

Delft University of Technology

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

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.

2 3 1 Classification: DIOLOGICAL SOLENCES: Microbiology					
3 4	1	1 Classification: BIOLOGICAL SCIENCES: Microbiology			
5	2	Research article			
6 7 8 9 10 11 2 13 14 15 16 17 8 9 21 22 34 25 26 27 28 9 0 31 22 33 4 56 27 28 9 0 31 22 33 45 36 37 8 9 40 41 42 43 44 5 6 7 8 9 0 12 23 45 56 57 8 9 50 57 8 9 57 57 8 9 57 57 8 9 57 57 57 8 9 57 57 57 57 57 57 57 57 57 57 57 57 57	3 4	Cas4-Cas1 is a PAM-processing factor mediating half-site spacer integration			
	5	during CRISPR adaptation			
	5				
	7	Sobaction N. Kiopor ^{1,2} Crictobal Almondros ^{1,2} Anna C. Haagsma ^{1,2} Arian			
	/	Derendrogt ^{3,4} Albert J.D. Heek ^{3,4} Sten J.J. Bround ^{1,28}			
	8	Barendregt ^{e, 4} , Albert J.R. Heck ^{e, 4} , Stan J.J. Brouns ^{1,23}			
	9				
	10	¹ Department of Bionanoscience, Delft University of Technology, Delft, Netherlands			
	11	- Kavii institute of Nanoscience, Delitt, Netherlands.			
	12	³ Biomolecular Mass Spectrometry and Proteomics, Bijvoet Center for Biomolecular			
	13	Research, Utrecht Institute of Pharmaceutical Sciences, Utrecht University, Utrecht,			
	14	Netherlands.			
	15	⁴ 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			
59 60					

22 Abstract

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

Introduction

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

against HHPV-2 (Li et al., 2014). Biochemical evidence was provided by Plagens et al. showing a protein-protein interaction in vitro between Cas4 and the type I-A adaptation fusion protein Cas1/2 and Csa1 demonstrating that Cas4 directly interacts with the adaptation machinery (Plagens et al., 2012). Recently, several studies defined the role of Cas4 in more detail, finding that the presence of Cas4 increases the fidelity of spacer integration. Specifically, the Cas4 protein of a cyanobacterial type I-D CRISPR-Cas system facilitated the integration of interference-proficient spacers that carry the consensus PAM of the type I-D CRISPR-Cas system. Spacers acquired in the presence of Cas4 displayed shorter lengths compared to those acquired in the absence of Cas4 or in the presence of catalytically inactive variant of the protein (Kieper et al., 2018). Additionally, two Cas4 variants (Cas4-1 and Cas4-2) encoded in the type I-A system of Pyrococcus furiosus were shown to define the upstream protospacer adjacent motif (PAM) and a downstream NW motif in vivo (Shiimori et al., 2018). Deletion of cas4-1 and cas4-2 resulted in incorrect processing of prespacers with respect to up- and downstream motifs, random orientation of the integrated spacer as well as large deviations from the consensus spacer length (Shiimori et al., 2018). Previously, biochemical studies of the Cas4 protein, containing an iron-sulfur cluster and a RecB domain, found several nuclease activities, demonstrating endo- and exonuclease activities (Lemak et al., 2013; Lemak et al., 2014; Zhang et al., 2012). The requirement of Cas4 for prespacer processing was therefore in accordance with the previously described biochemical activities. Indeed, Lee et al. provided the first mechanistic details of how Cas4 proteins ensure PAM-processing and correct spacer orientation (Lee et al., 2019; Lee et al., 2018). It was shown that Cas4 tightly interacts with the Cas1 integrase, forming a heterohexameric complex composed of two Cas1 dimers and two Cas4 subunits (Lee et al., 2018). This complex would interact with

double-stranded prespacer substrates and endonucleolytically cleave PAM sequences in long 3' overhangs, ensuring that only PAM-processed spacers would be eventually integrated into the CRISPR array (Lee et al., 2018). Interestingly, the authors did not see interaction of the Cas4-Cas1 complex with Cas2 in their initial experiments, suggesting the possibility that a prespacer is required for assembly of the full complex. After supplying a dsDNA substrate to the three adaptation proteins, Lee et al. could demonstrate the assembly of the full Cas4-Cas1-Cas2 complex (Lee et al., 2019). This complex was shown to assemble in a mixture of symmetric and asymmetric architectures as shown by negative-staining Electron Microscopy, in which the asymmetric complex would contain only a Cas4 monomer associated with one of the Cas1 dimers (Lee et al., 2019). The authors suggested that this asymmetry might aid in the differential processing of prespacer substrates in which the Cas4 containing half of the complex would interact with the PAM-containing overhang. This hypothesis is supported by the findings in the type I-A system, in which two independent Cas4 homologs are dedicated processing factors for the PAM- and NW motif containing prespacer overhangs (Shiimori et al., 2018). However, how this asymmetric processing is orchestrated in CRISPR systems containing only a single cas4 gene is currently unknown. In this work we provide mechanistic insights of an asymmetric complex, specifically how the Cas4-Cas1 complex is able to recognize and sequence specifically process the PAM sequence of the type I-D CRISPR-Cas system. Previously we have shown that the type I-D Cas4 protein facilitates the integration of PAM-compliant spacers in vivo (Kieper et al., 2018). We demonstrate that Cas4 strongly interacts with the Cas1 integrase forming a heteromeric Cas4₁-Cas1₂ complex. This heteromeric complex does not require the Cas2 protein for processing and half-site integration of PAM-containing prespacer substrates. The catalytic activity of Cas4 is required for

prespacer cleavage and is crucially dependent on the presence of Cas1 in order to recognize and process the PAM overhang. We show that this Cas4-Cas1 complex does not cleave the non-PAM containing overhang. Processing of the non-PAM containing overhang potentially relies on the Cas1-Cas2 complex and likely requires host-factor nucleases. We provide a model in which an asymmetric adaptation complex differentially acts on PAM and non-PAM containing overhangs, providing cues for the correct orientation of spacer integration. This correct PAM processing as well as a functional orientation explains the importance and hence the strong conservation of the cas4 gene, increasing the integration of interference-proficient spacers.

Material & Methods

Bacterial strains and growth conditions

E. coli strains DH5α and BL21 were grown in Lysogeny Broth (LB) at 37°C and continuous shaking at 180 rpm or grown on LB agar plates (LBA) containing 1.5% (wt/vol) agar. When required, the media were supplemented with 100 µg ml⁻¹ ampicillin, 50 µg ml⁻¹ spectinomycin, 25 µg ml⁻¹ chloramphenicol (see Table S1 for plasmids and corresponding selection markers).

Plasmid construction and transformation

Plasmids used in this study are listed in Table S1. All cloning steps were performed in *E. coli* DH5a. Primers described in Table S2 were used for PCR amplification of the type I-D CRISPR-Cas locus (cas4, cas1, cas2 and leader-repeat-spacer1) from Synechocystis cell material using the Q5 high-fidelity Polymerase (New England Biolabs). PCR amplicons were subsequently cloned into Berkeley MacroLab LIC (https://gb3.berkeley.edu/facility/gb3-macrolab/) using vectors either ligation-independent cloning (LIC), or into the pACYCDuet-1 vector system (Novagen (EMD Millipore) using conventional restriction-ligation cloning. The cas4^{D76A+K91A} mutant was obtained using a PCR-based mutagenesis of pCas4^{D76A} using primers listed in Table S2. All plasmids were verified by Sanger-sequencing (Macrogen Europe, Amsterdam, The Netherlands). Bacterial transformations were either carried out by electroporation (2.5 kV, 25 mF, 200 V) using an ECM 630 electroporator (BTX Harvard Apparatus) or chemically competent cells prepared according to manufacturer's manual (Mix&Go, Zymo research). Electrocompetent cells were prepared following a protocol adapted from (Gonzales et al., 2013). Transformants were selected on LBA supplemented with appropriate antibiotics.

³ 152 **Protein expression and purification**

Plasmid encoded cas genes were either co-expressed or expressed individually in E. coli BL21 AI cells (Invitrogen). Pre-cultures were grown from individual colonies and used for inoculation of pre-warmed (37°C) LB medium at an initial OD₆₀₀=0.05. Protein expression was induced at OD₆₀₀=0.5 by addition of IPTG and L-arabinose preceded by a 30-minute cold-shock. Cultures were subsequently grown overnight at 20°C with continuous shaking. Cells were harvested by centrifugation (10 min, 4°C, 2400xg) and subsequently resuspended in lysis buffer (50 mM HEPES pH 7.5, 300 mM KCl, 5% Glycerol, 1 mM DTT, 25 mM Imidazole, 0.1% Triton-X 100) supplemented with cOmplete[™], EDTA-free Protease Inhibitor Cocktail (Roche). Cells were lysed by two passages through a CF1 cell disruptor (Constant Systems Ltd.) equilibrated with lysis buffer at a constant pressure of 1 kbar. Lysates were cleared by centrifugation (45 min, 4°C, 25000xg) and filtered through 0.45 µm filter. Protein was bound in batch to HIS-Select (Sigma Aldrich) IMAC resin for 30 min at 4°C and rotary shaking. IMAC resin was then loaded onto Pierce gravity-flow columns (Thermo Scientific) and washed with 10 CV wash buffer (50 mM HEPES pH 7.5, 300 mM KCl, 5% Glycerol, 1 mM DTT, 50 mM Imidazole). Proteins were subsequently block eluted in 0.5 ml elution buffer (50 mM HEPES pH 7.5, 300 mM KCl, 5% Glycerol, 1 mM DTT, 250 mM Imidazole). Protein concentration and purity was determined by NanoDrop A280 spectroscopy and SDS PAGE analysis. Protein elution fractions were pooled and subjected to size exclusion chromatography using Superdex 200 10/300 GL (GE Healthcare) column with 0.5 ml/min flow rate using elution buffer as mobile phase. Cas1-Cas2 complex IMAC elution fractions used for integration assays were prepared for ion-exchange chromatography by adjusting the KCL concentration to 30 mM and subsequently loaded onto HiTrap Heparin HP column (GE Healthcare). Cas1-Cas2 complexes were eluted by gradually increasing KCL concentration to 1 M. Resulting fractions

1		
2 3 4	178	were analyzed by SDS-PAGE and appropriate fractions pooled, snap frozen and
5 6 7	179	stored at -80°C.
8	180	Native Mass Spectrometry
9 10 11 12	181	Cas4-Cas1 and Cas1-Cas2 complexes were buffer exchanged into 500 mM
13 14 15	182	ammonium acetate (pH 7.5) using seven sequential steps on a centrifugal filter with
16 17 18	183	a molecular weight cut-off of 10 kDa (Sartorius) at 4°C. MS measurements were
20 21 22	184	performed in positive mode by directly infusing the individual complexes at a
23 24 25 26	185	concentration of 1 μM using an LCT electrospray time-of-flight (Waters, United
20 27 28 29	186	Kingdom) adjusted for optimal performance in high mass detection (Tahallah et al.,
30 31 32 33	187	2001; van den Heuvel et al., 2006). The needles used for electrospray were
34 35 36 37	188	prepared in house from borosilicate capillaries (Kwik-Fil, World Precision
38 39 40	189	Instruments, Sarasota, FL) on a P97 puller (SutterInstruments, Novato, USA) and
41 42 43 44	190	gold coated by using an Edwards Scancoat Six Pirani 501 Sputter Coater (Edwards
45 46 47 48	191	Laboratories, Milpitas, USA). During the measurement the capillary voltage was
49 50 51	192	kept at 1200V, cone voltage between 80-150V and the source pressure was
52 53 54 55	193	increased to \approx 8mbar. Exact mass measurements of the individual Cas proteins
56 57 58 59 60	194	were acquired under denaturing conditions by adding formic acid to a final

concentration of 5%. All spectra were mass calibrated by using an aqueous solution of cesium iodide (25 mg/ml). Mass spectra were accumulated, averaged, smoothened and centered, using the software MassLynx 4.1 (Waters, United Kingdom). Nuclease assays Oligo-nucleotide sequences used in this study are indicated in Table S2. Oligo-nucleotides with C6-Amino modifications on the 5' terminus were obtained from ELLA Biotech (Planegg, Germany). Cy5 or Cy3 (GE Healthcare) labelling of 5' termini was done in 100 mM Sodium-bicarbonate buffer as described by (Joo and Ha, 2012). Unlabeled oligo-nucleotides were obtained from Integrated DNA Technologies (IDT). Nuclease assays were performed in buffer R (5 mM HEPES pH 7.5, 100 mM Sodium-Glutamate supplemented with 2 mM MnCl₂ and 10 mM MgCl₂. Annealed and Cy5 and Cy3 labelled oligo-nucleotides (Cy3-BN1829+Cy5-BN1830) were added to a final concentration of 125 nM and purified protein complexes to a final concentration of 500 nM. Reactions were incubated for 1 hour at 30°C after which reactions were quenched by addition of Proteinase K (Thermo Fischer) and incubation for 1 hour at 37°C. The resulting products were analyzed on denaturing PAGE (10% acrylamide, 8M Urea) and analyzed with Amersham Typhoon fluorescence gel scanner (GE Healthcare). In vitro spacer integration assays Oligo-nucleotide integrations with either Cy5 labeled or unlabeled oligo-nucleotides

were performed by pre-incubating indicated protein complexes (500 nM) with oligo-nucleotides (250 nM) on ice for 15 min. Following pre-incubation, either linear CRISPR substrate (obtained by Q5 high-fidelity PCR from pCRISPR using primers BN015+BN1398) or supercoiled pCRISPR were added to a final concentration of 7.5

nM. Reaction mixtures were incubated at 30°C for 1 hour after which reactions were guenched by addition of Proteinase K (Thermo Fischer) and incubation for 1 hour at 37°C. Reactions were run on 1% native agarose gels for 45 min and gels subsequently stained with SYBR gold (Sigma Aldrich). Gels were scanned for Cy5 and SYBR gold using Amersham Typhoon fluorescence gel scanner (GE Healthcare). For PCR analysis of *in vitro* integration, unlabeled oligo-nucleotides were used in the reaction. Open-circular plasmid DNA was gel isolated and DNA purified using Zymoclean gel recovery kit (ZymoResearch) after which integration was assessed by PCR using primers BN1711+BN1713 (leader distal integration; correct spacer orientation), BN1711+BN1714 (leader distal integration; incorrect spacer orientation). BN1712+BN1713 (leader proximal integration; incorrect spacer orientation) and BN1712+BN1714 (leader proximal integration; correct spacer orientation). Purified PCR amplicons were subjected to MiSeq sequencing (Illumina).

Next generation sequencing and statistical analysis

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

Page 12 of 39

the in vitro integration assay. Overhang processing was analyzed by annotating the
 primers used for amplification and comparing the overhangs post-integration to the
 initial oligo-nucleotide sequence.

10 248

248 In vivo spacer integration assays

E. coli BL21 AI cells were co-transformed with either pCas1-Cas2 and pEmpty or pCas1-2 and pCas4 (wild-type Cas4 or Cas4D76A+K91A). One transformant for each combination was grown in LB at 37°C and continuous shaking (180 rpm) to OD₆₀₀=0.3 and made electrocompetent after which pCRISPR was transformed. For each treatment three individual colonies were grown in SOB medium (LB supplemented with 10 mM MgSO4 and 10 mM MgCl2) at 37°C and continuous shaking (180 rpm) to OD₆₀₀=0.3 after which protein expression was induced by addition of 0.2% L-arabinose and 0.5 mM IPTG. Induced cultures were grown for additional 2 hours at 37°C and continuous shaking (180 rpm). Cells were made electrocompetent and annealed pre-spacer oligo-nucleotides (BN1763+BN1768) electroporated at a final concentration of 1 µM. After 30 min recovery cells were harvested and plasmid DNA extracted using GeneJET plasmid miniprep kit (Thermo Scientific). Extracted plasmid DNA was normalized to 0.5 ng µl⁻¹ and subsequently 2 µl used in half-site integration PCRs using primers BN1711+BN1713 (leader distal integration; correct spacer orientation), BN1711+BN1714 (leader distal integration; incorrect spacer orientation). BN1712+BN1713 (leader proximal integration; incorrect spacer orientation) and BN1712+BN1714 (leader proximal integration; correct spacer orientation). PCR amplicons were validated by agarose gel electrophoresis and purified with the GeneJET PCR Purification kit (Thermo Scientific). Purified PCR amplicons were subjected to MiSeq sequencing (Illumina).

Results

PAM-containing overhang processing depends on orientation of spacer integration

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

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

Page 14 of 39

typically broader and shifted towards longer overhang lengths. Analyzing the non-PAM containing overhangs of correctly oriented spacers did not display any differences between the conditions (with and without Cas4), suggesting that prespacer overhangs without PAM are not processed by Cas4, but rather by endogenous *E. coli* nucleases. Spacers integrated in the incorrect orientation (Fig. 1E) showed similar 3' overhangs under all conditions. We observed most accurate processing when Cas4 was present in its active form. In those samples the presence of Cas4 led to an increase in the shortening of overhangs, with a predominant overhang length of 6 nucleotides. In the E. coli model system, host factor nucleases potentially act on both PAM and non-PAM containing 3' overhangs that remain unprotected by the core Cas1-Cas2 complex holding the prespacer. Cas4 does not specifically cleave non-PAM containing overhangs, but requires the presence of the PAM in order to engage in sequence-specific processing.

Cas4 forms a strong heteromeric complex with Cas1

In order to assess whether the overhang processing connected to the presence of Cas4 was a result of Cas4 specifically interacting with the Cas1-Cas2 integration complex, we first investigated the formation of a Cas1-Cas2 complex. The Cas1 protein was N-terminally His₆-tagged and co-expressed with Cas2 in *E. coli* BL21-AI cells. After the initial nickel-affinity pull-down from cleared cell lysate, the elution fraction was subjected to size exclusion chromatography (SEC), which resulted in one peak containing aggregated protein and another peak species (Fig. 2A). This peak contained three proteins (Fig. 2A) for which the tagged-Cas1, untagged Cas1 (due to autoproteolysis of the tag) and Cas2 protein identity was confirmed by mass spectrometry. Next, we co-expressed the His6-tagged Cas1 with untagged Cas4 and observed strong co-purification of both proteins (Fig. 2B). In order to verify the Cas4-

Cas1 interaction, a reverse tagging strategy was used (His6-tagged Cas4 co-expressed with untagged Cas1), which again confirmed the presence of a Cas4-Cas1 complex (Fig. S1). When tagged-Cas1 and Cas2 were co-expressed along with Cas4, we observed a strong co-purification of Cas4 and Cas1 that abolished formation of the Cas1-Cas2 complex since Cas2 eluted separately as a low molecular weight species. This fraction also contained minor amounts of Cas1 and Cas4 that did not assemble into higher order complexes. Our results demonstrate that under these conditions Cas1 can form complexes with Cas4 or Cas2, and that these complexes appear to be mutually exclusive. In the presence of both Cas4 and Cas2, Cas1 strongly favors the interaction with Cas4 over the interaction with Cas2.

Cas4 associates with Cas1 in a 1:2 ratio

Next, we determined the stoichiometry of the formed Cas4-Cas1 complex. Previously, Lee et al. demonstrated that the heteromeric complex consists of two Cas1 dimers that each associate with a single Cas4 monomer (Lee et al., 2018). To gain insight into the composition of the untagged Cas4-Cas1 complex (Fig. S2) native protein mass spectrometry analysis was performed (van den Heuvel et al., 2006). The mass spectrum (Fig. 2D) revealed a distribution of different complex species, with the most abundant mass-over-charge (m/z) peaks consisting of either Cas1 dimers (73.6 \pm 1.6 kDa) or the Cas4-Cas1 complex consisting of a single Cas1 dimer and a Cas4 monomer resulting in a Cas4₁-Cas1₂ complex of 96.3 kDa (Fig. 2D). Even though we observed co-purification of Cas1 and Cas2 in the SEC analysis, the native mass spectrum of the Cas1-Cas2 complex resulted in mainly Cas1 dimers (Fig. S3) with Cas2 likely being lost during the native MS sample preparation.

Page 16 of 39

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

The acquisition of functional spacers not only requires appropriate prespacer selection, but also PAM-compliant processing. We have previously shown that the presence of Cas4 in addition to the core adaptation proteins Cas1 and Cas2 significantly increases the integration of spacers with a correctly processed PAM in vivo (Kieper et al., 2018). Due to the strong interaction of Cas4 and Cas1 we aimed to test whether PAM processing is mediated only by Cas4, or if the heteromeric Cas4-Cas1 complex is required. In order to address this question, we performed prespacer cleavage assays with a dual-labelled model prespacer (Fig. 3A). This model prespacer consisted of a 25 bp duplex flanked by 13 nucleotide 3' overhangs on each side and fluorescent labels at their 5' ends. The top strand was labelled with Cv3 and did not contain a PAM sequence in its 3' overhang of the, while the bottom strand was labelled with Cy5 and contained the I-D consensus PAM. We found that neither free Cas4 nor Cas1-Cas2 was able to catalyze 3' overhang cleavage. However, the addition of the Cas4-Cas1 complex resulted in a defined band corresponding to processing of the PAM sequence within the PAM-containing overhang (Fig. 3A). This result suggests that PAM recognition is mediated by the interactions within the Cas4-Cas1 complex, where Cas4 acts as the catalytic subunit of the complex.

We have previously shown that mutating D76 in the conserved RecB domain of Cas4 abolished integration of PAM-proficient spacers in vivo (Kieper et al. (2018). Using the D76A mutant in our in vitro cleavage assay fully abolished processing activity of the Cas4-Cas1 complex, demonstrating that the RecB domain of Cas4 is indeed the catalytically active site required for PAM processing. Interestingly, although Cas4 did not show processing activity on its own, combining Cas1-Cas2 and Cas4 fully restored processing to similar levels as the Cas4-Cas1 complex. The addition of both, the Cas1-

Cas2 and the Cas4-Cas1 complex, resulted in processing of the PAM overhang as observed with the Cas4-Cas1 complex alone or combination of Cas4 and Cas1-2 complex. All conditions that showed cleavage of the substrate resulted in a single defined band, suggesting that cleavage occurred via an endonuclease mechanism which is in line with previous studies (Lee et al., 2018). The processing of the non-PAM containing overhang was not observed in any of the conditions, indicating that the processing of the non-PAM site presumably relies on host factor nucleases such as DnaQ-like exonucleases or Exonuclease T as recently found in in the I-E system (Kim et al., 2020; Ramachandran et al., 2020). Our results demonstrate that sequence specific Cas4 activity requires the presence of Cas1 and that the Cas4-Cas1 complex is the core processing complex that sequence-specifically recognizes and processes the PAM sequence before integration.

386 The Cas4-Cas1 complex integrates new spacers into both linear and 387 supercoiled DNA 388 Next we tested whether the Oce4 Cos4 complex set only set only set.

Next, we tested whether the Cas4-Cas1 complex not only processes prespacer substrates but also catalyzes their integration into the CRISPR array. We performed adaptation assays using supercoiled plasmid DNA containing the type I-D leader and a single repeat (pCRISPR; Fig. 1A) as well as linear CRISPR array substrates generated by PCR (Fig. 3B). Both linear and plasmid CRISPR loci were incubated with a Cy5-labelled prespacer, the Cas4-Cas1 and Cas1-2 complexes. Intriguingly, we observed coupling of the labelled prespacer by the Cas4-Cas1 complex to both CRISPR substrates, showing that this sub complex is proficient in catalyzing at least half-site spacer integration (Fig. 3B&C). Spacer integration into plasmid DNA resulted in the formation of open-circular (OC) plasmid conformations. Merging the Cy5 signal of the prespacer and the plasmid DNA signal confirmed that the prespacer was indeed coupled to the OC form of the plasmid. The Cas1-2 complex was able to integrate the

Page 18 of 39

prespacer into both linear and supercoiled arrays similar to the Cas4-Cas1 complex. Our observation demonstrates that at least two different sub-complexes exist, which are both capable of catalyzing half-site spacer integration. Taken together, based on the selective PAM-overhang processing of the Cas4-Cas1 complex, we hypothesize that the Cas4-Cas1 complex processes and integrates the PAM containing overhang and the Cas1-Cas2 complex the non-PAM containing overhang.

Correct spacer orientation requires overhang processing prior to integration

In order to analyze the accuracy of spacer integration by the Cas4-Cas1 complex in more detail, OC plasmid resulting from the integration reaction was gel purified and subjected to half-site integration PCRs as described previously (Fig. 1C). PCR products were subjected to Illumina MiSeg sequencing and prespacer sequences were extracted. This approach allowed us to assess 3' overhang processing before integration at the leader-proximal or leader-distal integration site. Interestingly, PAM-containing overhangs only showed sequence specific processing when the spacer was correctly oriented with respect to the PAM (Fig. 3D). The Cas4-Cas1 complex cleaved 65% of correctly oriented spacers exactly downstream of the PAM, however, we also observed incorrect removal of a single nucleotide in 25% of sequences and removal of 2 or more nucleotides in 10% of the sequences. Surprisingly, incorrectly oriented spacers did not show any processing of the PAM-containing overhang (Fig. 3E), indicating that integration in the correct orientation is preceded by the processing of the overhang. As predicted from the bulk cleavage assays, we did not observe any processing of the non-PAM containing overhangs regardless of the spacer orientation. Our data show that integration of new spacers in the correct orientation by Cas4-Cas1 requires PAM recognition and processing before a spacer can be integrated.

Page 19 of 39

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

448 Discussion

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

spacers has been revealed only in the recent years (Kieper et al., 2018; Lee et al., 2019; Lee et al., 2018; Shiimori et al., 2018; Zhang et al., 2019). Here we provide a new mechanistic understanding of how Cas4-dependent PAM selection is achieved during CRISPR adaptation, and specifically how asymmetry of the adaptation complex drives the selection, processing and integration of PAM-compatible spacers. We present a model in which two independent subcomplexes. Cas4-Cas1 and Cas1-Cas2. selectively process the two 3' overhangs of a prespacer (Fig. 5). The interaction of Cas4 with the Cas1 integrase protein is central to the recognition and processing of PAM-containing prespacer substrates. Formation of this Cas4-Cas1 complex is mutually exclusive with formation of the Cas1-Cas2 complex, which may suggest distinct roles of both subcomplexes. We found that the Cas4-Cas1 subcomplex displays prespacer cleavage activity only on PAM-containing 3' overhangs. Cas4-Cas1 removes the PAM via endonuclease cleavage while Cas1-2 defines overhang trimming likely through host factor nucleases. Subsequently, Cas1-Cas2 initiates coupling of the non-PAM overhang to the leader-repeat junction followed by integration of the processed PAM-site at the repeat-spacer junction.

Our findings are consistent with previous studies that established the existence of two mutually exclusive Cas4-Cas1 and Cas1-2 complexes in type I-C CRISPR-Cas systems (Lee et al., 2019; Lee et al., 2018), and expand our understanding of the roles of these subcomplexes. Moreover, the RecB-domain mediated activity of Cas4 is dependent on the presence of Cas1, since Cas4 alone is not able to recognize and process the PAM sequence. This observation suggests that the Cas4-Cas1 interaction is essential for sequence specific recognition of the PAM. It remains to be determined whether the PAM sequence recognition domain is located within Cas4 or Cas1. Interestingly, PAM selection in the type I-E system is mediated by the C-terminal tail of

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

Cas4-Cas1 did not display activity on non-PAM containing overhangs in vitro, however, processing of the non-PAM overhang was observed in our *in vivo* setup, suggesting that processing involves other non-Cas proteins. This finding is in line with the Cas4-deficient type I-E system in which host factors such as the ExoT and DnaQ-like exonucleases are required for processing both, PAM-containing and non-PAM overhangs. (Kim et al., 2020; Ramachandran et al., 2020). We propose that, in analogy to the E. coli type I-E system, 3'-5' exonucleases act as trimming factors for non-PAM 3' overhangs in the native Synechocystis PCC6803 host. The Cas1 protein of the type I-E system recognizes and protects the PAM from premature trimming, causing a delayed processing of the PAM end that ensures correct orientation (Kim et al., 2020). Upon activation, the Cas4-Cas1 complex sequence specifically removes the type I-D PAM, although incorrect processing was observed in vivo and in vitro that would result in single-nucleotide slipped spacers. Recently, it was observed in the type I-F system that slipped spacers increase primed adaptation which enhances the spacer diversity of the population (Jackson et al., 2019). Our results suggest the possibility that such erroneous PAM processing could promote the integration of slipped spacers and by extension, primed adaptation as found in other type I CRISPR-Cas systems (Nussenzweig and Marraffini, 2020).

Cryo-EM structures of the type I-C Cas4-Cas1-Cas2 complex revealed that the complex might undergo a conformational change (e.g. causing dissociation of Cas4 from the complex) in order to allow for Cas1 mediated integration of the PAM-end site of the spacer. Lee et al. showed that 50% of their Cas4-Cas1-Cas2 complex structures

lacked the Cas4 density on one site of the complex, resulting in an asymmetric complex. Our observation of two integrase complexes (Cas4-Cas1 and Cas1-Cas2) that are independently capable of at least half-site spacer integration points towards a similar asymmetrical organization of the full adaptation complex, in which Cas4-Cas1 is involved in PAM-site and Cas1-Cas2 in non-PAM site integration. By testing asymmetric spacer precursors, we demonstrate that integrase activity of the type I-D Cas4-Cas1 complex is potentially halted until integration of the non-PAM overhang has occurred. We propose a model for the type I-D system that relies on a delayed PAM-site integration by Cas4 in order to result in a correctly oriented spacer. In summary, we propose a mechanism in which two functionally independent complexes, Cas4-Cas1 and Cas1-Cas2, sequentially process and integrate prespacer substrates. This mechanism ensures correct spacer orientation as well as correct PAM-processing, thereby resulting in interference-proficient CRISPR adaptation.

Acknowledgements

We would like to thank Marre Niessen for early contributions. We thank Dr. Viktorija Globyte for critical reading of the manuscript. We thank Rob B.M. Koehorst for providing the Synechocystis PCC6803 strain.

Author contributions: S.N.K, C.A. and S.J.J.B. designed the experiments. S.N.K, 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 analyzed the data. S.N.K, C.A. and S.J.J.B. wrote the paper with input from all authors.

2 3	526	
5 6	527	FUNDING: FOM (Projectruimte 15PR3188-2); Netherlands Organisation for Scientific
7 8	528	Research [VICI VI.C.182.027]
9 10 11	529	
12 13	530	Conflict of interest statement. None declared.
14 15 16	531	
10 17 18	532	
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		

535 References

Arslan, Z., Hermanns, V., Wurm, R., Wagner, R., and Pul, Ü. (2014). Detection and characterization of spacer integration intermediates in type I-E CRISPR-Cas system. Nucleic Acids Research 42, 7884-7893. Barrangou, R. (2013). CRISPR-Cas systems and RNA-guided interference. WIREs RNA 4, 267-278. Brouns, S.J.J., Jore, M.M., Lundgren, M., Westra, E.R., Slijkhuis, R.J.H., Snijders, A.P.L., 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. Fagerlund, R.D., Wilkinson, M.E., Klykov, O., Barendregt, A., Pearce, F.G., Kieper, S.N., Maxwell, H.W.R., Capolupo, A., Heck, A.J.R., Krause, K.L., et al. (2017). Spacer capture and integration by a type I-F Cas1-Cas2-3 CRISPR adaptation complex. Proc Natl Acad Sci U S A 114, E5122-E5128. Gonzales, M.F., Brooks, T., Pukatzki, S.U., and Provenzano, D. (2013). Rapid Protocol for Preparation of Electrocompetent Escherichia coli and Vibrio cholerae. Journal of Visualized Experiments : JoVE, 50684. Hou, Z., and Zhang, Y. (2018). Insights into a Mysterious CRISPR Adaptation Factor, Cas4. Mol Cell 70, 757-758. Hudaiberdiev, S., Shmakov, S., Wolf, Y.I., Terns, M.P., Makarova, K.S., and Koonin, E.V. (2017). Phylogenomics of Cas4 family nucleases. BMC Evolutionary Biology 17, 232-232. Jackson, S.A., Birkholz, N., Malone, L.M., and Fineran, P.C. (2019). Imprecise Spacer Acquisition Generates CRISPR-Cas Immune Diversity through Primed Adaptation. Cell Host and Microbe

³⁸ 563 25, 250-260.e254.
 ⁴⁰ 564
 ⁴¹ 565 Jackson, S.A., McKenzie, R.E., Fagerlund, R.D., Kieper, S.N., Fineran, P.C., and Brouns, S.J.

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

Jansen, R., Embden, J.D.A.V., Gaastra, W., and Schouls, L.M. (2002). Identification of genes
that are associated with DNA repeats in prokaryotes. Molecular Microbiology *43*, 1565-1575.
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.

574 Kieper, S.N., Almendros, C., Behler, J., McKenzie, R.E., Nobrega, F.L., Haagsma, A.C., Vink,
 575 J.N.A., Hess, W.R., and Brouns, S.J.J. (2018). Cas4 Facilitates PAM-Compatible Spacer Selection
 576 during CRISPR Adaptation. Cell Reports 22, 3377-3384.

57 578 Kim, S., Loeff, L., Colombo, S., Jergic, S., Brouns, S.J.J., and Joo, C. (2020). Selective loading and
 57 579 processing of prespacers for precise CRISPR adaptation. Nature.

⁵⁹ 580 ₆₀

1				
2				
5 4	581	Lee, H., Dhingra, Y., and Sashital, D.G. (2019). The Cas4-Cas1-Cas2 complex mediates precise		
5	582	prespacer processing during CRISPR adaptation. eLife 8, 1-84.		
6	583			
7	584	Lee, H., Zhou, Y., Taylor, D.W., and Sashital, D.G. (2018). Cas4-Dependent Prespacer Processing		
8	585	Ensures High-Fidelity Programming of CRISPR Arrays. Molecular Cell 70, 48-59.e45.		
9 10	586			
11	587	Lemak, S., Beloglazova, N., Nocek, B., Skarina, T., Flick, R., Brown, G., Popovic, A., Joachimiak,		
12	588	A., Savchenko, A., and Yakunin, A.F. (2013). Toroidal Structure and DNA Cleavage by the		
13	589	CRISPR-Associated [4Fe-4S] Cluster Containing Cas4 Nuclease SSO0001 from Sulfolobus		
14	590	solfataricus. Journal of the American Chemical Society 135, 17476-17487.		
15 16	591			
17	592	Lemak, S., Nocek, B., Beloglazova, N., Skarina, T., Flick, R., Brown, G., Joachimiak, A.,		
18	593	Savchenko, A., and Yakunin, A.F. (2014). The CRISPR-associated Cas4 protein Pcal_0546 from		
19	594	Pyrobaculum calidifontis contains a [2Fe-2S] cluster: crystal structure and nuclease activity.		
20	595	Nucleic Acids Research 42, 11144-11155.		
21 22	596			
23	597	Li, M., Wang, R., Zhao, D., and Xiang, H. (2014). Adaptation of the Haloarcula hispanica CRISPR-		
24	598	Cas system to a purified virus strictly requires a priming process. Nucleic Acids Research 42,		
25	599	2483-2492.		
26 27	600			
27 28	601	Liu, T., Liu, Z., Ye, Q., Pan, S., Wang, X., Li, Y., Peng, W., Liang, Y., She, Q., and Peng, N. (2017).		
29	602	Coupling transcriptional activation of CRISPR–Cas system and DNA repair genes by Csa3a in		
30	603	Sulfolobus islandicus. Nucleic Acids Research 45, 8978-8992.		
31	604			
32	605	McGinn, J., and Marraffini, L.A. (2019). Molecular mechanisms of CRISPR–Cas spacer		
33 34	606	acquisition. Nature Reviews Microbiology 17, 7-12.		
35	607			
36	608	Nuñez, J.K., Harrington, L.B., Kranzusch, P.J., Engelman, A.N., and Doudna, J.A. (2015a).		
37	609	Foreign DNA capture during CRISPR–Cas adaptive immunity. Nature 527, 535-538.		
38	610			
39 40	611	Nuñez, J.K., Kranzusch, P.J., Noeske, J., Wright, A.V., Davies, C.W., and Doudna, J.A. (2014).		
41	612	Cas1-Cas2 complex formation mediates spacer acquisition during CRISPR-Cas adaptive		
42	613	immunity. Nature Structural & Molecular Biology 21, 528-534.		
43	614			
44 45	615	Nuñez, J.K., Lee, A.S.Y., Engelman, A., and Doudna, J.A. (2015b), Integrase-mediated spacer		
45 46	616	acquisition during CRISPR-Cas adaptive immunity. Nature 519, 193-198.		
47	617			
48	618	Nussenzweig, P.M., and Marraffini, I.A. (2020). Molecular Mechanisms of CRISPR-Cas		
49	619	Immunity in Bacteria, Annu Rev Genet 54, 93-120.		
50 51	620			
52	621	Plagens A Tiaden B Hagemann A Randau I and Hensel R (2012) Characterization of		
53	622	the CRISPR/Cas Subtype I-A System of the Hyperthermonbilic Crenarchaeon Thermonroteus		
54	623	tenax Journal of Bacteriology 194, 2491-2500		
55	624			
50 57	625	Ramachandran A Summerville Learn R A DeRell and Railey S (2020) Processing		
58	626	and integration of functionally oriented prespacers in the Escherichia coli CRISPR system		
59	627	depends on hacterial host exonucleases. I Riol Chem 295, 3403-3414		
60	027			

Rollie, C., Schneider, S., Brinkmann, A.S., Bolt, E.L., and White, M.F. (2015). Intrinsic sequence specificity of the Cas1 integrase directs new spacer acquisition. eLife 4, e08716. Shiimori, M., Garrett, S.C., Graveley, B.R., and Terns, M.P. (2018). Cas4 Nucleases Define the PAM, Length, and Orientation of DNA Fragments Integrated at CRISPR Loci. Molecular cell 70, 814-824.e816. Tahallah, N., Pinkse, M., Maier, C.S., and Heck, A.J. (2001). The effect of the source pressure on the abundance of ions of noncovalent protein assemblies in an electrospray ionization orthogonal time-of-flight instrument. Rapid Commun Mass Spectrom 15, 596-601. van den Heuvel, R.H., van Duijn, E., Mazon, H., Synowsky, S.A., Lorenzen, K., Versluis, C., Brouns, S.J., Langridge, D., van der Oost, J., Hoyes, J., et al. (2006). Improving the performance of a quadrupole time-of-flight instrument for macromolecular mass spectrometry. Anal Chem 78, 7473-7483. van der Oost, J., Westra, E.R., Jackson, R.N., and Wiedenheft, B. (2014). Unravelling the structural and mechanistic basis of CRISPR–Cas systems. Nature Reviews Microbiology 12, 479-492. Yosef, I., Goren, M.G., and Qimron, U. (2012). Proteins and DNA elements essential for the CRISPR adaptation process in Escherichia coli. Nucleic Acids Research 40, 5569-5576. Zhang, J., Kasciukovic, T., and White, M.F. (2012). The CRISPR Associated Protein Cas4 Is a 5' to 3' DNA Exonuclease with an Iron-Sulfur Cluster. PLoS ONE 7, e47232-e47232. Zhang, Z., Pan, S., Liu, T., Li, Y., and Peng, N. (2019). Cas4 Nucleases Can Effect Specific Integration of CRISPR Spacers. Journal of Bacteriology 201.

Page 27 of 39



Figure 1



Figure 2





Figure 4







Figure S1







Figure S3

1 Figure Legends

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

A Genetic organization of the type I-D CRISPR-locus. Genes constituting the interference
machinery are located upstream of the adaptation complex. The adaptation complex consisting
of *cas4*, *cas1* and *cas2* is highlighted in purple. Downstream of *cas2* is the leader sequence
followed by the type I-D array.

8 B Experimental design of spacer processing in vivo assay. Idealized pre-spacer substrates were
9 electroporated into *E. coli* cells carrying the plasmid encoded minimalized type I-D array and
10 expressing the adaptation genes *cas1* and *cas2*. The Cas4 protein was either omitted or co11 expressed as the wild-type or D76A+K91A mutant protein. Half-site integration was assessed by
12 PCR as depicted in C.

13 C PCR scheme allowing differentiation of spacer orientation, integration site and prespacer
 14 processing. PCR amplicons of correct (PAM overhang integrated at Repeat-Spacer junction) and
 15 incorrect spacer orientations and Leader-Repeat (L-R) or Repeat-Spacer (R-S) integration site
 16 were subjected to high-throughput sequencing.

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

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

A N-terminally tagged *Cas1 associates with untagged Cas2 in the absence of Cas4. Unnumbered
 peak contains protein aggregates.

B Complex formation of N-terminally tagged *Cas1 and untagged Cas4.

C Co-expression of *Cas1, Cas2 and Cas4. *Cas1 elutes separately (peak 1) from the Cas4-Cas1
 complex (peak 2). Cas2 together with dissociated Cas1 and Cas4 elutes as a low molecular weight
 peak (peak 3). Unnumbered peak contains protein aggregates.

D Native Mass Spectrometry of Cas4-Cas1 complex as shown in B (native spectrum obtained after removal of His-SUMO tag from Cas1 by TEV protease cleavage (Fig. S2)). Cas4 monomers assemble with Cas1 dimers into a Cas4₁-Cas1₂ complex. Cas1 dimers are also frequently observed. Free monomers of Cas1 and Cas4 are less frequent in the measured sample.

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

A Pre-spacer model substrate containing a 25 bp duplex region flanked by 13 nt 3' overhangs incubated with different protein combinations. The non-PAM strand is 5' Cy3 labelled and the PAM-containing strand 5' Cy5 labelled. Cleavage of PAM containing 3' overhang results in Cy5 labelled fragment of 30 nt. Protein samples consist of co-purified Cas4-Cas1 and Cas4D76A-Cas1 complexes, combined Cas4-Cas1 and Cas1-Cas2 complex and individually purified Cas4 in the presence of Cas1-Cas2.

B Integration of labelled pre-spacer (PS) into linear CRISPR DNA consisting of the type I-D leader
 sequence (L) and a single Repeat (R). Labelled spacer imaged via Cy5, total DNA via SYBR gold
 stain. Merge of Cy5 and SYBR gold channels indicates integration of Cy5 labelled pre-spacer
 resulting in a higher molecular weight band.

Page 37 of 39

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

D-E High-throughput sequencing of Cas4-Cas1 integrated pre-spacers. Reaction was performed similar to the assay shown in **C** using unlabeled pre-spacer DNA. OC plasmids were gel extracted followed by PCRs specific for the leader-repeat (L-R) and repeat-spacer (R-S) integration as well as correct and incorrect spacer orientation. Bar graphs represent the percentage of spacers with specific overhang length depending on integration site and orientation (D-correct; E-incorrect). n = number of sequenced integration events.

58 Figure 4 - Spacer overhang preferences of Cas1-Cas2 and Cas4-Cas1 complexes.

A - Integration activity with respect to 3' overhang requirements. Pre-spacer substrates were 5'
Cy5 labelled in order to follow coupling to pCRISPR. Phosphorylated (P) 3' overhangs were used
to block integration of one of the DNA strands.

Figure 5 – Model of Cas4-Cas1 and Cas1-Cas2 assisted spacer selection, processing and
integration. Prespacers with long PAM- and non-PAM containing 3' overhangs are bound by
Cas4-Cas1 (PAM overhang) and Cas1-Cas2 (non-PAM overhang). Following processing by hostfactor nucleases, the non-PAM site of the spacer is integrated at the leader-repeat site (first halfsite integration). Subsequently, the second half-site integration (spacer-site integration) of the

1		
2 3 4	68	PAM-s
5 6	69	site of
7 8 9	70	full-sit
10 11	71	
12 13 14	72	
15 16	73	
17 18	74	
19 20 21	75	Supp
22 23	76	Figur
24 25	77	Cas1
26 27 28	78	PAGE
28 29 30	79	
31 32	80	Figur
33 34	81	after ⁻
35 36 37	82	Cas4
38 39	83	of SE
40 41	84	Figur
42 43 44	85	show
45 46	86	
47 48		
49 50		
51 52		
52 53		
54 55		
56		
57		
50 59		
60		

68	PAM-site occurs, likely orchestrated by release of the Cas4-processed overhang into the integrase
69	site of Cas1. Unwinding of the repeat followed by gap-repair completes repeat duplication and
70	full-site spacer integration.
71	
72	
73	
74	
75	Supplementary Figures Legends
76	Figure S1 – Related to Fig. 2B – Co-purification of His ₆ -SUMO-Cas4 and untagged
77	Cas1. A Size exclusion chromatogram of His_6 -SUMO-Cas4 and untagged Cas1. B SDS
78	PAGE analysis of SEC purified proteins.
79	
80	Figure S2 – Related to Fig. 2D – Co-purification of untagged Cas4 and untagged Cas1
81	after TEV-protease cleavage of His_6 -SUMO-TEV tag. A Size exclusion chromatogram of
82	Cas4 and Cas1 (A280 – total protein; A400 – Cas4 FeS-cluster). B SDS PAGE analysis
83	of SEC purified proteins.
84	Figure S3 – Related to Fig 2A – Native Mass Spectrometry of Cas1-Cas2 complex as
85	shown in Fig. 2A.
86	

Name in this study	Name	Insert	Vector	Resistance	Source
pCas2	pTU084	Synechocystis PCC6803 Type I-D cas2 (deltaCas1)	pET-T7	Amp	Kieper et al. (2018)
pCas1	pTU085	Synechocystis PCC6803 Type I-D <i>cas1</i> (deltaCas2)	pET-T7	Amp	Kieper et al. (2018)
pCas1	pTU092	Synechocystis PCC6803 Type I-D cas1	pET-T7	Spec	This study
pCas4 ^{D76A}	pTU086	Synechocystis PCC6803 Type I-D cas4 (D76A)	pET-T7	Spec	Kieper et al. (2018)
pCas4 ^{D76A+K91A}	pTU411	Synechocystis PCC6803 Type I-D cas4(D76A+K91A)	pET-T7	Spec	This study
pCas4	pTU130	Synechocystis PCC6803 Type I-D cas4	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 cas1-cas2	pET-T7	Amp	Kieper et al. (2018)
pEmp	pTU116	NA	pET-T7	Spec	Addgene Plasmid #48329

Table S1 - Plasmids used in this study

Table S2 - Oligonucleotides used in this study

Name	Sequence	Description
BN015	CGTCCATGGGAAGTCATTCTTCAAATTTTGGC	Leader Fw
BN277	GTGGAATACGCAAAAGGC	cas4 mutagenesis K91A Fw
BN278	AGGAATTAATAAGCCATCACTTTC	cas4 mutagenesis K91A Rv
BN1398	GCTAGTTATTGCTCAGCGG	pCRISPR bb Rv
BN1711	GGAAGGTTTGCCAAAGTC	Leader Distal Half-Site Integration
BN1712	CTGTTCGACTTAAGCATTATGC	Leader Proximal Half-Site Integration
BN1713	ATCGACACCACCA	OligoSpecific Primer Fw (PAM overhang)
BN1714	CGTGGTGGTGTCGAT	OligoSpecific Primer Rv (non-PAM overhang)
BN1763	CTACCATCGACACCACCACGCTGGCTTTTTAACTTTT	25 nt duplex 13nt PAM 3' ovhng
BN1768	GCCAGCGTGGTGGTGTCGATGGTAGTTTTTTGTTTT	25 nt duplex 13nt RvC PAM 3' ovhng
BN1829	CTACCATCGACACCACCACGCTGGCTTTTTTGTTTTT	25 nt duplex 13nt RvC PAM 3' ovhng (5' C6-Amino)
BN1830	GCCAGCGTGGTGGTGTCGATGGTAGTTTTTAACTTTT	25 nt duplex 13nt PAM 3' ovhng (5' C6-Amino)