

# Role of nucleotide identity in effective CRISPR target escape mutations

Tim Künne<sup>1,2</sup>, Yifan Zhu<sup>1</sup>, Fausia da Silva<sup>1</sup>, Nico Konstantinides<sup>1</sup>, Rebecca E. McKenzie<sup>3</sup>, Ryan N. Jackson<sup>4</sup> and Stan J.J. Brouns<sup>1,3,\*</sup>

<sup>1</sup>Laboratory of Microbiology, Department of Agrotechnology and Food Sciences, Wageningen University, Stippeneng 4, 6708 WE Wageningen, The Netherlands, <sup>2</sup>Laboratory of Food Microbiology, Department of Agrotechnology and Food Sciences, Wageningen University, Bornse Weiland 9, 6708 WG Wageningen, The Netherlands, <sup>3</sup>Kavli Institute of Nanoscience, Department of Bionanoscience, Delft University of Technology, Van der Maasweg 9, 2629 HZ Delft, The Netherlands and <sup>4</sup>Department of Chemistry and Biochemistry, Utah State University, 0300 Old Main Hill, Logan, UT, USA

Received May 18, 2018; Revised July 16, 2018; Editorial Decision July 17, 2018; Accepted August 10, 2018

## ABSTRACT

**Prokaryotes use primed CRISPR adaptation to update their memory bank of spacers against invading genetic elements that have escaped CRISPR interference through mutations in their protospacer target site. We previously observed a trend that nucleotide-dependent mismatches between crRNA and the protospacer strongly influence the efficiency of primed CRISPR adaptation. Here we show that guanine-substitutions in the target strand of the protospacer are highly detrimental to CRISPR interference and interference-dependent priming, while cytosine-substitutions are more readily tolerated. Furthermore, we show that this effect is based on strongly decreased binding affinity of the effector complex Cascade for guanine-mismatched targets, while cytosine-mismatched targets only minimally affect target DNA binding. Structural modeling of Cascade-bound targets with mismatches shows that steric clashes of mismatched guanines lead to unfavorable conformations of the RNA-DNA duplex. This effect has strong implications for the natural selection of target site mutations that lead to effective escape from type I CRISPR–Cas systems.**

## INTRODUCTION

Clustered regularly interspaced short palindromic repeats (CRISPR) together with CRISPR-associated (Cas) proteins provide immunity against foreign nucleic acids in prokaryotes (1,2). The constant battle between prokaryotes and their viruses is one of the oldest and most prominent predator-prey interactions on our planet (3,4). The CRISPR array consists of identical repeat units separated

by unique spacers. In many cases spacer sequences are derived from foreign genetic elements although ‘self’-derived spacers can also be found (5–9). CRISPR–Cas systems are currently divided into class 1 and class 2 systems, encoding multi-subunit or single-subunit crRNA effector complexes, respectively (10–12). Class 1 systems encompass type I, III and IV, while class 2 consists of types II, V and VI. Type I systems are the most widely distributed CRISPR type making up approximately 50% of all CRISPR systems in both Bacteria and Archaea (10). Type I CRISPR–Cas systems contain the universally conserved *cas1* and *cas2* genes, the hallmark *cas3* helicase-nuclease and a set of genes encoding for the Cascade-like effector complexes. The mechanism of CRISPR–Cas defence is divided into three stages: adaptation, expression and interference (13). First, a new spacer is acquired from an invader DNA that has not previously been encountered and is incorporated into the CRISPR array by the Cas1-2 complex (adaptation) (1,14). Next, the whole array is transcribed from the AT-rich leader sequence into long pre-CRISPR RNA (pre-crRNA) and subsequently processed into mature crRNAs that each carry one spacer (expression) (2,15). The crRNA assembles with Cas proteins to form surveillance complexes that make up the core of all CRISPR systems (10). In the last stage (interference), these surveillance complexes scan the cell for complementary targets and flag them for destruction, leading to immunity (16–19). Invaders can escape immunity by acquiring mutations in their recognition sequence (protospacer) or protospacer adjacent motif (PAM), which implies that the host has to acquire a new spacer in order to regain immunity. Several type I systems possess a primed acquisition mechanism that leads to rapid acquisition of new spacers when escape protospacers are detected and interference levels are insufficient to clear a threat (20–24). Unlike naïve acquisition, which requires only *cas1* and *cas2*, primed acquisition requires all *cas* genes and a targeting spacer (21,25).

\*To whom correspondence should be addressed. Tel: +31 15 278 3920; Email: stanbrouns@gmail.com

A number of studies have described the effect of mutations on interference and priming in the type I-E system of *Escherichia coli*. Two early studies have shown, on a small scale, that interference tolerates only few mutations in the protospacer, and no mutations in the seed and PAM, while priming is slightly more tolerant (21,26). Our previous work has extended this knowledge on a large scale, showing that interference tolerates mutations in the seed to a low degree and that priming is extremely robust against mutations in the entire protospacer (27). More recently, it was shown that mutation tolerance of the different immune responses is dependent on the spacer choice (28), and that the binding affinity of Cascade to target DNA is a major determinant for interference and priming (29).

While the number and position of the mutations is clearly important, we previously observed that the identity of an individual protospacer mutation (A, C, G or T) has a large effect on the efficiency of primed spacer acquisition as well (27). Cytosine substitutions in the target strand of the protospacer appeared to positively affect priming, while guanine substitutions negatively affected priming. In contrast, adenine and thymine did not show that trend, suggesting something more complex than a purine or pyrimidine effect.

Here, we show that C and G mutations affect priming by altering the rate of target degradation. We show that, while the overall effect is strongly dependent on the position of the mutations, C mutations repress interference only moderately compared to G mutations at the same protospacer positions. Furthermore, we show that this property is caused by a higher mismatch penalty for G mutations in the target strand of the protospacer compared to C mutations, resulting in lower Cascade binding affinities for mutant targets containing G substitutions. Finally, we use structural modeling to reveal that the molecular basis of this nucleotide bias resides in steric hindrance of mismatched guanines in the target strand of the protospacer.

## MATERIALS AND METHODS

### Bacterial strains and growth conditions

*Escherichia coli* strain KD263 was obtained from (30). *Escherichia coli* strains were grown at 37°C in Luria Broth (LB; 5 g/l NaCl, 5 g/l yeast extract, and 10 g/l tryptone) at 180 rpm or on LB-agar plates containing 1.5% (w/v) agar. When required, medium was supplemented with the following: ampicillin (Amp; 100 µg/ml), chloramphenicol (Cm; 34 µg/ml), or kanamycin (Km; 50 µg/ml). Bacterial growth was measured at 600 nm (OD<sub>600</sub>).

### Molecular biology and DNA sequencing

All oligonucleotides are listed in Supplementary Table S3. All plasmids are listed in Supplementary Table S4. All strains and plasmids were confirmed by PCR and sequencing (GATC-Biotech). Plasmids were prepared using GeneJET Plasmid Miniprep Kits (Thermo Scientific). DNA from PCR was purified using the DNA Clean and Concentrator and Gel DNA Recovery Kit (Zymo Research). The protospacer plasmid set was constructed by cutting pWUR925 with XbaI and SacI, removing the

kanamycin resistance marker, and ligating a PCR product containing the streptomycin resistance marker and the desired protospacer (primers: BG7167/7395-7 for controls, BG7167/8393-8410 for mutant set).

### Plasmid loss assay

The assay was carried out in *E. coli* KD263 cells, which have inducible *cas* gene expression. Expression was induced with 0.2% L-arabinose and 0.5 mM IPTG where appropriate. *Escherichia coli* KD263 cells were transformed with the target plasmids (pWUR926-946) by heat shock. Individual colonies were picked in duplicate and grown overnight in 5 ml LB supplemented with 2% glucose to repress *cas* gene expression. The next day, cultures were transferred 1:100 into induced medium (0.2% L-Arabinose, 0.5 mM IPTG) and plasmid loss was monitored. Samples were taken at the time of induction and 1.5, 3, 4.5, 6, 7, 24 and 48 h post induction (HPI). Dilutions were plated on non-selective plates containing 0.2% rhamnose and plasmid loss was quantified based on loss of red color. Liquid culture samples were screened for spacer integration by colony PCR using OneTaq (NEB). Acquisition of spacers was detected by PCR using primers BG5301 and BG5302. PCR products were visualized on 2% agarose gels and stained with SYBR-safe (Invitrogen). PCR products were sequenced using Sanger sequencing at GATC (Konstantz, Germany) using primer BG5301.

### Protein purification

All proteins were expressed in BL21-AI cells. Cascade was purified as described earlier (31). MBP-Cas3 was purified as described in (32).

### Oligo annealing and labelling

Complementary oligo nucleotides (BG9069-9074) were mixed (1:1) in a Tris-sodium buffer, heated to 95°C and slowly cooled to room temperature. Duplexes were checked on a native 20% acrylamide gel for residual single stranded oligo. The non-target substrate was PCR amplified from pWUR928 using BG9141/2. Duplexes were then labeled with  $\gamma$ -<sup>32</sup>P-ATP using T4 PNK (NEB) and free label was removed using a G25 column.

### EMSA assays

Purified Cascade complex with spacer8 crRNA was incubated with plasmid or oligos at a range of molar ratios (1:1-96:1, Cascade:DNA) in buffer A (20 mM HEPES pH 7.5, 75 mM NaCl, 1 mM DTT) for 30 min (33). Plasmid reactions were run on 1% native agarose gels for 18 h at 22 mA in 8 mM sodium-borate buffer. Gels were post stained with SYBR Safe (Invitrogen). Oligo reactions were run on 5% native acrylamide gels at 4 mA for 18 h. Gels were exposed to a phosphor screen (GE Healthcare) and scanned using a phosphor imager (Bio-Rad PMI). Shifted (Cascade bound DNA) and unshifted (free DNA) bands were quantified using the GeneTools software (Syngene) or ImageJ and free Cascade concentration (*X*) was plotted against the fraction

of bound DNA ( $Y$ ). The curves were fitted with the following formula:  $Y = (\text{amplitude} * X) / (K_d + X)$  (34). The amplitude is the maximum fraction of bound DNA. The affinity ratio is determined as  $\text{amplitude} / K_d$  to correct for the variable amplitudes (35).

### Cas3 DNA degradation assays

Plasmid-based assays were performed by incubating 70 nM Cas3 with 100 nM Cascade and 3.5 nM plasmid DNA. Reactions were incubated in buffer R (5 mM HEPES, pH 8, 60 mM KCl) supplemented with 10  $\mu$ M CoCl<sub>2</sub>, 10 mM MgCl<sub>2</sub>, 2 mM ATP at 37°C for the indicated amount of time. Reactions were quenched on ice with 6 $\times$  DNA loading dye (Thermo scientific). Reactions were run on 0.8% agarose gels at 100 V for 40 min and supercoiled plasmid bands were quantified using the GeneTools software (Syngene).

### Structural modeling of target mismatches

Atomic models of the Cascade complex bound to mismatched DNA targets were made with the molecular modeling program Coot (36). To visualize how G and C mismatches would affect target binding, the G and C mismatches of the C7/G7 target were modeled into the crRNA spacer sequence of dsDNA bound Cascade, (PDB: 5H9E), using the simple mutate tool in Coot. To model wobble basepairs the Rotate Translate Zone/Chain/Molecule tool was used to move nucleotides of the target strand as a rigid bodies into wobble positions. Rendering of atomic model images was performed with PyMOL Molecular Graphics System, Version 2.0 Schödinger, LLC.

## RESULTS

### Statistical scoring of C and G mutants

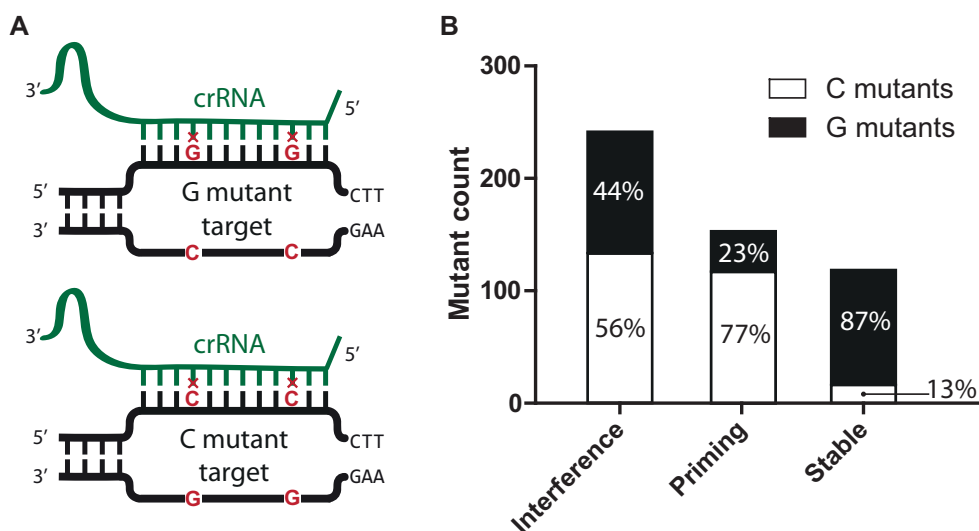
Previously we performed a high throughput plasmid loss assay with a large library of PAM/protospacer mutants, which lead to their classification as causing either (i) interference, (ii) priming or (iii) stable plasmid maintenance (27). Data analysis revealed a nucleotide bias, where C or G mutations have a positive or negative effect on priming, respectively. However, mutants were scored for the number of C or G mutations irrespective of the presence of any other mutations. To verify that the observed effect is purely based on the C and G mutations, we re-analyzed the original dataset, but this time we selected mutants with only C or G mutations to exclude the influence of other mutations (Figure 1A). All analyses were done using effective mutations, thus excluding positions 6, 12, 18, 24, 30 (i.e. kinks) in the crRNA–DNA duplex, which do not participate in base pairing (27,37–39). We included mutants with at least two effective mutations, since single mutants do not show a nucleotide specific effect and most single mutants lead to direct interference (27). Strictly C<sub>n</sub> or G<sub>n</sub> mutants ( $n \geq 2$ , 500 total) were grouped according to their classifications and counted (Figure 1B). The priming group contains mainly C mutants (117 out of 152, 77%), while the stable group contains mostly G mutants (94 out of 108, 87%). This confirms that C mutations generally stimulate priming, while G mutations generally repress priming. Interference has only a

slight preference for C mutants over G mutants (56%/44% for C/G respectively), which suggests either that interference is largely unaffected by the type of mutation or that the nucleotide bias is not pronounced enough in interference because the majority of mutants in this group carry only two effective mutations. We consider it very likely that the type of mutations indeed also affects direct interference, because priming is directly dependent on interference (20,35,40–42).

To address the question to what extent priming and interference are influenced by the type of mutations, and to analyze the effect of the mutations in more detail, we selected four priming protospacers from the dataset with only C mutations (1C–4C) and five stable protospacers with only G mutations (5G–9G) (Supplementary Table S1). The mutants were selected based on two criteria: (i) the mutations had to be effective mutations (i.e. not at kink positions), and (ii) the original nucleotide must be A or T, so that we can switch the mutations from C to G or vice versa without reverting to WT. After selecting the mutants, we designed the respective conversion mutants (1G–4G, 5C–9C).

### G mutants strongly inhibit direct interference

First we performed plasmid loss assays to accurately determine and quantify the ability of the mutant protospacers to trigger direct interference and priming. No plasmid loss and no spacer acquisition was observed with a non-target plasmid after 48 h, showing that CRISPR-independent plasmid loss and naïve acquisition do not occur at detectable levels in this timeframe. When comparing the respective pairs of C and G mutants, we observed that the C mutants consistently showed more rapid plasmid loss than the G mutants (Figure 2A and B, Supplementary Figure S1). Especially, some mutants that were switched from G to C (5C, 7C, 8C) drastically increased their speed of plasmid loss to almost WT levels. Two pairs of mutants show only small differences between the C and G version (6C/G, 9G/C). The original 9G mutant already shows rapid plasmid loss, which simply cannot be increased much more in the 9C mutant. The original 6G mutant is stable and the 6C mutant is only able to show strongly delayed priming. This is likely an additive effect of the individual positions of the mutations (15,20,21,32), which might be more detrimental for interference/priming regardless of nucleotide identity. Position 15 and 21 are in fact the middle positions of their respective segments, which we have previously shown to be more sensitive to mismatches (35). In many of the mutants, significant plasmid loss was observed within 5 hours, indicating that this was caused by direct interference rather than priming (Figure 2A and B). This is supported by the analysis of spacer acquisition showing that spacer acquisition initiated after the onset of plasmid loss (Figure 2A–C). Spacer acquisition also initiated earlier in most C mutants compared to their respective G mutants. The extent of priming, i.e. the fraction of the population that acquired new spacers, on the other hand is not consistent with the type of mutations. For example, the 9G/C mutant pair where the interference is already very high in the G mutant (and even higher in the C mutant) shows opposite behavior with respect to their priming response. Here, the G mutant shows a low level of early priming, while the C mutant shows no



**Figure 1.** Statistics of G/C bias. (A) Schematic representation of G and C mutants in an R-loop with crRNA. The mutation refers to the nucleotide substitution in the basepairing strand of the dsDNA target in the context of the R-loop. (B) Statistical analysis of the high throughput plasmid loss dataset from (27). Only effective mutations, thus excluding positions 6, 12, 18, 24, 30, are considered. Mutants with only C or G mutations (>2) are counted for each group of immune responses (Interference, Priming, Stable).

priming. We observe that the extent of priming is the highest when plasmid loss is occurring at intermediate speeds, while rapid or slow plasmid loss leads to a low extent of priming (Figure 2D). This is very well in line with the model proposed in previous work, i.e. that priming is dependent on inefficient interference that leads to persistence of the invader in the host cell, providing sufficient time for spacer acquisition to take place (35,40,41,43). In summary, mutant protospacers with mismatched guanines in the target strand inhibit CRISPR interference and priming much more than cytosines.

### G mutants inhibit Cas3 degradation rate

To elucidate the molecular basis of this difference in interference of C and G mutants, we first performed Cas3 activity assays with the set of target plasmids as an indicator for the level of direct interference and interference-dependent priming (35). Cas3 assays show a higher average activity of C mutants over G mutants (Figure 3, Supplementary Figure S2). Furthermore, looking at the individual C/G mutant pairs, we see consistently higher activity of the C mutants compared to their corresponding G mutants. This confirms that the more rapid plasmid loss of C mutants is indeed caused by a more efficient plasmid degradation by Cas3.

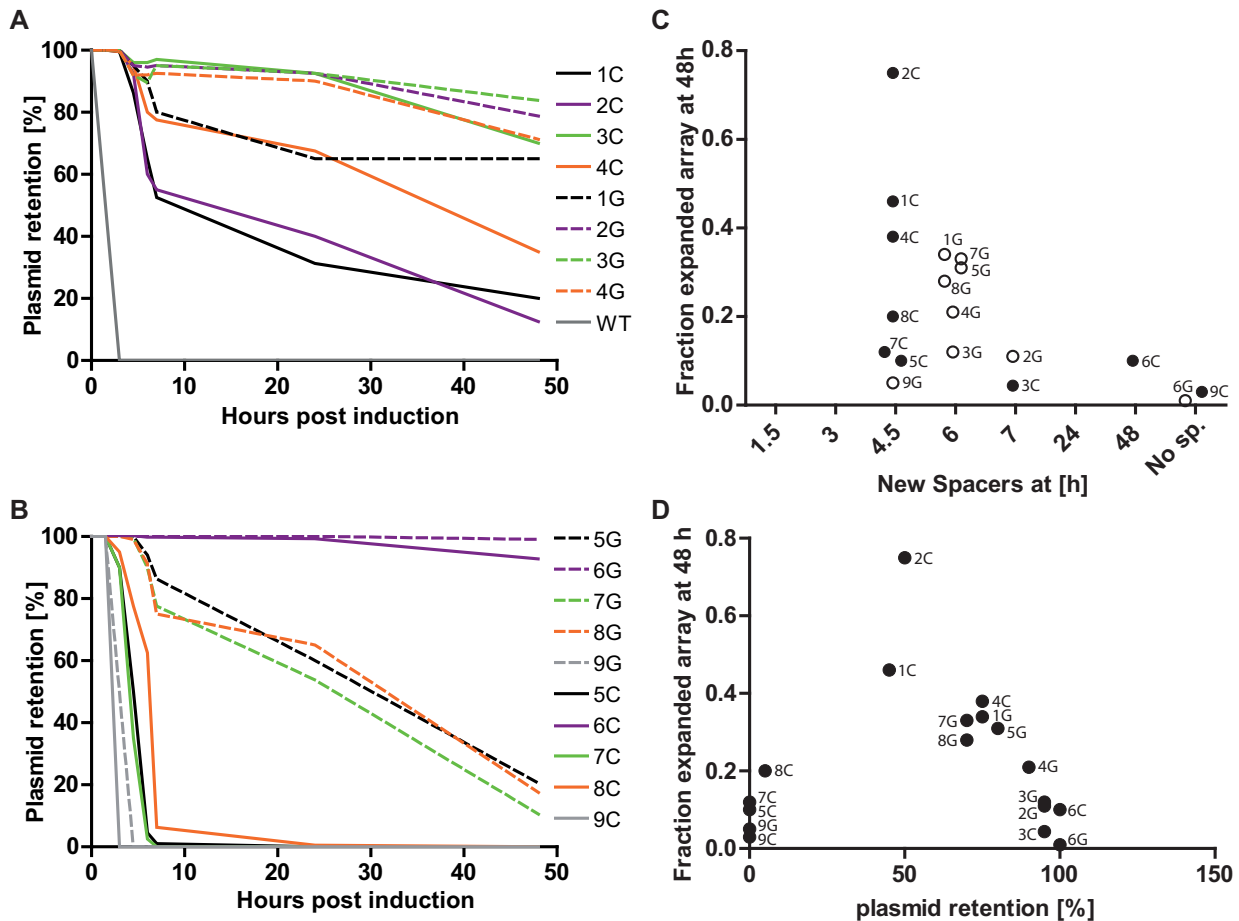
### G mutations in the target strand disrupt Cascade binding

To determine whether the difference in Cas3 activity is caused by affinity differences of Cascade for the mutant protospacers, we performed sensitive DNA binding assays. These assays were set up to assess whether DNA binding affinity is affected only by the mismatches between the crRNA and the DNA target strand, or also by nucleotide preferences of Cascade in the non-target strand. The non-target strand has been proposed to make interactions with the Cse2 dimer and might therefore have an effect on overall

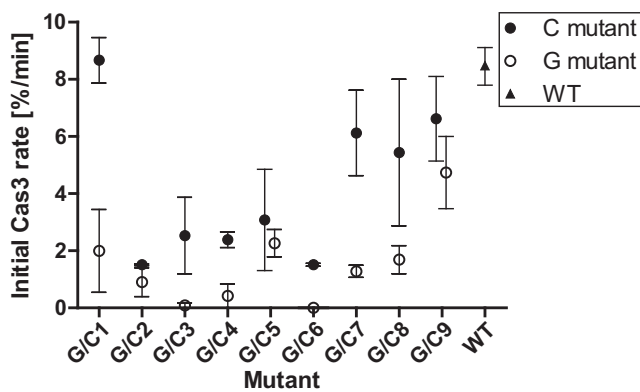
R-loop stability (39,44–47). To address these questions, we designed oligonucleotides for each strand carrying a protospacer with the WT sequence, C mutations or G mutations. The sequences were chosen based on the mutant pair 7C/G. The oligonucleotides were annealed in certain combinations that allow investigation of the effect of the mutations in either strand separately (Figure 4A). We then performed electrophoretic mobility shift assays (EMSA) to measure the binding affinity of Cascade with the five different oligo duplexes (Figure 4B). This revealed that the mutations in the non-target strand have little effect on the affinity and that C mutations on the target strand are readily tolerated. However, G mutations in the target strand disrupt Cascade binding to a much greater extent than C mutations. The affinity ratio (amplitude/ $K_d$ , higher is better) is used as a proxy for binding affinity, since both the amplitude (maximum binding) and the  $K_d$  are variable (35). We observed very similar binding affinity for the C and G mutants of the non-target strand ( $WT_T/C_{NT}$ ,  $WT_T/G_{NT}$ ), the full WT duplex ( $WT_T/WT_{NT}$ ) and the C mutant of the target strand ( $C_T/WT_{NT}$ ) (Figure 4B, Supplementary Table S2). Although the C mutant in the target strand ( $C_T/WT_{NT}$ ) shows a lower amplitude than the others (0.7 versus  $\sim 1$ ), it also has a very low  $K_d$  (13 nM versus 23–46 nM) (Supplementary Table S2). This shows that binding in the  $C_T/WT_{NT}$  mutant is only slightly impaired at higher protein concentrations, while affinity seems unaffected at low protein concentrations. The G mutant in the target strand ( $G_T/WT_{NT}$ ), in sharp contrast, shows almost no detectable binding.

### Steric clashes of mismatched G nucleotides likely distort target conformation

To better understand the molecular basis of why G mutants disrupt Cascade binding more than C mutants, we modeled the G and C mismatches of the C7/G7 target onto the

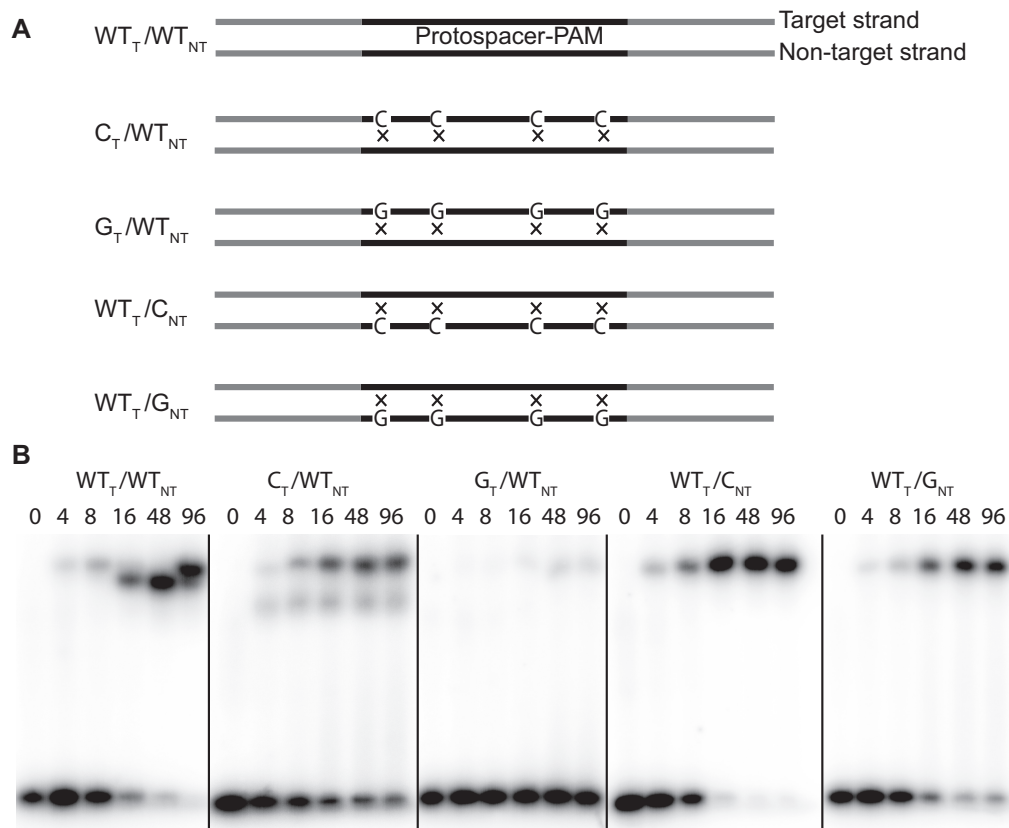


**Figure 2.** Plasmid loss and spacer acquisition. Plasmid loss and spacer acquisition assays of individual mutant plasmids. Two independent assays were carried out, each in duplicate. The standard error of the mean of the measurements can be found in Supplementary Figure S1. **(A)** Plasmid loss curves of C/G pairs 1–4 and the WT. **(B)** Plasmid loss curves of C/G pairs 5–9. **(C)** Analysis of priming during plasmid loss assays, indicating the first occurrence of visible CRISPR array expansion (X-axis) and the extent of priming (fraction of population with expanded array after 48 h, Y-axis). **(D)** Plot showing the relation of speed of plasmid loss (represented by the remaining plasmid after 10 h) with the extent of priming (represented by the fraction of the population with an expanded array after 48 h).



**Figure 3.** *In vitro* Cas3 activity assays. The initial Cas3 reaction rate (over first 10 min) is plotted for each mutant plasmid. The corresponding G/C mutant pairs are shown together. The rate is measured as the percentage of plasmid degradation per minute. The rates are the average of two independent experiments, error bars indicate the individual measurements. C mutants are indicated by full circles, G mutants are indicated by empty circles. Individual Cas3 activity graphs can be found in Supplementary Figure S2.

crystal structure of Cascade bound to a dsDNA (Figure 5A and B) (44). The crRNA of Cascade pairs with complementary DNA with five pseudo-A-form segments of five nucleotides that superimpose with an RMSD of  $\sim 0.45$  Å (34,37–39,44). Thumb domains of adjacent Cas7 subunits flank each base paired segment at six nucleotide intervals, making direct contacts with the crRNA and bound target. Simple replacement of T's for G's to form dG:rA mismatches at spacer positions 4 and 27 reveals steric clashes between the G and A bases, while modeled dC:rA mismatches show no clashes (Figure 5C). A Watson Crick alignment of the dC:rA pair suggests the mismatch cannot form any hydrogen bonds, but lateral motion across the Watson–Crick Face of the dC:rA mismatch may allow a dC:rA+ wobble to form that could possibly stabilize the mismatch. Simple mutation of A nucleotides to G's at spacer positions 13 and 31 initially suggested room for the mismatch. However, electrostatic repulsion between the O-6 position of G and the O-2 position of U likely pushes the G's towards a more thermodynamically stable dG:rU wobble conformation. Our model suggests dG:rU wobbles



**Figure 4.** EMSA with mismatched oligo duplexes. **(A)** Overview of the five oligonucleotide duplex combinations tested. Each combination contains one WT strand and one mutated strand. **(B)** EMSA of Cascade with five different  $^{32}\text{P}$ -labelled oligo duplexes. Cascade to DNA molar ratios are indicated above each lane. Top bands are Cascade bound DNA and bottom bands are free DNA. Vertical lines separate individual gels that have been adjusted to show comparable signal intensity. The affinity ratio (amplitude/ $K_d$ ,  $\text{nM}^{-1}$ ) is given below each gel.

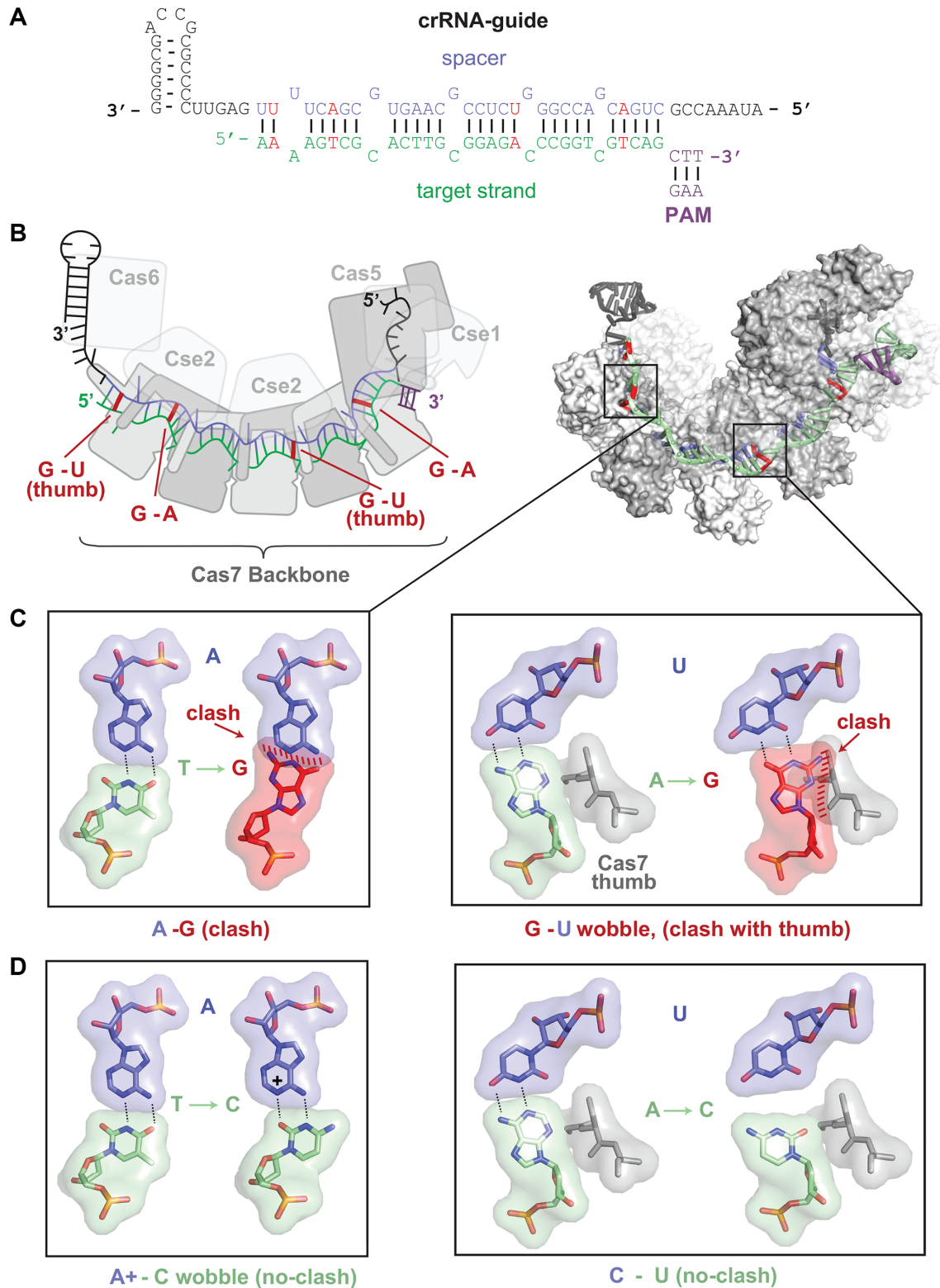
may be allowed through most positions of the duplex, however when the dG:rU wobble conformation is modeled at the first position of the five nucleotide segment a large clash between the G and the adjacent thumb is observed (Figure 5D). Because of electrostatic repulsion between G and U in a Watson–Crick conformation, and the inability to form a stable wobble because of clashes with the Cas7 thumb, dG:rU mismatches adjacent to Cas7 thumbs on the PAM distal side (positions 7, 13, 19, 25 and 30 of the spacer) likely cause distortions in the phosphate backbone of the target. In contrast, dC:rU mismatches at these thumb adjacent positions appear to be able to adopt a normal backbone conformation.

In conclusion, the strict structure and extensive interaction of protein backbone and guide RNA results in a complex range of nucleotide and position specific effects, where a mismatched dG is often detrimental to the stability and alternative basepairing potential of a complex with target DNA.

## DISCUSSION

In this study, we have shown that the type I-E CRISPR–Cas system more readily tolerates cytosine mutations in the target strand of the protospacer, while guanine mutations severely reduce the efficiency of direct interference. This difference is caused by a strong reduction of Cascade bind-

ing affinity of targets with guanine substitutions in the target strand of the protospacer, while the binding affinity is hardly affected in case of cytosine substitutions. The decreased binding affinity results in lower target degradation rates and consequently affects the interference-dependent priming process. The direct effect on interference and priming of the G mutants is also shown by the fact that, although C mutants in all cases lead to earlier priming than G mutants, the extent of priming in the whole population is not directly related to the type of mutation. Instead, the extent of priming follows the model that was conceived in an early study (20) and established in recent studies. These studies showed that the priming process is directly dependent on direct interference (40–42) and that in fact direct interference produces the precursor molecules for new spacers that fuel the priming process (35). However, next to requiring interference for the production of precursors, priming is also dependent on sufficient time of persistence of invader DNA in the cell. Only prolonged persistence gives sufficient opportunity for spacer capture and integration. Thus, a very high rate of direct interference, such as for a WT target or the C9/G9 mutant pair in this study, on the one hand leads to a very early onset of priming, but on the other hand to a very low extent of priming. Mutants with low rates of direct interference lead to late onset of priming and a low extent of priming due to the lack of precursor generation. Mutants



**Figure 5.** Guanine mismatches clash with adenosine bases and the PAM distal side of Cas7 thumbs. (A) Schematic of the crRNA with repeats colored black, spacer colored blue, and complementary DNA sequence colored green. Mismatch positions with the seven C/G target DNA are highlighted in red. (B) Schematic (left) and surface rendition of the crystal structure of Cascade bound to a dsDNA target (pdb 5H9E). Protein subunits are colored grey. The crRNA and target DNA are colored as in A. Cse1, Cse2 and Cas5 subunits are shown as transparent to better visualize the mismatch positions (red). (C) G and C mismatches were modeled into the Cascade structure at 7 C/G mismatch positions. Hydrogen bonds are indicated as dashed lines. Representatives of dG:rA (left) and dG:rU (right) mismatches are shown. Clashes are depicted as red lines. (D). Representatives of dC:rA+ (left) and dC:rU (right) mismatches are shown. No clashes are expected.

with an intermediate rate of direct interference, however, result in relatively early and a high extent of priming, due to prolonged persistence and simultaneous degradation of the invader.

The crRNA-target DNA hybrid in Cascade makes intensive contact with the Cascade protein distorting the duplex away from a classic A-form helix conformation. Therefore, we cannot use available thermodynamic models for determining mismatch energies of A-form RNA-DNA hybrids (48–52). Instead, by modelling mismatched C and G nucleotides into an existing crystal structure of the target bound Cascade complex, we show steric clashes for dG:rA mismatch pairs, while dC:rA pairs fit well in the existing structure. The dC:rA mismatch pair, in addition, might even be stabilized by a dC:rA+ wobble if the bases shift out of the standard Watson Crick alignment. Our model suggests dG:rU mismatch pairs likely form electrostatically favorable wobble pairs. There appears to be room enough around each base pair for such a wobble to form except for positions adjacent to the Cas7 thumbs but distal to the PAM (positions 7, 13, 19, 25, 31). At these positions the G base clashes with the Cas7 thumb. Thus, dG:rA mismatch pairs at every position, and dG:rU pairs at thumb adjacent positions may result in clashes that likely cause a distortion of the phosphate backbone of the DNA target strand to accommodate the G bases. This distortion in the DNA backbone might additionally disrupt basepairing at neighboring positions. Target DNA binding has been shown to initiate at the PAM and progress in a directional zipping process, starting from the PAM proximal end of the protospacer (seed) (53–55). Usually, after complete basepairing of the target DNA, the Cascade complex undergoes a conformational change by shifting its Cse2 subunits, thereby locking the target (54,56,57). This locked state leads to very high affinity binding and is required for interference. Mutations have been shown to prevent Cascade from entering the locked state, but this strongly depends on the amount and position of the mutations (56,57). It is possible that the distortion of the DNA backbone by mismatched G bases either strongly interrupts the zipping process or prevents the Cascade complex to change into the locked conformation, therefore causing the greatly reduced binding affinity that we observed.

The detrimental effect of G mismatches also has consequences for the success of viral escape mutants. Although viruses have been shown in a number of systems to preferably mutate the PAM or seed to escape CRISPR–Cas immunity (26,58–62), systems capable of priming can rapidly regain immunity against these mutants (21–23,27). Thus, only mutants with sufficient mutations to completely escape immunity have increased long-term survival. For these escape viruses it would therefore be beneficial to accumulate G mutations in the targeted strand of their protospacers to maximize the chances of escape. Since the CRISPR–Cas system can target either strand, viruses cannot simply prefer G over C substitutions in general. Instead, the selective pressure on the viruses by type I systems should lead to an overrepresentation of G mutations in protospacers of viruses in natural ecosystems. However, detection of this in natural systems would require very deep sequencing of metagenomes. Alternatively, laboratory co-evolution exper-

iments with host and phage should reveal a biased mutation strategy of viral escapers.

No comparable nucleotide bias has yet been described for other CRISPR–Cas systems, however, it is possible that other type I Cascade-like complexes exhibit similar biases due to rigid crRNA structure, extensive protein–RNA interactions and distorted helix conformation. Type II effector complexes (e.g. Cas9 and Cas12) make intensive protein–RNA interactions with the repeat part of the crRNA and the tracrRNA (63–65). However, the spacer part of the crRNA is flexible and not anchored at the 5' end, allowing uninterrupted basepairing of guide and target by winding of the guide around the target. Although the guide–target duplex in Cas9 also does not form a regular A-form helix, the less rigid nature of the guide RNA likely incorporates mismatched bases more easily. Consequentially, Cas9–target interactions are much more determined by thermodynamics and kinetics rather than structural implications (66,67). Effector complexes in type III systems have a similar overall architecture to Cascade, with a tightly bound crRNA that makes extensive contact with the protein backbone (68). However, type III effector complexes are mechanistically very distinct from Cascade, containing nuclease functionality, targeting RNA, and being able to degrade single stranded DNA non-specifically during transcription (17,69–73). Furthermore, it was recently shown that type III systems have a very high targeting flexibility, requiring viral escapers to lose the entire protospacer (74). Thus, despite structural similarities between type I and type III systems, type III effector complexes can bind to highly mismatched targets efficiently, while type I effector complexes display severe penalties in binding affinity. This apparent disadvantage of type I systems is resolved by the priming mechanism, which can be triggered even by highly mismatched, low affinity targets and restore immunity. In summary, we have demonstrated that, besides well-known positional effects of mutations in protospacers, the identity of a mismatch pair strongly affects CRISPR interference and priming.

## SUPPLEMENTARY DATA

Supplementary Data are available at NAR Online.

## FUNDING

European Research Council [639707]; Netherlands Organization for Scientific Research VIDI grant [864.11.005 to S.J.J.B.]; Technical University Delft start up grant (to S.J.J.B.). Funding for open access charge: ERC starting grant [639707].

*Conflict of interest statement.* None declared.

## REFERENCES

1. Barrangou, R., Fremaux, C., Deveau, H., Richards, M., Boyaval, P., Moineau, S., Romero, D.A. and Horvath, P. (2007) CRISPR provides acquired resistance against viruses in prokaryotes. *Science*, **315**, 1709–1712.
2. Brouns, S.J., Jore, M.M., Lundgren, M., Westra, E.R., Slijkhuys, R.J., Snijders, A.P., Dickman, M.J., Makarova, K.S., Koonin, E.V. and van der Oost, J. (2008) Small CRISPR RNAs guide antiviral defense in prokaryotes. *Science*, **321**, 960–964.



3. Breitbart, M. and Rohwer, F. (2005) Here a virus, there a virus, everywhere the same virus? *Trends Microbiol.*, **13**, 278–284.
4. Rohwer, F. and Thurber, R.V. (2009) Viruses manipulate the marine environment. *Nature*, **459**, 207–212.
5. Bolotin, A., Quinquis, B., Sorokin, A. and Ehrlich, S.D. (2005) Clustered regularly interspaced short palindromic repeats (CRISPRs) have spacers of extrachromosomal origin. *Microbiology*, **151**, 2551–2561.
6. Makarova, K., Grishin, N., Shabalina, S., Wolf, Y. and Koonin, E. (2006) A putative RNA-interference-based immune system in prokaryotes: computational analysis of the predicted enzymatic machinery, functional analogies with eukaryotic RNAi, and hypothetical mechanisms of action. *Biol. Direct*, **1**, 7.
7. Mojica, F.J.M., Díez-Villaseñor, C.S., García-Martínez, J. and Soria, E. (2005) Intervening sequences of regularly spaced prokaryotic repeats derive from foreign genetic elements. *J. Mol. Evol.*, **60**, 174–182.
8. Pourcel, C., Salvignol, G. and Vergnaud, G. (2005) CRISPR elements in *Yersinia pestis* acquire new repeats by preferential uptake of bacteriophage DNA, and provide additional tools for evolutionary studies. *Microbiology*, **151**, 653–663.
9. Stern, A., Keren, L., Wurtzel, O., Amitai, G. and Sorek, R. (2010) Self-targeting by CRISPR: gene regulation or autoimmunity? *Trends Genet.: TIG*, **26**, 335–340.
10. Makarova, K.S., Wolf, Y.I., Alkhnbashi, O.S., Costa, F., Shah, S.A., Saunders, S.J., Barrangou, R., Brouns, S.J., Charpentier, E., Haft, D.H. *et al.* (2015) An updated evolutionary classification of CRISPR–Cas systems. *Nat. Rev. Microbiol.*, **13**, 722–736.
11. Makarova, K.S., Zhang, F. and Koonin, E.V. (2017) SnapShot: class 1 CRISPR–Cas systems. *Cell*, **168**, 946.
12. Makarova, K.S., Zhang, F. and Koonin, E.V. (2017) SnapShot: class 2 CRISPR–Cas systems. *Cell*, **168**, 328.
13. van der Oost, J., Westra, E.R., Jackson, R.N. and Wiedenheft, B. (2014) Unravelling the structural and mechanistic basis of CRISPR–Cas systems. *Nat. Rev. Microbiol.*, **12**, 479–492.
14. Jackson, S.A., McKenzie, R.E., Fagerlund, R.D., Kieper, S.N., Fineran, P.C. and Brouns, S.J. (2017) CRISPR–Cas: adapting to change. *Science*, **356**, eaal5056.
15. Charpentier, E., Richter, H., van der Oost, J. and White, M.F. (2015) Biogenesis pathways of RNA guides in archaeal and bacterial CRISPR–Cas adaptive immunity. *FEMS Microbiol. Rev.*, **39**, 428–441.
16. Garneau, J.E., Dupuis, M.-E., Villion, M., Romero, D.A., Barrangou, R., Boyaval, P., Fremaux, C., Horvath, P., Magadan, A.H. and Moineau, S. (2010) The CRISPR/Cas bacterial immune system cleaves bacteriophage and plasmid DNA. *Nature*, **468**, 67–71.
17. Hale, C.R., Zhao, P., Olson, S., Duff, M.O., Graveley, B.R., Wells, L., Terns, R.M. and Terns, M.P. (2009) RNA-guided RNA cleavage by a CRISPR RNA–Cas protein complex. *Cell*, **139**, 945–956.
18. Marraffini, L.A. and Sontheimer, E.J. (2008) CRISPR interference limits horizontal gene transfer in staphylococci by targeting DNA. *Science*, **322**, 1843–1845.
19. Westra, E.R., van Erp, P.B., Künne, T., Wong, S.P., Staals, R.H., Seegers, C.L., 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.
20. Swarts, D.C., Mosterd, C., van Passel, M.W. and Brouns, S.J. (2012) CRISPR interference directs strand specific spacer acquisition. *PLoS One*, **7**, e35888.
21. Datsenko, K.A., Pougach, K., Tikhonov, A., Wanner, B.L., Severinov, K. and Semenova, E. (2012) Molecular memory of prior infections activates the CRISPR/Cas adaptive bacterial immunity system. *Nat. Commun.*, **3**, 945.
22. Li, M., Wang, R., Zhao, D. and Xiang, H. (2014) Adaptation of the *Haloarcula hispanica* CRISPR–Cas system to a purified virus strictly requires a priming process. *Nucleic Acids Res.*, **42**, 2483–2492.
23. Richter, C., Dy, R.L., McKenzie, R.E., Watson, B.N., Taylor, C., Chang, J.T., McNeil, M.B., Staals, R.H. and Fineran, P.C. (2014) Priming in the Type I-F CRISPR–Cas system triggers strand-independent spacer acquisition, bi-directionally from the primed protospacer. *Nucleic Acids Res.*, **42**, 8516–8526.
24. Vorontsova, D., Datsenko, K.A., Medvedeva, S., Bondy-Denomy, J., Savitskaya, E.E., Pougach, K., Logacheva, M., Wiedenheft, B., Davidson, A.R., Severinov, K. *et al.* (2015) Foreign DNA acquisition by the I-F CRISPR–Cas system requires all components of the interference machinery. *Nucleic Acids Res.*, **43**, 10848–10860.
25. Yosef, I., Goren, M.G. and Qimron, U. (2012) Proteins and DNA elements essential for the CRISPR adaptation process in *Escherichia coli*. *Nucleic Acids Res.*, **40**, 5569–5576.
26. Semenova, E., Jore, M.M., Datsenko, K.A., Semenova, A., Westra, E.R., Wanner, B., van der Oost, J., Brouns, S.J.J. and Severinov, K. (2011) Interference by clustered regularly interspaced short palindromic repeat (CRISPR) RNA is governed by a seed sequence. *Proc. Natl. Acad. Sci. U.S.A.*, **108**, 10098–10103.
27. Fineran, P.C., Gerritzen, M.J., Suarez-Diez, M., Künne, T., Boekhorst, J., van Hijum, S.A., Staals, R.H. and Brouns, S.J. (2014) Degenerate target sites mediate rapid primed CRISPR adaptation. *Proc. Natl. Acad. Sci. U.S.A.*, **111**, E1629–E1638.
28. Xue, C., Seetharam, A.S., Musharova, O., Severinov, K., SJ, J.B., Severin, A.J. and Sashital, D.G. (2015) CRISPR interference and priming varies with individual spacer sequences. *Nucleic Acids Res.*, **43**, 10831–10847.
29. Cooper, L.A., Stringer, A.M. and Wade, J.T. (2018) Determining the specificity of cascade binding, interference, and primed adaptation in vivo in the *Escherichia coli* type I-E CRISPR–Cas system. *mBio*, **9**, e02100-17.
30. Shmakov, S., Savitskaya, E., Semenova, E., Logacheva, M.D., Datsenko, K.A. and Severinov, K. (2014) Pervasive generation of oppositely oriented spacers during CRISPR adaptation. *Nucleic Acids Res.*, **42**, 5907–5916.
31. Jore, M.M., Lundgren, M., van Duijn, E., Bultema, J.B., Westra, E.R., Waghmare, S.P., Wiedenheft, B., Pul, U., Wurm, R., Wagner, R. *et al.* (2011) Structural basis for CRISPR RNA-guided DNA recognition by Cascade. *Nat. Struct. Mol. Biol.*, **18**, 529–536.
32. Mulepati, S. and Bailey, S. (2013) In vitro reconstitution of an *Escherichia coli* RNA-guided immune system reveals unidirectional, ATP-dependent degradation of DNA target. *J. Biol. Chem.*, **288**, 22184–22192.
33. Künne, T., Westra, E.R. and Brouns, S.J. (2015) Electrophoretic mobility shift assay of DNA and CRISPR–Cas ribonucleoprotein complexes. *Methods Mol. Biol.*, **1311**, 171–184.
34. van Erp, P.B., Jackson, R.N., Carter, J., Golden, S.M., Bailey, S. and Wiedenheft, B. (2015) Mechanism of CRISPR–RNA guided recognition of DNA targets in *Escherichia coli*. *Nucleic Acids Res.*, **43**, 8381–8391.
35. Künne, T., Kieper, S.N., Bannenberg, J.W., Vogel, A.I., Mielliet, W.R., Klein, M., Depken, M., Suarez-Diez, M. and Brouns, S.J. (2016) Cas3-derived target DNA degradation fragments fuel primed CRISPR adaptation. *Mol. Cell*, **63**, 852–864.
36. Emsley, P., Lohkamp, B., Scott, W.G. and Cowtan, K. (2010) Features and development of Coot. *Acta Crystallogr. D, Biol. Crystallogr.*, **66**, 486–501.
37. Zhao, H., Sheng, G., Wang, J., Wang, M., Bunkoczi, G., Gong, W., Wei, Z. and Wang, Y. (2014) Crystal structure of the RNA-guided immune surveillance Cascade complex in *Escherichia coli*. *Nature*, **515**, 147–150.
38. Mulepati, S., Heroux, A. and Bailey, S. (2014) Structural biology. Crystal structure of a CRISPR RNA-guided surveillance complex bound to a ssDNA target. *Science*, **345**, 1479–1484.
39. Jackson, R.N., Golden, S.M., van Erp, P.B., Carter, J., Westra, E.R., Brouns, S.J., van der Oost, J., Terwilliger, T.C., Read, R.J. and Wiedenheft, B. (2014) Structural biology. Crystal structure of the CRISPR RNA-guided surveillance complex from *Escherichia coli*. *Science*, **345**, 1473–1479.
40. Semenova, E., Savitskaya, E., Musharova, O., Strotskaya, A., Vorontsova, D., Datsenko, K.A., Logacheva, M.D. and Severinov, K. (2016) Highly efficient primed spacer acquisition from targets destroyed by the *Escherichia coli* type I-E CRISPR–Cas interfering complex. *Proc. Natl. Acad. Sci. U.S.A.*, **113**, 7626–7631.
41. Staals, R.H., Jackson, S.A., Biswas, A., Brouns, S.J., Brown, C.M. and Fineran, P.C. (2016) Interference-driven spacer acquisition is dominant over naive and primed adaptation in a native CRISPR–Cas system. *Nat. Commun.*, **7**, 12853.
42. Krivoy, A., Rutkauskas, M., Kuznedelov, K., Musharova, O., Rouillon, C., Severinov, K. and Seidel, R. (2018) Primed CRISPR adaptation in *Escherichia coli* cells does not depend on conformational changes in the Cascade effector complex detected in Vitro. *Nucleic Acids Res.*, **46**, 4087–4098.

43. Severinov, K., Isolatov, I. and Semenova, E. (2016) The influence of Copy-Number of targeted extrachromosomal genetic elements on the outcome of CRISPR–Cas defense. *Front. Mol. Biosci.*, **3**, 45.
44. Hayes, R.P., Xiao, Y., Ding, F., van Erp, P.B., Rajashankar, K., Bailey, S., Wiedenheft, B. and Ke, A. (2016) Structural basis for promiscuous PAM recognition in type I-E Cascade from *E. coli*. *Nature*, **530**, 499–503.
45. Nam, K.H., Huang, Q. and Ke, A. (2012) Nucleic acid binding surface and dimer interface revealed by CRISPR-associated CasB protein structures. *FEBS Lett.*, **586**, 3956–3961.
46. Wiedenheft, B., Lander, G.C., Zhou, K., Jore, M.M., Brouns, S.J., van der Oost, J., Doudna, J.A. and Nogales, E. (2011) Structures of the RNA-guided surveillance complex from a bacterial immune system. *Nature*, **477**, 486–489.
47. Xiao, Y., Luo, M., Hayes, R.P., Kim, J., Ng, S., Ding, F., Liao, M. and Ke, A. (2017) Structure basis for directional R-loop formation and substrate handover mechanisms in type I CRISPR–Cas system. *Cell*, **170**, 48–60.
48. Huang, Y., Chen, C. and Russu, I.M. (2009) Dynamics and stability of individual base pairs in two homologous RNA–DNA hybrids. *Biochemistry*, **48**, 3988–3997.
49. Sugimoto, N., Nakano, M. and Nakano, S. (2000) Thermodynamics–structure relationship of single mismatches in RNA/DNA duplexes. *Biochemistry*, **39**, 11270–11281.
50. Watkins, J.N.E., Kennelly, W.J., Tsay, M.J., Tuin, A., Swenson, L., Lee, H.-R., Morosyuk, S., Hicks, D.A. and SantaLucia, J.J. (2011) Thermodynamic contributions of single internal rA–dA, rC–dC, rG–dG and rU–dT mismatches in RNA/DNA duplexes. *Nucleic Acids Res.*, **39**, 1894–1902.
51. Wu, P., Nakano, S. and Sugimoto, N. (2002) Temperature dependence of thermodynamic properties for DNA/DNA and RNA/DNA duplex formation. *Eur. J. Biochem.*, **269**, 2821–2830.
52. Zhu, J. and Wartell, R.M. (1999) The effect of base sequence on the stability of RNA and DNA single base bulges. *Biochemistry*, **38**, 15986–15993.
53. Xue, C., Zhu, Y., Zhang, X., Shin, Y.K. and Sashital, D.G. (2017) Real-Time observation of target search by the CRISPR surveillance complex cascade. *Cell Rep.*, **21**, 3717–3727.
54. Rutkauskas, M., Sinkunas, T., Songailiene, I., Tikhomirova, M.S., Siksnys, V. and Seidel, R. (2015) Directional R-Loop formation by the CRISPR–Cas surveillance complex cascade provides efficient off-target site rejection. *Cell Rep.*, **10**, 1534–1543.
55. Szczelkun, M.D., Tikhomirova, M.S., Sinkunas, T., Gasiunas, G., Karvelis, T., Pschera, P., Siksnys, V. and Seidel, R. (2014) Direct observation of R-loop formation by single RNA-guided Cas9 and Cascade effector complexes. *Proc. Natl. Acad. Sci. U.S.A.*, **111**, 9798–9803.
56. Xue, C., Whitis, N.R. and Sashital, D.G. (2016) Conformational control of cascade interference and priming activities in CRISPR immunity. *Mol. Cell*, **64**, 826–834.
57. Blosser, T.R., Loeff, L., Westra, E.R., Vlot, M., Künne, T., Sobota, M., Dekker, C., Brouns, S.J. and Joo, C. (2015) Two distinct DNA binding modes guide dual roles of a CRISPR–Cas protein complex. *Mol. Cell*, **58**, 60–70.
58. Deveau, H., Barrangou, R., Garneau, J.E., Labonte, J., Fremaux, C., Boyaval, P., Romero, D.A., Horvath, P. and Moineau, S. (2008) Phage response to CRISPR-encoded resistance in *Streptococcus thermophilus*. *J. Bacteriol.*, **190**, 1390–1400.
59. Kupczok, A. and Bollback, J.P. (2014) Motif depletion in bacteriophages infecting hosts with CRISPR systems. *BMC Genomics*, **15**, 663.
60. Box, A.M., McGuffie, M.J., O'Hara, B.J. and Seed, K.D. (2015) Functional analysis of bacteriophage immunity through a type I-E CRISPR–Cas system in *Vibrio cholerae* and its application in bacteriophage genome engineering. *J. Bacteriol.*, **198**, 578–590.
61. Paez-Espino, D., Sharon, I., Morovic, W., Stahl, B., Thomas, B.C., Barrangou, R. and Banfield, J.F. (2015) CRISPR immunity drives rapid phage genome evolution in *Streptococcus thermophilus*. *mBio*, **6**, e00262-15.
62. Künne, T., Swarts, D.C. and Brouns, S.J. (2014) Planting the seed: target recognition of short guide RNAs. *Trends Microbiol.*, **22**, 74–83.
63. Anders, C., Niewoehner, O., Duerst, A. and Jinek, M. (2014) Structural basis of PAM-dependent target DNA recognition by the Cas9 endonuclease. *Nature*, **513**, 569–573.
64. Nishimasu, H., Ran, F.A., Hsu, P.D., Konermann, S., Shehata, S.I., Dohmae, N., Ishitani, R., Zhang, F. and Nureki, O. (2014) Crystal structure of Cas9 in complex with guide RNA and target DNA. *Cell*, **156**, 935–949.
65. Swarts, D.C., van der Oost, J. and Jinek, M. (2017) Structural basis for guide RNA processing and seed-dependent DNA targeting by CRISPR–Cas12a. *Mol. Cell*, **66**, 221–233.
66. Farasat, I. and Salis, H.M. (2016) A biophysical model of CRISPR/Cas9 activity for rational design of genome editing and gene regulation. *PLoS Comput. Biol.*, **12**, e1004724.
67. Klein, M., Eslami-Mossallam, B., Arroyo, D.G. and Depken, M. (2018) Hybridization kinetics explains CRISPR–Cas off-targeting rules. *Cell Rep.*, **22**, 1413–1423.
68. Jackson, R.N. and Wiedenheft, B. (2015) A conserved structural chassis for mounting versatile CRISPR RNA-Guided immune responses. *Mol. Cell*, **58**, 722–728.
69. Elmore, J.R., Sheppard, N.F., Ramia, N., Deighan, T., Li, H., Terns, R.M. and Terns, M.P. (2016) Bipartite recognition of target RNAs activates DNA cleavage by the Type III-B CRISPR–Cas system. *Genes Dev.*, **30**, 447–459.
70. Estrella, M.A., Kuo, F.T. and Bailey, S. (2016) RNA-activated DNA cleavage by the Type III-B CRISPR–Cas effector complex. *Genes Dev.*, **30**, 460–470.
71. Jiang, W., Samai, P. and Marraffini, L.A. (2016) Degradation of phage transcripts by CRISPR-Associated RNases enables type III CRISPR–Cas immunity. *Cell*, **164**, 710–721.
72. Kazlauskienė, M., Tamulaitis, G., Kostiuk, G., Venclovas, C. and Siksnys, V. (2016) Spatiotemporal control of type III-A CRISPR–Cas immunity: coupling DNA degradation with the target RNA recognition. *Mol. Cell*, **62**, 295–306.
73. Samai, P., Pyenson, N., Jiang, W., Goldberg, G.W., Hatoum-Aslan, A. and Marraffini, L.A. (2015) Co-transcriptional DNA and RNA cleavage during type III CRISPR–Cas immunity. *Cell*, **161**, 1164–1174.
74. Pyenson, N.C., Gayvert, K., Varble, A., Elemento, O. and Marraffini, L.A. (2017) Broad targeting specificity during bacterial type III CRISPR–Cas immunity constrains viral escape. *Cell Host Microbe*, **22**, 343–353.