

# Identification of Essential Protein Domains From High-density Transposon Insertion Sequencing

**A.S.M. Zisanur Rahman**

University of Manitoba

**Lukas Timmerman**

University of Manitoba

**Flyn Gallardo**

University of Manitoba

**Silvia T. Cardona** (✉ [silvia.cardona@umanitoba.ca](mailto:silvia.cardona@umanitoba.ca))

University of Manitoba

---

## Research Article

**Keywords:** DUFs, Essential genes, Tn-seq, Protein Domain

**Posted Date:** June 17th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-589027/v1>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

1 **Identification of Essential Protein Domains from High-density Transposon**  
2 **Insertion Sequencing**

3 A. S. M. Zisanur Rahman<sup>1</sup>, Lukas Timmerman<sup>2</sup>, Flyn Gallardo<sup>1</sup>, Silvia T. Cardona<sup>1,3\*</sup>.

4

5 <sup>1</sup>Department of Microbiology, University of Manitoba, Winnipeg, MB, Canada.

6 <sup>2</sup>Department of Computer Science, University of Manitoba, Winnipeg, MB, Canada.

7 <sup>3</sup>Department of Medical Microbiology & Infectious Diseases, University of Manitoba,  
8 Winnipeg, Canada.

9 \*To whom correspondence should be addressed: [silvia.cardona@umanitoba.ca](mailto:silvia.cardona@umanitoba.ca)

10

11

12

13

14

15 **Abstract**

16           A first clue to gene function can be obtained by examining whether a gene is required  
17 for life in certain standard conditions, that is, whether a gene is essential. In bacteria, essential  
18 genes are usually identified by high-density transposon mutagenesis followed by sequencing  
19 of insertion sites (Tn-seq). These studies assign the term “essential” to whole genes rather than  
20 the protein domain sequences that confer the essential functions. However, genes can code for  
21 multiple protein domains that evolve their functions independently. Therefore, when essential  
22 genes code for more than one protein domain, only one of them could be essential. In this study,  
23 we defined this subset of genes as “essential domain-containing” (EDC) genes. Using a Tn-seq  
24 data set built-in *Burkholderia cenocepacia* K56-2, we developed an *in silico* pipeline to identify  
25 EDC genes and the essential protein domains they encode. We found forty candidate EDC  
26 genes and demonstrated growth defect phenotypes using CRISPR interference (CRISPRi).  
27 This analysis included two knockdowns of genes encoding the protein domains of unknown  
28 function DUF2213 and DUF4148. These essential domains are conserved in more than two  
29 hundred bacterial species, including human and plant pathogens. Together, our study suggests  
30 that essentiality should be assigned to individual protein domains rather than genes,  
31 contributing to a first functional characterization of protein domains of unknown function.

32

33 **Keywords**

34 DUFs, Essential genes, Tn-seq, Protein Domain

35 **Author contributions**

36 ASMZR- performed the majority of the experiments and wrote the manuscript

37 LT- created the python script and contributed to manuscript editing

38 FG- created CRISPRi mutants and contributed to manuscript editing

39 STC- conceived the idea, supervised the work, provided financial support, and edited the final  
40 version of the manuscript.

41

## 42 **Introduction**

43           A first step when characterizing gene function should be asking whether a given gene  
44 encodes an essential cellular function, whether the gene is necessary for the survival of the  
45 organism. A widely accepted method to identify essential genes in bacteria is high-density  
46 transposon mutagenesis, followed by Illumina-sequencing of the transposon insertion junctions  
47 (Tn-seq)<sup>1</sup>. During Tn-seq, transposon mutant cells are pooled and grown in optimal conditions,  
48 allowing cells with a transposon insertion located in a non-essential element to survive. Cells  
49 with a transposon insertion in an essential element should be lost or depleted from the  
50 population. When transposon insertions are identified by Illumina sequencing, read counts per  
51 gene in the central 70-90% of the open reading frame (disruptive insertions) are normalized by  
52 gene length and used to predict essentiality. 5-15% sequences from the 3' and 5' ends are  
53 usually removed from the analysis, as insertions within the terminal regions are likely non-  
54 disruptive<sup>2-5</sup>. While disrupted genes are regarded as “non-essential,” the method yields a list  
55 of putative essential genes as those with zero or very few mapped reads (**Figure 1a and b**)<sup>3</sup>.

56           Another step towards identifying gene function is the annotation of the protein domains  
57 encoded by genes. Protein domains are functional or structural units that can fold, evolve, and  
58 function independently. Homology-based protein domain prediction and function assignment  
59 are effective starting points for understanding protein function, even when diverse protein  
60 architectures add complexity to functional annotations<sup>6,7</sup>. While domain databases such as  
61 Pfam<sup>8</sup> and InterPro<sup>9</sup> aim to provide maximum sequence coverage to predict protein domain  
62 identity, approximately 30% of all domains listed in these databases (Pfam 33.1 and InterPro  
63 81.0) are ‘domains of unknown function (DUFs).’ Single DUFs are usually predicted to span  
64 through functionally uncharacterized proteins. However, studies suggest that at least some of  
65 these proteins may contain more than one domain<sup>10,11</sup>.

66 While robust and comprehensive, Tn-seq studies do not consider that genes may encode  
67 for more than one protein domain. Tn-seq analysis may classify a gene as “non-essential” due  
68 to the presence of transposon insertions in a non-essential coding region, despite the gene  
69 coding for a second domain not spanning through the whole gene length that might be essential  
70 <sup>3,12,13</sup>. We operationally defined this subclass of essential genes as “essential domain-  
71 containing” (EDC) genes (**Figure 1c and d**) and set out to identify them in a Tn-seq dataset  
72 built-in *Burkholderia cenocepacia* K56-2 <sup>14</sup>. By analyzing biases in transposon density in  
73 genes previously identified as “non-essential”, we found 40 genes where the encoded proteins  
74 contained putative essential and non-essential domains. Using a CRISPR Interference  
75 (CRISPRi) <sup>15</sup> platform we developed for *Burkholderia* <sup>16</sup>, we experimentally confirmed growth  
76 defects, representing the loss of an essential function, in 27 EDC gene knockdowns. The  
77 identified EDC genes include ten encoding known multidomain proteins and two entirely  
78 uncharacterized genes encoding different N-terminal DUFs, demonstrating the utility of the  
79 approach. This study highlights that gene essentiality depends on the function of individual  
80 protein domains rather than entire proteins.

## 81 **Results**

### 82 *Identification of EDC Genes from Tn-seq Data*

83 To identify EDC genes in *B. cenocepacia* K56-2, we built a custom script that used our  
84 previous Tn-seq data <sup>14</sup> to select genes that i) were not previously found to be essential in *B.*  
85 *cenocepacia* K56-2 <sup>14</sup>, and ii) had an asymmetric distribution of transposon insertions (**Figure**  
86 **2**). The script split each gene into two equal parts and selected genes with reads in only one  
87 region to identify genes with transposon insertion biases. We worked under the assumption that  
88 (i) each half could represent one functional domain and (ii) one of the domains may be essential  
89 while the other may not. We arbitrarily set the parameters “min ratio” and “min reads” to 0 and  
90 0.14, respectively (see Material and Methods and **Supplementary Figure 1**). These settings

91 looked for genes that had zero reads at one end, while the number of reads in the non-empty  
92 end was at least 14% of that region's length. For example, if a section of a gene was 100 bp in  
93 length, it would require at least 14 reads mapped to that section to be considered non-essential.  
94 With these settings, the script produced an extensive list of 178 candidate EDC genes  
95 (**Supplementary Table 1**).

#### 96 *Bioinformatic Analysis of the Candidate EDC Genes*

97 We reasoned that if EDC genes contained essential protein domains, then the essential  
98 protein domains may be encoded by essential genes in at least some other bacteria. We then  
99 searched for essential ortholog genes of the 178 candidate EDC genes by BLASTx searches  
100 against the 'Database of Essential Genes (DEG)'<sup>17</sup> using 50% sequence alignment and 30%  
101 sequence identity as the cut-off. We found that 40 of the 178 genes had orthologs annotated as  
102 'essential' in other bacterial species. We wished to interrogate the domains encoded by these  
103 40 genes using UniProt<sup>18</sup> based on InterPro domains<sup>9</sup>. InterPro predicts the domain  
104 information by matching the protein or nucleic acid sequences against the member databases  
105 (collectively known as InterPro consortium) to identify 'signatures' associated with known  
106 domains. Thus, the InterPro prediction relies on the availability of sequence characterization  
107 and annotation. This analysis showed that from the 40 candidate EDC genes predicted to be  
108 essential by homology with other essential genes, 10 genes encoded multidomain proteins, and  
109 7 of them were well-characterized, such as the N-terminal domain of DnaK and NusA (Figure  
110 3a). The remaining genes were predicted to have one single annotated domain (19 genes) that  
111 did not span the whole gene-length or encoded uncharacterized proteins (11 genes)  
112 (**Supplementary Table 2**). All 40 genes had transposon insertions located in one half of the  
113 gene, showing that the script was able to identify genes with biased transposon insertions  
114 (**Supplementary Figure 2**). Taken together, these results suggest that the identified genes

115 could be essential due to the presence of essential protein domains orthologues. Notably, 17

116 DNA regions were identified as coding for new putative essential protein domains (**Table 1**).

117 **Table 1:** Essential genes and domains identified based on biased transposon insertions.

<b>K56-2 Locus Tag</b>	<b>Homolog J2315 Locus Tag</b>	<b>Product Name</b>	<b>Function</b>	<b>Reads at 5' Half</b>	<b>Reads at 3' Half</b>	<b>Identified Essential Domain</b>
WQ49_RS0 0050	BCAL3469	cell division protein FtsL	Essential cell division protein	0	23	Domain (FtsL)
WQ49_RS0 0770	BCAL3328	NUDIX hydrolase	Nucleoside- diphosphatase	0	49	Domain (Nudix hydrolase)
WQ49_RS0 0885	BCAL3305	preprotein translocase subunit YajC	Secretase/inse rtase	21	0	new
WQ49_RS0 1035	BCAL3270	DnaK	Chaperone	0	227	N-terminal Domain
WQ49_RS0 2920	BCAM1451	hypothetical protein	Unknown	43	0	new
WQ49_RS0 3160	BCAM1502	hypothetical protein	Unknown	59	0	new
WQ49_RS0 3550	QU43_RS6 2245	hypothetical protein	Unknown	33	0	new
WQ49_RS0 3805	BCAM1624	MaoC family dehydratase	MaoC-like dehydratase	46	0	new
WQ49_RS0 4450	BCAM1749	hypothetical protein	Unknown	17	0	new
WQ49_RS0 7360	BCAM2338	glycosyl transferase family 1	UDP- glycosyltransf erase	0	152	Domain (Glyco_transf_ 28)
WQ49_RS0 7395	QU43_RS6 6100	hypothetical protein	Unknown	0	58	new
WQ49_RS0 9185	BCAS0417	cytochrome biogenesis protein CcdA	electron transfer	0	38	new
WQ49_RS1 0495	BCAS0158	hypothetical protein	Unknown	0	34	Domain (DUF4148)
WQ49_RS1 1915	BCAL0324	TatB	protein transmembran e transporter	0	57	Domain (TatA_B_E)
WQ49_RS1 2045	BCAL0298	thiamine biosynthesis protein ThiS	thiamine biosynthesis protein ThiS	0	50	Domain (ThiS)

WQ49_RS1 2280	BCAL0250	50S ribosomal protein L18	structural constituent of ribosome	0	65	Domain (Ribosomal_L 18p)
WQ49_RS1 2305	BCAL0245	RplX	structural constituent of ribosome	20	0	Domain (L24- Pfam)
WQ49_RS1 2315	BCAL0243	30S ribosomal protein S17	structural constituent of ribosome	0	64	new
WQ49_RS1 2365	BCAL0233	RpsJ	structural constituent of ribosome	0	25	new
WQ49_RS1 6145	BCAM1066	hypothetical protein	Unknown	0	425	Domain (DUF2213)
WQ49_RS1 8705	BCAM0549	molecular chaperone GroES	Chaperone	0	21	Domain (Cpn10)
WQ49_RS2 2170	BCAM2699	alpha/beta hydrolase	Putative hydrolase	120	0	Domain (Abhydrolase_ 3)
WQ49_RS2 3945	BCAL0558	Cca	3'-cytidine- cytidine- tRNA adenylyltransf erase	0	79	Domain (PolyA Polymerase)/D omain (Binding)
WQ49_RS2 4070	BCAL0585	hypothetical protein	Unknown	0	23	new
WQ49_RS2 5525	BCAL0878	FmdB family transcriptional regulator	Regulatory activity	0	30	Domain (CxxC_CXXC _SSSS)
WQ49_RS2 5680	BCAL0909	16S rRNA maturation RNase YbeY	Endoribonucl ease activity	68	0	Domain (UPF0054)
WQ49_RS2 6625	BCAL2715	RpmG	structural constituent of ribosome	0	31	Domain (Ribosomal_L 33)
WQ49_RS2 7920	BCAL2334	NADH-quinone oxidoreductase subunit K	NADH dehydrogenas e	0	21	Domain (Oxidored_q2)
WQ49_RS2 8635	BCAL2199	Fe-S cluster assembly transcriptional regulator IscR	DNA-binding transcription factor	39	0	Domain (Rrf2)
WQ49_RS2 9230	BCAL2091	30S ribosomal protein S2	structural constituent of ribosome	0	86	Domain (Ribosomal_S2 )

WQ49_RS3 0770	BCAL1788	biopolymer transporter ExbD	Transmembrane transporter	0	47	Domain (ExbD)
WQ49_RS3 1735	NA	hypothetical protein	Unknown	0	42	new
WQ49_RS3 1805	BCAL1585	transcriptional regulator	DNA binding	44	0	new
WQ49_RS3 2210	BCAL1506	NusA	DNA-binding transcription factor	0	93	Domain (NusA_N)
WQ49_RS3 2225	BCAL1503	SMC-Scp complex	Cell Division/chromosome separation	0	94	Domain (SMC)
WQ49_RS3 2625	BCAL1424	ABC transporter	ATPase	63	0	new
WQ49_RS3 4660	BCAL0990	50S ribosomal protein L32	structural constituent of ribosome	27	0	new
WQ49_RS3 4895	BCAL2925	50S ribosomal protein L19	structural constituent of ribosome	0	26	Domain (Ribosomal_L 19)
WQ49_RS3 5060	BCAL2958	membrane protein	Porin activity	43	0	Domain (OmpA)
WQ49_RS0 3390	BCAM1545	LuxR family transcriptional regulator	DNA binding	251	0	Domain (HTH luxR-type)

118

119 *CRISPRi Knockdowns of EDC Genes Show Growth Defects*

120 To phenotypically characterize the effect of knocking down EDC genes, we used  
121 CRISPR interference or CRISPRi<sup>16</sup> to create knockdown mutants of the genes of interest.  
122 CRISPRi comprises a chromosomally integrated dCas9 under the control of a rhamnose-  
123 inducible promoter and plasmid-borne sgRNA driven by a constitutively active synthetic  
124 promoter, P<sub>J23119</sub><sup>16</sup>. Simultaneous expression of dCas9 and a target-specific sgRNA allows the  
125 dCas9 to bind the target DNA region and, thus, sterically interfere with transcription by RNA  
126 polymerase<sup>15,16</sup>. To inhibit the expression of the candidate genes, we designed two sgRNAs  
127 against each of the candidate genes targeting the start codon and adjacent region on the non-  
128 template strand (**Supplementary Figure 3a and c**). For phenotypic characterization, we grew

129 the cells in LB with and without rhamnose. Upon induction of dCas9 with rhamnose, 27 out of  
130 the 40 candidate genes showed at least 25% growth inhibition relative to the uninduced  
131 condition (**Supplementary Figure 3b-d**).

### 132 *DUF2213 and DUF4148 Appear to be Essential Domains*

133 The presence of DUFs is a common feature of hypothetical or uncharacterized proteins.  
134 To initiate functional characterization of DUFs, we focused on two genes containing DUF-  
135 coding sequences, which their respective CRISPRi mutants demonstrated a conditional growth  
136 defect (**Figure 3b**). WQ49\_RS16145 (BCAM1066) and WQ49\_RS10495 (BCAS0158)  
137 contain DUF2213 (Pfam accession PF09979) and DUF4148 (Pfam accession PF13663),  
138 respectively at the N-terminal end of the proteins (**Figure 3b**). BLAST searches of BCAM1066  
139 and BCAS0158 genes as a query against the DEG <sup>17</sup> showed that BCAM1066  
140 (WQ49\_RS16145) had 30% sequence similarity with *lysK* (B8GXH3) from *Caulobacter*  
141 *crescentus*, and BCAS0158 (WQ49\_RS10495) had a 52% sequence identity with a predicted  
142 amino acid permease (BPSS1112) from *Burkholderia pseudomallei* K96243 (data not shown).  
143 Mining of the Pfam database (<https://pfam.xfam.org/>) showed that these DUFs are well  
144 conserved across the bacterial species: DUF2213 is present in 209 bacterial species, including  
145 bacterial pathogens (*Acinetobacter baumannii*, *Enterobacter cloacae*, *Haemophilus influenzae*,  
146 *Burkholderia cepacia*, *Shigella flexneri*), plant pathogens (*Agrobacterium tumefaciens*), and  
147 biotechnologically relevant species (*Pseudomonas putida*) (**Figure 4a and Supplementary**  
148 **Table 5**). DUF4148 is found in 204 bacterial species, primarily in *Burkholderia* species (i. e.  
149 *Burkholderia cepacia*, *Burkholderia mallei*, *Burkholderia vietnamiensis*) and plant pathogens  
150 such as *Ralstonia solanacearum* (**Figure 4b and Supplementary Table 5**). DUF2213 is also  
151 present in many phage-related proteins (**Figure 4a**). Eight unique domain architectures were  
152 observed for proteins containing DUF2213 and five for DUF4148 (**Figure 4c-d**). DUF2213 is  
153 associated with another essential domain PF00293, a NUDIX hydrolase (**Figure 4c**). In other

154 proteins, DUF2213 is associated with the LPD3 domain (PF18798) and DUF1073 (PF06381)  
155 which is also conserved across bacterial species <sup>11</sup> (**Figure 4c**). On the other hand, Pfam  
156 analysis of DUF4148 shows that DUF4148 differs in domain length among species and is  
157 associated with the Pfam domain PF00144, known to confer resistance against  $\beta$ -lactams  
158 (**Figure 4d**) <sup>19</sup>. Nonetheless, the N-terminus was highly conserved, suggesting it is functionally  
159 significant. The Pfam-based analysis of species distribution also revealed that DUF2213 is  
160 present in six eukaryotic species (five metazoans and one fungal species), whereas DUF4148  
161 is present in five eukaryotic species (three viridiplantae species and two metazoan species).  
162 The widespread distribution of these DUFs indicates the functional importance of these  
163 essential domains, creating an impetus for further characterization.

## 164 **Discussion**

165 A first step in the functional characterization of proteins is performed through protein  
166 depletion and growth phenotype characterization. As multidomain proteins can perform  
167 multiple functions driven by the activity of their individual domains <sup>20</sup>, the function assigned  
168 to a gene product could indeed correspond to one of its domains and not to the whole protein.  
169 That is the case of essential genes identified by Tn-seq <sup>1</sup>. In standard Tn-seq analysis the  
170 condition of essentiality is assigned to genes and not to domains, resulting in incorrect  
171 classification of many essential genes as non-essential. Rather, essentiality assignment pipeline  
172 should be revised to analyze the essentiality of individual protein domains <sup>21</sup>. Indeed,  
173 essentiality can be assigned to individual domains of a multidomain protein rather than the  
174 entire protein <sup>12,13</sup>. In this work, we defined as essential-domain-containing (EDC) genes those  
175 genes that encode more than one protein domain, with one of the domains coding for an  
176 essential function. By analyzing a Tn-seq dataset <sup>14</sup> for transposon insertion biases, we show  
177 that standard Tn-seq analysis pipelines may miss EDC genes, whose detection often requires  
178 either manual curation or additional considerations <sup>22</sup>.

179 We validated our approach by identifying previously characterized multidomain  
180 essential proteins in which the essential function is assigned to one single domain. For instance,  
181 our analysis of biases in the Tn-seq dataset showed that the N-terminal domain of NusA<sup>23</sup> is  
182 sufficient to mediate the essential function, in agreement with previous work<sup>24</sup>. Similarly, the  
183 *B. cenocepacia* K56-2 *dnaK* gene was previously defined as non-essential<sup>14</sup>; however, we  
184 found that the Tn-seq reads mapped onto *dnaK* were biased toward the C-terminal domain  
185 (CTD), suggesting that only the NTD is necessary for its essential function. (**Figure 3b,**  
186 **Supplementary Figure 2**). DnaK is a multidomain protein and a master regulator of the  
187 chaperone network<sup>25</sup>. DnaK comprises an N-terminal ATPase domain (NTD) and a C-terminal  
188 substrate-binding domain (CTD)<sup>25</sup>. Perturbations either within the NTD that leads to the  
189 abrogation of the ATPase activity or within the conserved linker peptide that impairs the  
190 interdomain mechanistic interaction abrogate the *in vivo* activity of DnaK<sup>26,27</sup>.

191 While 14 EDC genes that demonstrated a growth defect when knocked down code for  
192 proteins annotated to have a single domain, none of these domains span the entire gene, and  
193 transposon insertions are only mapped to the annotated domain (**Supplementary Figure 2**).  
194 Thus, it is possible that the remaining regions code for novel domains that perform the essential  
195 biological functions independently of the adjacent sequences. Indeed, multidomain proteins  
196 that are involved in direct protein-protein interactions are more often detected as essential than  
197 proteins with a single domain<sup>12</sup>, hinting towards the functional contribution of individual  
198 domains within a protein complex.

199 We demonstrated a conditional growth defect in 27 out of 40 CRISPRi mutants of EDC  
200 genes. It remains a possibility that the sgRNAs designed for CRISPRi-mediated gene silencing  
201 of the remaining 13 genes were not efficient in target binding, thus yielding no growth defect.  
202 CRISPRi is more effective in blocking transcription initiation than elongation, and is the most  
203 efficient in silencing gene expression when promoter regions are targeted with gRNAs<sup>15,28-30</sup>.

204 However, as promoter regions for *B. cenocepacia* genomes remained largely unannotated we  
205 targeted translation start sites. It remains to be investigated whether targeting the promoter  
206 region to block the transcription initiation rather than elongation might yield conditional a  
207 growth phenotype in the remaining 13 genes.

208 A large portion of the protein domains that lack functional assignment can be grouped  
209 within the DUF category. DUFs are members of ever-increasing uncharacterized protein  
210 families; they are the object of experimental and computational efforts towards their functional  
211 characterization<sup>10,31–33</sup>. Determining if a DUF is essential is among the first steps in functional  
212 characterization. In this study, we focused on two EDC genes that encode putative essential  
213 DUFs: DUF2213 and DUF4148. Both domains have a high degree of conservation across  
214 diverse phyla, which highlights their biological relevance. DUF2213, a phage-associated  
215 domain (PF09979), is well distributed across bacteria and phages. Interestingly, we found that  
216 DUF4148 (PF13663) is putatively essential and associated with  $\beta$ -lactamase (PF00144)  
217 (**Figure 4**).

218 In summary, our study identified 27 EDC genes whose knockdown produced a growth  
219 defect, highlighting the essential nature of one of their protein domains. By leveraging a Tn-  
220 Seq dataset in *B. cenocepacia* K56-2<sup>14</sup>, we demonstrate that the essential nature of protein-  
221 coding genes is a function of the individual protein domains they encode. We propose that  
222 determining essentiality of a domain of unknown function should be the first step in the process  
223 to define their function.

## 224 **Methods**

### 225 *Bacterial Strains and Growth Conditions*

226 The list of bacterial strains and plasmids used in this study is provided in  
227 **Supplementary Table 3**. Bacterial strains were grown in LB-Lennox medium (Difco) at 37°C.  
228 *E. coli* strain MM290 carrying the helper plasmid pRK2013 was selected in kanamycin

229 40µg/mL (Fisher Scientific). Donor strains of *E. coli* DH5α and *B. cenocepacia* K56-2 carrying  
230 the sgRNA plasmids were selected in trimethoprim 50µg/mL and 100µg/mL (Sigma),  
231 respectively.

### 232 *Identification of EDC Genes from Tn-Seq Dataset*

233 Candidate EDC genes were identified with a custom python script using the Tn-seq  
234 dataset<sup>14</sup>. The script analyzed every gene previously classified as “non-essential” by splitting  
235 it into two equal halves and counting the number of reads mapped to each half-gene. The script  
236 then used the “min ratio” and “min reads” as filtering criteria to call EDC genes. “Min ratio”  
237 was defined as the desired ratio of reads between the halves of the gene. “Min reads” was  
238 defined as the minimum number of reads in the non-empty end that is equal to a 14% of that  
239 half's length. Min reads was set to 0.14, while min ratio was set as 0. For each gene, 10% from  
240 each end of the gene was discarded from the analysis. The parameters can be changed to yield  
241 either more stringent or more general results. The script is available at  
242 <https://github.com/cardonalab/EssentialDomains>

### 243 *Bioinformatic Analysis*

244 Orthologous essential genes were identified using BLASTx against DEG<sup>15</sup>.  
245 Multidomain information was fetched from the UniProt database based on Pfam<sup>8</sup> and InterPro<sup>9</sup>  
246 domain features. DUF containing genes were characterized using the Pfam tool available on  
247 the Pfam website (<https://pfam.xfam.org/>). Domain sequences were retrieved in FASTA format  
248 from the Pfam database<sup>8</sup> and aligned by Clustal Ω<sup>34</sup>. Maximum-likelihood phylogenetic trees  
249 were generated with MEGA-X<sup>35</sup> using a Jones-Taylor-Thornton (JTT)-based model<sup>36</sup>  
250 applying 100 bootstrap values. Phylogenetic trees were visualized, edited and taxonomic labels  
251 were assigned using Interactive Tree Of Life (i-TOL)<sup>37</sup>. Bootstrap values are represented on a  
252 scale of 0 to 1. Taxonomic annotations were labelled based on the NCBI taxonomy database  
253 using UniProt identifiers.

254 *Creating Knockdown Mutants of the Candidate EDC Genes with CRISPRi*

255 CRISPRi mutants of the EDC genes were created as previously described <sup>16</sup>.  
256 Briefly, pSCB2-sgRNA<sub>v2</sub>, a modified plasmid from pSCB2-sgRNA <sup>16</sup>, was used as the  
257 template for inverse PCR to insert 20bp target-specific sgRNA sequence. Inverse PCR was  
258 performed using Q5 high-fidelity polymerase (NEB), forward primers with individual sgRNAs  
259 as 5' tail, and 1092 as the reverse primer. The resultant fragments were ligated to create circular  
260 plasmids by incubating 0.5µL of the respective PCR products with quick ligation buffer (NEB),  
261 0.25 µL *DpnI*, 0.25 µL T4 polynucleotide kinase (NEB), and 0.25 µL T4 ligase (NEB) for 30  
262 minutes at 37°C. Resultant plasmids were transformed into *E. coli* DH5α, recovered for 2h and  
263 selected in LB supplemented with trimethoprim 50µg/mL (Sigma). The transformants were  
264 further confirmed by colony PCR using primers 1409 and 848. *E. coli* strains carrying the  
265 sgRNA plasmids were used as donors, and *E. coli* MM290/pRK2013 as the helper for  
266 triparental mating to introduce the sgRNA plasmids into *B. cenocepacia* K56-2 containing the  
267 chromosomally integrated dCas9 under the control of a rhamnose inducible promoter, as  
268 described previously <sup>38</sup>. Trimethoprim resistant colonies (100µg/mL) were selected and  
269 screened by colony PCR using the primers 1409 and 848. The list of all the primers used in this  
270 study is provided in **Supplementary Table 4**.

271 *Conditional Growth Phenotype Analysis of the CRISPRi Mutants*

272 To determine the conditional growth phenotype of the candidate genes, overnight  
273 cultures of the CRISPRi mutants were back diluted to OD<sub>600nm</sub> 0.01. The cultures were grown  
274 at 37°C for 20-24 hours with continuous shaking in a 384-well plate containing LB broth  
275 supplemented with trimethoprim 100µg/mL and with/without 1% rhamnose. OD<sub>600nm</sub> readings  
276 were taken at 1 hour intervals using BioTek Synergy 2 microplate reader.

277 **Acknowledgments**

278 This work was supported by grants from the Canadian Institutes of Health Research  
279 (CIHR), Cystic Fibrosis Foundation, Cystic Fibrosis Canada to STC; ASMZR was supported  
280 by a University of Manitoba Graduate Fellowship (UMGF). The authors thank Dr. Georg  
281 Hausner, Andrew Hogan, Dustin Maydaniuk and rest of the Cardona lab members for critically  
282 reading the manuscript.

### 283 **Conflict of Interest Statement**

284 The authors declare no conflict of interest.

### 285 **References**

- 286 1. Opijnen, T. van, Bodi, K. L. & Camilli, A. Tn-seq: high-throughput parallel sequencing  
287 for fitness and genetic interaction studies in microorganisms. *Nat. Methods* **6**, 767–772  
288 (2009).
- 289 2. Akerley, B. J. *et al.* Systematic identification of essential genes by in vitro mariner  
290 mutagenesis. *Proc Natl Acad Sci U S A* **95**, 8927–8932 (1998).
- 291 3. Chao, M. C., Abel, S., Davis, B. M. & Waldor, M. K. The design and analysis of  
292 transposon insertion sequencing experiments. *Nat. Rev. Microbiol.* **14**, 119–128 (2016).
- 293 4. Langridge, G. C. *et al.* Simultaneous assay of every *Salmonella Typhi* gene using one  
294 million transposon mutants. *Genome Res.* **19**, 2308–2316 (2009).
- 295 5. Shields, R. C., Zeng, L., Culp, D. J. & Burne, R. A. Genomewide Identification of  
296 Essential Genes and Fitness Determinants of *Streptococcus mutans* UA159. *mSphere* **3**,  
297 e00031-18 (2018).
- 298 6. Forslund, S. K., Kaduk, M. & Sonnhammer, E. L. L. Evolution of Protein Domain  
299 Architectures. in *Evolutionary Genomics* (ed. Anisimova, M.) vol. 1910 469–504  
300 (Springer New York, 2019).

- 301 7. Schnoes, A. M., Brown, S. D., Dodevski, I. & Babbitt, P. C. Annotation Error in Public  
302 Databases: Misannotation of Molecular Function in Enzyme Superfamilies. *PLoS*  
303 *Comput. Biol.* **5**, e1000605 (2009).
- 304 8. El-Gebali, S. *et al.* The Pfam protein families database in 2019. *Nucleic Acids Res.* **47**,  
305 D427–D432 (2019).
- 306 9. Mitchell, A. L. *et al.* InterPro in 2019: improving coverage, classification and access to  
307 protein sequence annotations. *Nucleic Acids Res.* **47**, D351–D360 (2019).
- 308 10. Bateman, A., Coggill, P. & Finn, R. D. DUFs: families in search of function. *Acta*  
309 *Crystallograph. Sect. F Struct. Biol. Cryst. Commun.* **66**, 1148–1152 (2010).
- 310 11. Goodacre, N. F., Gerloff, D. L. & Uetz, P. Protein Domains of Unknown Function Are  
311 Essential in Bacteria. *mBio* **5**, e00744-13 (2014).
- 312 12. Lluch-Senar, M. *et al.* Defining a minimal cell: essentiality of small ORFs and ncRNAs  
313 in a genome-reduced bacterium. *Mol. Syst. Biol.* **11**, 780 (2015).
- 314 13. Lu, Y. *et al.* A novel essential domain perspective for exploring gene essentiality.  
315 *Bioinformatics* **31**, 2921–2929 (2015).
- 316 14. Gislason, A. S., Turner, K., Domaratzki, M. & Cardona, S. T. Comparative analysis of  
317 the *Burkholderia cenocepacia* K56-2 essential genome reveals cell envelope functions  
318 that are uniquely required for survival in species of the genus *Burkholderia*. *Microb.*  
319 *Genomics* **3**, e000140 (2017).
- 320 15. Qi, L. S. *et al.* Repurposing CRISPR as an RNA-guided platform for sequence-specific  
321 control of gene expression. *Cell* **152**, 1173–1183 (2013).
- 322 16. Hogan, A. M., Rahman, A. S. M. Z., Lightly, T. J. & Cardona, S. T. A Broad-Host-Range  
323 CRISPRi Toolkit for Silencing Gene Expression in *Burkholderia*. *ACS Synth. Biol.* **8**,  
324 2372–2384 (2019).

- 325 17. Luo, H. *et al.* DEG 15, an update of the Database of Essential Genes that includes built-in  
326 analysis tools. *Nucleic Acids Res.* **49**, D677–D686 (2021).
- 327 18. The UniProt Consortium. UniProt: a worldwide hub of protein knowledge. *Nucleic Acids*  
328 *Res.* **47**, D506–D515 (2019).
- 329 19. Gao, M., Glenn, A. E., Blacutt, A. A. & Gold, S. E. Fungal Lactamases: Their  
330 Occurrence and Function. *Front. Microbiol.* **8**, 1775 (2017).
- 331 20. Kanaan, S. P., Huang, C., Wuchty, S., Chen, D. Z. & Izaguirre, J. A. Inferring Protein–  
332 Protein Interactions from Multiple Protein Domain Combinations. in *Computational*  
333 *Systems Biology* (eds. Ireton, R., Montgomery, K., Bumgarner, R., Samudrala, R. &  
334 McDermott, J.) vol. 541 43–59 (Humana Press, 2009).
- 335 21. Miravet-Verde, S., Burgos, R., Delgado, J., Lluch-Senar, M. & Serrano, L. FASTQINS  
336 and ANUBIS: two bioinformatic tools to explore facts and artifacts in transposon  
337 sequencing and essentiality studies. *Nucleic Acids Res.* **48**, e102 (2020).
- 338 22. Goodall, E. C. A. *et al.* The Essential Genome of *Escherichia coli* K-12. *mBio* **9**, e02096-  
339 17 (2018).
- 340 23. Qayyum, M. Z., Dey, D. & Sen, R. Transcription Elongation Factor NusA Is a General  
341 Antagonist of Rho-dependent Termination in *Escherichia coli*. *J. Biol. Chem.* **291**, 8090–  
342 8108 (2016).
- 343 24. Ha, K. S., Touloukhonov, I., Vassylyev, D. G. & Landick, R. The NusA N-Terminal  
344 Domain Is Necessary and Sufficient for Enhancement of Transcriptional Pausing via  
345 Interaction with the RNA Exit Channel of RNA Polymerase. *J. Mol. Biol.* **401**, 708–725  
346 (2010).
- 347 25. Wu, C.-C., Naveen, V., Chien, C.-H., Chang, Y.-W. & Hsiao, C.-D. Crystal Structure of  
348 DnaK Protein Complexed with Nucleotide Exchange Factor GrpE in DnaK Chaperone

- 349 System: INSIGHT INTO INTERMOLECULAR COMMUNICATION. *J. Biol. Chem.*  
350 **287**, 21461–21470 (2012).
- 351 26. Barthel, T. K., Zhang, J. & Walker, G. C. ATPase-Defective Derivatives of *Escherichia*  
352 *coli*DnaK That Behave Differently with Respect to ATP-Induced Conformational Change  
353 and Peptide Release. *J. Bacteriol.* **183**, 5482–5490 (2001).
- 354 27. Vogel, M., Mayer, M. P. & Bukau, B. Allosteric Regulation of Hsp70 Chaperones  
355 Involves a Conserved Interdomain Linker. *J. Biol. Chem.* **281**, 38705–38711 (2006).
- 356 28. Bikard, D. *et al.* Programmable repression and activation of bacterial gene expression  
357 using an engineered CRISPR-Cas system. *Nucleic Acids Res.* **41**, 7429–7437 (2013).
- 358 29. Hawkins, J. S., Wong, S., Peters, J. M., Almeida, R. & Qi, L. S. Targeted Transcriptional  
359 Repression in Bacteria Using CRISPR Interference (CRISPRi). *Methods Mol. Biol.*  
360 *Clifton NJ* **1311**, 349–362 (2015).
- 361 30. Vigouroux, A., Oldewurtel, E., Cui, L., Bikard, D. & Teeffelen, S. van. Tuning dCas9's  
362 ability to block transcription enables robust, noiseless knockdown of bacterial genes.  
363 *Mol. Syst. Biol.* **14**, e7899 (2018).
- 364 31. Bastard, K. *et al.* Revealing the hidden functional diversity of an enzyme family. *Nat.*  
365 *Chem. Biol.* **10**, 42–49 (2014).
- 366 32. Dessailly, B. H. *et al.* PSI-2: Structural Genomics to Cover Protein Domain Family  
367 Space. *Structure* **17**, 869–881 (2009).
- 368 33. Zhang, X. *et al.* Assignment of function to a domain of unknown function: DUF1537 is a  
369 new kinase family in catabolic pathways for acid sugars. *Proc. Natl. Acad. Sci.* **113**,  
370 E4161–E4169 (2016).
- 371 34. Sievers, F. *et al.* Fast, scalable generation of high-quality protein multiple sequence  
372 alignments using Clustal Omega. *Mol. Syst. Biol.* **7**, 539 (2011).

- 373 35. Kumar, S., Stecher, G., Li, M., Knyaz, C. & Tamura, K. MEGA X: Molecular  
374 Evolutionary Genetics Analysis across Computing Platforms. *Mol. Biol. Evol.* **35**, 1547–  
375 1549 (2018).
- 376 36. Jones, D. T., Taylor, W. R. & Thornton, J. M. The rapid generation of mutation data  
377 matrices from protein sequences. *Bioinformatics* **8**, 275–282 (1992).
- 378 37. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v4: recent updates and new  
379 developments. *Nucleic Acids Res.* **47**, W256–W259 (2019).
- 380 38. Hogan, A. M. *et al.* Competitive Fitness of Essential Gene Knockdowns Reveals a  
381 Broad-Spectrum Antibacterial Inhibitor of the Cell Division Protein FtsZ. *Antimicrob.*  
382 *Agents Chemother.* **62**, e01231-18 (2018).
- 383

# Figures

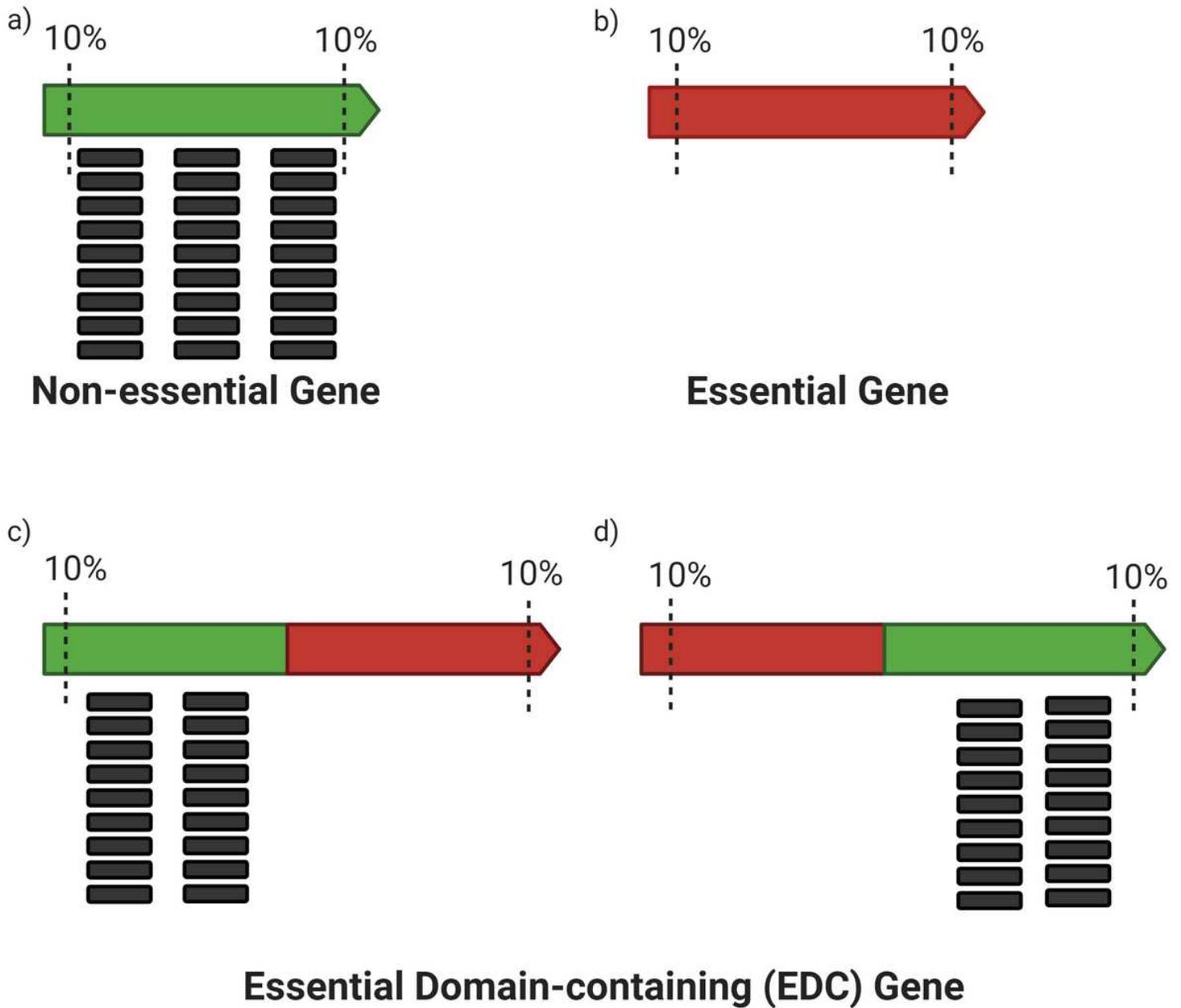
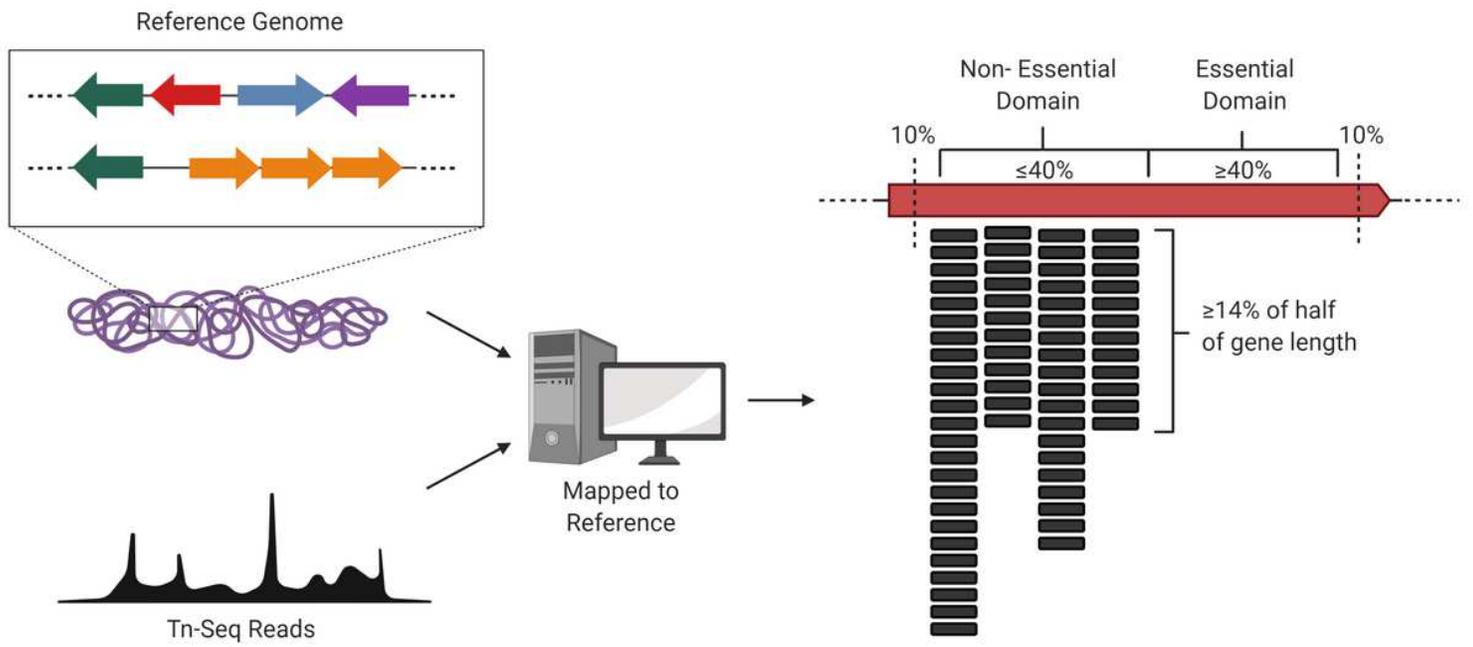


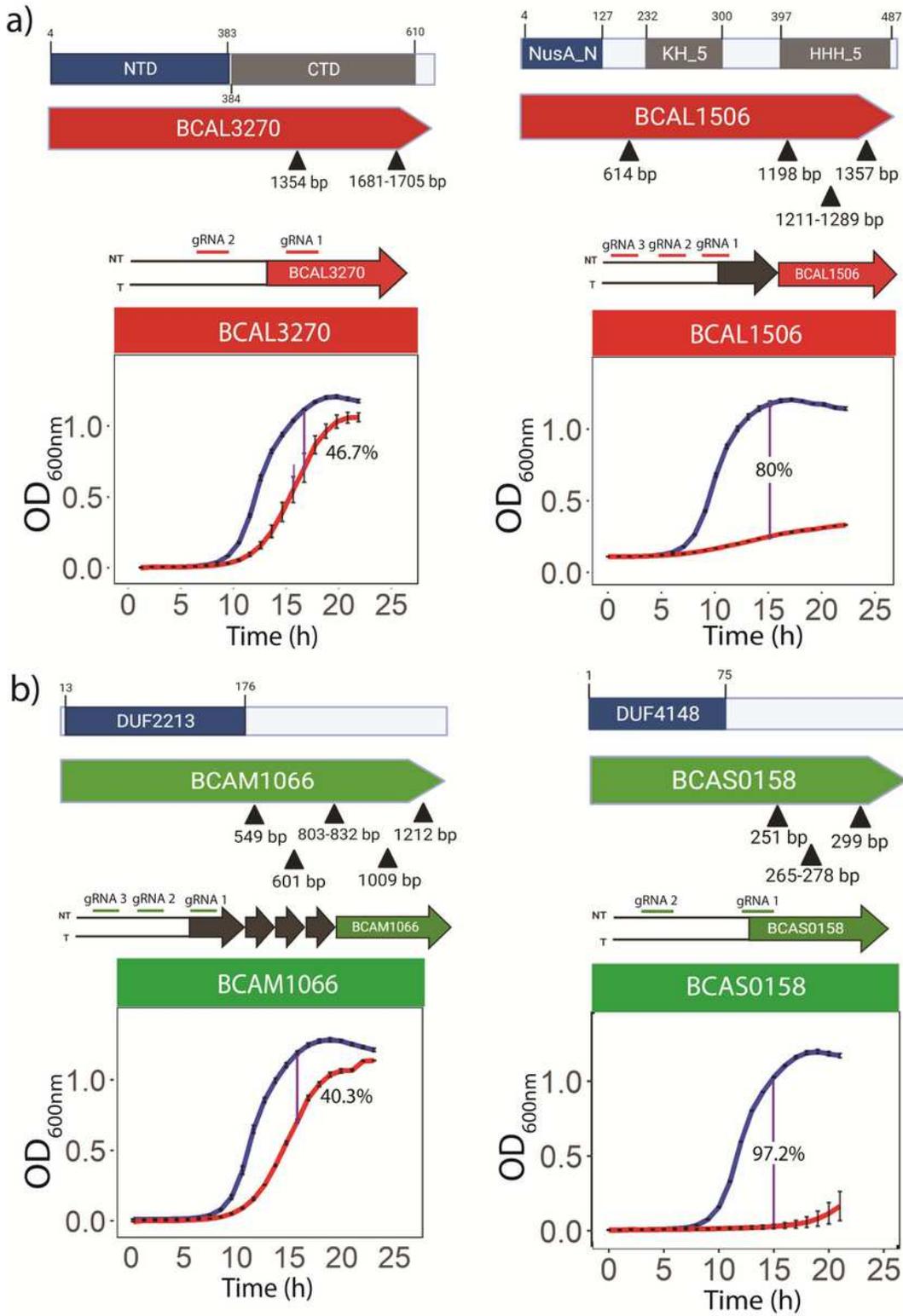
Figure 1

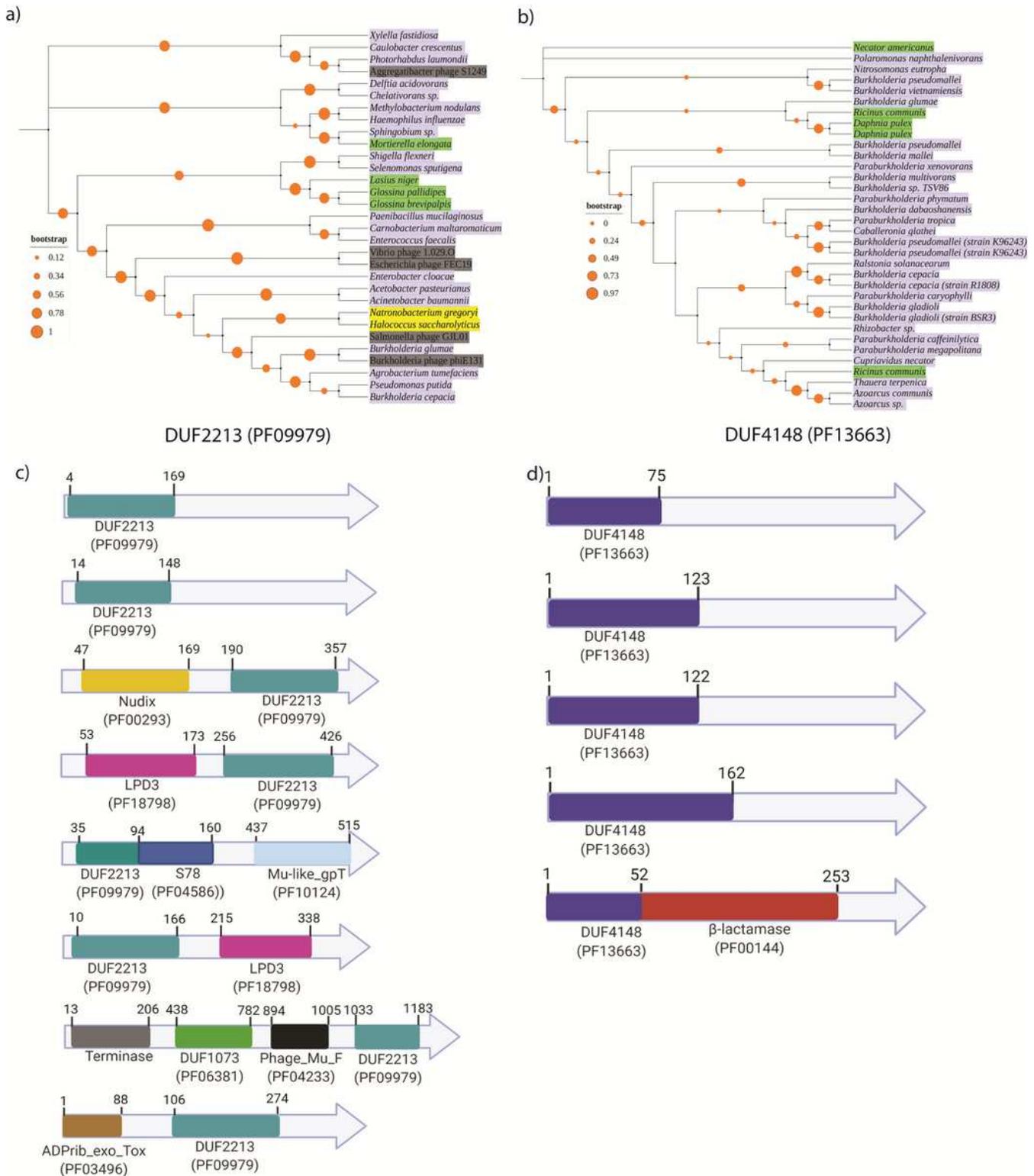
While disrupted genes are regarded as “non-essential,” the method yields a list of putative essential genes as those with zero or very few mapped reads (1a and 1b) We operationally defined this subclass of essential genes as “essential domain-containing” (EDC) genes (1c and 1d)



**Figure 2**

Asymmetric distribution of transposon insertions





**Figure 4**

DUF2213 is also present in many phage-related proteins (4a). DUF4148 is found in 204 bacterial species, primarily in Burkholderia species(4b) Eight unique domain architectures were observed for proteins containing DUF2213 and five for DUF4148 (4c-d).

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementalFiles.pdf](#)
- [SupplementaryTable1.csv](#)
- [SupplementaryTable2.xlsx](#)
- [SupplementaryTable5.xlsx](#)