1    **Discovery of integrons in Archaea: platforms for cross-domain gene transfer**

2

3    Timothy M. Ghaly[1]*, Sasha G. Tetu[1,2], Anahit Penesyan[1,2], Qin Qi[1], Vaheesan Rajabal[1,2] and

4    Michael R. Gillings[1,2]

5

6    *Corresponding author: timothy.ghaly@mq.edu.au

7

8    [1]School of Natural Sciences, Macquarie University, New South Wales, 2109

9    [2]ARC Centre of Excellence in Synthetic Biology, Macquarie University, New South Wales,

10   2109

11

12   Running title: Integrons in Archaea

13

14   Keywords: Horizontal gene transfer; prokaryotes; metagenome-assembled genomes;

15   evolution; site-specific DNA recombination

16

17

18 **Summary**

19

20 Horizontal gene transfer between different domains of life is increasingly being recognised as

21 an important driver of evolution, with the potential to provide the recipient with new gene

22 functionality and assist niche adaptation[1-3]. However, the molecular mechanisms underlying

23 the integration of exogenous genes from foreign domains are mostly unknown. Integrons are

24 a family of genetic elements that facilitate this process within Bacteria via site-specific DNA

25 recombination[4-7]. Integrons, however, have not been reported outside Bacteria, and thus their

26 potential role in cross-domain gene transfer has not been investigated. Here we show that

27 integrons are also present among diverse phyla within the domain Archaea. Further, we

28 provide experimental evidence that integron-mediated recombination can facilitate the

29 recruitment of archaeal genes by bacteria. Our findings establish a new mechanism that can

30 facilitate horizontal gene transfer between the two domains of prokaryotes, which has

31 important implications for prokaryotic evolution in both clinical and environmental contexts.

32

33 **Main**

34

35 Horizontal gene transfer between different domains of life can be a major driver in species

36 evolution[8]. There are now numerous examples of genes that have been transferred between

37 Archaea, Bacteria and Eukarya[3,9-13]. Among the consequences of such gene transfers are the

38 gain of novel biochemical functions and the ability to colonise specific environmental

39 niches[1-3]. However, the molecular mechanisms for most of these transfer events are unknown.

40 Integrons are genetic elements known to facilitate horizontal gene transfer within

41 Bacteria[4-7]. Integrons can capture exogenous genes, known as gene cassettes, by site-specific

42 recombination. Gene cassette capture is mediated by an integron integrase (IntI), which

43    catalyses the recombination between the recombination site of the inserting cassette (*attC*)

44    and the endogenous integron attachment site (*attI*), immediately adjacent to the *intI* gene.

45    Multiple gene cassettes can be inserted within a single integron, forming cassette arrays that

46    range from 1 to 200+ sequential cassettes[4,6]. Integrons are mostly known for their role in

47    driving the global antibiotic resistance crisis by disseminating diverse resistance determinants

48    among bacterial pathogens[14,15]. However, it is now clear that integrons play a much broader

49    role in bacterial evolution and niche adaptation[16]. The functions encoded by integron gene

50    cassettes are extraordinarily diverse and extend far beyond those of clinical relevance[7,17,18].

51        To date, integrons have only been found within bacterial genomes, where they have

52    been detected within diverse phyla[19]. However, gene cassette amplicon sequencing has

53    yielded cassette-encoded proteins that share homology with archaeal proteins[20,21]. Without

54    broader genomic context, however, the taxonomic residence of such gene cassettes is

55    unknown.

56        Here, we screened all publicly available archaeal genomes to show for the first time

57    that integrons are not limited to Bacteria, but are also present in Archaea. Archaeal integrons

58    exhibit the same characteristics and functional components as bacterial integrons. Further, we

59    demonstrate experimentally that diverse archaeal gene cassettes can be successfully recruited

60    by a bacterial host, facilitated by integron-mediated recombination. Such a mechanism can

61    potentially facilitate a cross-domain highway of gene transfer between Archaea and Bacteria,

62    with important implications for prokaryotic evolution.

63

64    **Discovery of integrons in Archaea**

65

66    Here, we report the discovery of integrons in the domain Archaea. We screened 6,718

67    archaeal genomes for integrons using the standard criteria applied to integron surveys in

68    Bacteria[19,22,23]. These include the presence of integron integrase genes and/or clusters of gene

69    cassette *attC*s (defined as at least two *attC*s with less than 4 kb between each). We identified

70    integrons in 75 archaeal metagenome-assembled genomes (MAGs) from 9 phyla (Fig. 1 and

71    Supplementary Table 1). It is not surprising that integrons were detected only in MAGs,

72    given that they constituted ~95% of all available archaeal genomes. However, to ensure that

73    these integrons did not arise from contaminating bacterial contigs, incorrectly binned with

74    archaeal MAGs, we applied stringent MAG refinement and quality filtering (see Methods for

75    details). Additionally, we found that ~7% of integron-bearing MAGs had at least one

76    archaeal phylogenetic marker gene on the same contig as an integron (Supplementary Table

77    2), confirming these to be located on archaeal chromosomes. No integron was ever co-located

78    with a bacterial marker gene. The markers used for this analysis consisted of a

79    comprehensive set of 122 archaeal and 120 bacterial proteins identified as suitable for

80    phylogenetic inference[24].

81        Among the 75 archaeal genomes, we detected six IntIs and 539 *attC* sites (excluding

82    all singleton *attC*s). We found that archaeal *attC*s and IntIs are largely restricted to one clade

83    of Archaea (Fig. 1), with some outliers, suggesting that integron diversification, for the most

84    part, has likely occurred within one archaeal clade, with occasional horizontal movements to

85    other archaeal phyla. In particular, integrons were significantly enriched in the phylum

86    Asgardarchaeota ($\chi^2$ test, p < 0.00001) (Fig. 1), being detected in almost 8% of available

87    Asgard genomes. Asgardarchaeota contributed the most genomes with detectable integrons

88    (28%) and the greatest number of gene cassettes (24.9%), despite having relatively few

89    genomes among the dataset (comprising 4% of available archaeal genomes). We also

90    detected integrons in 3-4% of genomes from the phyla Hadarchaeota and

91    Hydrothermoarchaeota (Fig. 1), although these comprised few available genomes (n<50). A

92    skewed phylogenetic distribution of integrons has similarly been observed among Bacteria[19].

93    For example, in the phylum Proteobacteria, integrons are enriched within the class

94    Gammaproteobacteria (20% of genomes), while being entirely absent from its sister class

95    Alphaproteobacteria. This is intriguing given that integrons have been detected at widely

96    varying prevalence in more distantly related bacterial phyla such as Cyanobacteria,

97    Spirochaetota, Planctomycetota, Chloroflexota, Bacteroidota and Desulfobacterota[19,22].

98

99    *Genetic structure of archaeal integrons*

100       We found that archaeal integrons exhibit the same structure and functional

101    components as bacterial integron cassette arrays (Extended Data Fig. 1). That is, tandem

102    arrays of short open reading frames (ORFs), generally in the same orientation, interspersed

103    by *attC* recombination sites. Archaeal *attC*s exhibit the same single-stranded folding structure

104    as bacterial *attC*s, which is essential for them to act as structure-specific DNA recombination

105    sites[25-31]. We also note that archaeal IntIs exhibit the defining characteristics of bacterial

106    IntIs, being tyrosine recombinases that possess a unique IntI-specific additional domain

107    surrounding the patch III motif region necessary for integron-mediated recombination[32]. We

108    found examples of 'complete' integrons, these being cassette arrays adjacent to a detectable

109    *intI* gene (Extended Data Fig. 1). We also found examples of putative *attI* sites, which act as

110    insertion points for incoming gene cassettes. These *attI*s were immediately downstream of the

111    *intI* gene, semi-conserved across distinct archaeal phyla (Extended Data Fig. 2a,b), and

112    exhibited the same canonical insertion point as all known bacterial *attI*s (Extended Data Fig.

113    2c).

114       Most archaeal integrons that we identified were CALINs (clusters of *attC*s lacking

115    integron integrases; Supplementary Table 3). This is not surprising given the fragmented

116    nature of MAGs, and the high prevalence of CALINs also found in bacterial genomes.

117    Indeed, among Bacteria, CALINs are more abundant than complete integrons that possess an

118    *intI* gene, and exhibit a much wider taxonomic distribution[19]. Two In0 elements were also

119    detected among Archaea. These are integrons that have an *intI* gene without an adjacent *attC*

120    site (Extended Data Fig. 1). However, both archaeal genomes with an In0 also had clusters of

121    *attC* sites on other contigs. Among our dataset, the longest array of *attC*s on the same contig

122    was 12, however, we found as many as 107 *attC*s (over 18 contigs) within a single MAG

123    (Supplementary Table 1). The number of *attC*s within a single MAG ranged from 2 to 107,

124    with an average of 7 *attC*s.

125

126    **Platforms for cross-domain gene transfer**

127

128    Archaeal gene cassettes with *attC*s from diverse phyla can be recognised and recruited by

129    Bacteria (Fig. 2). We demonstrate that cassette insertion (*attC* x *attI* recombination) can

130    occur following the conjugation of circular DNA molecules with archaeal *attC*s into an

131    *Escherichia coli* recipient harbouring a bacterial class 1 integron (Fig. 2a). Insertion events

132    were confirmed with Sanger sequencing of the PCR-amplified *attC*/*attI* recombination

133    junctions (Fig. 2a, Extended Data Fig. 3). We found that recruitment of cassettes with

134    archaeal *attC*s occurred at similar frequencies to that of the paradigmatic bacterial *attC* site,

135    *attC*$_{aadA7}$, which we used as a positive control (Fig. 2b, Extended Data Table 1). We observed

136    an average recombination frequency of $2.5 \times 10^{-1}$ between *attI1* and *attC*$_{aadA7}$. Comparable

137    frequencies (ranging from $1.9 \times 10^{-4} - 3.2 \times 10^{-1}$, with an average of $5.1 \times 10^{-2}$) were observed for

138    eight out of nine archaeal *attC*s (Kruskal-Wallis test, p=0.488), which were selected from

139    multiple archaeal phyla. Further, we confirmed that cassette recruitment was mediated by

140    IntI1 activity, since no *attC* x *attI* recombination events were detected when *intI1* was absent

141    or when its expression was suppressed (Extended Data Table 1). We therefore show that

142    integron-mediated gene transfer can occur between the two domains of prokaryotes.

143    Importantly, we find that the most clinically significant class of integrons (class 1)

144    can recruit archaeal cassettes as efficiently as bacterial cassettes. Class 1 integrons are highly

145    promiscuous due to their association with diverse mobile genetic elements, facilitating their

146    spread into at least 100 bacterial species[7]. They collectively carry more than 130 different

147    resistance genes[14], most of which are of unknown taxonomic origin[22]. Our findings open the

148    possibility that Archaea could be an unexplored source of class 1 integron gene cassettes.

149    Regardless, our findings indicate that any bacterial strain with a class 1 integron has the

150    capacity to incorporate exogenous genes from diverse archaeal phyla, greatly expanding the

151    genetic pool that they have access to.

152    The cross-domain transfer of integron gene cassettes is possibly widespread. For

153    example, we detected 23 *attC*s from six archaeal genomes that exhibited 95-100% nucleotide

154    identity to *attC*s within sequenced bacterial integrons (Supplementary Table 4). The archaeal

155    *attC*s were from three phyla: Nanoarchaeota, Thermoproteota and Hadarchaeota. The

156    homologous *attC*s in Bacteria were found in 26 genomes from 5 phyla: Proteobacteria,

157    Spirochaetota, Myxococcota, Nitrospirota and Desulfobacterota. One of these *attC* sites was

158    associated with a class 1 integron gene cassette, encoding an NADPH-dependent

159    oxidoreductase found on five different Enterobacteriaceae plasmids (Supplementary Table 4).

160    In Archaea, however, this *attC* site was part of a cassette that encoded a ligand-binding

161    protein of unknown function. Nevertheless, since strong *attC* homology is a characteristic of

162    cassettes that share the same taxonomic origin[22,33,34], it is possible that some clinically

163    relevant gene cassettes now found on class 1 integrons might be of archaeal origin.

164

165    **Diversity of archaeal integrons**

166

167    *Diversity of integron integrases*

168　　　　　Archaeal IntIs are phylogenetically distinct from bacterial IntIs (Fig. 3). We detected

169　　six IntIs from four archaeal phyla (Fig. 1), however, three of these were excluded from

170　　further phylogenetic analysis based on either short sequence length (< 200 amino acids) or

171　　partial coverage of the IntI-specific domain (Extended Data Fig. 4). We found that archaeal

172　　IntIs form their own monophyletic clade separate from known bacterial IntIs[22]. This strongly

173　　suggests that IntI radiation has occurred within Archaea and that their distribution, at least

174　　among the archaeal genomes in our dataset, is not likely to be the result of multiple IntI

175　　acquisitions from Bacteria. Regardless, we show that IntIs from distinct archaeal phyla,

176　　isolated from different environments, are more closely related to each other than they are to

177　　any bacterial IntI.

178　　　　　The closest sister clade to the archaeal IntIs comprises two Spirochaetota IntIs (Fig.

179　　3). Intriguingly, these two IntIs are phylogenetically distinct from 'typical' Spirochaetota

180　　IntIs, which are generally in reverse orientation[5,35]. Further, the two Spirochaetota that

181　　harboured atypical IntIs were isolated from extreme environments: a brine layer within an

182　　alkaline lake and a hot spring, respectively; environments known to have a relatively high

183　　abundance of Archaea[36]. Thus, these atypical Spirochaetota IntIs might have been

184　　horizontally acquired from Archaea that share the same extreme environments.

185

186　　*Diversity of attC recombination sites*

187　　　　　Archaeal *attC*s exhibit broad sequence and structural diversity (Fig. 4a). We find that

188　　some archaeal phyla possess *attC*s with a restricted diversity (e.g., Hadarchaeota and

189　　Aenigmatarchaeota), while other phyla have extremely variable *attC*s distributed throughout

190　　the *attC* diversity space (e.g., Asgardarchaeota, Nanoarchaeota and Thermoproteota). This

191　　distribution could indicate that different taxa have different propensities for horizontal

192　　exchange of gene cassettes[7,22]. We show that archaeal *attC*s are significantly more similar

193     within a genome than between genomes (Fig. 4b). This characteristic is also a hallmark of

194     chromosomal bacterial integrons[19,34]. We also show that *attC*s are more similar between

195     different genomes from the same archaeal order than they are between genomes from

196     different orders (Fig. 4c). This order-level *attC* homology is also seen within Bacteria[22,33].

197     Thus, the ecological and evolutionary forces that promote and/or constrain *attC* diversity[7] are

198     likely to be similar for both Archaea and Bacteria.

199         There is a clear overlap in the sequence and structural diversity of *attC*s from Archaea

200     and Bacteria (Fig. 4a). This provides additional evidence that the mechanistic overlap

201     between archaeal and bacterial *attC*s is extensive, and thus, cross-domain transfer of cassettes

202     could be common in shared environments. It also suggests that the recruitment of extra-

203     domain gene cassettes can be facilitated by diverse classes of integrons, of which there are

204     thousands (based on IntI amino acid homology[37]). The broad distribution of integrons among

205     the two domains suggest that integron-mediated transfer plays an important role in

206     prokaryotic evolution.

207

208     *Functional diversity of gene cassettes*

209         We detected 549 cassette-encoded proteins among Archaea. Only 23.1% of these

210     could be classified into a known COG category (Extended Data Fig. 5). In contrast, 47.4% of

211     all proteins from the 75 integron-bearing archaeal genomes could be assigned a known COG

212     category. This underrepresentation ($\chi^2$ test, p < 0.00001) of known COGs among cassette

213     proteins has previously been reported for bacterial integrons[4,5,33]. To gain further insight into

214     possible cassette functions, eggNOG 5.0[38] and Pfam[39] database searches were performed,

215     assigning putative functions to 228 (41.5%) of the archaeal cassette-encoded proteins. Out of

216     those with functional predictions, proteins involved in toxin-antitoxin (TA) systems (10.5%);

217     phage resistance proteins via DNA methylation or restriction endonuclease activities (8.3%);

218     and acetyltransferases (4.4%) were particularly prevalent (Supplementary Table 5). These are

219     the functions most commonly reported for gene cassettes in Bacteria[5,7,33,34,40]. TA gene

220     cassettes are particularly common in bacterial integrons, where they can stabilise very large

221     cassette arrays[41,42]. The antitoxin modules of TA cassettes can also counteract the toxins of

222     homologous systems found on plasmids and phage, thus potentially protecting their host from

223     invading mobile elements[43,44].

224         In addition, 13.2% of archaeal cassette-encoded proteins had signal peptides, which

225     represents a significant enrichment relative to their broader genomic contexts (6.9%, $\chi^2$ test, p

226     < 0.00001). Signal peptides are short amino acid tag sequences that target proteins into, or

227     across, membranes. Again, transmembrane and secreted proteins are commonly encoded by

228     gene cassettes in Bacteria[33], and are hypothesised to help facilitate interactions with their

229     broader environment[7].

230         Indeed, we find that functions of archaeal cassettes are associated with their

231     environment (Fig. 5). Functional families cluster according to their specific environment, and

232     these environmental clusters, in turn, group according to their broader environmental type

233     (Fig. 5). This environmentally explicit clustering might be the result of local ecological and

234     evolutionary forces. That is, gene cassettes in Archaea confer niche-specific functional traits

235     and/or horizontal transfer of cassettes occurs between archaeal phyla co-located in the same

236     environment.

237

238     **Conclusion**

239

240     Here, we present the first evidence of integrons in the domain Archaea. We demonstrate that

241     they have the same functional characteristics as bacterial integrons. We also present

242     experimental evidence that bacteria can successfully recruit archaeal gene cassettes,

243  facilitated by integron-mediated DNA recombination. Our results thus establish a novel

244  mechanism for cross-domain gene transfer between Archaea and Bacteria. We also find that,

245  although archaeal IntIs are phylogenetically distinct from bacterial IntIs, their associated *attC*

246  recombination sites are shared with Bacteria. This suggests that integron-mediated cross-

247  domain gene transfer is widespread and plays an important role in prokaryotic evolution.

248

249  **Methods**

250

251  *Data acquisition and quality filtering*

252      All available archaeal genomes were downloaded from the NCBI Assembly Database

253  (n=8,160; last accessed 2021-Oct-5). Of these, ~ 95% were metagenome-assembled genomes

254  (MAGs). We applied stringent filtering criteria to remove low quality MAGs. First, to

255  improve MAG quality, we identified and removed contaminating contigs from each MAG

256  using MAGpurify v2.1.2[45] with the following modules: '*phylo-markers*', which finds

257  taxonomically discordant contigs using 100 archaeal and 88 bacterial single-copy taxonomic

258  marker genes from the PhyEco database[46]; '*clade-markers*', which finds contaminating

259  contigs using a database of 855,764 clade-specific prokaryotic marker genes (MetaPhlAn2

260  database[47]); '*tetra-freq*', which employs principal component analysis (PCA) to identify

261  contaminating contigs with outlier tetra-nucleotide frequency; and '*gc-content*', which uses

262  PCA to identify contaminating contigs with outlier GC content.

263      After refinement, the quality of the genomes was assessed using CheckM v1.1.3[48],

264  which uses single-copy lineage-specific marker genes to estimate genome completeness and

265  contamination. There is strong community consensus that high quality MAGs are those that

266  are more than 90% complete and have less than 5% contamination, while medium quality

267  MAGs have a completeness ≥50% and contamination <10%[24,45,49-52]. In this context,

268    however, we were more concerned with the level of contamination than completeness, and

269    thus removed all genomes with an estimated contamination ≥5%. The completeness of the

270    remaining genomes ranged from 15% – 100%, with a median of 81%. The estimated

271    contamination ranged from 0% – 4.98%, with a median of 0.93%.

272         Archaeal genomes were assigned taxonomic classifications based on the Genome

273    Taxonomy Database (GTDB)[49-51] using GTDB-Tk v1.6.0[53] with release 06-RS202 of the

274    GTDB. We employed the *classify_wf* command with default settings. This workflow

275    identifies and aligns 120 bacterial and 122 archaeal phylogenetic marker genes[24]. GTDB-Tk

276    then classifies each genome based on its placement into domain-specific reference trees (built

277    from 47,899 prokaryote genomes), its relative evolutionary divergence, and average

278    nucleotide identity to reference genomes in the GTDB. Any genomes not classified within

279    the domain Archaea were removed. This resulted in a final set of 6,718 archaeal genomes

280    retained for further analysis.

281         To infer the phylum-level phylogeny of Archaea, the highest quality representative

282    genome from each phylum was selected based on its genome quality score (defined by Parks

283    et al.[24] as the estimated completeness of a genome minus five times its estimated

284    contamination). From representative genomes, a concatenated multiple protein sequence

285    alignment of the 122 archaeal phylogenetic markers was generated using GTDB-Tk v1.6.0[53].

286    A maximum-likelihood tree was generated from the alignment using IQ-TREE v1.6.12[54] with

287    the best-suited protein model as determined by ModelFinder[55] and 1,000 bootstrap replicates

288    [parameters: -m MFP -bb 1000].

289

290    *Integron detection*

291         Due to faster processing speeds of large datasets, we initially screened all filtered

292    genomes for *attC* recombination sites using *attC*-screening.sh[37]

293    (https://github.com/timghaly/integron-filtering) with default parameters. This script uses the

294    HattCI[56] + Infernal[57] pipeline (first described by Pereira *et al.*[23]) to search for the conserved

295    sequence and structure of *attC* sites. Genomes that had at least one detectable *attC* site were

296    additionally screened using IntegronFinder v2.0rc6[19] [parameters: --local-max --cpu 24 --

297    gbk], which searches for integron integrases and gene cassette arrays. Only IntIs, *attCs* and

298    cassette-encoded proteins identified by IntegronFinder were included in downstream

299    analyses.

300         To ensure that these integrons were not from contaminating bacterial contigs that had

301    been incorrectly binned with archaeal MAGs, we screened all contigs containing an integron

302    for prokaryotic marker genes using GTDB-Tk v1.6.0[53]. These consisted of 122 archaeal and

303    120 bacterial proteins identified as suitable phylogenetic markers[24]. We found a total of nine

304    prokaryotic marker genes among seven integron-bearing contigs. All nine markers were

305    confirmed to be archaeal via a BLASTP search of the NCBI nr database (Supplementary

306    Table 2).

307

308    *Analysis of integron integrases, attC sites and cassette-encoded proteins*

309         IntegronFinder identifies IntIs using the overlap of two protein hidden Markov model

310    (HMM) profiles. The first is the Pfam profile PF00589 to identify tyrosine recombinases, and

311    the second is a protein profile built from the IntI-specific domain that separates IntIs from

312    other tyrosine recombinases[32]. Identified archaeal IntIs, with matches to both protein profiles,

313    were placed in a phylogeny alongside a set of previously identified bacterial IntIs[22]. IntIs

314    shorter than 200 amino acids or those that did not span the complete IntI-specific domain

315    were removed from phylogenetic analysis. The remaining IntIs were aligned using MAFFT

316    v7.271[58] [parameters: --localpair --maxiterate 1000] and trimmed using trimAl v1.2rev59

317    [parameters: -automated1]. A maximum-likelihood tree was generated from the alignment

318     using IQ-TREE v1.6.12[54] with the best-suited protein model as determined by ModelFinder[55]

319     and 1,000 bootstrap replicates [parameters: -m MFP -bb 1000].

320         The sequence and structural diversity of *attC*s was assessed using RNAclust v1.3[59] as

321     previously described[22]. RNAclust uses LocARNA[60,61] to generate pairwise structural

322     alignments (based on both sequence and folding structure) of input sequences. RNAclust then

323     calculates pairwise distances to create a hierarchical-clustering tree from a WPGMA analysis.

324     All archaeal *attC*s along with a set of previously identified *attC*s from representative bacterial

325     taxa[22] were clustered using RNAclust's default parameters.

326         Cassette-encoded proteins identified by IntegronFinder were functionally annotated

327     using InterProScan v5.44-79.0[62], with default parameters against the Pfam[39] database, and

328     eggNOG-mapper v2.0.1b[63,64], executed in DIAMOND[65] mode against the eggNOG 5.0

329     database[38]. To identify cassettes that encode transmembrane and secreted proteins, we

330     searched protein sequences for prokaryotic signal peptides using SignalP 5.0[66] with default

331     parameters. The correlation analysis of cassette functions was performed as described in

332     Penesyan et al[67]. Briefly, Pearson's correlations, based on co-occurrences between Pfam

333     functions, specific environments and archaeal phyla were calculated using the Hmisc v4.5-0

334     R package[68]. The network was generated from all positive correlations with p-values <0.05

335     using the ForceAtlas2 layout algorithm[69] within the Gephi software[70]. Specific correlations

336     and the description of Pfam functions are listed in Supplementary Table 6.

337

338     *Bacterial strains and plasmids for attC recombination assays*

339         The bacterial strains and plasmids used in this study are listed in Supplementary Table

340     7. LB medium (Lennox) was used to grow bacterial strains supplemented with appropriate

341     antimicrobial agents. The final concentrations of antimicrobial agents used were kanamycin

342     (Km) = 50 µg/mL, carbenicillin (Cb) = 75 µg/mL, and chloramphenicol (Cm) = 20 µg/mL.

343    LB medium was supplemented with 0.3 mM 2,6-diaminopimelic acid (DAP) to culture the

344    auxotrophic *E. coli* WM3064 λpir strain[71].

345

346    *Construction of attC donor strains*

347         Nine archaeal *attC*s, selected from diverse archaeal phyla (Supplementary Table 8)

348    along with one bacterial *attC* (*attC$_{aadA7}$* ) were used for the recombination assays. Two donor

349    strains were constructed for each *attC*, delivering either the *attC* top or bottom strands via

350    conjugation. Overlapping forward and reverse primers were designed to generate each *attC*

351    sequence flanked by *Xba*I and *Bam*HI overhangs respectively (e.g. primer pair *attC*-aadA7-

352    FW/REV for *attC$_{aadA7}$*). The annealed primer dimers were then ligated into the mobilisable

353    suicide vector pJP5603[72,73]. The *attC* top strand donor strains were generated by transforming

354    the ligation product into electrocompetent cells of the DAP auxotrophic *E. coli* strain

355    WM3064 λpir. Using the same procedures, all *attC* top strand donor plasmids and strains

356    were constructed using the pairs of long primers listed in Supplementary Table 9.

357         To deliver *attC* bottom strands, the pJP5603rev (pJPrev) vector was generated to

358    invert *oriT* orientation relative to that of the pJP5603 parental vector. The multiple-cloning

359    site and vector backbone of pJP5603 were PCR amplified using the primer pairs pJP-MCS-

360    FW/REV and pJP-Backbone-FW/REV respectively (with *Xho*I and *Mlu*I restriction sites

361    introduced) followed by restriction digest and ligation. The same primer pairs for generating

362    the top strand donor plasmids were used to create the bottom strand donor plasmids and

363    strains by cloning the same *attC* sequences into the *Xba*I/*Bam*HI sites of pJPrev.

364

365    *Construction of the recipient strain*

366    We generated a recipient strain using *E. coli* UB5201[74] that carried the *intI1* gene and the

367    *attI1* recombination site residing on the pBAD24[75] and pACYC184[76] backbones,

368     respectively. The *intI1* gene of the R388 plasmid[77] was PCR amplified using the primer pair

369     *intI1_Eco*RI-F/*intI1_Hin*dIII-R (Supplementary Table 9). The L-arabinose inducible

370     pBAD::*intI1* plasmid was generated by cloning *intI1* into the pBAD24 expression vector. The

371     pACYC184::*attI1* recipient plasmid was created by assembling the *attI1* sequence (from

372     R388) into the pACYC184 plasmid backbone using the NEBuilder HiFi DNA Assembly

373     Cloning Kit (New England Biolabs, United States). The PCR products required for the

374     assembly were generated using the *attI1*_fw/*attI1*_rev and

375     pACYC184_backbone_F/pACYC184_backbone_R primer pairs. *E. coli* UB5201 strain was

376     co-transformed with pBAD::*intI1* and pACYC184::*attI1* to generate the *E. coli* UB5201 +

377     pBAD::*intI1* + pACYC184::*attI1* recipient strain for *attC* x *attI* suicide conjugation assays. *E.*

378     *coli* UB5201 + pBAD24 + pACYC184::*attI1* was created as an *intI1*-negative control strain.

379     All plasmid constructs were confirmed by Sanger sequencing and restriction enzyme digests.

380

381     *attC x attI suicide conjugation assays*

382     The frequencies of recombination between the archaeal *attC* sequences and the class 1

383     integron *attI1* site were quantified using previously established *attC* x *attI* suicide

384     conjugation methods[25,29,31,78,79] with minor modifications. Briefly, the Cb-resistant UB5201 +

385     pBAD::*intI1* + pACYC184::*attI1* recipient strain was filter-mated with Km-resistant

386     WM3064 λpir *attC* donor strains in DAP-supplemented LB media. The expression of *intI1*

387     was either induced using L-arabinose (2 mg/mL) or suppressed with D-glucose (10 mg/mL).

388     After 6 hours of incubation at 37˚C, the recovered conjugation mix was plated on DAP-free

389     LB agar with Km, as well as on LB agar containing Cb. This method allowed for negative

390     selection of the donor strain, which cannot grow in the absence of DAP, and positive

391     selection of the recombinant recipient clones, which become Km-resistant following plasmid

392     co-integration (Fig. 2a). The recombination frequency was determined as the ratio of the

393     colony forming units (CFU) for Km-resistant recombinants to the CFU for the total number

394     of Cb-resistant recipients after two days of incubation. All assays were performed in three

395     biological replicates, and recombination frequencies were calculated as the mean of the three

396     independent experiments. To confirm the co-integrates, colony PCR was performed on eight

397     randomly chosen colonies per conjugation set for each biological replicate using the

398     following primer pairs pACYC_F/M13F and pACYC_R/M13R (Extended Data Fig. 3).

399     Sanger sequencing of PCR products was performed for four recombinant colonies per

400     conjugation set.

401

402     **Data availability**

403     All genome sequences were downloaded from the NCBI Assembly Database

404     (https://www.ncbi.nlm.nih.gov/assembly; last accessed 2021-Oct-5).

405

406     **Code availability**

407     All code and software used in this study are described within the manuscript.

408

409     **References**

410

411   1     Husnik, F. & McCutcheon, J. P. Functional horizontal gene transfer from bacteria to
412         eukaryotes. *Nature Reviews Microbiology* **16**, 67-79, doi:10.1038/nrmicro.2017.137
413         (2018).
414   2     Metcalf, J. A., Funkhouser-Jones, L. J., Brileya, K., Reysenbach, A.-L. &
415         Bordenstein, S. R. Antibacterial gene transfer across the tree of life. *eLife* **3**, e04266
416         (2014).
417   3     Schönknecht, G. *et al.* Gene transfer from Bacteria and Archaea facilitated evolution
418         of an extremophilic eukaryote. *Science* **339**, 1207-1210,
419         doi:doi:10.1126/science.1231707 (2013).
420   4     Mazel, D. Integrons: agents of bacterial evolution. *Nature Reviews Microbiology* **4**,
421         608-620, doi:10.1038/nrmicro1462 (2006).
422   5     Boucher, Y., Labbate, M., Koenig, J. E. & Stokes, H. W. Integrons: mobilizable
423         platforms that promote genetic diversity in bacteria. *Trends in Microbiology* **15**, 301-
424         309, doi:https://doi.org/10.1016/j.tim.2007.05.004 (2007).

425  6   Gillings, M. R. Integrons: past, present, and future. *Microbiology and Molecular Biology Reviews* **78**, 257-277 (2014).
427  7   Ghaly, T. M. *et al.* The natural history of integrons. *Microorganisms* **9**, 2212 (2021).
428  8   Bruto, M. *et al.* in *Evolutionary Biology: Exobiology and Evolutionary Mechanisms* (ed Pierre Pontarotti) 165-179 (Springer Berlin Heidelberg, 2013).
430  9   Sutherland, K. M., Ward, L. M., Colombero, C.-R. & Johnston, D. T. Inter-domain horizontal gene transfer of nickel-binding superoxide dismutase. *Geobiology* **19**, 450-459, doi:https://doi.org/10.1111/gbi.12448 (2021).
433  10  Frigaard, N.-U., Martinez, A., Mincer, T. J. & DeLong, E. F. Proteorhodopsin lateral gene transfer between marine planktonic Bacteria and Archaea. *Nature* **439**, 847-850, doi:10.1038/nature04435 (2006).
436  11  Dunning Hotopp, J. C. Horizontal gene transfer between bacteria and animals. *Trends in Genetics* **27**, 157-163, doi:https://doi.org/10.1016/j.tig.2011.01.005 (2011).
438  12  Bock, R. The give-and-take of DNA: horizontal gene transfer in plants. *Trends in Plant Science* **15**, 11-22, doi:https://doi.org/10.1016/j.tplants.2009.10.001 (2010).
440  13  Nelson, K. E. *et al.* Evidence for lateral gene transfer between Archaea and Bacteria from genome sequence of *Thermotoga maritima*. *Nature* **399**, 323-329 (1999).
442  14  Partridge, S. R., Tsafnat, G., Coiera, E. & Iredell, J. R. Gene cassettes and cassette arrays in mobile resistance integrons. *FEMS Microbiology Reviews* **33**, 757-784, doi:10.1111/j.1574-6976.2009.00175.x (2009).
445  15  Zhu, Y.-G. *et al.* Continental-scale pollution of estuaries with antibiotic resistance genes. *Nature Microbiology* **2**, 16270, doi:10.1038/nmicrobiol.2016.270 (2017).
447  16  Escudero, J. A., Loot, C., Nivina, A. & Mazel, D. The integron: adaptation on demand. *Microbiology Spectrum* **3**, 3.2.10, doi:doi:10.1128/microbiolspec.MDNA3-0019-2014 (2015).
450  17  Ghaly, T. M., Geoghegan, J. L., Tetu, S. G. & Gillings, M. R. The peril and promise of integrons: beyond antibiotic resistance. *Trends in Microbiology* **28**, 455-464, doi:https://doi.org/10.1016/j.tim.2019.12.002 (2020).
453  18  Ghaly, T. M., Geoghegan, J. L., Alroy, J. & Gillings, M. R. High diversity and rapid spatial turnover of integron gene cassettes in soil. *Environmental Microbiology* **21**, 1567-1574, doi:https://doi.org/10.1111/1462-2920.14551 (2019).
456  19  Cury, J., Jové, T., Touchon, M., Néron, B. & Rocha, E. P. Identification and analysis of integrons and cassette arrays in bacterial genomes. *Nucleic Acids Research* **44**, 4539-4550, doi:10.1093/nar/gkw319 (2016).
459  20  Elsaied, H., Stokes, H. W., Yoshioka, H., Mitani, Y. & Maruyama, A. Novel integrons and gene cassettes from a Cascadian submarine gas-hydrate-bearing core. *FEMS Microbiology Ecology* **87**, 343-356, doi:10.1111/1574-6941.12227 (2014).
462  21  Koenig, J. E. *et al.* Integron gene cassettes and degradation of compounds associated with industrial waste: The case of the Sydney Tar Ponds. *PLoS One* **4**, e5276, doi:10.1371/journal.pone.0005276 (2009).
465  22  Ghaly, T. M., Tetu, S. G. & Gillings, M. R. Predicting the taxonomic and environmental sources of integron gene cassettes using structural and sequence homology of *attC* sites. *Communications Biology* **4**, 946, doi:10.1038/s42003-021-02489-0 (2021).
469  23  Pereira, M. B. *et al.* A comprehensive survey of integron-associated genes present in metagenomes. *BMC Genomics* **21**, 495, doi:10.1186/s12864-020-06830-5 (2020).
471  24  Parks, D. H. *et al.* Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life. *Nature Microbiology* **2**, 1533-1542, doi:10.1038/s41564-017-0012-7 (2017).

25   Bouvier, M., Demarre, G. & Mazel, D. Integron cassette insertion: a recombination process involving a folded single strand substrate. *The EMBO Journal* **24**, 4356-4367, doi:https://doi.org/10.1038/sj.emboj.7600898 (2005).

26   MacDonald, D., Demarre, G., Bouvier, M., Mazel, D. & Gopaul, D. N. Structural basis for broad DNA-specificity in integron recombination. *Nature* **440**, 1157-1162, doi:10.1038/nature04643 (2006).

27   Bouvier, M., Ducos-Galand, M., Loot, C., Bikard, D. & Mazel, D. Structural features of single-stranded integron cassette *attC* sites and their role in strand selection. *PLoS Genetics* **5**, e1000632 (2009).

28   Demarre, G., Frumerie, C., Gopaul, D. N. & Mazel, D. Identification of key structural determinants of the IntI1 integron integrase that influence *attC×attI1* recombination efficiency. *Nucleic Acids Research* **35**, 6475-6489, doi:10.1093/nar/gkm709 (2007).

29   Nivina, A., Escudero, J. A., Vit, C., Mazel, D. & Loot, C. Efficiency of integron cassette insertion in correct orientation is ensured by the interplay of the three unpaired features of *attC* recombination sites. *Nucleic Acids Research* **44**, 7792-7803, doi:10.1093/nar/gkw646 (2016).

30   Mukhortava, A. *et al.* Structural heterogeneity of *attC* integron recombination sites revealed by optical tweezers. *Nucleic Acids Research* **47**, 1861-1870, doi:10.1093/nar/gky1258 (2018).

31   Nivina, A. *et al.* Structure-specific DNA recombination sites: Design, validation, and machine learning–based refinement. *Science Advances* **6**, eaay2922 (2020).

32   Messier, N. & Roy, P. H. Integron integrases possess a unique additional domain necessary for activity. *Journal of Bacteriology* **183**, 6699-6706 (2001).

33   Rowe-Magnus, D. A., Guerout, A.-M., Biskri, L., Bouige, P. & Mazel, D. Comparative analysis of superintegrons: engineering extensive genetic diversity in the Vibrionaceae. *Genome Research* **13**, 428-442 (2003).

34   Rowe-Magnus, D. A. *et al.* The evolutionary history of chromosomal super-integrons provides an ancestry for multiresistant integrons. *Proceedings of the National Academy of Sciences* **98**, 652-657 (2001).

35   Wu, Y.-W., Doak, T. G. & Ye, Y. The gain and loss of chromosomal integron systems in the *Treponema* species. *BMC Evolutionary Biology* **13**, 16, doi:10.1186/1471-2148-13-16 (2013).

36   Rampelotto, P. H. Extremophiles and extreme environments. *Life* **3**, 482-485 (2013).

37   Ghaly, T. M. *et al.* Methods for the targeted sequencing and analysis of integrons and their gene cassettes from complex microbial communities. *bioRxiv*, 2021.2009.2008.459516, doi:10.1101/2021.09.08.459516 (2021).

38   Huerta-Cepas, J. *et al.* eggNOG 5.0: a hierarchical, functionally and phylogenetically annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic Acids Research* **47**, D309-D314 (2019).

39   El-Gebali, S. *et al.* The Pfam protein families database in 2019. *Nucleic Acids Research* **47**, D427-D432 (2019).

40   Cambray, G., Guerout, A.-M. & Mazel, D. Integrons. *Annual Review of Genetics* **44**, 141-166 (2010).

41   Iqbal, N., Guérout, A.-M., Krin, E., Le Roux, F. & Mazel, D. Comprehensive functional analysis of the 18 *Vibrio cholerae* N16961 toxin-antitoxin systems substantiates their role in stabilizing the superintegron. *Journal of Bacteriology* **197**, 2150-2159 (2015).

42   Szekeres, S., Dauti, M., Wilde, C., Mazel, D. & Rowe-Magnus, D. A. Chromosomal toxin–antitoxin loci can diminish large-scale genome reductions in the absence of selection. *Molecular Microbiology* **63**, 1588-1605 (2007).

524  43  Wilbaux, M., Mine, N., Guérout, A.-M., Mazel, D. & Van Melderen, L. Functional
525      interactions between coexisting toxin-antitoxin systems of the *ccd* family in
526      *Escherichia coli* O157: H7. *Journal of Bacteriology* **189**, 2712-2719 (2007).
527  44  Guérout, A.-M. *et al.* Characterization of the *phd-doc* and *ccd* toxin-antitoxin
528      cassettes from *Vibrio* superintegrons. *Journal of Bacteriology* **195**, 2270-2283 (2013).
529  45  Nayfach, S., Shi, Z. J., Seshadri, R., Pollard, K. S. & Kyrpides, N. C. New insights
530      from uncultivated genomes of the global human gut microbiome. *Nature* **568**, 505-
531      510, doi:10.1038/s41586-019-1058-x (2019).
532  46  Wu, D., Jospin, G. & Eisen, J. A. Systematic identification of gene families for use as
533      "markers" for phylogenetic and phylogeny-driven ecological studies of bacteria and
534      archaea and their major subgroups. *PLoS One* **8**, e77033 (2013).
535  47  Truong, D. T. *et al.* MetaPhlAn2 for enhanced metagenomic taxonomic profiling.
536      *Nature Methods* **12**, 902-903 (2015).
537  48  Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W.
538      CheckM: assessing the quality of microbial genomes recovered from isolates, single
539      cells, and metagenomes. *Genome Research* **25**, 1043-1055,
540      doi:10.1101/gr.186072.114 (2015).
541  49  Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome phylogeny
542      substantially revises the tree of life. *Nature Biotechnology* **36**, 996-1004,
543      doi:10.1038/nbt.4229 (2018).
544  50  Parks, D. H. *et al.* A complete domain-to-species taxonomy for Bacteria and Archaea.
545      *Nature Biotechnology* **38**, 1079-1086, doi:10.1038/s41587-020-0501-8 (2020).
546  51  Rinke, C. *et al.* A standardized archaeal taxonomy for the Genome Taxonomy
547      Database. *Nature Microbiology* **6**, 946-959, doi:10.1038/s41564-021-00918-8 (2021).
548  52  Bowers, R. M. *et al.* Minimum information about a single amplified genome
549      (MISAG) and a metagenome-assembled genome (MIMAG) of bacteria and archaea.
550      *Nature Biotechnology* **35**, 725-731, doi:10.1038/nbt.3893 (2017).
551  53  Chaumeil, P.-A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-Tk: a toolkit to
552      classify genomes with the Genome Taxonomy Database. *Bioinformatics* **36**, 1925-
553      1927, doi:10.1093/bioinformatics/btz848 (2019).
554  54  Nguyen, L.-T., Schmidt, H. A., Von Haeseler, A. & Minh, B. Q. IQ-TREE: a fast and
555      effective stochastic algorithm for estimating maximum-likelihood phylogenies.
556      *Molecular Biology and Evolution* **32**, 268-274 (2015).
557  55  Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K., Von Haeseler, A. & Jermiin, L. S.
558      ModelFinder: fast model selection for accurate phylogenetic estimates. *Nature*
559      *Methods* **14**, 587-589 (2017).
560  56  Pereira, M. B., Wallroth, M., Kristiansson, E. & Axelson-Fisk, M. HattCI: fast and
561      accurate *attC* site identification using hidden Markov models. *Journal of*
562      *Computational Biology* **23**, 891-902 (2016).
563  57  Nawrocki, E. P. & Eddy, S. R. Infernal 1.1: 100-fold faster RNA homology searches.
564      *Bioinformatics* **29**, 2933-2935 (2013).
565  58  Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version
566      7: improvements in performance and usability. *Molecular Biology and Evolution* **30**,
567      772-780 (2013).
568  59  Engelhardt, J., Heyne, S., Will, S. & Reiche, K. *RNAclust Documentation*,
569      <http://www.bioinf.uni-leipzig.de/~kristin/Software/RNAclust/manual.pdf> (2010).
570  60  Will, S., Joshi, T., Hofacker, I. L., Stadler, P. F. & Backofen, R. LocARNA-P:
571      accurate boundary prediction and improved detection of structural RNAs. *RNA* **18**,
572      900-914 (2012).

573   61   Will, S., Reiche, K., Hofacker, I. L., Stadler, P. F. & Backofen, R. Inferring
574        noncoding RNA families and classes by means of genome-scale structure-based
575        clustering. *PLoS Computational Biology* **3**, e65 (2007).
576   62   Jones, P. *et al.* InterProScan 5: genome-scale protein function classification.
577        *Bioinformatics* **30**, 1236-1240 (2014).
578   63   Huerta-Cepas, J. *et al.* Fast genome-wide functional annotation through orthology
579        assignment by eggNOG-mapper. *Molecular Biology and Evolution* **34**, 2115-2122
580        (2017).
581   64   Cantalapiedra, C. P., Hernández-Plaza, A., Letunic, I., Bork, P. & Huerta-Cepas, J.
582        eggNOG-mapper v2: functional annotation, orthology assignments, and domain
583        prediction at the metagenomic scale. *Molecular Biology and Evolution*,
584        doi:10.1093/molbev/msab293 (2021).
585   65   Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using
586        DIAMOND. *Nature Methods* **12**, 59-60 (2015).
587   66   Almagro Armenteros, J. J. *et al.* SignalP 5.0 improves signal peptide predictions using
588        deep neural networks. *Nature Biotechnology* **37**, 420-423, doi:10.1038/s41587-019-
589        0036-z (2019).
590   67   Penesyan, A., Nagy, S. S., Kjelleberg, S., Gillings, M. R. & Paulsen, I. T. Rapid
591        microevolution of biofilm cells in response to antibiotics. *npj Biofilms and
592        Microbiomes* **5**, 34, doi:10.1038/s41522-019-0108-3 (2019).
593   68   Harrell, F. E. & Dupont, C. Hmisc: harrell miscellaneous. R package version 4.5-0.
594        https://CRAN.R-project.org/package=Hmisc.  (2021).
595   69   Jacomy, M., Venturini, T., Heymann, S. & Bastian, M. ForceAtlas2, a continuous
596        graph layout algorithm for handy network visualization designed for the Gephi
597        software. *PLoS One* **9**, e98679 (2014).
598   70   Bastian, M., Heymann, S. & Jacomy, M. in *Third International AAAI Conference on
599        Weblogs and Social Media*.
600   71   Dehio, C. & Meyer, M. Maintenance of broad-host-range incompatibility group P and
601        group Q plasmids and transposition of Tn*5* in *Bartonella henselae* following conjugal
602        plasmid transfer from *Escherichia coli*. *Journal of Bacteriology* **179**, 538-540,
603        doi:doi:10.1128/jb.179.2.538-540.1997 (1997).
604   72   Penfold, R. J. & Pemberton, J. M. An improved suicide vector for construction of
605        chromosomal insertion mutations in bacteria. *Gene* **118**, 145-146,
606        doi:https://doi.org/10.1016/0378-1119(92)90263-O (1992).
607   73   Riedel, T., Rohlfs, M., Buchholz, I., Wagner-Döbler, I. & Reck, M. Complete
608        sequence of the suicide vector pJP5603. *Plasmid* **69**, 104-107,
609        doi:10.1016/j.plasmid.2012.07.005 (2013).
610   74   Sanchez, J., Bennett, P. M. & Richmond, M. H. Expression of *elt*-B, the gene
611        encoding the B subunit of the heat-labile enterotoxin of *Escherichia coli*, when cloned
612        in pACYC184. *FEMS Microbiology Letters* **14**, 1-5, doi:10.1111/j.1574-
613        6968.1982.tb08623.x (1982).
614   75   Guzman, L. M., Belin, D., Carson, M. J. & Beckwith, J. Tight regulation, modulation,
615        and high-level expression by vectors containing the arabinose PBAD promoter.
616        *Journal of Bacteriology* **177**, 4121-4130, doi:10.1128/jb.177.14.4121-4130.1995
617        (1995).
618   76   Rose, R. E. The nucleotide sequence of pACYC184. *Nucleic Acids Research* **16**, 355,
619        doi:10.1093/nar/16.1.355 (1988).
620   77   Avila, P. & de la Cruz, F. Physical and genetic map of the IncW plasmid R388.
621        *Plasmid* **20**, 155-157, doi:10.1016/0147-619x(88)90019-4 (1988).

622  78    Vit, C., Loot, C., Escudero, J. A., Nivina, A. & Mazel, D. in *Horizontal Gene*
623        *Transfer: Methods and Protocols*   (ed Fernando de la Cruz)  189-208 (Springer US,
624        2020).
625  79    Vit, C. *et al.* Cassette recruitment in the chromosomal Integron of *Vibrio cholerae*.
626        *Nucleic Acids Research* **49**, 5654-5670, doi:10.1093/nar/gkab412 (2021).
627  80    Moura, A. *et al.* INTEGRALL: a database and search engine for integrons, integrases
628        and gene cassettes. *Bioinformatics* **25**, 1096-1098, doi:10.1093/bioinformatics/btp105
629        (2009).
630

631

## Acknowledgements

636

## Author contributions

TMG contributed to the conception of the study, performed all data analyses, wrote the

original draft of the paper, and contributed to the final editing of the paper. SGT contributed

to the conception of the study and the final editing of the paper. AP performed the correlation

analysis of cassette functions, and contributed to the final editing of the paper. QQ was

involved with the design and implementation of the experimental work, and contributed to

the final editing of the paper. VR was involved with the design and implementation of the

experimental work, and contributed to the final editing of the paper. MRG contributed to the

conception of the study and the final editing and revision of the paper. All authors

contributed to the article and approved the final submitted version.

647

## Competing interests

The authors declare no competing interests.

650

651 **Materials & Correspondence**

652 Correspondence and material requests should be addressed to TMG
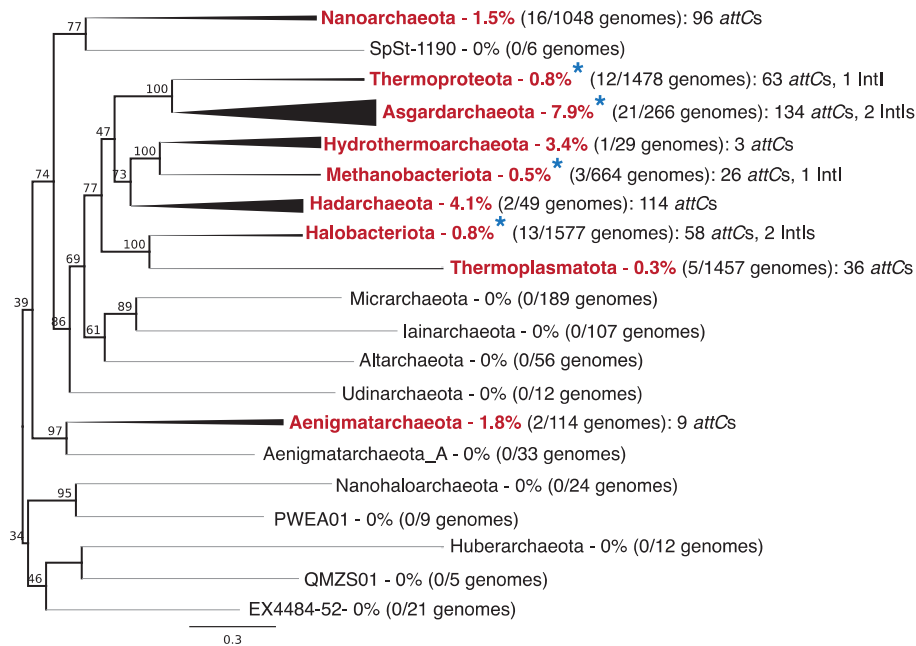
653 (timothy.ghaly@mq.edu.au).

654

655

656 **Figures**

657



658

659 **Fig. 1: Phylogenetic distribution of integrons among Archaea**. Archaeal phyla found to
660 carry integrons are labelled in red, and those found to have an integron integrase gene (*intI*)
661 are denoted with blue asterisks. Branch thickness indicates the proportion of genomes with
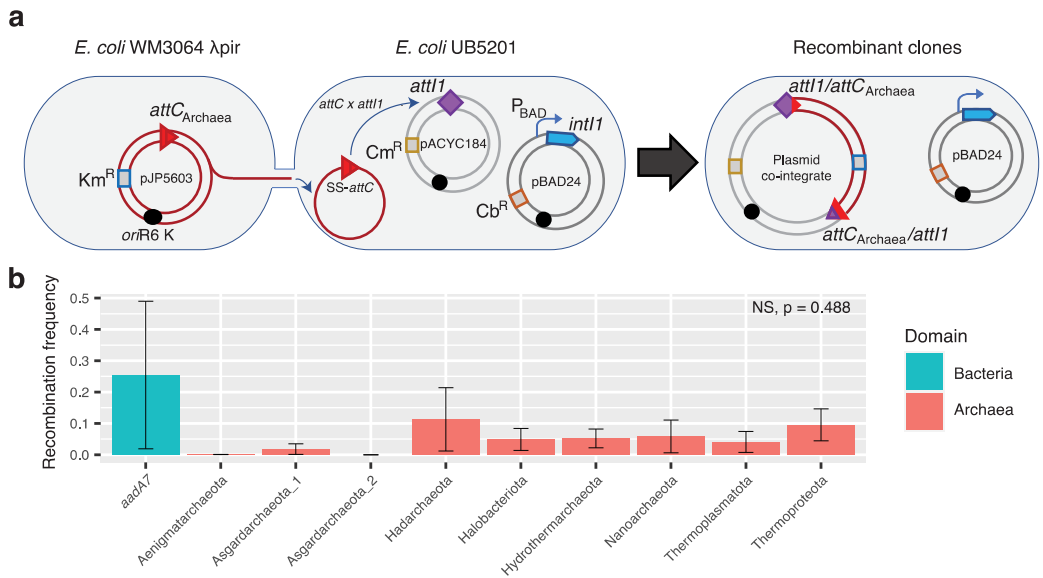662 integrons for each phylum.

663

664



665

666 **Fig. 2: Cassette recruitment (*attC* x *attI* recombination) assays. a,** schematic outlining the
667 experimental setup of the cassette insertion assays. The suicide vector pJP5603 with an *attC*
668 site is delivered into the recipient *E. coli* UB5201 strain via conjugation. The recipient strain
669 carries an *intI1* gene, expressed from the inducible P_BAD promoter, and an *attI1* site, residing
670 on the pBAD24 and pACYC184 backbones, respectively. The donor suicide vector cannot
671 replicate within the recipient host, and thus, can only persist following *attC* x *attI*
672 recombination to form a plasmid co-integrate. **b,** average recombination frequencies (±1 S.E.)
673 between *attI1* and nine archaeal *attC*s (with phyla of origin labelled along the X-axis) and the

674 paradigmatic bacterial *attC* site (*attC$_{aadA7}$*), used as positive control. Average frequencies
675 were calculated following three independent cassette insertion assays (see Methods for
676 details). No statistically significant difference in recombination frequencies were detected
677 among the tested *attC*s (Kruskal-Wallis test, n=27, df=8, p=0.488). Recombination
678 frequencies are shown for *attC* bottom strands only. See Extended Data Table 1 for *attC* top
679 strand recombination frequencies.
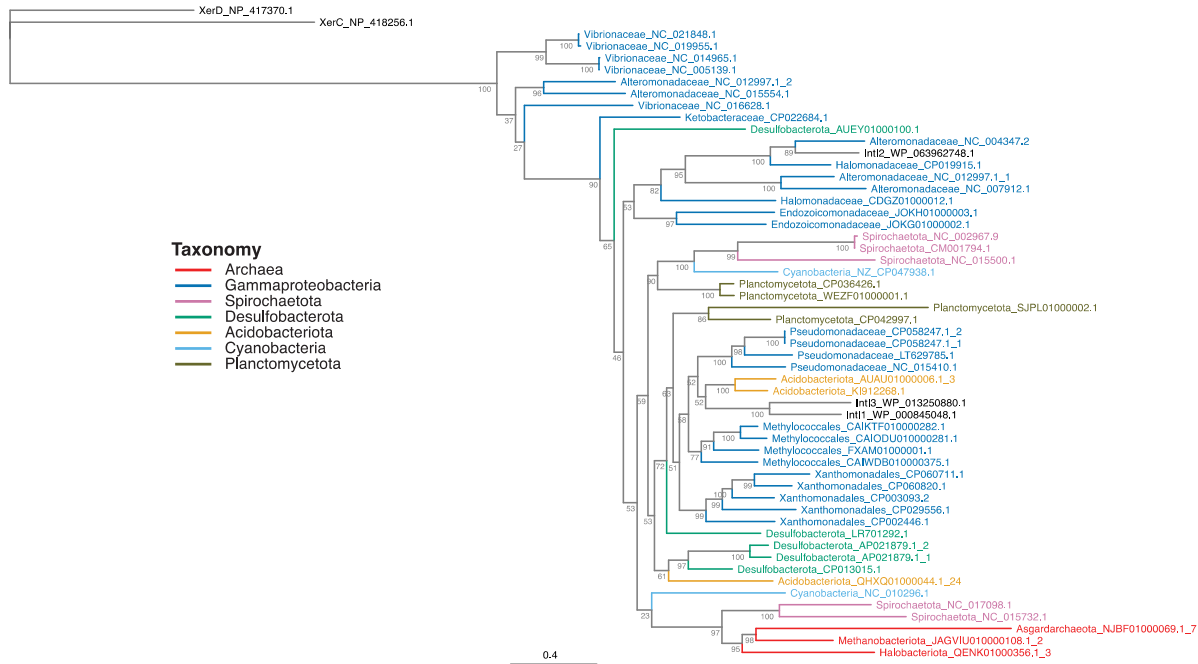
680
681
682
683
684



685
686 **Fig. 3: Phylogeny of integron integrases from Archaea and Bacteria.** To root the tree, the
687 tyrosine recombinases XerC and XerD from *Escherichia coli* were used as outgroups.
688 Integron integrases (IntIs) are coloured according to their taxonomy.
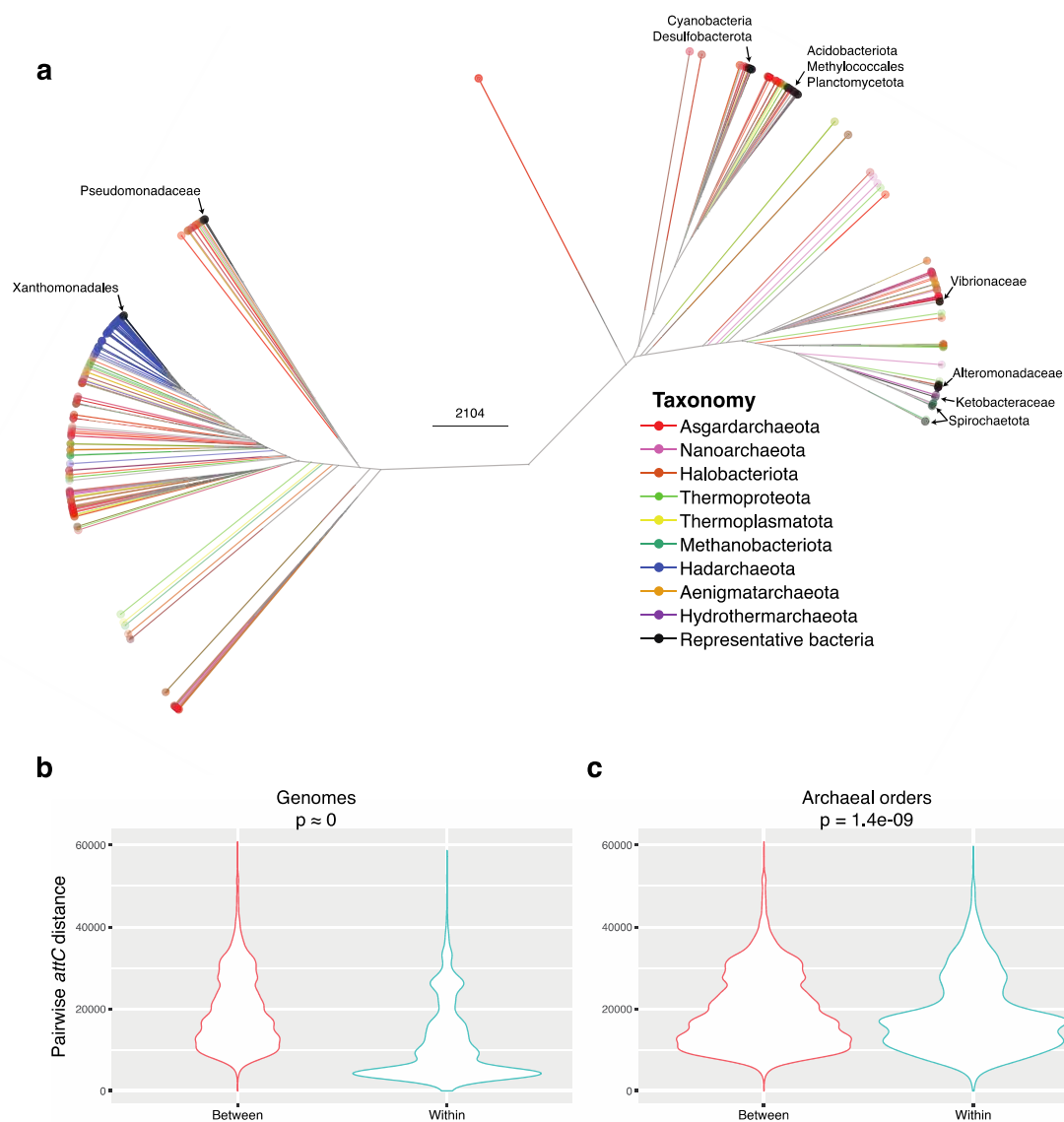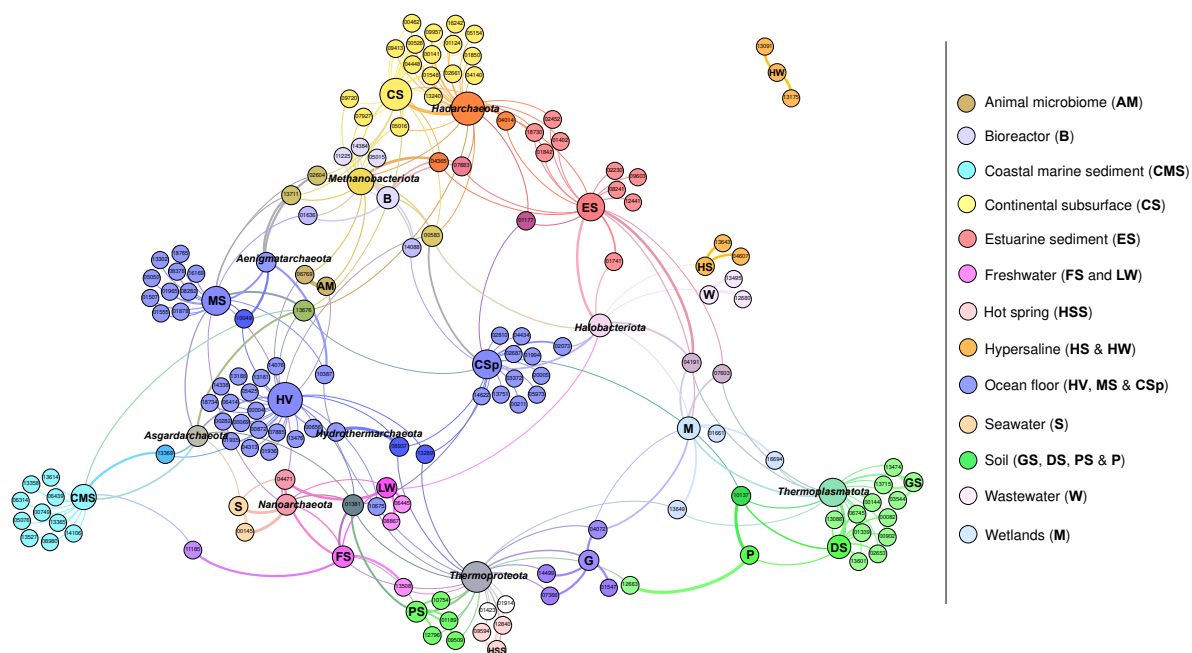689
690

**Fig. 4: Structural and sequence diversity of archaeal *attC* recombination sites. a,** structure-based clustering of all archaeal and representative bacterial *attC*s. Branches and tips are coloured according to archaeal phylum. The taxa of bacterial *attC*s are labelled with arrows. **b**, distribution of the sequence and structural distances calculated for all pairwise comparisons of *attC*s within and between genomes. **c**, distribution of distances for all pairwise comparisons of *attC*s from different genomes that are either from the same or different archaeal orders.

**AM**: Animal Microbiome; **B**: Bioreactor; **CMS**: Coastal Marine Sediment; **CSp**: Cold Seep; **CS**: Continental Subsurface; **DS**: Desert Soil; **ES**: Estuarine Sediment;
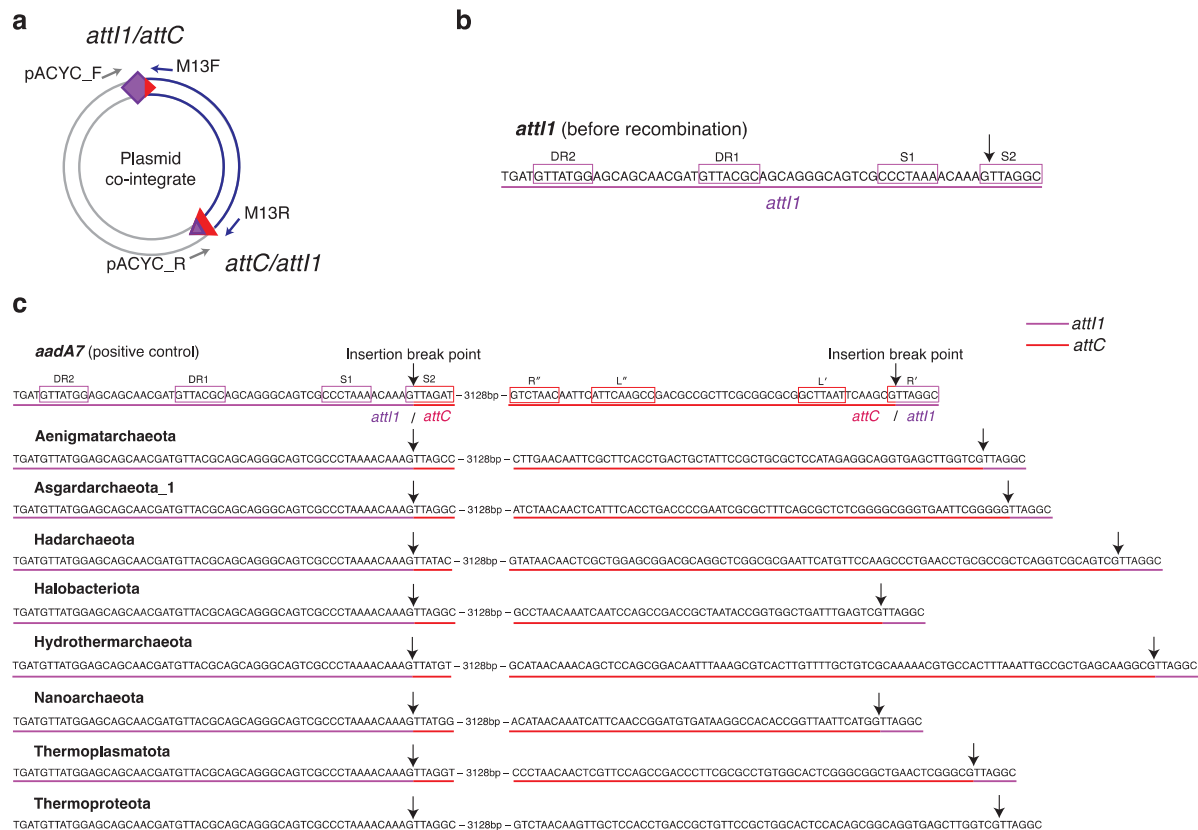**FS**: Freshwater Sediment; **G**: Groundwater; **GS**: Grassland Soil; **HSS**: Hot Spring Sediment; **HV**: Hydrothermal Vent; **HS**: Hypersaline Sediment; **HW**: Hypersaline Water;
**MS**: Marine Sediment; **LW**: Lake Water; **M**: Mangroves; **P**: Peatland; **PS**: Permafrost Soil; **S**: Seawater; **W**: Wastewater.

**Fig. 5: A network linking Pfam functions of archaeal integron gene cassettes with their taxonomic and environmental contexts.** The force-directed representation of the network is constructed based on co-occurrence patterns and correlations (p < 0.05) between Pfam functions, taxonomic groups, and specific environments from which the organisms were sampled. Nodes that represent taxonomic groups and specific environments are labelled accordingly. All other nodes denote Pfam functions and are labelled with a Pfam number preceded by 'PF'. Specific environments are grouped into broader environment types, each of which is coloured as per the panel. Pfams directly linked to specific environment types are coloured in corresponding colours. Pfams linked to more than one environment type are coloured in overlapping colours. The size of the node is relative to the node authority based on the degree of correlations. Edges (the lines connecting the nodes) represent correlations between nodes. Edge colour denotes the overlapping colour of the two nodes it connects. Edge thickness represents the strength of correlation. The full description of all correlations and Pfam functions is presented in Supplementary Table 6.
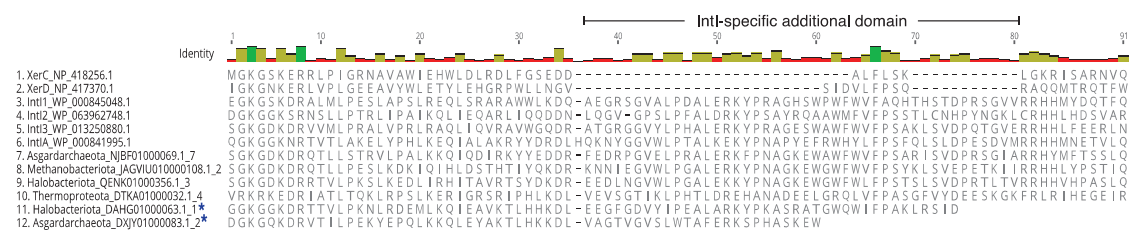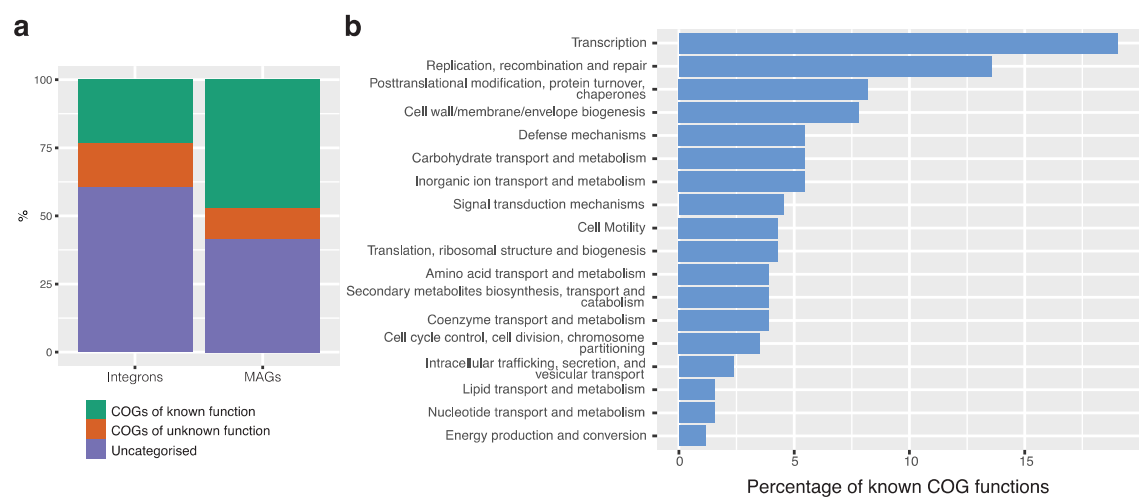
719 **Extended data**

720

721 **Extended Data Table 1. Average recombination frequencies for the *attC* x *attI* suicide**
722 **conjugation assays.**

| | *attC* bottom strand (*intI1* induced*) | *attC* bottom strand (*intI1* suppressed†) | *attC* top strand (*intI1* induced) | *attC* top strand (*intI1* suppressed) |
|---|---|---|---|---|
| $attC_{aadA7}$ | $2.54 \times 10^{-1}$ | ND†† | $2.48 \times 10^{-3}$ | ND |
| $attC_{Aenigmatarchaeota}$ | $5.46 \times 10^{-4}$ | ND | $8.66 \times 10^{-7}$ | ND |
| $attC_{Asgardarchaeota\_1}$ | $1.79 \times 10^{-2}$ | ND | $6.74 \times 10^{-4}$ | ND |
| $attC_{Asgardarchaeota\_2}$ | ND | ND | ND | ND |
| $attC_{Hadarchaeota}$ | $1.13 \times 10^{-1}$ | ND | $1.18 \times 10^{-3}$ | ND |
| $attC_{Halobacteriota}$ | $4.88 \times 10^{-2}$ | ND | $4.33 \times 10^{-4}$ | ND |
| $attC_{Hydrothermarchaeota}$ | $5.21 \times 10^{-2}$ | ND | $6.92 \times 10^{-3}$ | ND |
| $attC_{Nanoarchaeota}$ | $5.84 \times 10^{-2}$ | ND | $1.55 \times 10^{-3}$ | ND |
| $attC_{Thermoplasmatota}$ | $4.08 \times 10^{-2}$ | ND | $2.28 \times 10^{-3}$ | ND |
| $attC_{Thermoproteota}$ | $9.54 \times 10^{-2}$ | ND | $1.80 \times 10^{-3}$ | ND |

723 *induced using L-arabinose; †suppressed using D-glucose; ††ND = Not detected
724
725
726
727
728



729
730 **Extended Data Fig. 1: Example structure of archaeal integrons.** Maps of all 'complete
731 integrons', which are those that comprise an integron integrase gene (*intI*) and at least one
732 gene cassette recombination site (*attC*); all 'In0' elements, which are those with *intI* but no
733 detectable *attC* site; and three examples of 'CALINs' (clusters of *attC*s lacking integron
734 integrases).
735
736

737



**Extended Data Fig. 2: Putative archaeal integron recombination sites, *attI*s. a**, maps showing the location of putative archaeal *attI*s. **b**, sequence alignment of the two putative archaeal *attI*s. **c**, multiple sequence alignment of the two archaeal *attI*s and all annotated bacterial *attI*s from the INTEGRALL database[80]. Nucleotides are coloured if they match with at least 50% of the sequences. Vertical arrows indicate the canonical insertion point of an inserting gene cassette.

**Extended Data Fig. 3: Sanger sequencing of *attI1* x *attC* recombination junctions. a,** schematic of PCR primer pairs (grey and blue arrows) that amplify the recombination junctions following cassette insertion (*attI1* x *attC* recombination). **b,** *attI1* sequence before recombination. Boxes denoted with S1 and S2 indicate the core IntI1 binding sites, and the direct repeats signified by DR1 and DR2, are additional strong and weak IntI1 binding sites, respectively. The black arrow indicates the insertion break point where cleavage takes place during recombination. **c,** Sanger sequence data of the recombinant clones following *attI1* recombination with the paradigmatic bacterial *attC* site (*attC_aadA7*), used as positive control, and eight archaeal *attC*s. Black arrows indicate the insertion break points following recombination. For *attC_aadA7*, the two sets of paired inverted repeats are boxed (R′ to R″ and L′ to L″).



**Extended Data Fig. 4: A multiple protein sequence alignment of the additional domain unique to integron integrases.** Sequences (1) and (2) are tyrosine recombinases XerC and XerD that lack the IntI-specific domain. Sequences (3) to (6) are bacterial IntIs, and (7) to (12) are IntIs from Archaea. Blue asterisks indicate IntIs that did not span the full additional domain and were excluded from phylogenetic analysis.

**Extended Data Fig. 5: COG functional analysis of archaeal gene cassettes. a**, percentage of proteins assigned a COG category. 'Integrons' represent all cassette-encode proteins in Archaea, while 'MAGs' indicate all proteins from the 75 integron-bearing archaeal genomes. **b**, percentage of COGs with known functions assigned archaeal cassette-encoded proteins.