T1TAdb: the database of Type I Toxin-Antitoxin systems

- 3 Nicolas J. Tourasse* and Fabien Darfeuille*
- 4

2

- 5 Univ. Bordeaux, CNRS, INSERM, ARNA, UMR 5320, U1212, F-33000 Bordeaux, France.
- 6 *corresponding authors. E-mail: <u>nicolas.tourasse@inserm.fr;</u> <u>fabien.darfeuille@inserm.fr</u>

7

- 8 Running head: Database of type I toxin-antitoxin systems
- 9
- 10 Keywords: toxin-antitoxin system, database, genome, antisense RNA, RNA structure, bioinformatics

11

1 Abstract

2 Type I toxin-antitoxin (T1TA) systems constitute a large class of genetic modules with 3 antisense RNA (asRNA)-mediated regulation of gene expression. They are widespread in bacteria 4 and consist of an mRNA coding for a toxic protein and a noncoding asRNA that acts as an antitoxin 5 preventing the synthesis of the toxin by directly basepairing to its cognate mRNA. The co- and post-6 transcriptional regulation of T1TA systems is intimately linked to RNA sequence and structure, 7 therefore it is essential to have an accurate annotation of the mRNA and asRNA molecules to 8 understand this regulation. However, most T1TA systems have been identified by means of 9 bioinformatic analyses solely based on the toxin protein sequences, and there is no central repository 10 of information on their specific RNA features. Here we present the first database dedicated to type I 11 TA systems, named T1TAdb. It is an open-access web database (https://d-lab.arna.cnrs.fr/t1tadb) 12 with a collection of \sim 1,900 loci in \sim 500 bacterial strains in which a toxin-coding sequence has been 13 previously identified. RNA molecules were annotated with a bioinformatic procedure based on key 14 determinants of the mRNA structure and the genetic organization of the T1TA loci. Besides RNA 15 and protein secondary structure predictions, T1TAdb also identifies promoter, ribosome-binding, and 16 mRNA-asRNA interaction sites. It also includes tools for comparative analysis, such as sequence 17 similarity search and computation of structural multiple alignments, which are annotated with 18 covariation information. To our knowledge, T1TAdb represents the largest collection of features, 19 sequences, and structural annotations on this class of genetic modules.

20

21 Introduction

Toxin-antitoxin (TA) systems are encoded within small genetic loci found in most of bacterial genomes including those of pathogens. They are usually composed of two adjacent genes: a stable toxin and a labile antitoxin that inhibits the toxin's action or expression and whose depletion rapidly leads to cell death or growth arrest (Harms et al. 2018). Six types of TA systems have been described so far depending on the nature and mode of action of the antitoxin (reviewed in (Page and Peti 2016)),

1 and a seventh type has been recently proposed (Wang et al. 2021). While the toxin is always a protein,

2 the antitoxin can be either a protein (types II, IV, V, VI, and VII) or an RNA (types I and III).

3 Type II TA systems are by far the most well studied class of TA systems. Extensive data for 4 these systems are available in two databases, TADB 2.0 (Xie et al. 2018) and TASmania (Akarsu et 5 al. 2019), which also include limited data for other types of TA systems. These databases also provide 6 a tool (TAfinder and TASer, respectively) to scan and predict TA systems (Akarsu et al. 2019; Xie et 7 al. 2018). Other TA-specific resources include BtoxDB (Barbosa et al. 2015), a database of TA 8 protein structural data, and RASTA_Bacteria (Sevin and Barloy-Hubler 2007), a tool to scan for 9 toxins and antitoxins in bacterial genomes (unfortunately, both are no longer maintained). Overall, 10 none of these existing databases contains expanded data about type I TA (T1TA) systems. The RNA 11 families database Rfam (Kalvari et al. 2021) currently contains more than 600 antitoxin RNA entries 12 that belong to 13 known T1TA families. Rfam is a useful resource that provides, for all families, 13 sequences, alignments, structures, covariance models (that can be used to search genomes), 14 phylogenetic trees, and links to Wikipedia articles. However, it contains no data for the associated 15 toxin mRNAs of the T1TA systems, as Rfam is focused on non-coding RNAs.

16 A T1TA system consists of a relatively short mRNA (150 to 400 nucleotides long) coding for 17 a small protein (20-60 amino acids in length) whose expression is toxic to the host cell and an 18 antisense RNA (asRNA; 60-200 nt in length) that serves as a counteracting antitoxin to prevent the 19 synthesis of its cognate toxin by directly basepairing to the mRNA. Numerous aspects of T1TA 20 systems have been studied including RNA structure, toxin-antitoxin interaction, regulatory gene 21 expression mechanisms (transcription, translation, degradation, processing), mechanism of action of 22 the toxin, and function in cell physiology (Masachis and Darfeuille, 2018). While TA systems located 23 on plasmids have been demonstrated to contribute to plasmid maintenance (via the mechanism of 24 postsegregational killing, (Greenfield et al. 2000; Gerdes et al. 1986)), the roles of TA systems 25 encoded on the chromosome has only been addressed for a few TA loci (Brielle et al. 2016; Peltier et 26 al. 2020).

1 A few T1TA systems have been experimentally characterized, but hundreds have been 2 identified by bioinformatic analyses (Arnion et al. 2017; Fozo et al. 2010). By exhaustive amino acid 3 sequence homology searches using PSI-BLAST and TBLASTN run with customized parameters, 4 Fozo and coworkers identified ~900 sequences of ORFs coding for homologs of known type I toxin 5 peptides from 7 families in 95 bacterial species and 229 strains (Suppl. Table S5 in (Fozo et al. 2010)). 6 In addition, through searches based on characteristics of T1TA loci (such as tandem repeats and 7 hydrophobicity) they proposed more than 2,000 novel toxin ORFs in hundreds of genomes (Suppl. 8 Table S8 therein). These analyses were essentially protein-based and did not investigate the mRNA 9 features. Based on thermodynamic local free energies, they could detect the position of asRNA genes 10 of a few known families, but the coordinates spanned only a 100-nt window and the precise start and 11 end boundaries of the full-length molecules were not obtained (Suppl. Table S9 therein). Most of the 12 identified T1TA loci are not yet annotated in genome records, and it is necessary to have an accurate 13 annotation of the complete mRNA and asRNA molecules to understand the co- and post-14 transcriptional regulation of T1TA systems, which is determined by RNA sequence and structure 15 (Masachis and Darfeuille 2018). Thus, there is a deep need for a central repository of T1TA systems 16 that would include RNA information. In this work, we have built a database, named T1TAdb, which 17 gathers all described and predicted loci of T1TA systems and where the mRNA and asRNA 18 coordinates are annotated. We describe below the content and main features of T1TAdb, along with 19 the procedure used to identify mRNAs and asRNAs.

20

21 Results and Discussion

We have developed T1TAdb, the first database dedicated to T1TA systems. It can be accessed by users through a graphical web interface. A preliminary, development version was put on-line in January 2019 and has been regularly updated, corrected, and improved (<u>https://d-</u> <u>lab.arna.cnrs.fr/t1tadb</u>). The database provides sequence, secondary structure, and genomic information on T1TA loci. In its current, initial version, it is limited to bacterial genomes that were

1 reported in the literature to carry T1TA systems. It contains 1894 loci belonging to 24 families from 2 218 bacterial species and 493 strains. Data on toxin mRNA, antitoxin asRNA, and toxin peptides 3 were taken from the current literature. Only a small number of loci have been experimentally 4 characterized and the bulk of the data in T1TAdb are mainly based on the results of genome-wide 5 bioinformatic studies. This includes the AapA/IsoA family that has been extensively curated in about 6 100 genomes of *Helicobacter* and *Campylobacter* (Arnion et al. 2017) and the large number of loci 7 found in hundreds of genomes in the study of (Fozo et al. 2010). The analyses by (Fozo et al. 2010) 8 did not investigate the mRNA genes and recovered only partially the asRNA genes. Therefore, in our 9 work, a major effort to build the T1TAdb database has been devoted to the annotation of the genomic 10 locations of the toxin mRNA and antitoxin asRNA corresponding to each toxin ORF of known family 11 reported by (Fozo et al. 2010), in order to reconstruct the complete loci.

12

Annotation of toxin mRNA and antitoxin asRNA genomic localization based on RNA structural
features

15 The genomic regions defining the mRNAs and asRNAs were predicted using RNAMotif 16 (Macke et al. 2001) and RNASurface (Soldatov et al. 2014), respectively. RNAMotif is used to 17 identify regions that can adopt a predefined secondary structure, while RNASurface predicts regions 18 that are structurally more stable than the rest of the genome. A comparison of the asRNA coordinates 19 predicted by RNASurface with those experimentally determined by transcriptome analyses revealed 20 a good agreement. RNASurface predictions were usually within 10 bases of the experimental 21 coordinates, at both the 5' and 3' ends, for diverse T1TA families in Helicobacter pylori and 22 Escherichia coli, but larger differences were more often observed for Enterococcus faecalis (Supp. 23 Tables S2-S6). Moreover, the 5' ends of mRNAs predicted by our RNAMotif structural descriptors 24 showed a similar accuracy. In the case of T1TA systems, a complete TA locus was obtained when a 25 pair of mRNA and asRNA could be predicted, with lengths, orientations, complementary interaction 26 regions, and relative positions matching the genetic organization of a given known TA family (Figure

1 1 and Suppl. Table S1; see Materials and Methods for details). As can be seen in Table 1, using this 2 procedure, we were able to recover RNA regions of the expected family for 83 to 89% (depending 3 on the family) of the toxin ORFs belonging to five of the seven known families detected by (Fozo et 4 al. 2010) (Suppl. Table S5 therein). The success rate was a bit lower for the TxpA/RatA family (75%), 5 and particularly low for the Fst/RNAII family (29%). For most antitoxin asRNAs of the Fst/RNAII 6 family, RNASurface predicts a sequence that folds into a secondary structure different from the 7 typical RNAII structure and IntaRNA does not identify an interaction region with the toxin mRNA 8 at the expected location. Whether this is a prediction issue specific to this family or if this actually 9 represents structural heterogeneity among Fst/RNAII loci requires further study to be determined. In 10 a few additional cases (less than 10% for each family), a locus could be recovered, but was assigned 11 to a different family. For example, locus TA05747 from *Streptococcus thermophilus* classified in the 12 Fst/RNAII family (based on peptide features) was identified by our procedure (based on RNA and 13 genetic organization features) as a member of the TxpA/RatA family. This discrepancy could be due 14 to a wrong family assignment by either bioinformatic method or peculiarities in the locus that make 15 it fortuitously look more similar to another family. Important discrepancies occurred for 17 loci (from 16 known families) where our procedure assigned a locus from a Gram-positive bacterium to a family 17 found in Gram-negative organisms (or vice-versa) exhibiting a completely different organization. 18 These were clearly erroneous hits and half of them (e.g., TA07190 from E. coli and TA07194 from 19 Salmonella enterica) were manually corrected and found to match the family predicted by (Fozo et 20 al. 2010). Interestingly, in many of the remaining cases of discrepancy the families show similarities. 21 For instance, locus TA06047 from E. coli classified in the Ldr/Rdl family (based on peptide features) 22 was identified as a member of the Hok/Sok family (based on RNA and genetic organization features). 23 However, Ldr/Rdl systems exhibit an organization (including a regulatory leader ORF) similar to that 24 of Hok/Sok systems and may be regulated in the same manner (Kawano 2012), thus the RNA-based 25 classification reflects these similarities. Locus TA05436 from Staphylococcus aureus is classified as 26 Fst/RNAII according to peptide features while it was identified as SprA1/SprA1as according to RNA

1 features, but SprA1/SprA1as systems are in fact homologs of Fst/RNAII systems (Weaver et al. 2009; 2 Kwong et al. 2010), and similarly for several loci (e.g., TA05467 from S. aureus) classified as either 3 TxpA/RatA or SprG/SprF (Pinel-Marie et al. 2014). The structural descriptors that we used for 4 RNAMotif searches were mainly based on key determinants of the toxin-encoding mRNA secondary 5 structure, in particular Shine-Dalgarno (SD) and complementary anti-SD sequences, as well as 6 terminator stem-loops in the 3' end, stem-loops in the 5' end, and regions involved in 5'-3' long-7 distance interactions. Because transcription and translation are coupled in bacteria, *cis*-encoded 8 regulatory elements prevent premature translation of the toxin mRNA by sequestering the SD 9 sequence, either co- or post-transcriptionally. Depending on the T1TA system, the anti-SD motif is 10 located either at a short distance upstream of the SD sequence, thus creating a sequestering stem-11 loop, or near the end of the mRNA occluding the SD sequence via a long-distance base-pairing 12 interaction between the 5' and 3' extremities of the mRNA (Masachis and Darfeuille 2018). The 13 overall success of our results demonstrates that taking into account these elements, which are essential 14 for the expression and regulation of T1TA systems, is critical for identifying type I toxin mRNAs.

15

16 *Overview of T1TAdb features and tools*

17 Information in T1TAdb can be searched by entering one or several keywords via the 18 "Keyword Search" tab. A given keyword may be matched against any field in the database or against 19 a specific field such as toxin or antitoxin name, species, strain, taxonomy, locus ID, or genome 20 accession number. Submitting the search form will return a table listing all loci that match the search 21 criteria. The "Select loci" button in the upper left corner of the table allows to further refine the 22 selection of loci by replicon (chromosome or plasmid), TA family or taxonomy, and to select/deselect 23 all loci. The table indicates the locus identifier (TA*nnnn*), TA family, host strain and taxonomy, and 24 reference publication of each locus. For loci that were predicted using structural descriptors of a TA 25 family different from that reported in the literature (i.e., in (Fozo et al. 2010)), the predicted family is 26 indicated in the "putative family" column, and groups of loci from the same genomic region matching

1 multiple family characteristics appear under a specific identifier (TAGnnnn; "Overlapping Locus 2 Group", see Materials and Methods for details). The leader ORF is also indicated for loci that have 3 been described in the literature to carry a regulatory leader ORF or for which a leader ORF is 4 annotated (for Hok/Sok and Ldr/Rdl families). The loci table provides explanatory notes that appear 5 when hovering over the red asterisk next to specific items. Clicking on a locus identifier leads to a 6 "Locus Details" page that gives detailed information about the locus. The page is divided into four 7 panels (Figure 2). The "Locus" panel provides genomic information with an interactive full genome 8 map (drawn with CGView; (Stothard and Wishart 2005)) of the host strain, showing the location of 9 all TA loci harbored by the strain, and an embedded interactive genome browser (IGV; (Robinson et 10 al. 2011); https://github.com/igvteam/igv.js/wiki) for viewing the genomic context around the locus. 11 The "mRNA" and "sRNA" panels provide the sequence and secondary structure diagram of the toxin 12 mRNA and antitoxin asRNA, respectively. The location of relevant motifs (promoter -10 box, 13 start/stop codon, SD sequence, leader ORF, mRNA-asRNA interaction region) are annotated on the 14 sequences and /or diagrams. The "Peptide" panel shows the sequence, secondary structure model and 15 hydrophobicity plot of the toxin peptide. The "mRNA", "sRNA", and "Peptide" panels include a 16 "BLAST Sequence" link for launching a sequence homology search (using BLAST+; (Camacho et 17 al. 2009)) of the selected query against every mRNA, asRNA, ORF, peptide, or genome sequence in 18 T1TAdb. The database also allows the user to input or upload one or several query sequences to 19 perform a BLAST search via the "Sequence Search" tab.

The T1TAdb home page contains links to a full list of organisms ("view organism list"; "T1TAdb Organisms" page) and TA families ("view family list"; "T1TAdb Families" page) included in the database. Selecting a particular organism name or family name will return a table listing all loci that belong to that organism or family. The family table on the "T1TAdb Families" page provides links to the corresponding Rfam entries for antitoxin RNA families that are included in Rfam. In addition, the table is organized to indicate groups of T1TA families that share some degree of homology (in sequence, structure, and/or genetic organization).

1 A notable implementation of T1TAdb is that, in addition to homology searches by BLAST, 2 users can perform multiple sequence alignments (using MAFFT; (Katoh and Standley 2013)) to 3 compare mRNA, asRNA, or peptide sequences. Alignments can be launched via the "Align selected 4 sequences" button on top of the table that is returned after selecting loci by organism or family, or 5 following a keyword search. An embedded interactive alignment viewer (MSAViewer; (Yachdav et 6 al. 2016)) is provided to browse and download the alignments. For RNA sequences, structural 7 information is taken into account in MAFFT to produce a structural alignment (Katoh and Toh 2008), 8 and a consensus 2D structure is predicted (using RNAalifold; (Bernhart et al. 2008)) from the 9 alignment. Furthermore, a covariation analysis is performed (using R-chie; (Lai et al. 2012)) to 10 annotate the alignment with covariation information ((Tourasse and Darfeuille 2020); dedicated links 11 are provided to download the consensus and covariation data). With these programs (and BLAST) 12 T1TAdb aims to provide a selected set of tools (run with most default parameters) for comparative 13 sequence and structure analysis, although users may wish to use custom settings or some of the 14 numerous other tools that exist for doing this kind of analysis.

In addition to alignment results, the other data in T1TAdb are also available for download and further custom analysis. Once a set of loci has been selected, the corresponding data can be downloaded using the "Download selected loci" menu at the top of the loci table. Users can obtain a spreadsheet file containing the detailed information about the loci (including host organism, genomic coordinates and sequences of the RNA, ORF, peptide, and promoter features) or sequence files in FASTA format. In the "Locus Details" page for a given locus the structure diagrams can be saved in various formats (SVG, PDF, PNG, or GIF).

The list of publications from which sequences, structures, and other information about the T1TA systems have been obtained is given in the "References" page. Links to other resources on TA systems, as well as the various software and tools used to build T1TAdb, are provided in the "Links" page.

26

9

1 Example of insights gained from the use of T1TAdb

2 We present below an example illustrating the usefulness and added value provided by 3 T1TAdb. In particular, we would like to highlight the type of sequence and structure analysis that can 4 be done with the database. Examination of the structural diagrams of individual Sok antitoxin 5 asRNAs from the Hok/Sok family suggested that there may be different structural types. This 6 hypothesis could be tested in T1TAdb by performing structural alignments of different subsets of Sok 7 RNAs which revealed two subgroups: one subgroup had a 2D structure conforming to that reported 8 by (Franch et al. 1997) (Figure 3A), while the other carries an additional stem-loop element at the 5' 9 end (Figure 3B). The covariation analysis that is run on the alignments in T1TAdb confirmed that 10 this extra element is well supported by covarying positions that reveal compensatory substitutions in 11 a number of sequences. Furthermore, the annotations of the Sok sequences provided in T1TAdb 12 indicated the presence of a -10 box promoter motif (TANNNT, where N means any nucleotide) at an 13 appropriate location upstream of the stem-loop. Although these observations should be verified 14 experimentally, this example shows the potential of T1TAdb to reveal new features in RNA-mediated 15 regulation of T1TA systems by comparative analysis that could not be identified during the study of 16 single T1TA loci.

17

18 Conclusions and perspectives

19 In this work, we have built the first web database totally dedicated to T1TA systems, named 20 T1TAdb. T1TA systems are widespread in bacteria, including pathogenic species, and are studied for 21 numerous aspects including RNA and protein structure and function, regulation of gene expression 22 and cell physiology. In addition to giving access to the collection of loci that have been reported in 23 the literature, the database brings an added value by the annotation of toxin mRNAs and antitoxin 24 asRNAs of loci that were predicted solely by analysis of the toxin peptide sequence, and by providing 25 tools to compare their sequences and structures. This information, together with a wealth of sequence, 26 structure, and genomic data provided on T1TA systems, along with the other databases and tools

dedicated to TA systems (such as TADB and TASmania), will certainly serve the scientific community to gain deeper insights of the distribution, evolution, structure, and function of TA systems. This could reveal to be particularly important for the TA systems located in bacterial chromosomes, whose function remains largely unknown. We also believe that the knowledge contained in T1TAdb, and more specifically the RNA structures and loci organization, is a strong asset to help design new tools for identifying new T1TA systems.

7 In future developments of T1TAdb, we plan to expand the collection of T1TA systems by 8 predicting TA loci in all bacterial genomes available. In addition, the procedure described in this 9 work for the annotation of mRNAs and asRNAs of T1TA systems, which is based on structural 10 features, can be turned into an automatic *de novo* prediction tool, most likely with additional steps to 11 control the false-positive rate. The RNA data in T1TAdb may also be used to generate alignments 12 and covariance models in order to search genome databases for conserved RNAs, as exemplified in 13 the Rfam database (Kalvari et al. 2021). We also anticipate to incorporate transcriptome (RNA-Seq) 14 data into T1TAdb. As the number of sequenced bacterial transcriptomes is increasing, such 15 information would be valuable to check whether TA loci are expressed and to verify the genomic 16 coordinates of the predicted RNAs. Further improvements to T1TAdb will also include tools for 17 simultaneously visualizing and comparing the genomic context (gene neighborhoods) across multiple 18 strains or for the reconstruction of phylogenetic trees, and will allow users to submit data for inclusion 19 in T1TAdb.

20

21 Availability

22

T1TAdb is freely available on-line at <u>https://d-lab.arna.cnrs.fr/t1tadb</u>.

23

24 Materials and Methods

25 Data collection

1 Sequence, secondary structure, and genomic localization of T1TA loci (including mRNA, 2 asRNA, ORF, and promoter -10 boxes) from the various families were taken from the publications 3 where they were initially discovered and/or characterized (Arnion et al. 2017; Durand et al. 2012; 4 Fozo et al. 2010; Weaver et al. 2009; Maikova et al. 2018; Folli et al. 2017; Kristiansen et al. 2016; 5 Wen and Fozo 2014; Pinel-Marie et al. 2014; Wen et al. 2014; Jahn and Brantl 2013; Weaver 2012; 6 Fozo 2012; Sayed et al. 2012; Han et al. 2010; Sharma et al. 2010; Darfeuille et al. 2007; Pichon and 7 Felden 2005; Kawano et al. 2002; Pedersen and Gerdes 1999; Franch et al. 1997; Masachis and 8 Darfeuille 2018; Meißner et al. 2016; Peltier et al. 2020; Germain-Amiot et al. 2019; Guo et al. 2014; 9 Kawano et al. 2007; Andresen et al. 2020). In most cases, reports were on one or few copies of a 10 given TA system, identified or studied in one or a few bacterial strains. Information on TA systems 11 belonging to the AapA/IsoA family identified by bioinformatic analyses and manual curation in ~100 12 genomes of *Helicobacter* and *Campylobacter* were taken from (Arnion et al. 2017). (Sharma et al. 13 2010) reported members of the AapB/IsoB, AapC/IsoC, and AapD/IsoD families in six *Helicobacter* 14 strains. For the sake of consistency and completeness, we located the homologues of these loci in all 15 other Helicobacter genomes screened by (Arnion et al. 2017) through sequence homology searches 16 (using BLAST+ 2.2.31; (Camacho et al. 2009)) and multiple sequence alignment. (Kristiansen et al. 17 2016) identified mRNAs and ORFs belonging to the DinQ/AgrB family in 15 Gram-negative bacteria 18 and we devised a procedure to predict the corresponding asRNAs (see below). The majority of the 19 data in T1TAdb are based on the genome-wide bioinformatic searches by (Fozo et al. 2010) who 20 identified sequences of ORFs coding for type I toxin peptides of known and novel families in 21 hundreds of bacterial strains. For toxins belonging to known TA families, we performed 22 computational analyses to locate the coordinates of the toxin mRNA and antitoxin asRNA in order to 23 identify the complete TA locus corresponding to each reported toxin ORF.

24

25 Prediction of toxin mRNA and antitoxin asRNA genomic localization

1 Identification of mRNAs and asRNAs was based on the characteristics of known examples. 2 The prediction pipeline is summarized in Fig. 1. Genomic locations of toxin mRNAs were determined 3 using RNAMotif 3.1.1 (Macke et al. 2001). RNAMotif takes as input a descriptor file that contains 4 parameters describing the various elements that make up the secondary structure of a particular RNA 5 (helices, loops, etc. and their respective lengths). Using this descriptor RNAMotif scans genome 6 sequences to find regions that can adopt (i.e., can be folded into) the specified structure. RNAMotif 7 descriptors for mRNAs belonging to 16 known T1TA families were written based on sequences and 8 secondary structures reported in the literature cited above. We did not include every structural 9 element in the descriptors, but rather focused on key determinants such as SD and anti-SD sequences 10 (that sequester the SD motif), location of the ORF within the mRNA, stem-loops in the 5' end or 11 terminator stem-loops in the 3' end, and regions implicated in 5'-3' long-distance interactions (textual 12 summary of RNAMotif parameters given in Suppl. Table S1). The remaining parts of the structure 13 were specified as undefined regions. To avoid descriptors being too specific to a given mRNA 14 instance, lengths and distances between the various elements were usually not set to a defined value 15 but to a min./max. value or a relatively broad range of values. All genomes surveyed by (Fozo et al. 16 2010) (genome sequences downloaded from NCBI RefSeq; (Haft et al. 2018)) were scanned by 17 RNAMotif with descriptors of all 16 TA families. Hits that spanned the coordinates of the ORFs 18 reported by Fozo *et al.* on the same DNA strand were extracted (by means of the "intersectbed" utility 19 from the BEDTools 2.24.0 package (Quinlan and Hall 2010)). Among those, hits that contained ORFs 20 that were in-frame with the ORFs of Fozo et al. (ORFs may not always be identical and could differ 21 slightly in length in some cases depending on the start codon used) and whose ORF and total mRNA 22 length matched the typical range of ORF and mRNA length for a given family were retrieved, whether 23 or not the family matched that reported by Fozo and coworkers.

Antitoxin asRNAs were localized using RNASurface 1.1 (Soldatov et al. 2014), which predicts sequence regions that are structurally more stable than the rest of the genome by local folding of segments up to a predefined size. Genome sequences were scanned with RNASurface (run with

options "--winmin 50 --winmax 300 -z -2 -d 500") to identify structured segments of length 50-300
nt that have a z-score lower than -2 (corresponding to a false-positive rate of 5%, (Soldatov et al.
2014)). For a given TA family, RNASurface segments whose length corresponded to the typical
asRNA length for that family (+/- 20%) were retrieved.

Putative promoter -10 boxes were predicted by searching (by means of the PERL module
Regexp::Exhaustive) for the presence of the sequence motif TANNNT (where N means any
nucleotide) in a window covering the region -20 to -6 upstream of the identified mRNAs and asRNAs.
In case where multiple motifs were present the 3'-most (i.e., the closest to the RNA start) was
selected.

10 In order to find the pair of mRNA and asRNA corresponding to the same TA locus, RNAMotif 11 and RNASurface results were intersected (using "intersectbed" from BEDTools) based on the defined 12 genetic organization of the various T1TA families. In the ShoB/OhsC, TisB/IstR1, Zor/Orz, and 13 DinQ/AgrB families, found in Gram-negative bacteria, the asRNA is located 5' to the mRNA and 14 does not overlap it, whereas in the Fst/RNAII, TxpA/RatA, YonT/as-YonT, SprA1/SprA1as, 15 SprG/SprF, and BsrE-G-H/as-BsrE-G-H families, all found in Gram-positive bacteria, the asRNA is 16 located on the 3' side of the mRNA and partially overlaps it (overlaps also the ORF except in the 17 Fst/RNAII family); the asRNA is fully overlapped by the mRNA in the AapA/IsoA, Ibs/Sib, 18 Hok/Sok, and Ldr/Rdl families from Gram-negative organisms, but overlaps the ORF only in the 19 AapA/IsoA and Ibs/Sib families (spans the entire ORF in Ibs/Sib), whereas it is located 5' of the ORF 20 in the Hok/Sok and Ldr/Rdl families ((Wen and Fozo 2014; Arnion et al. 2017; Pinel-Marie et al. 21 2014; Durand et al. 2012; Sayed et al. 2012); Suppl. Table S1). For families where the two RNAs do 22 not overlap, a max. distance of 300 nt between the RNAs was allowed (except for the DinQ/AgrB 23 family, see below), and for families where the two RNAs do overlap, min. and max. limits were set 24 for the length of the overlap region based on known examples (Suppl. Table S1). In T1TA systems 25 the mRNA and asRNA share a complementary region of interaction. In families with overlapping 26 RNAs this region generally corresponds to the overlap region, but for the Fst/RNAII and

1 SprA1/SprA1as families interaction occurs outside the overlap sequence (Weaver et al. 2009; Sayed 2 et al. 2012). The software IntaRNA 3.1.0.2 (Mann et al. 2017) was used to find mRNA-asRNA pairs 3 that share a complementary region at the expected location for these two families and to identify 4 regions of interaction for families with non-overlapping RNAs (the interaction region was set to be 5 at least 15 nt long). A specific processing was carried out for the DinQ/AgrB family because in some 6 bacterial strains the locus includes an AgrA ncRNA gene that is homologous to the AgrB asRNA and 7 that can be located in-between DinQ and AgrB. In Escherichia coli K-12 substr. MG1655 the 8 sequence of the interaction region in AgrA contains a number of mismatches and it has been shown 9 that AgrA does not bind the DinQ mRNA and that only AgrB, which shares a region almost fully 10 complementary to DinQ, acts as an antitoxin (Kristiansen et al. 2016). Therefore, to accommodate 11 for the possible presence of the AgrA and AgrB paralogs, the max. distance between the mRNA and 12 asRNA was extended from 300 to 500 nt. In cases where two asRNAs were predicted in this region 13 by RNASurface and for both of them a possible sequence of interaction (**1** 15 nt) with the mRNA 14 was identified by IntaRNA, we selected the one for which the interaction was predicted to be the most 15 stable (i.e., had the lowest free energy as computed by IntaRNA). The same procedure was followed 16 to predict the asRNAs corresponding to the DinQ-like loci reported in (Kristiansen et al. 2016) where 17 only mRNAs and ORFs were annotated.

18 Following the characteristics described above, all pairs of mRNAs and asRNAs that were in 19 the correct orientation and distance and that shared a complementary interaction region were 20 identified for each family. This set of pairs was then intersected with the set of mRNAs that were 21 found to encode ORFs corresponding to those reported by (Fozo et al. 2010), to obtain the pairs that 22 include a known ORF (Fig. 1). For a given locus, there were usually several mRNAs and asRNAs 23 matching these criteria because we used a relatively broad range of values for the parameters 24 describing the length, structure, and organization of the RNAs. Due to the flexibility in the RNAMotif 25 structural descriptors multiple overlapping mRNAs spanning the same ORF were always found, 26 differing by their lengths and a few bases in their start/end coordinates. Among those, the one for

1 which a promoter -10 box could be predicted was selected as the mRNA for the given locus. If a 2 promoter was predicted for multiple mRNAs, or if no promoter was found for any of the mRNAs, 3 then RNA sequences were folded using MFOLD 3.6 (Zuker 1989, 2003) and the minimum free 4 energy (MFE) of the most stable structure of each RNA was normalized by sequence length (adjusted 5 MFE; AMFE) to compare the stability among RNAs. The mRNA that had the smallest AMFE was 6 taken as the mRNA for the locus. The use of AMFE was warranted by the fact that mRNA lengths of 7 the different T1TA families are in the range 200-400 nt where AMFE is almost length-independent 8 (Trotta 2014). For some loci, there were also several possible asRNAs matching the TA family 9 characteristics as RNASurface often predicts multiple RNA segments of different lengths overlapping 10 the same genomic region. Among those, the one for which a promoter -10 box was predicted was 11 selected as the asRNA for the locus. If a promoter was predicted for multiple asRNAs, or if no 12 promoter was found for any of the asRNAs, then the one with the lowest RNASurface z-score was 13 chosen as the putative asRNA for the given locus. Z-score was used here as normalized folding 14 stability measure because asRNA lengths are in the range 60-200 nt where AMFE is length-dependent 15 and thus cannot be reliably used to compare stabilities among RNAs (Trotta 2014).

To verify that the correct asRNA was associated with a given mRNA, an alignment of the antitoxin sequences was made for each family to check the homology and completeness of the sequences. The alignment included only the RNAs whose family corresponded to the family of the ORF in (Fozo et al. 2010), in which case one should expect these sequences to be homologous to each other. Incomplete sequences or sequences whose coordinates were shifted relative to the others were manually corrected. No alignment was done for the mRNAs because they are constrained by the location of the ORFs and thus are normally homologous.

If no suitable mRNA-asRNA pair spanning an ORF identified by (Fozo et al. 2010) could be found, that particular ORF was not included in T1TAdb. However, in cases where no mRNA-asRNA pair corresponding to the same family as that assigned to the ORF in (Fozo et al. 2010) could be identified but the ORF was included in an RNA pair of a different family, the locus was retained and

1 a "putative_family" flag was set to indicate that it matched an alternative family. This flag was also 2 set when ORFs from unknown or novel families were part of loci corresponding to known families. 3 If RNA pairs of multiple TA families spanned the same ORF, only the one matching the family of 4 the ORF was retained. If there were no such pair, then alternative loci were filtered according to the 5 type of organism in which the ORF was encoded, i.e., for a Gram-negative bacterium only the RNA 6 pairs from families found in Gram-negative bacteria were retained, and similarly for Gram-positive 7 organisms. If there were no such pairs, then all alternative mRNA-asRNA pairs covering the ORF 8 were retained. The multiple loci spanning a given ORF were organized into an "Overlapping Locus 9 Group", which represents a group of loci that have been identified using characteristics of different 10 TA families and that overlap the same genomic region, but we could not determine which one is the 11 real locus in this region.

12

13 Database implementation

14 T1TAdb is implemented as relational database in PostgreSOL 10.13 а 15 (https://www.postgresql.org/). The database schema was designed following the five rules of 16 normalization (http://www.barrywise.com/2008/01/database-normalization-and-design-techniques/) 17 to avoid data redundancy and inconsistent dependency among tables. The graphical web interface 18 was developed using the PERL Catalyst framework (http://www.catalystframework.org/) to control 19 and manage connections and SQL requests to the database. Graphical design of the web pages was 20 done in dynamic and responsive HTML, JavaScript, and Cascading Style Sheets (CSS) 21 (https://www.w3schools.com/), along with the PERL Template Toolkit 2.26 templating system 22 (http://www.template-toolkit.org/). The T1TAdb website is run via the Apache HTTP 2.4.6 server 23 (https://httpd.apache.org/) under the Linux CentOS 7.8 operating system. Secondary structures of 24 RNA were predicted using MFOLD 3.6 (Zuker 2003, 1989) and annotated diagrams highlighting the 25 location of specific motifs (start/stop codon, SD sequence, interaction region)were generated with 26 VARNA 3.93 (Darty et al. 2009). Secondary structures of toxin peptides were predicted using

1 PSIPRED 4.02 (McGuffin et al. 2000) run with PSI-BLAST 2.2.26 against the UniRef90 protein 2 sequence database (https://www.uniprot.org/help/uniref) and drawn with POLYVIEW-2D (Porollo 3 et al. 2004). Hydrophobicity plots were computed with ProtScale (Gasteiger et al. 2005). Interactive 4 genomic maps in SVG format showing the localizations of TA loci were drawn using CGView 5 (Stothard and Wishart 2005). The embeddable JavaScript/CSS version of the IGV browser 6 ((Robinson et al. 2011); https://github.com/igvteam/igv.js/wiki) was used to provide an interactive 7 visualization of the genomic context flanking TA loci, and SVG images of the genomic context were 8 generated by of Gviz the means package 9 (https://bioconductor.org/packages/release/bioc/html/Gviz.html) in R 3.5.3 ((R Development Core 10 Team 2019); https://www.r-project.org/). Sequence similarity searches in T1TAdb are done with 11 BLAST+ 2.2.31 (Camacho et al. 2009) and multiple sequence alignments are computed by MAFFT 12 7.407 (Katoh and Standley 2013). For peptide sequences MAFFT is run with the method "mafft-13 linsi" and the option "--localpair", whereas for RNA sequences MAFFT is run with the method 14 "mafft-xinsi" and the option "--scarnapair" to incorporate structure information and produce a 15 structural alignment (Katoh and Toh 2008). Alignments are visualized using the embeddable 16 MSAViewer (Yachdav et al. 2016), which is part of the BioJS JavaScript tools (Yachdav et al. 2015). 17 In addition, for RNA alignments, a consensus structure is predicted by means of RNA alignments. 18 the ViennaRNA 2.1.9 package (Lorenz et al. 2011; Bernhart et al. 2008) and a covariation analysis is 19 conducted using R-chie (Lai et al. 2012). Other bioinformatic software such as EMBOSS 6.6.0 (Rice 20 et al. 2000), **BioPERL** (Stajich et al. 2002), and FASTX-Toolkit 0.0.14 21 (http://hannonlab.cshl.edu/fastx toolkit/) were also employed to generate and/or process the data 22 included in T1TAdb.

23

24 Acknowledgements

This work was supported by Institut National de la Santé et de la Recherche Médicale
[INSERM U1212], Centre National de la Recherche Scientifique [CNRS UMR5320], Bordeaux

| 1 | University, and Agence | Nationale de la Recherche | [ANR-12-BSV5-0025-Bactox | 1, ANR-12-BSV6- |
|---|------------------------|---------------------------|--------------------------|-----------------|
|---|------------------------|---------------------------|--------------------------|-----------------|

2 0007-asSUPYCO]. The web server hosting T1TAdb is provided by the ODS Web Hosting service of

3 CNRS. We thank Dr. Stéphane Thore and Dr. Sébastien Fribourg for providing additional computing

4 power. Finally, we thank all present and past members of the Darfeuille laboratory for fruitful

5 discussions on this project and Anthony Bugaut, Isabelle Iost, Anaïs Le Rhun, Simon Bonabal, and

- 6 Olga Soutourina for critical comments on the manuscript.
- 7

8 **References**

9 Akarsu H, Bordes P, Mansour M, Bigot DJ, Genevaux P, Falquet L. 2019. TASmania: A bacterial

10 toxin-antitoxin systems database. *PLoS Comput Biol* **15**: e1006946.

- 11 Andresen L, Martínez-Burgo Y, Zangelin JN, Rizvanovic A, Holmqvist E. 2020. The small toxic
- salmonella protein timp targets the cytoplasmic membrane and is repressed by the small rna
 timr. *MBio* 11: e01659-20.
- 14 Arnion H, Korkut DN, Gelo SM, Chabas S, Reignier J, Iost I, Darfeuille F. 2017. Mechanistic
- 15 insights into type I toxin antitoxin systems in Helicobacter pylori: The importance of mRNA
- 16 folding in controlling toxin expression. *Nucleic Acids Res* **45**: 4782–4795.
- 17 Barbosa LCB, Garrido SS, Marchetto R. 2015. BtoxDB: A comprehensive database of protein
- 18 structural data on toxin-antitoxin systems. *Comput Biol Med* **58**: 146–153.

19 Bernhart SH, Hofacker IL, Will S, Gruber AR, Stadler PF. 2008. RNAalifold: Improved consensus

- 20 structure prediction for RNA alignments. *BMC Bioinformatics* **9**: 474.
- Brielle R, Pinel-Marie M-L, Felden B. 2016. Linking bacterial type I toxins with their actions. *Curr Opin Microbiol* 30: 114–121.
- 23 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009.
- 24 BLAST+: Architecture and applications. *BMC Bioinformatics* **10**: 1–9.
- 25 Darfeuille F, Unoson C, Vogel J, Wagner EGH. 2007. An Antisense RNA Inhibits Translation by
- 26 Competing with Standby Ribosomes. *Mol Cell* **26**: 381–392.

- 1 Darty K, Denise A, Ponty Y. 2009. VARNA: Interactive drawing and editing of the RNA secondary
- 2 structure. *Bioinformatics* **25**: 1974–1975.
- 3 Durand S, Jahn N, Condon C, Brantl S. 2012. Type I toxin-antitoxin systems in Bacillus subtilis.
- 4 *RNA Biol* **9**: 1491–1497.
- 5 Folli C, Levante A, Percudani R, Amidani D, Bottazzi S, Ferrari A, Rivetti C, Neviani E, Lazzi C.
- 6 2017. Toward the identification of a type i toxin-antitoxin system in the plasmid DNA of dairy
- 7 Lactobacillus rhamnosus. *Sci Rep* **7**: 1–13.
- 8 Fozo EM. 2012. New type I toxin-antitoxin families from "wild" and laboratory strains of E. coli:
- 9 Ibs-Sib, ShoB-OhsC and Zor-Orz. *RNA Biol* **9**: 1504–1512.
- 10 Fozo EM, Makarova KS, Shabalina SA, Yutin N, Koonin E V., Storz G. 2010. Abundance of type I
- 11 toxin-antitoxin systems in bacteria: Searches for new candidates and discovery of novel
- 12 families. *Nucleic Acids Res* **38**: 3743–3759.
- 13 Franch T, Gultyaev AP, Gerdes K. 1997. Programmed cell death by hok/sok of plasmid R1:
- 14 processing at the hok mRNA 3'-end triggers structural rearrangements that allow translation
- and antisense RNA binding. *J Mol Biol* **273**: 38–51.
- 16 Gasteiger E, Hoogland C, Gattiker A, Duvaud S, Wilkins MR, Appel RD, Bairoch A. 2005. The
- 17 Proteomics Protocols Handbook. 571–608.
- 18 Gerdes K, Rasmussen PB, Molin S. 1986. Unique type of plasmid maintenance function:
- 19 postsegregational killing of plasmid-free cells. *Proc Natl Acad Sci* 83: 3116–3120.
- 20 Germain-Amiot N, Augagneur Y, Camberlein E, Nicolas I, Lecureur V, Rouillon A, Felden B.
- 21 2019. A novel Staphylococcus aureus cis-trans type I toxin-antitoxin module with dual effects
- 22 on bacteria and host cells. *Nucleic Acids Res* **47**: 1759–1773.
- 23 Greenfield TJ, Ehli E, Kirshenmann T, Franch T, Gerdes K, Weaver KE. 2000. The antisense RNA
- of the par locus of pAD1 regulates the expression of a 33-amino-acid toxic peptide by an
- 25 unusual mechanism. *Mol Microbiol* **37**: 652–660.
- 26 Guo Y, Quiroga C, Chen Q, McAnulty MJ, Benedik MJ, Wood TK, Wang X. 2014. RalR (a

- 1 DNase) and RalA (a small RNA) form a type I toxin-antitoxin system in Escherichia coli.
- 2 *Nucleic Acids Res* **42**: 6448–6462.
- 3 Haft DH, DiCuccio M, Badretdin A, Brover V, Chetvernin V, O'Neill K, Li W, Chitsaz F,
- 4 Derbyshire MK, Gonzales NR, et al. 2018. RefSeq: An update on prokaryotic genome
- 5 annotation and curation. *Nucleic Acids Res* **46**: D851–D860.
- 6 Han K, Kim KS, Bak G, Park H, Lee Y. 2010. Recognition and discrimination of target mRNAs by
- 7 Sib RNAs, a cis-encoded sRNA family. *Nucleic Acids Res* **38**: 5851–5866.
- Harms A, Brodersen DE, Mitarai N, Gerdes K. 2018. Toxins, Targets, and Triggers: An Overview
 of Toxin-Antitoxin Biology. *Mol Cell* 70: 768–784.
- 10 Jahn N, Brantl S. 2013. One antitoxin-two functions: SR4 controls toxin mRNA decay and
- 11 translation. *Nucleic Acids Res* **41**: 9870–9880.
- 12 Kalvari I, Nawrocki EP, Ontiveros-Palacios N, Argasinska J, Lamkiewicz K, Marz M, Griffiths-
- 13 Jones S, Toffano-Nioche C, Gautheret D, Weinberg Z, et al. 2021. Rfam 14: expanded
- 14 coverage of metagenomic, viral and microRNA families. *Nucleic Acids Res* **49**: D192–D200.

15 Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7:

- 16 Improvements in performance and usability. *Mol Biol Evol* **30**: 772–780.
- 17 Katoh K, Toh H. 2008. Improved accuracy of multiple ncRNA alignment by incorporating
- 18 structural information into a MAFFT-based framework. *BMC Bioinformatics* **9**: 212.
- 19 Kawano M. 2012. RNA Biology Divergently overlapping cis-encoded antisense RNA regulating
- 20 toxin-antitoxin systems from E. coli hok/sok, ldr/rdl, symE/symR. 9: 1520–1527.
- 21 Kawano M, Aravind L, Storz G. 2007. An antisense RNA controls synthesis of an SOS-induced
- 22 toxin evolved from an antitoxin. *Mol Microbiol* **64**: 738–754.
- 23 Kawano M, Oshima T, Kasai H, Mori H. 2002. Molecular characterization of long direct repeat
- 24 (LDR) sequences expressing a stable mRNA encoding for a 35-amino-acid cell-killing peptide
- and a cis-encoded small antisense RNA in Escherichia coli. *Mol Microbiol* **45**: 333–349.
- 26 Kristiansen KI, Weel-Sneve R, Booth JA, Bjørås M. 2016. Mutually exclusive RNA secondary

| 1 | structures regulate translation initiation of DinQ in Escherichia coli. RNA 22: 1739–1749. |
|----|--|
| 2 | Kwong SM, Jensen SO, Firth N. 2010. Prevalence of Fst-like toxin-antitoxin systems. Microbiology |
| 3 | 156 : 975–977. |
| 4 | Lai D, Proctor JR, Zhu JYA, Meyer IM. 2012. R-CHIE: A web server and R package for |
| 5 | visualizing RNA secondary structures. Nucleic Acids Res 40: e95. |
| 6 | Lorenz R, Bernhart SH, Höner zu Siederdissen C, Tafer H, Flamm C, Stadler PF, Hofacker IL. |
| 7 | 2011. ViennaRNA Package 2.0. Algorithms Mol Biol 6: 26. |
| 8 | Macke TJ, Ecker DJ, Gutell RR, Gautheret D, Case DA, Sampath R. 2001. RNAMotif, an RNA |
| 9 | secondary structure definition and search algorithm. Nucleic Acids Res 29: 4724–4735. |
| 10 | Maikova A, Peltier J, Boudry P, Hajnsdorf E, Kint N, Monot M, Poquet I, Martin-Verstraete I, |
| 11 | Dupuy B, Soutourina O. 2018. Discovery of new type I toxin-antitoxin systems adjacent to |
| 12 | CRISPR arrays in Clostridium difficile. Nucleic Acids Res 46: 4733-4751. |
| 13 | Mann M, Wright PR, Backofen R. 2017. IntaRNA 2.0: Enhanced and customizable prediction of |
| 14 | RNA-RNA interactions. Nucleic Acids Res 45: W435–W439. |
| 15 | Masachis S, Darfeuille F. 2018. Type I Toxin-Antitoxin Systems: Regulating Toxin Expression via |
| 16 | Shine-Dalgarno Sequence Sequestration and Small RNA Binding. Regul with RNA Bact |
| 17 | Archaea 173–190. |
| 18 | McGuffin LJ, Bryson K, Jones DT. 2000. The PSIPRED protein structure prediction server. |
| 19 | Bioinformatics 16: 404–405. |
| 20 | Meißner C, Jahn N, Brantl S. 2016. In Vitro Characterization of the Type I Toxin-Antitoxin System |
| 21 | bsrE/SR5 from Bacillus subtilis. J Biol Chem 291: 560–71. |
| 22 | Page R, Peti W. 2016. Toxin-antitoxin systems in bacterial growth arrest and persistence. Nat Chem |
| 23 | <i>Biol</i> 12 : 208–214. |
| 24 | Pedersen K, Gerdes K. 1999. Multiple hok genes on the chromosome of Escherichia coli. Mol |
| 25 | Microbiol 32 : 1090–1102. |
| | |

26 Peltier J, Hamiot A, Garneau JR, Boudry P, Maikova A, Hajnsdorf E, Fortier LC, Dupuy B,

Downloaded from rnajournal.cshlp.org on April 28, 2024 - Published by Cold Spring Harbor Laboratory Press

Tourasse and Darfeuille

| 1 | Soutourina O. 2020. Type | e I toxin-antitoxin systems | contribute to | mobile genetic elements |
|---|--------------------------|-----------------------------|---------------|-------------------------|
|---|--------------------------|-----------------------------|---------------|-------------------------|

2 maintenance in Clostridioides difficile and can be used as a counter-selectable marker for

3 chromosomal manipulation. *Commun Biol* **3**: 718.

- 4 Pichon C, Felden B. 2005. Small RNA genes expressed from Staphylococcus aureus genomic and
- 5 pathogenicity islands with specific expression among pathogenic strains. *Proc Natl Acad Sci U*
- 6 *S A* **102**: 14249–14254.
- 7 Pinel-Marie M-L, Brielle R, Felden B. 2014. Dual Toxic-Peptide-Coding Staphylococcus aureus
- 8 RNA under Antisense Regulation Targets Host Cells and Bacterial Rivals Unequally. *Cell Rep*9 7: 424–435.
- 10 Porollo AA, Adamczak R, Meller J. 2004. POLYVIEW: A flexible visualization tool for structural

11 and functional annotations of proteins. *Bioinformatics* **20**: 2460–2462.

Quinlan AR, Hall IM. 2010. BEDTools: A flexible suite of utilities for comparing genomic
 features. *Bioinformatics* 26: 841–842.

14 R Development Core Team. 2019. R: A language and environment for statistical computing. R

15 Found Stat Comput. https://www.r-project.org/.

16 Rice P, Longden I, Bleasby A. 2000. EMBOSS: the European Molecular Biology Open Software

- 17 Suite. *Trends Genet* **16**: 276–7.
- 18 Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, Mesirov JP. 2011.
- 19 Integrative genomics viewer. *Nat Biotechnol* **29**: 24–6.

20 Sayed N, Jousselin A, Felden B. 2012. A cis-antisense RNA acts in trans in Staphylococcus aureus

- 21 to control translation of a human cytolytic peptide. *Nat Struct Mol Biol* **19**: 105–113.
- 22 Sevin EW, Barloy-Hubler F. 2007. RASTA-Bacteria: A web-based tool for identifying toxin-
- antitoxin loci in prokaryotes. *Genome Biol* 8.
- 24 Sharma CM, Hoffmann S, Darfeuille F, Reignier J, Findeiß S, Sittka A, Chabas S, Reiche K,
- 25 Hackermüller J, Reinhardt R, et al. 2010. The primary transcriptome of the major human
- 26 pathogen Helicobacter pylori. *Nature* **464**: 250–255.

| 1 | Soldatov RA, Vinogradova S V., Mironov AA. 2014. RNASurface: Fast and accurate detection of |
|----|---|
| 2 | locally optimal potentially structured RNA segments. Bioinformatics 30: 457-463. |
| 3 | Stajich JE, Block D, Boulez K, Brenner SE, Chervitz SA, Dagdigian C, Fuellen G, Gilbert JGR, |
| 4 | Korf I, Lapp H, et al. 2002. The Bioperl toolkit: Perl modules for the life sciences. Genome |
| 5 | <i>Res</i> 12 : 1611–1618. |
| 6 | Stothard P, Wishart DS. 2005. Circular genome visualization and exploration using CGView. |
| 7 | <i>Bioinformatics</i> 21 : 537–539. |
| 8 | Tourasse NJ, Darfeuille F. 2020. Structural Alignment and Covariation Analysis of RNA |
| 9 | Sequences. Bio-Protocol 10: e3511. |
| 10 | Trotta E. 2014. On the normalization of the minimum free energy of RNAs by sequence length. |
| 11 | PLoS One 9. |
| 12 | Wang X, Yao J, Sun YC, Wood TK. 2021. Type VII Toxin/Antitoxin Classification System for |
| 13 | Antitoxins that Enzymatically Neutralize Toxins. Trends Microbiol 29: 388–393. |
| 14 | Weaver KE. 2012. The par toxin-antitoxin system from Enterococcus faecalis plasmid pAD1 and its |
| 15 | chromosomal homologs. RNA Biol 9: 1498–1503. |
| 16 | Weaver KE, Reddy SG, Brinkman CL, Patel S, Bayles KW, Endres JL. 2009. Identification and |
| 17 | characterization of a family of toxin-antitoxin systems related to the Enterococcus faecalis |
| 18 | plasmid pAD1 par addiction module. <i>Microbiology</i> 155: 2930–2940. |
| 19 | Wen J, Fozo EM. 2014. sRNA antitoxins: More than one way to repress a toxin. Toxins (Basel) 6: |
| 20 | 2310–2335. |
| 21 | Wen J, Won D, Fozo EM. 2014. The ZorO-OrzO type I toxin-Antitoxin locus: Repression by the |
| 22 | OrzO antitoxin. Nucleic Acids Res 42: 1930–1946. |
| 23 | Xie Y, Wei Y, Shen Y, Li X, Zhou H, Tai C, Deng Z, Ou HY. 2018. TADB 2.0: An updated |
| 24 | database of bacterial type II toxin-antitoxin loci. Nucleic Acids Res 46: D749–D753. |
| 25 | Yachdav G, Goldberg T, Wilzbach S, Dao D, Shih I, Choudhary S, Crouch S, Franz M, García A, |
| 26 | García LJ, et al. 2015. Anatomy of BioJS, an open source community for the life sciences. |

- 1 Elife **4**: 1–7.
- 2 Yachdav G, Wilzbach S, Rauscher B, Sheridan R, Sillitoe I, Procter J, Lewis SE, Rost B, Goldberg
- 3 T. 2016. MSAViewer: Interactive JavaScript visualization of multiple sequence alignments.

4 *Bioinformatics* **32**: 3501–3503.

5 Zuker M. 2003. Mfold web server for nucleic acid folding and hybridization prediction. *Nucleic*

6 *Acids Res* **31**: 3406–3415.

7 Zuker M. 1989. On finding all suboptimal foldings of an RNA molecule. *Science* **244**: 48–52.

8

9

1 Tables

- 2 Table 1. Identification by RNAMotif and RNASurface of toxin-antitoxin mRNA-asRNA pairs
- 3 corresponding to toxin ORF loci predicted by (Fozo et al. 2010).
- 4

| TA family | # of ORFs predicted | # of mRNA-asRNA pairs identified at the corresponding loci and with | # of mRNA-asRNA pairs identified at the corresponding loci but with |
|------------------------|---------------------|--|--|
| | ~y (=, | the corresponding family | an alternative family |
| Ldr/Rdl, | 298 (162, 136) | 175 (136, 39) | 16 (5, 11) |
| Fst/RNAII ^a | | | |
| Hok/Sok | 182 | 160 | 0 |
| Ibs/Sib | 210 | 182 | 0 |
| ShoB/OhsC | 31 | 26 | 3 |
| TisB/IstR1 | 45 | 39 | 0 |
| TxpA/RatA | 122 | 91 | 12 |

5

^aIn (Fozo et al. 2010) members of the Ldr/Rdl and Fst/RNAII families were classified together due
to remote similarity between toxin protein sequences of the two families (numbers in parentheses
correspond to counts that would be obtained if systems from Gram-negative and Gram-positive
organisms are assigned to the Ldr/Rdl and Fst/RNAII family, respectively).

10

1 Figure legends

- 2 **Fig. 1.** Automatic annotation of mRNAs and asRNAs of known type I TA systems.
- 3
- 4 Fig. 2. Screenshot of a "Locus Details" page in T1TAdb showing the various panels with sequence,
- 5 structure, and genomic information.
- 6

Fig. 3. Structural alignments of two subgroups of Sok asRNAs from various Gram-negative bacteria. The structure of the subgroup shown in panel A matches that reported by (Franch et al. 1997) while the subgroup shown in panel B exhibits an additional stem-loop element at the 5' end. Alignments were computed using MAFFT (and slightly corrected manually) and annotated with covariation information using R-chie. A region of 30 nt upstream of the RNAs was added that includes matches to the -10 promoter box motif (UANNNU, where N means any nucleotide) highlighted in dark orange.



T1TAdb Locus Details

| | Locu | s TAO | 0126 | | | | |
|-------------------------------|---|----------|------|------------|------------------------|------------|----------|
| Locus | | | | × | | | |
| ID | TA00126 | | | | genomic context (Gviz | (IGVis) | |
| Family | AapA/IsoA | | | 1 E00E mb | genomic context (OVIZ) | 1 5025 mb | |
| Organism | Helicobacter pylori 26695 chromosome | | - | 1.5025 110 | | 1.5035 110 | |
| Taxonomy | Bacteria; Proteobacteria; Epsilonproteobacteria; Campylobacterales; Helicobacteraceae; Helicobacter; Helicobacter pylori 26695 | 50 | | | 1.503 mb | | 1.504 mb |
| Genomic location | chromosome | 9 | ksaA | | | | HP1433 |
| Genomic coordinates | 1503082-1503303 | ozin ORF | | | | | |
| Genome accession (GenBank) | NC_000915.1 | adin T | | | АарАз-1 | | |
| Genome map | CGView | ÷ 8 | | | AapA3 | | |
| Genomic context | IGV.is | 5 S | | | | | |
| Upstream gene | ksgA : rRNA small subunit methyltransferase A : | Antit | | | 18043 | | |
| Downstream gene | HP1433 : hypothetical protein : | | | | 13043 | | |

Toxin mRNA

| | | predicted 2D structure (mfold / VARNA) |
|---|---|--|
| Name | AapA3 | |
| Locus_tag | | ů _. υ⊘ů ^{, ψ} `110 100− ∪⊙ů u⊝ů α⊂ů x − ∪ |
| Genomic coordinates | 1503082-1503303 (strand +) | $\begin{array}{cccccccccccccccccccccccccccccccccccc$ |
| Length | 222 nt | |
| Sequence (with 50 nt upstream, and -10 lov in red; mRNA8NA mRNA8NA madefilined, 078 in underlined, 078 in Underlined, 078 in Underlined, 078 in UPPERCASE) | ttgatgcggtctaggggtgtttaagggggtttgttgtaggatttcatca CGCCCCAIAGTIGGAAAGTGCAAGCGTTGCTTAGAT TGCCTTACTTGGCAAGCGTTGCTTGGGTGGGGGGGGGG | 80 - a - U 80 - a - U 90 - b |
| Predicted 2D structure | PDE PNG | (A) → U (A) → U (A) → U (C) = 0 (C) = 0 (C |

Antitoxin asRNA

| Name | IsoA3 |
|--|--|
| Locus_tag | |
| Genomic coordinates | 1503082-1503160 (strand -) |
| Length | 79 nt |
| Sequence (with 50 nt upstream, and -10 box in red; transcript sequence in UPPERCASE) | aaacgatcactttaagccccaaaaagcaaaatcaaagtataatgtttc CAAGAGCGTTGCCACTITGTGTTCATGGCATGCTCCTTTGACATAG GATTGCCCCATATTGCACATAGGGGGG → <u>BLAST sequence</u> |
| Predicted 2D structure | PDF, PNG |



predicted 2D structure (PSIPRED / POLYVIEW-2D)

hydrophobicity plot (ProtScale)

ProtScale output for user_sequence Hphob. / Kyte & Doolittle

 \cong

15 20 25 Position $H - \alpha$ and other helices $E - \beta$ -strand or bridge

MKHKSGKRSWKTLYFEFAFLGLKVIVSVKR

Toxin Peptide

ORF length (peptide length) 93 nt (30 aa)

 $\begin{array}{l} \text{Atgaaacacaaaagtggcaaacgctcttggaaaacattatactttgggt}\\ ORF sequence & \texttt{Tgctttttggggcttaaagtggtaagtggtaaagtggt}\\ \rightarrow \underline{BLAST} sequence \\ \end{array}$

Peptide sequence $\xrightarrow{MKHKSGKRSWKTLYFEFAFLGLKVIVSVKR*} \rightarrow \frac{BLAST sequence}{BLAST sequence}$

Predicted 2D structure PDF, GIF

Hydrophobicity plot PDF, <u>GIF</u>

⇒**₽**lab

T TA

INSERM U1212, CNRS UMR5320, University of Bordeaux



Fig. 2

Home Keyword Search Sequence Search Contact References

Sponsors Links

Fig. 3A

TA06056 A salmonicida A449 TA06101 P mirabilis HT4320 TA06105 S boydii CDC 3083-94 TA06046 E coli UTI89 TA06063 C sakazakii ATCC BAA-894 TA05975 E coli str. K-12 substr. MG1655 TA05975 E coli str. K-12 substr. MG1655 TA05970 E coli str. K-12 substr. W3110 K-12 TA05070 E coli str. K12 substr. DH10B TA06088 E coli str. K12 substr. DH10B TA06098 E coli Str. K12 substr. DH10B TA06098 E coli Str. K12 substr. DH10B TA06098 E coli Str. K12 substr. DH10B TA06078 E coli Str. K12 substr. DH10B TA06088 E coli Str. K12 substr. DH10B TA06108 S hoydIi CDC 3083-94 TA06108 E coli ATCC 8739 TA0613 E coli ATCC 8739 TA06140 E coli D157:H7 str. Sakai TA06140 E coli D157:H7 str. Sub93 genome. TA06012 E coli D157:H7 str. Sub93 genome. TA06012 E coli O157:H7 str. Str. Sakai TA06110 S hoydIi CDC 3083-94 TA06110 S hoydIi CDC 3083-94 TA06107 E albertIi TW07627 scf 1109867046288 TA06157 Y kristensēnii ATCC 33638 contig00012 TA05987 E coli O157:H7 str. SDL933 genome. TA06985 E coli O157:H7 str. Sub933 genome. TA06985 E coli O157:H7 str. Sub1933 genome. TA06985 E coli O157:H7 str. SDL933 genome. TA06985 E coli O157:H7 str. SDL933 genome. TA06985 E coli O157:H7 str. SDL933 genome. TA06000 E coli O157:H7 str. SDL933 genome. TA06000 E coli O157:H7 str. Str. BDL933 genome. TA06145 E coli O157:H7 str. Str. SAkai TA06145 E coli O157:H7 str. Str. SAkai TA06145 E coli O157:H7 str. Str. SAkai TA06145 E coli O157:H7 str. EC113 gcontig 1105762721293 TA06146 E coli O157:H7 str. EC369 gcontig 1106613679340 TA06147 E coli O157:H7 str. EC4486 gcontig 1106613679340 TA06149 E coli O157:H7 str. EC4486 gcontig 1106603634700 TA07122 E coli TA06145 E coli STR7 str. EC4486 gcontig 1106603634700 TA07122 E coli TA06031 E coli STR7 str. EC4486 gcontig 1106603634700 TA07031 E coli STR7 str. ST TA06071 E coli HS TA07122 E coli TA06031 E coli EH41 TA06114 S enterica serovar Kentucky str. CVM29188 TA06115 S enterica serovar Heidelberg str. SL476 TA06026 S enterica serovar Typhimurium TA05977 S sonnei TA05377 S Sonnel TA07128 E coli ECOR24 ECOR24 contig 19 consensus TA06126 E coli 0127:H6 E2348769 E2348/69. TA06017 S enterica serovar Typhi str. CT18 TA05981_S_typhi TA06158_Y_aldovae_ATCC_35236_contig00244 TA06144_Y_bercovieri_ATCC_43970_contig01191 TA06123_E_coli_SE11_DNA TA06142_E_coli_E110019_gcontig_1112495928642 TA06142⁻E⁻coli⁻E110019 gcontig_1112495928642 TA06072⁻E⁻coli⁻E24377A TA06040⁻S⁻boydIi Sb227 TA07126⁻E⁻coli⁻C⁻ TA06106⁻S⁻boydIi⁻CDC 3083-94 TA06034⁻S⁻sonnei⁻Ss046 TA06074⁻S⁻proteamaculans 568 TA06153⁻C⁻youngae ATCC_29220⁻C_sp-1.0_Cont2.17 TA07125⁻E⁻coli⁻ TA07123⁻E⁻coli⁻ TA07123⁻E⁻coli⁻ TA07129⁻H⁻alveI TA06171⁻P⁻mirabIlis ATCC_29906_SCAFFOLD3 TA06100⁻P⁻mirabIlis⁻HI4320 TA06100 P_mirabilis HI4320 TA06100 P milabilis ATCC 29906 SCAFFOLD31 TA06160 P rettgeri DSM 1T31 P rettgeri-1.0_Cont16.13 TA06005 E coli 0157:H7 str. Sakai TA06154 E cancerogenus ATCC_35316 E cancerogenus-1.0_Cont1.4 TA06134 E cancelogenus AICC_33316 E cancel TA06102 E coli TA06076 E coli APEC 01 TA05980 E coli B171 TA06141 E coli E22 gcontig_1112495653200 TA06078_E_coli____ TA07124_E_coli_K-12 TA06008 E coli 0157:H7 str. Sakai TA06007 E coli 0157:H7 str. Sakai TA05992 E coli 0157:H7 str. EDL933 genome. TA06173 E coli 83972 SCAFFOLD3 TA06044 E coli UTI89 TA06044 E COIL UT189 TA05998 E COLL 017:H7 str. Sakai TA06125 E COLL 0127:H6 E2348/69 E2348/69. TA06002 E COLL 0157:H7 str. Sakai TA06170 E COLL 0157:H7 str. EC4024 scf_1109799330226_genomic_scaffold TA06033 S sonnei Ss046 TA06031 S flourozzi 20 str. 24575 TA06021 S flexneri 2a str. 2457T TA06041 S boydii Sb227 TA05989 E coli OI57:H7 str. EDL933 genome. TA05986 E coli OI57:H7 str. EDL933 genome. TA05112_E_coli_ TA06112_E_coli_ TA06165_P_angustum_S14_1099604003227 TA06116_K_pneumoniae_342



Conservation Covariation One-sided

Covariation [-2,-1] (-1,0] (0,1] (1,2]

Invalid

Unpaired

Gap

Ambiguous

Fig. 3B

Downloaded from rnajournal.cshlp.org on April 28, 2024 - Published by Cold Spring Harbor Laboratory Press

TA05974_E_coli_str._K-12_substr._MG1655 TA06070 E coli Str. K-12 Substr. MG105 TA06070 E coli HS TA06139 E sp. T 1 43 supercont1.22 TA06087 E coli Str. K12 substr. DH10B TA06097 E coli SMS-3-5 TA05968 E coli str. K-12 substr. W3110 K-12 TA06143 E coli 53638 gcontig 1105238512145 TA06084 E coli str. K12 substr. DH10B TA06084 E coli str. K12 substr. DH10B TA06020 S flexneri Za str. 2457T TA05966 E coli str. K-12 substr. W3110 K-12 TA05972 E coli str. K-12 substr. MG1655 TA06082 E coli ATCC 8739 TA06147 E coli O157:H7 str. EC4113 gcontig_1105762727453 TA06037 S dysenteriae Sd197 chromosome TA06001 E coli O157:H7 str. Sakai TA05988 E coli 0157:H7 str. Sakai TA05988 E coli 0157: H7 str. EDL933 genome. TA06039 S boydii Sb227 TA06132_E_fergusonii_ATCC_35469_chromosome TA06132 E Tergusonii Arcc 35469 Cr TA06066 E coli HS TA06032 S sonnei Ss046 TA06045 E coli UTI89 TA06138 E sp. I 1 43 supercont1.1 TA06077 E coli APEC 01 TA06009 E coli 0157:H7 str. Sakai TA05993 E coli 0157:H7 str. EDL933 genome. TA05993 E COIL 0157:H7 Str. EDL933 genome. TA06137 K pneumoniae NTUH-KZ044 DNA TA06060 K pneumoniae MGH 78578 TA06117 K pneumoniae 342 TA05965 E coli str. K-12 substr. W3110 K-12 TA06083 E coli str. K12 substr. DH10B TA05971 E coli str. K-12 substr. MG1655 TA057127 E coli str. K-12 substr. MG1655 TA07127 E coli strain B TA05996 E coli 0157:H7 str. EDL933 genome. TA05997 E coli 0157:H7 str. EDL933 genome. TA05997 E COIL 0157:H7 Str. EDL933 TA06131 S baltica OS223 TA06130 S baltica OS223 TA06122 E coli SEIL DNA TA05976 E phage 933W TA05999 E coli 0157:H7 str. Sakai TA06006 E coli 0157:H7 str. Sakai TA06164 E coli BL21(DE3) ctg71 TA06165 enterica serovar Twohis TA06164 E coli BL21(DE3) ctg71 TA06016 S enterica serovar Typhi str. CT18 TA06120 S enterica serovar Enteritidis str. P125109 TA06036 S dysenteriae Sd197 chromosome TA05991 E coli 0157:H7 str. EDL933 genome. TA06004 E coli 0157:H7 str. Sakai TA06057 E sp. 638 TA06035 E coli 0157:H7 str. EDL933 TA05978 E coli 0157:H7 str. Sakai TA06043 E coli 0157:H7 str. Sakai TA06043 E coli UTI89 TA05984 E coli K-12 TA06136 E coli chi7122 TA06090 E coli SMS-3-5 TA06042 E coli TA06118 K pneumoniae 342 TA07121 E coli TA06075_S_proteamaculans_568 TA05979_S_flexneri_2b TA06073 S proteamaculans 568 TA06156 P_rustigianii_DSM_4541 P_rustigianii-1.0_Cont4.1



Conservation Covariation

Invalid

Unpaired



T1TAdb: the database of Type I Toxin-Antitoxin systems

Nicolas J Tourasse and Fabien Darfeuille

RNA published online September 16, 2021

| Supplemental Material | http://rnajournal.cshlp.org/content/suppl/2021/09/16/rna.078802.121.DC1 | |
|--|--|--|
| P <p< th=""><th>Published online September 16, 2021 in advance of the print journal.</th></p<> | Published online September 16, 2021 in advance of the print journal. | |
| Accepted Manuscript | Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version. | |
| Open Access | Freely available online through the RNA Open Access option. | |
| Creative Commons License | This article, published in <i>RNA</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at http://creativecommons.org/licenses/by-nc/4.0/ . | |
| Email Alerting Service | Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here. | |



To subscribe to RNA go to: http://rnajournal.cshlp.org/subscriptions