

1 Evolution and classification of the CRISPR/Cas systems

2
3
4 Kira S. Makarova¹, Daniel H. Haft², Rodolphe Barrangou³, Stan J. J. Brouns⁴,
5 Emmanuelle Charpentier⁵, Philippe Horvath⁶, Sylvain Moineau⁷, Francisco J. M.
6 Mojica⁸, Yuri I. Wolf¹, Alexander F. Yakunin⁹, John van der Oost⁴, and Eugene V.
7 Koonin^{1,*}

8
9 ¹National Center for Biotechnology Information, National Library of Medicine, National Institutes of
10 Health, 8600 Rockville Pike, Bethesda, MD 20894

11
12 ²The J Craig Venter Institute, 9704 Medical Center Drive, Rockville, MD 20850

13
14 ³Danisco USA Inc., 3329 Agriculture Drive, Madison, WI 53716

15
16 ⁴Laboratory of Microbiology, Wageningen University, Dreijenplein 10, 6703 HB Wageningen, The
17 Netherlands

18
19 ⁵Laboratory for Molecular Infection Medicine Sweden, Umeå Centre for Microbial Research, Department
20 of Molecular Biology, Umeå University, S-90187 Umeå, Sweden

21
22 ⁶Danisco France SAS, BP10, 86220 Dangé-Saint-Romain, France

23
24 ⁷Département de Biochimie, Microbiologie et Bio-informatique, Faculté des Sciences et de Génie, Groupe
25 de Recherche en Ecologie Buccale, Université Laval, Quebec City, Quebec, G1V 0A6, Canada

26
27 ⁸Departamento de Fisiología, Genética y Microbiología, Universidad de Alicante, 03080-Alicante, Spain

28
29 ⁹Banting and Best Department of Medical Research, University of Toronto, Toronto, Ontario, Canada

30
31 *for correspondence: koonin@ncbi.nlm.nih.gov

32 **Abstract**

33

34 The CRISPR/Cas modules are adaptive immunity systems that are widespread in
35 Archaea and Bacteria. These defense systems show extraordinary diversity in operon
36 architecture along with high rates of *cas* gene evolution. Several classifications of
37 CRISPR/Cas systems and Cas proteins have been proposed but it appears difficult to
38 capture the full complexity of these systems in a coherent scheme. Here, we provide an
39 updated analysis of the evolutionary relationships between CRISPR/Cas systems and Cas
40 proteins. Three major types of CRISPR/Cas systems are delineated, with a further
41 division into several subtypes and a few chimeric variants. Given the complexity of
42 genomic architectures and the extremely dynamic evolution of the CRISPR/Cas systems,
43 a unified classification of these systems should be based on multiple criteria.
44 Accordingly, a “polythetic” classification is proposed that integrates the phylogenies of
45 the most common *cas* genes, the sequence and arrangements of the CRISPR repeats, and
46 the architecture of the CRISPR/Cas loci.

47

48 The CRISPR/Cas (clustered regularly interspaced short palindromic repeats /
49 CRISPR-associated proteins) modules are adaptive immunity systems that act against
50 invading genetic elements and are encoded by most Archaea and many Bacteria ¹⁻⁶
51 (Supplementary Table 1). Distinct arrays of short repeats interspersed with unique
52 spacers (CRISPR) have been recognized in bacterial and archaeal genomes for years, and
53 although it has been proposed that these repeat arrays could have an important common
54 function ⁷, the nature of that function has been elucidated only much later.
55 Independently, Cas protein sequences encoded by putative operons adjacent to CRISPR
56 were analyzed in detail with computational methods and found to contain domains
57 characteristic of several nucleases, a helicase, a polymerase and RNA-binding proteins⁸.
58 Initially, it was speculated that these proteins might jointly constitute a novel DNA repair
59 system ⁹. However, the observation that some of the unique CRISPR spacers were
60 (nearly) identical to fragments of virus and plasmid genes led to the hypothesis that
61 CRISPR/Cas might be involved in defense against selfish elements ¹⁰⁻¹². On the basis of
62 these findings combined with a comprehensive computational re-analysis of the Cas
63 proteins ^{13,14}, a model was proposed ¹⁴ that drew an analogy between the putative novel
64 prokaryotic defense system and the eukaryotic RNA interference (RNAi) mechanisms¹⁵.
65 However, unlike the eukaryotic RNAi systems, the CRISPR/Cas system integrates a
66 small piece of DNA derived from foreign nucleic acid into the CRISPR locus of the host
67 genome as the first step in the series of events that leads to immunity to the given agent
68 ¹⁴. The hypothesis that the CRISPR/Cas system plays a role in defense against invading
69 DNA has been validated by the demonstration that integration of a short bacteriophage-
70 specific sequence into the CRISPR locus of the lactic acid bacterium *Streptococcus*

71 *thermophilus* conferred resistance to the cognate phage¹⁶. The phage resistance was
72 abrogated by even a single mismatch between the CRISPR insert (referred to as spacer)
73 and the target phage sequence (referred to as proto-spacer)¹⁶.

74

75 The CRISPR/Cas systems mediate immunity to invading genetic elements via
76 three stages (Figure 1)^{1-3,6}. The first stage is adaptation which involves the integration of
77 short pieces of DNA homologous to virus or plasmid sequences into the CRISPR loci¹⁶⁻
78 ¹⁸. Viral challenge typically triggers insertion of a single virus-derived, resistance-
79 conferring spacer with a characteristic length of approximately 30 bp; acquisition of
80 multiple spacers from the same phage is less frequent. The insertion of new spacers has
81 been reported to depend on short (few nucleotides) proto-spacer adjacent motifs (PAMs)
82 which differ between variants of the CRISPR/Cas system and appear to determine the
83 selection of the inserted spacer^{19,20}.

84

85 The second stage is the expression and processing, during which the long primary
86 transcript of a CRISPR locus (pre-crRNA) is processed into short crRNAs . The latter
87 processing step is catalyzed by endoribonucleases that either operate as a subunit of a
88 larger complex (e.g. Cascade, CRISPR-associated complex for antiviral defense in
89 *Escherichia coli*) (Figure 1) or as a stand alone enzyme (e.g., Cas6 in the archaeon
90 *Pyrococcus furiosus*). In the case of the Cascade complex²¹, the mature crRNA remains
91 associated with the complex after initial endonuclease cleavage (Figure 1), whereas in *P.*
92 *furiosus* the crRNA, processed by Cas6, is passed on to a distinct Cas protein complex
93 (Cmr-type, see below) where it is processed further^{22,23}.

94 The third and final step is interference, during which the alien nucleic acid
95 (foreign DNA or RNA) is targeted and cleaved within the proto-spacer sequence^{6,17,18}.
96 The crRNAs guide the respective complexes of Cas proteins to the complementary target
97 sequences of invading viruses or plasmids that match the spacer. In *S. thermophilus* and
98 *E. coli* targeting either strand of the phage DNA confers immunity to the cognate phage,
99 an observation that is best compatible with DNA being the target^{16,24}. Furthermore,
100 insertion of a self-splicing intron into the proto-spacer sequence of the target gene
101 rendered the respective plasmid resistant to the CRISPR-mediated immunity in
102 *Staphylococcus epidermidis*, indicating that the invading DNA rather than mRNA is
103 targeted²⁵. However, *in vitro* experiments with the CRISPR/Cas system from *P. furiosus*
104 showed that in this case the crRNA rather targets the viral mRNA²³. These findings
105 emphasize the remarkable mechanistic and functional diversity of the CRISPR/Cas
106 systems, although the full range of their activities remains to be determined. Various Cas
107 proteins might participate in either one or multiple stages of the CRISPR/Cas system
108 action, most likely, as protein complexes⁶.

109 In agreement with the bioinformatic predictions, nuclease activities, RNase
110 and/or DNase, have been demonstrated for several Cas proteins, including the two
111 universal core Cas proteins: Cas1 (a metal-dependent DNase with no sequence
112 specificity that has been proposed to be involved in the integration of the alien DNA
113 (spacer) into the CRISPR cassettes²⁶) and Cas2 (a metal-dependent endoribonuclease
114 whose role in the CRISPR/Cas mechanism remains unclear²⁷). The Cas proteins known
115 as RAMPs (Repeat-Associated Mysterious Proteins) contain a double ferredoxin-fold

116 domain; some of the RAMPs have been shown to possess sequence- or structure-specific
117 RNase activity that is involved in the processing of pre-crRNA transcripts^{21,22,24}.

118

119 The CRISPR/Cas systems can be divided into two distinct, quasi-independent
120 subsystems. The highly conserved “information processing” subsystem includes the Cas1
121 and Cas2 proteins in its core, and is thought to be involved in the maintenance and
122 replenishment of the spacer-repeat library. The “executive” subsystem, which is highly
123 variable in content, typically includes multiple RAMP proteins, and is involved in
124 processing of the CRISPR transcript (Cascade complex and its analogs) and the crRNA-
125 directed interference of invading genetic elements.

126

127 Extensive bioinformatic analyses have shown that the genomes of various
128 CRISPR-containing organisms encode approximately 65 distinct Cas proteins which can
129 be classified into 23 to 45 families depending on the classification criteria (granularity of
130 clustering)^{13,14}. Furthermore, 8 distinct subtypes of the CRISPR/Cas systems (CASS1 to
131 CASS8) have been delineated on the basis of the composition and architecture of the *cas*
132 operons and Cas1 phylogeny^{13,14}.

133 The diversity of CRISPR systems identified in newly sequenced genomes (in a
134 representative set of 703 archaeal and bacterial genomes, 310 (44%) encode one or more
135 CRISPR/Cas modules; Table 2 and Supplementary Table 1) is rapidly increasing^{1,4},
136 hence an urgent need exists for a rational and unified classification and nomenclature of
137 the *cas* genes. In this article, we summarize the shortcomings of the existing
138 classifications and nomenclatures of the CRISPR/Cas systems, and propose a new,

139 “polythetic” classification which combines information from comparative-genomic and
140 phylogenetic analyses.

141

142 **Problem with the existing CRISPR/Cas classification**

143 The original, widely used classification proposed by Haft et al. was based on the
144 topology of the Cas1 phylogenetic tree and *cas* operon organization in eight organisms¹³.

145 The names of four core *cas* genes were kept as they were originally proposed by Jansen
146 et al. in 2002⁸. Two other core genes, *cas5* and *cas6*, were then added using the same
147 principle¹³. In addition, gene names for proteins specific to each of the eight CRISPR
148 systems were proposed. For example, the unique genes found in the *E. coli* system were
149 denoted *cse1* (CRISPR system of *E. coli* gene number 1), *cse2*, *cse3*, *cse4*, and *cse5*
150 (elsewhere, these *E. coli* genes were also labeled as *casA*, *casB*, *casC*, *casD*, and *casE*)²¹.

151 Although the original approach¹³ offers attractive simplicity, its major
152 shortcoming is the failure to identify distant relationships between many proteins. For
153 example, the proteins that belong to the COG1857 family, that are present in the majority
154 of the CRISPR/Cas systems and are obviously orthologous, have been given at least 5
155 different names: *Cse3*, *Csd2*, *Csh2*, *Cst2*, and *Csa2* (Table 1). The other shortcoming of
156 this classification is that it does not take into account the complexity of the relationships
157 between the CRISPR/Cas systems and the respective Bacteria and Archaea. In particular,
158 the classification ignores the apparent relatedness between several CRISPR/Cas systems:
159 for example, the *E. coli* and *Y. pestis* systems are definitely related, and so are the *A. pernix*,
160 *T. neapoliensis*/Hmari and *D. vulgaris* systems which share a common signature gene of the BH0338
161 family. Conversely, extensive recombination within CRISPR operons has resulted in

162 hybrid CRISPR/Cas systems that cannot be assigned to any of the proposed groups
163 although they contain typical *cas* genes. The linkage between CRISPR/Cas groups and
164 particular organisms can be misleading due to the presence of multiple CRISPR/Cas
165 systems in the same organism, the presence of different CRISPR/Cas systems in different
166 strains of same species, and again, the rather wide spread of hybrid systems. The
167 inconsistencies between the nomenclatures of the CRISPR/Cas systems and the names of
168 Cas proteins are rapidly growing. In particular, many of these proteins are currently
169 classified into families that do not have systematic names pointing to their involvement
170 with CRISPR/Cas (e.g., CXXC_CXXC family protein, GSU0053 family protein, etc.).

171

172 **Evolutionary relationships as a basis for a new classification of CRISPR/Cas**
173 **systems**

174 Two Cas proteins (Cas1 and Cas2) are present in all CRISPR/Cas systems that are
175 predicted to be functionally active and are thought to be involved in spacer integration
176 (the adaptation stage) as the “information processing” subsystem. The *cas1* and *cas2*
177 genes comprise the cores of three distinct types (I, II and III) of CRISPR/Cas systems that
178 form the basis of a new, polythetic (based on multiple criteria) classification we propose
179 here (Figure 2 and Table 1).

180

181 ***Type I CRISPR/Cas system.*** Typical Type I loci contain a gene for a predicted
182 helicase/nuclease (Cas3) as well as several other proteins that probably form Cascade-
183 like complexes with different compositions^{21,24}. These complexes include multiple
184 proteins of the RAMP superfamily, in particular, the widespread COG1857 (Cas7) family

185 recently proposed to adopt the RAMP fold (KSM, unpublished observations), and
186 BH0338-like families; in addition, the complexes may contain other, less conserved
187 subunits. In the Cascade complex, a RAMP protein has been demonstrated to be the
188 major enzyme (RNA endonuclease) that catalyzes the processing of the long
189 spacer/repeat-containing transcript into a mature crRNA^{21,24}. In most cases, the catalytic
190 RAMP protein (Cas6, Cas6e and Cas6f; see Table 1) does not belong to the most
191 prevalent Cas5 or Cas7 family of RAMPs and is often encoded in the periphery of the
192 respective operon. However, an exception might be the subtype I-C system (also known
193 as Dvulg/CASS1, Table 1 and Figure 2) in which either Cas5 or Cas7 are likely
194 candidates to possess RNase activity. The Type I CRISPR systems appear to target
195 DNA, and the cleavage might be catalyzed by the HD nuclease domains of Cas3 and/or
196 by the RecB-family nucleases (Cas4). However, the fact that in several Type I
197 CRISPR/Cas systems the RecB nuclease domain is fused to Cas1, may suggest a role for
198 Cas4 in spacer acquisition.

199

200 ***Type II CRISPR/Cas system.*** The Type II systems include the “HNH” type
201 system (*Streptococcus*-like, also known as Nmeni subtype or CASS4, Table 3) in which
202 (in addition to the ubiquitous Cas1 and Cas2) a single, very large protein, Cas9
203 (COG3513), appears to be sufficient for both generating crRNA and cleaving the target
204 DNA. The Cas9 protein (~1,000 amino acids) contains at least two nuclease domains,
205 namely, the N-terminal RuvC-like nuclease (RNase H fold) and the HNH (McrA-like)
206 nuclease domain that is located in the middle of the protein. The functional specialization
207 of these nuclease domains remains to be elucidated. However, the HNH nuclease domain

208 is abundant in restriction enzymes and possesses endonuclease activity^{28,29}, so it is likely
209 to be responsible for target cleavage.

210

211 ***Type III CRISPR/Cas system.*** The Type III systems contain polymerase-RAMP
212 modules in which at least some of the RAMPs appear to be involved in the processing of
213 the CRISPR-spacer transcripts analogously to the Cascade complex. Targeting of plasmid
214 DNA by this system (subtype III-B) has been demonstrated *in vivo* in *S. epidermidis*²⁵,
215 and it seems plausible that the HD domain of the polymerase-like protein (COG1353) is
216 involved in the target DNA cleavage. There is also strong evidence that at least *in vitro*
217 the Type III-A CRISPR/Cas system from *P. furiosus* can target RNA²³. Apart from the
218 universal Cas2 protein, the only identified ribonucleases in the Type III CRISPR/Cas
219 systems are RAMP proteins. Beside Cas6 that is involved in CRISPR transcript
220 processing, Type III systems contain at least two additional RAMPs. These RAMPs
221 appear to be the most likely candidate enzymes for the subsequent trimming of the
222 crRNA. In many organisms, Type III CRISPR/Cas operons lack the *cas1-cas2* gene pair;
223 in all these cases, an additional CRISPR locus (Type I or II) is present in the respective
224 genome, so the polymerase-RAMP module probably interacts with Cas1-Cas2 *in trans*.
225 In other organisms, the polymerase-RAMP modules are present in a single operon with
226 *cas1* and *cas2*, e.g., a module with the typical architecture in *Staphylococcus epidermidis*
227 and *Mycobacterium tuberculosis* (Type III-A) and a distinct version in *Halorhodospira*
228 *halophila* (Type III-B). In these organisms the Type III operon is the only CRISPR/Cas
229 locus, suggesting that, combined with Cas1-Cas2, the polymerase-RAMP module forms a
230 fully functional, autonomous Type III system.

231

232 The three types of CRISPR systems show a distinctly non-uniform distribution
233 among the major lineages of Archaea and Bacteria (Table 2). In particular, the Type II
234 systems so far have been found exclusively in bacteria whereas Type III systems are
235 more common in Archaea. The previously observed trend of overrepresentation of
236 CRISPR in Archaea compared to Bacteria³⁰ still holds (Table 2). Moreover, the majority
237 of archaeal genomes carry more than one CRISPR/Cas system; typically, different
238 modules within the same genome are unrelated.

239

240 **Subtypes of the CRISPR/Cas systems and their evolution**

241 Based on the gene composition and architecture of the respective *cas* operons, the
242 three basic types of CRISPR/Cas systems can be further classified into subtypes that
243 largely agree with the previously delineated variants^{13, 14}. Each of the subtypes contains a
244 distinct signature gene (Figure 2 and Table 1). The ubiquitous, highly conserved Cas1
245 protein can be used as a scaffold to investigate the evolution of the CRISPR/Cas system
246 (the other universal protein, Cas2, is too small to yield a well-resolved tree). The
247 phylogenetic tree of Cas1 includes several well-resolved branches that generally agree
248 with the classification of CRISPR/Cas systems into subtypes (I-A, I-B, I-C, I-E, I-F, II
249 and III-A)¹⁴, with a few notable exceptions (Figure 3). In particular, Cas1 proteins
250 associated with the polymerase-RAMP module (type III) appear in several unrelated
251 positions in the tree (Figure 3), suggesting that this module can operate with a variety of
252 *cas1-cas2* genes both *in cis* and *in trans*.

253

254 The CRISPR repeats can be classified into at least 12 groups based on sequence
255 similarity³¹. Four groups of CRISPR repeats clearly correspond to distinct CRISPR/Cas
256 subtypes (all of Type I): group 2 – I-E (Ecoli/CASS2), group 3 – I-C (Dvulg/CASS1),
257 group 4 – I-F (Ypest/CASS3) and group 10 – type II (Nmeni/CASS4). These four
258 variants of CRISPR/Cas systems have the most stable operon organizations; on the other
259 hand, subtypes I-A, I-B, I-D, and Type III appear to be prone to recombination between
260 different subtypes. Structural characteristics of CRISPR repeats of these four groups can
261 be potentially employed for classification, in addition to phylogenetic data and signature
262 genes. Other groups of repeats cannot be unequivocally associated with particular
263 CRISPR/Cas system subtypes.

264 Integration of all the above considerations in a dendrogram, reflects our present
265 understanding of the evolutionary history of CRISPR/Cas systems (Figure 2). Subtypes
266 of the Type I system are grouped according to their operon organizations and the
267 phylogeny of their Cas1 proteins.

268

269 **Proposals for the CRISPR/Cas system nomenclature**

270 Most of the CRISPR/Cas loci can be readily classified into the proposed types and
271 subtypes based on the presence of type- and subtype-specific signature genes (Table 1
272 and Table 3). In addition, we introduce the catch-all subtypes I-U, II-U and III-U (U for
273 unclassified) for systems that lack currently defined subtype-specific signature genes but
274 might fit one of the established subtypes based on further structure and sequence analysis
275 or potentially could become founders of new subtypes. In the same vein, we propose

276 Type U for the loci that cannot be classified even at the type level (e.g., the CRISPR/Cas
277 system in *Acidithiobacillus ferrooxidans* ATCC 23270 discussed below).

278

279 We propose to retain the well-established names for core genes of the information
280 processing (ubiquitous *cas1* and *cas2*) and executive (*cas3-6*, characteristic for Type I
281 system) components of the CRISPR/Cas systems. In several cases for which orthology
282 could be confidently traced, we extend the usage of these six *cas* gene names (thus, *cmx5*
283 of I-C is renamed *cas5* and *cmx6* is renamed *cas6*). In cases when significant sequence
284 similarity between Cas proteins is observed but orthologous relationships cannot be
285 definitively assigned, we use an additional letter derived from the subtype label (hence
286 *cas6e* and *cas6f*, former *cse3* and *csy4*, respectively, which are likely to be extremely
287 divergent derivatives of *cas6*; Table 1).

288

289 In Type I systems, there are two additional genes for which orthology is readily
290 detectable between different subtypes. We refer to these genes as *cas7* and *cas8(abc)*;
291 both encode subunits of the Cascade complex (Table 1). The *cas8a*, *cas8b* and *cas8c*
292 genes are the signature genes for subtypes I-A, I-B and I-C, respectively. In type II and
293 type III systems, the respective signature genes are designated *cas9* (formerly *csn1* and
294 *csx12*) and *cas10* (*cmr2*, *csm1* and *csx11*).

295

296 When a gene clearly is a fusion or fission of established genes, we propose an *ad*
297 *hoc* nomenclature indicating the relationship of this variant to the “canonical” forms:

298 thus, *cas2/cas3* in I-F systems denotes a fusion of *cas2* and *cas3*, whereas *cas3'* and
299 *cas3''* denote the solo helicase domain and the solo HD domain, respectively.

300

301 For other, less common genes that have been named previously¹³, the “legacy”
302 nomenclature can be used along with the family classification given in Table 1 (eg. *csx9*).
303 Due to the fact that many of the Cas protein sequences are highly diverged, it is expected
304 that with the increasing representation of sequences and structures, many of these genes
305 eventually will be included into existing families. We propose to continue assigning
306 further “numerical” names to newly merged orthologous families in the future (*cas11*,
307 *cas12* etc.).

308

309 For the remaining CRISPR-associated genes, we propose to assign interim gene
310 names (*csx1*, where “x” stands for unclassified family), with an indication of the
311 (super)family where known (e.g. *csx1*, COG1517 family or *csx10*, RAMP superfamily).

312

313 **Outstanding problems**

314 ***Subtype assignment.*** As pointed out above, the phylogenetic tree of Cas1
315 reproduces most of the previously established groups fairly well, with the exception of
316 the Type III systems (Figure 3). However, for the deep branches, assigning a subtype can
317 be problematic. In many cases, detailed analysis of the gene orders reveals a more
318 complicated picture, with different arrangements of *cas* genes in the operons, apparently,
319 due to frequent horizontal gene transfer (HGT) and recombination that involve the
320 CRISPR loci. In particular, a notable recombinant CRISPR/Cas system is present at least

321 in certain cyanobacteria (e.g., *Synechocystis* sp. PCC 6803: slr7010-ssr7072) and
322 Archaea. In this case, the Type I-C system combines with a distinct polymerase/RAMP
323 module genes in the following arrangement: *cas3*, *cas10* (predicted inactivated
324 polymerase with a HD-domain), *csc2* (COG1337 family, RAMP superfamily), *csc2*
325 (RAMP subfamily), *cas6*, *cas4*, *cas1*, *cas2*. This peculiar hybrid system containing
326 signature genes for both type I and type III systems is represented in a variety of genomes
327 (and thus likely functional), so we introduce it as a new subtype I-D (Figure1).

328

329 Another interesting CRISPR/Cas system typified by *Acidithiobacillus*
330 *ferrooxidans* ATCC 23270 (AFE_1037-AFE_1040) was found so far only in four
331 genomes. This locus appears to possess a distinct gene content and potentially could
332 contribute to our understanding of the functions and evolution of CRISPR/Cas systems in
333 general. This system contains neither of the two ubiquitous core genes (*cas1* or *cas2*) nor
334 any other signature genes of the three CRISPR/Cas types or the 10 subtypes. The *A.*
335 *ferrooxidans* system consists of four genes denoted *csf1*, *csf2*, *csf3* and *csf4*
336 (TIGRFAMS: TIGR03114, TIGR03115, TIGR03116, TIGR03117, respectively). These
337 genes encode, respectively, a Zn-finger domain containing protein, a protein containing
338 two RAMP domains, another distinct RAMP protein and a DinG-like helicase of the
339 XPD family³⁰. According to the CRISPRdb database³², a CRISPR array is present in the
340 vicinity of the above four genes in all of the respective genomes; the architecture of these
341 arrays is unique in each genome. Thus, this system might function in conjunction with
342 different CRISPR arrays and does not require a distinct repeat signature. Indeed, three of
343 the four genomes containing this system additionally possess *cas1* and *cas2* genes that

344 are located in other parts of these genomes and are associated with Type I CRISPR/Cas
345 systems. It remains unclear whether this is a self-sufficient system or rather a defective
346 system that captures and utilizes preexisting CRISPR arrays generated by other, Cas1-
347 containing CRISPR/Cas systems. More data are needed to classify this novel system as a
348 separate CRISPR/Cas type, but this finding illustrates the diversity of CRISPR systems
349 and the challenges associated with their classification.

350

351 ***Gene name assignments.*** Many *cas* genes, in particular genes that encode RAMP
352 proteins, evolve at exceptionally high rates. Gene (protein) family assignment becomes
353 increasingly complicated with the appearance of CRISPR/Cas systems containing genes
354 that encode highly divergent proteins that, after the structure is solved, might (or might
355 not) fall into a known Cas protein family. For example, a CRISPR system very similar to
356 that of subtype I-F (Ypest/CASS3 as determined by Cas1 similarity) is present in
357 *Photobacterium profundum* and several other bacteria. This system includes two proteins,
358 PBPRB1992 and PBPRB1993, that show no significant sequence similarity to any Cas
359 proteins. However, analysis of the sequence motifs that are conserved in these proteins,
360 the predicted secondary structure, as well as the length and position of the corresponding
361 genes in the operon, strongly suggest that they belong to the Cas7 and Cas5 families of
362 RAMPs, respectively. Another example includes the CRISPR/Cas system of *Geobacter*
363 *sulfurreducens* that, according to the phylogeny of Cas1, should be assigned to the I-C
364 subtype. This operon encodes three uncharacterized proteins GSU0052, GSU0053,
365 GSU0054; the last two proteins contain several motifs similar to the characteristic motifs
366 of the RAMP superfamily and thus might be RAMP homologs (Table 1). However, none

367 of these proteins could be linked to known Cas families, even using the most sensitive of
368 the available methods for remote sequence similarity detection³³⁻³⁵. Thus, only
369 comparison of solved structures will shed light on the relationships between these and
370 other highly diverged Cas proteins and known Cas families. In such cases, assignment of
371 new gene names appears to be premature because these proteins are likely to eventually
372 assume already existing names. Therefore, it is proposed that these genes are given
373 temporary “*csx*” names.

374
375 Many CRISPR loci belong to “islands” that contain various “high-mobility”
376 genes such as components of other defense systems, toxin-antitoxins, and transposases
377³⁶. Some of these genes can be erroneously linked to CRISPR/Cas systems, so caution
378 should be exercised in the classification and naming of genes as “*cas*” or even “*csx*”
379 before functional connections with CRISPR/Cas systems are convincingly established.

380
381 An additional challenge to the nomenclature is presented by the variable domain
382 architectures of some of the Cas proteins including domain fusions and fissions as
383 discussed above for the Cas3 protein. Other notable fusions include *cas2-cas3* in the
384 Ypest/CASS3 system, *cas1-cas4* (eg. GSU0057 from *Geobacter sulfurreducens*), *cas1*
385 and DEDDh family exonuclease (eg. LBUL_0800 from *Lactobacillus delbrueckii subsp.*
386 *bulgaricus*), *cas1* and reverse transcriptase fusion (eg. VVA1544 from *Vibrio vulnificus*),
387 and more.

388 In several genomes, homologs of some *cas* genes also appear in contexts
389 different from CRISPR/Cas systems; these proteins might either represent (components
390 of) distinct antiviral defense systems, or they could be involved in other functions such

391 as DNA repair. These proteins include RAMPs of the COG5551 subfamily, COG1517,
392 COG1468 and COG3513 families. In cases like this, classification and labeling of such
393 genes as “*cas*” should be avoided.

394

395 The CRISPR arrays contain few stop codons and accordingly are often
396 erroneously translated into bogus “hypothetical proteins”. Unfortunately, these artifacts
397 then enter the databases and tend to be amplified during the analysis of new genomes, so
398 currently there are at least two Pfam entries each of which consists of non-existent
399 “pseudo-Cas proteins” (pfam11194; pfam11664). Care should be taken during the
400 annotation of new genome sequences to avoid further proliferation of these errors.

401

402 **Conclusion**

403 Given the complexity and the highly dynamic mode of evolution of the CRISPR/Cas
404 systems, it would be counter-productive to attempt classification on the basis of any
405 single criterion, for instance, the phylogeny of Cas1. Thus, we propose here a
406 “polythetic” classification that integrates the phylogenies of the conserved *cas* genes,
407 sequences and structural similarities between other Cas proteins, and the composition and
408 organization of the (putative) operons. It should be emphasized that a robust family
409 classification of the Cas proteins, many of which diverge rapidly, is not only a matter of
410 convenient description, but also a basis for experimental validation of the respective
411 functional predictions. Therefore it is important that this classification is continuously
412 updated and revised when necessary, using new sequence and structure information
413 combined with state of the art computational methods. The classification described in this

414 article is available at the NCBI website along with tools for identification of Cas proteins
415 (ftp://ftp.ncbi.nih.gov/pub/wolf/_suppl/CRISPRclass/index.html). In the future, a fine-
416 grain classification of the CRISPR/Cas systems should become feasible on the basis of
417 phylogenies and structures of Cas proteins, *cas* operon organizations of *cas* genes, and
418 CRISPR repeat architectures.

419

420 **Acknowledgements**

421 The authors thank Michael Terns for critical reading of the manuscript and useful
422 discussions. KSM, YIW and EVK are supported by the intramural funds of the US
423 Department of Health and Human Services (National Library of Medicine); E.C.
424 acknowledges funding from Umeå University and the Swedish Research Councils.S.M.
425 acknowledges funding from NSERC of Canada (Discovery program); FJMM
426 acknowledges support from the University of Alicante (Vicerrectorado de Investigacion,
427 Desarrollo e Innovacion) for the use of its research technical services; AFY is supported
428 by Genome Canada (through the Ontario Genomics Institute; 2009-OGI-ABC-1405).

429

430

431 **References**

432

433

- 434 1 Deveau, H., Garneau, J. E. & Moineau, S. CRISPR/Cas System and Its Role in
435 Phage-Bacteria Interactions. *Annu Rev Microbiol* **64**, 475-493 (2010).
- 436 2 Horvath, P. & Barrangou, R. CRISPR/Cas, the immune system of bacteria and
437 archaea. *Science* **327**, 167-170 (2010).
- 438 3 Karginov, F. V. & Hannon, G. J. The CRISPR system: small RNA-guided
439 defense in bacteria and archaea. *Mol Cell* **37**, 7-19 (2010).
- 440 4 Koonin, E. V. & Makarova, K. S. CRISPR-Cas: an adaptive immunity system in
441 prokaryotes. *F1000 Biol Rep* **1**, 95 (2009).
- 442 5 Sorek, R., Kunin, V. & Hugenholz, P. CRISPR--a widespread system that
443 provides acquired resistance against phages in bacteria and archaea. *Nat Rev*
444 *Microbiol* **6**, 181-186 (2008).
- 445 6 van der Oost, J., Jore, M. M., Westra, E. R., Lundgren, M. & Brouns, S. J.
446 CRISPR-based adaptive and heritable immunity in prokaryotes. *Trends Biochem*
447 *Sci* **34**, 401-407 (2009).
- 448 7 Mojica, F. J., Diez-Villasenor, C., Soria, E. & Juez, G. Biological significance of
449 a family of regularly spaced repeats in the genomes of Archaea, Bacteria and
450 mitochondria. *Mol Microbiol* **36**, 244-246 (2000).
- 451 8 Jansen, R., Embden, J. D., Gastra, W. & Schouls, L. M. Identification of genes
452 that are associated with DNA repeats in prokaryotes. *Mol Microbiol* **43**, 1565-
453 1575 (2002).
- 454 9 Makarova, K. S., Aravind, L., Grishin, N. V., Rogozin, I. B. & Koonin, E. V. A
455 DNA repair system specific for thermophilic Archaea and bacteria predicted by
456 genomic context analysis. *Nucleic Acids Res* **30**, 482-496 (2002).
- 457 10 Mojica, F. J., Diez-Villasenor, C., Garcia-Martinez, J. & Soria, E. Intervening
458 sequences of regularly spaced prokaryotic repeats derive from foreign genetic
459 elements. *J Mol Evol* **60**, 174-182 (2005).
- 460 11 Bolotin, A., Quinquis, B., Sorokin, A. & Ehrlich, S. D. Clustered regularly
461 interspaced short palindrome repeats (CRISPRs) have spacers of
462 extrachromosomal origin. *Microbiology* **151**, 2551-2561 (2005).
- 463 12 Pourcel, C., Salvignol, G. & Vergnaud, G. CRISPR elements in *Yersinia pestis*
464 acquire new repeats by preferential uptake of bacteriophage DNA, and provide
465 additional tools for evolutionary studies. *Microbiology* **151**, 653-663 (2005).
- 466 13 Haft, D. H., Selengut, J., Mongodin, E. F. & Nelson, K. E. A guild of 45
467 CRISPR-associated (Cas) protein families and multiple CRISPR/Cas subtypes
468 exist in prokaryotic genomes. *PLoS Comput Biol* **1**, e60 (2005).
- 469 14 Makarova, K. S., Grishin, N. V., Shabalina, S. A., Wolf, Y. I. & Koonin, E. V. A
470 putative RNA-interference-based immune system in prokaryotes: computational
471 analysis of the predicted enzymatic machinery, functional analogies with
472 eukaryotic RNAi, and hypothetical mechanisms of action. *Biol Direct* **1**, 7 (2006).
- 473 15 Carthew, R. W. & Sontheimer, E. J. Origins and Mechanisms of miRNAs and
474 siRNAs. *Cell* **136**, 642-655 (2009).

- 475 16 Barrangou, R. *et al.* CRISPR provides acquired resistance against viruses in
476 prokaryotes. *Science* **315**, 1709-1712 (2007).
- 477 17 Garneau, J. E. *et al.* The CRISPR/Cas bacterial immune system cleaves
478 bacteriophage and plasmid DNA. *Nature* **468**, 67-71 (2010).
- 479 18 Sontheimer, E. J. & Marraffini, L. A. Microbiology: slicer for DNA. *Nature* **468**,
480 45-46 (2010).
- 481 19 Mojica, F. J., Diez-Villasenor, C., Garcia-Martinez, J. & Almendros, C. Short
482 motif sequences determine the targets of the prokaryotic CRISPR defence system.
483 *Microbiology* **155**, 733-740 (2009).
- 484 20 Marraffini, L. A. & Sontheimer, E. J. Self versus non-self discrimination during
485 CRISPR RNA-directed immunity. *Nature* **463**, 568-571 (2010).
- 486 21 Brouns, S. J. *et al.* Small CRISPR RNAs guide antiviral defense in prokaryotes.
487 *Science* **321**, 960-964 (2008).
- 488 22 Carte, J., Wang, R., Li, H., Terns, R. M. & Terns, M. P. Cas6 is an
489 endoribonuclease that generates guide RNAs for invader defense in prokaryotes.
490 *Genes Dev* **22**, 3489-3496 (2008).
- 491 23 Hale, C. R. *et al.* RNA-guided RNA cleavage by a CRISPR RNA-Cas protein
492 complex. *Cell* **139**, 945-956 (2009).
- 493 24 Haurwitz, R. E., Jinek, M., Wiedenheft, B., Zhou, K. & Doudna, J. A. Sequence-
494 and structure-specific RNA processing by a CRISPR endonuclease. *Science* **329**,
495 1355-1358 (2010).
- 496 25 Marraffini, L. A. & Sontheimer, E. J. CRISPR interference limits horizontal gene
497 transfer in staphylococci by targeting DNA. *Science* **322**, 1843-1845 (2008).
- 498 26 Wiedenheft, B. *et al.* Structural basis for DNase activity of a conserved protein
499 implicated in CRISPR-mediated genome defense. *Structure* **17**, 904-912 (2009).
- 500 27 Beloglazova, N. *et al.* A novel family of sequence-specific endoribonucleases
501 associated with the clustered regularly interspaced short palindromic repeats. *J*
502 *Biol Chem* **283**, 20361-20371 (2008).
- 503 28 Kleanthous, C. *et al.* Structural and mechanistic basis of immunity toward
504 endonuclease colicins. *Nat Struct Biol* **6**, 243-252 (1999).
- 505 29 Jakubauskas, A., Giedriene, J., Bujnicki, J. M. & Janulaitis, A. Identification of a
506 single HNH active site in type IIS restriction endonuclease Eco31I. *J Mol Biol*
507 **370**, 157-169 (2007).
- 508 30 White, M. F. Structure, function and evolution of the XPD family of iron-sulfur-
509 containing 5'-->3' DNA helicases. *Biochem Soc Trans* **37**, 547-551 (2009).
- 510 31 Kunin, V., Sorek, R. & Hugenholtz, P. Evolutionary conservation of sequence and
511 secondary structures in CRISPR repeats. *Genome Biol* **8**, R61 (2007).
- 512 32 Grissa, I., Vergnaud, G. & Pourcel, C. The CRISPRdb database and tools to
513 display CRISPRs and to generate dictionaries of spacers and repeats. *BMC*
514 *Bioinformatics* **8**, 172 (2007).
- 515 33 Altschul, S. F. & Koonin, E. V. PSI-BLAST - a tool for making discoveries in
516 sequence databases. *Trends Biochem Sci.* **23**, 444-447 (1998).
- 517 34 Marchler-Bauer, A. & Bryant, S. H. CD-Search: protein domain annotations on
518 the fly. *Nucleic Acids Res* **32**, W327-331 (2004).

519 35 Soding, J., Remmert, M., Biegert, A. & Lupas, A. N. HHsenser: exhaustive
520 transitive profile search using HMM-HMM comparison. *Nucleic Acids Res* **34**,
521 W374-378 (2006).

522 36 Makarova, K. S., Wolf, Y. I., van der Oost, J. & Koonin, E. V. Prokaryotic
523 homologs of Argonaute proteins are predicted to function as key components of a
524 novel system of defense against mobile genetic elements. *Biol Direct* **4**, 29,
525 (2009).

526 37 Deveau, H. *et al.* Phage response to CRISPR-encoded resistance in *Streptococcus*
527 *thermophilus*. *J Bacteriol* **190**, 1390-1400 (2008).

528 38 Guindon, S. & Gascuel, O. A simple, fast, and accurate algorithm to estimate
529 large phylogenies by maximum likelihood. *Syst Biol* **52**, 696-704 (2003).

530 39 Han, D., Lehmann, K. & Krauss, G. SSO1450--a CAS1 protein from *Sulfolobus*
531 *solfatarius* P2 with high affinity for RNA and DNA. *FEBS Lett* **583**, 1928-1932
532 (2009).

533 40 Han, D. & Krauss, G. Characterization of the endonuclease SSO2001 from
534 *Sulfolobus solfataricus* P2. *FEBS Lett* **583**, 771-776 (2009).

535 41 Guy, C. P., Majernik, A. I., Chong, J. P. & Bolt, E. L. A novel nuclease-ATPase
536 (Nar71) from archaea is part of a proposed thermophilic DNA repair system.
537 *Nucleic Acids Res* **32**, 6176-6186 (2004).

538 42 Selengut, J. D. *et al.* TIGRFAMs and Genome Properties: tools for the assignment
539 of molecular function and biological process in prokaryotic genomes. *Nucleic*
540 *Acids Res* **35**, D260-264 (2007).

541
542
543
544
545

546 Figure legends

547

548 Figure 1. **The current model of CRISPR/Cas mechanism.**

549 The figure shows a schematic representation of the three stages of the CRISPR/Cas

550 action based on a summary of biochemical and genetic data for the CRISPR/Cas systems

551 from different organisms and some specific details for the *Escherichia coli* system (Type

552 IE, see Fig.2).

553 (1) Adaptation. DNA fragments of an invading virus/plasmid are integrated as a new

554 unit at the leader side of a CRISPR in the host chromosome, with the

555 simultaneous duplication of a repeat. Selection of sequences to be integrated

556 (proto-spacers) is likely to depend on a proto-spacer adjacent motif (* PAM)¹⁹.

557 Although no direct evidence is presently available on the mechanism of spacer

558 acquisition, the core of the CRISPR/Cas system (Cas1/Cas2), possibly with

559 additional proteins, are prime candidates^{16,21}. In Type II and Type III systems,

560 substantial variations have been demonstrated (see text).

561

562 (2) Expression and Processing. In the Type IE system, a Cascade complex is

563 responsible for processing of pre-crRNA to crRNA via a single cut in the repeat

564 part²¹.

565 (3) Interference. In *E. coli*, the Cascade complex with the crRNA guide target the

566 complementary DNA of an invading virus/plasmid, and Cas3 is required, most

567 likely, to cleave the alien DNA through the endonuclease activity of the HD

568 domain²¹. The PAM (*) appears to play an important role in the interference

569 process^{20,37}. In Type II and Type III systems, no Cas3 ortholog is involved.

570

571

572 **Figure 2. The three major types and 10 subtypes of CRISPR systems and their**
573 **relationships**

574 The operon organization cartoons and color scheme were generally adopted from Ref. 14,
575 with the addition of the identification of Cas7 (COG1857) as a member of the RAMP
576 superfamily. Orthologous genes are color-coded and identified by a family name as in
577 Table 1: bold for definitive proposals and regular font for “legacy names”. The signature
578 genes for CRISPR/cas types are shown within green *boxes*, and for subtypes within red
579 boxes. The letters above the genes show major categories of Cas proteins: L, large
580 CASCADE subunit; S, small CASCADE subunit; R, RAMP CASCADE subunit; RE,
581 RAMP family RNase involved in crRNA processing (experimentally characterized
582 nucleases shown be asterisks); T, transcriptional regulator. Genes coding for inactivated
583 (putative) polymerases are shown by crosses. Question marks denote tentative predictions
584 based on weak sequence similarity. For subtype I-A, the *cas8a1* and *cas8a2* genes are
585 typically mutually exclusive but both can be considered signature genes for the subtype.
586 For type III systems, the *cas1* and *cas2* genes are shown in a dashed box that indicates the
587 loose association of these genes with the type III polymerase/RAMP modules.

588

589

590 **Figure 3. A schematic representation of the phylogenetic tree for Cas1/COG1518**
591 **proteins.**

592 The maximum likelihood tree was constructed using the PHYML program³⁸ from 182
593 informative positions in the multiple alignment of a representative set of 228 Cas1
594 proteins from 442 complete genomes. Six major CRISPR/Cas system subtypes of Type I
595 systems as well as type II and type III systems are color-coded. Dashed lines show *cas1*
596 genes found in “hybrid” CRISPR loci containing genes from both Type I and Type III
597 CRISPR/Cas systems (see text). The subtypes of CRISPR/Cas systems are denoted as in
598 Figure 2. Subtypes I-U, II-U and III-U (U for unclassified) denote CRISPR/Cas systems
599 that lack currently defined subtype-specific signature genes (see text).

Table 1. Classification and nomenclature of CRISPR genes

Proposed gene name	Type or subtype	^a Name from ¹³	Name from ²¹	Structure (PDB code)	^b Family (superfamily)	Representatives	Comment and references
<i>cas1</i>	Type I, Type II, Type III	<i>cas1</i>	<i>cas1</i>	3GOD, 3LFX, 2YZS	COG1518	ygbT, SPy1047, SERP2463	metal-dependent deoxyribonuclease; a unique fold consisting of a N-terminal β strand domain and a C-terminal α -helical domain ^{26,39} ; also binds RNA ^{26,39}
<i>cas2</i>	Type I, Type II, Type III	<i>cas2</i>	<i>cas2</i>	2IVY, 2I8E, 3EXC	COG1343, COG3512	ygbF, SPy1047, SERP2462, y1723 (N-term. domain)	small protein related to VapD; shown to be a RNase specific to U-rich regions ²⁷
<i>cas3'</i>	Type I	<i>cas3</i>	<i>cas3</i>	-	COG1203	ygcB, APE1232	DNA helicase; most proteins have fusion to HD nuclease ²¹ (<i>cas3'</i>)
<i>cas3''</i>	I-A, I-B	-	-	-	COG2254	APE1231, BH0336	HD-like nuclease, specifically digesting double-stranded oligonucleotides and preferably cleaving at G:C pairs ⁴⁰
<i>cas4</i>	I-A, I-B, I-C, I-D, II-B	<i>cas4, csal</i>	-	-	COG1468	APE1239, BH0340	RecB-like nuclease with three-cysteine C-terminal cluster
<i>cas5</i>	I-A, I-B, I-C, I-E	<i>cas5e, d, a, t, h, p, cmx5</i>	<i>casD</i>	3KG4	COG1688 (RAMP)	APE1234, BH0337, DevS, ygcI	predicted subunit of the Cascade complex ²¹ ; in subtype I-C this protein might be the endoribonuclease that generates crRNAs
<i>cas6</i>	I-A, I-B, I-D, III-A, III-B	<i>cas6, cmx6</i>	-	3I4H	COG1583, COG5551 (RAMP)	PF1131, slr7014	Cas6 is an endoribonuclease that generates crRNAs ^{22,23} , predicted subunit of Cascade complex
<i>cas6e</i>	I-E	<i>cse3</i>	<i>casE</i>	1WJ9	(RAMP)	ygcH	homologous to Cas6, but distinct family
<i>cas6f</i>	I-F	<i>csy4</i>	-	2XLJ	(RAMP)	y1727	homologous to Cas6, but distinct family; shown to be is an endoribonuclease that generates crRNAs ²¹
<i>cas7</i>	I-A, I-B, I-C, I-E	<i>cse4, csd2, csh2, cst2, csa2, csp1</i>	<i>casC</i>	-	COG1857, COG3649 (RAMP)	YgcJ, DevR	α/β protein; subunit of Cascade complex ²¹
<i>cas8a1</i>	I-A	<i>csx8, cmx1, csp2, cst1, CxxC-CxxC, csx13</i>	-	-	BH0338-like	BH0338, MTH1090, TM1802, LA3191 ^c	large proteins, some contain Zn-finger domain; nuclease activity has been reported for MTH1090 ⁴¹

<i>Cas8a2</i>	I-A	<i>csa4, csx9</i>	-	-	PH0918	MJ0385, PF0637, AF0070, PH0918, SSO1401, AF1873	see cas8a
<i>cas8b</i>	I-B	<i>csH1,</i> <i>TM1802,</i>	-	-	-		see cas8a
<i>cas8c</i>	I-C	<i>csd1</i>	-	-	-		see cas8a
<i>cas9</i>	Type II	<i>csn1, csx12</i>	-	-	COG3513	SPy1046, FTN_0757	very large protein containing McrA/HNH-nuclease related domain and a RuvC-like nuclease domain;
<i>cas10</i>	Type III	<i>cmr2, csm1,</i> <i>csx11</i>	-	-	COG1353	MTH326, alr1562 ^c , slr7011 ^c	multidomain protein with permuted HD nuclease domain, palm domain, polymerase-thumb-like domain and Zn-ribbon; MTH326-like has inactivated polymerase catalytic domain; alr1562 and slr7011 – predicted only on the basis of size, presence of HD domain, and location with RAMPs in one operon; subunit of Cmr complex ²³
<i>cas10d</i>	I-D	<i>csc3</i>	-	-	COG1353	slr7011	inactivated homolog of Cas10, contains N-terminal HD domain
<i>csy1</i>	I-F	<i>csy1</i>	-	-	y1724-like	y1724	~450 aa protein, predicted to be a subunit of Cascade complex
<i>csy2</i>	I-F	<i>csy2</i>	-	-	(RAMP)	y1725	predicted Cas7 ortholog
<i>csy3</i>	I-F	<i>csy3</i>	-	-	(RAMP)	y1726	predicted Cas5 ortholog
<i>cse1</i>	I-E	<i>cse1</i>	<i>casA</i>	-	ygcL-like	ygcL	large Zn-finger containing proteins; a subunit of the Cascade complex ²¹ ; signature gene for Ecoli/CASS2 subtype
<i>cse2</i>	I-E	<i>cse2</i>	<i>casB</i>	2ZCA	ygcK-like	ygcK	~180 aa protein; a subunit of the Cascade complex ²¹
<i>csc1</i>	I-D	<i>csc1</i>	-	-	alr1563-like (RAMP)	alr1563	
<i>csc2</i>	I-D	<i>csc2, csc1</i>	-	-	COG1337 (RAMP)	slr7012	
<i>csa5</i>	I-A	<i>csa5</i>	-	-	AF1870	AF1870, SSO1398, PF0643, MJ0380	~150 aa protein; alpha-helical

<i>csn2</i>	II-A	<i>csn2</i>	-	-	SPy1049-like	SPy1049	~220 aa protein; predicted to be a functional analog of Cas4 ⁶ based on anti-correlated phyletic patterns
<i>csm2</i>	III-A	<i>csm2</i>	-	-	COG1421	MTH1081, SERP2460	~150 aa protein; mostly α -helical protein; part of distinct polymerase cassette
<i>csm3</i>	III-A	<i>csm3, csc2</i>	-	-	COG1337 (RAMP)	SERP2459, MTH1080	
<i>csm4</i>	III-A	<i>csm4</i>	-	-	COG1567 (RAMP)	SERP2458, MTH1079	
<i>csm5</i>	III-A	<i>csm5</i>	-	-	COG1332 (RAMP)	SERP2457, MTH1078	
<i>csm6</i>	III-A	<i>csm6, APE2256</i>	-	2WTE	COG1517	SSO1445, APE2256	HTH-type transcriptional regulator; often fused to COG1517-like domain
<i>cmr1</i>	III-B	<i>cmr1</i>	-	-	COG1367, (RAMP)	PF1130	subunit of Cmr complex ²³
<i>cmr3</i>	III-B	<i>cmr3</i>	-	-	COG1769 (RAMP)	PF1128	subunit of Cmr complex ²³
<i>cmr4</i>	III-B	<i>cmr4</i>	-	-	COG1336 (RAMP)	PF1126	subunit of Cmr complex ²³
<i>cmr5</i>	III-B	<i>cmr5</i>	-	2ZOP, 2OEB	COG3337	PF1125, MTH324	subunit of Cmr complex ²³
<i>cmr6</i>	III-B	<i>cmr6</i>	-	-	COG1604 (RAMP)	PF1124	subunit of Cmr complex ²³
<i>csb2</i>	I-U ^d	-	-	-	(RAMP?)	GSU0054, Balac_1305	Contains RAMP superfamily motif (G-rich loop) and conserved histidine at the C-terminus
<i>csb1</i>	I-U	<i>GSU0053</i>	-	-	(RAMP?)	GSU0053, Balac_1306	Contains several motifs similar to Cas7 family
<i>csb3</i>	I-U	-	-	-	(RAMP?)	Balac_1303	Contains RAMP superfamily motif (G-rich loop) and conserved histidine at the C-terminus
<i>csx17</i>	I-U	-	-	-	-	Btus_2683	
<i>csx14</i>	I-U	-	-	-	-	GSU0052	
<i>csx10</i>	I-U	<i>csx10</i>	-	-	(RAMP)	Caur_2274	
<i>csp2</i>	I-U	<i>csp2</i>	-	-	PG2018	PG2018	predicted Cas8 ortholog
<i>csx16</i>	III-U	<i>VVA1548</i>	-	-	-	VVA1548	~100 aa protein; often seen in proximity to COG1517
<i>csaX</i>	III-U	<i>csaX</i>	-	-	-	SSO1438	~350 aa protein, no prediction
<i>csx3</i>	III-U	<i>csx3</i>	-	-	-	AF1864	~ 100 aa domain, in some cases fused to COG1517 family domains

<i>csx1</i>	III-U	<i>csa3</i> , <i>csx1</i> (DxTHG motif), <i>csx2</i> , <i>NE0113</i> , <i>TIGR02710</i>	-	1XMX, 2I71	COG1517, COG4006	PF1127, MJ1666, TM1812, NE0113	some are fused to HTH domain (see COG1517/HTH); some proteins have the domain duplication; some have a fusion with HTH and RecB-family nuclease domain; domain appears to have a Rossmann-like fold.
<i>csx15</i>	???	-	-	-	TTE2665	TTE2665	~130 aa protein, no prediction; some are fused to AAA ATPase domain
<i>csf4</i>	Type U	<i>csf4</i>	-	-	-	AFE_1037	DinG family helicase
<i>csf3</i>	Type U	<i>csf3</i>	-	-	(RAMP)	AFE_1040	
<i>csf2</i>	Type U	<i>csf2</i>	-	-	(RAMP)	AFE_1039	
<i>csf1</i>	Type U	<i>csf1</i>	-	-	-	AFE_1038	Zn-finger domain

Note:

^a Subsequent to the original publication¹³, Haft et al. introduced a number of new types of the CRISPR system and gene names that are included in the TIGRFAM database⁴² but mostly fit into previously described gene/protein families. Most of these new names are included in this column.

^b Most of the families correspond to those proposed by Makarova et al.¹⁴ with a few changes and additions.

^c Tentative predictions based on weak sequence similarity, sequence length and gene order in an operon.

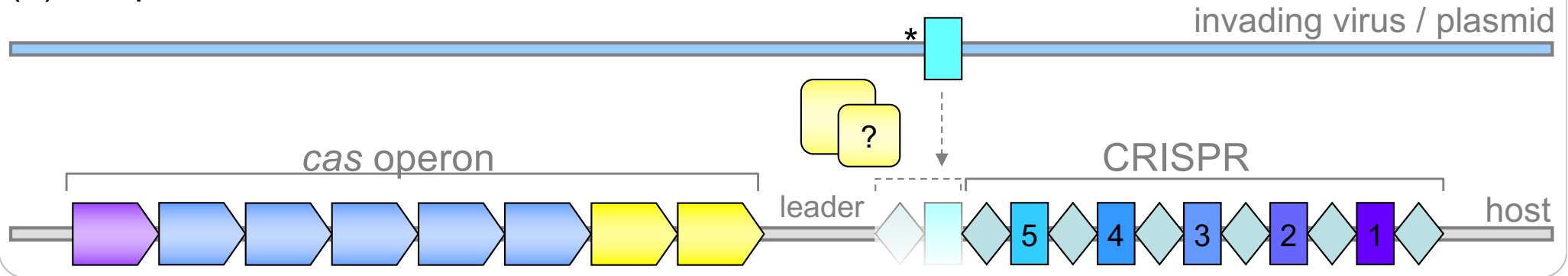
^d Unclassified.

Signature genes for CRISPR/Cas system types and subtypes are shown in the second column in bold.

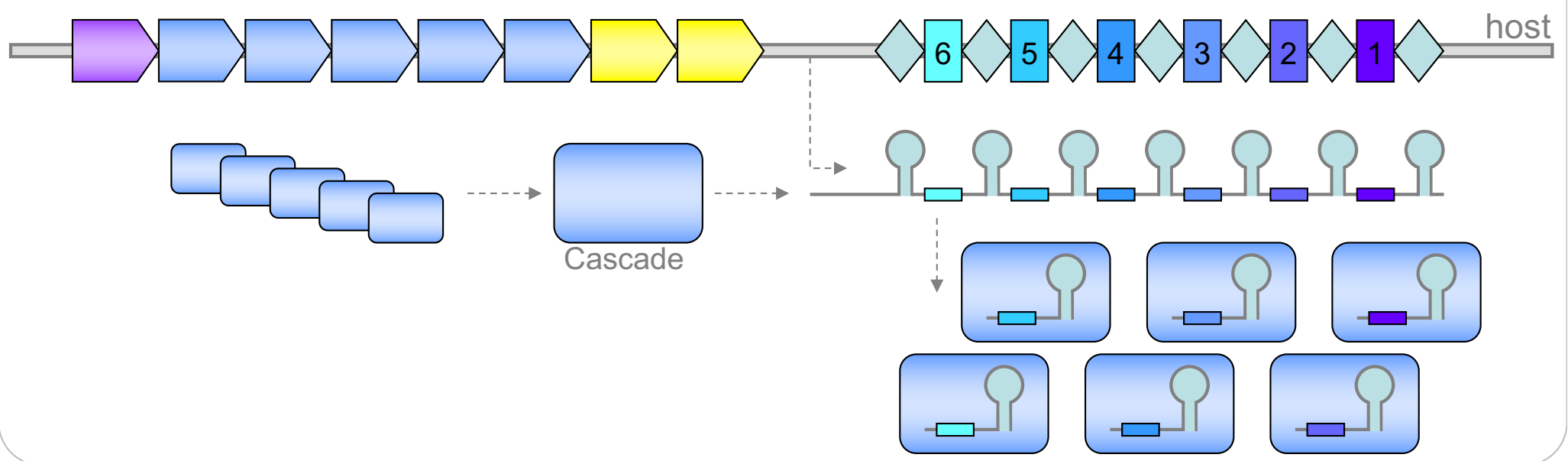
Table 2. Taxonomic distribution of three CRISPR system types

	Genomes analyzed	Cas1	Fraction of genomes with cas1	Type I (Cas7 and Cas3)	Type II (Cas9)	Type III (Cas10)
Archaea	67	54	0.81	50	0	40
Bacteria	639	256	0.40	245	65	99
Crenarchaeota	17	15	0.88	15	0	16
Euryarchaeota	47	37	0.79	33	0	23
Actinobacteria	72	26	0.36	28	15	8
Aquificae	7	5	0.71	7	1	4
Bacteroidetes/ Chlorobi group	32	16	0.50	14	2	6
Chlamydiae/ Verrucomicrobia group	10	2	0.20	0	1	1
Chloroflexi	10	9	0.90	9	2	7
Cyanobacteria	14	7	0.50	7	1	7
Firmicutes	126	56	0.44	40	17	23
Proteobacteria	318	107	0.34	117	20	22
Spirochaetes	13	3	0.23	2	1	0
Thermotogae	11	10	0.91	10	0	9

(1) Adaptation



(2) Expression & Processing



(3) Interference

