

Supporting Information

Salomon et al. 10.1073/pnas.1406110111

SI Materials and Methods

Bacterial Strains and Media. The *Vibrio parahaemolyticus* RIMD 2210633 derivative strain POR1 (RIMD 2210633 Δ tdhAS) (1) and its derivatives, as well as *Vibrio alginolyticus* 12G01 and its derivatives, were routinely cultured in marine LB (MLB) broth [LB broth containing 3% (wt/vol) sodium chloride] or on marine minimal media agar (2) at 30 °C. *Escherichia coli* DH5 α and *E. coli* S17 (λ pir) were routinely cultured in 2 \times yeast-extract and tryptone broth at 37 °C. The medium was supplemented with kanamycin (30 μ g/mL for *E. coli* and 250 μ g/mL for *Vibriosis* species) or chloramphenicol (25 μ g/mL), when necessary, to maintain a plasmid.

Plasmids. For arabinose-inducible expression of *vp1389*, *vp1416*, and *vti2*, the genes coding the sequences were amplified and cloned into the pBAD/Myc-His vector (Invitrogen), in which the antibiotic resistance was changed from ampicillin to kanamycin, or into the pBAD33 vector (Addgene). Empty pBAD/Myc-His and pBAD33 vectors were used to provide resistance to kanamycin and chloramphenicol, respectively.

Construction of Deletion Strains. Gene deletions were performed as previously described using the pDM4 vector (3, 4) for *vp1388*, *vp1388-vp1389*, *vp1415*, *vp1415-vp1416*, *vpa1263*, *vpa1263-vti2*, *v12G01_02265*, *v12G01_02265-v12G01_02260*, and *v12G01_01540* (*V. alginolyticus hcp1*). Construction of POR1 Δ *hcp1* was described previously (4).

Bacterial Killing Assays. Bacterial killing assays were performed as previously described (4). Cells were cocultured at an OD₆₀₀ ratio of 1:4 (prey/attacker) on MLB (for *V. parahaemolyticus* attacker) or LB (for *V. alginolyticus* attacker) agar plates for 4 h at 30 °C in triplicate. Where appropriate, plates were supplemented with 0.1% arabinose to induce expression from plasmids. Each experiment was done in triplicate and repeated at least twice, with similar results. A representative experiment is presented. A two-tailed Student *t* test was used to determine significance between indicated sample groups unless stated otherwise.

Proteomics and MS. *V. parahaemolyticus* cultures of POR1 Δ *opaR* [type 6 secretion system (T6SS1)⁺] and POR1 Δ *opaR* Δ *hcp1* (T6SS1⁻) were grown in triplicate in 50 mL of M9 minimal media with 3% (wt/vol) NaCl media containing 20 μ M phenamil for 5 h at 30 °C. Media were collected, and protein was precipitated as previously described (4). Protein samples were run 10 mm into the top of an SDS/PAGE gel, stained with Coomassie Blue, and excised. Overnight digestion with trypsin (Promega) was performed after reduction and alkylation with DTT and iodoacetamide (Sigma-Aldrich). The resulting samples were analyzed by tandem MS using a QExactive mass spectrometer (Thermo Electron) coupled to an Ultimate 3000 RSLC-Nano liquid chromatography system (Dionex). Peptides were loaded onto a 180- μ m i.d., 15-cm long, self-packed column containing 1.9 μ m C18 resin (Dr. Maisch, Ammerbuch, Germany) and eluted with a gradient of 0–40% buffer B for 60 min. Buffer A consisted of 2% (vol/vol) acetonitrile (ACN) and 0.1% formic acid in water. Buffer B consisted of 80% (vol/vol) ACN, 10% (vol/vol) trifluoroethanol, and 0.08% formic acid in water. The QExactive mass spectrometer acquired up to 10 high-energy, collision-induced dissociation fragment spectra for each full spectrum acquired.

Raw MS data files were converted to peak list format using ProteoWizard msconvert (version 3.0.3535) (5). The resulting files were analyzed using the central proteomics facilities pipeline (CPFP), version 2.0.3 (6, 7). Peptide identification was performed using the X!Tandem (8) and open MS search algorithm (OMSSA) (9) search engines against a database consisting of *V. parahaemolyticus* RIMD 2210633 sequences from the National Center for Biotechnology Information, with common contaminants and reversed decoy sequences appended (10). Fragment and precursor tolerances of 20 ppm and 0.1 Da were specified, and three missed cleavages were allowed. Carbamidomethylation of Cys was specified as a fixed modification, and oxidation of Met was specified as a variable modification. Label-free quantitation of proteins across samples was performed using SING normalized spectral index software (11).

To identify statistically significant differences in protein amount between T6SS1⁺ and T6SS1⁻ strains, SING quantitation results for three biological replicates per strain were processed using the power law global error model (PLGEM) package in R (12, 13). Protein identifications were filtered to an estimated 1% protein false discovery rate (FDR) using the concatenated target-decoy method (10). An additional requirement of two unique peptide sequences per protein was imposed, resulting in a final protein FDR <1%. Spectral index quantitation was performed using peptide-to-spectrum matches (PSMs) with a *q*-value \leq 0.01, corresponding to a 1% FDR rate for PSMs.

The dataset of tandem MS results was uploaded to the PeptideAtlas repository (www.peptideatlas.org/PASS/PASS00442; accession no. PASS00442).

Bioinformatics. The N-terminal sequence from VP1388 (GI|28898162, range 144–297) identifies conserved sequence motifs (PhhPhR and GhhYhhh) near the N terminus of numerous bacterial sequences. One of these sequences (GI|520945140, range 1–249) was queried against the nonredundant (NR) database using the position-specific iterative (PSI)-BLAST application of BLAST+ (14) [10 iterations, E-value cutoff of 0.005, collecting all subject GIs (NCBI protein sequence identifier)] to collect all potential MIX (marker for type six effectors) sequences. Collected sequences were filtered to include only those from complete genomes in the Reference Sequences (RefSeq) database (15) and were clustered using CLANS (16). Sequences from five resulting broad groups were independently aligned using the MAFFT server (17), and were colored according to conservation using Jalview (18) to highlight hydrophobicity patterns surrounding the conserved motifs. The corresponding domain sequences from each of the five multiple sequence alignments were aligned manually using the motifs and surrounding hydrophobicity patterns and PROMALS3D (19) alignments of representatives as a guide. Motifs were highlighted with WebLogo (20).

The domain organization of representative sequences was defined using Batch Entrez and conserved domains (CD)-search (21) with default cutoffs. Identified domains were sorted according to E-value, and incomplete domains with E-values >1.6e-4 were excluded from counts (excluded partial domains include mainly low-complexity sequence, such as coiled coils and transmembrane helices). Two low-complexity coil domains with E-values better than the cutoff were also excluded. Domains identified within this lower E-value range (0.01–1.6e⁻⁴) that were also identified in other sequences with higher confidence were included in the counts (i.e., 38 LysM domains were identified with an E-value range from 0.00015 to 6.8e⁻¹⁴). Unknown

sequences represent 855 sequences that identified no additional domains and 58 sequences that identified domains with low-complexity regions that were excluded. Given the presence of a number of excluded transmembrane helix (TMH)-containing domain predictions, Phobius (22) was used to predict TMHs in the C-terminal sequence of all unknown MIX-containing proteins. Several of the sequences with TMH-containing domains predicted by CD-search that were excluded due to low confidence were submitted to the HHPred server (23) for domain validation. All of the submitted sequences identified various colicin pore-forming toxins (ranging from 9–93.3% probability) in their predicted TMH-containing regions.

T6SS components VrgG (GI|28898168, range 1–525) and haemolysin coregulated proteins (Hcp) (GI|256599595, range 1–163) were queried against the NR database (two iterations with an E-value cutoff of 0.0001 and 10 iterations with an E-value cutoff of 0.001, respectively, collecting all subject GIs) using PSI-BLAST (14). Subject sequences were filtered to include only those from complete genomes in the RefSeq database (15), and VrgG sequences were clustered using CLANS (16) to purge false-positive hits. Species containing the resulting VrgG and Hcp sequence representatives were sorted and compared with those containing MIX using GeneVenn (24).

Genomic neighborhoods encompassing the T6SS machinery were explored using the Microbial Genome Database for Comparative Analysis (25) and the Search Tool for the Retrieval of Interacting Genes (STRING) (26). The STRING database v9.1 describing protein association networks was used to generate genome neighborhood information for (i) core T6SS components and (ii) MIX-containing groups. For the core T6SS components, a network was generated in the clusters of orthologous groups (COG) mode from Hcp (COG3157). Only neighborhood and gene fusion scores with a high-confidence cutoff (>0.7) were considered for generating the T6SS network. MIX sequences for each group were submitted to the STRING database using “multiple sequences” submission in the COG mode to identify relevant groups for neighborhood analysis. The nonorthologous group populated with the most MIX sequences was selected for neighborhood analysis, which was limited to neighborhood and gene fusion scores with a medium-confidence cutoff for network generation. For each of the five MIX groups, combined STRING scores were ranked, and all links with high confidence scores are shown in Table 1. For links of lower confidence, only those scores up to and including the top core T6SS component were reported. For the MIX II group, a low-confidence link to a T6SS component was also reported.

- Park KS, et al. (2004) Functional characterization of two type III secretion systems of *Vibrio parahaemolyticus*. *Infect Immun* 72(11):6659–6665.
- Eagon RG (1962) *Pseudomonas natriegens*, a marine bacterium with a generation time of less than 10 minutes. *Journal of Bacteriology* 83:736–737.
- Salomon D, et al. (2013) Effectors of animal and plant pathogens use a common domain to bind host phosphoinositides. *Nat Commun* 4:2973.
- Salomon D, Gonzalez H, Updegraff BL, Orth K (2013) *Vibrio parahaemolyticus* type VI secretion system 1 is activated in marine conditions to target bacteria, and is differentially regulated from system 2. *PLoS ONE* 8(4):e61086.
- Kessner D, Chambers M, Burke R, Agus D, Mallick P (2008) ProteoWizard: Open source software for rapid proteomics tools development. *Bioinformatics* 24(21):2534–2536.
- Trudgian DC, Mirzaei H (2012) Cloud CFP: A shotgun proteomics data analysis pipeline using cloud and high performance computing. *J Proteome Res* 11(12): 6282–6290.
- Trudgian DC, et al. (2010) CFP: A central proteomics facilities pipeline. *Bioinformatics* 26(8):1131–1132.
- Craig R, Beavis RC (2004) TANDEM: Matching proteins with tandem mass spectra. *Bioinformatics* 20(9):1466–1467.
- Geer LY, et al. (2004) Open mass spectrometry search algorithm. *J Proteome Res* 3(5): 958–964.
- Elias JE, Gygi SP (2007) Target-decoy search strategy for increased confidence in large-scale protein identifications by mass spectrometry. *Nat Methods* 4(3):207–214.
- Trudgian DC, et al. (2011) Comparative evaluation of label-free SINQ normalized spectral index quantitation in the central proteomics facilities pipeline. *Proteomics* 11(14):2790–2797.
- Pavelka N, et al. (2008) Statistical similarities between transcriptomics and quantitative shotgun proteomics data. *Mol Cell Proteomics* 7(4):631–644.
- Pavelka N, et al. (2004) A power law global error model for the identification of differentially expressed genes in microarray data. *BMC Bioinformatics* 5:203.
- Camacho C, et al. (2009) BLAST+: Architecture and applications. *BMC Bioinformatics* 10:421.
- Pruitt KD, Tatusova T, Brown GR, Maglott DR (2012) NCBI Reference Sequences (RefSeq): Current status, new features and genome annotation policy. *Nucleic Acids Res* 40(Database issue):D130–D135.
- Frickey T, Lupas A (2004) CLANS: A Java application for visualizing protein families based on pairwise similarity. *Bioinformatics* 20(18):3702–3704.
- Katoh K, Standley DM (2013) MAFFT multiple sequence alignment software version 7: Improvements in performance and usability. *Mol Biol Evol* 30(4):772–780.
- Waterhouse AM, Procter JB, Martin DM, Clamp M, Barton GJ (2009) Jalview Version 2—A multiple sequence alignment editor and analysis workbench. *Bioinformatics* 25(9): 1189–1191.
- Pei J, Tang M, Grishin NV (2008) PROMALS3D web server for accurate multiple protein sequence and structure alignments. *Nucleic Acids Res* 36(Web Server issue):W30–W34.
- Crooks GE, Hon G, Chandonia JM, Brenner SE (2004) WebLogo: A sequence logo generator. *Genome Res* 14(6):1188–1190.
- Marchler-Bauer A, et al. (2011) CDD: A Conserved Domain Database for the functional annotation of proteins. *Nucleic Acids Res* 39(Database issue):D225–D229.
- Käll L, Krogh A, Sonnhammer EL (2004) A combined transmembrane topology and signal peptide prediction method. *J Mol Biol* 338(5):1027–1036.
- Hildebrand A, Remmert M, Biegert A, Söding J (2009) Fast and accurate automatic structure prediction with HHPred. *Proteins* 77(Suppl 9):128–132.
- Pirooznia M, Nagarajan V, Deng Y (2007) GeneVenn—A web application for comparing gene lists using Venn diagrams. *Bioinformatics* 1(10):420–422.
- Uchiyama I, Mihara M, Nishide H, Chiba H (2013) MGD update 2013: The microbial genome database for exploring the diversity of microbial world. *Nucleic Acids Res* 41(Database issue):D631–D635.
- Franceschini A, et al. (2013) STRING v9.1: Protein-protein interaction networks, with increased coverage and integration. *Nucleic Acids Res* 41(Database issue):D808–D815.

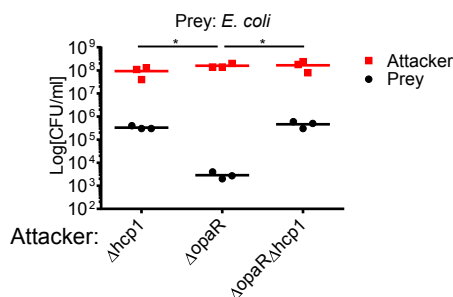


Fig. S1. *V. parahaemolyticus* $\Delta opaR$ strain kills *E. coli* in a T6SS1-dependent manner. Viability counts of *E. coli* prey 4 h after coculturing with *V. parahaemolyticus* POR1 $\Delta hcp1$, POR1 $\Delta opaR$, or POR1 $\Delta hcp1\Delta opaR$ attacker strains. Mixed cultures at an OD₆₀₀ ratio of 1:4 (prey/attacker) were spotted on MLB plates at 30 °C for 4 h. Asterisks mark statistical significance between prey sample groups at $t = 4$ h ($P < 0.05$).

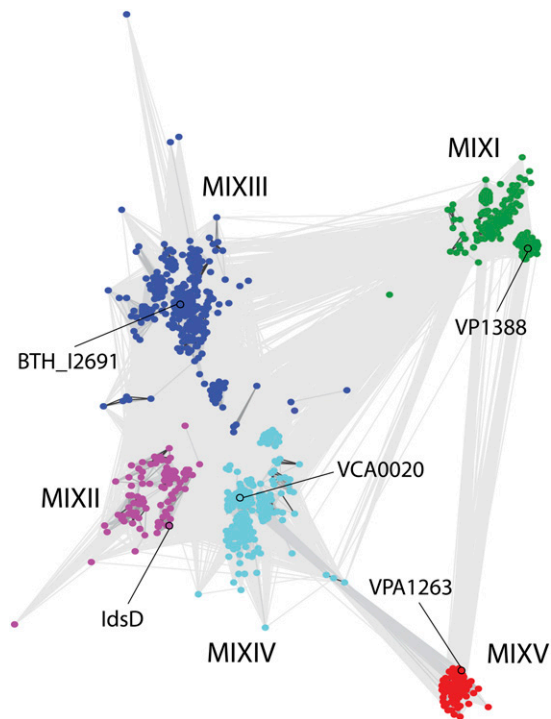


Fig. S4. MIX regions are grouped into five clans. Pairwise BLAST scores of all identified MIX sequences were used to cluster similar sequences into groups depicted in a 2D representation. Nodes are colored and labeled according to visual grouping: MIX I (green), MIX II (magenta), MIX III (blue), MIX IV (cyan), and MIX V (red), with connections between nodes representing pairwise BLAST scores better than $E = 0.0001$. MIX sequences of known T6SS effectors are labeled.

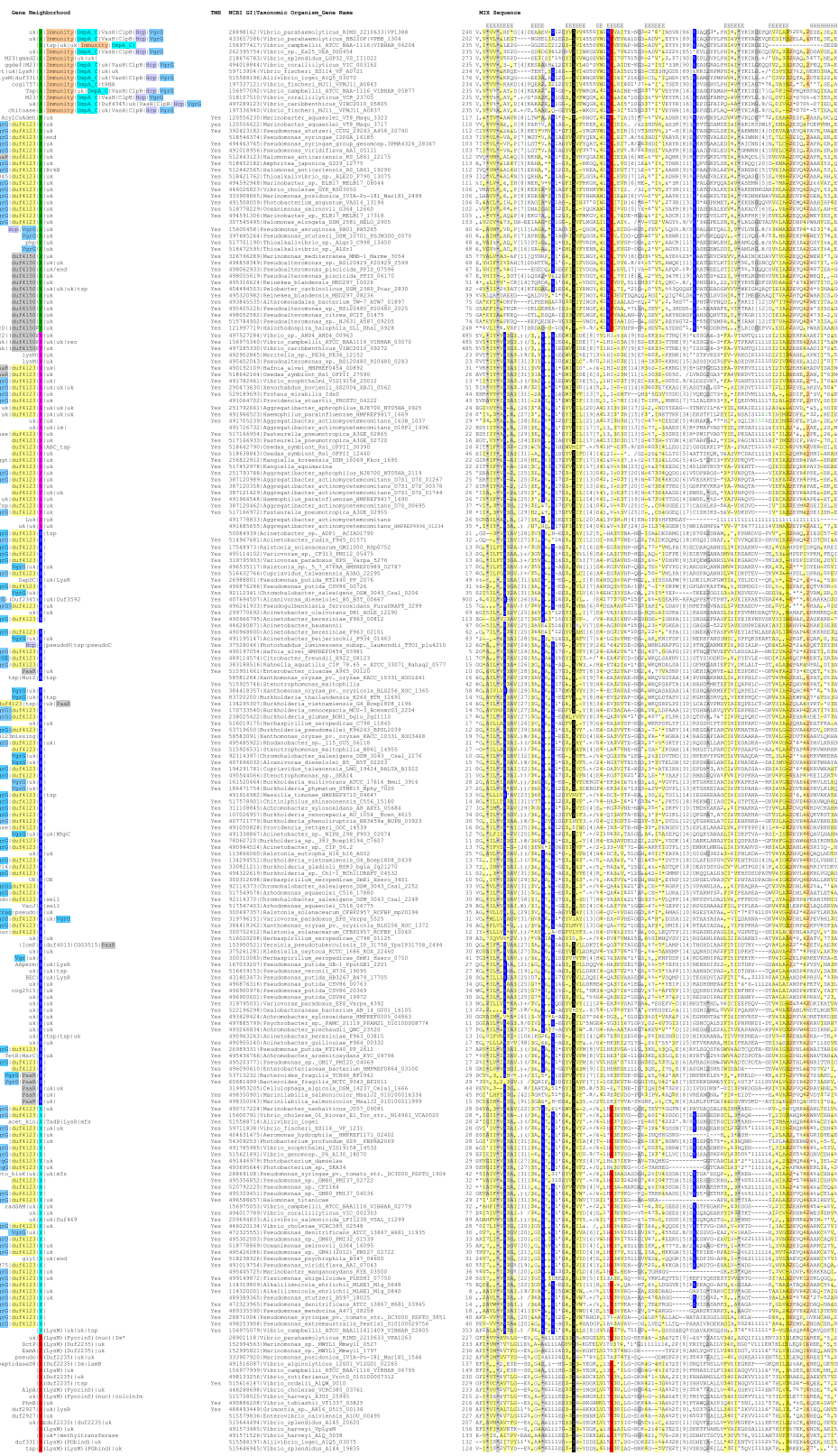


Fig. S5. Alignment of MIX sequences. Multiple sequence alignment of MIX representatives (partial sequences and redundant sequences removed at 90% identity cutoff) is shown. Sequences are labeled to the left by National Center for Biotechnology Information GI number, species, and gene name, with residue numbers corresponding to the first and last alignable residues indicated to the left and right of the sequence, respectively. Domain neighborhoods surrounding the MIX-encoding gene (represented by X, colored according to group as in Fig. S4) are indicated on the left, with T6SS components highlighted in blue (VgrG), lavender (Hcp), and gray (proline-alanine-alanine-arginine), and linked Duf4123 highlighted in yellow. uk, genes with unknown domains. The TMH column represents MIX-encoding genes that have predicted TMHs. Residue positions are highlighted according to group-wise conservation: mainly hydrophobic (yellow), mainly small (gray), mainly basic (blue), mainly acidic (red), mainly S/T (orange), and MIX motif Y (black). Unalignable regions between conserved core secondary structures (indicated above alignment) are omitted, with the number of deleted residues shown in brackets.

