1
2

# Structure-guided engineering of type I-F CASTs for targeted gene insertion in human cells

5

6

7   George D. Lampe[1]*, Ashley R. Liang[1,5]*, Dennis J. Zhang[1,6], Israel S. Fernández[2,3,#], Samuel H.
8   Sternberg[1,4#]

9
10
11
12

13   [1]Department of Biochemistry and Molecular Biophysics, Columbia University, New York, NY,
14   USA.
15   [2]Ikerbasque, Basque Foundation for Science, Bilbao, Spain.
16   [3]Instituto Biofisika (UPV/EHU, CSIC), University of the Basque Country, Leioa, Spain.
17   [4]Howard Hughes Medical Institute, Columbia University, New York, NY, USA.
18   [5]Present Address: Tornado Bio, San Francisco, CA, USA.
19   [6]Present Address: Section of Microbiology, Department of Biology, University of Copenhagen,
20   Copenhagen, Denmark.
21
22   *These authors contributed equally to this work.
23   [#]Co-corresponding authors. E-mail: shsternberg@gmail.com; israel.s.fernandez@gmail.com
24

## ABSTRACT

Conventional genome editing tools rely on DNA double-strand breaks (DSBs) and host recombination proteins to achieve large insertions, resulting in a heterogeneous mixture of undesirable editing outcomes. We recently leveraged a type I-F CRISPR-associated transposase (CAST) from the *Pseudoalteromonas* Tn*7016* transposon (*Pse*CAST) for DSB-free, RNA-guided DNA integration in human cells, taking advantage of its programmability and large payload capacity. *Pse*CAST is the only characterized CAST system that has achieved human genomic DNA insertions, but multiple lines of evidence suggest that DNA binding may be a critical bottleneck that limits high-efficiency activity. Here we report structural determinants of target DNA recognition by the *Pse*CAST QCascade complex using single-particle cryogenic electron microscopy (cryoEM), which revealed novel subtype-specific interactions and RNA-DNA heteroduplex features. By combining our structural data with target DNA library screens and rationally engineered protein mutations, we uncovered CAST variants that exhibit increased integration efficiency and modified PAM stringency. Structure predictions of key interfaces in the transpososome holoenzyme also revealed opportunities for the design of hybrid CASTs, which we leveraged to build chimeric systems that combine high-activity DNA binding and DNA integration modules. Collectively, our work provides unique structural insights into type I-F CAST systems while showcasing multiple diverse strategies to investigate and engineer new RNA-guided transposase architectures for human genome editing applications.

## INTRODUCTION

Canonical CRISPR-Cas systems that have been leveraged for programmable gene editing, such as Cas9 nucleases, cause targeted DNA double-strand breaks (DSBs) that provoke the cell to activate DNA repair mechanisms[1,2]. Non-homologous end joining (NHEJ) is the most efficient repair pathway in human cells, which leads to indel mutations, and although homology-directed repair (HDR) offers the ability to generate precise modifications or insertions, it is inefficient in most cell types, inaccessible in non-dividing cells, and requires large homology arms for each new insertion site[3,4]. Furthermore, HDR efficiencies decrease drastically with insertion size, and aberrant editing pathways that occur at non-negligible frequencies can cause large chromosomal truncations and/or rearrangements[5–10]. Second generation editors, including base and prime editors, employ nickase-variant Cas proteins to bypass DSB intermediates, but indel byproducts still arise and edits are generally restricted to single-base pair (bp) changes or small insertions (<50 bp)[11–14], thus failing to address the need for large DNA insertion technology. CRISPR-associated transposases (CASTs), on the other hand, leverage a CRISPR-associated DNA targeting module and a transposase effector module that allow for highly specific and programmable insertions, which are both DSB-free and multi-kilobases in size[15–17].

To date, four CAST subtypes have been characterized in bacteria: type I-B, I-D, I-F, and V-K[15,16,18,19]. These subtypes encode unique architectures for both the targeting and integration steps of the transposition pathway: type I CASTs rely on TnsABC proteins for integration and a multi-subunit complex for DNA targeting that includes TniQ and Cascade components (TniQ-

65  Cascade, hereafter simply QCascade), with Cascade itself comprising 3-5 unique protein
66  components in varying oligomeric states[20–22]; whereas type V-K CASTs rely on only TnsBC for
67  integration[16,23,24] and a simpler Cas12k-TniQ-S15 co-complex for DNA targeting[25]. Individual
68  homologs within each of these CAST subtypes also vary in sequence identity[26,27], subunit
69  composition and fusion connectivity[18,24,28], DNA targeting modules, crRNA guide sequence[18,26,29],
70  and host factor requirements[17,25,30], thus representing a diverse pool of potential starting points
71  for tool development. Although type V-K CASTs are more compact systems in terms of coding
72  size, they exhibit multiple undesirable biochemical properties — including reduced specificity[31–
73  33], low overall editing efficiencies[16,31], and poor product purity[24,34,35] — that would necessitate
74  extensive optimization for potential research and therapeutic genome engineering applications.
75  In contrast, type I-F CASTs exhibit highly specific and homogeneous integration products, with
76  demonstrably greater efficiencies than types I-B, I-D, and V-K[15–19,24].

77      CAST systems have been the focus of extensive structural efforts using cryogenic electron
78  microscopy (cryoEM) in recent years. The type V-K ShCAST system from *Scytonema hoffmannii*
79  has been systematically investigated[25,36–39], with a recent report of the holo transpososome
80  architecture that revealed intricacies of the megadalton complex containing Cas12k, TniQ, TnsB,
81  TnsC, single-guide RNA, partial donor and target DNA substrates, and the bacterial host factor
82  S15[39]. Structural studies of type I-B and I-F CASTs have largely focused on the QCascade DNA
83  targeting module and the accessory TnsC ATPase[20,21,40–43], with no structures of the
84  endonuclease-transposase TnsAB module described to date. Intriguingly, QCascade structures
85  exhibit distinct conformations across different systems: type I-B CASTs feature a single TniQ
86  monomer that recruits TnsC to the Cascade-bound target DNA[21], whereas type I-F CASTs feature
87  a TniQ homodimer that is stably associated with Cascade[20]. Thus far, two I-F CAST systems from
88  subtypes I-F3a and I-F3b have been deeply characterized — *Vch*CAST (Tn*6677*) and *Asa*CAST
89  (Tn*6900*), respectively —  both of which are only distantly related to *Pse*CAST (Tn*7016*), a system
90  that we recently exploited for targeted DNA integration in human cells[17].

91      The *Pse*CAST RNA-guided transposase was identified as a lead candidate for human
92  genome engineering applications through a systematic screen of diverse type I-F CAST systems[17]
93  (**Fig. 1a**). Although our first study reported editing activities that reached single-digit efficiencies
94  at genomic target sites, representing a ~100-fold improvement over our original candidate,
95  *Vch*CAST, these efficiencies remain limiting for downstream applications. We hypothesized that
96  identifying bottlenecks in the system would inform more targeted rationally engineering,
97  developed several assays to investigate intermediate events and overall integration efficiencies
98  in human cells[17], and then applied these assays to *Vch*CAST and *Pse*CAST, the only type I-F
99  CASTs shown to successfully perform RNA-guided integration in human cells. Intriguingly, while
100 *Pse*CAST promoted comparatively robust DNA integration, it exhibited markedly weaker DNA
101 binding activity relative to *Vch*CAST. We therefore hypothesized that, alongside parallel efforts to
102 engineer and evolve hyperactive transposase variants, the *Pse*CAST QCascade module would
103 represent a promising focus area to improve DNA targeting and thus editing efficiencies.

104     Towards that goal, here we report the cryoEM structure of *Pse*CAST QCascade and the
105 effect of targeted mutations in the PAM- and crRNA-interacting regions on DNA integration.

106    Separately, we leveraged AlphaFold-Multimer to predict protein-protein interactions within the
107    TnsABC co-complex, inspiring the rational design of novel chimeric CAST systems that enable
108    divergent DNA targeting and DNA integration modules to be combined into a single functional
109    system. Collectively, this work establishes multiple biochemically- and structurally-guided
110    approaches to engineer CAST systems for improved editing efficiencies in human cells.

111

## RESULTS

### CryoEM structure of *Pse*CAST QCascade complex

114    We previously demonstrated that *Vch*CAST and *Pse*CAST, two distantly related type I-F
115    CASTs[17,26], exhibit distinct DNA binding and integration efficiencies (**Fig. 1a-c**). Given our
116    previous mechanistic and structural studies of the QCascade complex from *Vch*CAST[20,40], we
117    hypothesized that structure-guided engineering of the *Pse*CAST QCascade complex might reveal
118    novel interactions and open a path to improve overall integration efficiencies. We therefore
119    purified recombinant *Pse*QCascade after carefully optimizing the expression vector design
120    (**Supplementary Fig. 1**) and set out to determine the cryoEM structure.

121    We incubated the purified *Pse*QCascade complex, which is expected to comprise a
122    1:6:1:2:1 stoichiometry of Cas8:Cas7:Cas6:TniQ:crRNA components (**Fig. 1d**), with a double-
123    stranded DNA (dsDNA) substrate containing a 32-bp target sequence and 5′-CC-3′ PAM, and
124    then subjected the sample to electron microscopy. Preliminary cryoEM experiments revealed a
125    homogeneous behavior with multiple views and no apparent disassembly (**Supplementary Fig.
126    2a**), and the overall architecture was consistent with other type I-F QCascade complexes,
127    comprising six Cas7 monomers (named hereafter Cas7.1 to Cas7.6) that form a pseudo-helical
128    assembly coating the crRNA molecule (**Fig. 1e**). The Cas8 protein contains two domains: a bulky
129    domain that interacts with Cas7.1 and binds the crRNA 5′ end and PAM sequence, and a second
130    α-helical domain that exhibited a dynamic behavior (**Fig. 1f**). Towards the crRNA 3′ end (hereafter
131    PAM-distal region), the RNA hairpin is stabilized by Cas6, which also binds the TniQ dimer.
132    Preliminary maps exhibited greater mobility for the TniQ dimer compared to other QCascade
133    components (**Supplementary Fig. 2b,c**). The quality of the maps approaching the TniQ dimer
134    region degrades rapidly, contrasting the excellent map quality for the PAM-adjacent region
135    (**Supplementary Fig. 2d**). Multibody approaches in Relion4 improved the overall resolution, with
136    approximately 2.6 Å and 3.0 Å resolution estimates in the PAM-proximal and PAM-distal regions,
137    respectively (**Methods**).

138    To further characterize the dynamics of the system and confirm the existence of novel
139    interactions, we complemented our multibody analysis in Relion4 with cryoDRGN[44], a machine-
140    learning approach for cryoEM analysis (**Supplementary Fig. 3**). CryoDRGN revealed multiple
141    populations of the complex, with the TniQ dimer populating a wide range of positions relative to
142    the rest of the complex that pivot around Cas6 and Cas7.6. The dimer adopts an 'open'
143    conformation that lacks any direct interactions with Cas8, as well as multiple intermediate, 'closed'
144    conformations that approach the tip of the Cas8 α-helical domain (**Supplementary Fig. 3b**). In a
145    recent structure of a homologous QCascade complex bound to target DNA, the Cas8 α-helical

146    domain exhibits a different conformation, almost perpendicular to the inner face of the TniQ dimer
147    and aligned with the bulky domain of Cas8[22]. In our dataset, we did not observe this extended
148    conformation and instead detected alternative TniQ-Cas8 interactions that are established
149    between the most distal end of the TniQ dimer and the apical part of the Cas8 α-helical domain,
150    which were revealed through low-pass filtered maps (**Supplementary Fig. 3c**). Both the TniQ
151    dimer and the Cas8 α-helical domains remain in parallel configurations, with only marginal
152    contacts at the periphery of the complex. Despite the apparent flexibility in this interaction
153    (**Supplementary movie 1 and 2**), the Cas8 α-helical domain is essential for RNA-guided DNA
154    integration activity, as revealed by the complete loss of human cell activity when we replaced the
155    domain with a flexible glycine-serine linker (**Supplementary Fig. 4**).
156
157    **Stabilizing protein-RNA and protein-protein interactions**
158    The overall architecture of the TniQ dimer is similar to the *Vch*CAST QCascade dimer[20],
159    with an antiparallel head-to-tail configuration, forming a compact unit that laterally approaches the
160    interface formed by Cas6 and Cas7.6 (**Fig. 2a**). The C-terminal domain of one TniQ monomer
161    interacts with Cas6, and the N-terminal domain of the other TniQ monomer interacts with Cas7.6.
162    At the core of this four-fold interface, the crRNA appears to play a critical role, with residues 40–
163    45 establishing multiple RNA-protein stacking interactions (**Fig. 2b,c**).
164    We hypothesized that crRNA interactions with Cas6, Cas7.1, TniQ.1, and TniQ.2 are
165    crucial for robust QCascade complex formation, and that disrupting them would prevent
166    transposase recruitment and abolish integration activity. We therefore introduced alanine point
167    mutations to disrupt nucleobase-side chain stacking interactions and investigated the resulting
168    effects in human genomic DNA integration assays. Alanine substitutions to Cas6 and TniQ
169    residues contacting the crRNA were well tolerated, whereas a Cas7 R143A mutation (Cas7$^{R143A}$)
170    abolished integration activity (**Fig. 2d**). The crRNA trajectory in the hinge region between Cas7.6
171    and Cas6 differs in *Pse*CAST and *Vch*CAST (**Fig. 2e**), and *Pse*CAST crRNA residue G41 seems
172    to play a key role as an interaction "hub," establishing coincident contact with TniQ.1, TniQ.2, and
173    Cas7.6 by adopting a unique, extruded conformation.
174    We next explored protein-protein interactions that we similarly hypothesized would
175    contribute to QCascade function, in part by playing a role in downstream transposase recruitment
176    to the target site. The first of these interactions involved a hydrophobic patch on Cas6 cradling
177    hydrophobic residues in the loop connecting TniQ.1 α-helices W262–K275 and F312–S327 (**Fig.
178    3a,b**), which is conserved across homologous QCascade complexes, with minor variations.
179    Specifically, a hydrophobic residue in the TniQ.1 connecting loop (I282 in *Pse*CAST, V270 in
180    *Vch*CAST) inserts deeply into the Cas6 hydrophobic patch to anchor the TniQ monomer to the
181    Cascade module (**Fig. 3c**). The cradle structure of this interaction potentially acts as a pivot point,
182    facilitating dynamic TniQ movement. Disruption of these hydrophobic interactions via introduction
183    of charged arginine residues in either TniQ or Cas6 led to a marked reduction in integration
184    efficiencies (**Fig. 3d**). The other TniQ monomer (TniQ.2) interacts electrostatically with Cas7.6 via
185    α-helix Y33–L47 and adjacent residues (**Fig. 3e**). Given the multimeric assembly of Cas7
186    monomers along the crRNA, loop regions observed to interact with TniQ.2 may have pleiotropic

187    functions, possibly participating in Cas7 monomer-monomer interactions (**Supplementary Fig.**
188    **5**). With the goal of selectively perturbing Cas7.6-TniQ.2 interactions to investigate its importance,
189    we avoided mutagenizing residues that might affect the Cas7 monomer-monomer contacts and
190    thus focused on loops A and B (**Supplementary Fig. 5b**). Alanine mutations within the TniQ-
191    interacting regions abolished DNA integration, whereas several mutations within Cas7 had
192    surprisingly little to no impact on overall DNA integration activity (**Fig. 3f**).

193

194    **Protein engineering modulates PAM stringency and improves DNA integration**
195    In comparison to other type I-F CASTs, *Pse*CAST exhibits a remarkably flexible PAM
196    preference, with almost no sequence preference at both the -1 and -2 positions in *E. coli*
197    transposition assays[26]; this property may lead to a dramatic increase in the effective search space
198    for the 32-bp guide. Inspired by previous work investigating CRISPR-Cas9 activity and PAM
199    search space[45], we hypothesized that inefficient DNA targeting due to a flexible PAM preference
200    may represent a rate-limiting step in RNA-guided DNA integration, especially within the cellular
201    milieu of human cells, whose genome is ~1000× larger than *E. coli*. We therefore set out to
202    specifically engineer QCascade variants that might exhibit altered PAM specificity and thus direct
203    altered DNA integration efficiencies.
204    After leveraging the excellent quality of our cryoEM map in the area surrounding Cas8, we
205    identified two hydrophobic alanine residues at the center of the PAM-interacting region. In
206    contrast, systems with stricter PAM preferences — *Vch*CAST, *Asa*CAST, and *Pae*Cascade from
207    a *Pseudomonas aeruginosa* type I-F1 CRISPR-Cas system[26,46] — feature polar residues at the
208    equivalent positions, which allow for hydrogen bonding with specific PAM nucleotides (**Fig. 4a,b,**
209    **Supplementary Fig. 6a**). Based on these observations, we reasoned that mutation of A143 and
210    A144 to residues with greater hydrogen bonding potential might improve PAM stringency, reduce
211    the effective search space, and result in more efficient DNA targeting. We also chose to
212    mutagenize residues 125–127, as this region also interacts with the PAM (**Fig. 4b,**
213    **Supplementary Fig. 6a**). We analyzed the sequence conservation at these PAM-interacting
214    regions and compared *Pse*CAST to other Cascade homologs that have previously exhibited either
215    robust DNA integration activity or stringent PAM preferences (**Supplementary Fig. 6b,c**).
216    Collectively, we designed fifteen Cas8 variants with PAM-interacting mutations, varying from
217    single point mutations at A243 or A244 to larger mutations in which the entire PAM-interacting
218    region was grafted from a homolog.
219    We quantified changes in PAM preference by performing an episomal PAM library screen
220    in HEK293T cells, in which a target plasmid (pTarget) contained an *AAVS1* target site directly
221    downstream of a randomized 4-bp PAM library (**Supplementary Fig. 6d**). After transiently
222    transfecting cells with pTarget, pDonor, and all the necessary protein-RNA expression vectors,
223    we isolated plasmid DNA, sequenced the PAM motifs from all successful integration products,
224    and constructed a consensus motif for each Cas8 variant; in parallel, we also quantified absolute
225    integration efficiencies at the genomic AAVS1 site, which contains a 5′-CC-3′ PAM (**Fig. 4c**). The
226    results revealed that certain mutations led to improvements in integration efficiencies by as much
227    as 3.5-fold, but without a clear correlation between PAM stringency and overall genomic

228   integration activity (**Fig. 4c**). For example, the variant with the greatest improvement in integration
229   activity, Cas8$^{R241K,A244S}$, actually exhibited a reduced PAM preference, compared to the stronger
230   preference for cytidine in the -2 position with WT Cas8 (**Fig. 4c, Supplementary Fig. 6e**).
231   Interestingly, Cas8$^{A243Q,A244N}$ exhibited decreased PAM preference, whereas when we grafted the
232   entire PAM region from a type I-F1 system ($_{241}$RPAAV$_{245}$>KPQNI), the resulting mutant restored
233   a strong preference for cytidine at the -1. Mutations within the upstream PAM-interacting region
234   (residues 125-127) showed moderate improvements on integration activity, with either unchanged
235   or moderately reduced PAM stringency (**Fig. 4c**). A Cas8$^{R241A}$ mutant with disrupted 'R-wedge,'
236   which normally forms stacking interactions with the -1 PAM position to help unwind dsDNA[47,48],
237   unexpectedly exhibited both WT integration efficiencies and PAM stringency (**Fig. 4c**).
238         Together, mutational profiling of the PAM-interacting region revealed key residues whose
239   mutation improved integration efficiencies, but the combination of PAM specificity and integration
240   activity results failed to support the hypothesis that PAM promiscuity is a key bottleneck towards
241   achieving higher efficiency *Pse*CAST integration activity in human cells (**Fig. 4c, Supplementary**
242   **Fig. 6e**).
243         We also focused on PAM-proximal interactions with the upstream double-stranded DNA
244   region as another potential point of engineering and optimization. Previous work on canonical
245   type I-F1 defense systems revealed key interactions between dsDNA and the N-terminal region
246   of Cas8[47–49], with a positively charged vise domain undergoing a conformational change to 'clamp'
247   onto the PAM-adjacent sequence in a non-specific fashion. When comparing *Pse*Cas8 (from type
248   I-F3 *Pse*CAST) to *Pae*Cas8 (from type I-F1 *Pae*Cascade; **Supplementary Fig. 7a**), we observed
249   a markedly different conformation of the N-terminus, with the vise domain absent. Given this
250   potential deficiency, we hypothesized that substituting the *Pae*Cas8 vise domain in *Pse*Cas8
251   could improve DNA binding affinity and thus CAST activity. However, a thorough screening of
252   chimeric Cas8 constructs for human cell integration activity revealed a clear intolerance of
253   *Pse*Cas8 to sequence perturbations in this region (**Supplementary Fig. 7b**). We pursued
254   additional synthetic strategies to improve DNA binding of *Pse*QCascade by fusing a variety of
255   DNA-binding domains to the *Pse*Cas8 N-terminus of *Pse*Cas8 (**Supplementary Fig. 7c**), inspired
256   by engineering strategies previously applied to polymerases[50,51], reverse transcriptases[52], and
257   ligases[53]. However, these fusions exhibited no improvement relative to WT, and in some cases
258   reduced overall integration efficiencies (**Supplementary Fig. 7c**). Collectively, these experiments
259   suggest that either the DNA binding affinity of *Pse*Cas8 is not a critical bottleneck in the overall
260   transposition pathway, or that the tested variants fail to improve upon the WT activity in this
261   regard.
262
263   **Unfavorable nucleobase positioning along the RNA-DNA heteroduplex**
264         Cascade complexes bind the target DNA by forming a discontinuous RNA-DNA
265   heteroduplex in 6-bp segments[47,54], and we could clearly resolve RNA-DNA base pairs for the first
266   4 segments engaged by Cas7 monomers within the *Pse*QCascade complex, but the remaining
267   two segments featured weaker RNA density and no DNA density. Density for the RNA-DNA
268   heteroduplex across the first 3 segments (crRNA residues 9 to 26) was exceptionally good, with

269    clear separation within base pairs and features compatible with a local resolution beyond 3 Å. We
270    were therefore able to accurately model RNA-DNA interactions to a high level of confidence in
271    these regions of the map. The resulting view revealed peculiarities in the base-pair geometry, with
272    acute divergence from ideal values in some base pairs. The third and fourth base pair within each
273    segment exhibited severe deviation from ideal planarity values (buckling), while the first and fifth
274    base pair exhibited exacerbated propeller twist deviations. Only the second base pair across
275    distinct segments exhibited geometric and hydrogen-bonding distance values closer to
276    energetically favored conditions (**Fig. 4d-h**).

277    Type I-F Cascade complexes bind the target DNA, such that the two-stranded β-sheet
278    'finger' motif of each Cas7 monomer engages the crRNA to flip out every sixth nucleotide of the
279    32-nt spacer, thereby preventing RNA-DNA basepairing[20,47]. We hypothesized that finger motif
280    residues involved in this nucleotide dislocation might promote the consistent distortion of adjacent
281    base pairs, and to explore this effect, we introduced Cas7 mutations intended to relax this
282    distortion, hoping to promote energetically favorable hydrogen-bonding geometries and stabilize
283    the RNA-DNA heteroduplex. Taking advantage of the high local resolution around this region
284    (**Supplementary Fig. 8a,b**), we identified numerous bulky hydrophobic residues —including I69,
285    L70, and L224 — that were not highly conserved across nearby homologs (**Supplementary Fig.
286    8c**), and subjected them to site-directed mutagenesis.

287    After generating the desired Cas7 mutations, we performed genomic DNA integration
288    experiments in HEK293T cells at the AAVS1 locus (**Fig. 4i**). Intriguingly, the Cas7 heteroduplex-
289    interacting residues, though not highly conserved, appeared to have low tolerance for mutations.
290    While Cas7$^{L224F}$ and various valine mutations exhibited near-WT integration efficiencies, all other
291    mutations, including Cas7$^{I69P}$, resulted in detrimental impacts on DNA integration (**Fig. 4i**).
292    Intriguingly, L70H, which would theoretically recapitulate a stacking interaction observed in our
293    previous *Vch*CAST structure[20], completely ablated integration activity (**Fig. 4i**). Together, the
294    intolerance to perturbations in the Cas7 finger domain suggests these nucleobase kinking
295    interactions may in fact be necessary for proper successful DNA integration.

296

**Structure-based engineering of chimeric CAST systems.**

298    Rational engineering of *Pse*QCascade yielded only moderate improvements in integration
299    activity, suggesting a non-trivial path forward to overcome the apparently weak DNA binding
300    activity in human cells[17]. Although recent studies shed light on the kinetics of Cascade target
301    search and recognition[55,56], the intermediate steps of Cascade complex formation, TniQ-Cascade
302    association, and 3D-diffusion remain poorly understood, particularly in human cells. *Pse*CAST
303    was originally identified through a homolog screen that investigated both overall integration
304    activity and several subunit-specific properties: crRNA processing, TnsB-donor DNA interactions,
305    and QCascade and TnsC-mediated transcriptional activation[17]. Through this screening process,
306    *Vch*CAST (Tn*6677*) and *Pse*CAST (Tn*7016*) were the only two systems that yielded detectable
307    DNA integration in human cells, despite exhibiting distinct subunit-specific activities. Based on
308    these results, we hypothesized that natural CAST systems may be unlikely to possess optimal
309    human cell properties across all recombinant components, and we therefore set out to design

310    chimeric CAST systems that would enable 'crosstalk' between otherwise orthogonal components.
311    Our specific goal was to combine the most active DNA targeting and DNA integration machineries
312    derived from divergent CASTs (**Fig. 5a**).

313    To identify robust DNA targeting homologs, we tested DNA binding activity across 20 type
314    I-F CASTs via transcriptional repression in *E. coli*[40,57] (**Supplementary Fig. 9a**). Surprisingly,
315    QCascade complexes from only two systems — *Vch*CAST and Tn*7005* — exhibited RFP
316    repression under the tested conditions, with only weak activity from *Pse*CAST and Tn*7000*
317    (**Supplementary Fig. 9b**). Yet when we tested the overall DNA integration activity of *Vch*CAST
318    and *Pse*CAST at the exact same sites used for transcriptional repression, we again observed
319    greater integration activity for *Pse*CAST, mirroring our results in human cells[17] (**Supplementary**
320    **Fig. 9c**). This reinforced the conclusion that the weak DNA targeting activity of *Pse*CAST may
321    impose a lower ceiling on achievable DNA integration efficiencies in diverse cell types, despite
322    having co-evolved with a highly active transposition (TnsABC) module.

323    We sought to address this potential bottleneck by combining the TnsABC machinery from
324    *Pse*CAST with the QCascade machinery from *Vch*CAST. We previously demonstrated that
325    intrinsic CAST modularity precludes simply mixing and matching components from evolutionary
326    diverse systems[26], but we were emboldened to attempt a more nuanced approach by taking
327    advantage of recent high-resolution structures[21,39], predicted structures via structural
328    alignments[58], and AlphaFold-multimer[59] predicted structures. (**Fig. 5b, Supplementary Fig. 10**).
329    In particular, a model for the putative TnsABC co-complex from *Pse*CAST featured the expected
330    heptameric arrangement of TnsC, similar to our empirical structures for *Vch*CAST[40], while also
331    revealing predicted interactions between *Pse*TnsC and the C-terminus of *Pse*TnsB that were
332    reminiscent of the TnsB 'hook' described for type V-K ShCAST[37,39] (**Fig. 5b, Supplementary Fig.**
333    **10a**). This model, in conjunction with experimentally determined type V-K structures and
334    biochemical studies of Tn*7*[60], led us to speculate that the C-terminal tail of TnsB functions a key
335    mediator of TnsC interactions, and that the specificity of CAST transpososome assembly would
336    be dictated in part by cognate TnsB-TnsC interactions. Importantly, we hypothesized that
337    reengineering this interaction would enable the TnsAB and donor DNA components from one
338    CAST system to be combined with the QCascade and TnsC components from an orthogonal
339    CAST system.

340    To test this hypothesis, we designed 16 chimeric TnsAB constructs in which different
341    lengths of the *Pse*TnsB C-terminus were substituted with corresponding residues from the
342    *Vch*TnsB C-terminus (**Fig. 5c**). These variants were then screened for RNA-guided DNA
343    integration activity in *E. coli*, in conjunction with *Vch*QCascade and *Vch*TnsC, but with a pDonor
344    containing transposon ends compatible with *Pse*TnsB (**Fig. 5d**). As expected, given our previous
345    work[26], WT *Pse*TnsAB, lacking any chimeric substitutions, showed undetectable activity when
346    combined with *Vch*CAST DNA targeting machinery (**Fig. 5e**). Remarkably, however, several
347    chimeric TnsAB designs were able to robustly rescue activity, showing up to ~10% integration
348    efficiencies (**Fig. 5e**). These designs, which only reprogrammed 20 – 29 amino acids in the C-
349    terminus of *Pse*TnsAB exhibited graft points between the *Pse* and *Vch*TnsB sequence in an
350    unstructured region that links the "hook" region of the C-terminus to the remainder of the protein

351   sequence (**Fig. 5c**); furthermore, when comparing this region to solved type V-K complexes, it is
352   located in a similar region as the 52-residue long "flexible linker" that was unresolved[39]. Together,
353   we conclude that substitutions in this region minimize disruptions to the overall protein fold, while
354   nonetheless providing a chimeric hook that is compatible for cognate interactions with *Vch*TnsC.
355            We next set out to investigate if this chimeric design is reciprocal; that is, can we rescue
356   DNA integration activity when combining *Pse*QCascade and *Pse*TnsC with a chimeric *Vch*TnsAB
357   design? After designing and cloning similar constructs, we were indeed able to detect integration
358   activity with the converse combination (**Supplementary Fig. 11a**). Furthermore, when we applied
359   these chimeric designs to a broader range of homologous TnsAB variants and their cognate mini-
360   Tn donor substrates, we also observed integration activity for chimeric designs derived from
361   additional transposon variants, denoted Tn*7005* and Tn*7015*[26]. Intriguingly, TnsAB chimeras
362   derived from Tn*7010* and Tn*7011* showed no evidence of activity (**Supplementary Fig. 11b**),
363   suggesting that some CASTs may require targeted screening to identify tolerable chimeric graft
364   points. Next, we explored whether this engineering approach could also generate compatible
365   chimeras between divergent CRISPR-associated transposons, candidate Type I-F (*Vch*CAST)
366   and Type V-K (*Sh*CAST) systems, each of which comprise distinct transposase architectures and
367   likely arose from unique domestication events[23]. TnsB variants derived from *Sh*CAST exhibited
368   low, but detectable levels of activity as well (**Supplementary Fig. 11b,c**); when we investigated
369   the transposon insertion orientation preference for type I/V CAST chimeras, we observed that
370   chimeras in which the TnsB was derived from *Sh*CAST exhibited a "T-LR" insertion preference,
371   as typically observed in previous *Sh*CAST studies[16,35], while type I-F CASTs exhibit a "T-RL"
372   preference[15,26] (**Supplementary Fig. 11d**). Together, these results reveal that rational, structure-
373   guided engineering of precise regions of CAST systems can overcome the natural orthogonality
374   of diverse systems, enabling novel genome editing designs.
375

## DISCUSSION

377            The unexpected paradox of poor DNA binding and strong overall integration activity of
378   *Pse*CAST (**Fig. 1b,c, Supplementary Fig. 9**), inspired us to determine cryoEM structures of
379   *Pse*QCascade and pursue rational engineering methods to improve DNA targeting. Given the
380   unique phenomenon among CAST systems to harbor 'homing' crRNAs that target conserved,
381   often essential, genes within the host genome[18,26,28,29], CAST-derived CRISPR modules may have
382   been naturally selected for weak DNA binding relative to their defense-associated CRISPR-Cas
383   counterparts, thereby reducing transcriptional repression of these essential genes. This possibility
384   underscores the need to develop a comprehensive understanding of all molecular requirements
385   and intermediate steps within the CAST transposition pathway.
386            The structure of *Pse*QCascade resembles previously determined DNA-bound type I-F
387   CAST structures[20,22], but several knowledge gaps still limit a complete understanding of the
388   mechanistic requirements for RNA-guided transposition. First, the functional relevance of the
389   Cas8 helical bundle remains a mystery. When comparing between three distinct, DNA-bound
390   QCascade structures[20,22], three different conformational states of the helical bundle have been
391   observed: a state in which the domain is unresolved, suggesting a conformationally dynamic

392  mode related to the open versus closed state of the overall QCascade complex[20]; a state in which
393  the domain is resolved, with close contact to the PAM-distal DNA[22]; and a state in which the helical
394  bundle is resolved but does not contact TniQ or the PAM-distal DNA (**Fig. 1e,f**). Despite this
395  heterogeneity, our deletions experiments clearly indicate that the helical bundle is crucial for
396  overall DNA integration to occur (**Supplementary Fig. 4**). Another area that will require future
397  study is the manner in which the QCascade complex binds TnsC, since these interactions have
398  not yet been captured for a type I-F CAST system; mutations in Cas7 that theoretically disrupt
399  Cas7.6 interactions with TniQ.2 appear to be tolerated (**Fig. 3e,f**). Although unexpected, this lends
400  credence to the possibility that only one of the two TniQ monomers present in type I-F CAST
401  complexes interacts with TnsC, which is supported by similar CAST structures from type I-B and
402  type V-K systems in which only one TniQ is present with TnsC at the target site (**Supplementary**
403  **Fig. 10**)[21,25,39]. Further *in vitro* biochemical studies, combined with structural insights into the holo
404  transpososome, will be necessary to shed light on these mechanistic aspects, including the extent
405  to which the helical bundle may regulate TnsC recruitment, and thus the targeting discrimination
406  between on- and off-target sites during CAST transposition[40].

407       Beyond defining structural requirements for transposition, our QCascade structure
408  revealed potential targets for rational engineering, most notably within the PAM interacting regions
409  of Cas8. The presence of alanine residues at this interface, rather than polar residues,
410  differentiates *Pse*CAST from homologous type I-F CAST systems (**Supplementary Fig. 6a**).
411  Interestingly, one of these homologous systems — *Vch*CAST — exhibited higher DNA binding
412  activity than *Pse*CAST in both human cells and *E. coli* (**Fig. 1c, Supplementary Fig. 9**), leading
413  us to hypothesize that reinstating polar residues might stabilize DNA-protein interactions, thereby
414  increasing DNA binding activity and integration efficiency. Mutation of even one of these alanine
415  residues yielded QCascade variants with integration efficiencies 2- to 3-fold above wild-type, but
416  interestingly, these changes did not accompany concomitant increases in PAM stringency (**Fig.**
417  **4b**). On the other hand, our episomal PAM screen in human cells revealed a wild-type 'CN'
418  preference that had not previously been observed in *E. coli*, and we hypothesize that this
419  difference may result from the larger DNA search space in the human cell milieu. The quality of
420  our cryoEM maps also provided a detailed view of RNA-DNA base-pairing interactions, enabling
421  visualization of energetically unfavorable nucleobase positioning along the heteroduplex (**Fig. 4d-**
422  **h**). Close analysis of the surrounding Cas7 residues implicated several hydrophobic side chains
423  in enforcing this positioning (**Supplementary Fig. 8**), and we therefore introduced mutations with
424  less bulky side chains to potentially stabilize heteroduplex formation. Interestingly, however, most
425  Cas7 variants complete abolished integration activity (**Fig. 4i**), suggesting that these mutations
426  adversely affected DNA binding and/or QCascade complex formation.

427       Alongside our efforts at engineering specific *Pse*CAST components for DNA integration
428  activity improvements, we considered a parallel path that would instead leverage pre-existing
429  components from homologous CAST systems. Our previous experiments revealed the orthogonal
430  properties of diverse type I-F CAST systems, which precluded mixing and matching of
431  homologous components into single systems[26], but we hypothesized that a more nuanced,
432  structure-guided approach could reveal unique opportunities for the construction of synthetic

433 chimeric designs that would retain key protein-protein interactions necessary for transposition. To
434 this end, we leveraged AlphaFold[59] to generate predicted structures of TnsA-TnsB interacting with
435 a heptameric TnsC ring (**Fig. 5b**), and based on the resemblance to previously determined type
436 V-K transpososome structures (**Supplementary Fig. 10a**)[39], we envisioned that reprogramming
437 the TnsB C-terminus could uncover functional chimeric CASTs. This hypothesis was borne out
438 with data demonstrating that chimeric CASTs, in which the DNA targeting module of *Vch*CAST
439 was combined with the DNA integration module of *Pse*CAST, functioned robustly for RNA-guided
440 DNA integration (**Fig. 5**). Next, we further extended these chimeric designs to a variety of type I-
441 F systems, and we demonstrated the first example of coordinated activity between type I-F and
442 type V-K CAST machineries (**Supplementary Fig. 11**). Based on these promising results, we
443 expect that future modifications will enable additional chimeric starting points for future
444 engineering, such as at the TniQ–TnsC interface (**Supplementary Fig. 10b,c**).

445 The ability to coordinate targeted integration with transposase proteins derived from
446 unique families[23] opens the door to novel, diverse chimeric CAST designs that can sample
447 combinatorial sequence spaces unexplored by evolution. With growing evidence that additional
448 CAST subtypes can be leveraged for genome editing applications in human cells[61–63], the ability
449 to exchange modules with ease may be key for future CAST engineering efforts. Collectively, our
450 work showcases diverse, structure-guided approaches to understand and improve CAST
451 function, and opens the door to a far greater combinatorial space for leveraging CASTs systems
452 as genome editing tools.

453

## METHODS

### Protein purification

454   
455

456   The TniQ-Cascade complex from *Pse*CAST (*Pse*QCascade) was overexpressed and
457   purified as previously described[20], with the following modifications. All proteins were codon
458   optimized and placed downstream of consensus RBS sequences, and TniQ contained an N-
459   terminal 10xHis-TEV tag. The minimal CRISPR array was encode upstream of *cas7* and
460   contained a 32 bp spacer targeting the AAVS1 locus (see **Supplementary Table 1** for detailed
461   plasmid sequences). After overnight expression at 0.5 mM IPTG, cell pellets were resuspended
462   in QCascade lysis buffer (50 mM Tris-Cl, pH 7.5, 700 mM NaCl, 0.5 mM PMSF, EDTA-free
463   Protease Inhibitor Cocktail tablets (Roche), 1 mM dithiothreitol (DTT), 5% glycerol) and lysed by
464   sonication. Lysates were clarified by centrifugation at 15,000 x g for 30 min at 4 °C. Initial
465   purification was performed by immobilized metal-ion affinity chromatography with NiNTA Agarose
466   (Qiagen) using NiNTA wash buffer (50 mM Tris-Cl, pH 7.5, 700 mM NaCl, 10 mM imidazole, 1
467   mM DTT, 5% glycerol) and NiNTA elution buffer (50 mM Tris-Cl pH 7.5, 700 mM NaCl, 300 mM
468   imidazole, 1 mM DTT, 5% glycerol). The sample was further purified by size exclusion
469   chromatography over a Superose 6 Increase 10/300 column (GE Healthcare) equilibrated with
470   QCascade storage buffer (20 mM Tris-Cl, pH 7.5, 700 mM NaCl, 1 mM DTT, 5% glycerol).
471   Fractions were pooled, concentrated, snap frozen in liquid nitrogen, and stored at −80 °C. TEV
472   cleavage was not performed.

473

### Plasmid construction

474   
475   Bacterial expression plasmids for *Pse*QCascade were codon-optimized for *E. coli* and
476   synthesized by GenScript. For human cell transfections, genetic components encoding *Pse*CAST
477   proteins were codon-optimized for human cells, synthesized by GenScript, and cloned into
478   pcDNA3.1 expression vectors. All CAST constructs were cloned into plasmids using a
479   combination of restriction digestion, ligation, Gibson assembly, and Golden Gate assembly. All
480   PCR fragments for cloning were generated in-house using Q5 DNA Polymerase (New England
481   Biolabs (NEB)) and gel purified using Qiagen Gel Extraction.

482   To clone the 4N PAM library used for HEK293T cell episomal integration assays, two
483   overlapping oligos containing 'NNNN' were phosphorylated with T4 PNK (NEB) and hybridized at
484   95 °C for 2 min before cooling to room temperature. The resulting oligoduplex was ligated into a
485   target plasmid vector predigested with BsmBI (55 °C for 2 h) using T4 DNA ligase (NEB). Cloning
486   reactions were transformed into chemically competent NEB Turbo *E. coli*, plated on agar plates
487   with the appropriate antibiotic to grow overnight, and inoculated in 5 uL LB media and antibiotic
488   for approximately 7 h. Colony counting was then performed to ensure sufficient library diversity.
489   Plasmids were then purified using Qiagen Miniprep columns verified by a combination of Sanger
490   sequencing (Azenta/Genewiz) and whole-plasmid nanopore sequencing (Plasmidsaurus), and
491   ultimately characterized by high-throughput sequencing (Illumina).

492

### CryoEM structure determination.

494       Purified *Pse*QCascade was serially diluted in a modified buffer (20 mM Tris-Cl, pH 7.5,
495   200 mM NaCl, 1 mM DTT) for initial imaging experiments. Target DNA (NTS: 5′-
496   TTCATCAAGCCATTGGACCGCCACAGTGGGGCCACTAGGGACAGGATTGGTGACCTTCGC
497   CTTGACGGCCAAAA-3′, TS: 5′-TTTTGGCCGTCAAGGCGAAGCTGAAAAGCAATGAAGCCAA
498   AGCGTCCTGTAAGGCGGTCCAATGGCTTGATGAA-3′) was duplexed by mixing the NTS and
499   TS in equimolar concentrations, heated to 95° C, and then cooled to room temperature. 50 µM
500   aliquots were then snap frozen. Purified *Pse*QCascade aliquots were incubated with a 5X molar
501   excess of target DNA for 10 min at room temperature with a total reaction volume of 50 µL. The
502   complex (2-4 µM range) was initially imaged in a Talos L120C (Thermo Fisher) electron
503   microscope equipped with a $LaB_6$ electron source and a Ceta-M camera. Negative staining
504   experiments were carried out using uranyl-formate solution at 0.75% (w/v) in water. CF-400
505   (EMS) continuous carbon grids were activated for 30 s using a $Ar/O_2$ gas mix plasma at 25 W
506   using a Solarus2 plasma cleaner (Gatan). Immediately after plasma activation, 3 µL of the
507   *Pse*QCascade/DNA complex at concentrations of 1, 2 and 4 µM were applied to the activated
508   grids. After 1 min incubation, the excess solution was gently blotted away, and 3 µL of 0.75%
509   uranyl-formate solution was added for an additional 1 min incubation. Excess staining solution
510   was blotted away and the grids were left on the bench drying for 5 min. Grid screening revealed
511   well stained, homogeneous, and dispersed particles with a circular shape compatible in
512   dimensions and shape with the estimated molecular size of the complex, as well as showing
513   similarities with previously reported images of other Cascade complexes (**Supplementary Fig.
514   2a**).

515       We chose the 1 µM concentration grid for manual collection of 10 negative staining images
516   (pixel size 2.5 Å/pixel, 1 s exposure, -2 to -3 µm defocus) for exploratory class-2D analysis in
517   Relion4. The resulting negative staining C2D averages confirmed the homogeneity of the sample
518   and its potential for high-resolution (**Supplementary Fig. 2a, left**). Next, we explored the behavior
519   of the complex under cryogenic conditions using the negative stain conditions as a reference
520   starting point. We vitrified UltraAu foil 1.2/1.3 'Gold' grids (Quantifoil) using a VitroBot Mark IV
521   (Thermo Fisher) set up to 100% humidity and 4 °C. The sample concentration was in the 2–4 µM
522   range. Grids were plasma cleaned with the same protocol described for the negative staining
523   grids, and after application of 3 µL solution, the grids were blotted and plunged frozen in liquid
524   ethane. Vitrobot settings were: blot force -5, drain and waiting time 0 with blotting times variating
525   between 2.5 to 3.5 s. Following these parameters, we froze 8 grids, 4 grids at 2 µM concentration
526   and 4 grids at 4 µM concentration. 2 grids, one at 2 µM and another at 4 µM concentration were
527   transferred to a cooled 910 side entry holder (Gatan) for screening under cryogenic conditions in
528   the same Talos L120C microscope used for negative staining using similar imaging conditions.
529   Both grids showed good ice distribution, with the 2 µM grid showing better particle distribution and
530   contrast in ice. Using SerialEM, we collected 10 images with similar settings as in negative
531   staining experiments for exploratory reference-free C2D analysis in Relion4 under cryogenic
532   conditions (**Supplementary Fig. 2a, middle**). The resulting C2D averages were promising, with
533   distinctive and multiple views of the complex. The grid was recovered and stored for high

534  resolution data collection in a Titan Krios G3i electron microscope equipped with a
535  BioQuantum/K3 energy filter and direct detection.

536      High resolution data was collected at high magnification with 2x hardware binning in the
537  K3 detector (0.6485 Å/pixel size after binning) at a fluence of ~20e$^-$/pixel/s and 1 s exposure time
538  for a total dose of ~50 e$^-$/Å$^2$. Defocus range was adjusted to vary between -0.8 to -2 μm, and the
539  total number of K3 fractions was adjusted to 50. 24 h collection on the recovered grid yielded
540  ~22,000 images which were on-the-fly motion corrected in Relion4 with ctf estimation in ctffind4.
541  Image processing was integrally done in Relion 4 and cryoDRGN. First, we manually selected
542  100 images for Laplacian picking, which yielded ~4,000 particles that were normalized and
543  extracted with 8 times binning. Fast C2D analysis using the VDAM algorithm generated C2D
544  averages in multiple orientations that were selected and used as training set for Topaz, used
545  through the Relion wrapper. Using the optimized trained model from Topaz, the full dataset of
546  ~22,000 images yielded ~1.5 million particles that after two C2D steps using T parameters of 3
547  and then 6 was reduced to ~667,000 particles. ArnA contamination accounted for the bulk of the
548  eliminated particles. Next, we refined the reduced dataset using a filtered map of *Vch*QCascade
549  as reference. We did not perform alignments with this initial classification (K20, tau fudge T=6).

550      We identified multiple classes with damaged or poorly aligned particles, a class without
551  the TniQ dimer, and a dominating class with better features. A re-extraction step was then
552  performed with the recenter option activated and at 4x binning (2.594 Å/pixel). After selection of
553  2D class averages showing secondary structure features, an ab-initio 3D model was
554  reconstructed using the Stochastic Gradient Descent (SGD) algorithm with all selected particles
555  from the class 2D job (K4, tau fudge T=3). A second 3D refinement produced a consensus
556  refinement in the 5 Å range that upon inspection showed clear secondary features and substantial
557  heterogeneity at the PAM distal region hosting the TniQ dimer. A soft-mask (10 pixel extension,
558  8 pixel soft edge and initial threshold of 0.002) was used for 3D classification without alignment
559  using 20 classes and T parameters 3, 6 and 8. A minor population (~8% of the particles) of
560  Cascade without TniQ was identified and removed from the dataset, together with poorly aligned
561  or damaged particles, reducing the total dataset to ~128,000 particles. Re-refinement of this
562  dataset after re-extraction to binning 2 (~1.2 Å/pixel) produced a sub-3Å map, but exacerbated
563  heterogeneity of the TniQ dimer region was evident.

564      Using focused classification of this region of the map produced multiple classes without
565  clear discrete states, suggesting continuous heterogeneity. Before applying a multibody
566  approach, we re-refined the ~128,000 particle dataset after refining the ctf parameters (defocus
567  values per particle and astigmatism per micrograph) followed by Bayesian particle polishing for
568  signal decay and local particle movement correction. We defined via soft masking (6 pixel mask
569  extension, 6 pixel soft edge decay, initial threshold 0.002) three rigid body groups: the first body
570  included Cas8, and the first Cas7 monomer (Cas7.1), the second body contained Cas7
571  monomers 2 to 5, and the third body included the TniQ dimer, Cas6, Cas7.6, and the crRNA 3′-
572  proximal hairpin. Residual rotation priors were defined to 10 degrees with translation offset of 2
573  pixels. We designed two wide masks: one (body 1) covering the best part of the map and including
574  Cas8, the first five Cas7 proteins, and surrounding densities including the corresponding sections

575    of the crRNA-DNA heteroduplex; and a second soft mask (body 2) covering Cas7.6, Cas6, and

576    the TniQ dimer. Multibody refinement produced maps with exceptional quality for each body, with

577    clear sub 3Å features for the Cas8 and the Cas7 regions. The maps for the PAM-distal body,

578    including the TniQ dimer, improved substantially, but residual heterogeneity remained, especially

579    at the distal end of the TniQ dimer.

580    We used ModelAngelo[64] for initial model building using the improved maps from the

581    multibody analysis. With default options and sequence information from the cloned constructs,

582    ModelAngelo correctly built approximately 90% of the residues. Manual inspection of the built

583    model corrected limited errors and completed areas where the resolution did not allow accurate

584    placement of side chains. The built models were refined against the multibody maps

585    independently, first with phenix refine (secondary structure restrain activated) and then with

586    Refmac5, adjusting the experimental/ideal geometry weights manually to avoid overfitting.

587    CryoDRGN analysis was performed with the final set of ~128,000 particles used for multibody

588    analysis in Relion. This set of particles was re-extracted to a box size of 128 pixels and an initial

589    training in 1 dimension (Zdim=1) was performed. After assessing the homogeneity of this set of

590    particles, 3 different training were performed with 2, 4 and 8 dimensions (Zdim=2, 4 and 8).

591    Principal component analysis (PCA), UMAP, and K-means clustering dimensionality reduction

592    techniques were used to explore the derived latent spaces, producing similar results irrespective

593    of the Zdim used. We perform a final training with particle re-extracted to 256 pixels size and Zdim

594    2 and 8. Exploration of the latent space derived from these training revealed multiple

595    conformations of the TniQ dimer, as shown in **Supplementary Figure 3**.

596

597    **Mammalian cell culture and transfections**

598    HEK293T cells were cultured at 37 °C and 5% $CO_2$ and maintained in DMEM media with

599    10% FBS and 100 U/mL of penicillin and streptomycin (Thermo Fisher Scientific). 24 h before

600    transfection, a 48-well plate was coated with poly-D-lysine (Thermo Fisher Scientific) and seeded

601    with 10,000 cells per well. Cells were transfected with DNA mixtures and 1 µL of Lipofectamine

602    2000 (Thermo Fisher Scientific) per the manufacturer's instructions. Transcriptional activation and

603    integration assays were performed as previously described[17]. For plasmid-based PAM library

604    assays, cells were co-transfected with the following *Pse*CAST CAST plasmids: 200 ng pTnsAB,

605    50 ng pTnsC, 75 ng pQCascade, 100 ng pCRISPR (crRNA), 200 ng pDonor, and 100 ng pTarget

606    (4N PAM library). Cells were harvested 4 days after transfection using previously described

607    methods[17].

608

609    **Analysis of HEK293T integration assays**

610    Genomic integration assays were analyzed as previously described[17]. In brief, 5 µL of

611    genomic lysate (10% of total lysate volume) was used for 2 rounds of PCR. In the first PCR, a

612    forward primer was used that anneals to the AAVS1 locus, and a reverse primer was used that

613    anneals to both the AAVS1 locus and a primer binding site in the donor DNA (see **Supplementary**

614    **Table 3** for oligonucleotide sequences). These oligos included 5′ overhangs encoding read 1 and

615    read 2 Illumina adapters. In the second PCR, 'universal' primers were used, which anneal to the

616    read 1 and read 2 sequences and append unique index sequences and the remaining Illumina
617    adapter sequences for next generation sequencing. Samples were then pooled, gel purified, and
618    sequenced on a NextSeq 500/550 with at least 75 cycles in read 1. The relative abundance of
619    reads that contain a *Pse*CAST transposon end sequence (representing an integration read) vs.
620    downstream AAVS1 sequence (unintegrated read) was calculated.
621        For the episomal PAM library assay, samples were prepared as above except a different
622    forward oligo was used that anneals directly upstream of the degenerate PAM library in PCR 1,
623    such that we would capture both the PAM sequence and the presence of the transposon end
624    sequence with the forward read (see **Supplementary Table 3** for oligonucleotide sequences).
625    PCR 1 cycles were reduced to 15 cycles. After Illumina sequencing, reads were filtered to have
626    a transposon end sequence, thus representing a PAM library member which was successfully
627    targeted by *Pse*CAST for DNA integration. The input library was sequenced as well, to calculate
628    enrichment and depletion scores. Library members were then ranked by their enrichment values
629    (proportion of output library / proportion of input library). The top 10% of library members were
630    used to generate a consensus WebLogo (Version 2.8.2, 2005-09-08, weblogo.berkeley.edu) for
631    the PAM preference of each Cas8 variant. All library members and their associated enrichment
632    values were used to generate PAM wheels using Krona[65].

633

634    ***E. coli* repression and integration assays**

635        *E. coli* transcriptional repression assays were performed as previously described[40,57], with
636    some minor modifications. In brief, an *E. coli* strain expressing mRFP from the chromosome, a
637    gift from L. S. Qi, was transformed with pQCascade. We initially attempted to use pQCascade
638    plasmids with a strong J23119 promoter, but due to toxicity associated with strong *Pse*QCascade
639    expression, we switched to a weaker J23101 promoter for all pQCascade constructs. We
640    designed crRNA sequences to target the template strand of mRFP proximal to the 5′ end of the
641    coding region (60 bp downstream of the mRFP start codon). Two replicates were performed for
642    each unique transformation, and relative mRFP repression was analyzed as previously
643    described[40].
644        Integration assays were performed as previously described[15,40], with the following
645    modifications. Although J23101 promoters were used for QCascade, J23119 promoters were still
646    used for constitutive expression of all TnsABC cassettes, as there was no observed toxicity. In
647    brief, TnsABC expression vectors harboring donor DNA (pDonor-TnsABC) encoded a *tnsA-tnsB-*
648    *tnsC* operon downstream of a strong constitutive promoter (J23119), as well as a mini-transposon
649    donor DNA of 0.9 and 1.2 kb in length for *Vch*CAST and *Pse*CAST, respectively, all on a pUC19
650    backbone. Strains harboring medium-strength J23101 promoter-controlled pQCascade
651    constructs were first made chemically competent, followed by duplicate transformations with
652    pDonor-TnsABC and lysate generation for qPCR after an 18 h incubation at 37 °C. Lysates were
653    analyzed via qPCR, as previously performed[15,40].

654

655    **Data availability**

656    Cryo-EM maps and models will be deposited on EMDB and PDB and released upon publication.

657    Source data for protein gels are included as **Supplementary Fig. 12**.

658

659    **Author Contributions**

660    G.D.L., A.R.L., and S.H.S. conceived of and designed the project. G.D.L. purified *Pse*QCascade.

661    G.D.L. and A.R.L. performed all cellular experiments and cellular experimental analyses, with the

662    exception of *E. coli* repression and integration assays, which were performed by D.J.Z and A.R.L.

663    I.S.F. collected cryoEM data and performed structure determination. G.D.L., A.R.L., I.S.F., and

664    S.H.S. discussed the data and wrote the manuscript, with input from D.J.Z.

665

679

## FIGURES



**Figure 1 | CryoEM structure of the TniQ-Cascade (QCascade) complex from *Pse*CAST. a,** Phylogenetic tree of type I-F CRISPR-associated transposons (CASTs), adapted from a previous publication[26]. Systems with previously solved QCascade structures are marked with red arrows, while *Pse*CAST is marked with a green arrow. Phylogenetic clades are colored. **b,** Experimental design to investigate both DNA binding and overall integration activities for CAST systems in human cells[17]. DNA binding is extrapolated from two different transcriptional activation assays, one in which VP64 is fused to Cas7 (left), and one in which VP64 is fused to TnsC (right). Overall integration efficiencies are measured via amplicon sequencing. **c,** Comparison of *Vch*CAST and *Pse*CAST across different assays in human cells. Although *Pse*CAST exhibits consistently weak transcriptional activation compared to *Vch*CAST, its absolute integration activity is approximately two orders of magnitude greater. DNA integration data is adapted from a previous publication[17]. **d,** Operonic architecture of *Pse*CAST components from the *Pse*CAST transposon, with genes encoding the QCascade complex labeled accordingly. **e,** Left, dominant reference-free 2D cryoEM class averages. Right, cryoEM densities with colored map regions corresponding to Cas8 (blue), Cas7 monomers 1-6 (light blue), Cas6 (purple), TniQ monomers 1-2 (orange, yellow), crRNA (gray), and target DNA (red) indicated. **f,** Refined model for the Cas8 α-helical domain and its positioning relative to the TniQ dimer interface.

**Figure 2 | The role of crRNA in the PAM-distal region of *Pse*QCascade. a,** Overall view of the cryoEM reconstruction of the *Pse*CAST QCascade complex. **b,** Magnified view of the dashed region in **a**, highlighting the cryoEM density (colored and semi-transparent) for interactions between the indicated crRNA nucleotides and protein subunits. **c,** Magnified view of the dashed regions in **b**, highlighting interactions between the crRNA and Cas6 (left), TniQ.1 (middle), and both TniQ.2 and Cas7.6 (right). Key interacting residues are labeled. **d,** Normalized RNA-guided DNA integration efficiency at *AAVS1* in HEK293T cells, as measured by amplicon sequencing. The indicated alanine mutations were designed to perturb specific RNA-protein interactions highlighted in **c**, and were compared to WT. NT, non-targeting crRNA. Data are shown as mean ± s.d. for n=3 biologically independent samples. **e,** Comparison of the crRNA conformation within the PAM-distal region, adjacent to the site of RNA hairpin stabilization by Cas6, for *Vch*CAST (PDB: 6PIJ) and *Pse*CAST (this study). The region around nucleotide G41 exhibits a distinct configuration for *Pse*CAST, likely affecting the behavior of the adjacent TniQ dimer.

**Figure 3 | TniQ recruitment to the Cas6-Cas7.6 interface of Cascade requires hydrophobic and electrostatic interactions. a,** Overall view of the *Pse*CAST QCascade complex, oriented to highlight the TniQ dimer (dark/light orange). **b,** Magnified view of the region indicated in **a**, showing how TniQ.1 (dark orange) interacts with a hydrophobic cavity on Cas6. The two visual renderings are colored either by Cas6 surface (purple, top) or hydrophobicity (bottom). **c,** Comparison of the hydrophobic interactions between TniQ.1 and Cas6 in *Pse*CAST (left) and *Vch*CAST (right, PDB: 6PIJ), with residues labeled. **d,** Normalized RNA-guided DNA integration efficiency at *AAVS1* in HEK293T cells, as measured by amplicon sequencing. The indicated arginine point mutations were designed to perturb TniQ.1-Cas6 hydrophobic interactions. NT, non-targeting crRNA. **e,** Magnified views of hydrogen bonding (top) and electrostatic (bottom) interactions between Cas7.6 (blue) and TniQ.2 helix (yellow). **f,** Normalized RNA-guided DNA integration efficiency at *AAVS1* in HEK293T cells, as measured by amplicon sequencing. Alanine mutations perturbing Cas7.6-TniQ interactions are generally tolerated. Data in **d, f** are shown as mean ± s.d. for n=3 biologically independent samples.

**Figure 4 | Structural and functional consequences of PAM and target DNA recognition by**
***Pse*QCascade. a, Top,** overall view of the *Pse*CAST QCascade complex, oriented to highlight
the target DNA recognition. **Bottom,** magnified views of the PAM binding pocket, with Cas8 and
DNA shown in blue and red, respectively. Residues A243 and A244 lack any base-specific,
hydrogen-bonding interactions with the DNA. **b,** Normalized genomic integration efficiencies at
AAVS1 for the indicated Cas8 mutants (top), plotted above the WebLogo for PAM preferences in
the -1 and -2 positions (bottom) derived from integration into pTarget. (For additional PAM
specificity data, see **Supplementary Fig. 6e**.) Integration efficiency data are shown as mean ±
s.d. for n=3 biologically independent samples. **c,** Magnified view of the experimental cryoEM
density map around Cas7.1 and Cas7.2, showing interactions with the crRNA (gray) and DNA
target strand (TS, red). NTS, DNA non-target strand. **d,** Overlay of the refined atomic model and
cryoEM density (semi-transparent) for the seed region of QCascade bound to the DNA target
strand. **e,** Schematic representation showing angles for the first five RNA-DNA base pairs (BP 1–
5) within the R-loop. **f,** View of the RNA-DNA heterduplex at right, highlighting the unfavorable
base-pairing surrounding flipped out nucleobases within the first 18 base pairs of the R-loop. **g,**
Magnified view of the RNA-DNA heteroduplex segments aligned at the flipped out base pair,
revealing consistent unfavorable angles at the adjacent base pairs. **h,** Normalized RNA-guided
DNA integration efficiency at *AAVS1* in HEK293T cells for the indicated Cas7 mutations, as
measured by amplicon sequencing. Data are shown as mean ± s.d. for n=3 biologically
independent samples.

**Figure 5 | AlphaFold-guided engineering of TnsABC to generate chimeric CAST systems. a,** Schematic showing the approach to generate a chimeric CAST system by combining optimal DNA targeting and DNA integration machineries from distinct CAST systems. **b,** AlphaFold-generated structure prediction of the TnsABC co-complex from *Pse*CAST. The C-terminal "hook" region of TnsB that putatively interact with TnsC is marked. **c,** Visualization of select TnsB graft points within the predicted *Pse*TnsABC structure. Residues where *Pse-Vch* chimerism was introduced are colored in blue, and the three top performing graft points (V585, S589, Q594; *Pse*TnsB numbering) from panel **e** are labeled. **d,** Experimental workflow to test chimeric TnsAB constructs for RNA-guided DNA integration activity. *E. coli* BL21(DE3) cells containing a pEffector encoding *Vch*QCascade and *Vch*TnsC were transformed with a plasmid encoding a mini-transposon (mini-Tn) and TnsAB, with TnsAB derived from either *Vch*CAST, *Pse*CAST, or a chimeric combination thereof. Integration efficiency was measured by qPCR (bottom). **e,** DNA integration efficiencies for each tested TnsAB chimera. The amino acid listed represents the position at which the reading frame was grafted from *Pse*TnsB (red) to *Vch*TnsB (blue). "Custom" denotes a variant in which 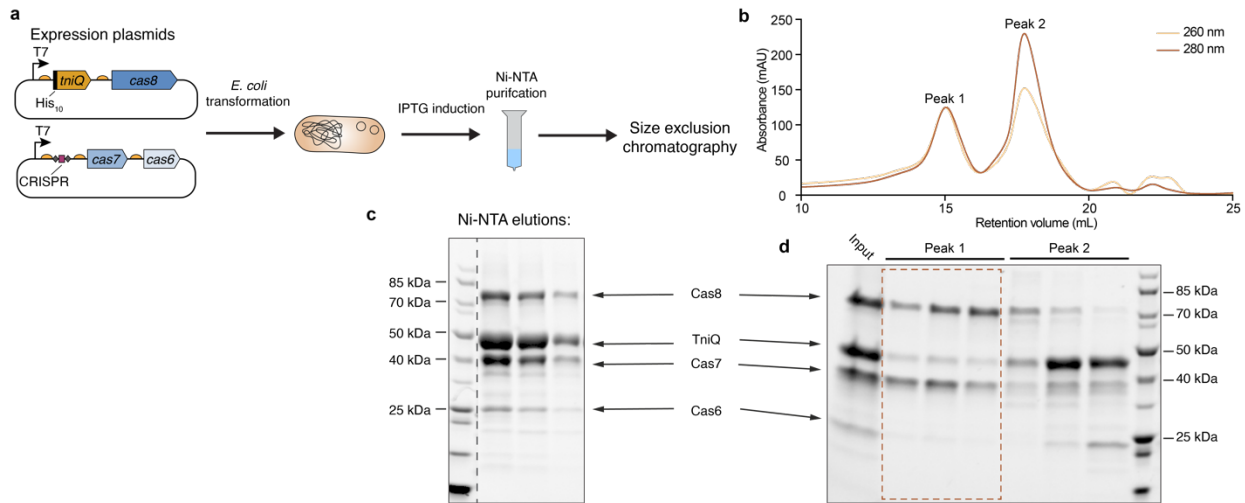multiple different *Vch*TnsB sequences were substituted (see **Supplementary Table 2** for details). Data are shown as mean for n=2 biologically independent samples.

## SUPPLEMENTARY FIGURES



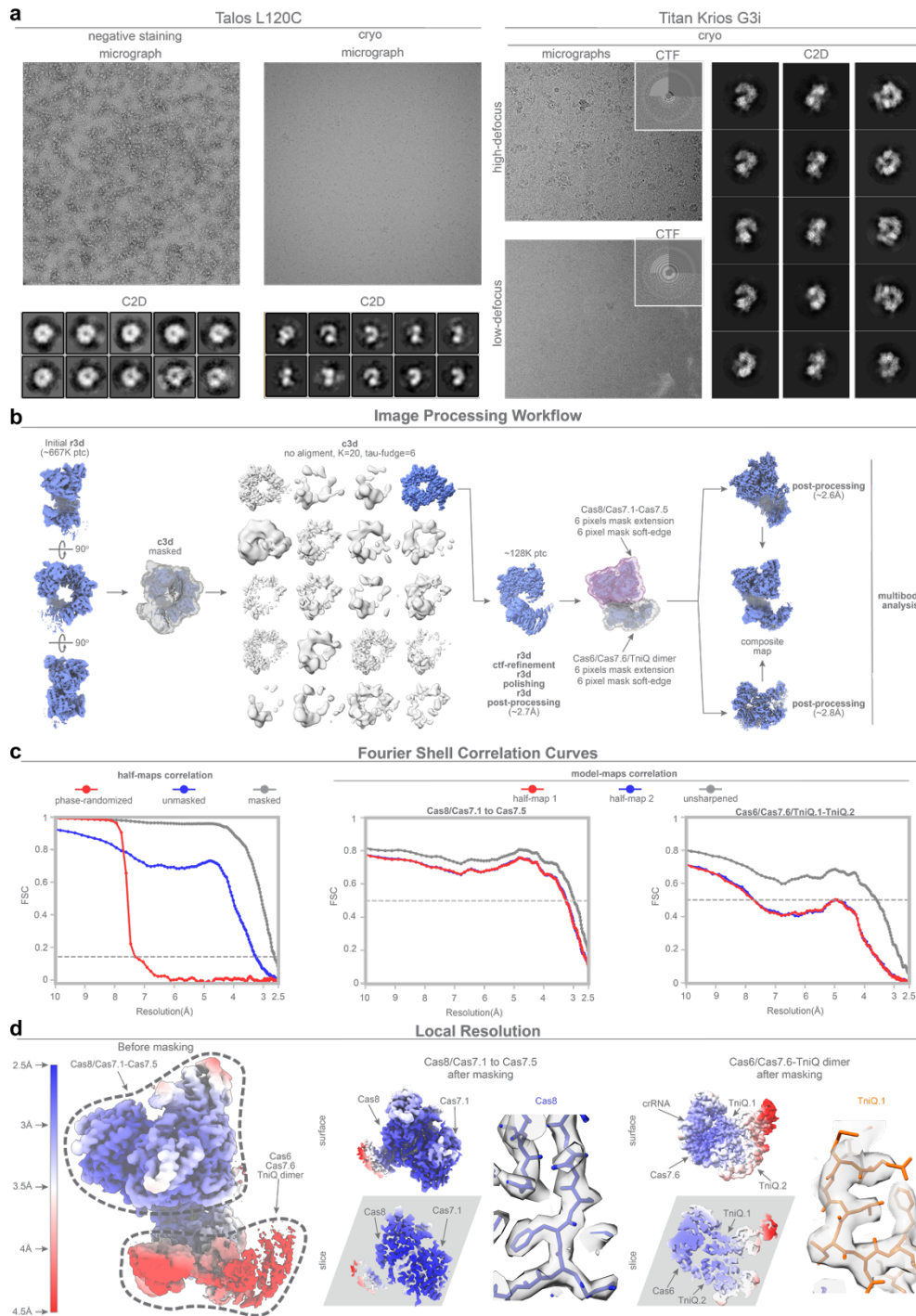**Supplementary Figure 1 | Purification of Qcascade from *Pse*CAST (Tn*7016*). a,** The two schematized expression plasmids (left) encode *E. coli* codon-optimized *Pse*CAST QCascade genes and a crRNA cassette, with a. strong ribosome binding site (half-circle) upstream of each protein-coding gene. After transformation of BL21(DE3) cells and IPTG induction, *Pse*QCascade was purified via Ni-NTA affinity chromotography and size exclusion chromatography (SEC). Codon-optimized expression plasmids were used after the native operon failed to generate detectable QCascade complexes after SEC. **b,** SEC chromatogram of *Pse*QCascade showing 2 distinct peaks. **c,** SDS-PAGE gel of representative Ni-NTA elution fractions that were pooled and used for SEC. QCascade subunits are labeled. **d,** SDS-PAGE gel of both peaks from SEC. Elutions from peak 1, marked with a red dashed box, were pooled and used for cryoEM.

787



788

789  **Supplementary Figure 2 | CryoEM imaging, data processing, and model refinement. a**,
790  Preliminary sample characterization for cryoEM grid optimization. Left, Talos L120C microscope
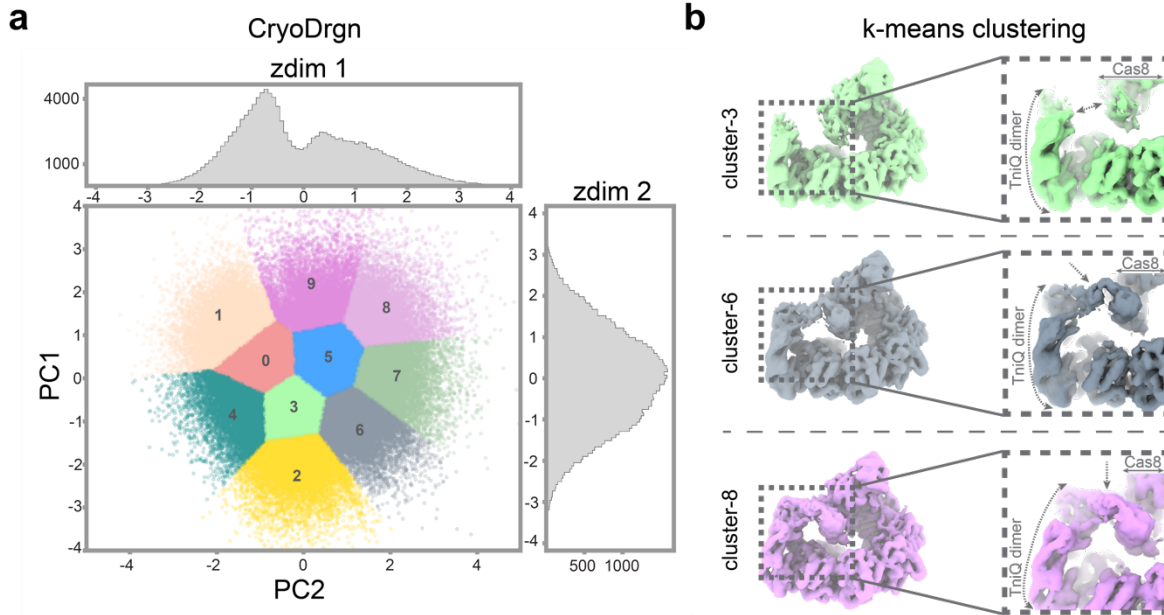791  analysis showing exemplary negative staining micrograph (left) and cryogenic micrograph (right).
792  Corresponding reference-free 2D class averages from particles obtained from 10 images are
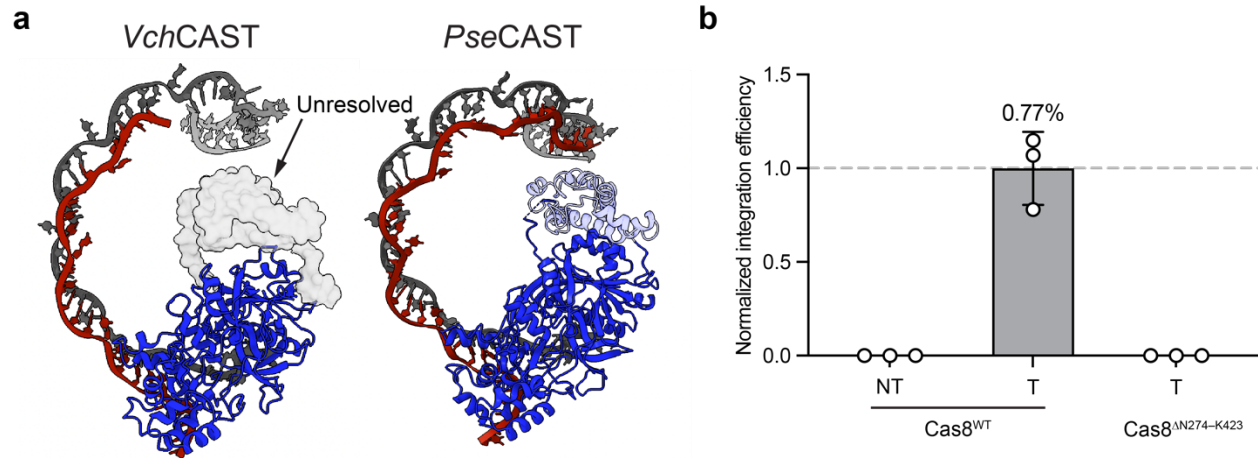793  shown below each image, with a calibrated pixel size of 2.5 Å. Right, two grids from the Talos
794  L120C screening were recovered and loaded into a Titan Krios G3i microscope, and a large

795    dataset was collected at a pixel size of 0.644 Å. Two images at different defoci are shown with
796    their corresponding CTF images (inset). Reference-free 2D class averages are shown on the
797    right, with multiple different views revealing details compatible with protein secondary structure.
798    **b**, Image processing workflow implemented in Relion4 for high-resolution structure determination.
799    Briefly, from left to right: an ab-initio 3D model was reconstructed after selection of 2D class
800    averages; using this as a reference, a consensus refinement was generated; inspection of this
801    preliminary map revealed heterogeneity, especially in the region of the TniQ dimer (**Methods**).
802    However, after unbinning and multiple rounds of 3D refinements, the map still exhibited residual
803    heterogeneity in the region adjacent to the Cas6 protein, suggesting mobility of the TniQ dimer
804    with respect to Cascade. To improve the maps and to analyze TniQ dynamics, two masks were
805    designed (**Methods**), yielding improved the densities and B-factors for the first body, but the
806    second body exhibited a significant improvement in terms of resolution and general density
807    quality. **c,** Fourier Shell Correlation (FSC) curves for the half-maps and model-maps. **d,** Local
808    resolution depictions of the final map before and after the multibody approach.
809

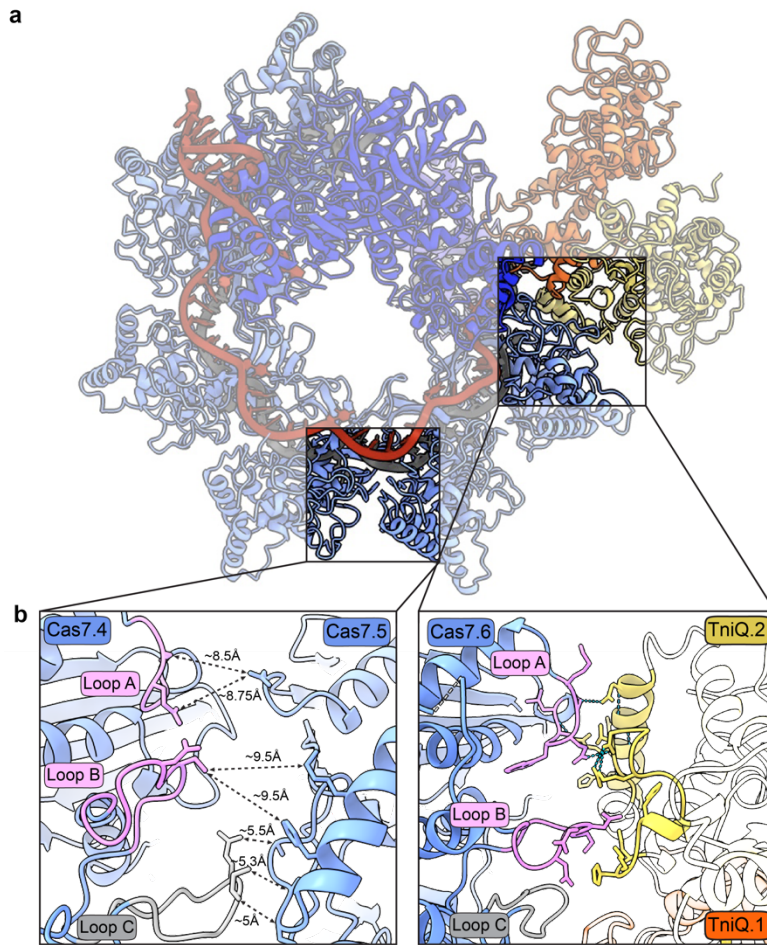**Supplementary Figure 3 | Visualization of TniQ dimer dynamics with cryoDRGN. a**, cryoDRGN analysis, using the same set of particles (~128,000) identified using Relion4 classifications, revealed dynamics of the TniQ dimer and uncovered multiple conformational states. We trained cryoDRGN on our dataset with multiple values of the zdim (2, 4 and 8), and found that the derived latent space for different runs was similar. Shown is a principal component analysis of the latent space derived from the run at zdim = 2. **b**, Segmentation of this latent space via k-mean clustering reveled multiple TniQ dimer conformations: an 'open' position, in which the TniQ dimer is distant from the Cas8 α-helical domain (cluster 3, green); an intermediate position, where the distal end of the TniQ dimer marginally contacts the Cas8 α-helical domain (cluster 6, grey); and a compact conformation, in which the TniQ dimer closely approaches the Cas8 α-helical domain  (cluster 8, pink). In all cryoDRGN-generated maps, the Cas8 α-helical domain remains in a similar position and conformation. with only the TniQ dimer exhibiting pronounced fluctuations

**Supplementary Figure 4 | Cas8 α-helical domain deletion abolishes RNA-guided DNA integration. a,** Comparison of select regions of the DNA-bound QCascade complex from *Vch*CAST (left, PDB: 6PIJ) and *Pse*CAST (right), including the crRNA (grey), target DNA (red), and Cas8 (blue). The Cas8 α-helical domain from *Pse*CAST (residues 274–423) is shown in light blue, and was replaced with a flexible, 10-amino acid GS linker in subsequent integration assays. **b,** Normalized efficiency of RNA-guided DNA integration at *AAVS1*, tested in HEK293T cells and measured by amplicon sequencing (**Methods**). Experiments used WT Cas8 and either a non-targeting (NT) or targeting (T) crRNA, or a targeting crRNA and Cas8 mutant, in which residues N274–K243 were replaced with a 10-amino acid GS linker. Data are shown as mean ± s.d. for n=3 independent biological samples.
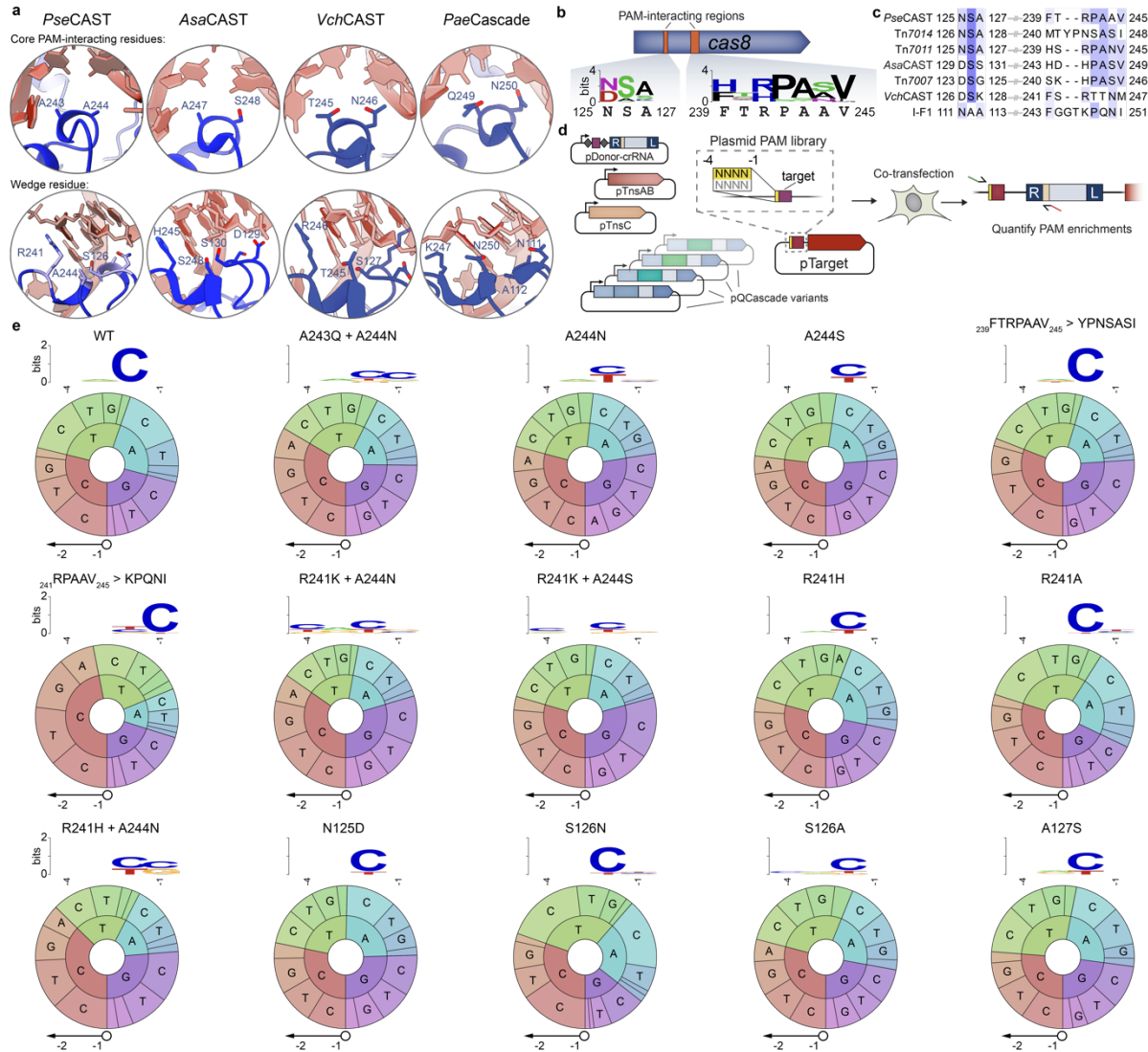
**Supplementary Figure 5 | Cas7 loops chosen to selectively perturb Cas7-TniQ.2 interactions. a,** View of overall *Pse*CAST QCascade complex, with specific regions for panel **b** highlighted. **b,** Magnified view of the different Cas7 loop interactions. Loop C participates in interactions at the interface between Cas7 monomers (left) and was therefore left intact. Amino acid sidechains in loops A and B (pink) that interact more closely with TniQ.2 were selected for mutagenesis, as detailed in **Figure 3e,f**.

**Supplementary Figure 6 | Experimental design and results for PAM library screening with Cas8. a,** Visualization of PAM binding pockets for diverse type I-F Cascade complexes (from left to right): *Pse*CAST (this study), *Asa*CAST (PDB: 7U5D), *Vch*CAST (PDB: 6PIJ), and *Pae*Cascade (PDB: 6NE0)[47]. The top inset shows core PAM-interacting residues; the bottom inset shows the wedge residue and additional interacting residues. **b,** Amino acid sequence conservation within PAM-interacting regions of *Pse*Cas8, with the WebLogo derived from a multiple sequence alignment (MSA) of 66 homologs; the *Pae*Cas8 WT sequence is shown below. **c,** MSA of the same regions from **b**, shown for diverse type I-F Cas8 homologs from both CAST and canonical type I-F1 CRISPR-Cas systems. Conserved residues are colored in blue. **d,** Mammalian PAM library assay workflow. A target pla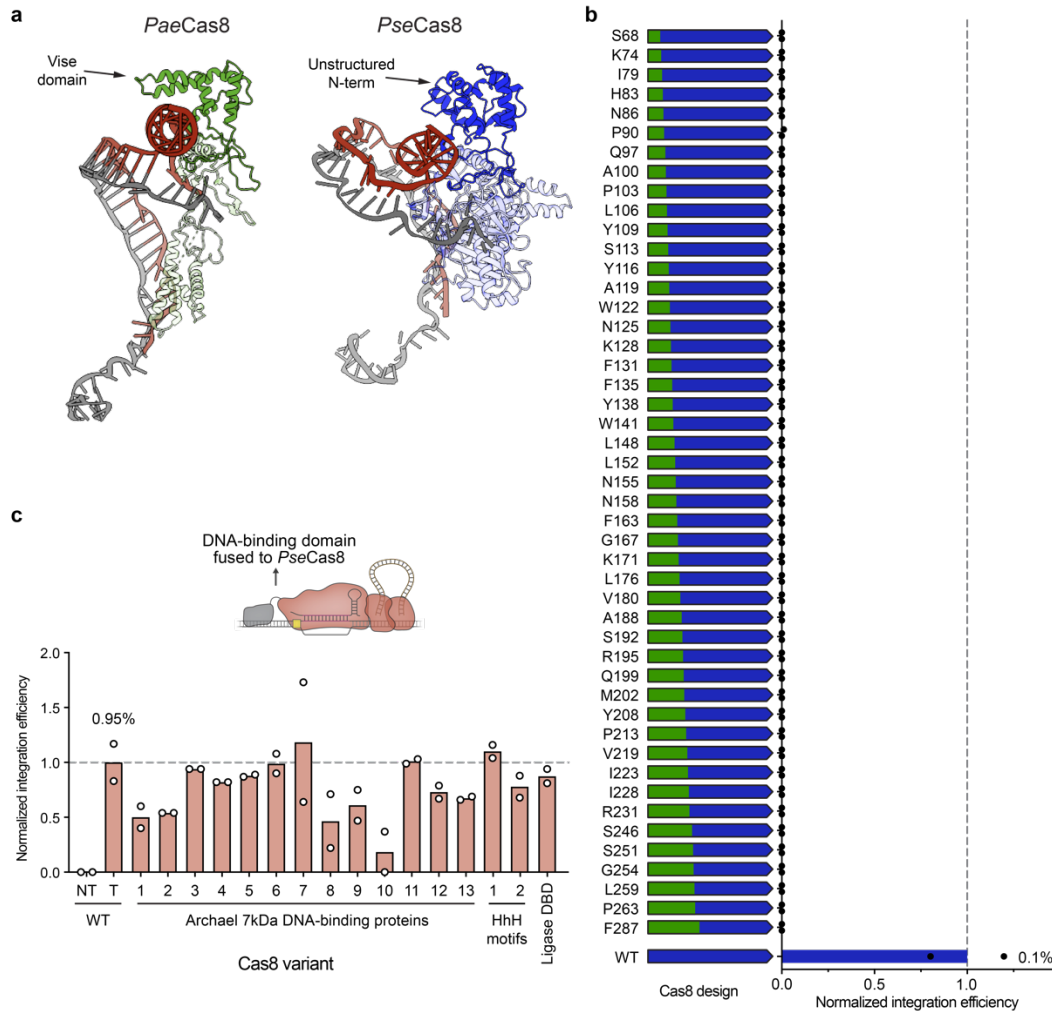smid (pTarget) was generated that contains an *AAVS1* target flanked by a 4-bp randomized PAM library. Individual Cas8 mutants were screened in each transfection via a plasmid-based integration assay, in which junction PCR and next-generation sequencing revealed PAM sequences enriched within integration products (**Methods**). **e,** Detailed PAM library data for all active Cas8 variants, showing the identity of the mutation(s) (top),
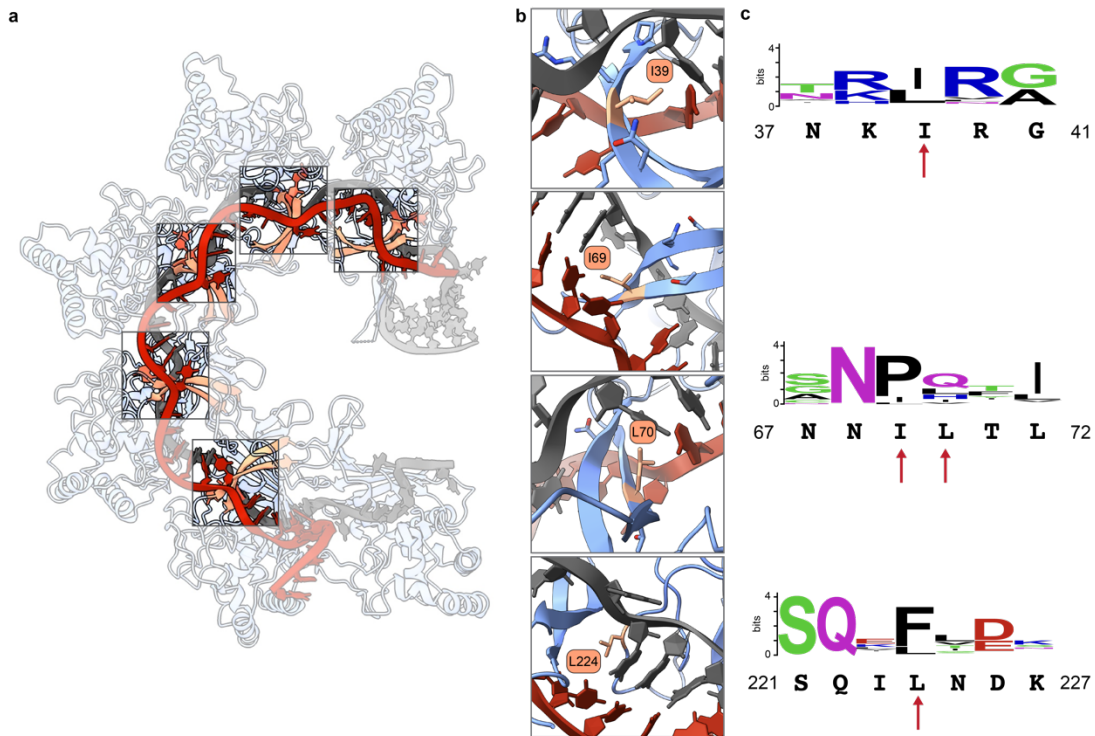
860    WebLogo of the top 10% of enriched library members (middle), and PAM wheel[65] of all library
861    members (bottom)[65]. The PAM wheel is displayed with the inner and outer rings representing the
862    -1 and -2 PAM positions, respectively.
863

**Supplementary Figure 7 | Investigating integration activity via engineering DNA-binding ability of Cas8. a,** Comparing N-terminal regions of *Pae*Cas8 (PDB: 6NE0) and *Pse*Cas8. While *Pae*Cas8 (left) shows a vise domain clamped around the dsDNA backbone, *Pse*Cas8 shows an unstructured region at the N-terminus that does not exhibit clear dsDNA backbone interactions. **b,** Normalized RNA-guided DNA integration efficiency at *AAVS1* in HEK293T cells as measured by amplicon sequencing, for a panel of chimeric Cas8 designs in which the N-terminus of *Pae*Cas8 (type I-F1 CRISPR-Cas system, green) was grafted onto the N-terminus of *Pse*Cas8 (type I-F3 *Pse*CAST, blue); the amino acid residue listed at left indicates the graft point (*Pse*Cas8 numbering). All chimeric designs tested were non-functional for DNA integration. **c,** Normalized RNA-guided DNA integration efficiency at *AAVS1* in HEK293T cells as measured by amplicon sequencing, for a panel of Cas8 fusions designed to improve DNA binding affinity. Thirteen unique archael 7 kDa DNA-binding proteins[66], two helix–hairpin–helix DNA binding motifs ('HhH')[67], and one binding domain from *Pyrococcus abyssi* DNA ligase[52] ('Ligase DBD') were tested as N-terminal *Pse*Cas8 fusions, compared to non-targeting (NT) and targeting (T) controls with WT Cas8. Data in **b** and **c** are shown as mean for n=2 biologically independent samples.

881



882

883 **Supplementary Figure 8 | Detailed view of Cas7 interactions with the RNA-DNA**
884 **heteroduplex. a,** View of overall *Pse*QCascade complex, with the five similar Cas7-crRNA
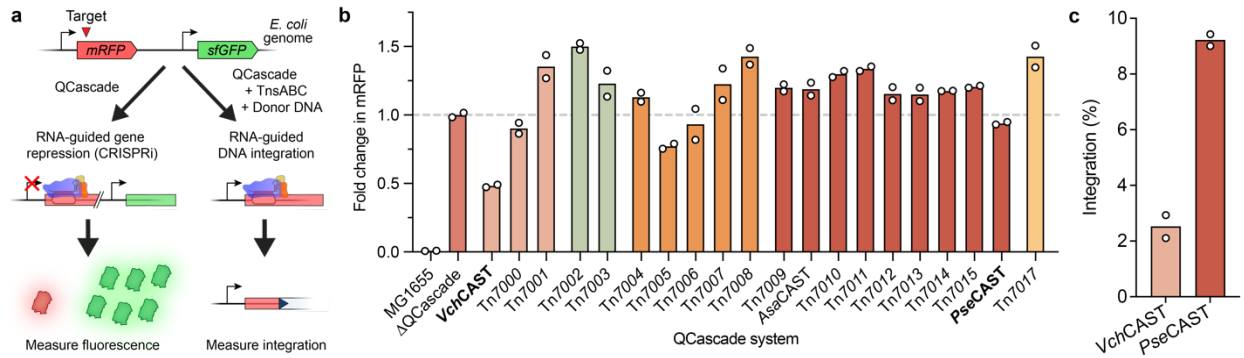885 interactions highlighted. **b,** Visualization of Cas7 residues that interact with the crRNA at each
886 flipped out nucleobase; residues with bulky and hydrophobic sidechains are highlighted and
887 labeled. **c,** *Pse*Cas7 sequence conservation at residues in panel **b**, from a multiple sequence
888 alignment of 98 homologs; the WT sequence is shown below the x-axis. Specific residues
889 selected for functional investigation are marked with red arrows.
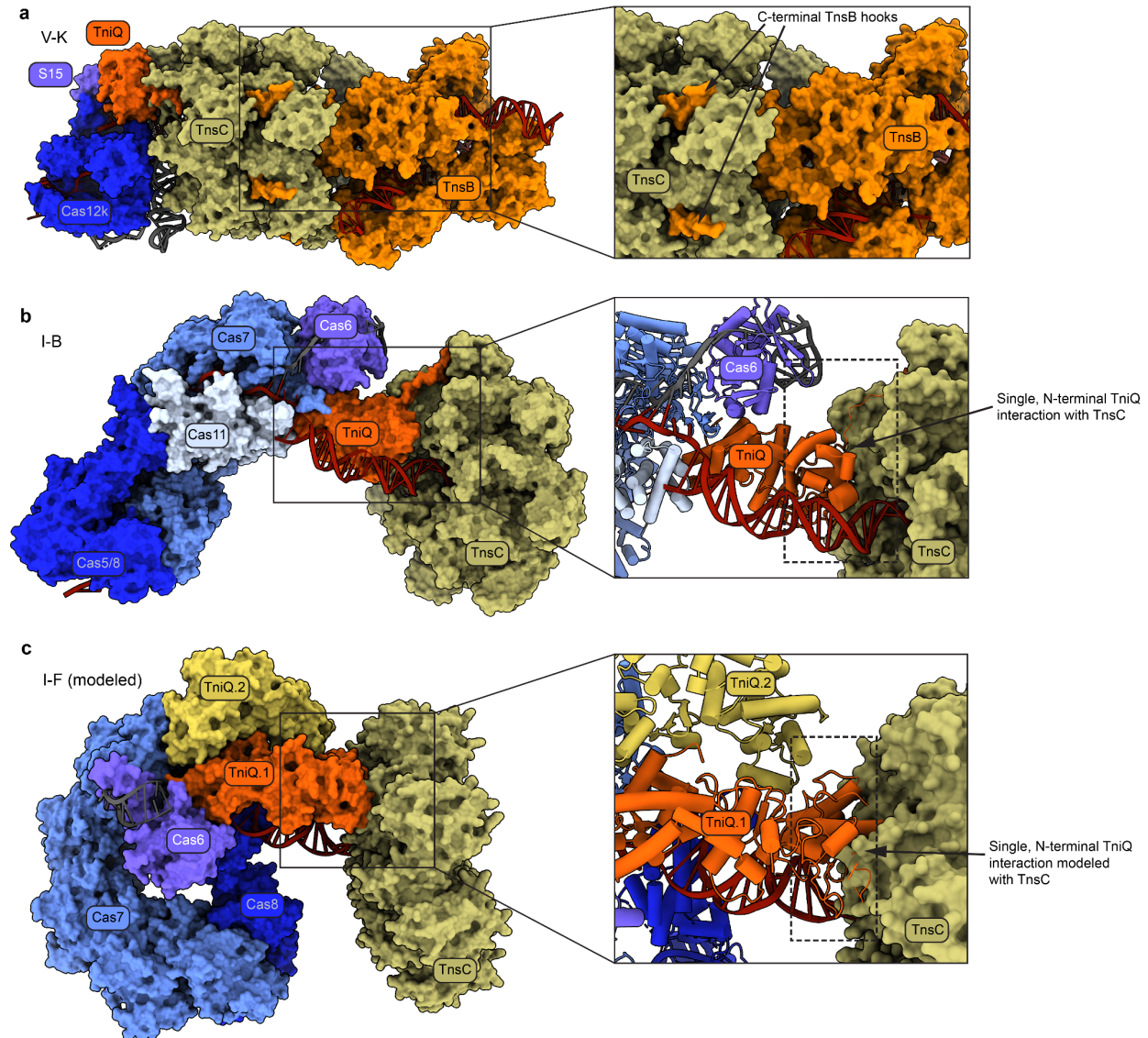
890

**Supplementary Figure 9 | DNA binding and integration activity of diverse CAST systems in** *E. coli*. **a,** Schematic of *E. coli* transcriptional repression and DNA integration assays to investigate CAST-encoded QCascade activity in bacteria. Using an engineered *E. coli* strain that constitutively expresses mRFP and sfGFP[57], transformation of[57]. a QCascade expression plasmid driven by a medium-strength J23101 promoter leads to target DNA binding (red triangle) and mRFP repression. Alternatively, when cells are co-transformed with QCascade, TnsABC, and pDonor, RNA-guided DNA integration occurs at the mRFP target site. **b,** Bar graph showing the fold change in mRFP fluorescence for each CAST-encoded QCascade system, relative to a control experiment lacking QCascade (ΔQCascade); *Vch*CAST and *Pse*CAST are highlighted in bold text. CAST systems are colored by phylogenetic clade, as shown in **Fig. 1a**. **c,** Bar graph comparing DNA integration activity for *Vch*CAST and *Pse*CAST at the same mRFP target site used for repression assays, as measured by qPCR. As observed in human cells, *Pse*CAST yields higher levels of DNA integration activity despite exhibiting apparent weaker QCascade-based DNA targeting and repression. Data in **b,c** are shown as mean for n=2 independent biological samples.
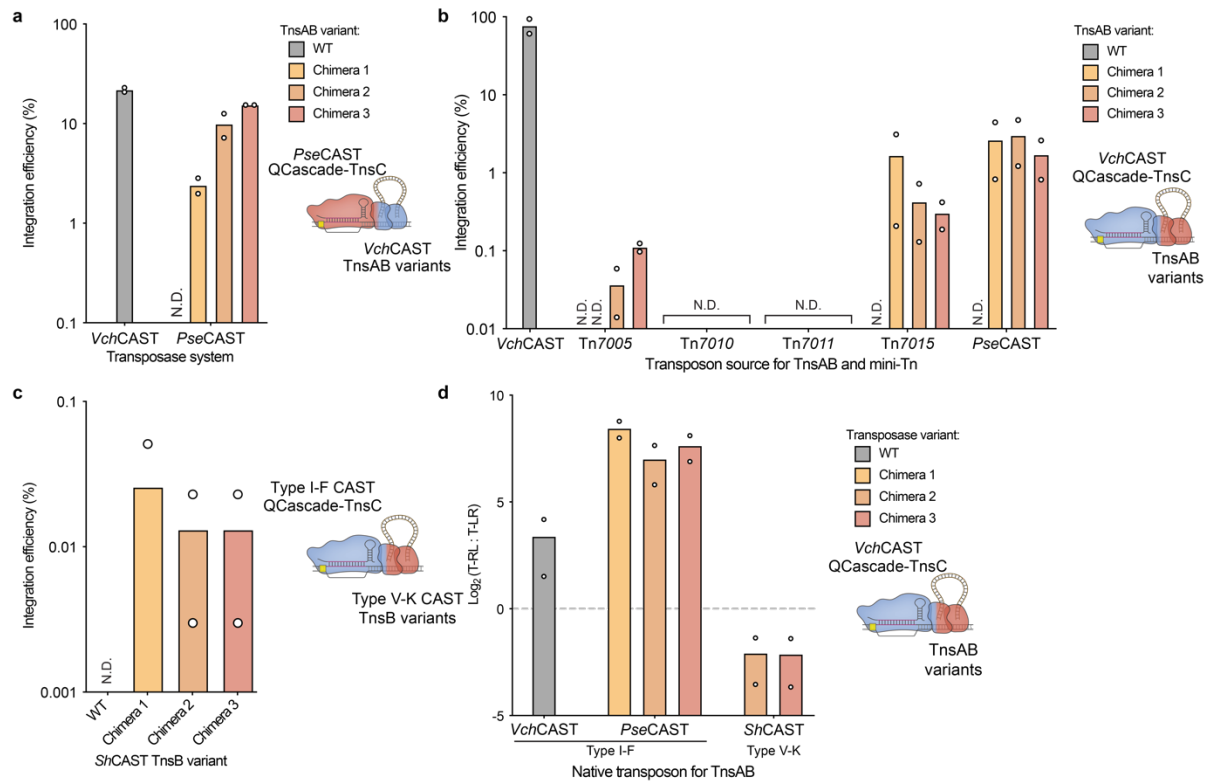
908



909

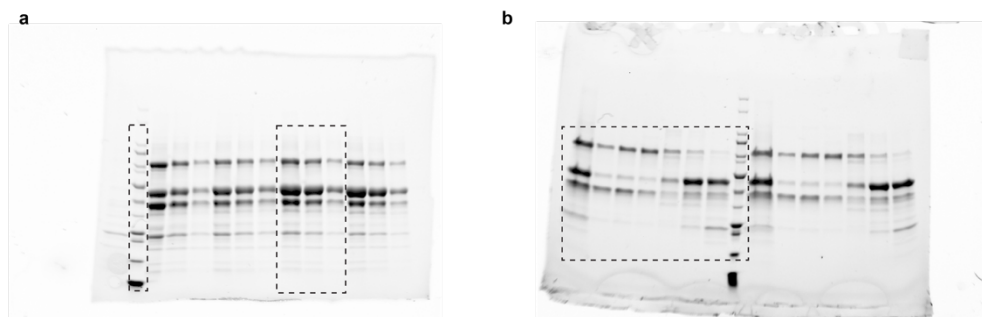**Supplementary Figure 10 | Structural inspiration for the rational design of chimeric CAST systems. a,** The holo transpososome structure from type V-K ShCAST system (PDB 8EA3), with a magnified view (right) showing how the C-terminal hook of TnsB docks into the TnsC ATPase. **b,** The QCascade-TnsC structure from type I-B *Pmc*CAST system (PDB 8FF4), with a magnified view (right) showing the N-terminus of a monomeric TniQ interacting with the TnsC ATPase. **c,** Predicted QCascade-TnsC structure from type I-F CAST, based on previous modelling[40] but with PDB ID: 7U5D, for which the PAM-distal DNA is better resolved. The magnified view (right) highlights the putative TniQ-TnsC interface, with the N-terminus of just one TniQ monomer within the dimeric arrangement interacting with the TnsC ATPase.

919

**Supplementary Figure 11 | Additional chimeric TnsB designs are functional for RNA-guided DNA integration. a,** Investigating reciprocal chimeric designs to coordinate transposition between *Pse*QCascade-TnsC and *Vch*TnsAB. WT TnsAB sequences for both *Vch*CAST and *Pse*CAST and three unique chimeric inspired by the most active variants in **Fig. 5e** (variants V585, S589, and Q594) were tested. Only chimeric TnsAB variants enabled coordinated DNA integration activity when combining *Pse*QCascade and *Pse*TnsC with *Vch*TnsAB and *Vch* mini-transposons. **b,** Exploring chimeric CASTs across multiple type I-F systems. Chimeric TnsAB variants enable coordinated transposition when combining *Vch*QCascade-TnsC with TnsAB constructs sourced from diverse Type I-F CASTs; Tn numbers were defined previously[26]. **c,** Designing chimeric CASTs across evolutionarily distinct CAST families. Chimeric *Sh*CAST TnsB constructs (inspired by functional chimeric *Pse*TnsABs) can coordinate low levels of transposition between type I-F and type V-K CAST systems. For chimera 1, only one of two biological replicates exhibited detectable integration. **d,** Insertion site orientation preference of *Vch*CAST, *Pse*CAST TnsAB chimeras, and *Sh*CAST TnsB chimeras. *Vch*CAST TnsAB and *Pse*CAST TnsAB chimeras adopt the common T-RL preference; *Sh*CAST TnsB chimeras invert the insertion site orientation preference, adopting the previously observed T-LR preference for *Sh*CAST systems[16]. Data shown as mean for n=2 independent biological samples. Chimeras 1, 2, and 3 for all homologs are listed in **Supplementary Table S2.**

**Supplementary Figure 12 | Uncropped protein gels. a,** Uncropped image used for **Supplementary Figure 1c**. The regions shown are marked with dashed boxes. **b,** Uncropped image used for **Supplementary Figure 1d**. The regions shown are marked with a dashed box.

945 **Supplementary Movie 1 | Conformational transitions revealed by Relion Multibody analysis**
946 **defining two bodies.** The first body included Cas8, the DNA-RNA duplex, and all Cas7
947 monomers; the second body included the TniQ dimer, Cas6, and the corresponding fragment of
948 the crRNA interacting with Cas6.
949
950 **Supplementary Movie 2 | CryoDRGN analysis of *Pse*QCascade flexibility visualized**
951 **through k-means clustering of the latent space.** Morphing between the two most populated
952 states after segmentations into 20 clusters is shown.
953

954  **SUPPLEMENTARY TABLES**

955

956  **Table S1. Description and sequence of plasmids used in this study.**

957  **Table S2. Sequence of chimeric TnsAB protein sequences used in this study.**

958  **Table S3: Oligonucleotides used for amplicon sequencing in this study.**

959

## REFERENCES

960   **REFERENCES**

961   1.   Branzei, D. & Foiani, M. Regulation of DNA repair throughout the cell cycle. *Nat. Rev. Mol.*

962        *Cell Biol.* **9**, 297–308 (2008).

963   2.   Heyer, W.-D., Ehmsen, K. T. & Liu, J. Regulation of Homologous Recombination in

964        Eukaryotes. *Annu. Rev. Genet.* **44**, 113–139 (2010).

965   3.   Pawelczak, K. S., Gavande, N. S., VanderVere-Carozza, P. S. & Turchi, J. J. Modulating

966        DNA Repair Pathways to Improve Precision Genome Engineering. *ACS Chem. Biol.* **13**,

967        389–396 (2018).

968   4.   Kanca, O. *et al.* An efficient CRISPR-based strategy to insert small and large fragments of

969        DNA using short homology arms. *eLife* **8**, e51539 (2019).

970   5.   Zuccaro, M. V. *et al.* Allele-Specific Chromosome Removal after Cas9 Cleavage in Human

971        Embryos. *Cell* **183**, 1–15 (2020).

972   6.   Adikusuma, F. *et al.* Large deletions induced by Cas9 cleavage. *Nature* **560**, E8–E9 (2018).

973   7.   Kosicki, M., Tomberg, K. & Bradley, A. Repair of double-strand breaks induced by CRISPR–

974        Cas9 leads to large deletions and complex rearrangements. *Nat. Biotechnol.* **36**, (2018).

975   8.   Leibowitz, M. L. *et al.* Chromothripsis as an on-target consequence of CRISPR–Cas9

976        genome editing. *Nat. Genet.* **53**, 895–905 (2021).

977   9.   Nahmad, A. D. *et al.* Frequent aneuploidy in primary human T cells after CRISPR–Cas9

978        cleavage. *Nat. Biotechnol.* **40**, 1807–1813 (2022).

979   10.  Tsuchida, C. A. *et al.* Mitigation of chromosome loss in clinical CRISPR-Cas9-engineered T

980        cells. *Cell* **186**, 4567-4582.e20 (2023).

981   11.  Komor, A. C., Kim, Y. B., Packer, M. S., Zuris, J. A. & Liu, D. R. Programmable editing of a

982        target base in genomic DNA without double-stranded DNA cleavage. *Nature* **533**, 420–424

983        (2016).

984   12.  Kim, Y. B. *et al.* Increasing the genome-targeting scope and precision of base editing with

985        engineered Cas9-cytidine deaminase fusions. *Nat. Biotechnol.* **35**, 371–376 (2017).

986    13. Gaudelli, N. M. et al. Programmable base editing of T to G C in genomic DNA without DNA

987          cleavage. Nature 551, 464–471 (2017).

988    14. Anzalone, A. V. et al. Search-and-replace genome editing without double-strand breaks or

989          donor DNA. Nature 576, 149–157 (2019).

990    15. Klompe, S. E., Vo, P. L. H., Halpin-Healy, T. S. & Sternberg, S. H. Transposon-encoded

991          CRISPR–Cas systems direct RNA-guided DNA integration. Nature 571, 219–225 (2019).

992    16. Strecker, J. et al. RNA-guided DNA insertion with CRISPR-associated transposases.

993          Science 364, 48–53 (2019).

994    17. Lampe, G. D. et al. Targeted DNA integration in human cells without double-strand breaks

995          using CRISPR-associated transposases. Nat. Biotechnol. 42, 87–98 (2023).

996    18. Saito, M. et al. Dual modes of CRISPR-associated transposon homing. Cell 184, 2441-

997          2453.e18 (2021).

998    19. Hsieh, S. & Peters, J. E. Discovery and characterization of novel type I-D CRISPR-guided

999          transposons identified among diverse Tn7-like elements in cyanobacteria. 51, 765–782

1000         (2023).

1001   20. Halpin-Healy, T. S., Klompe, S. E., Sternberg, S. H. & Fernández, I. S. Structural basis of

1002         DNA targeting by a transposon-encoded CRISPR–Cas system. Nature 577, 271–274

1003         (2020).

1004   21. Wang, S., Gabel, C., Siddique, R., Klose, T. & Chang, L. Molecular mechanism for Tn7-like

1005         transposon recruitment by a type I-B CRISPR effector. Cell 186, 4204-4215.e19 (2023).

1006   22. Park, J. U. et al. Multiple adaptations underly co-option of a CRISPR surveillance complex

1007         for RNA-guided DNA transposition. Mol. Cell 83, 1827-1838.e6 (2023).

1008   23. Faure, G. et al. CRISPR–Cas in mobile genetic elements: counter-defence and beyond.

1009         Nat. Rev. Microbiol. 17, 513–525 (2019).

1010   24. Vo, P. L. H., Acree, C., Smith, M. L. & Sternberg, S. H. Unbiased profiling of CRISPR RNA-

1011         guided transposition products by long-read sequencing. Mob. DNA 12, 1–17 (2021).

1012　25. Schmitz, M., Querques, I., Oberli, S., Chanez, C. & Jinek, M. Structural basis for the

1013　　　assembly of the type V CRISPR-associated transposon complex. *Cell* **185**, 4999-5010.e17

1014　　　(2022).

1015　26. Klompe, S. E. *et al.* Evolutionary and mechanistic diversity of Type I-F CRISPR-associated

1016　　　transposons. *Mol. Cell* **82**, 616-628.e5 (2022).

1017　27. Roberts, A., Nethery, M. A. & Barrangou, R. Functional characterization of diverse type I-F

1018　　　CRISPR-associated transposons. *Nucleic Acids Res.* **50**, 11670–11681 (2022).

1019　28. Rybarski, J. R., Hu, K., Hill, A. M., Wilke, C. O. & Finkelstein, I. J. Metagenomic discovery of

1020　　　CRISPR-associated transposons. *Proc. Natl. Acad. Sci. U. S. A.* **118**, (2021).

1021　29. Petassi, M. T., Hsieh, S. & Peters, J. E. Guide RNA Categorization Enables Target Site

1022　　　Choice in Tn7-CRISPR-Cas Transposons. *Cell* **183**, 1757-1771.e18 (2020).

1023　30. Walker, M. W. G., Klompe, S. E., Zhang, D. J. & Sternberg, S. H. Novel molecular

1024　　　requirements for CRISPR RNA-guided transposition. *Nucleic Acids Res.* **51**, 4519–4535

1025　　　(2023).

1026　31. Vo, P. L. H. *et al.* CRISPR RNA-guided integrases for high-efficiency, multiplexed bacterial

1027　　　genome engineering. *Nat. Biotechnol.* **39**, 480–489 (2021).

1028　32. Rubin, B. E. *et al.* Species- and site-specific genome editing in complex bacterial

1029　　　communities. *Nat. Microbiol.* **7**, 34–47 (2022).

1030　33. George, J. T. *et al.* Mechanism of target site selection by type V-K CRISPR-associated

1031　　　transposases. *Science* **382**, (2023).

1032　34. Strecker, J., Ladha, A., Makarova, K. S., Koonin, E. V. & Zhang, F. Response to Comment

1033　　　on "RNA-guided DNA insertion with CRISPR-associated transposases". *Science* **368**, 1–2

1034　　　(2020).

1035　35. Tou, C. J., Orr, B. & Kleinstiver, B. P. Precise cut-and-paste DNA insertion using engineered

1036　　　type V-K CRISPR-associated transposases. *Nat. Biotechnol.* (2023) doi:10.1038/s41587-

1037　　　022-01574-x.

1038   36. Park, J. U. *et al.* Structural basis for target site selection in RNA-guided DNA transposition

1039       systems. *Science* **373**, 768–774 (2021).

1040   37. Park, J. U., Tsai, A. W. L., Chen, T. H., Peters, J. E. & Kellogg, E. H. Mechanistic details of

1041       CRISPR-associated transposon recruitment and integration revealed by cryo-EM. *Proc.*

1042       *Natl. Acad. Sci. U. S. A.* **119**, 1–9 (2022).

1043   38. Querques, I., Schmitz, M., Oberli, S., Chanez, C. & Jinek, M. Target site selection and

1044       remodelling by type V CRISPR-transposon systems. *Nature* (2021) doi:10.1038/s41586-

1045       021-04030-z.

1046   39. Park, J. U. *et al.* Structures of the holo CRISPR RNA-guided transposon integration

1047       complex. *Nature* **613**, 775–782 (2023).

1048   40. Hoffmann, F. T. *et al.* Selective TnsC recruitment enhances the fidelity of RNA-guided

1049       transposition. *Nature* **609**, 384–393 (2022).

1050   41. Jia, N., Xie, W., de la Cruz, M. J., Eng, E. T. & Patel, D. J. Structure–function insights into

1051       the initial step of DNA integration by a CRISPR–Cas–Transposon complex. *Cell Res.* **30**,

1052       182–184 (2020).

1053   42. Wang, B., Xu, W. & Yang, H. Structural basis of a Tn7-like transposase recruitment and

1054       DNA loading to CRISPR-Cas surveillance complex. *Cell Res.* **30**, 185–187 (2020).

1055   43. Li, Z., Zhang, H., Xiao, R. & Chang, L. Cryo-EM structure of a type I-F CRISPR RNA guided

1056       surveillance complex bound to transposition protein TniQ. *Cell Res.* **30**, 179–181 (2020).

1057   44. Zhong, E. D., Bepler, T., Berger, B. & Davis, J. H. CryoDRGN: reconstruction of

1058       heterogeneous cryo-EM structures using neural networks. *Nat. Methods* **18**, 176–185

1059       (2021).

1060   45. Moreb, E. A., Hutmacher, M. & Lynch, M. D. CRISPR-Cas 'non-Target' Sites Inhibit On-

1061       Target Cutting Rates. *CRISPR J.* **3**, 550–561 (2020).

1062   46. Tuminauskaite, D. *et al.* DNA interference is controlled by R-loop length in a type I-F1

1063       CRISPR-Cas system. *BMC Biol.* **18**, 1–16 (2020).

1064    47. Rollins, M. C. F. et al. Structure Reveals a Mechanism of CRISPR-RNA-Guided Nuclease

1065         Recruitment and Anti-CRISPR Viral Mimicry. *Mol. Cell* **74**, 132-142.e5 (2019).

1066    48. Guo, T. W. et al. Cryo-EM Structures Reveal Mechanism and Inhibition of DNA Targeting by

1067         a CRISPR-Cas Surveillance Complex. *Cell* **171**, 414-426.e12 (2017).

1068    49. Chowdhury, S. et al. Structure Reveals Mechanisms of Viral Suppressors that Intercept a

1069         CRISPR RNA-Guided Surveillance Complex. *Cell* **169**, 47-57.e11 (2017).

1070    50. Wang, Y. et al. A novel strategy to engineer DNA polymerases for enhanced processivity

1071         and improved performance in vitro. *Nucleic Acids Res.* **32**, 1197–1207 (2004).

1072    51. de Vega, M., Lázaro, J. M., Mencía, M., Blanco, L. & Salas, M. Improvement of φ29 DNA

1073         polymerase amplification performance by fusion of DNA binding motifs. *Proc. Natl. Acad.*

1074         *Sci. U. S. A.* **107**, 16506–16511 (2010).

1075    52. Oscorbin, I. P., Wong, P. F., Boyarskikh, U. A., Khrapov, E. A. & Filipenko, M. L. The

1076         attachment of a DNA-binding Sso7d-like protein improves processivity and resistance to

1077         inhibitors of M-MuLV reverse transcriptase. *FEBS Lett.* **594**, 4338–4356 (2020).

1078    53. Tong, C. L., Kanwar, N., Morrone, D. J. & Seelig, B. Nature-inspired engineering of an

1079         artificial ligase enzyme by domain fusion. *Nucleic Acids Res.* **50**, 11175–11185 (2022).

1080    54. Jackson, R. N. et al. Crystal structure of the CRISPR RNA–guided surveillance complex

1081         from Escherichia coli. *Science* **345**, 1473–1479 (2014).

1082    55. Xue, C., Zhu, Y., Zhang, X., Shin, Y. K. & Sashital, D. G. Real-Time Observation of Target

1083         Search by the CRISPR Surveillance Complex Cascade. *Cell Rep.* **21**, 3717–3727 (2017).

1084    56. Aldag, P. et al. Dynamic interplay between target search and recognition for a Type I

1085         CRISPR-Cas system. *Nat. Commun.* **14**, (2023).

1086    57. Qi, L. S. et al. Repurposing CRISPR as an RNA-γuided platform for sequence-specific

1087         control of gene expression. *Cell* **152**, 1173–1183 (2013).

1088    58. Hoffmann, F. T. et al. Selective recruitment of the AAA + ATPase TnsC increases the fidelity

1089         of Type I-F CRISPR RNA-guided transposition. *Manuscr. Revis.* 1–60 (2021).

1090   59. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**,

1091        583–589 (2021).

1092   60. Skelding, Z., Sarnovsky, R. & Craig, N. L. Formation of a nucleoprotein complex containing

1093        Tn7 and its target DNA regulates transposition initiation. *EMBO J.* **21**, 3494–3504 (2002).

1094   61. Zhang, F., Saito, M. & Faure, G. Type I-B CRISPR-Associated Transposase Systems.

1095        (2024).

1096   62. Metagenomi Technologies. *S-1*. (2024).

1097   63. Strecker, J., Zhang, F. & Ladha, A. CRISPR-associated transposase systems and methods

1098        of use thereof.

1099   64. Jamali, K. *et al.* Automated model building and protein identification in cryo-EM maps.

1100        *Nature* (2024) doi:10.1038/s41586-024-07215-4.

1101   65. Ondov, B. D., Bergman, N. H. & Phillippy, A. M. Krona-385.pdf. *BMC Bioinformatics* **385**,

1102        (2011).

1103   66. Kalichuk, V. *et al.* The archaeal "7 kDa DNA-binding" proteins: extended characterization of

1104        an old gifted family. *Sci. Rep.* **6**, 37274 (2016).

1105   67. Shao, X. Common fold in helix-hairpin-helix proteins. *Nucleic Acids Res.* **28**, 2643–2650

1106        (2000).

1107