISPR–Cas systems were purged from the Csf2 tree. Following manual annotation, the protein sequences from the refined representative modules were pooled for another all-against-all sequence comparison. Protein sequences were clustered using the method previously described (16) and another aggregate module similarity tree was built for the refined representative type IV module set. The tree was overlaid with gene maps of the type IV modules marked with the obtained protein clustering information (Supplementary Figure S2). This was used for devising the subtypes/variants (Figure 1), which were then searched for in all downloaded genomes using the HMMs corresponding to all defined protein clusters. Subsequently, all resulting final type IV and non-type IV system proteins were subjected to another all-against-

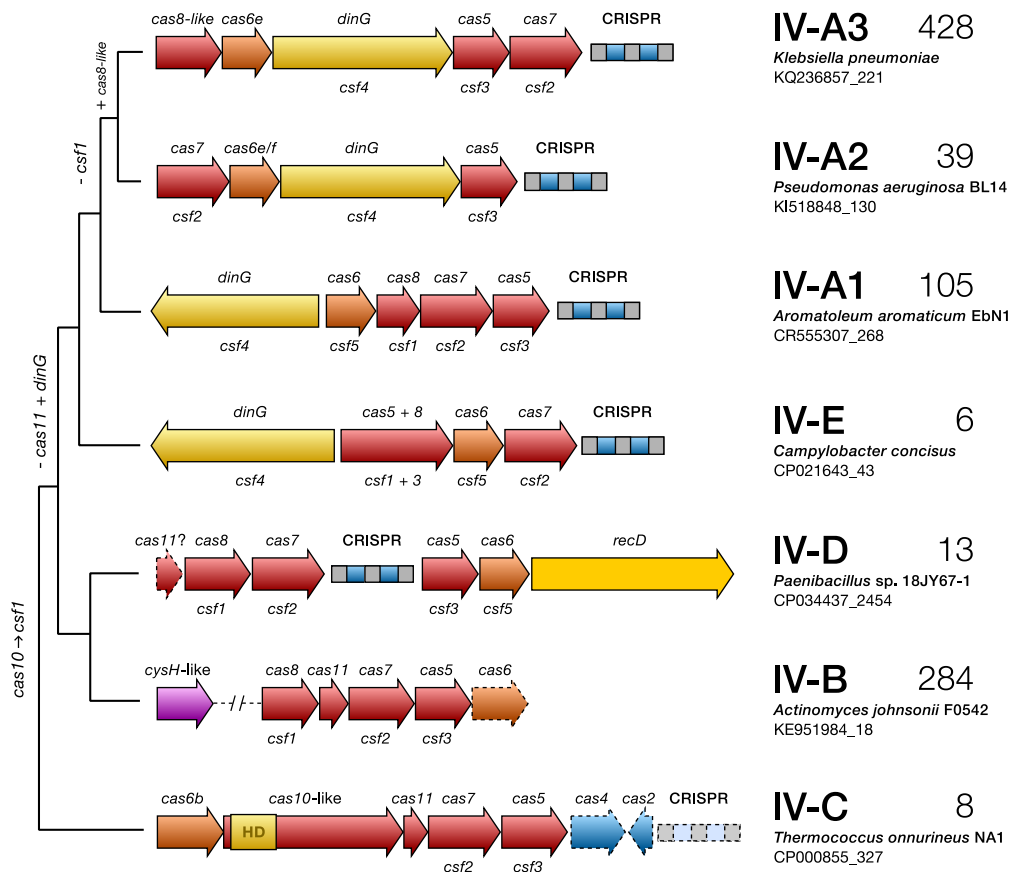
all sequence comparison using FASTA in order to build the final aggregate similarity subtype tree (Supplementary Figure S3) in addition to building the final Csf2 maximum-likelihood tree (Supplementary Figure S1).

### Spacer-protospacer match analysis

CRISPR arrays were detected with CRISPRCasFinder (4.2.17, (17)) and matched to a Type IV module if any predicted operon was within a 10 kb radius (distance to first gene in the operon, this cut-off was based on Supplementary Figure S13). Non-type IV CRISPR–Cas systems were also detected with CRISPRCasFinder in the same genome assemblies where type IV systems were found, and typing was manually corrected when necessary. Arrays were matched to an operon if it was within 10 kb (distance to first gene in the operon, see Supplementary Figure S13). Phage genomes were obtained from the April 2019 version of the millardlab.org phage database (<http://millardlab.org/bioinformatics/bacteriophage-genomes/>), and plasmid sequences were obtained from the PLSDB database (2019\_03.05, (18)). In order to rule out false positive matches to conserved spacers within undetected arrays on plasmids, the putative arrays in the plasmid database were masked when detected via CRISPRCasFinder (17,19), CRISPRdetect (2.2) (17,19) and CRT (1.2) (20). Furthermore, all unique repeats pertaining to arrays from the initial CRISPRCasFinder search were aligned (blastn -task blastn-short, (21)) against the masked PLSDB database, and putative arrays were defined if two or more matches (*E*-value < 0.1) were found within 100 bp, and these regions were masked as well. Spacers from the initial CRISPRCasFinder search yielded a total of 6021 type IV and 11 230 non-type IV spacers. These were collapsed into a unique spacer set using cd-hit-est (22) in order to avoid overrepresented (redundant) spacers from sequencing biases. The unique spacer set, consisting of 1051 type IV and 6778 non-type IV spacers (Supplementary Table S4), was aligned against the masked plasmid and phage databases using FASTA (13); a spacer match was considered significant when the *E*-value was <0.05.

### Targeted gene enrichment analysis

Enrichment in spacer targeting of certain functions was done by first predicting ORFs in all plasmid and phage genomes using Prodigal (23), and then clustering genes using the protein clustering algorithm previously described (16). The observed number of matches to each gene cluster was compared to  $10^5$  simulations of random draws from a binomial distribution with size *n* equal to the number of genes in the gene cluster and the probability, where the relative gene length was defined by dividing the length of each gene by the median gene length, and then finding the average length for each gene cluster. The above simulation only counted each spacer once, however, spacers usually match multiple genes in the same gene cluster. Therefore, each simulated match was multiplied by a random draw of the observed number of genes matched by a spacer matching that gene cluster.



**Figure 1.** A proposed classification of type IV CRISPR–Cas systems based on their genome loci architectures and evolutionary relationships. Phylogenetic tree depicting the typical operon organization of the identified subtype IV loci. A selected representative locus is shown for each clade wherein genes are colour-coded and labelled according to the protein families they encode, using both the cas (upper) and csf (lower) nomenclatures. Genes or CRISPR arrays that are not invariably present are represented with dashed lines on the gene maps. The number of loci identified for each clade is given on the right. Hypothesized gene gain/loss events over the course of evolution are shown on the left.

### PAM identification

From the 1016 CRISPR arrays (481 type IV and 535 non-type IV) detected in type IV containing complete and draft genomes, the consensus repeat for each array was aligned against corresponding consensus repeats for all other arrays using needleall (24). Consensus repeats that differed from each other by more than two mismatches were assigned to separate repeat clusters, resulting in 171 repeat clusters in total. The previous unique spacer matching output from FASTA was surveyed for protospacers pertaining to each of the 171 repeat clusters separately. Spacers matching several phages and plasmids in the database were only counted once to circumvent sequencing bias in the database. Logo plots were drawn from the ten nucleotides immediately flanking each side of each unique protospacer. Protospacers with alignment lengths smaller than the total spacer length had their coordinates adjusted so all flanks within a repeat cluster were properly aligned.

### CRISPR–Cas subtype co-occurrence analysis

Co-occurrence between type IV and non-type IV subtypes was analysed with phylogenetic logistic regression (phy-

loglm, maximum penalized likelihood estimation (13,25), with the non-type IV occurrence as the response and the type IV occurrence as the predictor. Besides the genomes with type IV systems, we supplemented the analysis with all complete genomes with at least one non-type IV operon (as defined by CRISPRCasFinder). The phylogenetic tree was based on 16S rRNA gene sequences detected by Barnmap (26, 27), aligned with mafft 7.307 (28), and tree made with FastTree2 (26), and was rooted by the archaeal clade. Edges of length zero were rescaled to the shortest non-zero branch length. Furthermore, outlier branches were pruned by removing tips for which the maximum phylogenetic variance-covariance was  $>2$ . Only non-type IV subtypes found in at least 100 genomes were included, and only the four most prevalent type IV subtypes were included. *P*-values were *fdr*-adjusted with the Benjamini–Hochberg method (29).

### CRISPR repeat heatmap

All unique CRISPR consensus repeats were aligned with the pairwise2 module from Biopython 1.73 (30). Repeats were globally aligned with globalxs with both open gap and extend gap penalties of 3, and no end gap penalties. Align-



ments were done on both strands, and the highest identity was used.

### Leader sequence analysis

Multiple sequence alignments were performed with the upstream regions of a series of representatives of co-occurring IV-A3 and I-E CRISPR arrays using MUSCLE (31). Alignments were analysed and visually displayed using Jalview (30,32). The corresponding leader sequence conservation profiles were generated using WebLogo 3 (33).

### Plasmid mobility prediction

The mobility of all plasmids (conjugative, mobilizable or non-mobilizable) in PLSDDB was predicted with mobtyper (34) with an *E*-value cut-off of  $1e-10$ . For calculating whether certain mobility types were enriched in targeted plasmids, the number of matches were scaled such that the sum for each spacer was 1, which ensured that each spacer only counted once, no matter how many matches it had.

### Plasmid/prophage prediction

To predict whether the CRISPR–Cas operons were located on chromosomes or MGEs, we used an iterative heuristic; first, contigs from complete genomes were annotated as plasmids or chromosomes as described in the NCBI name. Second, for draft genomes PlasFlow (35) was used to detect contigs that were putatively part of plasmids. Third, VirSorter (35, 36) was used to predict the presence of prophages, and all operons within a category 1, 2, 4 or 5 region were classified as putative prophages.

### Protein structure prediction

Protein homology models were generated with the Phyre2 protein structure prediction server (intensive mode) (37). Superimposition of protein structures were generated by the PyMol molecular visualization software (PyMOL Molecular Graphics System (C) Schrödinger, LLC).

## RESULTS

### Expanding the number of identified type IV CRISPR–Cas systems

In order to perform a comprehensive analysis of the diversity and distribution of type IV CRISPR–Cas systems, we first sought to expand the repertoire of currently identified loci. Although Csf1 has been proposed as a signature protein for type IV systems (1), we found that it was unsuitable, owing to its high level of sequence divergence between subtypes/variants and because of its absence from some loci. Instead, the Cas7-like (Csf2) protein was found to be the most conserved protein, and it was used as an initial query for searches against all publicly available complete and draft genomes (obtained from <ftp.ncbi.nih.gov>). Out of 883 detected Csf2 proteins (Supplementary Figure S1), 69 diverse representatives were selected for further analysis. The gene neighbourhoods of these representatives were

explored systematically and annotated manually via PSI-BLAST (15) searches, protein clustering (38) and profile-profile alignments (38, 39). An aggregate protein similarity tree was then generated including all proteins from the curated type IV modules. Finally, their corresponding gene maps were compared to gauge the diversity of their genetic compositions (Supplementary Figure S2).

### Type IV systems display a previously uncharted diversity of loci architectures

Our phylogenetic analysis outlines a hitherto unrecognized richness of type IV gene arrangements and reveals a complex evolutionary relationship between the different variants, pervaded by clear instances of horizontal gene transfer (Figure 1). The identified type IV loci are distributed across five major phylogenetically discrete groups that show consistent differences in their genetic compositions (Figure 1; Supplementary Figures S1–S3). Notably, a set of archaeal type IV modules were found to cluster as a clear outgroup and during the preparation of this work were proposed as a new subtype: IV-C (S.A.S. personal communication with K.S. Makarova). These distinctive loci share key organizational features with type III CRISPR–Cas systems, including the presence of a Cas10-like protein in place of Csf1, their common association with type I CRISPR–Cas systems, and the frequent absence of CRISPR arrays and adaptation modules (Supplementary Figure S4). Importantly, exhaustive protein domain searches with the IV-C Cas10-like protein revealed the typical Zn finger domain found in the middle section of other Class 1 CRISPR–Cas large subunits (Cas8, Csf1 and Cas10 families (40, 41)) and, similarly to type III Cas10 proteins, an N-terminal HD nuclease domain that is suggestive of DNA cleavage activity (Supplementary Table S1). However, no indication of the degenerate nucleotide cyclase palm domain motif ‘GGDD’ was found suggesting that, in contrast to *bona fide* Cas10s, this protein is not involved in oligoadenylate signaling.

Overall, subtypes IV-B and IV-A exhibit a high level of genetic diversity (Supplementary Figure S3). Subtype IV-B is composed of several phylogenetically divergent clades, merged here because of their similar genetic architectures, and subtype IV-A spans three major groups, hitherto referred as IV-A variants 1, 2 and 3. Notably, although subtypes IV-A2 and IV-A3 are closely related, they primarily differ in the absence (IV-A2), or presence (IV-A3), of a gene in their cas operons. We infer this gene encodes a Cas8-like protein due to its shared features with other Cas8 components, such as similar size and a zinc finger domain (Supplementary Table S2). Since Cas8 proteins often show little or no significant sequence similarity, even within subtypes (1, 14) (e.g. subtype I-B), and because all three IV-A variants cluster as a monophyletic group showing comparable modular architectures (Figure 1, Supplementary Figure S3), we maintain them within the same subtype. Notably, the common exchange of functional modules (42) between different CRISPR–Cas systems is particularly evident for IV-A2 and IV-A3, where Cas6 apparently has been recruited from subtype I-F and I-E systems, respectively (Figures 1 and 4b), highlighting a possible functional link between these subtypes.

Additionally, we identified a distinctive group of loci (named here subtype IV-D) which is unique in carrying a helicase of the RecD family in place of the archetypal DinG. This latter observation highlights the putatively central functional role of a dsDNA unwinding component in these systems. Moreover, while IV-B and IV-D appear to have diverged relatively recently, their classification into separate subtypes seems justified. Unlike subtype IV-D, IV-B loci are typically associated with a *cysH*-like gene (a member of the adenosine 5'-phosphosulfate reductase family (8, 43)) and they do not encode a helicase, a Cas6 (with rare exceptions: only 12 instances; Supplementary Data S1) or a CRISPR array. Finally, a few examples were found of an outgroup clade related to IV-A, labelled here as the putative subtype IV-E. Despite sharing similar modular architectures, their DinG components have diverged significantly (Supplementary Figure S5) and the Csf1 of subtype IV-E is fused to Csf3, as revealed by HHpred searches (Supplementary Table S3).

### Type IV systems are widely distributed across taxa and diverse MGEs

Our taxonomic analysis reveals a widespread, yet heterogeneous, distribution of type IV loci across a variety of prokaryotic genome backgrounds and they were primarily predicted to be encoded by MGEs (Figure 2, Supplementary Data S2). Subtypes IV-A and IV-B appear to be the most prevalent, contrasting with the sparse and relatively narrow taxonomic distribution of the other subtypes. While IV-A variants are mainly spread across proteobacterial plasmid-like conjugative elements, subtype IV-B is largely confined to predicted plasmids (and sometimes prophages) of Actinobacteria, and to a lesser extent Archaea, Firmicutes and Proteobacteria. The reduced group of subtype IV-C loci were found in Archaea, and no evidence for a preferential association with MGEs was found. Moreover, subtype IV-D occurs in some plasmids of Firmicutes and IV-E modules are present in Campylobacter and Bacteroides; some of the latter also residing in plasmid-like elements. Notably, in sharp contrast to the near-exclusive association of type IV systems with MGEs, we rarely found other CRISPR–Cas types to be encoded by plasmids or prophages, consistent with earlier reports and highlighting the uniqueness of type IV systems in this regard (1).

### Type IV spacer contents exhibit a strong bias towards plasmid protospacers

Statistical analyses of the distribution of spacer matches has proved a powerful tool for predicting functional properties of CRISPR–Cas systems and for understanding the ecology of the genomes carrying them (44,45). Given that type IV loci are primarily harboured by plasmids, semi-independent entities with selective pressures differing from those of their hosts (46), we sought to investigate whether type IV systems exhibit different targeting preferences from non-type IV systems. Therefore, we performed a comprehensive analysis of the spacer–protospacer matches for all the type IV-associated CRISPR arrays, and for the CRISPR arrays of

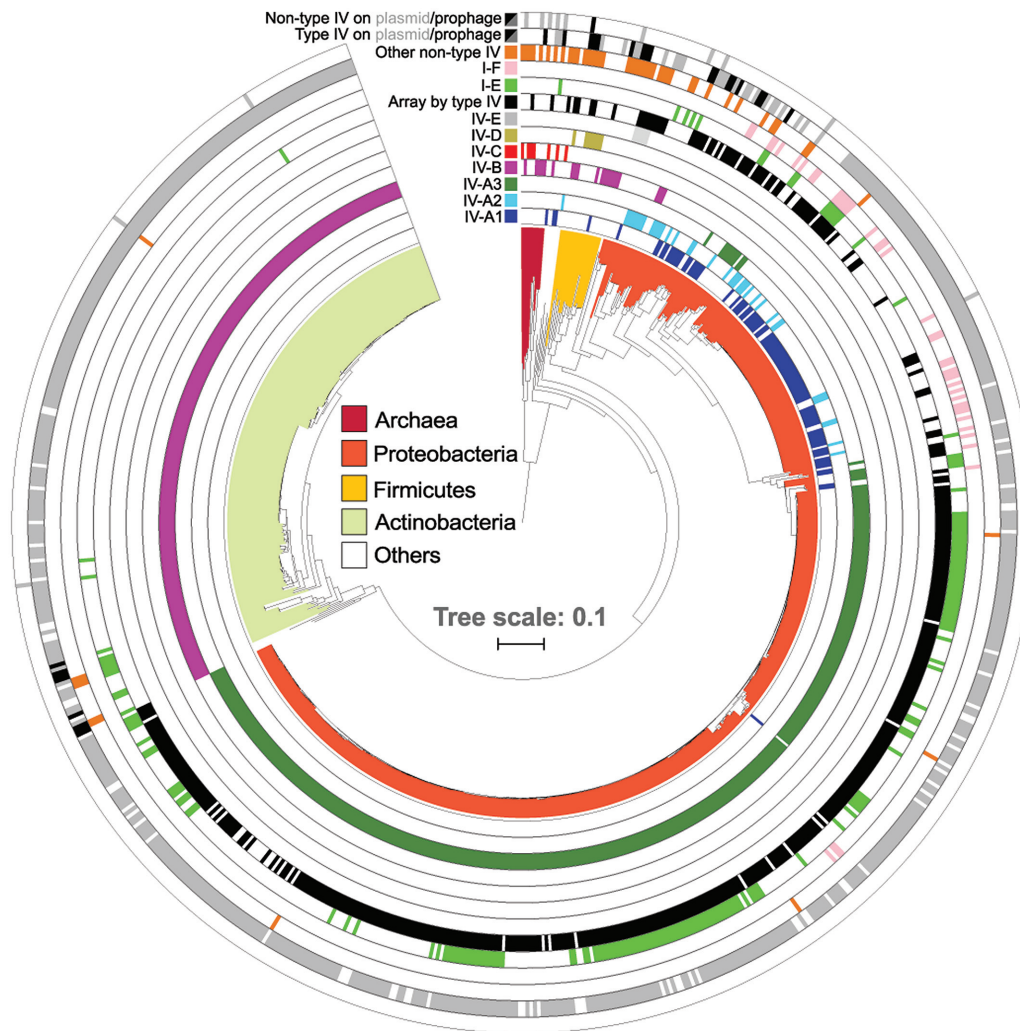
all other identified non-type IV systems present in the host genomes.

A total of 481 type IV and 535 non-type IV arrays were identified and 1051 and 6778 unique spacers, respectively, were extracted for spacer–protospacer match analyses (Supplementary Table S4). Consistent with earlier results (1, 8), only a small fraction of spacers yielded significant matches: ~12% and ~7%, for type IV and non-type IV, respectively (Figure 3A, Supplementary Table S4 and Data S3 and S4), which reflects the current undersampling of the microbiome and high evolutionary rates of MGEs (10). However, we observed that type IV systems displayed an exceptionally strong targeting bias towards plasmids, in contrast to the other co-occurring CRISPR–Cas systems (80% versus 26%, respectively). Importantly, this trend was valid for all DinG associated type IV subtypes and variants, whereas the remaining subtypes did not yield sufficient data. On the other hand, non-type IV subtypes overall exhibited the previously reported strong preference for viral targets (Figure 3A, Supplementary Table S4) (1, 10). Given that the type IV and non-type IV spacer contents investigated here originate from the same cellular environments, the results strongly underline an anti-plasmid function for type IV systems.

Next, in order to further explore the potential functional differences between type IV and non-type IV systems, we examined possible variations in their targeting preferences towards specific plasmid and viral gene families. Statistical analyses of the spacer match distributions revealed an enrichment of certain plasmid and virus-related genes, yet no consistent differences were observed between type IV and non-type IV targets (Supplementary Figure S6, Data S3 and S4). In agreement with previous reports (10), all CRISPR–Cas types revealed a targeting preference for conserved, and frequently plasmid-borne, genes; e.g. conjugative transfer machinery genes (Supplementary Figure S6). Although a similar pattern was observed for non-type IV viral gene matches, the corresponding analysis for type IV was inconclusive due to the low number of identified viral protospacers. Next, we investigated whether the plasmids targeted by type IV-derived spacers displayed any unifying biological features that could provide insights into the function of type IV systems. In general, we found that targeted plasmids tend to be relatively large (Figure 3C, targeted: 155 kb, PLSDB: 53 kb, median sizes,  $P < 2.2 \times 10^{-16}$ , Kolmogorov–Smirnov test), irrespective of their predicted mobility (Supplementary Figure S7), and there was a clear bias towards the targeting of conjugative plasmids (type IV: 48%, PLSDB: 30%, Figure 3B).

### Type IV associations with other CRISPR–Cas systems

The almost exclusive absence of adaptation module genes from type IV loci (Supplementary Data S5) raises the question as to the origin of CRISPR spacers. This aspect of type IV's biology is especially puzzling given the observed variability in spacer content between related type IV CRISPR loci. Notably, spacer acquisition invariably requires Cas1 and Cas2, the most conserved components of all CRISPR–Cas systems (47). This high conservation implies that there could be degrees of compatibility between adaptation mod-



**Figure 2.** Distribution of type IV loci across prokaryotic taxa and MGE types. Phylogenetic tree based on 16S rRNA gene sequences of all bacteria and archaea that carry type IV CRISPR–Cas systems. Concentric rings denote the presence or absence of type IV and other co-occurring non-type IV CRISPR–Cas loci in the same genomes, colour-coded according to the subtype/variant to which they belong. All non-type IV systems, except I-E (light green) and I-F (pink), were merged into one lane (orange) for visualization purposes. Type IV effector cas operons for which an associated CRISPR array was detected are shown (black). Based on genomic context analyses (Methods), CRISPR–Cas systems predicted to be encoded by plasmid-like elements (grey) or (pro)phages/viruses (black) are shown (Supplementary Data S2), for both type IV and non-type IV loci (two outermost ring lanes). 758 (of 883) identified type IV loci are displayed on the tree; for the remainder no 16S rRNA gene sequence was found in the genome.

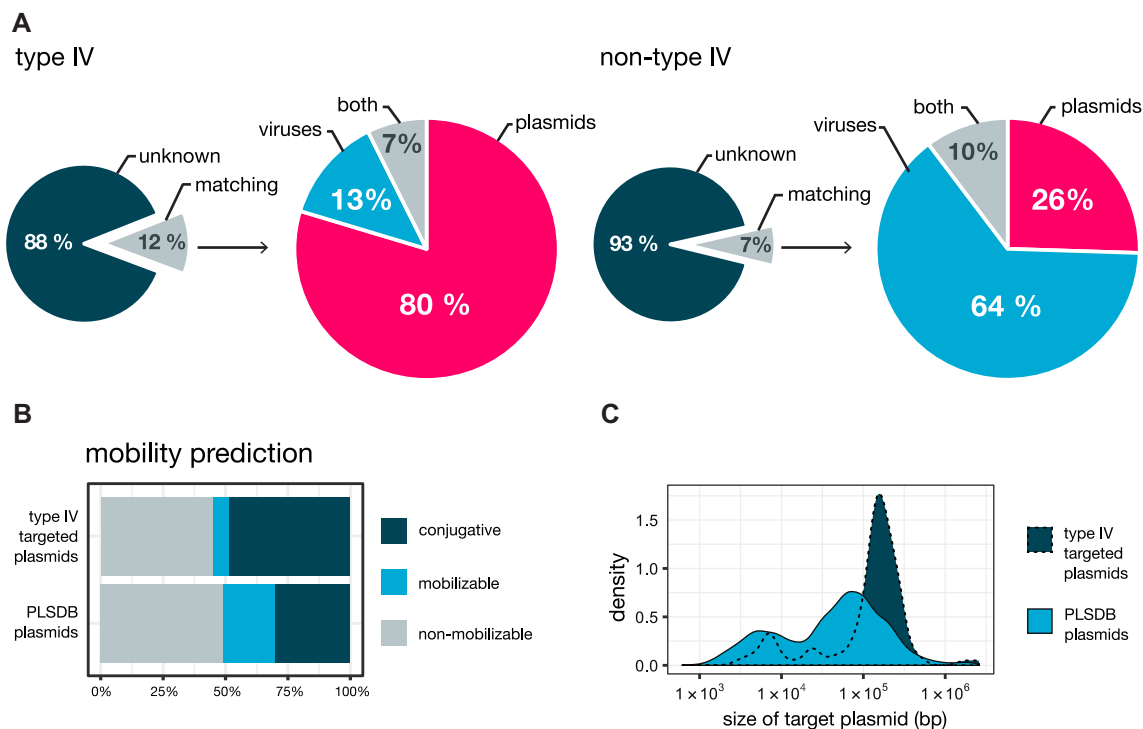
ules of different CRISPR–Cas types. Therefore, we reasoned that type IV loci could exploit this functional redundancy by co-opting Cas1/Cas2 adaptation modules from other CRISPR–Cas systems that coexist intracellularly.

To explore this hypothesis, we first searched for evidence of positive correlations between different type IV subtypes/variants and other CRISPR–Cas systems present within the same hosts (Figure 4A). Interestingly, significant positive correlations were found for subtypes IV-A1/2 and IV-A3, together with subtypes I-F and I-E, respectively (IV-A3 and I-E:  $P = 3.7 \times 10^{-13}$ , IV-A2 and I-F:  $P = 4.5 \times 10^{-10}$ , IV-A1 and I-F:  $P = 3.7 \times 10^{-17}$ , *fdr*-adjusted  $P$ -values from phylogenetic logistic regression). We also found significant negative correlations between several subtypes, including IV-A1 with I-B, I-E and I-C, IV-A3 with I-C, I-F and II-C, and IV-B with I-B, I-C, I-E and II-C, which could be due, at least partly, to the targeting of type IV-

carrying plasmids/MGEs by host-encoded CRISPR–Cas systems. Furthermore, co-clustering of CRISPR repeats demonstrated that type IV repeat sequences are similar to those from CRISPR loci with which they co-occur and/or correlate positively (IV-A1/2, IV-A3 and IV-D, with I-F, I-E and I-B, respectively) (Figure 4C), strengthening the notion of a potential functional connection between type IV and other co-encoded CRISPR–Cas systems.

PAM (protospacer adjacent motif) recognition is essential for Cas1/Cas2-dependent spacer acquisition and self/non-self discrimination in most CRISPR–Cas systems (48–50), yet such motifs have not yet been described for type IV systems. Therefore, we investigated whether PAMs could be identified and, if so, whether they are compatible with co-occurring non-type IV CRISPR–Cas systems. To test this, we predicted PAMs *in silico* by aligning protospacer flanking regions and a putative PAM was identi-





**Figure 3.** Spacers from type IV systems preferentially target plasmid-borne protospacers. (A). Comparison of spacer–protospacer matches detected for type IV systems (left) and the co-occurring non-type IV systems (right). A more detailed breakdown, by CRISPR–Cas subtype/variant, is presented in Supplementary Table S4. (B). Distribution of type IV spacer hits on plasmids as a function of predicted plasmid mobility. (C). Size distribution of the targeted plasmids. The mobility prediction and size for the collection of PLSDB plasmids are displayed as a reference in both ‘B’ and ‘C’ plots.

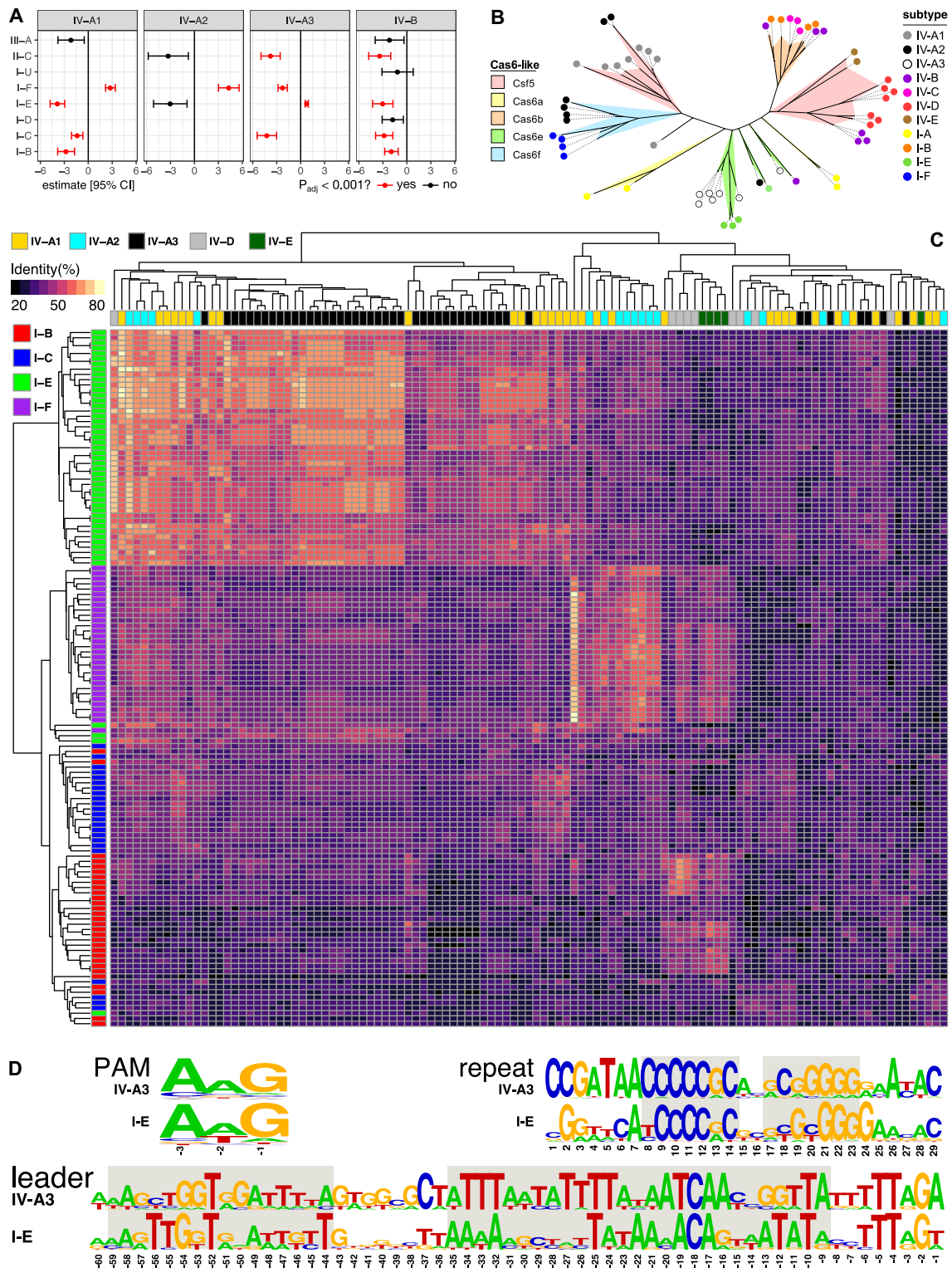
fied for subtype IV-A3 (Figure 4D). However, searches for other subtypes/variants were unsuccessful, likely due to the low number of spacer–protospacers matches (Supplementary Table S4). Importantly, the predicted subtype IV-A3 PAM (-AAG-) is identical to that of the positively correlating type I-E CRISPR–Cas system.

The higher numbers of detected subtype IV-A3 loci provided the basis for a case study involving more extensive comparative analyses. Alignments of consensus repeats of the positively correlating subtypes IV-A3 and I-E (Supplementary Figure S8) revealed the previously described recognition sites for the Cas1–Cas2e adaptation machinery (51) (Figure 4D). In addition, multiple sequence alignments of the upstream regions from co-occurring IV-A3 and I-E CRISPR arrays (Supplementary Figure S9A and B, respectively) showed similar conserved motifs in the leader region (Figure 4D). Importantly, these conserved sequences comprise the binding sites for the Cas1–Cas2e complex and the integration host factor (IHF), both of which are essential for uptake of new spacers into leader-repeat junctions of type I-E arrays (51,52). Next, we searched for evidence of preferential acquisition of spacers in the leader-end of IV-A3 CRISPR arrays, a phenomenon described for some CRISPR–Cas systems (53–56). However, clustering of related IV-A3 CRISPR loci and a comparison of their spacer contents did not reveal any clear support for this (Supplementary Figure S10). Finally, analysis of type IV-associated Cas6-like proteins yielded evidence for a polyphyletic origin, with independent acquisitions having occurred on multiple occasions (Figure 4B). For example, IV-A3 and IV-

A1/2 loci contain Cas6 variants that are more closely related to Cas6e and Cas6f, respectively, than to Cas6 (57), and co-occurring IV-C carries a Cas6b enzyme, further underlining the functional interrelations occurring between type IV and type I systems.

## DISCUSSION

In addition to the recognized adaptive immune functions of CRISPR–Cas systems, there is increasing evidence that diverse MGEs, including phages, giant viruses and transposons, have co-opted these systems for alternative functions (7, 58–61). The discovery of type IV systems, and of their frequent encoding on plasmids, is also relatively recent (1). To date, only two subtypes (IV-A and IV-B) are known and their biological functions and mechanisms of action remain obscure. In this work, we identify several novel type IV subtypes/variants and incorporate them into a revised type IV classification (Figure 1). In agreement with previous work, we found that the newly identified type IV loci are primarily encoded by prokaryotic MGEs, most of which are predicted to be plasmids (Figure 2). Notably, given the current limited sequence information covering the ‘dark matter’ of the mobilome (10, 59, 60), our findings likely underestimate the true diversity and distribution of these systems. Future comparative genomic characterizations will clearly benefit from including metagenomic sequence datasets and the continuing global effort to sample the meta-mobilome.



**Figure 4.** Interactions between type IV CRISPR–Cas systems and other co-encoded CRISPR–Cas systems in a host. **(A)** Co-occurrence analysis between type IV and non-type IV systems. Estimates are from phylogenetic logistic regressions, with *P*-values *fdr*-adjusted. Only estimates with standard errors <10 are shown. **(B)** Unrooted phylogenetic tree for Cas6/Csf3 built with representatives covering the diversity of type IV and type I subtypes/variants detected in this study. Each cluster is coloured according to the cas6-like family it corresponds to, and the coloured dot at the end of branches indicates the specific RISPR–Cas subtype/variant encoding such a Cas6-like protein. **(C)** Heat map depicting CRISPR repeat similarity of co-occurring CRISPR–Cas subtypes/variants clustered by average linkage hierarchical clustering. **(D)** PAM, consensus CRISPR repeat and leader sequence logos for the positively correlated subtypes IV-A3 and I-E. The short semi-palindromic repeats at the centre of the consensus repeat that are used as anchor sequences by the Cas1–Cas2e complex are highlighted in grey, as well as the conserved leader sequences comprising the binding sites for the Cas1–Cas2e complex (left) and the IHF (right).



Our analyses suggest an origin of type IV systems from a type III-like ancestor in archaea (most similar to the subtype IV-C described in this work), comparable to the evolutionary pathway proposed for the emergence of type I CRISPR–Cas systems (41). This is further supported by IV-C (a) carrying a Csf2 with structural similarities to Cmr4/Csm3, the Cas7-like helical backbone subunit in subtypes III-A/B (Supplementary Figure S11), (b) being exceptional for type IV in carrying an HD domain on its Cas10-like protein instead of a helicase and Cas8/Csf1 (Figure 1), (c) being most commonly found in archaeal hyperthermophile genomes and (d) being only intermittently associated with CRISPR arrays and type I systems (Supplementary Figure S4), all of which are characteristic properties of type III systems. With regard to the evolution of the remaining type IV subtypes, a parsimonious scenario involves streamlining of the Cas10-like protein into Csf1 generating a IV-B-like ancestor that acquired a RecD helicase which led, in turn, to the evolution of subtype IV-D. In a separate branch such a IV-B-like ancestor is speculated to have lost Cas11 and acquired DinG leading to the evolution of subtypes IV-A and IV-E. The fusion between Csf3 and Csf1 in IV-E is consistent with the proximity of these proteins in Class 1 effector complexes. As for IV-A, although the three variants are all closely related, IV-A2 appears to derive from IV-A1 after loss of Csf1, but retaining CRISPR–Cas functionality, while IV-A3 is a more recent variant of IV-A2 that seems to have gained a substitute for Csf1, the Cas8-like protein.

This evolutionary scenario is underpinned further by the overall taxonomic distribution of the identified type IV loci, ranging from the broad occurrence of IV-B, to derived variants such as IV-A3 being restricted to a few genera of Proteobacteria (Figure 2). Interestingly, the emergence of IV-D from an IV-B-like ancestor pool may have occurred more than once, as evidenced by the paraphyly of IV-Ds (Supplementary Figure S1). Subtypes IV-D, IV-A and IV-E are CRISPR array-associated, unlike IV-B which is more diverse and CysH associated. This likely reflects convergent evolution where type IV systems initiated as CRISPR–Cas immune systems and then evolved an altered functionality before reverting back into CRISPR–Cas systems via lateral acquisition of a DNA helicase.

Contrary to the strong viral targeting preference of all other known CRISPR–Cas types, our work reveals that type IV systems exhibit an exceptional targeting bias towards plasmid-like elements (Figure 3A, Supplementary Table S4). While preliminary work had reported that some spacers from type IV loci in *Klebsiella* plasmids matched plasmid genomes, no systematic analyses were carried out (62). The plasmid targeting bias described here was not detected in earlier systematic studies that employed lower numbers of non-redundant spacers and were primarily centred around the matching of spacers against (pro)virus databases (7,9,10,63). Our additional matching of spacers against PLSDB, a comprehensive database of >16 000 curated plasmid genomes (18), was key in determining the plasmid targeting bias. Importantly, since our analysis of the spacer contents from other non-type IV CRISPR–Cas systems coexisting intracellularly with type IV loci clearly yielded the established bias towards viral targets (Figure

3A, Supplementary Table S4), and both analyses were done matching spacers against the same databases, we conclude that the reported type IV plasmid bias cannot be an artefact.

Interestingly, a significant enrichment of certain targeted gene families was observed (Figure 3B), particularly those encoding components of complex molecular machineries including the conjugative transfer apparatus (Supplementary Figure S6, Data S3). An explanation for this phenomenon is that conserved genes are less prone to mutational escape from CRISPR–Cas targeting, and thus lead to a positive selection of their cognate spacers over time (64). Moreover, spacer retention may also be further enhanced when a targeted gene is shared by distinct MGEs, as also occurs, for example, with conjugative transfer genes.

The finding that type IV systems are carried by plasmid-like elements that primarily target other plasmids leads to the basic question as to how and why an anti-plasmid bias emerged. Our results indicate that type IV systems may have evolved to target plasmid-like elements more effectively than, for example, phages/viruses, although the mechanistic basis of such a bias remains unclear. Moreover, certain plasmids may provide strong competition for type IV CRISPR–Cas-carrying plasmids, leading to the selection of spacers against the former plasmids over time. Whereas phages/viruses can interfere with plasmid survival by killing the host, cells already carry potent defence systems against these fatal intruders (65). Therefore, plasmids may be more strongly challenged by other intracellular plasmid-like elements which, while not being especially detrimental to the host, may compete directly for common cellular resources (66,67). The latter argument receives support from the accepted community ecology view that similar entities compete more strongly for overlapping niches and resources (68–70). Notably, recent work has proposed that many (pro)phages readily engage in similar CRISPR-based inter-virus warfare dynamics, utilizing ‘mini-arrays’ with spacers targeting viruses to prevent host superinfection (7). In summary, our results imply that plasmid-like elements leverage type IV systems to eliminate other plasmids with similar properties and lifestyles, in order to monopolize the host environment.

In addition, our findings reveal an apparent functional cross-talk between type IV modules and other co-occurring CRISPR–Cas systems within a host, thereby providing a credible explanation for the minimal nature of type IV systems. Not only did some type IV subtypes correlate positively with specific type I subtypes (Figure 4A) but there were also additional parallels between some co-occurring pairs: PAM sequence sharing, high CRISPR repeat sequence similarity and a high similarity between the Cas6 processing enzymes (Figure 4d,c,b, respectively). Noteworthy, future experimental work is required to both validate the predicted IV-A3 PAM and establish whether the large subunit Csf1 facilitates PAM recognition on nucleic acid targets, as occurs for its homolog Cas8 in type I systems (71, 72). Furthermore, we also found shared conserved CRISPR leader motifs for the binding of the Cas1/2e adaptation machinery and IHF between co-occurring IV-A3 and I-E subtypes (Figure 4D). Although all these results are con-

sistent with the inference that type IV systems can rely on the Cas1–Cas2 adaptation machinery from co-occurring ‘helper’ type I systems, such co-functionality requires experimental validation. Nevertheless, this hypothesis receives support from the numerous accounts of type III systems lacking Cas1 and Cas2 which utilize CRISPR-arrays maintained by adaptation modules from neighbouring type I systems (73, 74), and is further reinforced by the evolutionary links demonstrated here between type IV and type III systems.

Nevertheless, alternative spacer acquisition strategies cannot be ruled out. These include, for example, the mechanism proposed for viral-derived orphan mini-arrays, where recombination with host CRISPRs seems most likely (7). Consistent with the latter hypothesis, we observe examples of spacer rearrangements between related IV-A3 CRISPR loci (Supplementary Figure S10). Interestingly, most type IV systems carry a Cas6-like component, suggesting that specific pre-crRNA processing may be necessary for exclusive crRNA coupling with type IV effector complexes. This extra level of specificity and the stark contrast in spacer targets between positively correlating type I and type IV systems indicates that, although there are functional ties at the adaptation stage, the crRNA utilization stage operates independently. Moreover, type IV systems may benefit from carrying their own Cas6 component by ensuring control over crRNA processing, especially in cells where no host-derived Cas6 is available.

Although elucidation of the specific targeting mechanism of type IV systems requires an experimental approach, it is likely that the associated helicase (DinG/RecD) is required and involved in the specificity towards plasmid targets. Although our exhaustive analyses did not locate a nucleolytic active site in these enzymes, the presence of a cryptic nuclease domain is possible. In such a case, RecD/DinG could function mechanistically similarly to Cas3, the helicase–nuclease effector component of type I CRISPR–Cas systems (75, 76). Interestingly, similarly to Cas3, some non-type IV-associated DinG helicases have evolved 3′→5′ exonuclease activity (7, 77). However, even in the absence of a nuclease component, type IV systems could still co-opt host-encoded restriction enzymes to cleave their targets, possibly by rendering them susceptible to degradation upon dsDNA unwinding. In support of this, type III systems have recently been shown to utilise host degradosome nucleases to ensure successful interference of diverse MGEs (78). Intriguingly, chromosomally derived RecD homologs are known to take part in the RecBCD complex (exonuclease V) which, in addition to playing a role in DNA repair, carries out defence functions through the degradation of invading genetic elements (79).

Binding of type IV effector complexes to DNA could also destabilize the target, especially if it constitutes a rapidly replicating element. The consequences of replication fork collisions with protein–nucleic acid complexes (e.g. the transcription machinery) on genome integrity are well documented and can include replication fork arrest, premature transcription termination, and double-strand DNA breaks (80). Notably, these physical conflicts are also known to destabilize plasmids, eventually leading to their extinction from within cell lineages (80, 81). The latter explana-

tion is compatible with the hypothesis that type IV systems function similarly to the artificially developed catalytically dead CRISPR–Cas systems, which bind DNA targets but lack cleavage activity (e.g. dCas9) (82). These so-called CRISPR interference (CRISPRi) systems, silence the expression of targeted genes by blocking transcription factor binding or RNA polymerase elongation (82). Moreover, type IV-mediated gene silencing could serve purposes beyond plasmid–plasmid warfare, such as altering host expression profiles to enhance plasmid propagation and/or stabilize maintenance, all piracy practices, which plasmids are known to invoke via diverse mechanisms (83–85). In the context of CRISPRi functionality, for which R-loop formation between the crRNA and the DNA target is key, the common association of DinG with type IV loci appears paradoxical, as it is well documented that the substrate for this helicase are R-loops that block replication fork advancement (86–88). Thus, it is tentative to speculate about the potential regulatory or antagonistic role of the helicase component in the removal of type IV crRNA–DNA hybrids, although the purpose of such a function remains unclear. Interestingly, *dinG* sometimes appears in the opposite orientation to the other genes in type IV loci (Figure 1), consistent with the notion that its expression might be controlled independently.

Subtype IV-B systems constitute the most reduced and enigmatic version of type IV systems, lacking identifiable CRISPR arrays, Cas6, and a helicase component. This exceptional combination of features led to the proposition that it performs a different function from the other type IV systems, e.g. similar to transposon-encoded CRISPR–Cas systems (7, 59, 85). Because type IV-B systems encode all the necessary components to generate a Cascade-like surveillance complex (Csf1, Csf2, Csf3), we hypothesized that it could accommodate pre-processed crRNAs originating from other co-occurring CRISPR–Cas systems. However, we found no evidence of neighbouring CRISPR arrays, mini-arrays, SRUs (7) or of palindromic sequences that could yield the characteristic stem-loop secondary structures of crRNAs. Interestingly, our data revealed significant negative correlations of IV-B with the presence of all other CRISPR–Cas systems in the hosts (Figure 4A). Taken together, it seems plausible that these systems could have been repurposed by plasmids/phages to bind and neutralise crRNAs that become available, thereby antagonizing other CRISPR–Cas functions in the intracellular milieu. Nonetheless, the complexity of such an anti-CRISPR (Acr) mechanism would greatly contrast that of all other Acrs described to date (89), thus rendering this explanation unlikely. The key to deciphering the function of subtype IV-B possibly resides in its obscure, nearly invariant, genomic association with *cysH*, a protein which seems to have co-diversified with this subtype (Supplementary Figure S12). Since *cysH* belongs to the phosphoadenosine phosphosulfate reductase family, to which DNA phosphorothioate modification enzymes also belong, these systems could be involved in epigenetic silencing, or either linked to or antagonizing, related RM functions (7, 89, 90).

Collectively, our results provide further evidence of the strong dynamic pairing between CRISPR–Cas systems and MGEs. This complex co-evolutionary interrelation fits the

described ‘guns for hire’ paradigm, where CRISPR–Cas components are recurrently co-opted by different genetic entities for myriad defence and offence functions (6). Noteworthy, repurposing the power and programmability of type IV systems for controlling plasmid propagation presents promising biotechnological applications, particularly in the face of the current growing concerns regarding the spread of virulence and antibiotic resistance determinants within and between microbiomes (91,92). Indeed, as the mysteries surrounding the biology of type IV systems continue to be unveiled further opportunities will arise for expanding the CRISPR–Cas molecular toolbox.

## SUPPLEMENTARY DATA

Supplementary Data are available at NAR Online.

## ACKNOWLEDGEMENTS

The authors would like to thank Jonas Stenlørkke Madsen and Joseph Nesme for helpful discussions and useful comments.

*Author contributions:* Conceptualization, R.P.-R., S.A.S., D.M.-M. and J.R.; direction and planning, R.P.-R., and S.A.S. with support from S.J.S.; investigation and data analysis, R.P.-R., J.R., S.A.S. and D.M.-M.; writing, R.P.-R. with support from S.A.S., J.R., D.M.-M. and R.A.G., and in consultation with S.J.S. and L.R.

## FUNDING

DFF Independent Research Fund Denmark (to R.P.-R.); Novo Nordisk Foundation Tandem programme (to J.R.); Lundbeckfonden (to S.J.S.); DFG SPP2141 and Heisenberg Programme (to L.R.); Capital Region of Denmark [A6291 to S.A.S.]. Funding for open access charge: Novonordisk Foundation Basic Bioscience programme.

*Conflict of interest statement.* None declared.

## REFERENCES

- Makarova, K.S., Wolf, Y.I., Alkhnbashi, O.S., Costa, F., Shah, S.A., Saunders, S.J., Barrangou, R., Brouns, S.J.J., Charpentier, E., Haft, D.H. *et al.* (2015) An updated evolutionary classification of CRISPR–Cas systems. *Nat. Rev. Microbiol.*, **13**, 722–736.
- Barrangou, R. and Doudna, J.A. (2016) Applications of CRISPR technologies in research and beyond. *Nat. Biotechnol.*, **34**, 933–941.
- Pickar-Oliver, A. and Gersbach, C.A. (2019) The next generation of CRISPR–Cas technologies and applications. *Nat. Rev. Mol. Cell Biol.*, **20**, 490–507.
- Hsu, P.D., Lander, E.S. and Zhang, F. (2014) Development and applications of CRISPR–Cas9 for genome engineering. *Cell*, **157**, 1262–1278.
- Koonin, E.V. and Makarova, K.S. (2017) Mobile genetic elements and evolution of CRISPR–Cas systems: all the way there and back. *Genome Biol. Evol.*, **9**, 2812–2825.
- Koonin, E.V. and Makarova, K.S. (2019) Origins and evolution of CRISPR–Cas systems. *Philos. Trans. R. Soc. Lond. B Biol. Sci.*, **374**, 20180087.
- Faure, G., Shmakov, S.A., Yan, W.X., Cheng, D.R., Scott, D.A., Peters, J.E., Makarova, K.S. and Koonin, E.V. (2019) CRISPR–Cas in mobile genetic elements: counter-defence and beyond. *Nat. Rev. Microbiol.*, **17**, 513–525.
- Shmakov, S.A., Makarova, K.S., Wolf, Y.I., Severinov, K.V. and Koonin, E.V. (2018) Systematic prediction of genes functionally linked to CRISPR–Cas systems by gene neighborhood analysis. *Proc. Natl. Acad. Sci. U.S.A.*, **115**, E5307–E5316.
- Özcan, A., Pausch, P., Linden, A., Wulf, A., Schühle, K., Heider, J., Urlaub, H., Heimerl, T., Bange, G. and Randau, L. (2019) Type IV CRISPR RNA processing and effector complex formation in *Aromatoleum aromaticum*. *Nat. Microbiol.*, **4**, 89–96.
- Shmakov, S.A., Sitnik, V., Makarova, K.S., Wolf, Y.I., Severinov, K.V. and Koonin, E.V. (2017) The CRISPR Spacer Space is Dominated by Sequences from Species-Specific Mobilomes. *MBio*, **8**, e01397-17.
- Haft, D.H., Selengut, J.D., Richter, R.A., Harkins, D., Basu, M.K. and Beck, E. (2013) TIGRFAMs and genome properties in 2013. *Nucleic Acids Res.*, **41**, D387–D395.
- Potter, S.C., Luciani, A., Eddy, S.R., Park, Y., Lopez, R. and Finn, R.D. (2018) HMMER web server: 2018 update. *Nucleic Acids Res.*, **46**, W200–W204.
- Pearson, W.R. (2016) Finding protein and nucleotide similarities with FASTA. *Curr. Protoc. Bioinformatics*, **53**, 3.9.1–3.9.25.
- Vestergaard, G., Garrett, R.A. and Shah, S.A. (2014) CRISPR adaptive immune systems of Archaea. *RNA Biol.*, **11**, 156–167.
- Altschul, S. (1997) Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.*, **25**, 3389–3402.
- Shah, S.A., Alkhnbashi, O.S., Behler, J., Han, W., She, Q., Hess, W.R., Garrett, R.A. and Backofen, R. (2019) Comprehensive search for accessory proteins encoded with archaeal and bacterial type III CRISPR–cas gene cassettes reveals 39 new cas gene families. *RNA Biol.*, **16**, 530–542.
- Couvin, D., Bernheim, A., Toffano-Nioche, C., Touchon, M., Michalik, J., Néron, B., Rocha, E.P.C., Vergnaud, G., Gautheret, D. and Pourcel, C. (2018) CRISPRCasFinder, an update of CRISPRFinder, includes a portable version, enhanced performance and integrates search for Cas proteins. *Nucleic Acids Res.*, **46**, W246–W251.
- Galata, V., Fehlmann, T., Backes, C. and Keller, A. (2019) PLSDB: a resource of complete bacterial plasmids. *Nucleic Acids Res.*, **47**, D195–D202.
- Biswas, A., Staals, R.H.J., Morales, S.E., Fineran, P.C. and Brown, C.M. (2016) CRISPRDetect: a flexible algorithm to define CRISPR arrays. *BMC Genomics*, **17**, 356.
- Bland, C., Ramsey, T.L., Sabree, F., Lowe, M., Brown, K., Kyrpides, N.C. and Hugenholtz, P. (2007) CRISPR recognition tool (CRT): a tool for automatic detection of clustered regularly interspaced palindromic repeats. *BMC Bioinformatics*, **8**, 209.
- Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K. and Madden, T.L. (2009) BLAST+: architecture and applications. *BMC Bioinformatics*, **10**, 421.
- Li, W. and Godzik, A. (2006) Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics*, **22**, 1658–1659.
- 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.
- Rice, P., Longden, I. and Bleasby, A. (2000) EMBOSS: the European molecular biology open software suite. *Trends Genet.*, **16**, 276–277.
- Ho, L. si T and Ané, C. (2014) A linear-time algorithm for Gaussian and non-Gaussian trait evolution models. *Syst. Biol.*, **63**, 397–408.
- Price, M.N., Dehal, P.S. and Arkin, A.P. (2010) FastTree 2 – approximately maximum-likelihood trees for large alignments. *PLoS One*, **5**, e9490.
- Seemann, T. (2014) Prokka: rapid prokaryotic genome annotation. *Bioinformatics*, **30**, 2068–2069.
- Katoh, K. and Standley, D.M. (2013) MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol. Biol. Evol.*, **30**, 772–780.
- Benjamini, Y. and Hochberg, Y. (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. R. Stat. Soc. B (Methodological)*, **57**, 289–300.
- Cock, P.J.A., Antao, T., Chang, J.T., Chapman, B.A., Cox, C.J., Dalke, A., Friedberg, I., Hamelryck, T., Kauff, F., Wilczynski, B. *et al.* (2009) Biopython: freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics*, **25**, 1422–1423.
- Edgar, R.C. (2004) MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.*, **32**, 1792–1797.



32. Waterhouse, A.M., Procter, J.B., Martin, D.M.A., Clamp, M. and Barton, G.J. (2009) Jalview Version 2—a multiple sequence alignment editor and analysis workbench. *Bioinformatics*, **25**, 1189–1191.
33. Crooks, G.E., Hon, G., Chandonia, J.-M. and Brenner, S.E. (2004) WebLogo: a sequence logo generator. *Genome Res.*, **14**, 1188–1190.
34. Robertson, J. and Nash, J.H.E. (2018) MOB-suite: software tools for clustering, reconstruction and typing of plasmids from draft assemblies. *Microb. Genom.*, **4**, 8.
35. Krawczyk, P.S., Lipinski, L. and Dziembowski, A. (2018) PlasFlow: predicting plasmid sequences in metagenomic data using genome signatures. *Nucleic Acids Res.*, **46**, e35.
36. Roux, S., Enault, F., Hurwitz, B.L. and Sullivan, M.B. (2015) VirSorter: mining viral signal from microbial genomic data. *PeerJ*, **3**, e985.
37. Kelley, L.A., Mezulis, S., Yates, C.M., Wang, M.N. and Sternberg, M.J.E. (2015) The Phyre2 web portal for protein modeling, prediction and analysis. *Nat. Protoc.*, **10**, 845–858.
38. Enright, A.J., Van Dongen, S. and Ouzounis, C.A. (2002) An efficient algorithm for large-scale detection of protein families. *Nucleic Acids Res.*, **30**, 1575–1584.
39. Steinegger, M., Meier, M., Mirdita, M., Voehringer, H., Haunsberger, S.J. and Soeding, J. (2019) HH-suite3 for fast remote homology detection and deep protein annotation. *BMC Bioinformatics*, **20**, 473.
40. Makarova, K.S. and Koonin, E.V. (2015) Annotation and classification of CRISPR–Cas systems. *Methods Mol. Biol.*, **1311**, 47–75.
41. Makarova, K.S., Aravind, L., Wolf, Y.I. and Koonin, E.V. (2011) Unification of Cas protein families and a simple scenario for the origin and evolution of CRISPR–Cas systems. *Biol. Direct*, **6**, 38.
42. Garrett, R.A., Vestergaard, G. and Shah, S.A. (2011) Archaeal CRISPR-based immune systems: exchangeable functional modules. *Trends Microbiol.*, **19**, 549–556.
43. Carroll, K.S., Gao, H., Chen, H., Stout, C.D., Leary, J.A. and Bertozzi, C.R. (2005) A conserved mechanism for sulfonucleotide reduction. *PLoS Biol.*, **3**, e250.
44. Shah, S.A., Hansen, N.R. and Garrett, R.A. (2009) Distribution of CRISPR spacer matches in viruses and plasmids of crenarchaeal acidothermophiles and implications for their inhibitory mechanism. *Biochem. Soc. Trans.*, **37**, 23–28.
45. Andersson, A.F. and Banfield, J.F. (2008) Virus population dynamics and acquired virus resistance in natural microbial communities. *Science*, **320**, 1047–1050.
46. Ghaly, T.M. and Gillings, M.R. (2018) Mobile DNAs as ecologically and evolutionarily independent units of life. *Trends Microbiol.*, **26**, 904–912.
47. Haft, D.H., Selengut, J., Mongodin, E.F. and Nelson, K.E. (2005) A guild of 45 CRISPR-associated (Cas) protein families and multiple CRISPR/Cas subtypes exist in prokaryotic genomes. *PLoS Comput. Biol.*, **1**, e60.
48. Mojica, F.J.M., Díez-Villaseñor, C., García-Martínez, J. and Almendros, C. (2009) Short motif sequences determine the targets of the prokaryotic CRISPR defence system. *Microbiology*, **155**, 733–740.
49. Shah, S.A., Erdmann, S., Mojica, F.J.M. and Garrett, R.A. (2013) Protospacer recognition motifs: mixed identities and functional diversity. *RNA Biol.*, **10**, 891–899.
50. Westra, E.R., Semenova, E., Datsenko, K.A., Jackson, R.N., Wiedenheft, B., Severinov, K. and Brouns, S.J.J. (2013) Type I-E CRISPR-cas systems discriminate target from non-target DNA through base pairing-independent PAM recognition. *PLoS Genet.*, **9**, e1003742.
51. Jackson, S.A., McKenzie, R.E., Fagerlund, R.D., Kieper, S.N., Fineran, P.C. and Brouns, S.J.J. (2017) CRISPR–Cas: adapting to change. *Science*, **356**, 6333.
52. Yoganand, K.N.R., Sivathanu, R., Nimkar, S. and Anand, B. (2017) Asymmetric positioning of Cas1-2 complex and Integration Host Factor induced DNA bending guide the unidirectional homing of protospacer in CRISPR–Cas type I-E system. *Nucleic Acids Res.*, **45**, 367–381.
53. Tyson, G.W. and Banfield, J.F. (2008) Rapidly evolving CRISPRs implicated in acquired resistance of microorganisms to viruses. *Environ. Microbiol.*, **10**, 200–207.
54. Weinberger, A.D., Sun, C.L., Pluciński, M.M., Denef, V.J., Thomas, B.C., Horvath, P., Barrangou, R., Gilmore, M.S., Getz, W.M. and Banfield, J.F. (2012) Persisting viral sequences shape microbial CRISPR-based immunity. *PLoS Comput. Biol.*, **8**, e1002475.
55. Thompson, C.P., Doak, A.N., Amirani, N., Schroeder, E.A., Wright, J., Kariyawasam, S., Lamendella, R. and Shariat, N.W. (2018) High-Resolution identification of multiple salmonella serovars in a single sample by using CRISPR-SeroSeq. *Appl. Environ. Microbiol.*, **84**, 21.
56. McGinn, J. and Marraffini, L.A. (2016) CRISPR–Cas systems optimize their immune response by specifying the site of spacer integration. *Mol. Cell*, **64**, 616–623.
57. Taylor, H.N., Warner, E.E., Armbrust, M.J., Crowley, V.M., Olsen, K.J. and Jackson, R.N. (2019) Structural basis of Type IV CRISPR RNA biogenesis by a Cas6 endoribonuclease. *RNA Biol.*, **16**, 1438–1447.
58. Seed, K.D., Lazinski, D.W., Calderwood, S.B. and Camilli, A. (2013) A bacteriophage encodes its own CRISPR/Cas adaptive response to evade host innate immunity. *Nature*, **494**, 489–491.
59. Klompe, S.E., Vo, P.L.H., Halpin-Healy, T.S. and Sternberg, S.H. (2019) Transposon-encoded CRISPR–Cas systems direct RNA-guided DNA integration. *Nature*, **571**, 219–225.
60. Strecker, J., Ladha, A., Gardner, Z., Schmid-Burgk, J.L., Makarova, K.S., Koonin, E.V. and Zhang, F. (2019) RNA-guided DNA insertion with CRISPR-associated transposases. *Science*, **365**, 48–53.
61. Al-Shayeb, B., Sachdeva, R., Chen, L.-X., Ward, F., Munk, P., Devoto, A., Castelle, C.J., Olm, M.R., Bouma-Gregson, K., Amano, Y. et al. (2019) Clades of huge phage from across earth's ecosystems. bioRxiv doi: <https://doi.org/10.1101/572362>, 11 March 2019, pre-print: not peer reviewed.
62. Newire, E., Aydin, A., Juma, S., Enne, V. and Roberts, A. (2019) Identification of a Type IV CRISPR–Cas system located exclusively on IncHI1B/ IncFIB plasmids in Enterobacteriaceae. bioRxiv doi: <https://doi.org/10.1101/536375>, 31 January 2019, pre-print: not peer reviewed.
63. McDonald, N.D., Regmi, A., Morreale, D.P., Borowski, J.D. and Boyd, E.F. (2019) CRISPR–Cas systems are present predominantly on mobile genetic elements in Vibrio species. *BMC Genomics*, **20**, 105.
64. Nasko, D.J., Ferrell, B.D., Moore, R.M., Bhavsar, J.D., Polson, S.W. and Wommack, K.E. (2019) CRISPR spacers indicate preferential matching of specific viroplankton genes. *MBio*, **10**, e02651-18.
65. Rostøl, J.T. and Marraffini, L. (2019) (Ph)ighting phages: how bacteria resist their parasites. *Cell Host Microbe*, **25**, 184–194.
66. Paulsson, J. (2002) Multileveled selection on plasmid replication. *Genetics*, **161**, 1373–1384.
67. Deane, S.M. and Rawlings, D.E. (2004) Plasmid evolution and interaction between the plasmid addiction stability systems of two related broad-host-range IncQ-like plasmids. *J. Bacteriol.*, **186**, 2123–2133.
68. Burns, J.H. and Strauss, S.Y. (2011) More closely related species are more ecologically similar in an experimental test. *Proc. Natl. Acad. Sci. U.S.A.*, **108**, 5302–5307.
69. Darwin, C. (1859) On the origin of species by means of natural selection, or the preservation of favoured races in the struggle for life. doi:10.5962/bhl.title.68064.
70. Russel, J., Røder, H.L., Madsen, J.S., Burmølle, M. and Sørensen, S.J. (2017) Antagonism correlates with metabolic similarity in diverse bacteria. *Proc. Natl. Acad. Sci. U.S.A.*, **114**, 10684–10688.
71. Cass, S.D.B., Haas, K.A., Stoll, B., Alkhnbashi, O.S., Sharma, K., Urlaub, H., Backofen, R., Marchfelder, A. and Bolt, E.L. (2015) The role of Cas8 in type I CRISPR interference. *Biosci. Rep.*, **35**, e00197.
72. Mulepati, S., Orr, A. and Bailey, S. (2012) Crystal structure of the largest subunit of a bacterial RNA-guided immune complex and its role in DNA target binding. *J. Biol. Chem.*, **287**, 22445–22449.
73. Deng, L., Garrett, R.A., Shah, S.A., Peng, X. and She, Q. (2013) A novel interference mechanism by a type IIIB CRISPR–Cmr module in *Sulfolobus*. *Mol. Microbiol.*, **87**, 1088–1099.
74. Silas, S., Lucas-Elio, P., Jackson, S.A., Aroca-Crevillén, A., Hansen, L.L., Fineran, P.C., Fire, A.Z. and Sánchez-Amat, A. (2017) Correction: type III CRISPR–Cas systems can provide redundancy to counteract viral escape from type I systems. *Elife*, **6**, e27601.
75. Sinkunas, T., Gasiunas, G., Fremaux, C., Barrangou, R., Horvath, P. and Siksnys, V. (2011) Cas3 is a single-stranded DNA nuclease and ATP-dependent helicase in the CRISPR/Cas immune system. *EMBO J.*, **30**, 1335–1342.
76. Westra, E.R., van Erp, P.B.G., Künne, T., Wong, S.P., Staals, R.H.J., Seegers, C.L.C., Bollen, S., Jore, M.M., Semenova, E., Severinov, K. et al. (2012) CRISPR immunity relies on the consecutive binding and

- degradation of negatively supercoiled invader DNA by Cascade and Cas3. *Mol. Cell*, **46**, 595–605.
77. McRobbie, A.-M., Meyer, B., Rouillon, C., Petrovic-Stojanovska, B., Liu, H. and White, M.F. (2012) Staphylococcus aureus DinG, a helicase that has evolved into a nuclease. *Biochem. J.*, **442**, 77–84.
78. Chou-Zheng, L. and Hatoum-Aslan, A. (2019) A type III-A CRISPR–Cas system employs degradosome nucleases to ensure robust immunity. *Elife*, **8**, e45393
79. Amundsen, S.K., Taylor, A.F. and Smith, G.R. (2002) A domain of RecC required for assembly of the regulatory RecD subunit into the Escherichia coli RecBCD holoenzyme. *Genetics*, **161**, 483–492.
80. Helmrich, A., Ballarino, M., Nudler, E. and Tora, L. (2013) Transcription-replication encounters, consequences and genomic instability. *Nat. Struct. Mol. Biol.*, **20**, 412–418.
81. Wein, T., Hülter, N.F., Mizrahi, I. and Dagan, T. (2019) Emergence of plasmid stability under non-selective conditions maintains antibiotic resistance. *Nat. Commun.*, **10**, 2595.
82. Qi, L.S., Larson, M.H., Gilbert, L.A., Doudna, J.A., Weissman, J.S., Arkin, A.P. and Lim, W.A. (2013) Repurposing CRISPR as an RNA-guided platform for sequence-specific control of gene expression. *Cell*, **152**, 1173–1183.
83. Madsen, J.S., Riber, L., Kot, W., Basfeld, A., Burmølle, M., Hansen, L.H. and Sørensen, S.J. (2016) Type 3 fimbriae encoded on plasmids are expressed from a unique promoter without affecting host motility, facilitating an exceptional phenotype that enhances conjugal plasmid transfer. *PLoS One*, **11**, e0162390.
84. San Millan, A., Toll-Riera, M., Qi, Q., Betts, A., Hopkinson, R.J., McCullagh, J. and MacLean, R.C. (2018) Integrative analysis of fitness and metabolic effects of plasmids in Pseudomonas aeruginosa PAO1. *ISME J.*, **12**, 3014–3024.
85. Venanzio, G.D., Di Venanzio, G., Moon, K.H., Weber, B.S., Lopez, J., Ly, P.M., Potter, R.F., Dantas, G. and Feldman, M.F. (2019) Multidrug-resistant plasmids repress chromosomally encoded T6SS to enable their dissemination. *Proc. Natl. Acad. Sci. U.S.A.*, **116**, 1378–1383.
86. Hawkins, M., Dimude, J.U., Howard, J.A.L., Smith, A.J., Dillingham, M.S., Savery, N.J., Rudolph, C.J. and McGlynn, P. (2019) Direct removal of RNA polymerase barriers to replication by accessory replicative helicases. *Nucleic Acids Res.*, **47**, 5100–5113.
87. Frye, S.A., Beyene, G.T., Namouchi, A., Gómez-Muñoz, M., Homberset, H., Kalayou, S., Riaz, T., Tønjum, T. and Balasingham, S.V. (2017) The helicase DinG responds to stress due to DNA double strand breaks. *PLoS One*, **12**, e0187900.
88. Boubakri, H., de Septenville, A.L., Viguera, E. and Michel, B. (2010) The helicases DinG, Rep and UvrD cooperate to promote replication across transcription units in vivo. *EMBO J.*, **29**, 145–157.
89. Borges, A.L., Davidson, A.R. and Bondy-Denomy, J. (2017) The discovery, mechanisms, and evolutionary impact of Anti-CRISPRs. *Annu. Rev. Virol.*, **4**, 37–59.
90. Wang, L., Jiang, S., Deng, Z., Dedon, P.C. and Chen, S. (2019) DNA phosphorothioate modification—a new multi-functional epigenetic system in bacteria. *FEMS Microbiol. Rev.*, **43**, 109–122.
91. World Health Organization (2014) Antimicrobial Resistance: Global Report on Surveillance.
92. Huddleston, J.R. (2014) Horizontal gene transfer in the human gastrointestinal tract: potential spread of antibiotic resistance genes. *Infect. Drug Resist.*, **7**, 167–176.