**High diversity and variability of pipolins among a wide range of pathogenic**
*Escherichia coli* **strains**

Saskia-Camille Flament-Simon[1*], María de Toro[2*], Liubov Chuprikova[4*], Miguel Blanco[1], Juan Moreno-González[3], Margarita Salas[3‡], Jorge Blanco[1], and Modesto Redrejo-Rodríguez[3,4#]

[1] Laboratorio de Referencia de *E. coli* (LREC), Departamento de Microbiología y Parasitología, Facultad de Veterinaria, Universidad de Santiago de Compostela (USC), Lugo, 27002,Spain.

[2] Plataforma de Genómica y Bioinformática, CIBIR (Centro de Investigación Biomédica de La Rioja), Logroño, 26006 La Rioja, Spain.

[3] Centro de Biología Molecular Severo Ochoa (Consejo Superior de Investigaciones Científicas-Universidad Autónoma de Madrid), 28049, Madrid, Spain.

[4] Departamento de Bioquímica, Universidad Autónoma de Madrid (UAM), Instituto de Investigaciones Biomédicas "Alberto Sols" CSIC-UAM, 28029, Madrid, Spain.

*Equal contribution

‡Deceased

#Correspondence to modesto.redrejo@uam.es

# SUPPLEMENTARY INFORMATION

**Supplementary Tables**

**Table S1 (XLS file). Characterization of pipolin-harboring LREC strains.**
Compilation of main features determined for LREC strains. Biosample and Enterobase
Uberstrain reference IDs are also indicated. D*: *Does not match with PCR data.
N.D. : non-detected.* Virulence factors detected by PCR screening and in silico
analysis are detailed.

**Table S2 (XLS file). Genbank genomes carrying pipolins.**
References for genomes and biosamples as well as main features from Enterobase
[21] are indicated.
*CRISPR/Cas cassettes and Integrons were surveyed as indicated in Methods.
**Integrity of detected integrons were analyzed with IntegronFinder [56] as indicated in
Materials and Methods. C, complete; In0, integron lacking attC site; CALIN, integron
lacking functional integrase gene.
***When available, PubMed ID of strain reporting publication

| Gene | Pipolins (%) | Annotation | UniProtKB HHPred best hit | eggNOG Description | KEGG KO | KEGG Description |
|------|-------------|------------|---------------------------|--------------------|---------|------------------|
| pipolB | 92 (100%) | Primer-independent DNA polymerase PolB | P03680 | | | |
| xerC_2 | 90 (98%) | Tyrosine recombinase XerC | P0A8P8 | Belongs to the 'phage' integrase family | | |
| group_1 | 87 (95%) | Uracil-DNA glycosylase | Q96YD0 | | | |
| group_6 | 86 (93%) | hypothetical protein | | | | |
| xerC_1 | 84 (85%) | Tyrosine recombinase XerC | P03700 | Belongs to the 'phage' integrase family | | |
| hisF | 78 (83%) | Type I restriction modification system methyltransferase (hsdM) | Q5M500 | HsdM N-terminal domain | K03427 | hsdM; type I restriction enzyme M protein [EC:2.1.1.72] |
| group_16 | 76 (82%) | metallohydrolase | Q57587 | Metal-dependent hydrolase | K07043 | uncharacterized protein |
| group_18 | 75 (82%) | hypothetical protein | | | | |
| group_5 | 75 (82%) | Type I site-specific deoxyribonuclease (hsdR) | P10486 | Type I restriction enzyme R protein N terminus (HSDR_N) | K01153 | hsdR; type I restriction enzyme, R subunit [EC:3.1.21.3] |
| group_10 | 75 (82%) | Protein of unknown function (DUF2787) | | Protein of unknown function (DUF2787) | | |
| group_13 | 74 (80%) | hypothetical protein | | Protein of unknown function (DUF726) | | |
| group_11 | 73 (79%) | hypothetical protein | | | | |
| group_24 | 73 (79%) | hypothetical protein | | | | |
| group_52 | 73 (79%) | Excisionase | A6T888 | | | |
| group_3 | 64 (70%) | hypothetical protein | | | | |
| group_8 | 59 (64%) | hypothetical protein | | | | |
| group_19 | 56 (61%) | WYL domain | A0A4Y3NDN0 | transcriptional regulator | | |
| group_58 | 53 (58%) | Znf/thioredoxin_put domain-containing protein | Q9A679 | | | |
| group_28 | 41 (45%) | Uncharacterized protein family (UPF0149) | P28366 | Uncharacterised protein family (UPF0149) | K07039 | uncharacterized protein |
| group_17 | 38 (41%) | IS1 family transposase IS1A | | cog cog3677 | | |
| group_23 | 31 (34%) | PD-(D/E)XK nuclease superfamily | | PD-(D/E)XK nuclease superfamily | | |
| group_31 | 30 (33%) | Restriction endonuclease | A0A0J9X157 | Restriction endonuclease | | |
| group_15 | 30 (33%) | IS1 family transposase IS1X2 | | cog cog1662 | K07480 | insertion element IS1 protein InsB |
| group_9 | 29 (32%) | Protein of unknown function (DUF4011) | | Protein of unknown function (DUF4011) | | |

| | | | | | | |
|---|---|---|---|---|---|---|
| group_34 | 24 (26%) | hypothetical protein | | type I restriction enzyme, R | | |
| group_53 | 23 (25%) | Protein of unknown function DUF262 | | Protein of unknown function (DUF1524) | | |
| group_14 | 22 (24%) | Uncharacterized protein family (UPF0149) | | | | |
| group_20 | 22 (24%) | WYL-domain containing protein | A0A4Y3NDN0 | | | |
| group_39 | 22 (24%) | Protein of unknown function DUF262 | | Protein of unknown function (DUF1524) | | |
| group_40 | 22 (24%) | hypothetical protein | | | | |
| group_25 | 22 (24%) | Protein of unknown function DUF262 | | Protein of unknown function (DUF1524) | | |
| group_55 | 21 (23%) | Protein of unknown function DUF262 | | Protein of unknown function DUF262 | | |
| group_32 | 19 (21%) | Znf/thioredoxin_put domain-containing protein | Q9A679 | | | |
| group_27 | 16 (17%) | Type I restriction modification enzyme | Q8R9Q6 | Type I restriction modification DNA specificity domain | K01154 | hsdS; type I restriction enzyme, S subunit [EC:3.1.21.3] |
| group_29 | 16 (17%) | IS3 family transposase ISEam1 | A0A0G3QIX7 | Transposase | K07483 | transposase |
| insK | 78 (16%) | IS3 family transposase ISEc14 | T0PD67 | silverDB | | |

**Table S3. Functional characterization of the most common pipolin genes.**

Annotation of genes from Roary shell-genome (present in more than 15% of pipolins) is indicated.
Functional groups in eggNOG and KEGG databases, as well as HHPred searches were also performed for
a detailed functional characterization.

**Supplementary Figures**

Figure S1

**Figure S1. Phylogeny of piPolB genes from analyzed pipolins.**
Nucleotide sequence of piPolB genes from new pipolins in LREC strains (cyan) and those retrieved from GenBank (orange) were aligned using Prank codon aware option and then used for maximum-likelihood phylogeny reconstruction with IQtree. Modelfinder Best-fit model was K3Pu+F+R2.
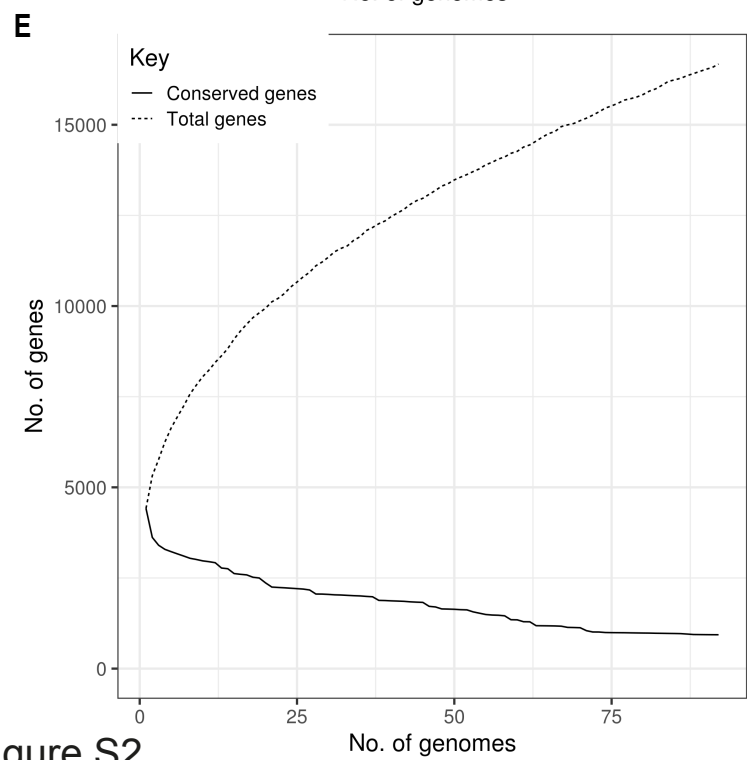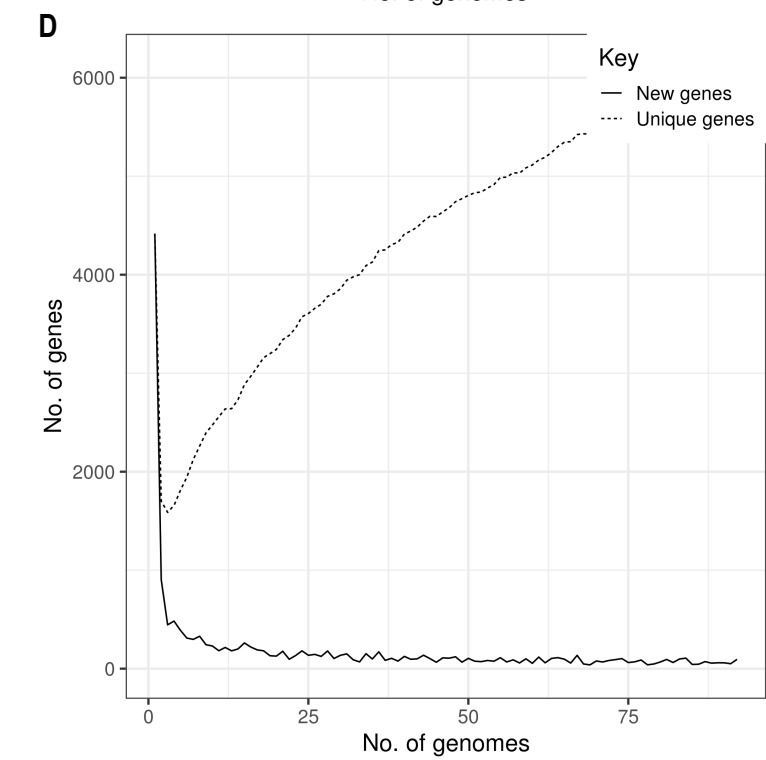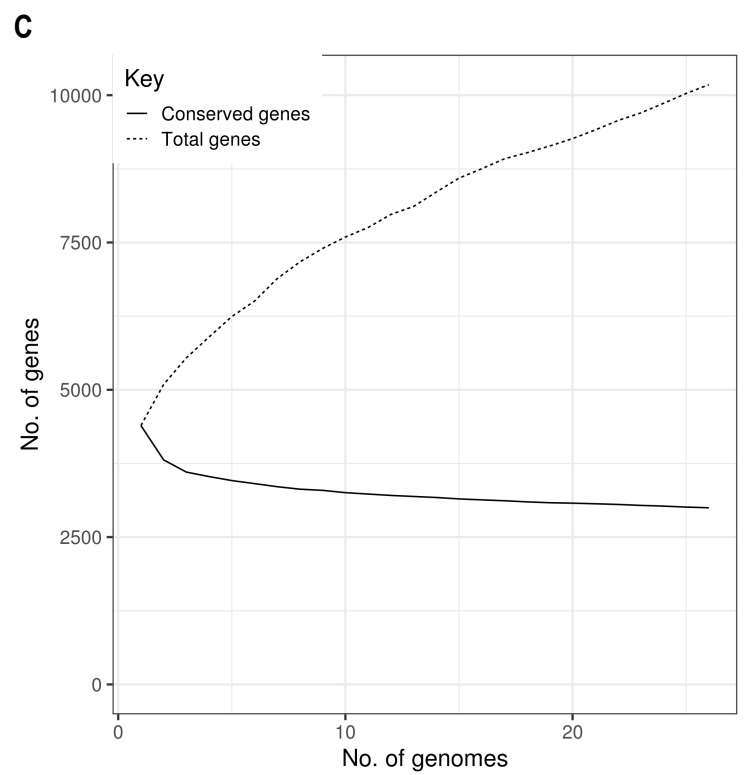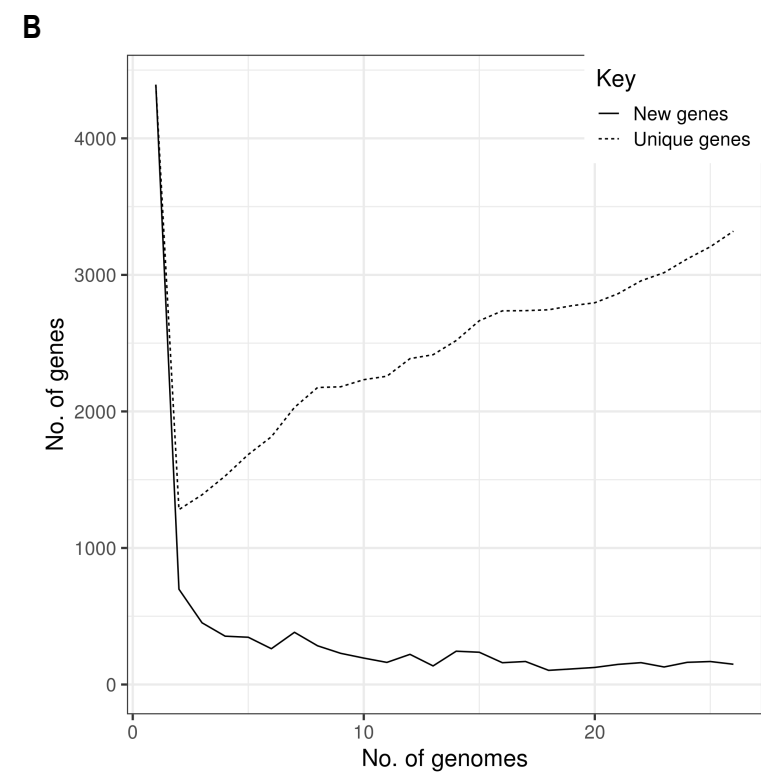
Figure S2

**Figure S2. Roary pangenome analysis of pipolins-harboring *E. coli* strains.**
A. Combined layout of hierarchical clustering of gene presence/absence for LREC *E. coli* genomes (left), along with color-coded markers (middle) and chromosome genetic structure (right). Representation of Roary output was rendered at Phandango website [67]. Lower panels show the accumulation of new vs. unique genes per isolate (B and D) and conserved vs. total genes per isolate (C and E) in the new (LREC, B and C) and all (LREC+GenBank, D and E) *E. coli* strains.
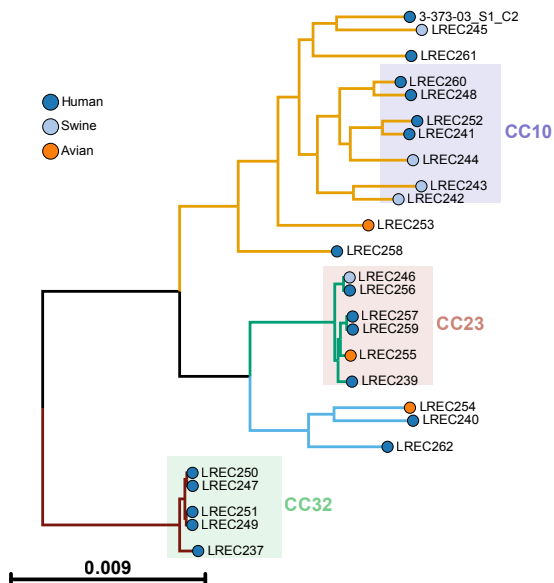
Figure S3

**Figure S3. Maximum-likelihood phylogeny of pipolins-harboring *E. coli* strains.** Codon aware core-genome alignment from Roary was used for the phylogeny reconstruction using IQtree. The best-fit model was GTR+F+R7 (according to ModelFinder). Scale bar indicates the substitution rate per site. The main features retrieved from Enterobase are indicated on the right: source, multilocus sequence type (ST) and clonal complexes (CC). Strains are colored based on the phylogenetic groups, with LREC strains highlighted in italics. Reference strain 3-373-03_S1_C2 is highlighted with a black asterisk (*).
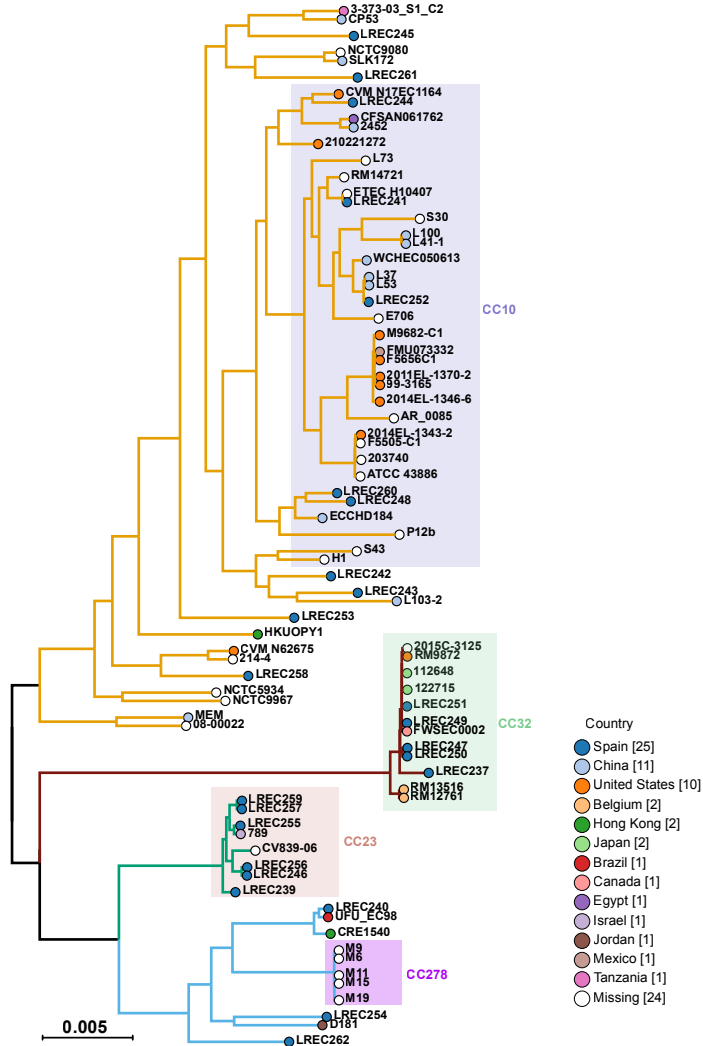
Figure S4

**Figure S4. Single Nucleotide Polymorphism (SNP) trees of the pipolin-harboring strains.**
A. LREC strains' SNP-based tree. Previously described pipolin-harboring isolate 3-373-03_S1_C2 was included as the reference genome. The SNP-tree was performed with EnteroBase [21] default parameters (min. 95% sites present) and included 132985 variant sites. Source and main clonal complexes are also indicated. B. SNP-based tree of the 58 pipolin-harboring strains available from Enterobase and the 25 LREC pipolin-harboring strains performed with EnteroBase. The tree included 170702 variant sites. The figure also includes the country of isolation, sequence types (ST) and clonal complexes (CC). Clade branches were colored by phylogenetic groups as in previous figures.
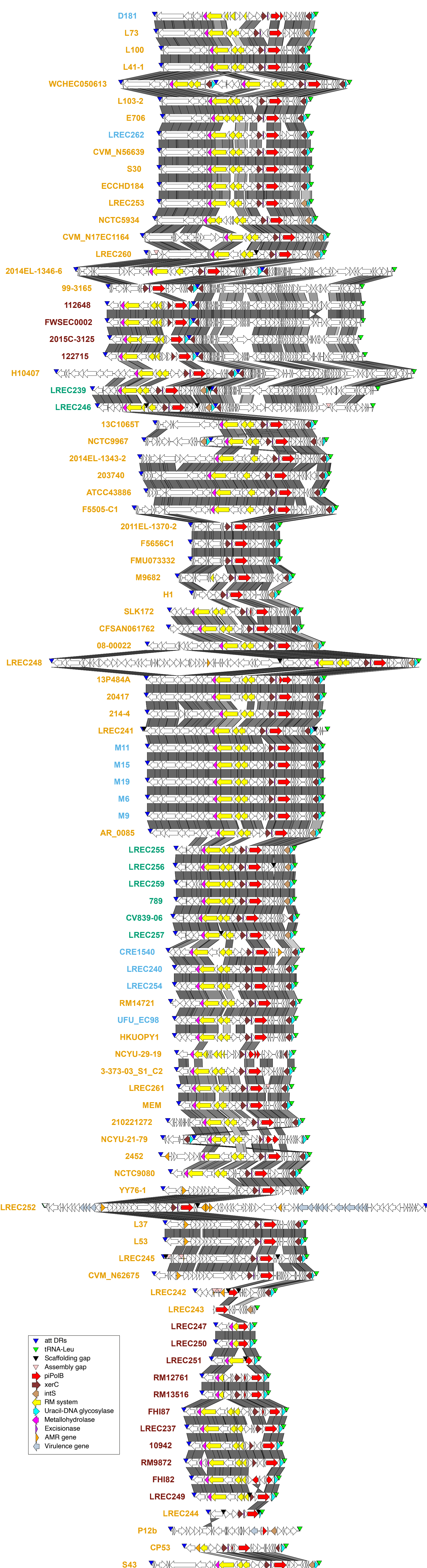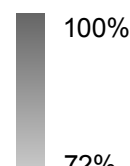
Figure S5

**Figure S5. Genetic structure of pipolins.**
Protein-coding genes are represented by arrows, indicating the direction of transcription, and colored as indicated in the legend. The image was generated by EasyFig software and re-annotated pipolins sorted according to the hierarchical clustering of the gene presence/absence. The greyscale on the right reflects the percent of amino acid identity between pairs of sequences. Names of pipolin-carrying strains are colored based in the phylogroups as in Figure S3.
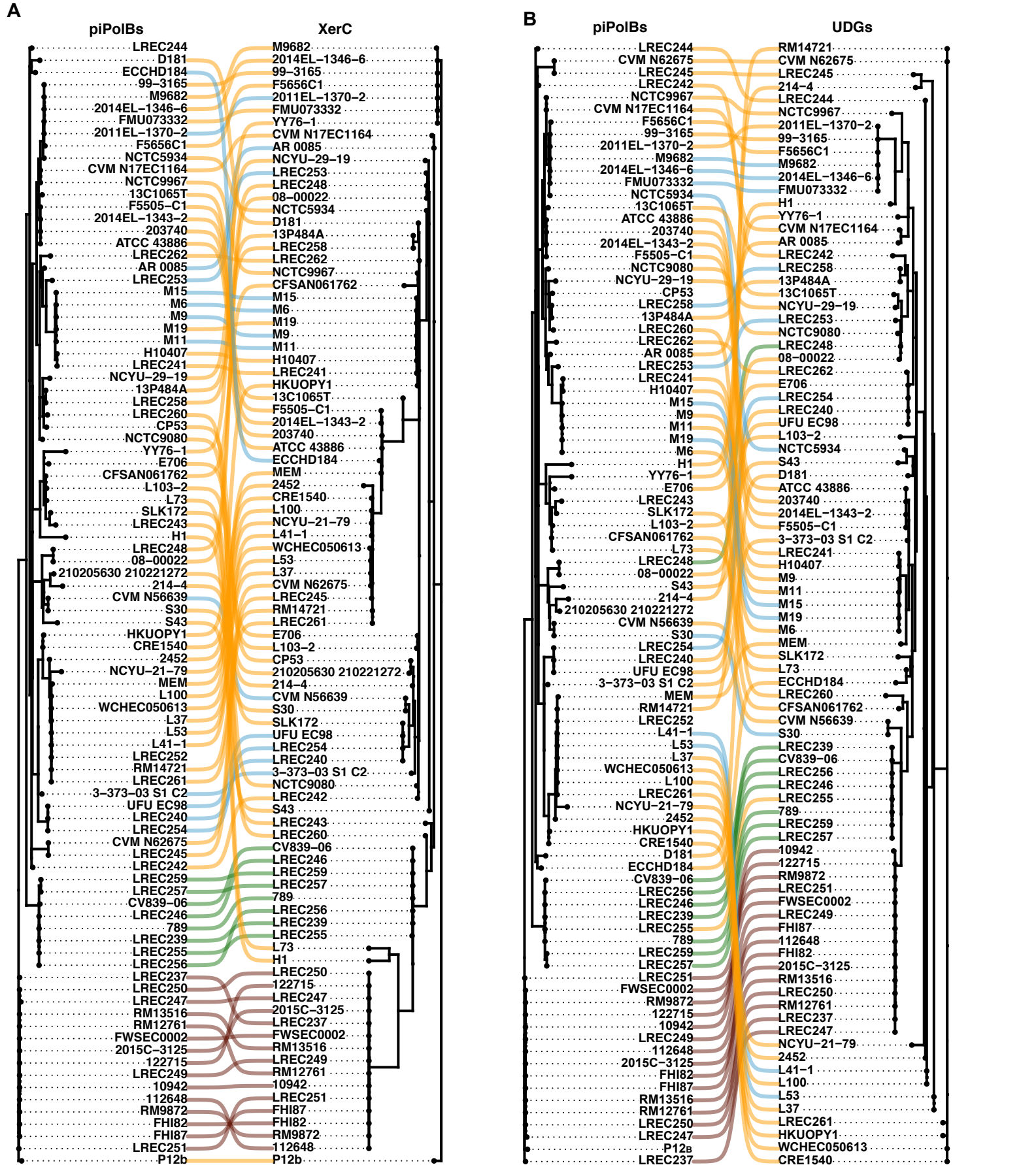
Figure S6

**Figure S6. Tanglegram of pipolins and host strains.**
Tanglegram representation of maximum-likelihood comparative phylogenies of piPolB gene (same as in Figure S1) with XerC (A) and UDG (B) genes of pipolins. Association lines are colored based on the phylogenetic groups as in Figure S3.

Figure S7

**Figure S7. Detailed cophylogeny analysis of pipolins and host strains with PACo.**
Squared Procrustes residues of each phylogenetic group were compared with the
remainder interactions (A, B, C, D). The p-value of Welch t-test is indicated. Panel E
shows the contribution of each pipolin-host association to the general pattern of
coevolution. Each bar represents a jack-knifed estimate of a squared residual. Error
bars represent the upper 95% confidence intervals from applying PACo to patristic
distances. The dashed line indicates the median squared residual value that can serve
as a threshold for congruent phylogenetic interactions. Bars are colored based on the
phylogenetic groups of the strains.