

## Cas4-Cas1 Is a Protospacer Adjacent Motif-Processing Factor Mediating Half-Site Spacer Integration during CRISPR Adaptation

Kieper, S.N.; Almendros Romero, C.; van Eijkeren-Haagsma, A.C.; Barendregt, Arjan; Heck, Albert J.R.; Brouns, S.J.J.

**DOI**

[10.1089/crispr.2021.0011](https://doi.org/10.1089/crispr.2021.0011)

**Publication date**

2021

**Document Version**

Accepted author manuscript

**Published in**

CRISPR Journal

**Citation (APA)**

Kieper, S. N., Almendros Romero, C., van Eijkeren-Haagsma, A. C., Barendregt, A., Heck, A. J. R., & Brouns, S. J. J. (2021). Cas4-Cas1 Is a Protospacer Adjacent Motif-Processing Factor Mediating Half-Site Spacer Integration during CRISPR Adaptation. *CRISPR Journal*, 4(4), 536-548. <https://doi.org/10.1089/crispr.2021.0011>

**Important note**

To cite this publication, please use the final published version (if applicable). Please check the document version above.

**Copyright**

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

**Takedown policy**

Please contact us and provide details if you believe this document breaches copyrights. We will remove access to the work immediately and investigate your claim.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1 Classification: BIOLOGICAL SCIENCES: Microbiology

2 Research article

3  
4 **Cas4-Cas1 is a PAM-processing factor mediating half-site spacer integration**  
5 **during CRISPR adaptation**

6  
7 Sebastian N. Kieper<sup>1,2</sup>, Cristobal Almendros<sup>1,2</sup>, Anna C. Haagsma<sup>1,2</sup>, Arjan  
8 Barendregt<sup>3,4</sup>, Albert J.R. Heck<sup>3,4</sup>, Stan J.J. Brouns<sup>1,2§</sup>

9  
10 <sup>1</sup> Department of Bionanoscience, Delft University of Technology, Delft, Netherlands

11 <sup>2</sup> Kavli Institute of Nanoscience, Delft, Netherlands.

12 <sup>3</sup> Biomolecular Mass Spectrometry and Proteomics, Bijvoet Center for Biomolecular  
13 Research, Utrecht Institute of Pharmaceutical Sciences, Utrecht University, Utrecht,  
14 Netherlands.

15 <sup>4</sup> Netherlands Proteomics Center, Utrecht, Netherlands.

16  
17 § Corresponding author: Brouns, S.J.J. (stanbrouns@gmail.com, Tel +31 15 278  
18 3920)

19  
20 Keywords: CRISPR adaptation, Cas4, Spacer acquisition, PAM selection

## 22 **Abstract**

23 The immunization of bacteria and archaea against invading viruses via CRISPR  
24 adaptation is critically reliant on the efficient capture, accurate processing and  
25 integration of CRISPR spacers into the host genome. The adaptation proteins Cas1  
26 and Cas2 are sufficient for successful spacer acquisition in some CRISPR-Cas  
27 systems. However, many CRISPR-Cas systems additionally require the Cas4 protein  
28 for efficient adaptation. Cas4 has been implied in selection and processing of spacer  
29 precursors, but the detailed mechanistic understanding of how Cas4 contributes to  
30 CRISPR adaptation is lacking. Here we biochemically reconstitute the CRISPR-Cas  
31 type I-D adaptation system and show two functionally distinct adaptation complexes,  
32 Cas4-Cas1 and Cas1-Cas2. The Cas4-Cas1 complex recognizes and cleaves PAM  
33 sequences in 3' overhangs in a sequence-specific manner, while the Cas1-Cas2  
34 complex defines the cleavage of non-PAM sites via host factor nucleases. Both sub-  
35 complexes are capable of mediating half-site integration, facilitating the integration of  
36 processed spacers in the correct, interference-proficient orientation. We provide a  
37 model in which an asymmetric adaptation complex differentially acts on PAM and non-  
38 PAM containing overhangs, providing cues for the correct orientation of spacer  
39 integration.

## 41 **Introduction**

42 Clustered Regularly Interspaced Short Palindromic Repeats (CRISPR) and their  
43 associated genes (*cas*) provide adaptive and inheritable immunity against mobile  
44 genetic elements (MGEs) in bacteria and archaea (Nussenzweig and Marraffini, 2020).  
45 The CRISPR array is composed of palindromic repeats interspersed by sequences  
46 derived from MGEs and serves as a template for the biogenesis of CRISPR RNAs  
47 (crRNAs) (Barrangou, 2013; van der Oost et al., 2014). Cas proteins subsequently  
48 assemble around the crRNA to form effector complexes that mediate the recognition  
49 and destruction of invading MGEs that have been recorded in the bacterial genome  
50 during previous infections (Brouns et al., 2008). Therefore, the main requirement for  
51 the establishment of immunity is the memorization of foreign genetic material in a step  
52 called CRISPR adaptation (Jackson et al., 2017; McGinn and Marraffini, 2019). The  
53 core machinery responsible for adaptation is composed of the Cas1 and Cas2 proteins  
54 that assemble into the adaptation complex (Nuñez et al., 2015a; Nuñez et al., 2014;  
55 Yosef et al., 2012). Among the first identified *cas* genes were the adaptation genes  
56 *cas1* and *cas2*, the *cas3* gene encoding the nuclease-helicase Cas3 and the *cas4*  
57 gene encoding a protein, the function of which has been unknown until recently (Hou  
58 and Zhang, 2018; Jansen et al., 2002). The *cas4* gene is widespread among several  
59 sub-types of type I, type II and type V systems and therefore present in the majority of  
60 CRISPR-Cas systems (Hudaiberdiev et al., 2017). Predictions of Cas4 function existed  
61 early on, based on the frequent colocalization of the CRISPR adaptation genes *cas1*  
62 and *cas2* and the *cas4* gene. This co-localization suggested that Cas4 could be  
63 contributing to the adaptation stage and the early studies indeed found supporting  
64 evidence for this hypothesis: Adaptation of the type I-A system of *Sulfolobus islandicus*  
65 was severely impaired upon deletion of *cas4* (Liu et al., 2017). Similarly, deleting *cas4*  
66 from the type I-B system of *Haloarcula hispanica* abrogated CRISPR adaptation

1  
2  
3 67 against HHPV-2 (Li et al., 2014). Biochemical evidence was provided by Plagens et al.  
4  
5 68 showing a protein-protein interaction *in vitro* between Cas4 and the type I-A adaptation  
6  
7 69 fusion protein Cas1/2 and Csa1 demonstrating that Cas4 directly interacts with the  
8  
9 70 adaptation machinery (Plagens et al., 2012). Recently, several studies defined the role  
10  
11 71 of Cas4 in more detail, finding that the presence of Cas4 increases the fidelity of spacer  
12  
13 72 integration. Specifically, the Cas4 protein of a cyanobacterial type I-D CRISPR-Cas  
14  
15 73 system facilitated the integration of interference-proficient spacers that carry the  
16  
17 74 consensus PAM of the type I-D CRISPR-Cas system. Spacers acquired in the  
18  
19 75 presence of Cas4 displayed shorter lengths compared to those acquired in the  
20  
21 76 absence of Cas4 or in the presence of catalytically inactive variant of the protein  
22  
23 77 (Kieper et al., 2018). Additionally, two Cas4 variants (Cas4-1 and Cas4-2) encoded in  
24  
25 78 the type I-A system of *Pyrococcus furiosus* were shown to define the upstream  
26  
27 79 protospacer adjacent motif (PAM) and a downstream NW motif *in vivo* (Shiimori et al.,  
28  
29 80 2018). Deletion of *cas4-1* and *cas4-2* resulted in incorrect processing of prespacers  
30  
31 81 with respect to up- and downstream motifs, random orientation of the integrated spacer  
32  
33 82 as well as large deviations from the consensus spacer length (Shiimori et al., 2018).  
34  
35 83 Previously, biochemical studies of the Cas4 protein, containing an iron-sulfur cluster  
36  
37 84 and a RecB domain, found several nuclease activities, demonstrating endo- and  
38  
39 85 exonuclease activities (Lemak et al., 2013; Lemak et al., 2014; Zhang et al., 2012).  
40  
41 86 The requirement of Cas4 for prespacer processing was therefore in accordance with  
42  
43 87 the previously described biochemical activities. Indeed, Lee et al. provided the first  
44  
45 88 mechanistic details of how Cas4 proteins ensure PAM-processing and correct spacer  
46  
47 89 orientation (Lee et al., 2019; Lee et al., 2018). It was shown that Cas4 tightly interacts  
48  
49 90 with the Cas1 integrase, forming a heterohexameric complex composed of two Cas1  
50  
51 91 dimers and two Cas4 subunits (Lee et al., 2018). This complex would interact with  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 92 double-stranded prespacer substrates and endonucleolytically cleave PAM sequences  
4  
5 93 in long 3' overhangs, ensuring that only PAM-processed spacers would be eventually  
6  
7 94 integrated into the CRISPR array (Lee et al., 2018). Interestingly, the authors did not  
8  
9  
10 95 see interaction of the Cas4-Cas1 complex with Cas2 in their initial experiments,  
11  
12 96 suggesting the possibility that a prespacer is required for assembly of the full complex.  
13  
14 97 After supplying a dsDNA substrate to the three adaptation proteins, Lee et al. could  
15  
16  
17 98 demonstrate the assembly of the full Cas4-Cas1-Cas2 complex (Lee et al., 2019). This  
18  
19 99 complex was shown to assemble in a mixture of symmetric and asymmetric  
20  
21 100 architectures as shown by negative-staining Electron Microscopy, in which the  
22  
23 101 asymmetric complex would contain only a Cas4 monomer associated with one of the  
24  
25  
26 102 Cas1 dimers (Lee et al., 2019). The authors suggested that this asymmetry might aid  
27  
28 103 in the differential processing of prespacer substrates in which the Cas4 containing half  
29  
30 104 of the complex would interact with the PAM-containing overhang. This hypothesis is  
31  
32  
33 105 supported by the findings in the type I-A system, in which two independent Cas4  
34  
35 106 homologs are dedicated processing factors for the PAM- and NW motif containing  
36  
37 107 prespacer overhangs (Shiimori et al., 2018). However, how this asymmetric processing  
38  
39 108 is orchestrated in CRISPR systems containing only a single *cas4* gene is currently  
40  
41  
42 109 unknown. In this work we provide mechanistic insights of an asymmetric complex,  
43  
44 110 specifically how the Cas4-Cas1 complex is able to recognize and sequence specifically  
45  
46  
47 111 process the PAM sequence of the type I-D CRISPR-Cas system. Previously we have  
48  
49 112 shown that the type I-D Cas4 protein facilitates the integration of PAM-compliant  
50  
51 113 spacers *in vivo* (Kieper et al., 2018). We demonstrate that Cas4 strongly interacts with  
52  
53 114 the Cas1 integrase forming a heteromeric Cas4<sub>1</sub>-Cas1<sub>2</sub> complex. This heteromeric  
54  
55  
56 115 complex does not require the Cas2 protein for processing and half-site integration of  
57  
58 116 PAM-containing prespacer substrates. The catalytic activity of Cas4 is required for  
59  
60

1  
2  
3 117 prespacer cleavage and is crucially dependent on the presence of Cas1 in order to  
4  
5 118 recognize and process the PAM overhang. We show that this Cas4-Cas1 complex  
6  
7 119 does not cleave the non-PAM containing overhang. Processing of the non-PAM  
8  
9  
10 120 containing overhang potentially relies on the Cas1-Cas2 complex and likely requires  
11  
12 121 host-factor nucleases. We provide a model in which an asymmetric adaptation  
13  
14 122 complex differentially acts on PAM and non-PAM containing overhangs, providing cues  
15  
16  
17 123 for the correct orientation of spacer integration. This correct PAM processing as well  
18  
19 124 as a functional orientation explains the importance and hence the strong conservation  
20  
21 125 of the *cas4* gene, increasing the integration of interference-proficient spacers.  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## 127 **Material & Methods**

128

### 129 **Bacterial strains and growth conditions**

130 *E. coli* strains DH5 $\alpha$  and BL21 were grown in Lysogeny Broth (LB) at 37°C and  
131 continuous shaking at 180 rpm or grown on LB agar plates (LBA) containing 1.5%  
132 (wt/vol) agar. When required, the media were supplemented with 100  $\mu\text{g ml}^{-1}$   
133 ampicillin, 50  $\mu\text{g ml}^{-1}$  spectinomycin, 25  $\mu\text{g ml}^{-1}$  chloramphenicol (see Table S1 for  
134 plasmids and corresponding selection markers).

### 135 **Plasmid construction and transformation**

136 Plasmids used in this study are listed in Table S1. All cloning steps were performed in  
137 *E. coli* DH5 $\alpha$ . Primers described in Table S2 were used for PCR amplification of the  
138 type I-D CRISPR-Cas locus (*cas4*, *cas1*, *cas2* and leader-repeat-spacer1) from  
139 *Synechocystis* cell material using the Q5 high-fidelity Polymerase (New England  
140 Biolabs). PCR amplicons were subsequently cloned into Berkeley MacroLab LIC  
141 vectors (<https://qb3.berkeley.edu/facility/qb3-macrolab/>) using either ligation-  
142 independent cloning (LIC), or into the pACYCDuet-1 vector system (Novagen (EMD  
143 Millipore) using conventional restriction-ligation cloning. The *cas4*<sup>D76A+K91A</sup> mutant was  
144 obtained using a PCR-based mutagenesis of pCas4<sup>D76A</sup> using primers listed in Table  
145 S2. All plasmids were verified by Sanger-sequencing (Macrogen Europe, Amsterdam,  
146 The Netherlands). Bacterial transformations were either carried out by electroporation  
147 (2.5 kV, 25 mF, 200 V) using an ECM 630 electroporator (BTX Harvard Apparatus) or  
148 chemically competent cells prepared according to manufacturer's manual (Mix&Go,  
149 Zymo research). Electrocompetent cells were prepared following a protocol adapted  
150 from (Gonzales et al., 2013). Transformants were selected on LBA supplemented with  
151 appropriate antibiotics.



### 152 **Protein expression and purification**

153 Plasmid encoded *cas* genes were either co-expressed or expressed individually in *E.*  
154 *coli* BL21 AI cells (Invitrogen). Pre-cultures were grown from individual colonies and  
155 used for inoculation of pre-warmed (37°C) LB medium at an initial OD<sub>600</sub>=0.05. Protein  
156 expression was induced at OD<sub>600</sub>=0.5 by addition of IPTG and L-arabinose preceded  
157 by a 30-minute cold-shock. Cultures were subsequently grown overnight at 20°C with  
158 continuous shaking. Cells were harvested by centrifugation (10 min, 4°C, 2400xg) and  
159 subsequently resuspended in lysis buffer (50 mM HEPES pH 7.5, 300 mM KCl, 5%  
160 Glycerol, 1 mM DTT, 25 mM Imidazole, 0.1% Triton-X 100) supplemented with  
161 cComplete™, EDTA-free Protease Inhibitor Cocktail (Roche). Cells were lysed by two  
162 passages through a CF1 cell disruptor (Constant Systems Ltd.) equilibrated with lysis  
163 buffer at a constant pressure of 1 kbar. Lysates were cleared by centrifugation (45 min,  
164 4°C, 25000xg) and filtered through 0.45 µm filter. Protein was bound in batch to HIS-  
165 Select (Sigma Aldrich) IMAC resin for 30 min at 4°C and rotary shaking. IMAC resin  
166 was then loaded onto Pierce gravity-flow columns (Thermo Scientific) and washed with  
167 10 CV wash buffer (50 mM HEPES pH 7.5, 300 mM KCl, 5% Glycerol, 1 mM DTT, 50  
168 mM Imidazole). Proteins were subsequently block eluted in 0.5 ml elution buffer (50  
169 mM HEPES pH 7.5, 300 mM KCl, 5% Glycerol, 1 mM DTT, 250 mM Imidazole). Protein  
170 concentration and purity was determined by NanoDrop A280 spectroscopy and SDS  
171 PAGE analysis. Protein elution fractions were pooled and subjected to size exclusion  
172 chromatography using Superdex 200 10/300 GL (GE Healthcare) column with 0.5  
173 ml/min flow rate using elution buffer as mobile phase. Cas1-Cas2 complex IMAC  
174 elution fractions used for integration assays were prepared for ion-exchange  
175 chromatography by adjusting the KCL concentration to 30 mM and subsequently  
176 loaded onto HiTrap Heparin HP column (GE Healthcare). Cas1-Cas2 complexes  
177 were eluted by gradually increasing KCL concentration to 1 M. Resulting fractions

1  
2  
3 178 were analyzed by SDS-PAGE and appropriate fractions pooled, snap frozen and  
4  
5 179 stored at -80°C.  
6  
7

### 8 180 **Native Mass Spectrometry**

9 181 Cas4-Cas1 and Cas1-Cas2 complexes were buffer exchanged into 500mM  
10  
11  
12  
13 182 ammonium acetate (pH 7.5) using seven sequential steps on a centrifugal filter with  
14  
15  
16  
17 183 a molecular weight cut-off of 10 kDa (Sartorius) at 4°C. MS measurements were  
18  
19  
20  
21 184 performed in positive mode by directly infusing the individual complexes at a  
22  
23  
24 185 concentration of 1  $\mu$ M using an LCT electrospray time-of-flight (Waters, United  
25  
26  
27  
28 186 Kingdom) adjusted for optimal performance in high mass detection (Tahallah et al.,  
29  
30  
31 187 2001; van den Heuvel et al., 2006). The needles used for electrospray were  
32  
33  
34  
35 188 prepared in house from borosilicate capillaries (Kwik-Fil, World Precision  
36  
37  
38  
39 189 Instruments, Sarasota, FL) on a P97 puller (SutterInstruments, Novato, USA) and  
40  
41  
42 190 gold coated by using an Edwards Scancoat Six Pirani 501 Sputter Coater (Edwards  
43  
44  
45  
46 191 Laboratories, Milpitas, USA). During the measurement the capillary voltage was  
47  
48  
49  
50 192 kept at 1200V, cone voltage between 80-150V and the source pressure was  
51  
52  
53 193 increased to  $\approx$ 8mbar. Exact mass measurements of the individual Cas proteins  
54  
55  
56  
57 194 were acquired under denaturing conditions by adding formic acid to a final  
58  
59  
60

1  
2  
3  
4 195 concentration of 5%. All spectra were mass calibrated by using an aqueous solution  
5  
6  
7 196 of cesium iodide (25mg/ml). Mass spectra were accumulated, averaged,  
8  
9  
10 197 smoothed and centered, using the software MassLynx 4.1 (Waters, United  
11  
12  
13  
14 198 Kingdom).

### 18 199 **Nuclease assays**

19 200 Oligo-nucleotide sequences used in this study are indicated in Table S2. Oligo-  
20  
21 201 nucleotides with C6-Amino modifications on the 5' terminus were obtained from ELLA  
22  
23 202 Biotech (Planegg, Germany). Cy5 or Cy3 (GE Healthcare) labelling of 5' termini was  
24  
25 203 done in 100 mM Sodium-bicarbonate buffer as described by (Joo and Ha, 2012).  
26  
27 204 Unlabeled oligo-nucleotides were obtained from Integrated DNA Technologies (IDT).  
28  
29 205 Nuclease assays were performed in buffer R (5 mM HEPES pH 7.5, 100 mM Sodium-  
30  
31 206 Glutamate supplemented with 2 mM  $MnCl_2$  and 10 mM  $MgCl_2$ . Annealed and Cy5 and  
32  
33 207 Cy3 labelled oligo-nucleotides (Cy3-BN1829+Cy5-BN1830) were added to a final  
34  
35 208 concentration of 125 nM and purified protein complexes to a final concentration of 500  
36  
37 209 nM. Reactions were incubated for 1 hour at 30°C after which reactions were quenched  
38  
39 210 by addition of Proteinase K (Thermo Fischer) and incubation for 1 hour at 37°C. The  
40  
41 211 resulting products were analyzed on denaturing PAGE (10% acrylamide, 8M Urea) and  
42  
43 212 analyzed with Amersham Typhoon fluorescence gel scanner (GE Healthcare).  
44  
45  
46  
47  
48

### 49 213 **In vitro spacer integration assays**

50 214 Oligo-nucleotide integrations with either Cy5 labeled or unlabeled oligo-nucleotides  
51  
52 215 were performed by pre-incubating indicated protein complexes (500 nM) with oligo-  
53  
54 216 nucleotides (250 nM) on ice for 15 min. Following pre-incubation, either linear CRISPR  
55  
56 217 substrate (obtained by Q5 high-fidelity PCR from pCRISPR using primers  
57  
58 218 BN015+BN1398) or supercoiled pCRISPR were added to a final concentration of 7.5  
59  
60

1  
2  
3 219 nM. Reaction mixtures were incubated at 30°C for 1 hour after which reactions were  
4  
5 220 quenched by addition of Proteinase K (Thermo Fischer) and incubation for 1 hour at  
6  
7 221 37°C. Reactions were run on 1% native agarose gels for 45 min and gels subsequently  
8  
9 222 stained with SYBR gold (Sigma Aldrich). Gels were scanned for Cy5 and SYBR gold  
10  
11 223 using Amersham Typhoon fluorescence gel scanner (GE Healthcare). For PCR  
12  
13 224 analysis of *in vitro* integration, unlabeled oligo-nucleotides were used in the reaction.  
14  
15 225 Open-circular plasmid DNA was gel isolated and DNA purified using Zymoclean gel  
16  
17 226 recovery kit (ZymoResearch) after which integration was assessed by PCR using  
18  
19 227 primers BN1711+BN1713 (leader distal integration; correct spacer orientation),  
20  
21 228 BN1711+BN1714 (leader distal integration; incorrect spacer orientation),  
22  
23 229 BN1712+BN1713 (leader proximal integration; incorrect spacer orientation) and  
24  
25 230 BN1712+BN1714 (leader proximal integration; correct spacer orientation). Purified  
26  
27 231 PCR amplicons were subjected to MiSeq sequencing (Illumina).  
28  
29  
30  
31  
32

33  
34  
35 232

### 233 **Next generation sequencing and statistical analysis**

234 After validation of PCR amplicons by gel electrophoresis and clean up with the  
235 GeneJET PCR Purification kit (Thermo Fisher Scientific) the samples were analyzed  
236 using Qubit fluorometric quantification (Invitrogen). Samples were prepared for  
237 sequencing with the Nextera XT DNA Library Preparation Kit (Illumina) and each library  
238 individually barcoded with the Nextera XT Index Kit v2 SetA (Illumina). Libraries were  
239 pooled equally and spiked with ~5% of the PhiX control library (Illumina) to artificially  
240 increase the genetic diversity before sequencing on a Nano flowcell (250 nt paired-  
241 end) with an Illumina MiSeq. Image analysis, base calling, de-multiplexing and data  
242 quality assessments were performed on the MiSeq instrument. FASTAQ files  
243 generated by the MiSeq were analyzed by pairing and merging the reads using  
244 Geneious 9.0.5 and subsequently extracting the oligo-nucleotide sequences used in

1  
2  
3 245 the in vitro integration assay. Overhang processing was analyzed by annotating the  
4  
5 246 primers used for amplification and comparing the overhangs post-integration to the  
6  
7  
8 247 initial oligo-nucleotide sequence.  
9

### 10 248 **In vivo spacer integration assays**

11 249 *E. coli* BL21 AI cells were co-transformed with either pCas1-Cas2 and pEmpty or  
12  
13 250 pCas1-2 and pCas4 (wild-type Cas4 or Cas4D76A+K91A). One transformant for each  
14  
15 251 combination was grown in LB at 37°C and continuous shaking (180 rpm) to OD<sub>600</sub>=0.3  
16  
17 252 and made electrocompetent after which pCRISPR was transformed. For each  
18  
19 253 treatment three individual colonies were grown in SOB medium (LB supplemented with  
20  
21 254 10 mM MgSO<sub>4</sub> and 10 mM MgCl<sub>2</sub>) at 37°C and continuous shaking (180 rpm) to  
22  
23 255 OD<sub>600</sub>=0.3 after which protein expression was induced by addition of 0.2% L-arabinose  
24  
25 256 and 0.5 mM IPTG. Induced cultures were grown for additional 2 hours at 37°C and  
26  
27 257 continuous shaking (180 rpm). Cells were made electrocompetent and annealed pre-  
28  
29 258 spacer oligo-nucleotides (BN1763+BN1768) electroporated at a final concentration of  
30  
31 259 1 μM. After 30 min recovery cells were harvested and plasmid DNA extracted using  
32  
33 260 GeneJET plasmid miniprep kit (Thermo Scientific). Extracted plasmid DNA was  
34  
35 261 normalized to 0.5 ng μl<sup>-1</sup> and subsequently 2 μl used in half-site integration PCRs using  
36  
37 262 primers BN1711+BN1713 (leader distal integration; correct spacer orientation),  
38  
39 263 BN1711+BN1714 (leader distal integration; incorrect spacer orientation),  
40  
41 264 BN1712+BN1713 (leader proximal integration; incorrect spacer orientation) and  
42  
43 265 BN1712+BN1714 (leader proximal integration; correct spacer orientation). PCR  
44  
45 266 amplicons were validated by agarose gel electrophoresis and purified with the  
46  
47 267 GeneJET PCR Purification kit (Thermo Scientific). Purified PCR amplicons were  
48  
49 268 subjected to MiSeq sequencing (Illumina).  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## 270 **Results**

### 271 **PAM-containing overhang processing depends on orientation of spacer** 272 **integration**

273  
274 We have previously demonstrated that the type I-D Cas4 protein facilitates the  
275 integration of PAM-compatible spacers *in vivo* (Kieper et al., 2018). In these  
276 experiments we looked at the total pool of spacers that were acquired from cytosolic  
277 DNA, obscuring the detailed mechanism that governs the processing of prespacer  
278 substrates. In order to obtain more detailed insights into processing of PAM and non-  
279 PAM containing substrates and how spacer orientation affects overhang processing,  
280 we electroporated an idealized prespacer substrate into *E. coli* cells overexpressing  
281 either Cas1-Cas2 or Cas4-Cas1-Cas2 proteins (Fig. 1A&B). Cas4 was either  
282 expressed as the wild-type protein or the catalytically inactive mutant (D76A, K91A).  
283 In addition to the adaptation genes, the cells were carrying a plasmid containing the  
284 type I-D leader and a single repeat (pCRISPR). By employing half-site integration  
285 PCRs followed by high-throughput sequencing, we analyzed the 3' overhangs after  
286 processing and integration *in vivo* (Fig. 1C). This approach allowed us to differentiate  
287 between correct and incorrect spacer orientation, as well as correct and incorrect PAM  
288 processing of their 3' end (Fig. 1D&E).

289 We observed that prespacer overhangs were trimmed by at least 5 nt in all cases,  
290 regardless of the presence or absence of Cas4. However, processing of PAM and non-  
291 PAM containing overhangs differed depending on the orientation in which the spacer  
292 was integrated. In particular, spacers integrated in the correct orientation (Fig. 1D)  
293 were more precisely processed in the PAM-containing overhang in the presence of  
294 Cas4. Although cells expressing only Cas1-Cas2 or a combination of Cas1-Cas2 with  
295 a catalytically inactive Cas4 double mutant (D76A, K91A) also displayed 30% to 35%  
296 of correct processing of the PAM overhang, their spacer size distributions were

1  
2  
3 297 typically broader and shifted towards longer overhang lengths. Analyzing the non-PAM  
4  
5 298 containing overhangs of correctly oriented spacers did not display any differences  
6  
7 299 between the conditions (with and without Cas4), suggesting that prespacer overhangs  
8  
9  
10 300 without PAM are not processed by Cas4, but rather by endogenous *E. coli* nucleases.  
11  
12 301 Spacers integrated in the incorrect orientation (Fig. 1E) showed similar 3' overhangs  
13  
14 302 under all conditions. We observed most accurate processing when Cas4 was present  
15  
16  
17 303 in its active form. In those samples the presence of Cas4 led to an increase in the  
18  
19 304 shortening of overhangs, with a predominant overhang length of 6 nucleotides. In the  
20  
21 305 *E. coli* model system, host factor nucleases potentially act on both PAM and non-PAM  
22  
23 306 containing 3' overhangs that remain unprotected by the core Cas1-Cas2 complex  
24  
25  
26 307 holding the prespacer. Cas4 does not specifically cleave non-PAM containing  
27  
28 308 overhangs, but requires the presence of the PAM in order to engage in sequence-  
29  
30  
31 309 specific processing.

32  
33 310

### 311 **Cas4 forms a strong heteromeric complex with Cas1**

312 In order to assess whether the overhang processing connected to the presence of  
33  
34  
35  
36  
37 313 Cas4 was a result of Cas4 specifically interacting with the Cas1-Cas2 integration  
38  
39  
40  
41 314 complex, we first investigated the formation of a Cas1-Cas2 complex. The Cas1 protein  
42  
43  
44 315 was N-terminally His<sub>6</sub>-tagged and co-expressed with Cas2 in *E. coli* BL21-AI cells.  
45  
46 316 After the initial nickel-affinity pull-down from cleared cell lysate, the elution fraction was  
47  
48 317 subjected to size exclusion chromatography (SEC), which resulted in one peak  
49  
50  
51 318 containing aggregated protein and another peak species (Fig. 2A). This peak  
52  
53 319 contained three proteins (Fig. 2A) for which the tagged-Cas1, untagged Cas1 (due to  
54  
55 320 autoproteolysis of the tag) and Cas2 protein identity was confirmed by mass  
56  
57 321 spectrometry. Next, we co-expressed the His<sub>6</sub>-tagged Cas1 with untagged Cas4 and  
58  
59 322 observed strong co-purification of both proteins (Fig. 2B). In order to verify the Cas4-

1  
2  
3 323 Cas1 interaction, a reverse tagging strategy was used (His6-tagged Cas4 co-  
4  
5 324 expressed with untagged Cas1), which again confirmed the presence of a Cas4-Cas1  
6  
7 325 complex (Fig. S1). When tagged-Cas1 and Cas2 were co-expressed along with Cas4,  
8  
9 326 we observed a strong co-purification of Cas4 and Cas1 that abolished formation of the  
10  
11 327 Cas1-Cas2 complex since Cas2 eluted separately as a low molecular weight species.  
12  
13 328 This fraction also contained minor amounts of Cas1 and Cas4 that did not assemble  
14  
15 329 into higher order complexes. Our results demonstrate that under these conditions Cas1  
16  
17 330 can form complexes with Cas4 or Cas2, and that these complexes appear to be  
18  
19 331 mutually exclusive. In the presence of both Cas4 and Cas2, Cas1 strongly favors the  
20  
21 332 interaction with Cas4 over the interaction with Cas2.  
22  
23  
24  
25

26 333

#### 27 28 334 **Cas4 associates with Cas1 in a 1:2 ratio**

29 335 Next, we determined the stoichiometry of the formed Cas4-Cas1 complex. Previously,  
30  
31 336 Lee et al. demonstrated that the heteromeric complex consists of two Cas1 dimers that  
32  
33 337 each associate with a single Cas4 monomer (Lee et al., 2018). To gain insight into the  
34  
35 338 composition of the untagged Cas4-Cas1 complex (Fig. S2) native protein mass  
36  
37 339 spectrometry analysis was performed (van den Heuvel et al., 2006). The mass  
38  
39 340 spectrum (Fig. 2D) revealed a distribution of different complex species, with the most  
40  
41 341 abundant mass-over-charge ( $m/z$ ) peaks consisting of either Cas1 dimers ( $73.6 \pm 1.6$   
42  
43 342 kDa) or the Cas4-Cas1 complex consisting of a single Cas1 dimer and a Cas4  
44  
45 343 monomer resulting in a Cas4<sub>1</sub>-Cas1<sub>2</sub> complex of 96.3 kDa (Fig. 2D). Even though we  
46  
47 344 observed co-purification of Cas1 and Cas2 in the SEC analysis, the native mass  
48  
49 345 spectrum of the Cas1-Cas2 complex resulted in mainly Cas1 dimers (Fig. S3) with  
50  
51 346 Cas2 likely being lost during the native MS sample preparation.  
52  
53  
54  
55  
56

57 347



1  
2  
3 348 **The Cas4-Cas1 complex sequence specifically processes PAM-**  
4 349 **containing 3' overhangs**

5 350 The acquisition of functional spacers not only requires appropriate prespacer selection,  
6  
7  
8 351 but also PAM-compliant processing. We have previously shown that the presence of  
9  
10 352 Cas4 in addition to the core adaptation proteins Cas1 and Cas2 significantly increases  
11  
12 353 the integration of spacers with a correctly processed PAM *in vivo* (Kieper et al., 2018).  
13  
14 354 Due to the strong interaction of Cas4 and Cas1 we aimed to test whether PAM  
15  
16 355 processing is mediated only by Cas4, or if the heteromeric Cas4-Cas1 complex is  
17  
18 356 required. In order to address this question, we performed prespacer cleavage assays  
19  
20 357 with a dual-labelled model prespacer (Fig. 3A). This model prespacer consisted of a  
21  
22 358 25 bp duplex flanked by 13 nucleotide 3' overhangs on each side and fluorescent labels  
23  
24 359 at their 5' ends. The top strand was labelled with Cy3 and did not contain a PAM  
25  
26 360 sequence in its 3' overhang of the, while the bottom strand was labelled with Cy5 and  
27  
28 361 contained the I-D consensus PAM. We found that neither free Cas4 nor Cas1-Cas2  
29  
30 362 was able to catalyze 3' overhang cleavage. However, the addition of the Cas4-Cas1  
31  
32 363 complex resulted in a defined band corresponding to processing of the PAM sequence  
33  
34 364 within the PAM-containing overhang (Fig. 3A). This result suggests that PAM  
35  
36 365 recognition is mediated by the interactions within the Cas4-Cas1 complex, where Cas4  
37  
38 366 acts as the catalytic subunit of the complex.

39  
40 367 We have previously shown that mutating D76 in the conserved RecB domain of Cas4  
41  
42 368 abolished integration of PAM-proficient spacers *in vivo* (Kieper et al. (2018). Using the  
43  
44 369 D76A mutant in our *in vitro* cleavage assay fully abolished processing activity of the  
45  
46 370 Cas4-Cas1 complex, demonstrating that the RecB domain of Cas4 is indeed the  
47  
48 371 catalytically active site required for PAM processing. Interestingly, although Cas4 did  
49  
50 372 not show processing activity on its own, combining Cas1-Cas2 and Cas4 fully restored  
51  
52 373 processing to similar levels as the Cas4-Cas1 complex. The addition of both, the Cas1-  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 374 Cas2 and the Cas4-Cas1 complex, resulted in processing of the PAM overhang as  
4  
5 375 observed with the Cas4-Cas1 complex alone or combination of Cas4 and Cas1-2  
6  
7  
8 376 complex. All conditions that showed cleavage of the substrate resulted in a single  
9  
10 377 defined band, suggesting that cleavage occurred via an endonuclease mechanism  
11  
12 378 which is in line with previous studies (Lee et al., 2018). The processing of the non-PAM  
13  
14 379 containing overhang was not observed in any of the conditions, indicating that the  
15  
16  
17 380 processing of the non-PAM site presumably relies on host factor nucleases such as  
18  
19 381 DnaQ-like exonucleases or Exonuclease T as recently found in in the I-E system (Kim  
20  
21 382 et al., 2020; Ramachandran et al., 2020). Our results demonstrate that sequence  
22  
23 383 specific Cas4 activity requires the presence of Cas1 and that the Cas4-Cas1 complex  
24  
25  
26 384 is the core processing complex that sequence-specifically recognizes and processes  
27  
28 385 the PAM sequence before integration.

30  
31 386 **The Cas4-Cas1 complex integrates new spacers into both linear and**  
32 387 **supercoiled DNA**

33 388 Next, we tested whether the Cas4-Cas1 complex not only processes prespacer  
34  
35 389 substrates but also catalyzes their integration into the CRISPR array. We performed  
36  
37  
38 390 adaptation assays using supercoiled plasmid DNA containing the type I-D leader and  
39  
40 391 a single repeat (pCRISPR; Fig. 1A) as well as linear CRISPR array substrates  
41  
42 392 generated by PCR (Fig. 3B). Both linear and plasmid CRISPR loci were incubated with  
43  
44 393 a Cy5-labelled prespacer, the Cas4-Cas1 and Cas1-2 complexes. Intriguingly, we  
45  
46 394 observed coupling of the labelled prespacer by the Cas4-Cas1 complex to both  
47  
48 395 CRISPR substrates, showing that this sub complex is proficient in catalyzing at least  
49  
50 396 half-site spacer integration (Fig. 3B&C). Spacer integration into plasmid DNA resulted  
51  
52 397 in the formation of open-circular (OC) plasmid conformations. Merging the Cy5 signal  
53  
54 398 of the prespacer and the plasmid DNA signal confirmed that the prespacer was indeed  
55  
56 399 coupled to the OC form of the plasmid. The Cas1-2 complex was able to integrate the  
57  
58  
59  
60

1  
2  
3 400 prespacer into both linear and supercoiled arrays similar to the Cas4-Cas1 complex.  
4  
5 401 Our observation demonstrates that at least two different sub-complexes exist, which  
6  
7 402 are both capable of catalyzing half-site spacer integration. Taken together, based on  
8  
9 403 the selective PAM-overhang processing of the Cas4-Cas1 complex, we hypothesize  
10  
11 404 that the Cas4-Cas1 complex processes and integrates the PAM containing overhang  
12  
13 405 and the Cas1-Cas2 complex the non-PAM containing overhang.  
14  
15  
16

### 17 406 **Correct spacer orientation requires overhang processing prior to** 18 407 **integration**

19 408 In order to analyze the accuracy of spacer integration by the Cas4-Cas1 complex in  
20  
21 409 more detail, OC plasmid resulting from the integration reaction was gel purified and  
22  
23 410 subjected to half-site integration PCRs as described previously (Fig. 1C). PCR  
24  
25 411 products were subjected to Illumina MiSeq sequencing and prespacer sequences were  
26  
27 412 extracted. This approach allowed us to assess 3' overhang processing before  
28  
29 413 integration at the leader-proximal or leader-distal integration site. Interestingly, PAM-  
30  
31 414 containing overhangs only showed sequence specific processing when the spacer was  
32  
33 415 correctly oriented with respect to the PAM (Fig. 3D). The Cas4-Cas1 complex cleaved  
34  
35 416 65% of correctly oriented spacers exactly downstream of the PAM, however, we also  
36  
37 417 observed incorrect removal of a single nucleotide in 25% of sequences and removal of  
38  
39 418 2 or more nucleotides in 10% of the sequences. Surprisingly, incorrectly oriented  
40  
41 419 spacers did not show any processing of the PAM-containing overhang (Fig. 3E),  
42  
43 420 indicating that integration in the correct orientation is preceded by the processing of  
44  
45 421 the overhang. As predicted from the bulk cleavage assays, we did not observe any  
46  
47 422 processing of the non-PAM containing overhangs regardless of the spacer orientation.  
48  
49 423 Our data show that integration of new spacers in the correct orientation by Cas4-Cas1  
50  
51 424 requires PAM recognition and processing before a spacer can be integrated.  
52  
53  
54  
55  
56  
57  
58  
59  
60

### 425 **Spacer integration preferentially initiates with the non-PAM overhang**

426 In type I CRISPR-Cas systems spacer integration initiates by first integrating the non-  
427 PAM end of the spacer at the leader repeat junction and proceeding with the coupling  
428 of the PAM-end of the spacer at the Repeat-Spacer boundary (Arslan et al., 2014;  
429 Nuñez et al., 2015b; Rollie et al., 2015). In the Type I-E CRISPR-Cas system that is  
430 lacking Cas4, directionality of spacer integration is dictated by the prespacer  
431 processing kinetics (Kim et al., 2020; Ramachandran et al., 2020). We therefore  
432 hypothesized that the prespacer processing of our Cas4 containing system could  
433 influence the orientation of the integrated spacer. To test the effect of the processed  
434 and unprocessed prespacer overhangs on integration, we assayed spacer integration  
435 using the Cy5-labelled prespacer substrates. Prespacers were either fully processed  
436 (5 nt 3' overhangs), with an unprocessed non-PAM overhang (13 nt) or with a  
437 processed, but integration-deficient (Fagerlund et al., 2017; Rollie et al., 2015) 3'  
438 phosphorylated non-PAM overhang (Fig. 4A). The fully processed substrate was  
439 efficiently coupled by Cas1-Cas2, Cas4-Cas1 as well as the combination of both.  
440 Similarly, the prespacer with an unprocessed non-PAM overhang was coupled  
441 efficiently in all three treatments. However, when the processed non-PAM overhang  
442 was blocked for integration by 3'-phosphorylation, neither of the protein complexes was  
443 able to efficiently couple the spacer, indicating that coupling of the PAM-overhang  
444 requires prior integration of the non-PAM overhang. Altogether, this mechanism  
445 ensures that integration of the PAM site of the spacer is halted until integration of the  
446 non-PAM site has occurred, resulting in the correct orientation of the spacer with  
447 respect to the PAM.

### 448 **Discussion**

449 Although Cas4 proteins have been recognized as part of the core cas gene machinery  
450 almost two decades ago (Jansen et al., 2002), its role in acquiring PAM-compatible

1  
2  
3 451 spacers has been revealed only in the recent years (Kieper et al., 2018; Lee et al.,  
4  
5 452 2019; Lee et al., 2018; Shiimori et al., 2018; Zhang et al., 2019). Here we provide a  
6  
7 453 new mechanistic understanding of how Cas4-dependent PAM selection is achieved  
8  
9 454 during CRISPR adaptation, and specifically how asymmetry of the adaptation complex  
10  
11 455 drives the selection, processing and integration of PAM-compatible spacers. We  
12  
13 456 present a model in which two independent subcomplexes, Cas4-Cas1 and Cas1-Cas2,  
14  
15 457 selectively process the two 3' overhangs of a prespacer (Fig. 5). The interaction of  
16  
17 458 Cas4 with the Cas1 integrase protein is central to the recognition and processing of  
18  
19 459 PAM-containing prespacer substrates. Formation of this Cas4-Cas1 complex is  
20  
21 460 mutually exclusive with formation of the Cas1-Cas2 complex, which may suggest  
22  
23 461 distinct roles of both subcomplexes. We found that the Cas4-Cas1 subcomplex  
24  
25 462 displays prespacer cleavage activity only on PAM-containing 3' overhangs. Cas4-Cas1  
26  
27 463 removes the PAM via endonuclease cleavage while Cas1-2 defines overhang trimming  
28  
29 464 likely through host factor nucleases. Subsequently, Cas1-Cas2 initiates coupling of the  
30  
31 465 non-PAM overhang to the leader-repeat junction followed by integration of the  
32  
33 466 processed PAM-site at the repeat-spacer junction.

34  
35 467 Our findings are consistent with previous studies that established the existence of two  
36  
37 468 mutually exclusive Cas4-Cas1 and Cas1-2 complexes in type I-C CRISPR-Cas  
38  
39 469 systems (Lee et al., 2019; Lee et al., 2018), and expand our understanding of the roles  
40  
41 470 of these subcomplexes. Moreover, the RecB-domain mediated activity of Cas4 is  
42  
43 471 dependent on the presence of Cas1, since Cas4 alone is not able to recognize and  
44  
45 472 process the PAM sequence. This observation suggests that the Cas4-Cas1 interaction  
46  
47 473 is essential for sequence specific recognition of the PAM. It remains to be determined  
48  
49 474 whether the PAM sequence recognition domain is located within Cas4 or Cas1.  
50  
51 475 Interestingly, PAM selection in the type I-E system is mediated by the C-terminal tail of  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 476 Cas1 (Kim et al., 2020), however, this C-terminal proportion is not conserved in the  
4  
5 477 type I-D Cas1 protein. Future structural and biochemical studies will have to address  
6  
7  
8 478 how PAM selection is achieved.

9  
10 479 Cas4-Cas1 did not display activity on non-PAM containing overhangs *in vitro*, however,  
11  
12 480 processing of the non-PAM overhang was observed in our *in vivo* setup, suggesting  
13  
14 481 that processing involves other non-Cas proteins. This finding is in line with the Cas4-  
15  
16 482 deficient type I-E system in which host factors such as the ExoT and DnaQ-like  
17  
18 483 exonucleases are required for processing both, PAM-containing and non-PAM  
19  
20 484 overhangs. (Kim et al., 2020; Ramachandran et al., 2020). We propose that, in analogy  
21  
22 485 to the *E. coli* type I-E system, 3'-5' exonucleases act as trimming factors for non-PAM  
23  
24 486 3' overhangs in the native *Synechocystis* PCC6803 host. The Cas1 protein of the type  
25  
26 487 I-E system recognizes and protects the PAM from premature trimming, causing a  
27  
28 488 delayed processing of the PAM end that ensures correct orientation (Kim et al., 2020).  
29  
30 489 Upon activation, the Cas4-Cas1 complex sequence specifically removes the type I-D  
31  
32 490 PAM, although incorrect processing was observed *in vivo* and *in vitro* that would result  
33  
34 491 in single-nucleotide slipped spacers. Recently, it was observed in the type I-F system  
35  
36 492 that slipped spacers increase primed adaptation which enhances the spacer diversity  
37  
38 493 of the population (Jackson et al., 2019). Our results suggest the possibility that such  
39  
40 494 erroneous PAM processing could promote the integration of slipped spacers and by  
41  
42 495 extension, primed adaptation as found in other type I CRISPR-Cas systems  
43  
44 496 (Nussenzweig and Marraffini, 2020).

45  
46 497 Cryo-EM structures of the type I-C Cas4-Cas1-Cas2 complex revealed that the  
47  
48 498 complex might undergo a conformational change (e.g. causing dissociation of Cas4  
49  
50 499 from the complex) in order to allow for Cas1 mediated integration of the PAM-end site  
51  
52 500 of the spacer. Lee et al. showed that 50% of their Cas4-Cas1-Cas2 complex structures  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 501 lacked the Cas4 density on one site of the complex, resulting in an asymmetric  
4  
5 502 complex. Our observation of two integrase complexes (Cas4-Cas1 and Cas1-Cas2)  
6  
7 503 that are independently capable of at least half-site spacer integration points towards a  
8  
9 504 similar asymmetrical organization of the full adaptation complex, in which Cas4-Cas1  
10  
11 505 is involved in PAM-site and Cas1-Cas2 in non-PAM site integration. By testing  
12  
13 506 asymmetric spacer precursors, we demonstrate that integrase activity of the type I-D  
14  
15 507 Cas4-Cas1 complex is potentially halted until integration of the non-PAM overhang has  
16  
17 508 occurred. We propose a model for the type I-D system that relies on a delayed PAM-  
18  
19 509 site integration by Cas4 in order to result in a correctly oriented spacer. In summary,  
20  
21 510 we propose a mechanism in which two functionally independent complexes, Cas4-  
22  
23 511 Cas1 and Cas1-Cas2, sequentially process and integrate prespacer substrates. This  
24  
25 512 mechanism ensures correct spacer orientation as well as correct PAM-processing,  
26  
27 513 thereby resulting in interference-proficient CRISPR adaptation.  
28  
29  
30  
31  
32

33 514

35 515 **Acknowledgements**

37 516 We would like to thank Marre Niessen for early contributions. We thank Dr. Viktorija  
38  
39 517 Globyte for critical reading of the manuscript. We thank Rob B.M. Koehorst for  
40  
41 518 providing the *Synechocystis* PCC6803 strain.  
42  
43  
44

45 519

46 520 Author contributions: S.N.K, C.A. and S.J.J.B. designed the experiments. S.N.K,  
47  
48 521 A.C.H. and A.B. performed the experiments. S.N.K., C.A., A.B., A.J.H.R and S.J.J.B  
49  
50 522 analyzed the data. S.N.K, C.A. and S.J.J.B. wrote the paper with input from all authors.  
51  
52

53 523

54 524

55 525

56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

526

527 FUNDING: FOM (Projectruimte 15PR3188-2); Netherlands Organisation for Scientific

528 Research [VICI VI.C.182.027]

529

530 Conflict of interest statement. None declared.

531

532

533



1  
2  
3 535 **References**  
4

- 5 536 Arslan, Z., Hermanns, V., Wurm, R., Wagner, R., and Pul, Ü. (2014). Detection and  
6 537 characterization of spacer integration intermediates in type I-E CRISPR–Cas system. *Nucleic*  
7 538 *Acids Research* *42*, 7884-7893.
- 9 539 Barrangou, R. (2013). CRISPR-Cas systems and RNA-guided interference. *WIREs RNA* *4*, 267-  
10 540 278.  
11 541
- 13 542 Brouns, S.J.J., Jore, M.M., Lundgren, M., Westra, E.R., Slijkhuys, R.J.H., Snijders, A.P.L.,  
14 543 Dickman, M.J., Makarova, K.S., Koonin, E.V., and van der Oost, J. (2008). Small CRISPR RNAs  
15 544 Guide Antiviral Defense in Prokaryotes. *Science* *321*, 960-964.  
16 545
- 18 546 Fagerlund, R.D., Wilkinson, M.E., Klykov, O., Barendregt, A., Pearce, F.G., Kieper, S.N.,  
19 547 Maxwell, H.W.R., Capolupo, A., Heck, A.J.R., Krause, K.L., *et al.* (2017). Spacer capture and  
20 548 integration by a type I-F Cas1-Cas2-3 CRISPR adaptation complex. *Proc Natl Acad Sci U S A* *114*,  
21 549 E5122-E5128.  
22 550
- 24 551 Gonzales, M.F., Brooks, T., Pukatzki, S.U., and Provenzano, D. (2013). Rapid Protocol for  
25 552 Preparation of Electrocompetent *Escherichia coli* and *Vibrio cholerae*. *Journal of Visualized*  
26 553 *Experiments : JoVE*, 50684.  
27 554
- 29 555 Hou, Z., and Zhang, Y. (2018). Insights into a Mysterious CRISPR Adaptation Factor, Cas4. *Mol*  
30 556 *Cell* *70*, 757-758.  
31 557
- 32 558 Hudaiberdiev, S., Shmakov, S., Wolf, Y.I., Terns, M.P., Makarova, K.S., and Koonin, E.V. (2017).  
33 559 Phylogenomics of Cas4 family nucleases. *BMC Evolutionary Biology* *17*, 232-232.  
34 560
- 36 561 Jackson, S.A., Birkholz, N., Malone, L.M., and Fineran, P.C. (2019). Imprecise Spacer Acquisition  
37 562 Generates CRISPR-Cas Immune Diversity through Primed Adaptation. *Cell Host and Microbe*  
38 563 *25*, 250-260.e254.  
39 564
- 41 565 Jackson, S.A., McKenzie, R.E., Fagerlund, R.D., Kieper, S.N., Fineran, P.C., and Brouns, S.J.  
42 566 (2017). CRISPR-Cas: Adapting to change. *Science* *356*.  
43 567
- 45 568 Jansen, R., Embden, J.D.A.V., Gaastra, W., and Schouls, L.M. (2002). Identification of genes  
46 569 that are associated with DNA repeats in prokaryotes. *Molecular Microbiology* *43*, 1565-1575.  
47 570
- 48 571 Joo, C., and Ha, T. (2012). Labeling DNA (or RNA) for single-molecule FRET. *Cold Spring Harb*  
49 572 *Protoc* *2012*, 1005-1008.  
50 573
- 52 574 Kieper, S.N., Almendros, C., Behler, J., McKenzie, R.E., Nobrega, F.L., Haagsma, A.C., Vink,  
53 575 J.N.A., Hess, W.R., and Brouns, S.J.J. (2018). Cas4 Facilitates PAM-Compatible Spacer Selection  
54 576 during CRISPR Adaptation. *Cell Reports* *22*, 3377-3384.  
55 577
- 57 578 Kim, S., Loeff, L., Colombo, S., Jergic, S., Brouns, S.J.J., and Joo, C. (2020). Selective loading and  
58 579 processing of pre-spacers for precise CRISPR adaptation. *Nature*.  
59 580  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

581 Lee, H., Dhingra, Y., and Sashital, D.G. (2019). The Cas4-Cas1-Cas2 complex mediates precise  
582 prespacer processing during CRISPR adaptation. *eLife* 8, 1-84.

583  
584 Lee, H., Zhou, Y., Taylor, D.W., and Sashital, D.G. (2018). Cas4-Dependent Prespacer Processing  
585 Ensures High-Fidelity Programming of CRISPR Arrays. *Molecular Cell* 70, 48-59.e45.

586  
587 Lemak, S., Beloglazova, N., Nocek, B., Skarina, T., Flick, R., Brown, G., Popovic, A., Joachimiak,  
588 A., Savchenko, A., and Yakunin, A.F. (2013). Toroidal Structure and DNA Cleavage by the  
589 CRISPR-Associated [4Fe-4S] Cluster Containing Cas4 Nuclease SSO0001 from *Sulfolobus*  
590 *solfataricus*. *Journal of the American Chemical Society* 135, 17476-17487.

591  
592 Lemak, S., Nocek, B., Beloglazova, N., Skarina, T., Flick, R., Brown, G., Joachimiak, A.,  
593 Savchenko, A., and Yakunin, A.F. (2014). The CRISPR-associated Cas4 protein Pcal\_0546 from  
594 *Pyrobaculum calidifontis* contains a [2Fe-2S] cluster: crystal structure and nuclease activity.  
595 *Nucleic Acids Research* 42, 11144-11155.

596  
597 Li, M., Wang, R., Zhao, D., and Xiang, H. (2014). Adaptation of the *Haloarcula hispanica* CRISPR-  
598 Cas system to a purified virus strictly requires a priming process. *Nucleic Acids Research* 42,  
599 2483-2492.

600  
601 Liu, T., Liu, Z., Ye, Q., Pan, S., Wang, X., Li, Y., Peng, W., Liang, Y., She, Q., and Peng, N. (2017).  
602 Coupling transcriptional activation of CRISPR-Cas system and DNA repair genes by Csa3a in  
603 *Sulfolobus islandicus*. *Nucleic Acids Research* 45, 8978-8992.

604  
605 McGinn, J., and Marraffini, L.A. (2019). Molecular mechanisms of CRISPR-Cas spacer  
606 acquisition. *Nature Reviews Microbiology* 17, 7-12.

607  
608 Nuñez, J.K., Harrington, L.B., Kranzusch, P.J., Engelman, A.N., and Doudna, J.A. (2015a).  
609 Foreign DNA capture during CRISPR-Cas adaptive immunity. *Nature* 527, 535-538.

610  
611 Nuñez, J.K., Kranzusch, P.J., Noeske, J., Wright, A.V., Davies, C.W., and Doudna, J.A. (2014).  
612 Cas1-Cas2 complex formation mediates spacer acquisition during CRISPR-Cas adaptive  
613 immunity. *Nature Structural & Molecular Biology* 21, 528-534.

614  
615 Nuñez, J.K., Lee, A.S.Y., Engelman, A., and Doudna, J.A. (2015b). Integrase-mediated spacer  
616 acquisition during CRISPR-Cas adaptive immunity. *Nature* 519, 193-198.

617  
618 Nussenzweig, P.M., and Marraffini, L.A. (2020). Molecular Mechanisms of CRISPR-Cas  
619 Immunity in Bacteria. *Annu Rev Genet* 54, 93-120.

620  
621 Plagens, A., Tjaden, B., Hagemann, A., Randau, L., and Hensel, R. (2012). Characterization of  
622 the CRISPR/Cas Subtype I-A System of the Hyperthermophilic Crenarchaeon *Thermoproteus*  
623 *tenax*. *Journal of Bacteriology* 194, 2491-2500.

624  
625 Ramachandran, A., Summerville, L., Learn, B.A., DeBell, L., and Bailey, S. (2020). Processing  
626 and integration of functionally oriented prespacers in the *Escherichia coli* CRISPR system  
627 depends on bacterial host exonucleases. *J Biol Chem* 295, 3403-3414.

- 1  
2  
3 628  
4 629 Rollie, C., Schneider, S., Brinkmann, A.S., Bolt, E.L., and White, M.F. (2015). Intrinsic sequence  
5 630 specificity of the Cas1 integrase directs new spacer acquisition. *eLife* 4, e08716.  
6 631  
7  
8 632 Shiimori, M., Garrett, S.C., Graveley, B.R., and Terns, M.P. (2018). Cas4 Nucleases Define the  
9 633 PAM, Length, and Orientation of DNA Fragments Integrated at CRISPR Loci. *Molecular cell* 70,  
10 634 814-824.e816.  
11 635  
12  
13 636 Tahallah, N., Pinkse, M., Maier, C.S., and Heck, A.J. (2001). The effect of the source pressure  
14 637 on the abundance of ions of noncovalent protein assemblies in an electrospray ionization  
15 638 orthogonal time-of-flight instrument. *Rapid Commun Mass Spectrom* 15, 596-601.  
16 639  
17  
18 640 van den Heuvel, R.H., van Duijn, E., Mazon, H., Synowsky, S.A., Lorenzen, K., Versluis, C.,  
19 641 Brouns, S.J., Langridge, D., van der Oost, J., Hoyes, J., *et al.* (2006). Improving the performance  
20 642 of a quadrupole time-of-flight instrument for macromolecular mass spectrometry. *Anal Chem*  
21 643 78, 7473-7483.  
22 644  
23  
24 645 van der Oost, J., Westra, E.R., Jackson, R.N., and Wiedenheft, B. (2014). Unravelling the  
25 646 structural and mechanistic basis of CRISPR–Cas systems. *Nature Reviews Microbiology* 12,  
26 647 479-492.  
27 648  
28  
29 649 Yosef, I., Goren, M.G., and Qimron, U. (2012). Proteins and DNA elements essential for the  
30 650 CRISPR adaptation process in *Escherichia coli*. *Nucleic Acids Research* 40, 5569-5576.  
31 651  
32 652 Zhang, J., Kasciukovic, T., and White, M.F. (2012). The CRISPR Associated Protein Cas4 Is a 5'  
33 653 to 3' DNA Exonuclease with an Iron-Sulfur Cluster. *PLoS ONE* 7, e47232-e47232.  
34 654  
35  
36 655 Zhang, Z., Pan, S., Liu, T., Li, Y., and Peng, N. (2019). Cas4 Nucleases Can Effect Specific  
37 656 Integration of CRISPR Spacers. *Journal of Bacteriology* 201.  
38 657  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

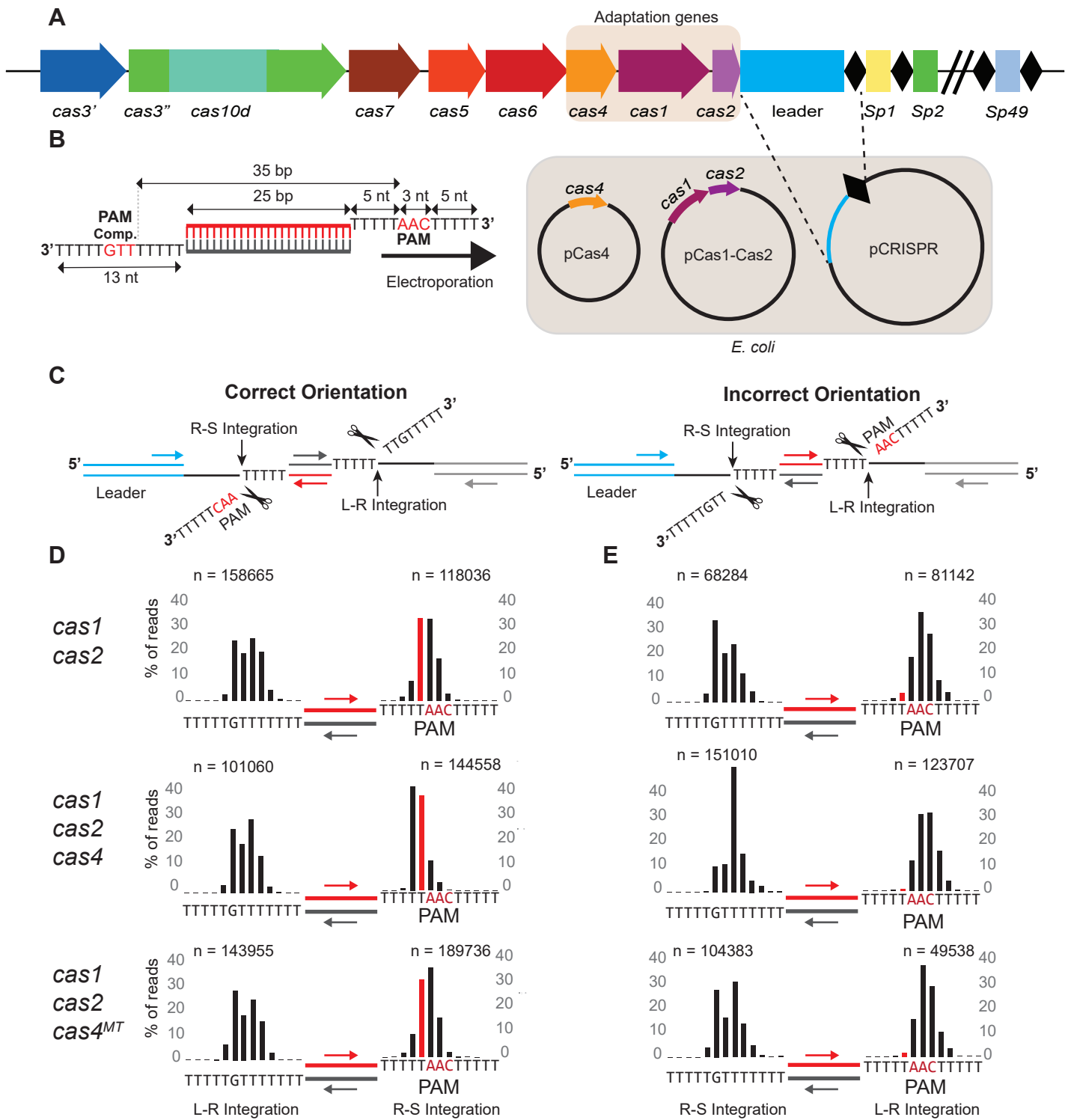


Figure 1

50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

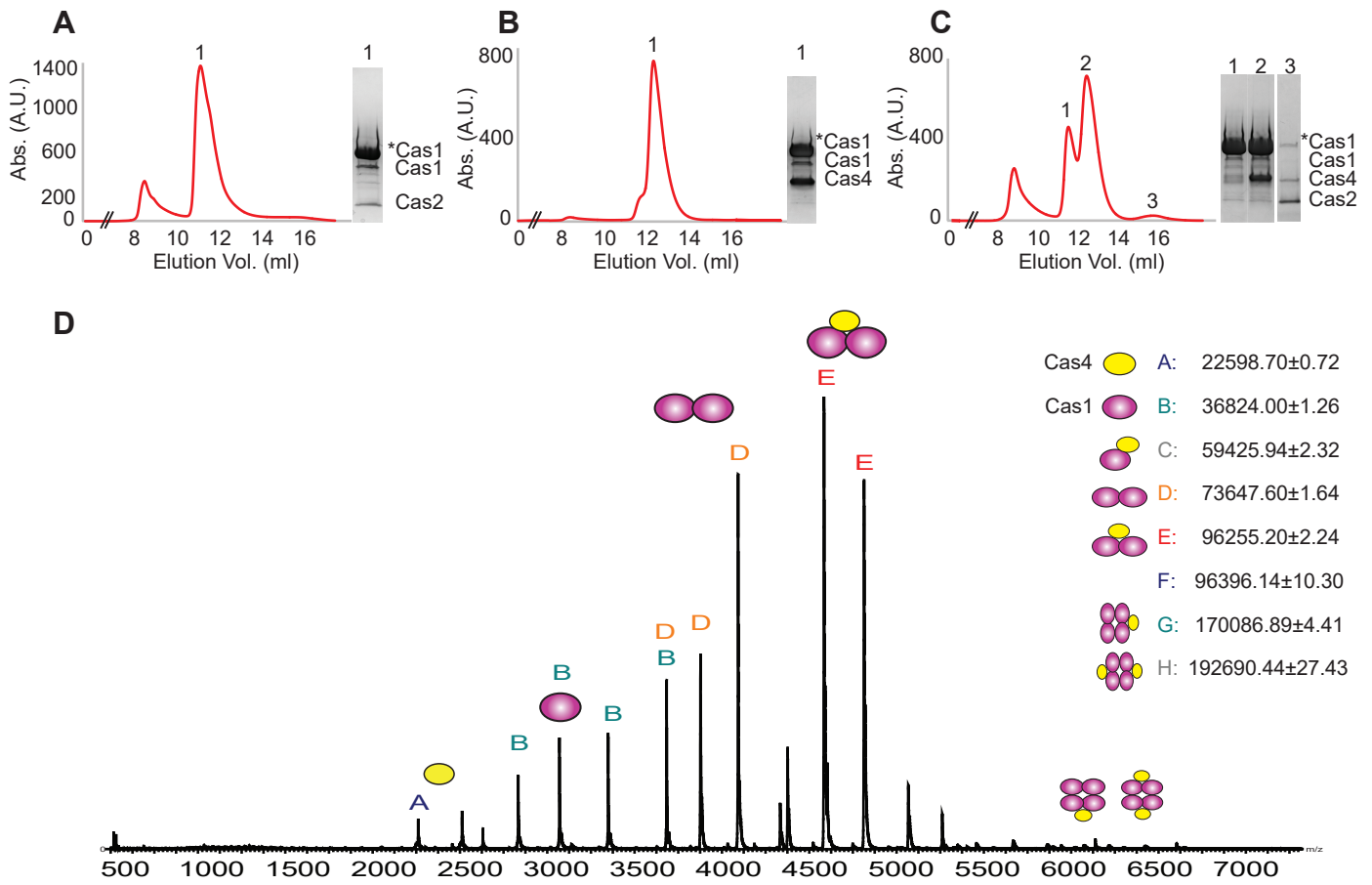
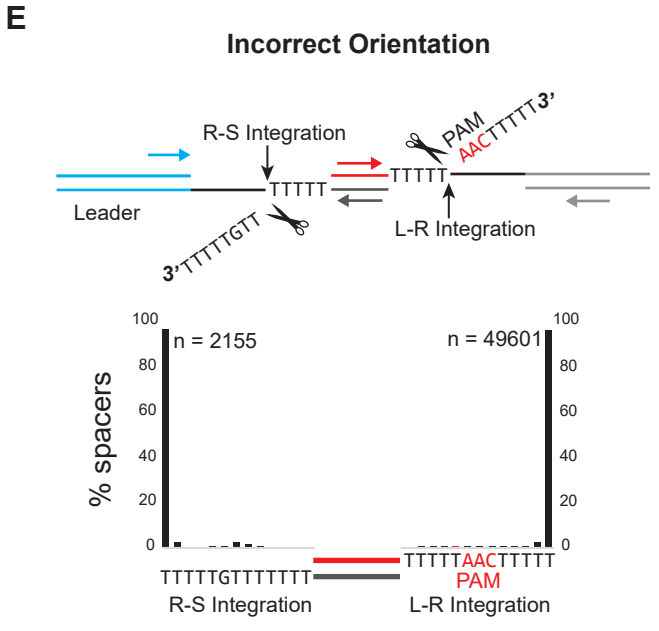
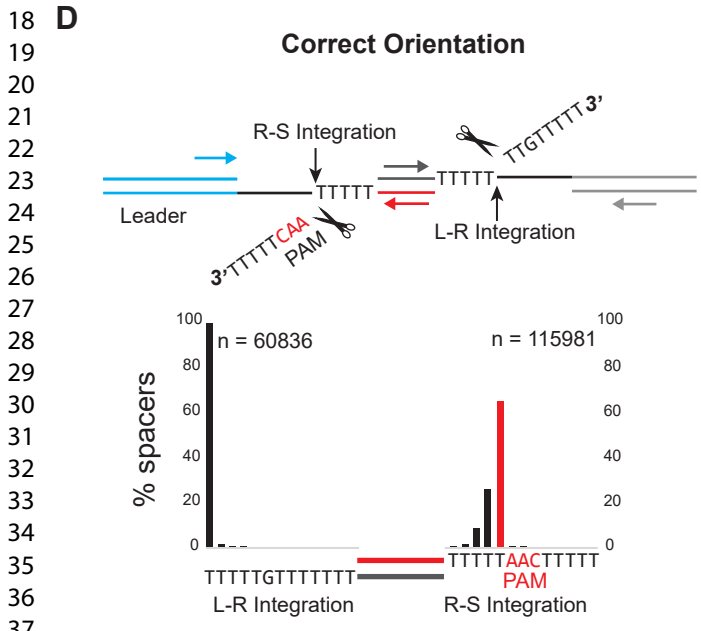
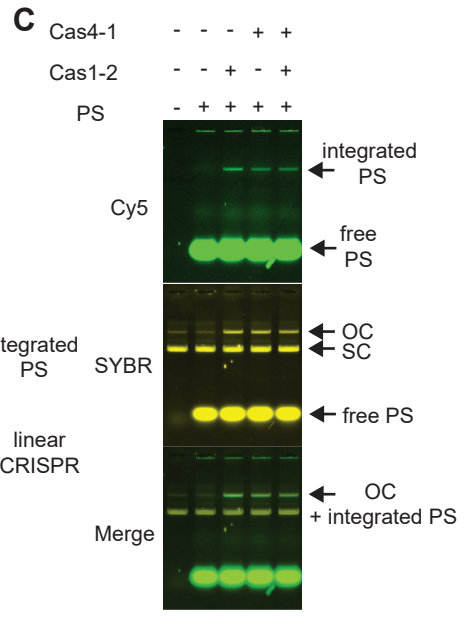
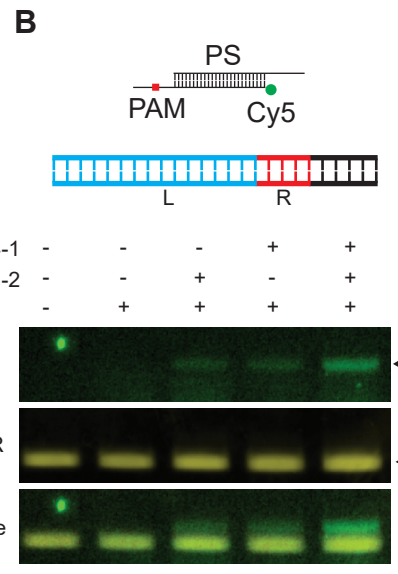
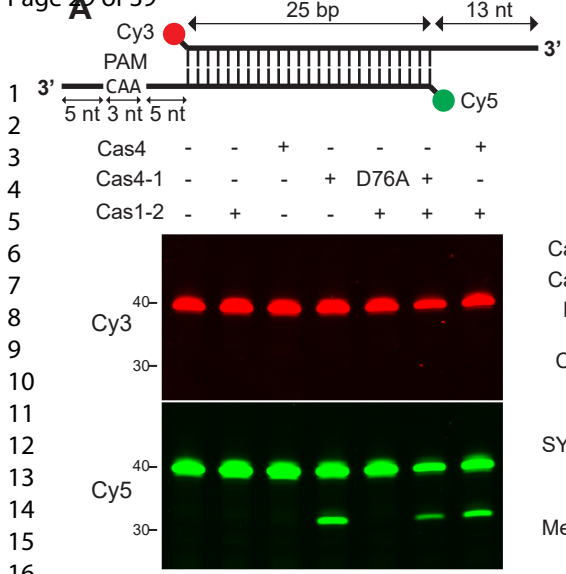


Figure 2



A

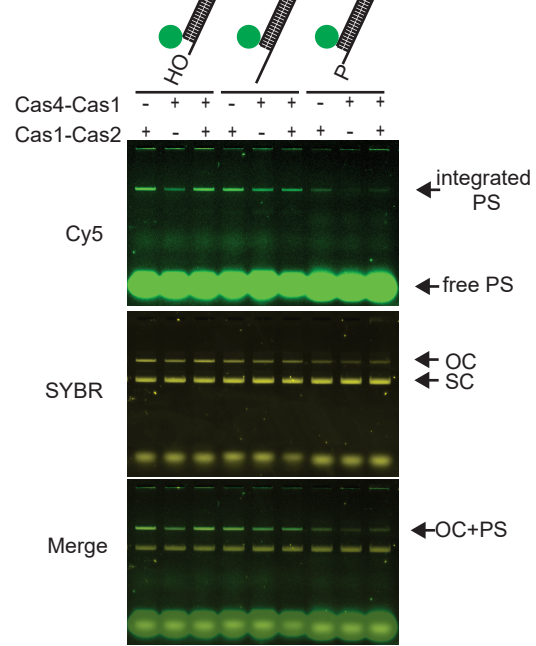


Figure 4

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

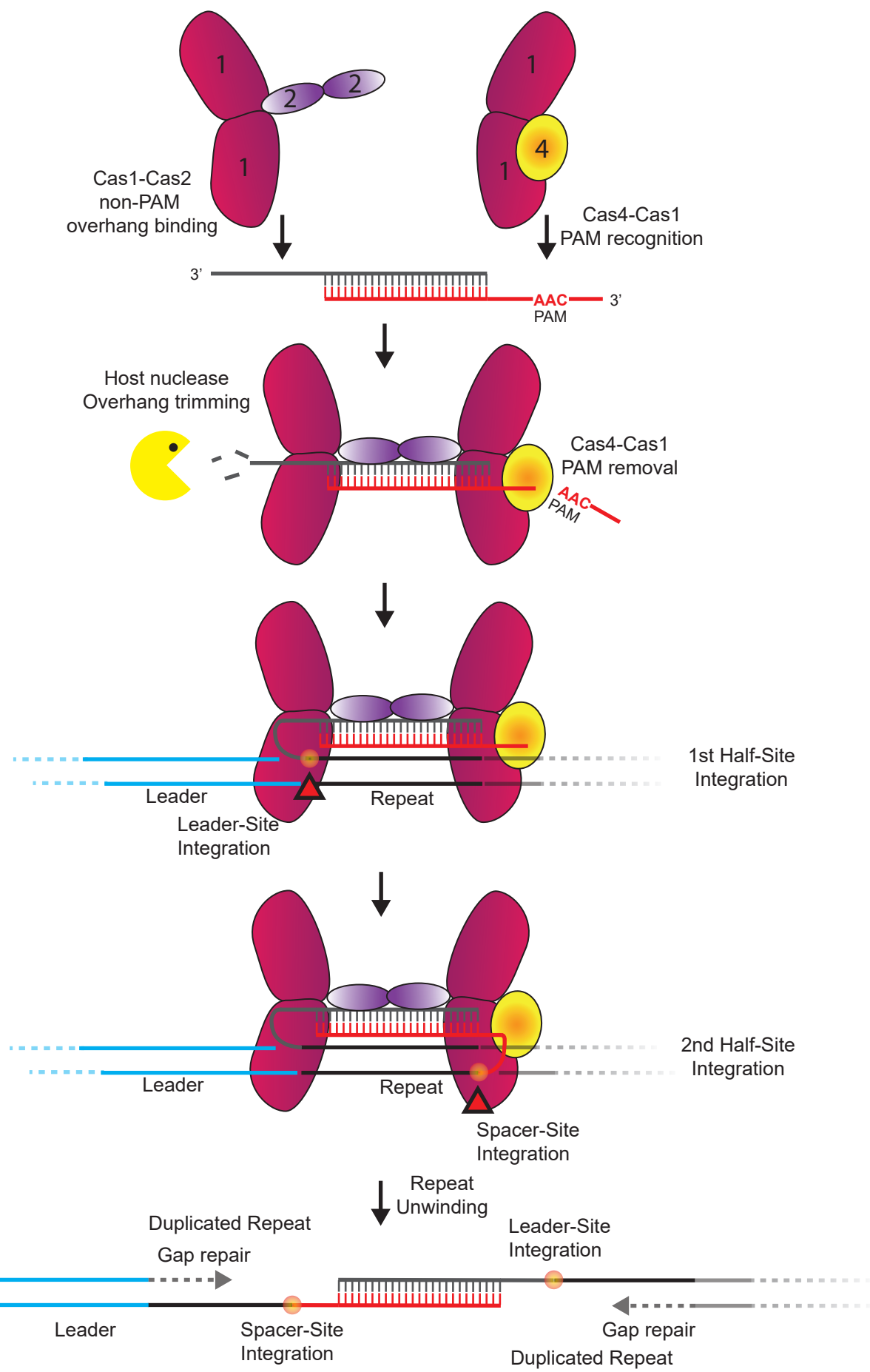


Figure 5



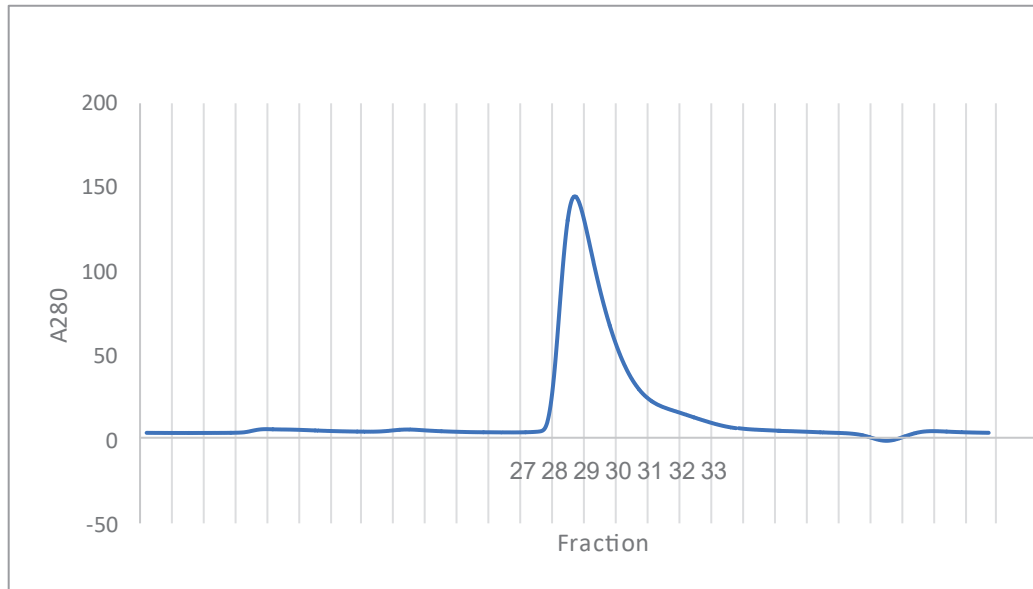
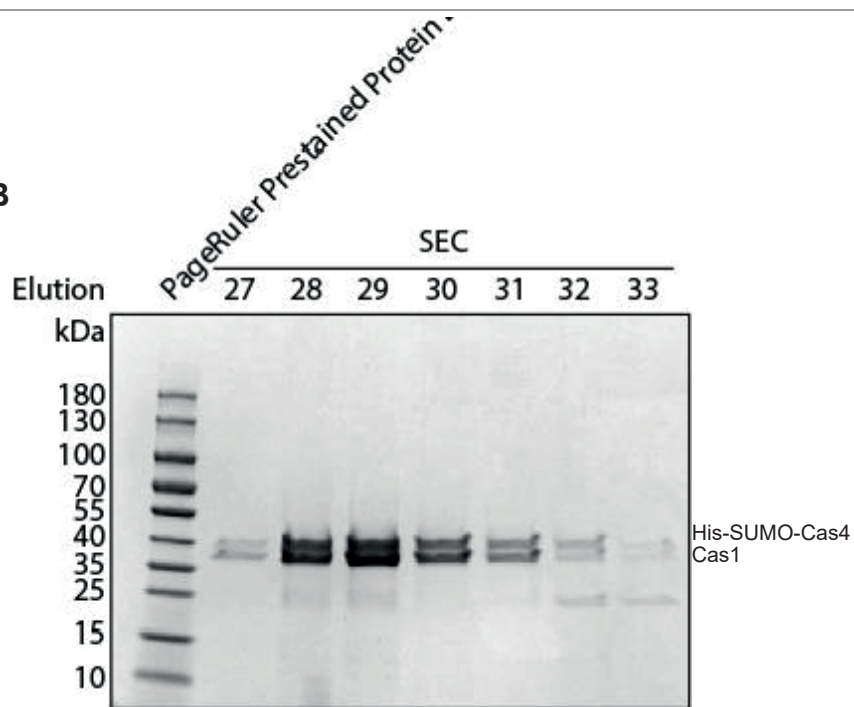
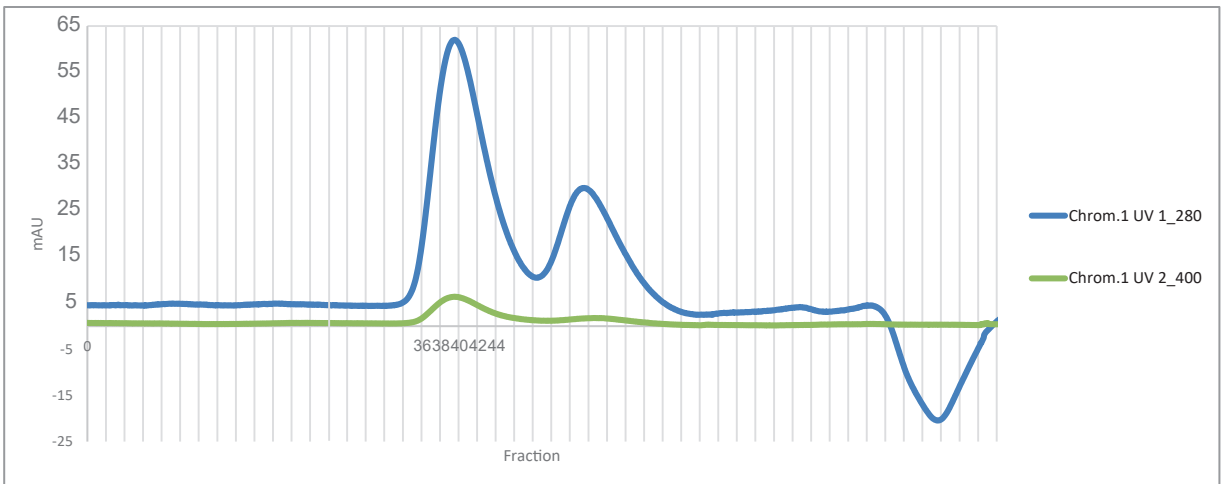
**A****B**

Figure S1

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**A**



**B**

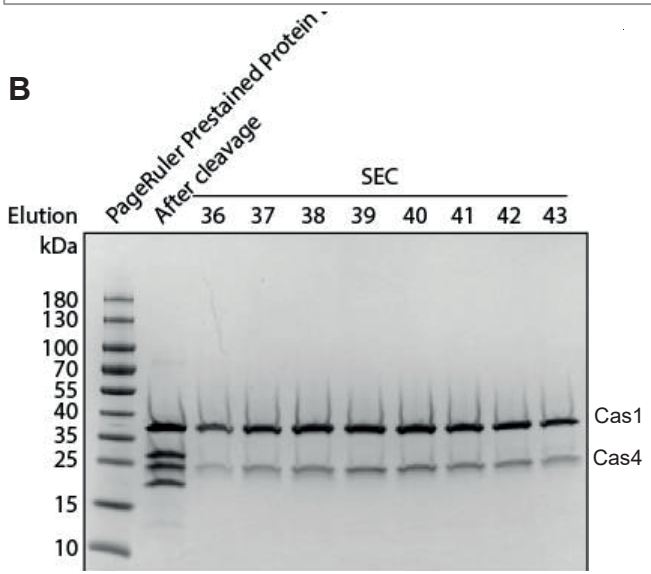


Figure S2

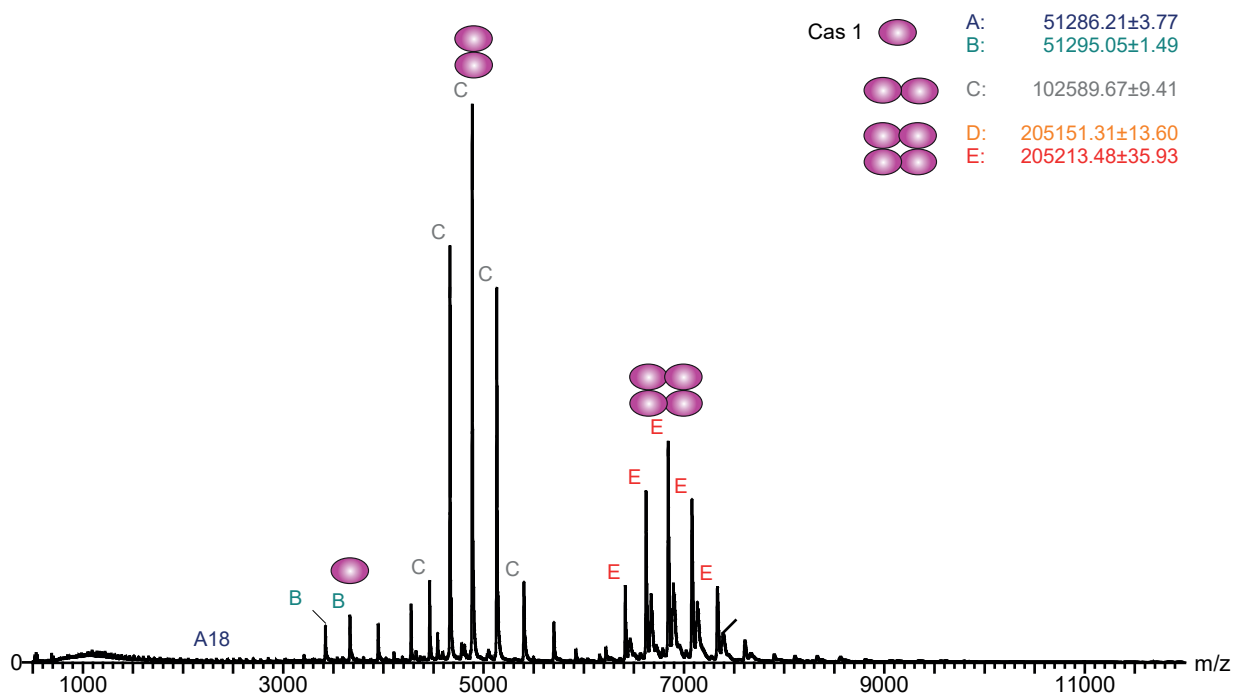


Figure S3

## 1 **Figure Legends**

### 2 **Figure 1 – Prespacer processing and half-site integration in vivo**

3  
4  
5  
6  
7  
8 **A** Genetic organization of the type I-D CRISPR-locus. Genes constituting the interference  
9  
10 machinery are located upstream of the adaptation complex. The adaptation complex consisting  
11  
12 of *cas4*, *cas1* and *cas2* is highlighted in purple. Downstream of *cas2* is the leader sequence  
13  
14 followed by the type I-D array.  
15  
16

17  
18 **B** Experimental design of spacer processing in vivo assay. Idealized pre-spacer substrates were  
19  
20 electroporated into *E. coli* cells carrying the plasmid encoded minimalized type I-D array and  
21  
22 expressing the adaptation genes *cas1* and *cas2*. The Cas4 protein was either omitted or co-  
23  
24 expressed as the wild-type or D76A+K91A mutant protein. Half-site integration was assessed by  
25  
26 PCR as depicted in **C**.  
27  
28  
29

30  
31 **C** PCR scheme allowing differentiation of spacer orientation, integration site and prespacer  
32  
33 processing. PCR amplicons of correct (PAM overhang integrated at Repeat-Spacer junction) and  
34  
35 incorrect spacer orientations and Leader-Repeat (L-R) or Repeat-Spacer (R-S) integration site  
36  
37 were subjected to high-throughput sequencing.  
38  
39

40  
41 **D-E** Overhang processing resulting from high-throughput sequencing of *in vivo* half-site  
42  
43 integration PCR. Bar charts indicate the overhang nucleotide left after processing and integration  
44  
45 as a percentage of the total number of sequenced spacer integration events. Red bars represent  
46  
47 correctly trimmed PAM-containing 3' overhangs. Integration was assessed in correct (D) or  
48  
49 incorrect (E) spacer orientation resulting from either Cas1-Cas2 or Cas1-Cas2 co-expressed with  
50  
51 Cas4 wild-type or Cas4 mutant (MT) background. n = number of sequenced integration events.  
52  
53

54  
55  
56  
57  
58  
59  
60  
23

1  
2  
3 **24 Figure 2 - Size exclusion chromatograms and SDS-PAGE analysis of peak fractions**

4  
5 **25 A** N-terminally tagged \*Cas1 associates with untagged Cas2 in the absence of Cas4. Unnumbered  
6  
7  
8 **26** peak contains protein aggregates.

9  
10 **27 B** Complex formation of N-terminally tagged \*Cas1 and untagged Cas4.

11  
12 **28 C** Co-expression of \*Cas1, Cas2 and Cas4. \*Cas1 elutes separately (peak 1) from the Cas4-Cas1  
13  
14  
15 **29** complex (peak 2). Cas2 together with dissociated Cas1 and Cas4 elutes as a low molecular weight  
16  
17  
18 **30** peak (peak 3). Unnumbered peak contains protein aggregates.

19  
20 **31 D** Native Mass Spectrometry of Cas4-Cas1 complex as shown in B (native spectrum obtained after  
21  
22  
23 **32** removal of His-SUMO tag from Cas1 by TEV protease cleavage (Fig. S2)). Cas4 monomers  
24  
25 **33** assemble with Cas1 dimers into a Cas4<sub>1</sub>-Cas1<sub>2</sub> complex. Cas1 dimers are also frequently  
26  
27  
28 **34** observed. Free monomers of Cas1 and Cas4 are less frequent in the measured sample.

29  
30 **35 Figure 3 – In vitro pre-spacer cleavage and integration**

31  
32 **36 A** Pre-spacer model substrate containing a 25 bp duplex region flanked by 13 nt 3' overhangs  
33  
34  
35 **37** incubated with different protein combinations. The non-PAM strand is 5' Cy3 labelled and the  
36  
37  
38 **38** PAM-containing strand 5' Cy5 labelled. Cleavage of PAM containing 3' overhang results in Cy5  
39  
40  
41 **39** labelled fragment of 30 nt. Protein samples consist of co-purified Cas4-Cas1 and Cas4D76A-Cas1  
42  
43  
44 **40** complexes, combined Cas4-Cas1 and Cas1-Cas2 complex and individually purified Cas4 in the  
45  
46  
47 **41** presence of Cas1-Cas2.

48  
49 **42 B** Integration of labelled pre-spacer (PS) into linear CRISPR DNA consisting of the type I-D leader  
50  
51  
52 **43** sequence (L) and a single Repeat (R). Labelled spacer imaged via Cy5, total DNA via SYBR gold  
53  
54  
55 **44** stain. Merge of Cy5 and SYBR gold channels indicates integration of Cy5 labelled pre-spacer  
56  
57  
58 **45** resulting in a higher molecular weight band.

59  
60

1  
2  
3 46 **C** Labelled pre-spacer (PS) integration into pCRISPR DNA. Similar to B, both adaption complexes  
4  
5  
6 47 facilitate integration into plasmid encoded CRISPR locus. Integration reaction is accompanied by  
7  
8 48 nicking of supercoiled (SC) plasmid DNA, resulting in formation of open-circular (OC) plasmid  
9  
10  
11 49 conformation. Merge image of Cy5 and SYBR channels shows co-localization of OC plasmid  
12  
13 50 species and Cy5 labelled spacer substrate.

15 51 **D-E** High-throughput sequencing of Cas4-Cas1 integrated pre-spacers. Reaction was performed  
16  
17  
18 52 similar to the assay shown in **C** using unlabeled pre-spacer DNA. OC plasmids were gel extracted  
19  
20  
21 53 followed by PCRs specific for the leader-repeat (L-R) and repeat-spacer (R-S) integration as well  
22  
23 54 as correct and incorrect spacer orientation. Bar graphs represent the percentage of spacers with  
24  
25 55 specific overhang length depending on integration site and orientation (D-correct; E-incorrect).  
26  
27  
28 56 n = number of sequenced integration events.

29  
30 57

31  
32  
33 58 **Figure 4 - Spacer overhang preferences of Cas1-Cas2 and Cas4-Cas1 complexes.**

34  
35  
36 59 **A** - Integration activity with respect to 3' overhang requirements. Pre-spacer substrates were 5'  
37  
38 60 Cy5 labelled in order to follow coupling to pCRISPR. Phosphorylated (P) 3' overhangs were used  
39  
40  
41 61 to block integration of one of the DNA strands.

42  
43 62

45 63 **Figure 5 – Model of Cas4-Cas1 and Cas1-Cas2 assisted spacer selection, processing and**  
46  
47  
48 64 **integration.** Prespacers with long PAM- and non-PAM containing 3' overhangs are bound by  
49  
50  
51 65 Cas4-Cas1 (PAM overhang) and Cas1-Cas2 (non-PAM overhang). Following processing by host-  
52  
53 66 factor nucleases, the non-PAM site of the spacer is integrated at the leader-repeat site (first half-  
54  
55 67 site integration). Subsequently, the second half-site integration (spacer-site integration) of the

1  
2  
3 68 PAM-site occurs, likely orchestrated by release of the Cas4-processed overhang into the integrase  
4  
5 69 site of Cas1. Unwinding of the repeat followed by gap-repair completes repeat duplication and  
6  
7  
8 70 full-site spacer integration.  
9

10 71

11  
12 72

13 73

14  
15 74

16  
17  
18  
19  
20 75 **Supplementary Figures Legends**

21  
22 76 **Figure S1 – Related to Fig. 2B** – Co-purification of His<sub>6</sub>-SUMO-Cas4 and untagged  
23  
24 77 Cas1. **A** Size exclusion chromatogram of His<sub>6</sub>-SUMO-Cas4 and untagged Cas1. **B** SDS  
25  
26 78 PAGE analysis of SEC purified proteins.  
27  
28

29 79

30  
31 80 **Figure S2 – Related to Fig. 2D** – Co-purification of untagged Cas4 and untagged Cas1  
32  
33 81 after TEV-protease cleavage of His<sub>6</sub>-SUMO-TEV tag. **A** Size exclusion chromatogram of  
34  
35 82 Cas4 and Cas1 (A280 – total protein; A400 – Cas4 FeS-cluster). **B** SDS PAGE analysis  
36  
37 83 of SEC purified proteins.  
38  
39

40 84 **Figure S3 – Related to Fig 2A** – Native Mass Spectrometry of Cas1-Cas2 complex as  
41  
42 85 shown in Fig. 2A.  
43  
44

45 86  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**Table S1** – Plasmids used in this study

Name in this study	Name	Insert	Vector	Resistance	Source
pCas2	pTU084	Synechocystis PCC6803 Type I-D <i>cas2</i> (delta <i>Cas1</i> )	pET-T7	Amp	Kieper et al. (2018)
pCas1	pTU085	Synechocystis PCC6803 Type I-D <i>cas1</i> (delta <i>Cas2</i> )	pET-T7	Amp	Kieper et al. (2018)
pCas1	pTU092	Synechocystis PCC6803 Type I-D <i>cas1</i>	pET-T7	Spec	This study
pCas4 <sup>D76A</sup>	pTU086	Synechocystis PCC6803 Type I-D <i>cas4</i> (D76A)	pET-T7	Spec	Kieper et al. (2018)
pCas4 <sup>D76A+K91A</sup>	pTU411	Synechocystis PCC6803 Type I-D <i>cas4</i> (D76A+K91A)	pET-T7	Spec	This study
pCas4	pTU130	Synechocystis PCC6803 Type I-D <i>cas4</i>	pET-T7	Spec	Kieper et al. (2018)
pCRISPR	pTU134	Synechocystis PCC6803 Type I-D Leader-R-S1	pACYCDuet1	Cm	Kieper et al. (2018)
pCas1-2	pTU70	Synechocystis PCC6803 Type I-D <i>cas1-cas2</i>	pET-T7	Amp	Kieper et al. (2018)
pEmp	pTU116	NA	pET-T7	Spec	Addgene Plasmid #48329

**Table S2** - Oligonucleotides used in this study

Name	Sequence	Description
BN015	CGTCCATGGGAAGTCATTCTTCAAATTTTGGC	Leader Fw
BN277	GTGGAATACGCAAAAGGC	<i>cas4</i> mutagenesis K91A Fw
BN278	AGGAATTAATAAGCCATCACTTTC	<i>cas4</i> mutagenesis K91A Rv
BN1398	GCTAGTTATTGCTCAGCGG	pCRISPR bb Rv
BN1711	GGAAGGTTTGCCAAAGTC	Leader Distal Half-Site Integration
BN1712	CTGTTGACTTAAGCATTATGC	Leader Proximal Half-Site Integration
BN1713	ATCGACACCACCACG	OligoSpecific Primer Fw (PAM overhang)
BN1714	CGTGGTGGTGTGCGAT	OligoSpecific Primer Rv (non-PAM overhang)
BN1763	CTACCATCGACACCACCACGCTGGCTTTTTAACTTTTT	25 nt duplex 13nt PAM 3' ovhng
BN1768	GCCAGCGTGGTGGTGTGCGATGGTAGTTTTTTGTTTTT	25 nt duplex 13nt RvC PAM 3' ovhng
BN1829	CTACCATCGACACCACCACGCTGGCTTTTTTTGTTTTT	25 nt duplex 13nt RvC PAM 3' ovhng (5' C6-Amino)
BN1830	GCCAGCGTGGTGGTGTGCGATGGTAGTTTTTTAACTTTTT	25 nt duplex 13nt PAM 3' ovhng (5' C6-Amino)