

Supplementary material

**Report and comparative genomics of NDM-5-producing *Escherichia coli* in a Portuguese hospital:
complex class 1 integrons as important players in *bla*_{NDM} spread**

Rafael D. S. Tavares ^{1,2}, Marta Tação ^{2,*}, Elmano Ramalheira ³, Sónia Ferreira ^{3,4,5}, Isabel Henriques ¹

¹ University of Coimbra, Centre for Functional Ecology (CFE), Department of Life Sciences, Calçada
Martim de Freitas, 3000-456 Coimbra, Portugal

² Centre for Environmental and Marine Studies (CESAM) and Department of Biology, University of
Aveiro, Campus Universitário Santiago, 3810-193 Aveiro, Portugal

³ Serviço de Patologia Clínica, Centro Hospitalar do Baixo Vouga-EPE, Avenida Artur Ravara, 3810-501
Aveiro, Portugal

⁴ Department of Medical Sciences, University of Aveiro, Campus Universitário Santiago, 3810-193
Aveiro, Portugal

⁵ Instituto de Educação e Cidadania, Largo da Igreja, 3770-033 Mamarrosa, Aveiro, Portugal

* Correspondence: martat@ua.pt;

Table S1. Primers and thermocycling conditions used in this study.

Primers	Target	Primer sequence (5' - 3')	Thermocycling program ¹	Reference	Positive control
NDM_fwd NDM_rev	<i>bla</i> _{NDM}	GGT TTG GCG ATC TGG TTT TC CGG AAT GGC TCA TCA CGA TC	10 min at 94 °C; (30 s at 94 °C, 40 s at 52 °C, 50 s at 72 °C) for 36 cycles; 7 min at 72 °C	[41]	<i>Enterobacter</i> CR8 ²
BOXA1R	BOX elements	CTA CGG CAA GGC GAC GCT GAC G	7 min at 94 °C; (1 min at 94 °C, 1 min at 53 °C, 8 min at 65 °C) for 30 cycles; 16 min at 65 °C	[42]	-
TEM_fwd TEM_rev	<i>bla</i> _{TEM}	CAT TTC CGT GTC GCC CTT ATT C CGT TCA TCC ATA GTT GCC TGA C	5 min at 94 °C; (40 s at 94 °C, 40 s at 60 °C, 1 min at 72 °C) for 30 cycles; 7 min at 72 °C	[43]	<i>E. coli</i> Ec355340 ³
tetB_fwd tetB_rev	<i>tetB</i>	TCA TTG CCG ATA CCA CCT CAG CCA ACC ATC ATG CTA TTC CAT CC	5 min at 94 °C; (30 s at 94 °C, 30 s at 53 °C, 30 s at 72 °C) for 30 cycles; 7 min at 72 °C	[44]	<i>E. coli</i> Ec355340 ³
intl1_fwd intl1_rev	<i>intl1</i>	CCT CCC GCA CGA TGA TC TCC ACG CAT CGT CAG GC	5 min at 94 °C; (30 s at 94 °C, 30 s at 55 °C, 30 s at 72 °C) for 30 cycles; 7 min at 72 °C	[45]	<i>E. coli</i> Ec355340 ³
R1_fwd ⁴ R1_rev ⁴	pEc355340_NDM-5	GCA CTG TTG CAA ATA GTC GGT GTC GGT AAC CTC GCG CAT A	5 min at 94 °C; (30 s at 94 °C, 30 s at 60 °C, 90 s at 72 °C) for 30 cycles; 7 min at 72 °C	This study	-
R2_fwd ⁴ R2_rev ⁴	pEc355340_NDM-5	TTT CGC GTC AGG GAT GGA AG TTT TCT GAA CCA GGT CGC CA	5 min at 94 °C; (30 s at 94 °C, 30 s at 60 °C, 90 s at 72 °C) for 30 cycles; 7 min at 72 °C	This study	-
R3_fwd ⁴ R3_rev ⁴	pEc355340_NDM-5	GGA AAA AGA GCG CGC TGA AA CAC ATA CCA GAA GCC GTC GT	5 min at 94 °C; (30 s at 94 °C, 30 s at 60 °C, 90 s at 72 °C) for 30 cycles; 7 min at 72 °C	This study	-

¹PCR reactions with a 25 µL-volume were prepared in autoclaved milli-Q water with a concentration of 0.05 U/µL of NZYtaq 2x Green Master Mix (NzyTech, Portugal) and 0.3 µM of each primer. The only exception was in BOX-PCR reaction where the concentration of primer BOXA1R was 0.4 µM.

²Bacterial isolate characterized by Teixeira *et al.* [14].

³Bacterial isolate characterized in this study. Used in this PCR assays as control since whole-genome sequencing data confirmed the presence of these genes and a comparative strategy to the original strain was used.

⁴Primers used for the assembly of pEc355340_NDM-5.

Table S2. Quality metrics of the sequenced genomes obtained in this study as determined by RAST tool (<https://rast.nmpdr.org/rast.cgi>).

Quality metrics	Ec355340	Ec355340ΔpNDM-5
Putative genome size (bp)	5,004,892	4,884,878
Number of contigs (>500 bp)	64	57
GC content (%)	50.9	50.9
Shortest contig size (bp)	597	597
Median sequence size (bp)	31,036	34,400
Mean sequence size (bp)	78,201.4	85,699.6
Longest contig size (bp)	491,699	491,699
N50	215,417	208,412
L50	8	9
Predicted coding sequences	4,867	4,724

Table S3. Virulence genes detected in strain Ec355340 by VirulenceFinder and/or VFAnalyzer. The number of genes identified per category is indicated within square brackets.

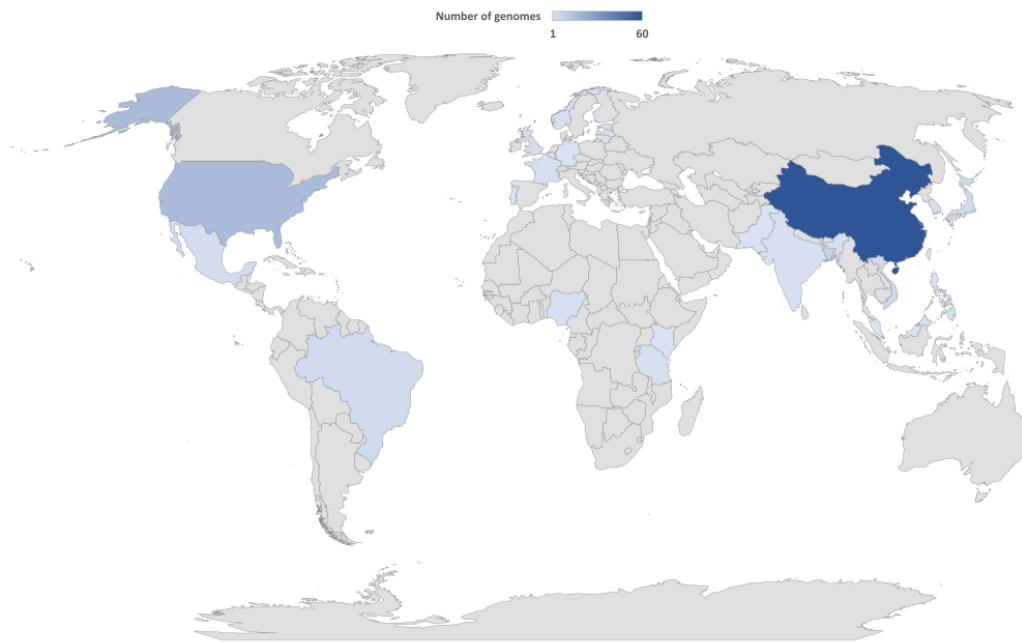
VF class	Virulence factors	Genes identified
Adherence [29]	CFA/I fimbriae [3]	<i>cfaA</i>
		<i>cfaB</i>
		<i>cfaC</i>
	<i>E. coli</i> common pilus (ECP) [6]	<i>ecpA</i>
		<i>ecpB</i>
		<i>ecpC</i>
		<i>ecpD</i>
		<i>ecpE</i>
		<i>ecpR</i>
	<i>E.coli</i> laminin-binding fimbriae (ELF) [4]	<i>elfA</i>
		<i>elfC</i>
		<i>elfD</i>
		<i>elfG</i>
	EaeH [1]	<i>eaeH</i>
	Hemorrhagic <i>E.coli</i> pilus (HCP) [3]	<i>hcpA</i>
		<i>hcpB</i>
		<i>hcpC</i>
	P fimbriae [3]	<i>papC</i>
		<i>papD</i>
		<i>papH</i>
	Type I fimbriae [8]	<i>fimA</i>
		<i>fimC</i>
		<i>fimD</i>
		<i>fimE</i>
		<i>fimF</i>
		<i>fimG</i>
		<i>fimH</i>
		<i>fimI</i>
	Long polar fimbriae [1]	<i>ipfA</i>
Autotransporter [2]	EhaB [1]	<i>ehaB</i>
	UpaG adhesin [1]	<i>upaG/ehaG</i>
Invasion [2]	Invasion of brain endothelial cells (Ibes) [2]	<i>ibeB</i>
		<i>ibeC</i>
Iron uptake [11]	Yersiniabactin siderophore [11]	<i>fyuA</i>
		<i>irp1</i>
		<i>irp2</i>
		<i>ybtA</i>
		<i>ybtE</i>
		<i>ybtP</i>

		<i>ybtQ</i>
		<i>ybtS</i>
		<i>ybtT</i>
		<i>ybtU</i>
		<i>ybtX</i>
Non-LEE encoded TTSS effectors [6]	EspL1 [1]	<i>espL1</i>
	EspL4 [1]	<i>espL4</i>
	EspR1 [1]	<i>espR1</i>
	EspX1 [1]	<i>espX1</i>
	EspX4 [1]	<i>espX4</i>
	EspX5 [1]	<i>espX5</i>
Secretion system [17]	ACE T6SS [17]	-
		<i>aec15</i>
		<i>aec16</i>
		<i>aec17</i>
		<i>aec18</i>
		<i>aec19</i>
		<i>aec22</i>
		<i>aec23</i>
		<i>aec24</i>
		<i>aec25</i>
		<i>aec26</i>
		<i>aec27/clpV</i>
		<i>aec28</i>
		<i>aec29</i>
		<i>aec30</i>
		<i>aec31</i>
		<i>aec32</i>
Toxins [1]	Hemolysin/cytolysin A [1]	<i>hlyE/clyA</i>
Serum resistance [2]	LPS rfb locus (<i>Klebsiella</i>) [1]	-
	Increased serum survival (Iss) protein [1]	<i>iss</i>
Immune evasion [1]	Outer membrane protein complement resistance [1]	<i>traT</i>
Others [2]	Glutamate decarboxylase (acid resistance) [1]	<i>gad</i>
	Tellurium ion resistance protein [1]	<i>terC</i>

Table S4. TOP100 most similar sequences to the plasmid pEc355340_NDM-5. The overall comparison was performed by a discontinuous megaBLAST using the pEc355340_NDM-5 full sequence as query. The existence of similar complex class 1 integrons associated with *bla*_{NDM} was assessed in this dataset by megaBLAST using the sequence of the complex integron identified in pEc355340_NDM-5 (14 kb). Integron arrays were manually annotated, and plasmid replicon types of relevant sequences were identified using PlasmidFinder 2.1. Metadata was also extracted from GenBank.

Table S5. *E. coli* genomes of the sequence type 156 retrieved from PATRIC database (<https://www.patricbrc.org/>). Metadata was obtained from PATRIC, the original GenBank records or from the reporting publication. Identification of antibiotic resistance genes, virulence genes, plasmid replicons and serotype were performed by tools at the Center for Genomic Epidemiology (<https://cge.cbs.dtu.dk/services/>). PathogenFinder was used to predict the potential pathogenicity of each genome. The *bla*_{NDM} genetic context was determined by manual inspection of the contig where the gene was identified in PATRIC using the tool “Genome Browser”. Class 1 integrons were screened by BLAST searches of *intI1* gene (accession NC_022375.1) against our genome dataset (<https://www.patricbrc.org/app/BLAST>). Variable regions were identified, when possible, by manual inspection of the *intI1*-carrying contigs.

A | *E. coli* ST156



B | NDM-5 positive *E. coli*

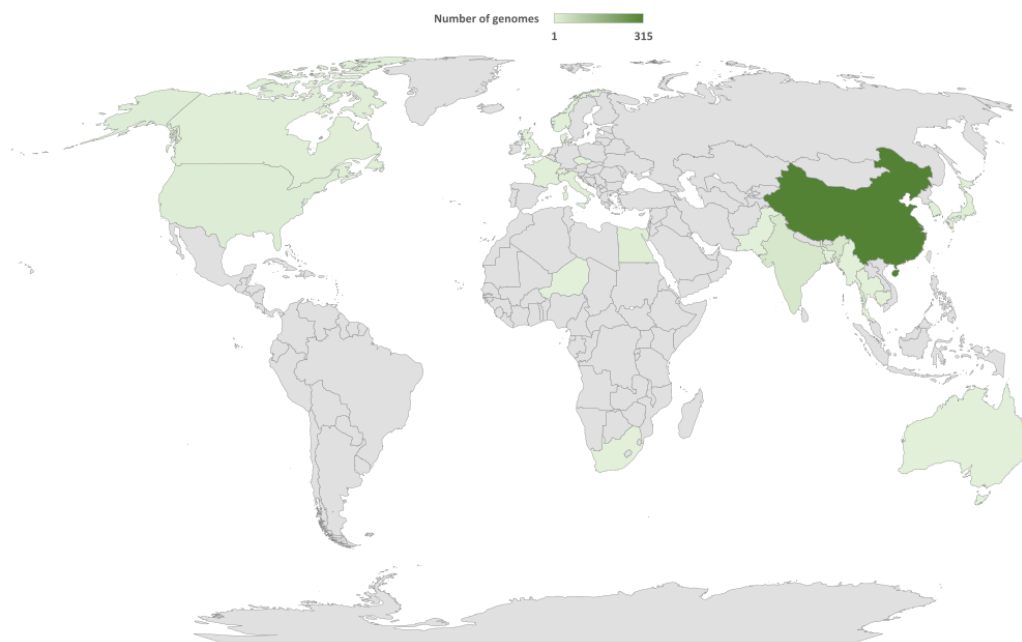


Figure S1. Geographical distribution of (A) *E. coli* ST156 genomes (n= 138) and (B) NDM-5 positive *E. coli* genomes (n=431) deposited in PATRIC database (<https://www.patricbrc.org/>) and analysed in this study. For 9 and 28 genomes of each respective dataset, information regarding country of isolation was not available.

Figure S2. Core SNP tree built from the *E. coli* ST156 genomes available in PATRIC, using REALPHY 1.13 with default settings (<https://realphy.unibas.ch/realphy/>). The genome of *E. coli* AR_452 was used as reference. Metadata related to NDM presence, isolation and geographical source of each genome is represented. The number of ARGs per genome and class of antibiotics to which it confers resistance is identified in a stack barplot while the number of VGs and plasmid replicon detected are displayed in simple barplots. The blue square indicates the cluster that includes Ec355340 and its closest genomes. This figure was built using iTOL version 6 (<https://itol.embl.de/>).

Table S6. Average nucleotide identity (ANI) matrix for the *E. coli* ST156 genomes closest to Ec355340.

Values are presented in percentage and in brackets are indicated the percentage of aligned nucleotides. Pairwise ANI calculations were performed at JSpeciesWS (Ribocon, <https://jspecies.ribohost.com/jspeciesws/>)

	Ec355340	Ec AR_452	Ec 89PenNDM	Ec 112
Ec355340	-	99.98 (99.03)	99.80 (92.72)	99.76 (91.63)
Ec AR_452	-	-	99.81 (94.13)	99.76 (93.44)
Ec 89PenNDM	-	-	-	99.63 (90.87)
Ec 112	-	-	-	-

Figure S3. Hierarchical clustering of *E. coli* ST156 genomes based on the presence/absence of ARGs, VGs and plasmid. The resemblance matrix was built using the Jaccard coefficient and the dendrogram was built using UPMGA as the clustering method in Primer v.6. The heatmap, plotted using R's ggplot2, represents presence/absence of all the genes identified in at least one of these genomes. NDM variants identified in each genome are indicated by colored squares behind each label. The 2 main clusters identified from this analysis are indicated as cluster I and cluster II.

Table S7. NDM-5 positive *E. coli* genomes available in PATRIC and corresponding metadata. Genomes were selected based on “PATRIC Local Family” feature. Specificity to NDM-5 variant was ensured by BLAST analysis (only perfect hits were considered). The sequence type affiliation for each genome was obtained using the MLST 2.0 tool (<https://cge.cbs.dtu.dk/services/MLST/>).