

 Open access • Posted Content • DOI:10.21203/RS.3.RS-199409/V1

Role of the mobilome in the global dissemination of the carbapenem resistance gene blaNDM — Source link

Mislav Acman, Ruobing Wang, Lucy van Dorp, Liam P. Shaw ...+7 more authors

Institutions: University College London, Peking University, John Radcliffe Hospital, University of Warwick

Published on: 17 Feb 2021 - bioRxiv (Cold Spring Harbor Laboratory)

Topics: Mobilome

Related papers:

- [Genomic characterisation and context of the blaNDM-1 carbapenemase in Escherichia coli ST101.](#)
- [Population structure and pangenome analysis of Enterobacter bugandensis uncover the presence of blaCTX-M-55, blaNDM-5 and blaIMI-1, along with sophisticated iron acquisition strategies.](#)
- [Integrated chromosomal and plasmid sequence analyses reveal diverse modes of carbapenemase gene spread among Klebsiella pneumoniae.](#)
- [Genomic analysis of carbapenemase-encoding plasmids from Klebsiella pneumoniae across Europe highlights three major patterns of dissemination](#)
- [Genomic evolution of the globally disseminated multidrug-resistant Klebsiella pneumoniae clonal group 147](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/role-of-the-mobilome-in-the-global-dissemination-of-the-3p80eakxgp>

Role of the mobilome in the global dissemination of the carbapenem resistance gene blaNDM

Mislav Acman (✉ mislav.acman@gmail.com)

University College London <https://orcid.org/0000-0003-2587-2836>

Ruobing Wang

Department of Clinical Laboratory, Peking University People's Hospital

Lucy van Dorp

University College London <https://orcid.org/0000-0002-6211-2310>

Liam Shaw

University of Oxford <https://orcid.org/0000-0001-7332-0820>

Qi Wang

Department of Clinical Laboratory, Peking University People's Hospital

Nina Luhmann

Warwick Medical School, University of Warwick

Yuyao Yin

Peking University People's Hospital

Shijun Sun

Peking University People's Hospital

Hongbin Chen

Department of Clinical Laboratory, Peking University People's Hospital

Hui Wang

Peking University People's Hospital

Francois Balloux

University College London <https://orcid.org/0000-0003-1978-7715>

Article

Keywords: blaNDM, resistance gene, host-adaptation

Posted Date: February 17th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-199409/v1>

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

[Read Full License](#)

Version of Record: A version of this preprint was published at Nature Communications on March 3rd, 2022. See the published version at <https://doi.org/10.1038/s41467-022-28819-2>.

1 Role of the mobilome in the global dissemination of the
2 carbapenem resistance gene *bla*_{NDM}

3 Mislav Acman^{1*}, Ruobing Wang², Lucy van Dorp¹, Liam P. Shaw³, Qi Wang², Nina Luhmann⁴, Yuyao Yin²,
4 Shijun Sun², Hongbin Chen², Hui Wang², Francois Balloux¹

5 1 UCL Genetics Institute, University College London, Gower Street, London, WC1E 6BT, UK

6 2 Department of Clinical Laboratory, Peking University People's Hospital, Beijing, 100044, China

7 3 Nuffield Department of Medicine, John Radcliffe Hospital, University of Oxford, Oxford OX3 9DU, UK

8 4 Warwick Medical School, University of Warwick, Coventry CV4 7AL, UK

9 * Corresponding Author

10 E-mail: mislav.acman.17@ucl.ac.uk

11 Abstract (249 words)

12 The mobile resistance gene *bla_{NDM}* encodes the NDM enzyme capable of hydrolysing carbapenems, a class of
13 antibiotics used to treat some of the most severe bacterial infections. *bla_{NDM}* is globally distributed across a variety
14 of Gram-negative bacteria and is typically located within a highly recombining transposon-rich genomic region
15 common to multiple plasmids types. As a result of this genomic complexity the dynamics underlying the
16 dissemination of *bla_{NDM}* remain poorly resolved. In this work, we compiled a dataset of over 2000 bacterial
17 genomes harbouring the *bla_{NDM}* gene including 112 new PacBio hybrid assemblies from clinical and livestock
18 associated isolates across China and developed a novel computational approach to track structural variants in
19 bacterial genomes. We were able to correlate specific structural variants with plasmid backbones, bacterial host
20 species and sampling locations, and identified multiple transposition events that occurred during the global
21 dissemination of *bla_{NDM}*. Our results highlight the most prominent transposons responsible for the global spread
22 of *bla_{NDM}* and suggest that genetic recombination, rather than mutation, was the dominant force driving the
23 evolution of the *bla_{NDM}* genomic region. By tracking the change in diversity among countries of collection of
24 *bla_{NDM}*-positive genomes, we estimate that *bla_{NDM}* reached global prevalence within 8-11 years after its initial
25 mobilization. Lastly, we observe notable correlation between plasmid backbones bearing *bla_{NDM}* and the sampling
26 location which suggests that the dissemination of resistance is mainly driven by successive between-plasmid
27 transposon jumps with plasmid exchange being largely restricted by the boundaries defined by bacterial host-
28 adaptation of individual plasmids.

29 Introduction

30 The increasing burden of antimicrobial resistance (AMR) poses a major challenge to human and veterinary health.
31 AMR can be conferred by vertically inherited point mutations or via the acquisition of horizontally transmitted
32 non-essential ‘accessory’ genes generally located in transposons and plasmids. The *bla_{NDM}* gene encoding the
33 NDM enzyme, a metallo- β -lactamase capable of hydrolysing most β -lactam antibiotics represents a typical
34 example of a mobile antibiotic resistance gene¹. Compounds belonging to the carbapenem class are commonly
35 employed to treat Gram-negative bacterial infections resistant to mainstay antibiotics and used as first-line
36 treatment for severe infections. The global prevalence of bacteria carrying *bla_{NDM}*, including carbapenem-resistant
37 *Acinetobacter baumannii* and *Enterobacteriaceae* in hospital settings, represents a major public health concern.

38 The *bla_{NDM}* gene was first identified in 2008 from a *Klebsiella pneumoniae* isolated from a urinary tract infection
39 in a Swedish patient returning from New Delhi, India². While *bla_{NDM}* now has a worldwide distribution, most of
40 the earliest cases have been linked to the Indian subcontinent, suggesting this region as a likely location for the
41 initial mobilisation event^{1,3-6}. Notably, NDM-positive *Acinetobacter baumannii* isolates have been retrospectively
42 identified from an Indian hospital in 2005⁷, which remain the earliest observations to date. However, an NDM-
43 positive *A. pittii* isolate was also isolated in 2006 from a Turkish patient with no history of travel outside Turkey⁸.

44 Although no complete genome sequences are publicly available from these earliest observations, the first NDM-
45 positive isolates from 2005 were shown to carry *bla_{NDM}* on multiple non-conjugative, but potentially mobilizable
46 plasmid backbones⁷. In addition, *bla_{NDM}* in these early isolates was positioned within a complete Tn125 transposon
47 with existing ISCR27 and IS26 insertion sequences (ISs), suggesting the possibility of complex patterns of
48 mobility since the gene’s initial integration. Subsequent NDM-positive isolates, spanning a range of species,
49 consistently harbour either a complete or fragmented IS*Aba125* (an IS constituting Tn125), which is always found
50 immediately upstream of *bla_{NDM}* providing a promoter region for the gene transcription^{1,5,9,10}. The presence of
51 IS*Aba125*, in some form, in all NDM-positive isolates to date, as well as the majority of the early observations
52 being in *A. baumannii*, has led to Tn125 being proposed as the likely transposon responsible for the initial
53 mobilization of *bla_{NDM}*, and *A. baumannii* as the ancestral host.

54 In addition, the NDM enzyme itself has been described as of possible chimeric origin^{10,11}, with the first six amino
55 acids in NDM matching to those in *aphA6*, a gene providing aminoglycoside resistance and also flanked by
56 IS*Aba125*. It is presumed that ISCR27, an IS which uses a rolling-circle (RC) transposition mechanism^{12,13},
57 initially mobilized a progenitor of *bla_{NDM}* in *Xanthomonas sp.* and placed it downstream of IS*Aba12*^{10,11,14,15}. The
58 NDM enzyme itself displays some polymorphism, with at least 29 distinct sequence variants having been
59 described to date. The most prevalent of these variants is the first to have been characterised, and is denoted NDM-
60 1¹⁶. Different NDM variants are mostly distinguished by a single amino-acid substitution, with the exception of
61 NDM-18 which carries a tandem repeat of five amino acids. None of the observed substitutions occur in the active
62 site and the functional impact of each of these substitutions remains under debate¹.

63 At present, NDM resistance is globally distributed and represents a major concern in healthcare settings. The gene
64 is found in at least 11 bacterial families and NDM-positive isolates have heterogeneous clonal backgrounds,
65 supporting multiple independent acquisitions of *bla_{NDM}*¹. The *bla_{NDM}* gene has been observed on bacterial

66 chromosomes^{17,18} but is most commonly harboured on plasmids, comprising multiple different backbones or
67 types. Furthermore, even within the same plasmid types, *bla_{NDM}* is found in a variety of genetic contexts, often
68 interspersed by multiple ISs and composite transposons^{1,11}. The immediate genetic environment of *bla_{NDM}* has
69 been reported to vary even in isolates from the same patient¹⁹. It is therefore clear that the emergence and
70 subsequent dissemination of NDM resistance, through a multitude of bacterial host species, is a dynamic and
71 multi-layer process involving multiple mobile genetic elements – ‘the mobilome’ – which abetted the mobility of
72 *bla_{NDM}* via a diverse set of processes, including genetic recombination, transposition, conjugation, transformation,
73 and transfer through outer-membrane vesicles (OMVs)²⁰⁻²³.

74 In this work, we reconstruct the individual roles of plasmids and ISs in the dissemination of NDM and provide a
75 comprehensive overview of the many genetic backgrounds harbouring the *bla_{NDM}* gene. To this end, we compiled
76 a global dataset of more than 2000 NDM-positive isolates including 112 newly generated hybrid PacBio
77 assemblies sampled from clinical and livestock settings across China. In order to decompose the high sequence
78 complexity of the immediate genomic contexts of *bla_{NDM}* in our large global dataset, we developed a novel
79 alignment-based method designed to uncover all structural variations flanking *bla_{NDM}*. This allowed us to pinpoint
80 individual insertion events for subsequent assessment. Correlating specific structural variants with plasmid
81 backbones, bacterial host genera and sampling locations, we are able to uncover transposition events underlying
82 the global spread of *bla_{NDM}*. We identify Tn125, Tn3000 and IS26 as the main contributors to *bla_{NDM}* mobility.
83 Furthermore, we provide evidence for genetic recombination being the main force driving evolution in this region.
84 We also identify plasmid backbones and bacterial hosts closely associated with specific sampling locations, as
85 well as an apparent plateau in the rate of spread of *bla_{NDM}* around 2014. Our findings position plasmids as the
86 main contributors to the local transmission of *bla_{NDM}*, while transposons seem to be more influential for spread at
87 a global scale.

88 Results

89

90 A global dataset of *bla*_{NDM} carriers

91 To study the genetic context and global spread of the *bla*_{NDM} resistance gene, a dataset of 2,148 bacterial genomes
92 (2,166 contigs) carrying at least one copy of *bla*_{NDM} were compiled from multiple sources (Figure 1). These
93 include: 795 bacterial genomes assembled using short read *de novo* assembly methods; 113 bacterial genomes
94 using hybrid PacBio-Illumina *de novo* assembly; and 1,240 RefSeq assemblies (See Methods, Supplementary
95 Table 1). Of the included *de novo* hybrid assemblies, 112 were newly generated for this study isolated from 87
96 hospitalized patients across China and 25 livestock farms. Overall, the dataset includes NDM-positive genomes
97 sampled across 67 states (Figure 1A). The majority of isolates were collected in East and South East Asian
98 countries with mainland China representing the predominant source of origin ($n=668$). A wide range of bacterial
99 species were represented with *Klebsiella* and *Escherichia* the primary genera each contributing 899 and 667
100 genomes, respectively (Figure 1B; Supplementary Data 1).

101 The majority of *bla*_{NDM} carriers in the global dataset were collected between 2014-2017 (74.41%, Figure 1C).
102 However, the dataset also includes 31 genomes from 2010 and earlier. These include the *K. pneumoniae* isolate
103 from 2008 in which *bla*_{NDM} was first characterized², as well as an earlier *A. baumannii* isolate from 2007 in an
104 individual of Balkan origin in Germany^{24,25} (Supplementary Data 1).

105 A substantial number of contigs isolated from our dataset were sufficient in length to enable identification of
106 putative plasmid backbones carrying *bla*_{NDM} (Supplementary Figure 1; See Methods). Within our filtered dataset
107 comprising 2,142 contigs (see Methods), we identified 482 replicon types using PlasmidFinder²⁶ and 194
108 circularized contigs in our dataset, of which 43 did not have a known replicon type. This resulted in a total of 525
109 putative plasmid sequences which also comprised 96 contigs (70 circularized) from our hybrid PacBio-Illumina
110 assemblies. Overall, 32 different plasmid replicon types were identified among *bla*_{NDM}-containing plasmid
111 sequences (Figure 1D). The most prevalent replicon type was IncX3, found in almost half (253/525, 48%) of the
112 included sequences. Nevertheless, the notable range in plasmid backbones harbouring *bla*_{NDM} indicates a high
113 recombination and/or transposition rate of the *bla*_{NDM} gene. At the same time, we observe some geographic
114 structure in plasmid replicon types (Supplementary Figure 2) signalling the importance of transposon movement
115 in the cross-continental spread of NDM-mediated resistance.

116

117 Resolving structural variants in the *bla*_{NDM} flanking regions

118 To gain a detailed overview of the transposition events and different genetic backgrounds harbouring *bla*_{NDM} we
119 developed an alignment-based approach to resolve structural variation in the genetic regions flanking *bla*_{NDM} (see
120 Methods, Figure 2). In brief, a pairwise discontinuous Mega BLAST search (v2.10.1+)^{27,28} was applied to all
121 *bla*_{NDM}-containing contigs in order to identify all possible homologous regions between each contig pair. Only
122 BLAST hits covering the complete *bla*_{NDM} gene were retained (Figure 2A). Next, starting from *bla*_{NDM}, a gradually
123 increasing ‘splitting threshold’ was introduced to monitor structural variants as they appeared upstream or
124 downstream of the gene. At each step, a network is constructed connecting contigs (nodes) that share a BLAST

125 hit with a minimum length as given by the ‘splitting threshold’ (Figure 2B). As we move upstream or downstream
126 and further away from the gene, the network starts to split into smaller clusters each carrying contigs that share
127 an uninterrupted stretch of homologous DNA. The splitting is visualized as a tree where branch lengths are scaled
128 to match the position within the sequence, and the thickness and the colour intensity of the branches corresponds
129 to the number of sequences which are homologous (Figure 2C). Given the approach uses the *bla_{NDM}* gene as an
130 anchor, it enables comparison between BLAST hits, but also limits the comparison to either upstream or
131 downstream flanking region and not both simultaneously.

132 The flanking region upstream of *bla_{NDM}* breaks down rather quickly: within a few hundred base pairs of the *bla_{NDM}*
133 start codon, the upstream flanking region splits into multiple structural variants, none of which dominates the
134 contig pool (Supplementary Figure 3). For instance, 99 different structural variants were identified only 1200 bp
135 from the *bla_{NDM}* start codon. This high variation in genome structure could be attributed to the many genetic
136 backgrounds in which *bla_{NDM}* is found as well as frequent genome rearrangements (Supplementary Figures 3).
137 The significance of the latter is also reflected by the number of fragments and complete insertion sequences present
138 in the region, including *ISAbal25* (132), *IS5* (385), *IS3000* (88), *ISKpn14* (44), and *ISEc33* (72), as well as almost
139 half the contigs (1,003, 46.93%) being excluded from the analysis for having too short an upstream flank
140 (Supplementary Figure 3). The transposition hotspot upstream of *bla_{NDM}* possibly hinders sequencing and genome
141 assembly efforts and enhances the presence of these short contig flanks. In agreement with previous work^{1,5,9,10},
142 more than 95% of sufficiently long contigs include a ~75 bp fraction of *ISAbal25*, supporting the notion of *TnI25*
143 as an ancestral transposon of the *bla_{NDM}* gene (Supplementary Figures 3 and 4).

144 The downstream flanking region exhibits more gradual structural diversification than the upstream region, with
145 one dominant putative ancestral background (Figure 3). As illustrated by the stem of the tree of structural
146 variations (Supplementary Figure 5), many of the 2,142 contigs analysed contain complete sequences of the same
147 genes: *ble* (2,047 contigs), *trpF* (1,770), *dsbD* (1,660), *cutA* (858), *groS* (673), *groL* (527). In total there are 1,229
148 contigs which are sufficiently long downstream of *bla_{NDM}* to harbour the full repertoire of the aforementioned
149 genes. When the analysis is restricted to those contigs of sufficient length, 42.9% of NDM-positive contigs carry
150 this full suite of genes downstream of *bla_{NDM}*.

151

152 Patterns of insertion events in *bla_{NDM}* flanking regions

153 Having reconstructed structural variation in the *bla_{NDM}* upstream (Supplementary Figure 3) and downstream
154 (Figure 3) flanking regions, we did not observe any strong overall signal in the distribution of associated plasmid
155 backbones, bacterial genera and sampling locations. However, closer examination of structural variants common
156 to sufficiently large pools of isolates allow distinct observations to be made. These more specific observations
157 appear to correlate to the events underlying the spread of *bla_{NDM}*. For instance, *IS3000* is found in 88 and 35
158 contigs on the upstream and downstream flanking regions respectively, almost exclusively in *Klebsiella* host
159 species and often on IncF plasmids (Figure 3 and Supplementary Figure 3). Thus, as previously suggested by
160 Campos et al., *Tn3000* likely re-mobilized *bla_{NDM}*, following the fossilization of *TnI25*²⁹; our analysis suggests
161 the secondary mobilization primarily happened in *Klebsiella* species.

162 Some structural variants appear geographically linked e.g., IS5 is predominantly found upstream of *bla*_{NDM} on
163 IncX3 plasmids from East Asia (Supplementary Figure 3), with none of these plasmids with IS5 having a matching
164 element on the downstream flanking region of *bla*_{NDM} to form a full composite transposon. IS5 is known to
165 enhance transcription of nearby promoters in *E. coli*³⁰ and its abundance and positioning just upstream of *bla*_{NDM}
166 suggests it may have assumed a similar role in this case. Interestingly, the NDM-5 variant has been increasing in
167 numbers in recent years (Supplementary Figure 6 A and B) and is mostly associated with both IncX3 plasmids
168 (Supplementary Figure 6 C and D) and isolates from East Asia (Supplementary Figure 6 G and H). Thus, an
169 increasing abundance of NDM-5 could be due to the aforementioned enhanced transcription caused by the
170 proximity to IS5. Other structural variants are observed across many global regions e.g., the *wapA* gene is found
171 truncating ISCR27 downstream on IncC plasmids (Figure 3).

172 One of the most commonly found transposable elements in the flanking regions (~30% prevalence) is an ISCR1-
173 like transposase (IS91 family transposase), hereafter referred to as ISCR1, coupled with the *folP* gene (Figure 3,
174 Supplementary Figure 5). This configuration is found at various positions downstream of *bla*_{NDM} and often
175 associated to IncF plasmids identified in *Escherichia* and *Klebsiella* species. In most cases, the orientation of
176 ISCR1 should prevent this element from mobilizing *bla*_{NDM}¹³, so it appears its role is to disrupt the surrounding
177 IS elements and transposons. Interestingly, ISCR1s are mainly found in complex class 1 integrons¹³, however, not
178 many annotated integrase genes are located within the vicinity of *bla*_{NDM}. In fact, only 11 contigs were found to
179 have an integrase <50 Kb away from *bla*_{NDM} and none showed any consistency in how the integrase is placed with
180 respect to *bla*_{NDM}. This may suggest integrons play at most a minor role in the dissemination of *bla*_{NDM}.

181 Another notable ISCR element is ISCR27 which is consistently found immediately downstream of the *groL* gene
182 (Figure 3, Supplementary Figure 5). The complete ISCR27 sequence is carried by 316 contigs, with another 211
183 contigs containing a fragmentary sequence. ISCR27 is found at high prevalence with 30.1% of sufficiently long
184 contigs harbouring this element. Contrary to its ISCR1-like relative, ISCR27 is correctly oriented to mobilize
185 *bla*_{NDM} as is presumed to have happened during the initial mobilization of the progenitor of *bla*_{NDM}¹⁰. However,
186 we find no evidence of subsequent ISCR27 mobility. The origin of rolling-circle replication of ISCR27 (*oriIS*;
187 GCGGTTGAACTTCCTATACC) is located 236 bp downstream of the ISCR27 transposase stop codon. The
188 region downstream of this stop codon in all structural variants bearing a complete ISCR27 is highly conserved for
189 at least 750 bp (Figure 3, Supplementary Figure 5). This suggests a reasonably conserved genetic background
190 surrounding ISCR27 as *bla*_{NDM} has been disseminated.

191 Surprisingly, only 58 contigs carried a complete IS*Aba125* downstream of *bla*_{NDM}, of which 53 carried an
192 IS*Aba125* sequence in proximity (<7886 bp) to the *bla*_{NDM} start codon. These account for a minority (7.4%) of
193 isolates when sufficiently long contigs are considered. Forty-five of these contigs contained a complete IS*Aba125*
194 both upstream and downstream of *bla*_{NDM} thus forming a complete Tn*I25* transposon. Even though the diversity
195 of bacterial genera carrying IS*Aba125* upstream is substantial (Supplementary Figure 3), the less preserved
196 downstream IS*Aba125* sequence is mostly found in the genera *Acinetobacter* and *Klebsiella* (Figure 3). This
197 supports the initial dissemination of *bla*_{NDM} by Tn*I25* to other plasmid backbones predominately being mediated
198 by these two genera, after which the transposon was disrupted by other rearrangements.

199 We note that more than 500 contigs were truncated around 3000 bp downstream of *bla*_{N_{DM}} (Figure 3). To
200 investigate the reasons behind this distinct cut-off point, we used 447 raw short-read sequencing samples from
201 our dataset (originally downloaded from SRA, see Methods) with contigs that carry *bla*_{N_{DM}} longer than 3000bp
202 (Supplementary Table 1). We compared the normalized number of reads with overhangs mapping to the end of
203 contigs ending 3000-3200 bp and longer contigs, ending >3200 bp downstream of *bla*_{N_{DM}} (Supplementary Figure
204 7A). On average, the normalized number of overhangs is two times higher in shorter contigs, which indicates that
205 a particular genetic region mapped by the overhanging reads is often present in more than one copy. Moreover,
206 when mapped back to the assembled contigs, the overhanging reads of shorter contigs are found on average on
207 three different contigs (>1000 bp) – twice as many as observed for longer contigs (Supplementary Figure 7B).
208 The presence of these overhanging reads on multiple contigs may point to within-isolate
209 transposition/rearrangement events between plasmids and/or bacterial genomes which seem to localise around
210 this region.

211 The shorter contigs (3000-3200 bp) are found across genera of *Enterobacteriaceae* including *Escherichia*,
212 *Klebsiella*, *Enterobacter*, *Citrobacter*, *Leclercia* and *Lelliottia*. What is more, the overhanging reads of shorter
213 contigs almost exclusively match the left inverted repeat (IRL) of IS26 sequence. In fact, over one third (157;
214 35.1%) of all analysed contigs' overhanging reads correspond to IS26 IRL. IS26, although often found in two
215 adjacent copies forming a seemingly composite transposon, is a so-called pseudo-composite (or pseudo-
216 compound) transposon³¹. In contrast to composite transposons, a fraction of DNA flanked by the two IS26 is
217 mobilized either via cointegrate formation or in the form of a translocatable unit (TU), which consists of a single
218 IS26 element and a mobilized fraction of DNA, and inserts preferentially next to another IS26^{31,32}. Interestingly,
219 no IS26 sequences were found upstream within contigs whose downstream overhanging reads match to IS26.
220 Assembly procedures are known to struggle with allele duplications which may explain the lack of IS26 sequences
221 upstream and the surge of truncated contigs³³. Nevertheless, the results above suggest an active within-isolate
222 movement of *bla*_{N_{DM}} via IS26 across *Enterobacteriaceae*.

223 In total, we identified 208 putative composite transposons (i.e., stretches of DNA flanked by at least two ISs
224 enclosing *bla*_{N_{DM}} <30 Kb apart) in 181 contigs. These comprised 18 different types with the five most frequent
225 being: IS26 (62 instances), IS*Aba125* (forming Tn*125*; 55 instances), IS3000 (forming Tn3000; 52), IS15 (13),
226 IS6100 (7). Interestingly, there are 38 cases where >2 of the same IS flank *bla*_{N_{DM}}. These are mostly IS26 (23).
227 Also, only 137 of the 208 putative transposons identified contained both complete flanking ISs, while others had
228 at least one IS partially truncated. Importantly, IS26, IS6100 and IS15, a known variant of IS26, are
229 phylogenetically related with all three falling into clade I of the IS6 family of insertion sequences whose members
230 are known to mobilize via cointegrate formation, as discussed above³⁴. The IS26s we identify are found at different
231 positions in the alignment, usually between 10-20 Kb apart, while other ISs are, for the most part, found at a fixed
232 position around *bla*_{N_{DM}}. This indicates increased activity and multiple independent acquisitions of IS26. As
233 expected, the transposons we identify are found on various plasmid backbones (Supplementary Figure 8C).
234 However, some trends can be identified in the distribution of associated bacterial genera and geographic region
235 of sampling (Supplementary Figure 8A and B). In particular, Tn3000 is almost exclusively found in *Klebsiella*
236 species and Tn125 predominantly in *Acinetobacter* and *Klebsiella*, while IS26 are found in *Escherichia* and
237 *Klebsiella*. In spite of these elements being present across the globe, some geographic structure is apparent. For

238 example, IS26 appears to dominate in East Asia while Tn3000 tends to occur in South Asia. Overall, the
239 distributions of various structural variants and transposons with respect to plasmid replicon types and bacterial
240 hosts suggest that most rearrangements in the *bla*_{NDM} flanking regions happened within *Escherichia* and *Klebsiella*
241 species where IS26, Tn125 and Tn3000 are the main contributors to *bla*_{NDM} mobility.

242

243 Mutations accumulated in *bla*_{NDM} transposons provide only weak evolutionary 244 signal

245 To further investigate the dynamics of spread of the *bla*_{NDM} gene, regression analyses and Bayesian molecular tip-
246 dating (implemented in BEAST2 v2.6.0)³⁵ were performed on full alignments of Tn125 (45 contigs) and Tn3000
247 (29 contigs) (Supplementary Figure 9). SNPs within each alignment were identified using a consensus sequence
248 approach (see Methods). Few SNPs are observed in the alignments of Tn125 (56 SNPs) and Tn3000 (14)
249 (Supplementary Figure 9A and B). In fact, a general observation was that relatively few SNPs are found in
250 alignments of any stretch of homologous sequence flanking *bla*_{NDM} relative to the number of structural variants.
251 For instance, only 80 SNPs are present in the 2,570 bp alignment of 1,711 contigs harbouring *bla*_{NDM}, *ble*, *trpF*,
252 and *dsbD* genes, while more than 50 different structural variations are found over the same distance downstream
253 of the *bla*_{NDM} start codon. Going downstream, the number of structural variants increases while the number of
254 newly accumulating SNPs plateaus, as fewer samples are available and the genetic background diversifies.

255 This restricted genetic diversity of the two transposon alignments results in only a weak temporal signal (see
256 Methods and Supplementary Figure 9A and B). While results should therefore be interpreted with appropriate
257 caution, we proceeded with Bayesian molecular tip-dating analyses to assess the relative timing of transposition
258 events involving Tn125 and Tn3000 (see Methods). All models converged well, though we note that both marginal
259 distributions of the most common recent ancestor (tMRCA) of Tn125 and Tn3000 (Supplementary Figure 9C and
260 D) overlap with the marginal distributions of the corresponding model priors (i.e., BEAST2 runs without SNP
261 data provided) (Supplementary Figure 9D) which is a likely consequence of the lack of genetic diversity.
262 Nevertheless, the tMRCA estimates of Tn125 and Tn3000 shift from the expectation under the priors. In particular,
263 the Tn3000 marginal distribution points to a later date indicating that the tMRCA of Tn3000 carrying *bla*_{NDM} gene
264 emerged after mid-2008, but still before the earliest sampling date at the end of 2011 (Supplementary Figure 9C).
265 In contrast, the marginal distribution of the Tn125 tMRCA shifts to an earlier date, suggesting this transposon
266 mobilized *bla*_{NDM} before 2009 and after 2004. This tMRCA distribution also includes the dates of the earliest
267 reported Tn125-*bla*_{NDM}-positive isolates from 2005⁷ which gives some credibility to these results.

268 The indications from molecular tip-dating fall into a wider narrative where *bla*_{NDM} spread was initially driven by
269 Tn125 mobilization before subsequent transposition by Tn3000, and others. However, the sparsity of SNPs within
270 the alignments, the weak temporal signal and the abundance of structural variants, plasmid backbones, transposons
271 and ISs argue in favour of genetic recombination, rather than *de novo* mutation, as the dominant mechanism
272 driving evolutionary change in the genetic region flanking *bla*_{NDM} gene.

273

274 Correlates with the global dissemination of *bla*_{NDM}

275 The earliest samples in our dataset span the years 2007 to 2010 and comprise 31 *bla*_{NDM}-positive genomes already
276 encompassing nine bacterial species, 13 countries, and three continents (23 confirmed clinical samples and 8 of
277 unknown origin from Asia, North America and Europe). Even though the exact time of emergence remains an
278 open question, such a wide host and geographic distribution, even in the earliest available samples, illustrates the
279 extraordinarily high mobility of *bla*_{NDM}. To track the spread of *bla*_{NDM} we estimated diversity over time for several
280 categorizations of *bla*_{NDM}-positive samples (Supplementary Figure 11, see Methods). In particular, for each year,
281 the diversity was estimated among samples' country of collection, associated bacterial genera, replicon types (i.e.,
282 plasmid backbones), SNP counts within 5000 bp alignment, and structural variants at positions 3000 bp and 5000
283 bp downstream of the *bla*_{NDM} gene. Shannon entropy was used as a measure of diversity and bootstrapping
284 implemented to provide confidence intervals around the entropy estimates. A strong sampling bias is present
285 among isolates from the same NCBI BioProject (Supplementary Figure 10). To account for this, we weighted
286 contigs during bootstrapping based on their BioProject affiliation (see Methods).

287 The change in diversity of the countries associated to *bla*_{NDM}-positive isolates was used to approximate the broad
288 patterns of global dissemination of NDM resistance (Supplementary Figure 11A). The diversity of sampling
289 countries through time plateaued between 2013-2015. In light of the earliest reports of NDM-positive samples in
290 2005, this indicates that it took eight to eleven years for NDM resistance to spread globally and is consistent with
291 our estimates based on phylogenetic tip-dating (Supplementary Figure 9C). Furthermore, the change in the
292 diversity of countries associated to *bla*_{NDM}-positive genomes was found to be positively correlated with all other
293 considered categories (Supplementary Figure 12) suggesting it holds information which can be leveraged to
294 reconstruct dissemination trends. The weakest correlation with the widest confidence interval was found between
295 the number of SNPs in the alignment and the diversity of countries of sample origin ($\rho = 0.407$ [0.119-0.753]),
296 followed by the bacterial genera ($\rho = 0.5$ [0.217-0.7]), then structural variants at 3000bp downstream of *bla*_{NDM}
297 ($\rho = 0.533$ [0.217-0.717], and 5000bp downstream ($\rho = 0.683$ [0.433-0.85]). Despite the overlap of confidence
298 intervals, this ordering again highlights the importance of genetic rearrangements and transposition in the
299 evolution of this genetic region.

300 The strongest correlation was found between the diversity of countries with NDM-positive isolates and the
301 replicon types of associated plasmid backbones ($\rho = 0.7$ [0.467-0.883]) supporting a strong dependence between
302 the two (Supplementary Figure 12B). To further investigate this relationship, we assessed the correlation between
303 genetic and geographical distance between pairs of contigs as a function of the distance downstream of *bla*_{NDM}
304 gene (Figure 4, see Methods). Starting from *bla*_{NDM} and moving downstream, we gradually extended the region
305 over which genetic distances were estimated. At each step, we estimated the correlation between genetic and
306 geographic distance.

307 Considering all contig sequences, a gradual increase in correlation between genetic and geographic distance was
308 observed as more of the sequence downstream of *bla*_{NDM} was included (Figure 4A). The same trend is observed
309 in an isolated case of "broad-range" IncF plasmids which have a wide geographical distribution (Figure 4B,
310 Supplementary Figure 2). However, no significant or sufficiently long consecutive correlations were found among
311 IncX3 and IncN plasmids (Supplementary Figure 13) likely due to the lack of longer plasmid sequences and more

312 restricted mean geographic distance between pairs of plasmids; both replicon types are mostly found in China and
313 India respectively (Supplementary Figure 2).

314 Nevertheless, considering *bla_{NDM}* is predominantly carried by plasmids¹, the trend identified in Figure 4 suggests
315 that plasmids carrying *bla_{NDM}* are geographically structured. Gene dissemination is a fundamentally spatial
316 process. Despite being theoretically mobile, in practice most plasmids may be both strongly host-constrained³⁶
317 and associated with particular locations or environmental niches³⁷. All in all, this could be hinting at the existence
318 of plasmid niches: settings to which particular plasmids are more adapted.

319 Discussion

320 Increasing levels of antimicrobial resistance in bacterial pathogens pose a major global health challenge, with
321 resistance to carbapenems a particularly concerning example. Understanding the main mechanisms by which
322 antibiotic resistance elements are disseminated is fundamental to our understanding of the spread of AMR, and
323 new methods are required to fully reconstruct the forces underlying the dynamic mobilome common to many
324 resistance elements. Here, we have compiled a global dataset of 2,148 bacterial genomes carrying *bla*_{NDM},
325 including 112 new hybrid assemblies from Chinese hospitals, to provide a comprehensive overview of the
326 different genetic backgrounds harbouring this resistance element and to gain insight into its mobility. In order to
327 do this, we developed a new alignment-based method to resolve the complex structural variations flanking this
328 major antibiotic resistance element.

329 Our results, summarized in Figure 3, highlight the vast diversity of genetic backgrounds and plasmids harbouring
330 *bla*_{NDM} and the predisposition of this region for genetic reshuffling. Moreover, we detected a markedly low SNP
331 prevalence and weak temporal signal, which points to the importance of genetic recombination and transposition
332 in driving the evolution of this region. In addition, we identified 18 different putative transposons within our
333 dataset, of which Tn125, Tn3000 and IS26 flanked pseudo-composite transposon are predominant and represent
334 the major contributors to plasmid jumps of *bla*_{NDM}. IS26 seems particularly promiscuous; it is often found inserted
335 at various positions around *bla*_{NDM} and with some indication of within-isolate activity. IS26 is known for its
336 increased activity and rearrangement of plasmids in clinical isolates³⁸ and has been observed to drive within-
337 plasmid heterogeneity even in a single *E. coli* isolate³⁹. Thus, it is a likely candidate driving *bla*_{NDM} gene
338 acquisition and extensive rearrangements found within *bla*_{NDM} region. Furthermore, IS5 was often and uniquely
339 found immediately upstream of *bla*_{NDM} and its peculiar positioning could foreshadow its role in increased
340 transcription of the gene³⁰. Little to no evidence was found for the involvement of integrons and RC transposition
341 of ISCR elements in spreading of *bla*_{NDM}. In fact, ISCR1 alongside other ISs, was mainly found disrupting the
342 *bla*_{NDM} region.

343 By assessing the change in entropy of countries where *bla*_{NDM}-positive isolates have been sequenced over time,
344 we traced the patterns underlying the spread of NDM resistance. Our assessment of diversity suggests that,
345 following a rapid dissemination, the spread of *bla*_{NDM} may have reached a plateau between 2013-2015, with
346 *bla*_{NDM} reaching a global prevalence 8-11 years after 2005. Such a rapid spread has also been suggested for other
347 significant mobile resistance genes: the *mcr-1* gene, mediating colistin resistance, is also estimated to have reached
348 global prevalence within a decade⁴⁰. The extent to which this model of 'rapid spread' applies to other transposon-
349 borne resistance elements remains to be determined.

350 We found a strong positive correlation between genetic distances between plasmid backbones bearing *bla*_{NDM} and
351 the geographic location in which they were sampled, suggesting the existence of a constraint on plasmid spread
352 i.e. plasmid niches. We presume plasmid niches exist thanks to local evolutionary pressures for which particular
353 plasmid backbones are optimized. Country boundaries limiting population movement, region-specific outbursts
354 of antibiotic usage and narrow host range of the majority of bacterial plasmids³⁶ all likely contribute to a restricted
355 geographical range. Thus, an introduction of another plasmid into a foreign plasmid niche may lead to plasmid

356 loss or fast adaptation by, for instance, acquisition of resistance and other accessory elements. This hypothetical
357 scenario also provides an opportunity for resistance to spread by transposition or recombination, by which a new
358 resistance gene is able to enter another plasmid niche. In the case of *bla_{NDM}*, this would also imply that after the
359 initial introduction of *bla_{NDM}* to a geographic region, dissemination and persistence of the gene could proceed
360 idiosyncratically - selection for carbapenem resistance being just one of many selective pressures acting on
361 plasmid diversity.

362 The importance of transposon movement has been previously demonstrated by our work on plasmid networks³⁶,
363 as well as several papers promoting a Russian-doll model of resistance mobility^{40,41}. In light of our results, we
364 suggest a conceptual framework of resistance gene dissemination where plasmid mobility is for the most part
365 restricted. Although plasmids can facilitate rapid spread within species and geographical regions, the momentum
366 of resistance dissemination is primarily reliant on between-plasmid transposon jumps and genetic recombination.

367 Methods

368

369 Compiling the dataset of NDM sequences

370 We compiled a global dataset of 2,148 bacterial genomes carrying the *bla_{NDM}* gene from several publicly available
371 databases. The vast majority of bacterial isolates were collected from patients (1,501), while 308 are of animal
372 origin (184 from chickens, 51 from other birds and 47 from flies), 244 are of an unknown origin, and 95 are
373 environmental samples (of which 36 are isolated from hospital environments). 1239 and 275 fully assembled
374 genomes were downloaded from NCBI Reference Sequence Database (RefSeq; accessed on 23rd of May 2019)^{42,43}
375 and EnteroBase⁴⁴ respectively. The EnteroBase repository was screened using BlastFrost (v1.0.0)⁴⁵ allowing for
376 one mismatch. In addition, we used the Bitsliced Genomic Signature Index (BIGSI) tool⁴⁶ to identify all Sequence
377 Read Archive (SRA) unassembled reads which carry the *bla_{NDM}* gene. At the time of writing, a publicly available
378 BIGSI demo did not include sequencing datasets from after December 2016. Therefore, we manually indexed and
379 screened an additional 355,375 SRA bacterial sequencing datasets starting from January 2017 to January 2019.
380 We required the presence of 95% of *bla_{NDM-1}* *k*-mers to identify NDM-positive samples from raw SRA reads. This
381 led to the inclusion of 522 isolates from reads downloaded from the SRA repository. Furthermore, we generated
382 112 new NDM-positive genomes using paired-end Illumina (Illumina HiSeq 2500) and PacBio (PacBio RS II)
383 sequencing of isolates from 87 hospitalized patients across China and 25 livestock farms. The sequenced isolates
384 were selected from two previous studies^{47,48}. The sequencing reads are available on the Short Read Archive (SRA)
385 under accession number **XXXXXXXX**. All reads were de novo assembled using Unicycler (v0.4.8)⁴⁹ using default
386 parameters while also specifying hybrid mode for those isolates for which we had both Illumina short-read and
387 PacBio long read sequencing data. Spades (v3.11.1)⁵⁰ was applied, without additional polishing, for cases where
388 Unicycler assemblies failed to resolve. Sequencing datasets without associated metadata on the date of sampling
389 were not included in the analysis.

390 In total, 2,165 contigs carrying the *bla_{NDM}* gene were identified using BLAST (v2.10.1+)²⁷. The full metadata
391 table of contigs containing *bla_{NDM}* is available as Supplementary Data 1. The table includes sample accession
392 numbers and information on host organism, collection date, sampling location, assembly status, and contig
393 plasmid type and circularity. Sixteen contigs (C165, C964, GCA_000764615, GCA_000814145,
394 GCA_001860505, GCA_002133365, GCA_002870165, GCA_003194305, GCA_003368345, GCA_003716765,
395 GCA_003860815, GCA_003950255, GCA_003991465, GCA_004795525, GCA_005155965, GCF_004357815)
396 were found to carry more than one copy of *bla_{NDM}* and were not included in our analyses. Two assemblies
397 (GCF_004358085 and GCF_004357805) had a single *bla_{NDM}* gene split into two contigs; these four contigs were
398 also excluded. Contigs GCA_00386065, C184 and C141 were removed due to poor assembly quality. This
399 filtering resulted in a dataset of 2,142 contigs (2,128 genomes) which were used in all subsequent analyses. Of
400 these, six genomes were found to contain *bla_{NDM}* on two contigs, each one harbouring a single copy of *bla_{NDM}*.

401

402 Annotating NDM-positive contigs

403 Coding sequences (CDS) of all NDM-positive contigs were annotated using the Prokka (v1.12)⁵¹ and Roary
404 (v3.12.0)⁵² pipelines run with default parameters. In addition, plasmid sequences were confirmed based on RefSeq
405 annotation (i.e., contigs labelled “plasmid”), contig circularity reported by Unicycler, or by the presence of a

406 plasmid replicon sequence⁵³. To identify plasmid replicon types, the contigs were screened against the
407 PlasmidFinder database (version 2020-02-25)²⁶ using BLAST (v2.10.1+)²⁷ where only BLAST hits with a
408 minimum coverage of 80% and percentage identity of >95% were retained. In cases where two or more replicon
409 hits were found at overlapping positions on a contig, the one with the higher percentage identity was retained. All
410 identified plasmid types are provided in Supplementary Data 1.

411

412 Resolving structural variants of NDM-positive contigs

413 Structural variations upstream and downstream of *bla*_{NDM} were resolved using a novel alignment-based approach,
414 as illustrated in Figure 2. First, contigs carrying *bla*_{NDM} were reoriented such that *bla*_{NDM} gene is on the positive-
415 sense DNA strand (i.e., facing 5' to 3' direction). A discontinuous Mega BLAST (v2.10.1+)²⁸ search with default
416 settings was applied against all pairs of retained contigs. This method was selected over the regular Mega BLAST
417 implementation as it is comparably fast, but more permissive towards dissimilar sequences with frequent gaps and
418 mismatches. BLAST hits including a complete *bla*_{NDM} gene on both contigs were selected and cropped to either
419 (i) the start of *bla*_{NDM} gene and the downstream sequence or (ii) the end of the *bla*_{NDM} gene and the upstream
420 sequence depending on the analysis at hand: the downstream or the upstream analysis respectively. This trimming
421 establishes *bla*_{NDM} as an anchor and forces the algorithm to consider only the region upstream or downstream of
422 the gene.

423 Next, the algorithm proceeds with a stepwise network analysis of BLAST hits. For this purpose, a ‘splitting
424 threshold’ was introduced. Starting from zero, the threshold is gradually increased by 10 bp. At each step, BLAST
425 hits with a length lower than the value given by the ‘splitting threshold’ are excluded. Then, a network is
426 constructed from the remaining BLAST hits such that contigs sharing a BLAST hit are connected with an edge.
427 The network is then broken down into components – groups of nodes (contigs) that share a common edge. It is
428 expected that contigs within each component share a homologous region downstream (or upstream) of *bla*_{NDM} at
429 least of the length given by the threshold. It is therefore not possible for a single contig to be assigned to multiple
430 components. Components of size <5 bp are labelled as ‘Other Structural Variants’ and are not considered in further
431 analyses. Also, contigs that are shorter than the defined ‘splitting threshold’ and share no edge with any other
432 contig are considered as ‘cutting short’.

433 By tracking the splitting of the network as the ‘splitting threshold’ is increased, one can determine clusters of
434 homologous contigs at any given position downstream or upstream from the anchor gene (here *bla*_{NDM}), as well
435 as the homology breakpoint. The precision of the algorithm is directly influenced by the step size which is, in this
436 case, 10 bp and the alignment algorithm, in this case discontinuous Mega BLAST. The described algorithm is
437 available at https://github.com/macman123/track_structural_variants

438

439 Date randomization, linear regression analyses and molecular tip-dating.

440 The 45 complete *Tn125* and 29 complete *Tn3000* contigs harbouring *bla*_{NDM} were sequentially aligned
441 (--pileup flag) using Clustal Omega (v1.2.3)⁵⁴ specifying *bla*_{NDM-1} as a profile. The consensus sequence over the
442 alignment was considered the closest match to a putative ancestral sequence and was hence used as a reference to

443 identify SNPs against. This approach was motivated by the fact that: (i) there is no appropriate outgroup sequence
444 available; (ii) the oldest contigs in the dataset can harbour non-ancestral SNPs; (iii) due to a short time span and
445 relatively few mutations present, it is unlikely that any one non-ancestral SNP has become dominant in the
446 population.

447 Date randomization and linear regression analyses considering the number of SNPs accumulated against the year
448 of sample collection provide an estimate of the strength of the temporal signal in the alignment⁵⁵⁻⁵⁷. We weighted
449 the linear regressions by the BioProject affiliation of the sequences in the alignments of the two transposons
450 (Supplementary Figure 8A and 8B). This was done to control the strong sampling biases present among samples
451 from the same NCBI BioProject, with contigs from the same BioProject tending to be genetically similar
452 irrespective of the sampling year (Supplementary Figure 9). While both Tn125 and Tn3000 showed positive
453 temporal signal (Supplementary Figure 8A and B), neither regression was significant ($p=0.1279$ and $p=0.1375$
454 respectively). The low sample size and the low genetic diversity in the two alignments may limit the statistical
455 power to detect temporal signal. Date-randomization analysis also showed that the estimated evolutionary rate for
456 both transposons fell within the distribution of slopes on randomized dates (Supplementary Figure 8A and B).

457 A further test of meaningful signal in the data is to consider the degree to which the dated alignment can drive the
458 posterior distribution away from the priors specified in Bayesian dating frameworks. BEAST2 (v2.6.0)³⁵ was run
459 on both transposon alignments specifying a strict molecular clock rate with a model averaging prior on the
460 substitution model⁵⁸ and a MCMC chain length of 5×10^8 (Supplementary Data 2). The long MCMC chain length
461 was chosen to ensure convergence. For both runs the Serial Birth-Death Skyline (BDSS) model was specified as
462 the tree prior. The BDSS model is commonly used for viral epidemics⁵⁹ which share many parallels with AMR
463 outbreaks. Similar to other birth-death models, the BDSS prior consists of three parameters: a rate of transmission
464 (an estimate transposon/plasmid mobility), recovery (an estimate of transposon fossilization or plasmid loss), and
465 sampling rate. Also, unlike coalescent models, BDSS does not attempt to estimate population sizes, which have
466 limited applicability to dating small genetic regions and mobile elements. We evaluated the prior and posterior
467 distributions across variables after discarding the first 20% of burn-in and after ensuring model convergence (an
468 effective sample size >200).

469

470 Estimating Shannon entropy among NDM-positive contigs

471 We estimated Shannon entropy ('diversity') for several categorizations of bla_{NDM} -containing contigs: country of
472 sampling, bacterial host genera, replicon type, SNP count within a 5000 bp alignment, and structural variants at
473 positions 3000 bp and 5000 bp downstream of the bla_{NDM} gene. The 5000 bp alignment consisted of 654 contigs
474 harbouring bla_{NDM} , ble , $trpF$, $dsbD$, $cutA$, $groS$ and $groL$ genes. To estimate entropy, we used a weighted
475 bootstrapping approach (1000 iterations) with the probability of pooling any one sample inversely proportional to
476 the number of samples contained in the corresponding BioProject. At each iteration, entropy was estimated for a
477 sampled set of contigs (X) classified into n unique categories according to the following formula:

$$478 \quad H(X) = -\sum_{i=1}^n P(x_i) \log P(x_i),$$

479 where probability $P(x_i)$ of any sample belonging to any particular category x_i (e.g., country or replicon type) is
480 approximated using the category's frequency. Accordingly, higher entropy values indicate an abundance of
481 equally likely categories, while lower entropy indicates a limited number of categories.

482

483 Estimating geographical and genetic distance between contigs

484 Geographical distance between pairs of selected contigs was determined using the *geodist*⁶⁰ R package and
485 reported sampling coordinates or centroids of countries of collection if the former was not available. The exact
486 Jaccard distance (JD) was used as a measure of the genetic distance. It was calculated using the tool Bindash⁶¹
487 with k -mer size equal to 21 bp. The JD is defined as the fraction of total k -mers not shared between two contigs.
488 For instance, JD=1 denotes no k -mers are shared. The two distance matrices (genetic and geographic) were
489 assessed using the *mantel* function from *vegan*⁶² package in R. To account for the sampling bias, pairs of contigs
490 belonging to the same BioProject were not considered while estimating the Spearman correlation and performing
491 the Mantel test between geographic and genetic distance.

492 Acknowledgements

493 M.A. was supported by a Ph.D. scholarship from University College London. H.W. is supported by National
494 Natural Science Foundation of China (81625014). L.v.D., H.W. and F.B. acknowledge financial support from the
495 Newton Fund UK-China NSFC initiative (MRC Grant MR/P007597/1 and 81661138006). L.v.D. and F.B. are
496 supported from a Wellcome Institutional Strategic Support Fund (ISSF3) – AI in Healthcare (19RX03). F.B.
497 additionally acknowledges support from the BBSRC GCRF scheme and the National Institute for Health Research
498 University College London Hospitals Biomedical Research Centre. L.v.D is supported by a UCL Excellence
499 Fellowship. M.A., L.v.D and F.B. acknowledge UCL Biosciences Big Data equipment grant from BBSRC
500 (BB/R01356X/1). L.P.S. acknowledges funding from the Antimicrobial Resistance Cross-council Initiative
501 supported by the seven UK research councils (NE/N019989/1). The funders had no role in study design, data
502 collection, interpretation of results, or the decision to submit the work for publication. Lastly, M.A. would like to
503 thank Nicola de Maio for informal discussions which led to the idea for the algorithm used to track structural
504 variants.

505 Competing interests

506 The authors declare no financial or non-financial competing interests.

507 Contributions

508 M.A., F.B., L.v.D. and H.W. conceived the project and designed the experiments. M.A., L.v.D., L.P.S., and N.L.
509 collected data from online repositories. R.W., Y.Y., Q.W., S.S, and H.C sequenced samples from Chinese
510 hospitals. M.A., L.v.D, and R.W. *de novo* assembled all the genomes. M.A. performed all the analyses under the
511 guidance of L.v.D and F.B. M.A., L.v.D. and F.B. take responsibility for the accuracy and availability of the
512 results. M.A. wrote the paper with contributions from L.v.D. and F.B. All authors read and commented on
513 successive drafts and all approved the content of the final version.

514 References

515

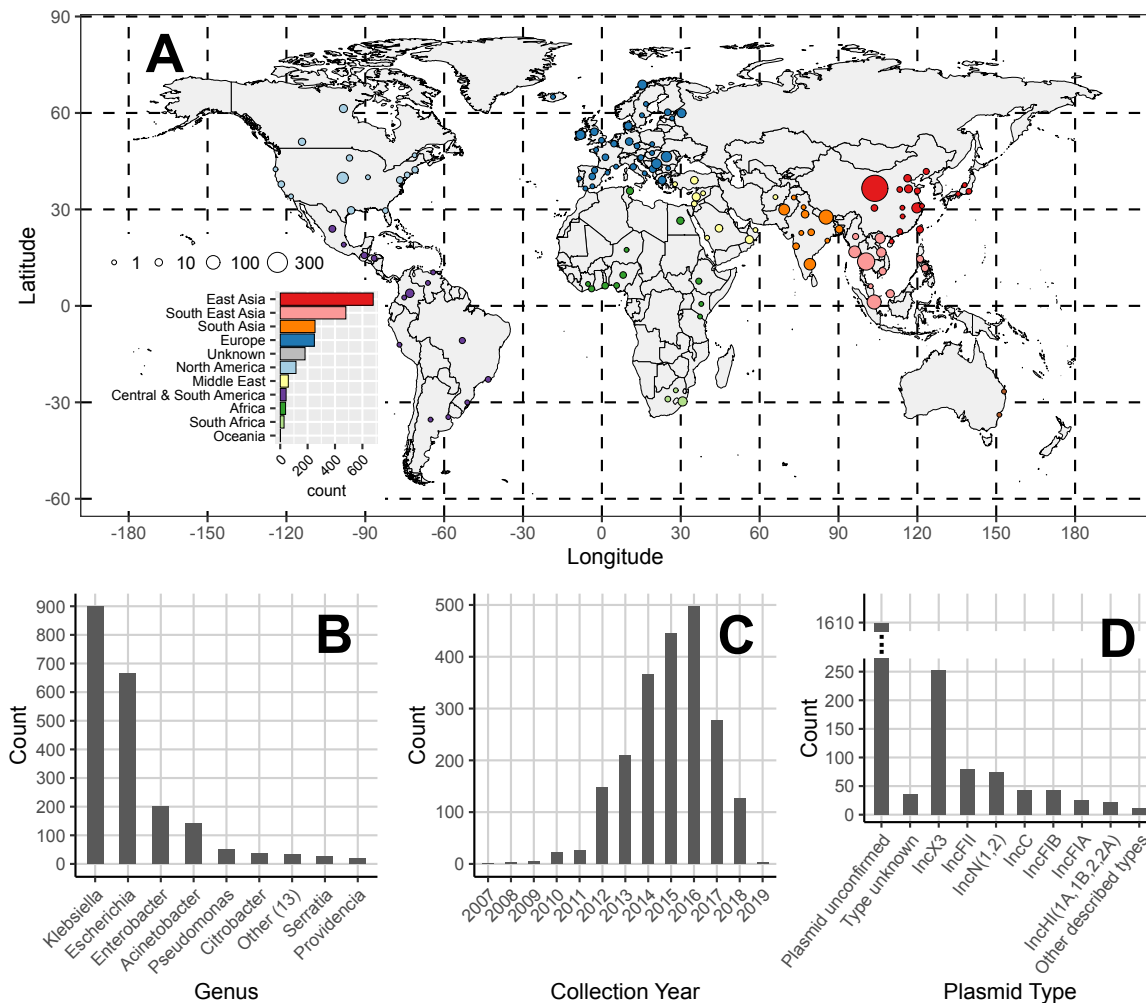
- 516 1. Wu, W. *et al.* NDM metallo- β -lactamases and their bacterial producers in health care settings. *Clinical*
517 *Microbiology Reviews* vol. 32 (2019).
- 518 2. Yong, D. *et al.* Characterization of a new metallo- β -lactamase gene, bla NDM-1, and a novel
519 erythromycin esterase gene carried on a unique genetic structure in *Klebsiella pneumoniae* sequence
520 type 14 from India. *Antimicrob. Agents Chemother.* **53**, 5046–5054 (2009).
- 521 3. Struelens, M. J. *et al.* New Delhi metallo-beta-lactamase 1–producing Enterobacteriaceae: emergence
522 and response in Europe. *Eurosurveillance* **15**, 19716 (2010).
- 523 4. Kumarasamy, K. K. *et al.* Emergence of a new antibiotic resistance mechanism in India, Pakistan, and
524 the UK: A molecular, biological, and epidemiological study. *Lancet Infect. Dis.* **10**, 597–602 (2010).
- 525 5. Poirel, L., Dortet, L., Bernabeu, S. & Nordmann, P. Genetic Features of bla NDM-1-Positive
526 Enterobacteriaceae. *Antimicrob. Agents Chemother.* **55**, 5403–5407 (2011).
- 527 6. Castanheira, M. *et al.* Early dissemination of NDM-1- and OXA-181-producing Enterobacteriaceae in
528 Indian hospitals: Report from the SENTRY Antimicrobial Surveillance Program, 2006–2007.
529 *Antimicrob. Agents Chemother.* **55**, 1274–1278 (2011).
- 530 7. Jones, L. S. *et al.* Plasmid carriage of blaNDM-1 in clinical *Acinetobacter baumannii* isolates from India.
531 *Antimicrob. Agents Chemother.* **58**, 4211–4213 (2014).
- 532 8. Roca, I. *et al.* Molecular characterization of NDM-1-producing *Acinetobacter pittii* isolated from
533 Turkey in 2006. *J. Antimicrob. Chemother.* **69**, 3437–3438 (2014).
- 534 9. Poirel, L., Bonnin, R. A. & Nordmann, P. Analysis of the resistome of a multidrug-resistant NDM-1-
535 producing *Escherichia coli* strain by high-throughput genome sequencing. *Antimicrob. Agents*
536 *Chemother.* **55**, 4224–4229 (2011).
- 537 10. Toleman, M. A., Spencer, J., Jones, L. & Walsh, T. R. bla NDM-1 is a chimera likely constructed in
538 *Acinetobacter baumannii*. *Antimicrob. Agents Chemother.* **56**, 2773–2776 (2012).
- 539 11. Partridge, S. R. & Iredell, J. R. Genetic Contexts of bla NDM-1. *Antimicrobial Agents and*
540 *Chemotherapy* vol. 56 6065–6067 (2012).
- 541 12. Toleman, M. A., Bennett, P. M. & Walsh, T. R. ISCR Elements: Novel Gene-Capturing Systems of the
542 21st Century? *Microbiol. Mol. Biol. Rev.* **70**, 296–316 (2006).
- 543 13. Ilyina, T. S. Mobile ISCR elements: Structure, functions, and role in emergence, increase, and spread of
544 blocks of bacterial multiple antibiotic resistance genes. *Molecular Genetics, Microbiology and Virology*
545 vol. 27 135–146 (2012).
- 546 14. Poirel, L. *et al.* Tn125-related acquisition of blaNDM-like genes in *Acinetobacter baumannii*.
547 *Antimicrob. Agents Chemother.* **56**, 1087–1089 (2012).
- 548 15. Sekizuka, T. *et al.* Complete Sequencing of the blaNDM-1-Positive IncA/C Plasmid from *Escherichia*
549 *coli* ST38 Isolate Suggests a Possible Origin from Plant Pathogens. *PLoS One* **6**, e25334 (2011).
- 550 16. Basu, S. Variants of the New Delhi metallo- β -lactamase: New kids on the block. *Future Microbiology*
551 vol. 15 465–467 (2020).
- 552 17. Baraniak, A. *et al.* NDM-producing Enterobacteriaceae in Poland, 2012–14: inter-regional outbreak of
553 *Klebsiella pneumoniae* ST11 and sporadic cases. *J. Antimicrob. Chemother.* **71**, 85–91 (2016).

- 554 18. Rahman, M. *et al.* Prevalence and Molecular Characterization of New Delhi Metallo-Beta-Lactamases
555 in Multidrug-Resistant *Pseudomonas aeruginosa* and *Acinetobacter baumannii* from India. *Microb.*
556 *Drug Resist.* **24**, 792–798 (2018).
- 557 19. Wailan, A. M. *et al.* Genetic contexts of blaNDM-1 in patients carrying multiple NDM-producing
558 strains. *Antimicrob. Agents Chemother.* **59**, 7405–7410 (2015).
- 559 20. González, L. J. *et al.* Membrane anchoring stabilizes and favors secretion of New Delhi metallo- β -
560 lactamase. *Nat. Chem. Biol.* **12**, 516–522 (2016).
- 561 21. Chatterjee, S., Mondal, A., Mitra, S. & Basu, S. *Acinetobacter baumannii* transfers the blaNDM-1 gene
562 via outer membrane vesicles. *J. Antimicrob. Chemother.* **72**, 2201–2207 (2017).
- 563 22. Lynch, T. *et al.* Molecular evolution of a klebsiella pneumoniae st278 isolate harboring blandm-7 and
564 involved in nosocomial transmission. *J. Infect. Dis.* **214**, 798–806 (2016).
- 565 23. Huang, T. W. *et al.* Copy Number Change of the NDM-1 Sequence in a Multidrug-Resistant Klebsiella
566 pneumoniae Clinical Isolate. *PLoS One* **8**, 1–12 (2013).
- 567 24. Sahl, J. W. *et al.* Phylogenetic and genomic diversity in isolates from the globally distributed
568 *Acinetobacter baumannii* ST25 lineage. *Sci. Rep.* **5**, (2015).
- 569 25. Bonnin, R. A. *et al.* Dissemination of New Delhi metallo- β -lactamase-1-producing *Acinetobacter*
570 *baumannii* in Europe. *Clin. Microbiol. Infect.* **18**, E362–E365 (2012).
- 571 26. Carattoli, A. *et al.* In Silico Detection and Typing of Plasmids using PlasmidFinder and Plasmid
572 Multilocus Sequence Typing. **58**, 3895–3903 (2014).
- 573 27. Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421 (2009).
- 574 28. Ma, B., Tromp, J. & Li, M. PatternHunter: Faster and more sensitive homology search. *Bioinformatics*
575 **18**, 440–445 (2002).
- 576 29. Campos, J. C. *et al.* Characterization of Tn3000, a Transposon Responsible for blaNDM-1
577 Dissemination among Enterobacteriaceae in Brazil, Nepal, Morocco, and India. *Antimicrob. Agents*
578 *Chemother.* **59**, 7387–95 (2015).
- 579 30. Schnetz, K. & Rak, B. IS5: A mobile enhancer of transcription in *Escherichia coli*. *Proc. Natl. Acad.*
580 *Sci. U. S. A.* **89**, 1244–1248 (1992).
- 581 31. Harmer, C. J., Pong, C. H. & Hall, R. M. Structures bounded by directly-oriented members of the IS26
582 family are pseudo-compound transposons. *Plasmid* vol. 111 102530 (2020).
- 583 32. Harmer, C. J., Moran, R. A. & Hall, R. M. Movement of IS26-Associated antibiotic resistance genes
584 occurs via a translocatable unit that includes a single IS26 and preferentially inserts adjacent to another
585 IS26. *MBio* **5**, (2014).
- 586 33. Sohn, J. II & Nam, J. W. The present and future of de novo whole-genome assembly. *Brief. Bioinform.*
587 **19**, 23–40 (2018).
- 588 34. Harmer, C. J. & Hall, R. M. An analysis of the IS6/IS26 family of insertion sequences: Is it a single
589 family? *Microb. Genomics* **5**, (2019).
- 590 35. Bouckaert, R. *et al.* BEAST 2.5: An advanced software platform for Bayesian evolutionary analysis.
591 *PLoS Comput. Biol.* **15**, e1006650 (2019).
- 592 36. Acman, M., van Dorp, L., Santini, J. M. & Balloux, F. Large-scale network analysis captures biological
593 features of bacterial plasmids. *Nat. Commun.* **11**, 1–11 (2020).

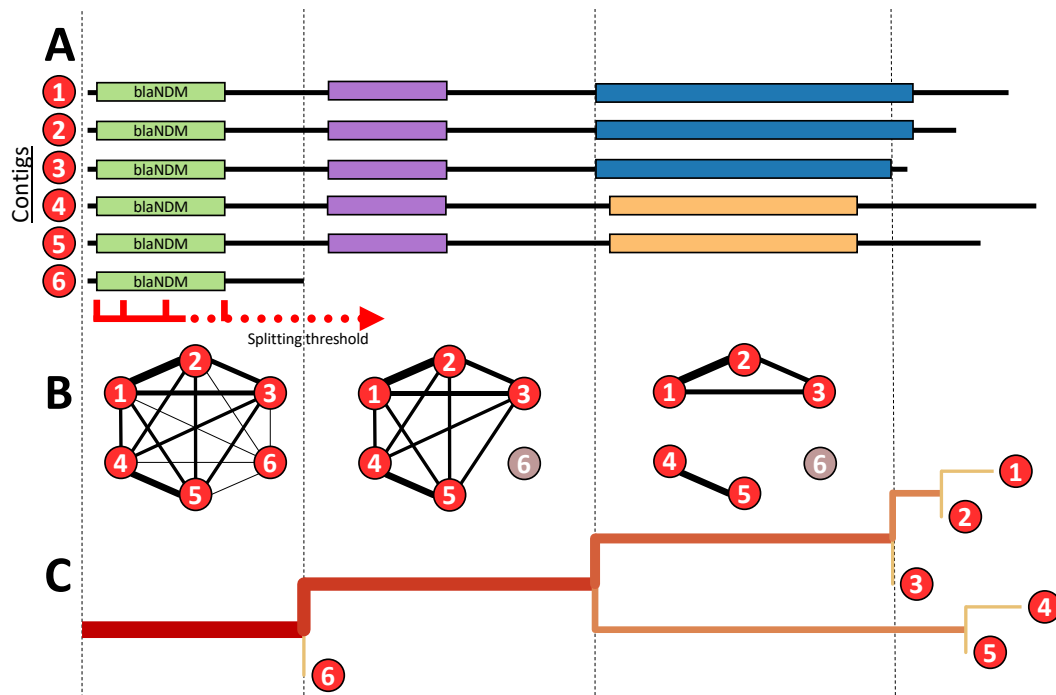
- 594 37. Shaw, L. *et al.* Niche and local geography shape the pangenome of wastewater- and livestock-associated
595 Enterobacteriaceae. 1–23 (2020) doi:10.1101/2020.07.23.215756.
- 596 38. He, S. *et al.* Insertion sequence IS26 reorganizes plasmids in clinically isolated multidrug-resistant
597 bacteria by replicative transposition. *MBio* **6**, 1–14 (2015).
- 598 39. He, D. D. *et al.* Antimicrobial resistance-encoding plasmid clusters with heterogeneous MDR regions
599 driven by IS26 in a single Escherichia coli isolate. *J. Antimicrob. Chemother.* **74**, 1511–1516 (2019).
- 600 40. Wang, R. *et al.* The global distribution and spread of the mobilized colistin resistance gene mcr-1. *Nat.*
601 *Commun.* **9**, 1–9 (2018).
- 602 41. Sheppard, A. E. *et al.* Nested Russian Doll-Like Genetic Mobility Drives Rapid Dissemination of the
603 Carbapenem Resistance Gene blaKPC. *Antimicrob. Agents Chemother.* **60**, 3767–3778 (2016).
- 604 42. Pruitt, K. D., Tatusova, T. & Maglott, D. R. NCBI reference sequences (RefSeq): A curated non-
605 redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res.* **35**, D61–D65
606 (2007).
- 607 43. O’Leary, N. A. *et al.* Reference sequence (RefSeq) database at NCBI: current status, taxonomic
608 expansion, and functional annotation. *Nucleic Acids Res.* **44**, D733–45 (2016).
- 609 44. Zhou, Z., Alikhan, N. F., Mohamed, K., Fan, Y. & Achtman, M. The EnteroBase user’s guide, with case
610 studies on Salmonella transmissions, Yersinia pestis phylogeny, and Escherichia core genomic diversity.
611 *Genome Res.* **30**, 138–152 (2020).
- 612 45. Luhmann, N., Holley, G. & Achtman, M. BlastFrost: Fast querying of 100,000s of bacterial genomes in
613 Bifrost graphs. *bioRxiv* 1–24 (2020) doi:10.1101/2020.01.21.914168.
- 614 46. Bradley, P., den Bakker, H. C., Rocha, E. P. C., McVean, G. & Iqbal, Z. Ultrafast search of all deposited
615 bacterial and viral genomic data. *Nat. Biotechnol.* **37**, 152–159 (2019).
- 616 47. Wang, R. *et al.* The prevalence of colistin resistance in Escherichia coli and Klebsiella pneumoniae
617 isolated from food animals in China: coexistence of mcr-1 and blaNDM with low fitness cost. *Int. J.*
618 *Antimicrob. Agents* **51**, 739–744 (2018).
- 619 48. Wang, Q. *et al.* Phenotypic and Genotypic Characterization of Carbapenem-resistant
620 Enterobacteriaceae: Data from a Longitudinal Large-scale CRE Study in China (2012–2016). *Clin.*
621 *Infect. Dis.* **67**, S196–S205 (2018).
- 622 49. Wick, R. R., Judd, L. M., Gorrie, C. L. & Holt, K. E. Unicycler: Resolving bacterial genome assemblies
623 from short and long sequencing reads. *PLoS Comput. Biol.* **13**, 1–22 (2017).
- 624 50. Bankevich, A. *et al.* SPAdes: a new genome assembly algorithm and its applications to single-cell
625 sequencing. *J. Comput. Biol.* **19**, 455–77 (2012).
- 626 51. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069 (2014).
- 627 52. Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics* **31**, 3691–
628 3693 (2015).
- 629 53. Orlek, A. *et al.* Plasmid classification in an era of whole-genome sequencing: Application in studies of
630 antibiotic resistance epidemiology. *Frontiers in Microbiology* vol. 8 1–10 (2017).
- 631 54. Sievers, F. *et al.* Fast, scalable generation of high-quality protein multiple sequence alignments using
632 Clustal Omega. *Mol. Syst. Biol.* **7**, 539 (2011).
- 633 55. Rieux, A. & Balloux, F. Inferences from tip-calibrated phylogenies: A review and a practical guide.

- 634 *Mol. Ecol.* **25**, 1911–1924 (2016).
- 635 56. Rambaut, A., Lam, T. T., Carvalho, L. M. & Pybus, O. G. Exploring the temporal structure of
636 heterochronous sequences using TempEst (formerly Path-O-Gen). *Virus Evol.* **2**, 1–7 (2016).
- 637 57. Duchene, S. *et al.* Bayesian Evaluation of Temporal Signal In Measurably Evolving Populations.
638 *bioRxiv* (2019) doi:10.1101/810697.
- 639 58. Bouckaert, R. R. & Drummond, A. J. bModelTest: Bayesian phylogenetic site model averaging and
640 model comparison. *BMC Evol. Biol.* **17**, 1–11 (2017).
- 641 59. Stadler, T., Kühnert, D., Bonhoeffer, S. & Drummond, A. J. Birth-death skyline plot reveals temporal
642 changes of epidemic spread in HIV and hepatitis C virus (HCV). *Proc. Natl. Acad. Sci. U. S. A.* **110**,
643 228–233 (2013).
- 644 60. Padgham, M. & Sumner, M. D. geodist: Fast, Dependency-Free Geodesic Distance Calculations.
645 (2020).
- 646 61. Zhao, X. BinDash, software for fast genome distance estimation on a typical personal laptop.
647 *Bioinformatics* **35**, 671–673 (2019).
- 648 62. Oksanen, J. *et al.* vegan: Community Ecology Package. (2019).
- 649

650 Figures



651
 652 **Figure 1. Composition of the global dataset of 2,148 NDM-positive samples.** (A) Geographic
 653 distribution of NDM-positive assemblies. Points are coloured by geographic region and their size reflects
 654 the number of samples they encompass. (B) Distribution of host bacterial genera of NDM-positive
 655 samples. (C) Distribution of sample collection years. (D) Identified plasmid types on contigs bearing the
 656 NDM-resistance. All uncircularized contigs with unknown plasmid type were labelled ‘plasmid
 657 unconfirmed’. On the other hand, all circularized contigs with an unknown plasmid type were still
 658 considered plasmids but labelled ‘type unknown’.



659

660

661

662

663

664

665

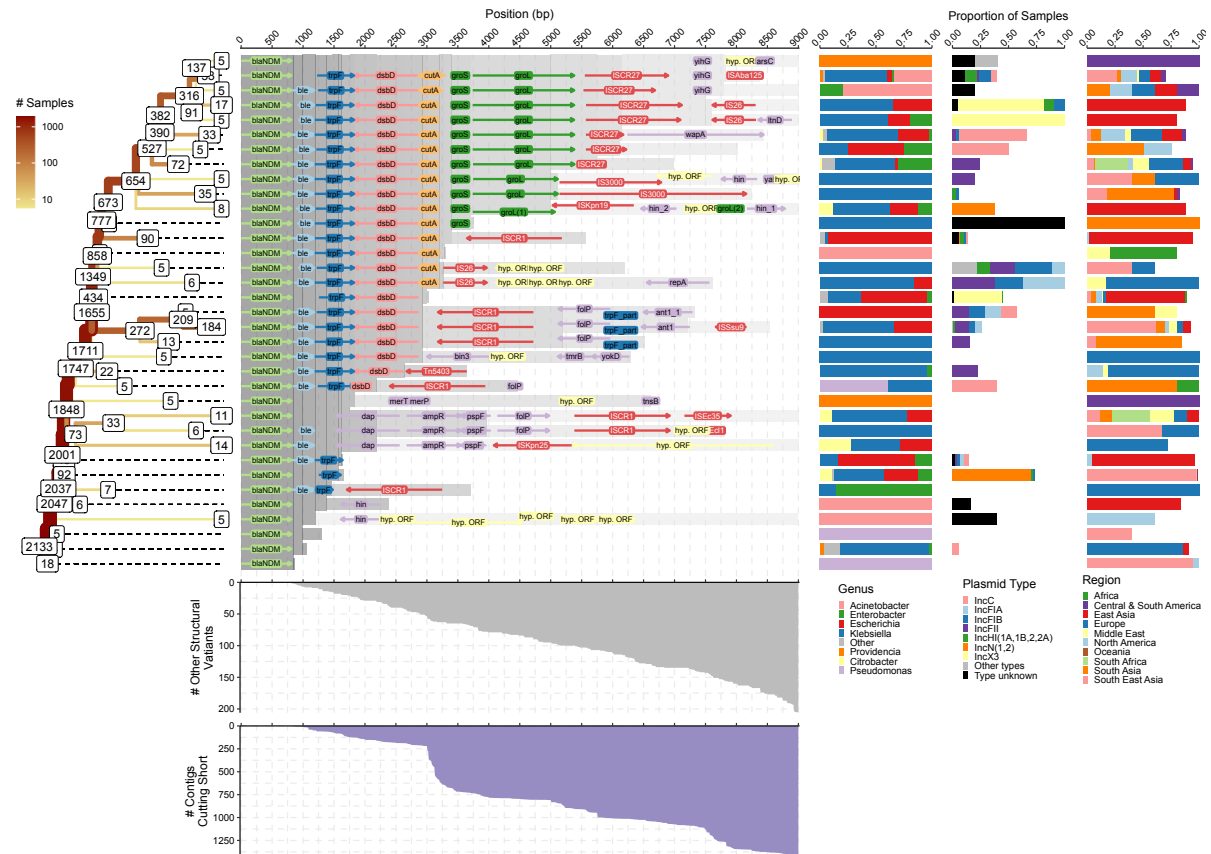
666

667

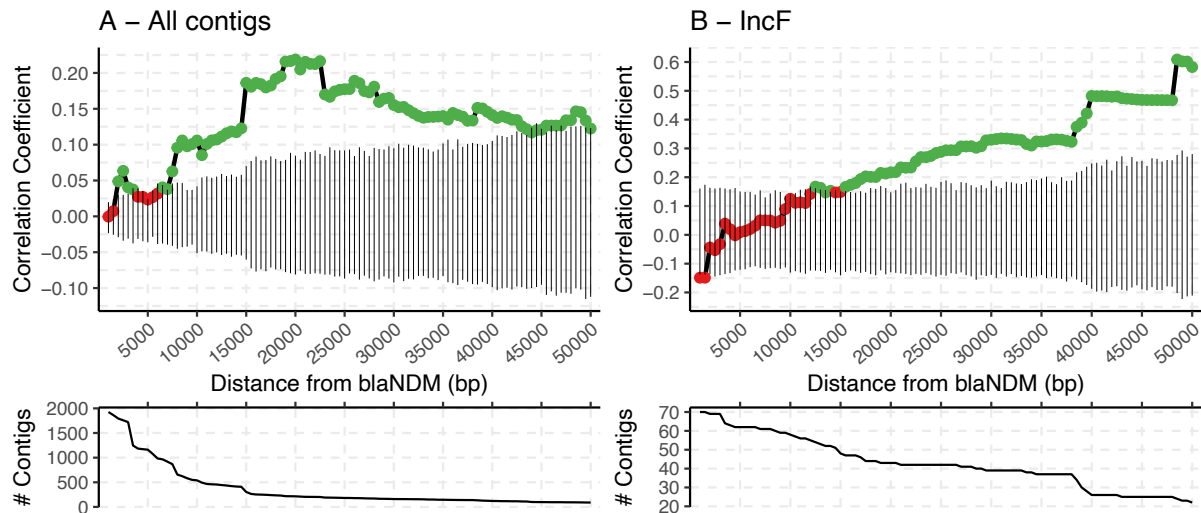
668

669

Figure 2. Schematic representation of the tracking algorithm splitting structural variant backgrounds upstream or downstream of *bla*_{NDM} gene. (A) A pairwise BLAST search is performed on all NDM-positive contigs. Starting from *bla*_{NDM} and continuing downstream or upstream, the inspected region is gradually increased using the 'splitting threshold'. (B) At each step, a graph is constructed connecting contigs (nodes) that share a BLAST hit with a minimum length as given by the 'splitting threshold'. Contigs which have the same structural variant at the certain position of the threshold belong to the same graph component, while the short contigs are singled out. (C) The splitting is visualized as a tree where branch lengths are scaled to match the position within the sequence, and the thickness and the colour intensity of the branches correspond to the number of sequences carrying the homology.



670
671 **Figure 3. Splitting of structural variants downstream of *bla*_{NDM}.** The 'splitting' tree for the most
672 common (i.e ≥ 5 contigs) structural variants is shown on the left-hand side. The labels on the nodes
673 indicate the number of contigs remaining on each branch. The other contigs either belong to other
674 structural variants or were removed due to being too short in length. The number of contigs cutting short
675 is indicated by the area chart at the bottom. Similarly, the number of contigs belonging to less common
676 structural variants is indicated by the upper area chart. The genome annotations of most common
677 structural variants are shown in the middle of the figure. The homologous regions are indicated by the
678 grey shading. Some of the structural variants and branches were intentionally cut short even though
679 their contigs were of sufficient size. This was done in order to prevent excessive bifurcation and to make
680 the tree easier to interpret. In particular, branches with percent change of contigs lost due to variation
681 and shortness above 10% were truncated. The distribution of genera, plasmid types and geographical
682 regions of samples that belong to a each of the common structural variant is shown on the right-hand
683 side.



684
685

Figure 4. The spearman correlation estimates between genetic and geographic distance of NDM-positive contigs as the DNA sequence upon which the genetic distance is measured is increased downstream of *bla*_{NDM} gene. The exact Jaccard index, an alignment-free metric, was used as a measure of genetic distance. Geographic distance between samples was estimated by the *geodist* (v0.0.6) R package using sampling coordinates or sampling country centroids if the former had not been provided. The analysis was performed on all contigs in the dataset that carry the *bla*_{NDM} gene (**A**) and the ones with confirmed IncF replicon type (**B**). In both cases, the genetic and geographic distance was measured between all pairs of contigs from a different BioProject which yielded two distance matrices: genetic and geographic. The Spearman correlation was then estimated between the two matrices and its significance evaluated using Mantel (randomization) test. Significant Spearman correlations (p-value <0.05) are indicated with green points and non-significant correlations with the red point, while the black vertical lines provide the 95% confidence interval of 1,000 Mantel test permutations. The genetic distance matrix and subsequent Spearman correlation were estimated multiple times by increasing the assessed DNA sequence starting from *bla*_{NDM} gene and continuing downstream. The two plots below the correlation graphs indicate the number of contigs used in the correlation analysis as the assessed DNA sequence is increased. See Supplementary Figure 12 for correlation analysis on IncX3 and IncN plasmids.

Figures

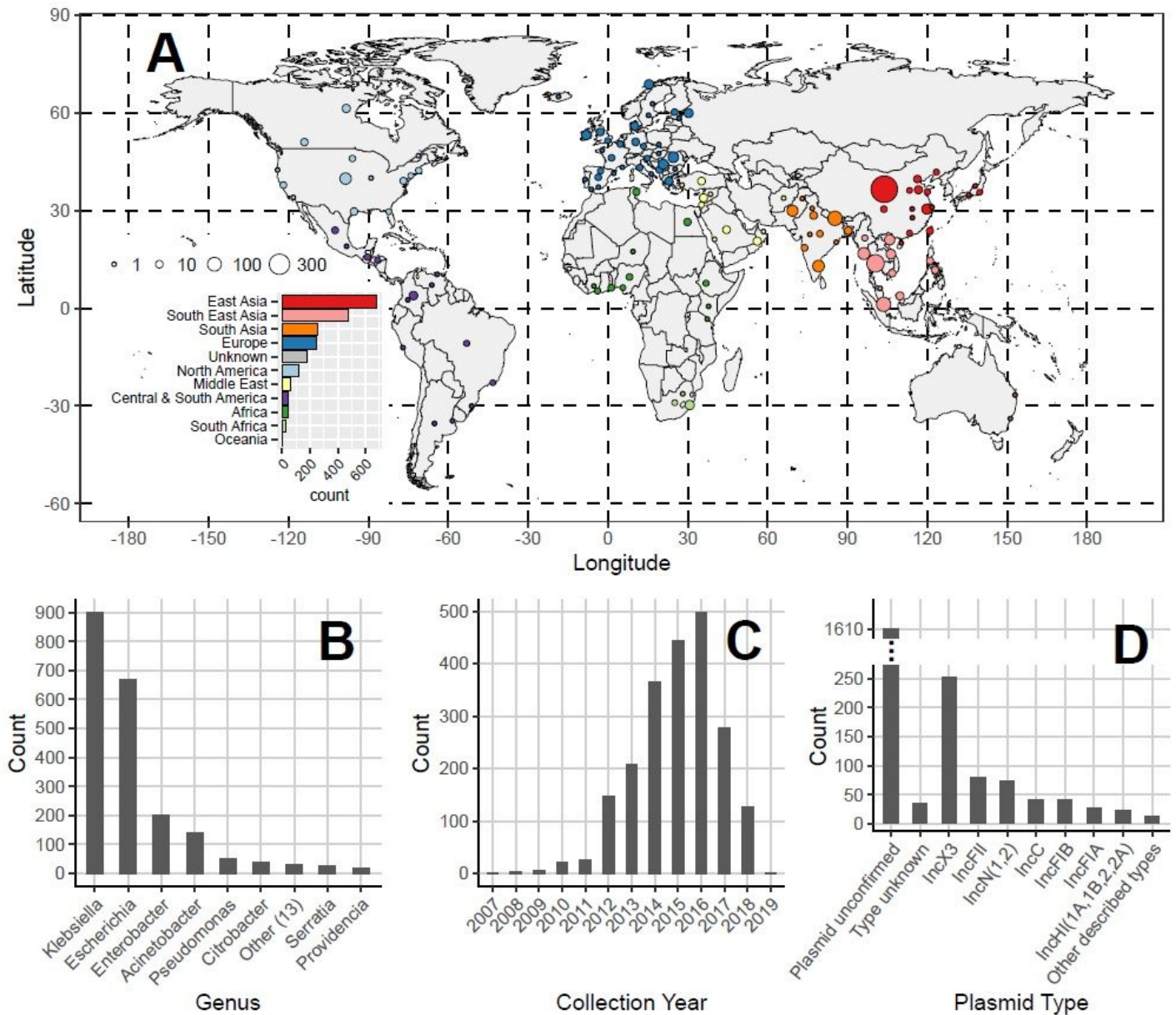


Figure 1

Composition of the global dataset of 2,148 NDM-positive samples. (A) Geographic distribution of NDM-positive assemblies. Points are coloured by geographic region and their size reflects the number of samples they encompass. (B) Distribution of host bacterial genera of NDM-positive samples. (C) Distribution of sample collection years. (D) Identified plasmid types on contigs bearing the NDM-resistance. All uncircularized contigs with unknown plasmid type were labelled 'plasmid unconfirmed'. On the other hand, all circularized contigs with an unknown plasmid type were still considered plasmids but labelled 'type unknown'. Note: The designations employed and the presentation of the material on this map do not imply the expression of any opinion whatsoever on the part of Research Square concerning

the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. This map has been provided by the authors.

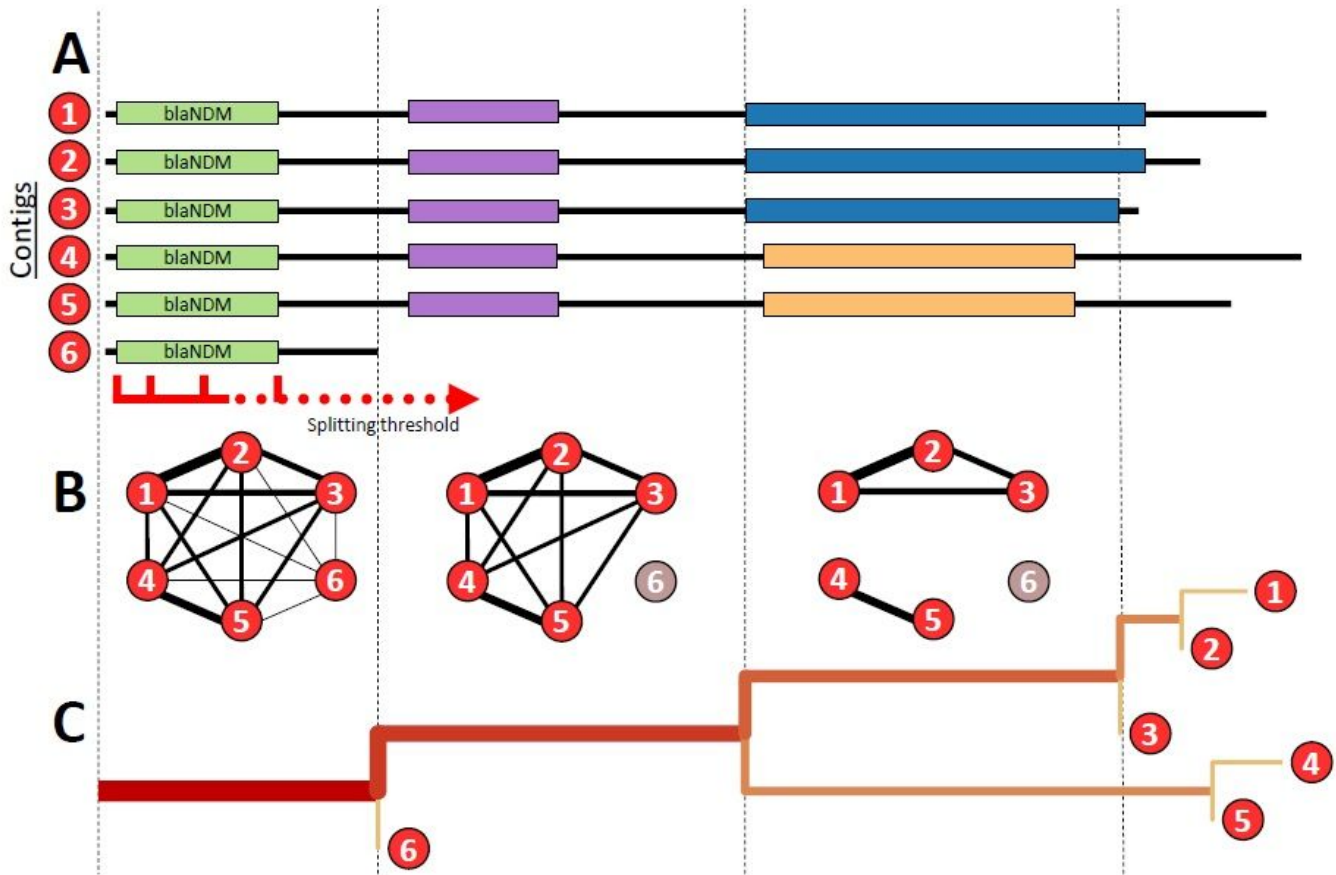


Figure 2

Schematic representation of the tracking algorithm splitting structural variant backgrounds upstream or downstream of blaNDM gene. (A) A pairwise BLAST search is performed on all NDM-positive contigs. Starting from blaNDM and continuing downstream or upstream, the inspected region is gradually increased using the 'splitting threshold'. (B) At each step, a graph is constructed connecting contigs (nodes) that share a BLAST hit with a minimum length as given by the 'splitting threshold'. Contigs which have the same structural variant at the certain position of the threshold belong to the same graph component, while the short contigs are singled out. (C) The splitting is visualized as a tree where branch lengths are scaled to match the position within the sequence, and the thickness and the colour intensity of the branches correspond to the number of sequences carrying the homology.

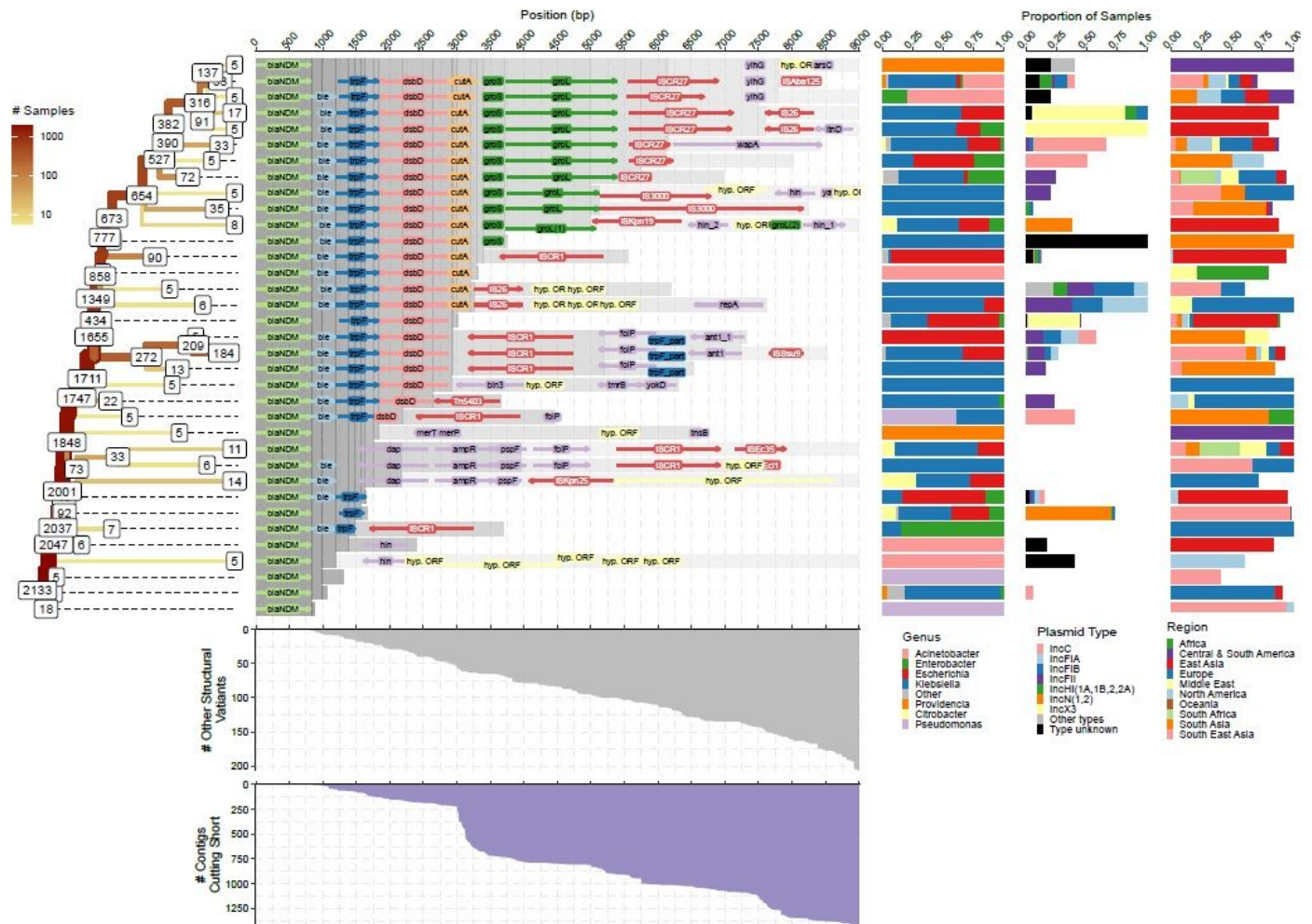


Figure 3

Splitting of structural variants downstream of *blaNDM*. The 'splitting' tree for the most common (i.e. ≥ 5 contigs) structural variants is shown on the left-hand side. The labels on the nodes indicate the number of contigs remaining on each branch. The other contigs either belong to other structural variants or were removed due to being too short in length. The number of contigs cutting short is indicated by the area chart at the bottom. Similarly, the number of contigs belonging to less common structural variants is indicated by the upper area chart. The genome annotations of most common structural variants are shown in the middle of the figure. The homologous regions are indicated by the grey shading. Some of the structural variants and branches were intentionally cut short even though their contigs were of sufficient size. This was done in order to prevent excessive bifurcation and to make the tree easier to interpret. In particular, branches with percent change of contigs lost due to variation and shortness above 10% were truncated. The distribution of genera, plasmid types and geographical regions of samples that belong to a each of the common structural variant is shown on the right-hand side.

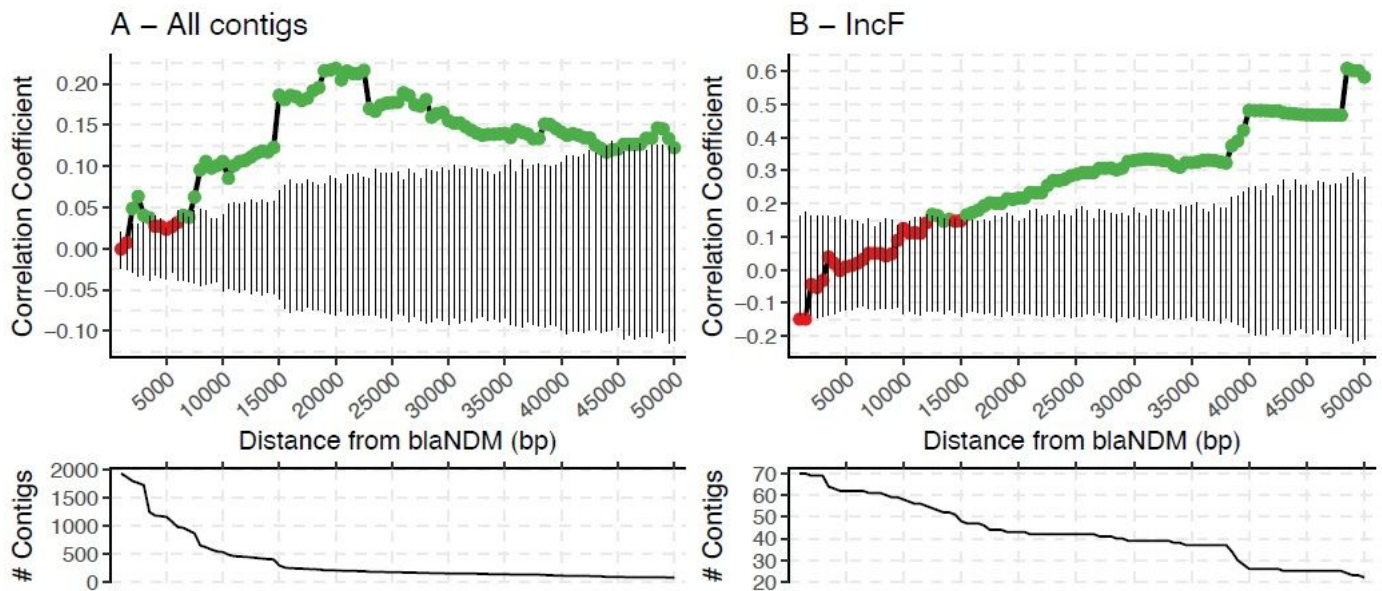


Figure 4

The spearman correlation estimates between genetic and geographic distance of NDM-positive contigs as the DNA sequence upon which the genetic distance is measured is increased downstream of blaNDM gene. The exact Jaccard index, an alignment-free metric, was used as a measure of genetic distance. Geographic distance between samples was estimated by the geodist (v0.0.6) R package using sampling coordinates or sampling country centroids if the former had not been provided. The analysis was performed on all contigs in the dataset that carry the blaNDM gene (A) and the ones with confirmed IncF replicon type (B). In both cases, the genetic and geographic distance was measured between all pairs of contigs from a different BioProject which yielded two distance matrices: genetic and geographic. The Spearman correlation was then estimated between the two matrices and its significance evaluated using Mantel (randomization) test. Significant Spearman correlations (p -value < 0.05) are indicated with green points and non-significant correlations with the red point, while the black vertical lines provide the 95% confidence interval of 1,000 Mantel test permutations. The genetic distance matrix and subsequent Spearman correlation were estimated multiple times by increasing the assessed DNA sequence starting from blaNDM gene and continuing downstream. The two plots below the correlation graphs indicate the number of contigs used in the correlation analysis as the assessed DNA sequence is increased. See Supplementary Figure 12 for correlation analysis on IncX3 and IncN plasmids.

Supplementary Files

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

- [Supplementarydata.zip](#)
- [Manuscriptsupplementary.pdf](#)