

# SCIENTIFIC REPORTS



OPEN

## Genome and plasmid diversity of Extended-Spectrum $\beta$ -Lactamase-producing *Escherichia coli* ST131 – tracking phylogenetic trajectories with Bayesian inference

Sofia Ny<sup>1,2</sup>, Linus Sandegren<sup>3</sup>, Marco Salemi<sup>4</sup> & Christian G. Giske<sup>1</sup>

Clonal lineages of ESBL (Extended-Spectrum  $\beta$ -Lactamase)-producing *E. coli* belonging to sequence type 131 (ST131) have disseminated globally during the last 30 years, leading to an increased prevalence of resistance to fluoroquinolones and extended-spectrum cephalosporins in clinical isolates of *E. coli*. We aimed to study if Swedish ESBL-producing ST131 isolates originated from single or multiple introductions to the population by assessing the amount of genetic variation, on chromosomal and plasmid level, between Swedish and international *E. coli* ST131. Bayesian inference of Swedish *E. coli* ST131 isolates ( $n = 29$ ), sequenced using PacBio RSII, together with an international ST131 dataset showed that the Swedish isolates were part of the international ST131 A, C1 and C2 clades. Highly conserved plasmids were identified in three clusters although they were separated by several years, which indicates a strong co-evolution between some ST131 lineages and specific plasmids. In conclusion, the tight clonal relationship observed within the ST131 clades, together with highly conserved plasmids, challenges investigation of strain transmission events. A combination of few SNPs on a genome-wide scale and an epidemiological temporospatial link, are needed to track the spread of the ST131 subclones.

Few antibiotic resistant clones have generated as much interest as *Escherichia coli* ST131. In recent years several papers have described its emergence, evolution, and molecular epidemiology globally<sup>1–9</sup>. The interest is justified since ST131 and especially its successful subclone C2/H30-Rx (hereafter referred to as C2) is the dominating multidrug-resistant (MDR) *E. coli* lineage in many parts of the world and is overrepresented among human infections. Its rapid dissemination has been suggested to be partially due to the advantageous combination of fluoroquinolone resistance as well as IncF plasmids encoding the extended-spectrum  $\beta$ -lactamase (ESBL) -gene *bla*<sub>CTX-M-15</sub><sup>2,10–12</sup>. The frequent association with IncF-plasmids encoding *bla*<sub>CTX-M-15</sub> adjacent to an active ISEcp1 transposase suggests that this clone also serves as a main driver of the spread of the *bla*<sub>CTX-M-15</sub> gene to other strains and species<sup>12</sup>.

The emergence of the successful ST131 C2 clade, carrying fluoroquinolone resistance and the *bla*<sub>CTX-M-15</sub> gene, started in the mid-1980s and was followed by a rapid international clonal expansion during the 1990s<sup>1–3,7</sup>. However, it was not until 2008, when MLST was broadly introduced, that the clone was simultaneously described in several places<sup>4,5,13</sup>. Since then the emergence of the ST131 subclone C1/H30-R (referred to as C1) with *bla*<sub>CTX-M-27</sub> has also been described<sup>14,15</sup>. Most genetic variation within the ST131 group is caused by recombinational events including traits such as virulence, antibiotic resistance and other accessory traits<sup>1,3,8</sup>. Studies on the

<sup>1</sup>Division of Clinical Microbiology, Department of Laboratory Medicine, Karolinska Institutet, Alfred Nobels allé 10, 141 52, Huddinge, Stockholm, Sweden. <sup>2</sup>Public Health Agency of Sweden, Nobels väg 18, 17182, Solna, Stockholm, Sweden. <sup>3</sup>Department of Medical Biochemistry and Microbiology, Uppsala University, BMC, Box 582, Husargatan 3, 75123, Uppsala, Sweden. <sup>4</sup>Department of Pathology, University of Florida. Emerging Pathogens Institute, University of Florida, P.O. Box 100009, Gainesville, Florida, 32610-0009, USA. Correspondence and requests for materials should be addressed to S.N. (email: [sofia.ny@folkhalsomyndigheten.se](mailto:sofia.ny@folkhalsomyndigheten.se))

arrangement of resistance genes in ST131 isolates could help to better understand the spread, selection due to resistance and the evolution of the accessory genome in this internationally circulating clone.

No ST131 isolates from Sweden have thus far been part of global phylogenetic studies. Sweden offers a different perspective, with low antibiotic pressure and low prevalence of resistant bacteria in animals and humans, compared to other countries represented in previous studies<sup>16</sup>. International travel to high prevalence regions is a known risk factor for being a community carrier of ESBL-producing *E. coli* in Sweden<sup>17–19</sup>. The ST131 subclone C2 was estimated to have a prevalence of 0.3% in Swedish community carriers in a nationwide study<sup>19</sup>.

We sought to elucidate if the Swedish ST131 population is mostly shaped by constant new influx of international isolates or if it is dominated by national spread from a single introduction. Given the low prevalence of ESBL-producing and low antibiotic selection pressure in Sweden, we hypothesised that ST131 isolates from Swedish carriers and patients would cluster with international clones in several different clades, if observed cases are mainly due to constant influx. This hypothesis was investigated by constructing a phylogeny of ESBL-producing ST131 isolates from i) patients with bloodstream infections and ii) non-symptomatic community carriers in Sweden in relation to previously published international ST131 isolates<sup>3</sup>. Plasmid similarities within clusters were used to study if there were specific plasmids associated with closely related isolates, which would provide an additional level of resolution during outbreaks. Swedish ST131 isolates (n = 29) were sequenced using long-read single molecule real-time sequencing (PacBio), allowing assembly of both chromosomes and plasmids. Bayesian phylogenetic inference was used to investigate the evolutionary patterns of Swedish ST131 isolates and their relationship to previously published genomes (n = 91) from strains circulating internationally<sup>3</sup>.

## Results

**Dataset generation and assembly.** The selected ESBL-producing *E. coli* ST131 isolates were from patients with bloodstream infection (n = 20) and from non-symptomatic community carriers (n = 9) in Sweden (Supplementary Table S1, deposited in NCBI under BioSample accession SAMN10839615 to 43). The isolates were selected based on heterogeneity in their phenotypic resistance pattern, *bla*<sub>CTX-M</sub> gene type and plasmid replicon type to achieve a broad representation of the ST131 population in Sweden.

The published ST131 dataset consisted of a mix between non-ESBL and ESBL-producing isolates from 6 countries on 4 continents<sup>3</sup>. SNP calling, on 120 isolates, was made on the 79% shared dynamic core genome with the reference strain (EC958), and the output alignment contained 13,351 SNPs. After removing recombinant regions 3,297 non-recombinant SNPs remained for Bayesian phylogenetic analysis

**Bayesian molecular clock calibration and mutation rate.** The correlation coefficient for the root-to-tip genetic divergence versus time in the TempEst analysis ( $R^2 = 0.25$ ) indicated a weak linear relationship between accumulated mutations and sampling time, but suggested that enough signal may be present to calibrate a relaxed molecular clock. Bayesian model comparison through Bayes factors confirmed the uncorrelated relaxed clock model and with Bayesian Skyline tree prior was the best fitting for the alignment (Supplementary Table S2). The median molecular clock rate was estimated to  $8.5 * 10^{-4}$  (95% HPD interval  $6.6 * 10^{-4}$  to  $1.03 * 10^{-3}$ ) mutations per site per year which translates to approximately 2.7 mutations per year per genome (Supplementary Table S2), with a median coefficient of variation for the clock rate of 0.41 (95% HPD interval 0.30 to 0.54).

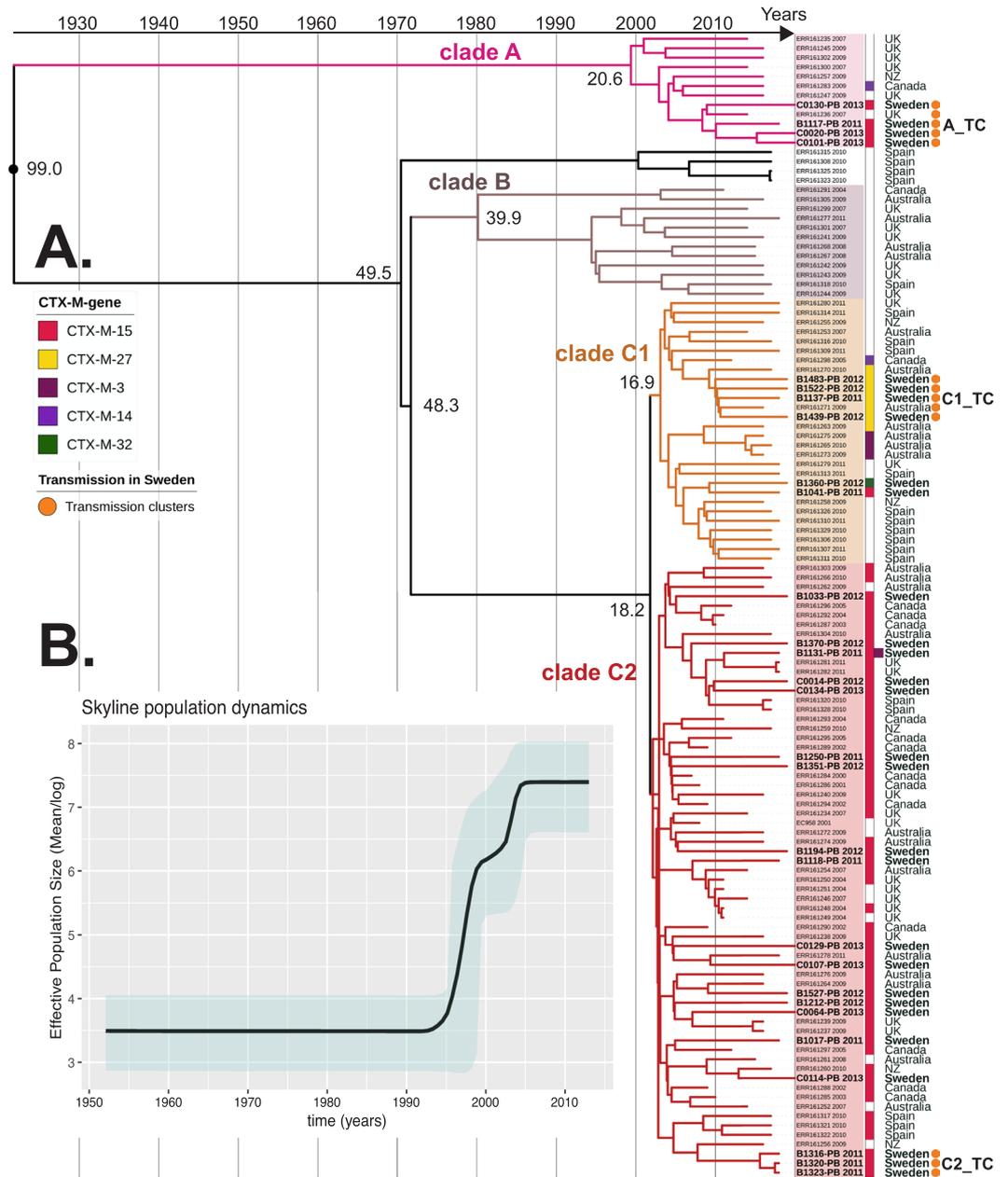
**Several introductions have shaped the Swedish ST131 population.** The combined phylogeny of the international and Swedish ST131 showed a scattered distribution of the Swedish isolates in clades A, C1 and C2 (Fig. 1). No isolates from Sweden belonged to clade B. A fifth clade, with isolates from Spain, which separated from the B/C clade in the beginning of the 1960s was also present in the tree branch, albeit with low statistical support (0.4 posterior probability). All other major clade separations had a branch support of 1.0. The major clade separations between A and B was suggested to have occurred in 1914 (95% HPD 1866 to 1952), between B and C 1964 (95% HPD 1950 to 1979) and between C1 and C2 in 1994 (95% HPD 1991 to 1997) (Fig. 1). The skyline plot showed a sharp increase in the estimated effective population size during the 1990s which levelled out after 2005 (Fig. 1).

The distribution of the Swedish isolates showed three suspected national transmission clusters, where Swedish isolates clustered together under a Most Recent Common Ancestor (MRCA), located in the clades A (five isolates), C1 (five isolates) and C2 (three isolates), marked with circles in Fig. 1.

The number of pairwise chromosomal SNPs separating the transmission cluster in clade A (A\_TC) were 18 to 47, in clade C1 (C1\_TC) 27 to 42 and in clade C2 (C2\_TC) 2 to 12 (Fig. 2, Supplementary Table S3). The median time estimations to MRCA was 12.0 years for the A\_TC, 10.0 years for the C1\_TC and 4.5 years for the C2\_TC. The five isolates in clade A\_TC were isolated at three different locations in Sweden; Stockholm, Malmö and Gothenburg and one was isolated in the UK. The five isolates in the C1\_TC, carrying *bla*<sub>CTX-M-27</sub>, were isolated at three different places in Sweden that are distantly separated geographically (Umeå, Stockholm and Lund) and one isolate was part of the international dataset and isolated in Australia. All isolates in the C2\_TC were isolated in the Swedish city Malmö (Fig. 2).

**Conserved plasmids identified in all suspected transmission clusters.** The three suspected Swedish transmission clusters, in clade A, clade C1 and clade C2, had similar arrangements of their resistance genes and a common MRCA close in time (<12 years) within each cluster (Figs 2 and 3, Supplementary Tables S5–S7).

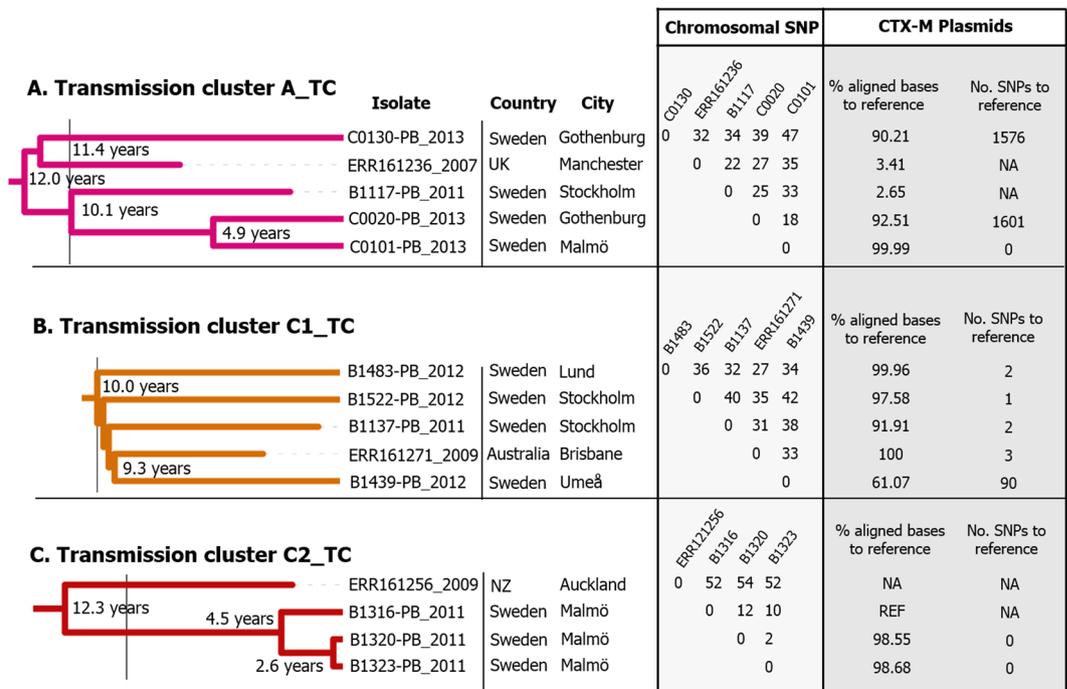
In clade A\_TC all four Swedish isolates included a plasmid with a single *bla*<sub>CTX-M-15</sub> gene (all around 127 kb) while the additional resistance genes were located on a second plasmid (Fig. 3, Fig. S1 Supplementary Tables S4 and S5). Alignment and SNP calling from shared regions displayed above 90% similarity between three of the clade A\_TC plasmids and a plasmid isolated in 2006 or 2007 in Germany (Fig. 2, reference plasmid HG530657.1 not included in the figure)<sup>20</sup>. One isolate from a Swedish community carrier (C0101-PB\_2013) only lacked 4



**Figure 1.** (A) BEAST ST131 phylogenetic tree with Swedish and international isolates, (n = 121) constructed from an alignment of 3,297 non-recombinant SNPs. (A) The ST131 tree consists of five clades A, B, C1 and C2 including an additional clade separation before the (A,C) clades. Visualized metadata includes *bla*<sub>CTX-M</sub> resistance gene (marked with coloured squares) and country of isolation. Three suspected transmission clusters clade A<sub>TC</sub>, clade C1<sub>TC</sub> and clade C2<sub>TC</sub> are marked with orange circles. Swedish isolates starting with B were isolated from patients with bloodstream infections while isolates starting with (C) were from community carriers. (B) Bayesian skyline plot of estimated changes in the effective population size (*N<sub>e</sub>*) over time with 95% CI displayed in blue.

nucleotides compared to the reference plasmid and had 0 SNPs in shared regions. Two Swedish plasmids in the A<sub>TC</sub> cluster differed from the reference with around 1600 SNPs respectively which were distributed over the entire plasmid (Fig. 2, Supplementary Table S5 and Fig. S2). Isolate B1117-PB\_2011 and the UK isolate in clade A<sub>TC</sub> contained completely different plasmids and the B1117-PB\_2011 only shared the region containing the *bla*<sub>CTX-M-15</sub> gene with the other plasmids in the cluster. Regarding the plasmid in isolate B1117-PB\_2011, NCBI BLAST identified two very similar plasmids, CP031903.1 and LT985295.1, both displaying 99% query coverage with 99.94% identities.

The four Swedish isolates in clade C1<sub>TC</sub> had plasmids with the same resistance gene arrangement, except for one strain (B1137-PB\_2011) that had lost part of the plasmid containing the resistance genes *mph(A)*, *sul2*, *strA* and *strB* (Fig. 3). Alignment and SNP calling from shared regions for the clade C1<sub>TC</sub> plasmids (97 to 150 kb in



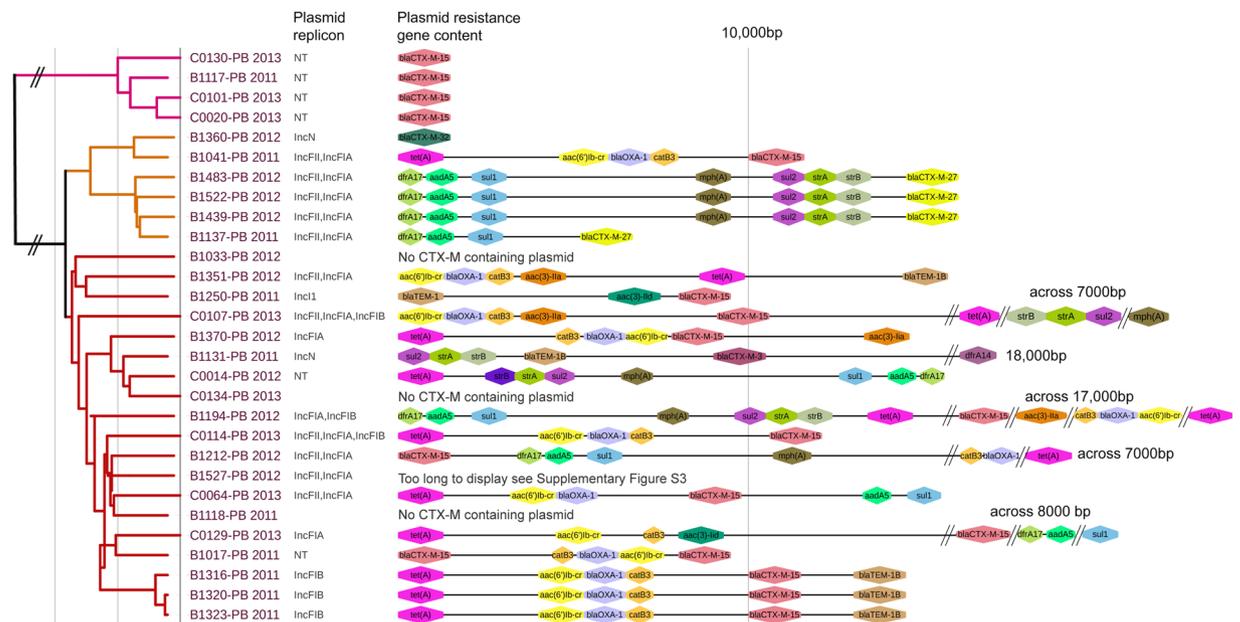
**Figure 2.** Chromosomal and plasmid differences in identified Swedish transmission clusters, in clade A (A\_TC), C1 (C1\_TC) and C2 (C2\_TC). Data included for each strain are; place of isolation (country and city), pairwise SNP differences on chromosomal level (chromosomal SNP), percentage aligned bases between *bla*<sub>CTX-M</sub> plasmids and a reference plasmid (% aligned bases to reference) and *bla*<sub>CTX-M</sub> plasmid SNP differences to reference plasmid (No. SNPs to reference). Reference plasmid for clade A\_TC GenBank: HG530657.1 (size:111,741 bp) and clade C1\_TC GenBank: CP021871.1 (size 134,899 bp). For C2\_TC an internal reference plasmid (B1316-PB\_2011, size: 101,922 bp) was used. NA = non applicable, UK = United Kingdom, NZ = New Zealand.

size) to a reference plasmid isolated in Germany in 2010 showed between 61 and 99% aligned bases (Fig. 2 (reference plasmid CP021871.1 not included in the figure), Supplementary Table S6)<sup>21</sup>. The plasmids differed from each other (including the reference plasmid) with 1 to 90 SNPs. The most diverse plasmid (B1439-PB\_2012) differed with approximately 90 SNPs that were all situated in a 3594 bp recombination event (Fig. 2, Supplementary Table S6 and Fig. S2). The B1439-PB\_2012 plasmid was smaller compared to the reference plasmid (87,000 bp vs 144,000 bp) and uncalled regions in the alignment corresponded to genes in the reference plasmid encoding different transposases, conjugative transfer systems and helicases.

The transmission clade C2\_TC had the same arrangement of resistance genes on their plasmids (98 to 102 kb) (Fig. 3). No closely related plasmid was identified in the NCBI search (all less than 50% query coverage) so B1316-PB\_2011 was used as reference. All plasmids had above 98% aligned bases and the SNP calling showed no variation in shared regions. Seven deletions were identified accounting for 244 to 1476 bp in total (Supplementary Table S7).

**Heterogeneous arrangement of resistance genes between plasmids and chromosome.** The distribution of *bla*<sub>CTX-M</sub> genes in the Swedish dataset showed that *bla*<sub>CTX-M-15</sub> was common in clade C2 but also appeared in the A and C1 clades while the *bla*<sub>CTX-M-27</sub> exclusively appeared in C1 (Fig. 1). The *fimH* types for the Swedish isolates were *fimH41* in clade A and *fimH30* in the rest of the isolates. The organisation of resistance genes was overall heterogeneous in their distribution between plasmids and chromosomal regions both within and between isolates. The majority of the resistance genes were located in the *bla*<sub>CTX-M</sub> encoding plasmid, except for two isolates in clade A (B1117-PB\_2011 and C0101-PB\_2013) (Fig. 3 and Supplementary Fig. S1). Eight isolates had more than one resistance plasmid (Fig. S1). The number of unique resistance genes were between 1 and 14 per isolate. One plasmid (B1527-PB\_2012) included five copies of the same resistance cassette (confirmed with a combination of long-read PacBio sequencing and coverage analysis using IonTorrent short-read data) Supplementary Fig. S3.

In total eight isolates had resistance genes located in the chromosome and two isolates had only chromosomally encoded resistance genes (C1118-PB\_2011, C0134-PB\_2013) (Fig. S4). One isolate (B1017-PB\_2011) had its resistance genes distributed over two plasmids and the chromosome. All chromosomally located ESBL-resistance was due to *bla*<sub>CTX-M-15</sub> insertions and in four cases it was the only resistance gene found on the chromosome (Fig. S4).



**Figure 3.** Comparison of resistance gene placement on *bla*<sub>CTX-M</sub> encoding plasmids among Swedish ST131 isolates (n = 29). For five of the isolates the sequences were cut to accommodate all resistance genes, the actual length of the region harbouring these genes is written next to the sequence. Isolates starting with B were from patients with bloodstream infections while isolates starting with C were from community carriers. NT = Non-Typable.

## Discussion

We describe the relationship between the Swedish ESBL-producing *E. coli* ST131 population and the international ST131 lineage. The results demonstrate that both national transmission clusters and introductions of international clones contribute to the ST131 population structure in Sweden. We also show close relationship between international clonal clusters, chromosomes and plasmids, in the ST131 A and C1 clades, without them being part of a national transmission chain. This finding illustrates the challenges when investigating outbreaks with conserved clonal lineages, as few methods can provide appropriate resolution to draw conclusions on strain transmission based on molecular typing data.

**Evolutionary history and clade separation in line with previous publications.** The presented ST131 tree (Fig. 1) including Swedish and international isolates showed the same topological distribution of clades as previously published phylogenetic trees<sup>1–3,7,9</sup>. Our results are also in line with a previous BEAST analysis favouring the same substitution model, clock model and tree prior<sup>1</sup>. Since we used an external dataset that was also included in the study by Zakour *et al.*, it is not surprising that our results are similar<sup>1,3</sup>. We saw slightly different values for time point clade separation (C1/C2 clade 1994 in our study versus 1986 for Zakour *et al.*) and mutation rate (SNP/site/year of  $8.4 \times 10^{-4}$  in our study vs  $4.39 \times 10^{-7}$  for Zakour *et al.*). These values are naturally connected since a higher mutation rate (our study) leads to more variations in a shorter time period and therefore a later divergence date. The large difference in mutation rate between the studies could be due to different methods in removing recombination, however very similar SNP alignment length after removing recombination suggests otherwise (3,297 bp in our study vs 3,779 bp for Zakour *et al.*). The mutation rate estimation, in our study, translates to approximately 2.7 fixated SNPs/year/genome which means that around 5 SNPs can be expected to differ between two isolates sharing a MRCA one year back in time. The number (2.7 SNPs) could be perceived as high, but importantly this is the number of SNPs in the entire genome (5 Mb), excluding recombinant SNPs and therefore an estimate to be used with caution. The estimated changes in effective population size (Skyline plot Fig. 1) show similar variations over time as previously suggested<sup>1</sup>. Considering that the first report of *bla*<sub>CTX-M-15</sub> was published in 2001 it is fascinating to see that the population increase was already ongoing at that time<sup>22</sup>. It also likely illustrates the powerful driving force of fluoroquinolones like ciprofloxacin that was introduced on the market 1987 and has been described by many as one of the most important factors for the success of the ST131 C2 lineage<sup>1,2,23</sup>. The skyline model suggests that a plateau was reached in 2005, after which the ST131 population entered a steady state where the population did neither decline nor increase (Fig. 1). This finding is in agreement with Swedish surveillance data, where the proportion of ST131 in ESBL-producing *E. coli* urinary tract infection has been stable since surveillance started in 2007<sup>24</sup>.

The observations herein suggest that a constant influx of ST131 has taken place to the Swedish ST131 population, since the Swedish isolates were mixed with the international sublineages in the phylogeny (Fig. 1). We also observed local national transmission clusters. One important aspect is that our dataset was not designed to quantify the burden of international import versus national spread of ESBL-producing ST131 in Sweden. To do

this a larger and differently selected dataset would be needed. National transmission could therefore be the largest contributor to infections with ESBL-producing *E. coli* ST131 in Sweden.

**High stability of plasmids within the major observed transmission clusters.** Because of the clonal nature of the isolates in the different ST131 clades it can be challenging to determine which isolates that have spread locally versus international influx. This was seen in both the A and C1 clade where strains clustering together was isolated at different places in Sweden and the world (UK and Australia) and therefore represent a tight clonal lineage rather than a direct transmission chain. Including more data from the international ST131 population could potentially have increased resolution within transmission cluster A and C1 which would be interesting clusters to examine in future analyses.

In clade A the cluster containing five isolates had limited chromosomal variation and three of the isolates had very similar plasmids (Fig. 2). These three Swedish plasmids showed high resemblance to a plasmid isolated from a patient in Germany in 2006 or 2007 (HG530657.1) that we used as reference for the plasmid SNP call (over 90% of sequence shared)<sup>20</sup>. One plasmid turned out to be almost identical to the German reference sharing 99.99% of its sequence and containing 0 SNPs (Clade A Fig. 2). Considering that the Swedish isolate was from a community carrier and isolated 5–6 years after the German isolate, it shows remarkable stability. The other two plasmids in the cluster had accumulated SNPs over time (around 1,600 SNPs each) but still shared over 90% of the sequence, and the variations seen were likely not due to recombination since an even SNP accumulation was seen across the genome (Supplementary Fig. S2). The Swedish B1117-PB\_2011 isolate had completely different plasmids that only shared the region containing the *bla*<sub>CTX-M-15</sub> gene itself with the other plasmids (Fig. 2). This gives insight into how unpredictable the plasmid composition of closely related strains can be. One likely explanation is that the highly mobile cassette harbouring the ISEcpI and *bla*<sub>CTX-M-15</sub> genes has moved between plasmids. This is also a likely explanation to why we saw high presence of *bla*<sub>CTX-M-15</sub> with chromosomal location in our Swedish dataset. Since one of the isolates in the cluster (Fig. 2) was isolated in the UK, and the Swedish isolates were from very different geographical regions, this clade A clonal lineage is likely circulating internationally among community carriers but also causes clinical infections.

The BEAST analysis showed that the C1\_TC had a MRCA around 10 years back in time (Figs 1 and 2). Since all isolates in this cluster, including one isolate from Australia, carried the same conserved plasmid, with only 1–3 SNPs or one recombination event, it must have been present already in the MRCA 10 years ago (Fig. 2, Supplementary Table S6)<sup>14,15,21</sup>. In addition, the Australian isolate just outside the cluster (ERR161270\_2010 Fig. 1) also carried the same conserved plasmid (Supplementary Table S6). The conserved plasmid together with the limited chromosomal SNP differences (around 30–40 SNPs) in the C1\_TC cluster indicate a strong co-evolution between the plasmid lineage carrying *bla*<sub>CTX-M-27</sub> and the specific C1 ST131 clade as has been described previously<sup>14,15</sup>. A dataset including more representatives of the C1 *bla*<sub>CTX-M-27</sub> clone might have given a better resolution and helped to separate the Swedish isolates into several clades.

The three Swedish C2\_TC isolates only differed by 2, 10 and 12 SNPs on a chromosomal level and no SNPs were seen in the 98% shared plasmid sequence. All three isolates came from the same Swedish city and likely represents local spread. The C2\_TC was the only transmission cluster identified as likely to be local for Sweden in this dataset.

## Conclusions

The evidence presented herein offers a deepened insight into the Swedish epidemiology of ESBL-producing ST131 and its close relationship to internationally disseminating ST131 clones. It is evident that the Swedish ST131 population is part of the international lineages and that several introductions combined with national transmission have formed the strain population. The clonal nature of the ST131 lineage, with highly conserved plasmids in some sub-lineages, complicates estimation of local circulation and transmission, and highlights the importance of temporospatial epidemiological links even in the genomic era. However, very close genetic relationships (a few chromosomal SNPs) could indicate a direct transmission even if the epidemiological link is unknown. Such small differences can only be detected by whole genome sequencing and not with traditional typing methods. Online tools and databases to analyse WGS data adapted for clinical microbiologists and public health workers could provide the possibility to detect outbreaks with clonal lineages like ST131, two examples are BacWGSTdb and Enterobase<sup>25,26</sup>. The presence of plasmids, which were highly conserved over many years, in a widely disseminated clonal lineage illustrates that identical plasmids sequences are not a certain evidence of plasmid transmission. If a suspected plasmid transmission takes place in a hospital and identical plasmids are identified it still might not be a direct spread if the plasmid originates from one of these conserved lineages. Therefore, despite the emergence of sequencing technology that allows for rapid plasmid sequencing in the clinical setting, more work is still to understand the role of conserved plasmids in certain conserved clonal lineages, and how this phenomenon can impact inference about plasmid transmission events.

## Material and Methods

**Swedish clinical isolates.** A total of 29 strains of *E. coli* ST131 isolates from patients with bloodstream infections (n = 20) and from non-symptomatic community carriers (n = 9) were selected from a larger dataset of ST131 (n = 177) from bloodstream infections and community carriers. The data collection and molecular typing was previously described<sup>19</sup>.

The isolates were selected from a Swedish nationwide collection of *E. coli* ST131 bloodstream isolates (n = 177) producing ESBL (2011–2012) and from a point prevalence study on community carriage (n = 16) of ESBL-producing *E. coli* (2013)<sup>19</sup>. We selected ST131 isolates from the community carriers and a subset of bloodstream isolates based on C2 status, phenotypic resistance profiles, *bla*<sub>CTX-M</sub> types, and plasmid replicon type<sup>2,19</sup>. Based on these criteria we included isolates to represent the diversity, as well as representative isolates from two

possible transmission clusters. The selected isolates had the ESBL genes *bla*<sub>CTX-M-15</sub> (n = 23), *bla*<sub>CTX-M-27</sub> (n = 4), *bla*<sub>CTX-M-3</sub> (n = 1), *bla*<sub>CTX-M-32</sub> (n = 1). The majority of isolates had different IncF plasmids (n = 25) (Supplementary Table S1).

**Sequencing and external international dataset.** Genomic DNA was extracted using Qiagen genomic tip 500/G kit and sequenced with PacBio RSII (Pacific Biosciences) using one SMRT<sup>™</sup> cell per isolate. For detailed data on the dataset see Supplementary Table S4. Hierarchical Genome Assembly Process (HGAP) was used to generate draft assemblies/genomes. Plasmids identified in transmission clusters in clade A, C1 and C2 were further closed using Flye assembler version 2.4.2 <https://github.com/fenderglass/Flye><sup>27,28</sup>. Assembled genomes were deposited at NCBI under BioProject PRJNA517648 with BioSample accessions SAMN10839615 to 43. Isolate B1527-PB\_2012 was also sequenced using Ion Torrent S5 XL (Thermo Fischer Scientific). A publicly available WGS Illumina dataset (FASTQ-files) of 91 *E. coli* ST131 was used as comparison and reference for the different subclones within ST131<sup>3</sup>. *FimH*, plasmid replicon and resistance genes were identified via the CGE website using FimType, PlasmidFinder and ResFinder respectively (<https://cge.cbs.dtu.dk/services/>) accessed on 2019-03-05<sup>29</sup>. All annotations were made in CLC Genomic Workbench v8.0.1 (Qiagen Bioinformatics, Redwood City, USA).

**SNP-call, alignments and recombination removal.** Variation calling, mapping and, *de novo* assembly of the 121 genomes was done using CLC Assembly Cell Version: 4.4.2.133896 (Qiagen Bioinformatics). Isolate EC958 was used as reference for the chromosome SNP call because of it being a well sequenced (long and short read), assembled and closed genome from the C2 clade<sup>30</sup>. For the external international dataset the deposited FASTQ files were used while for the Swedish PacBio dataset the HGAP assembled genomes in the SNP-call were used. Minimum coverage, for FASTQ, was set to 10 and the SNP ratio cut-off for SNP support was set to 90%.

Alignment and SNP call of suspected transmission cluster plasmids was done with Snippy v4.0 (<https://github.com/tseemann/snippy> accessed on 2019-03-05). SNP distance matrices were calculated using snp-dist v0.6 (<https://github.com/tseemann/snp-dists> accessed on 2019-03-05). Two reference plasmids were used for SNP-calling in transmission cluster A (HG530657.1, query coverage: 98%, identity 99.98%) and C1 (CP021871.1 query coverage: 100%, identity 99.99%) and identified by using NCBI nucleotide BLAST using the assembled plasmids as queries with the cut-offs >95% query coverage and >99% identity. For the transmission cluster in C2 an internal reference was used (B1316-PB\_2011) since the closest match at NCBI (CP036244.1) had only 67% query coverage. Re-assembly with IonTorrent data of the resistance plasmid from isolate B1527-PB\_2012 against the plasmid assembled from the PacBio data, was done using CLC Genomic Workbench v8.0.1 (Qiagen Bioinformatics). Gubbins version 2.3.1, with standard settings, was used to remove SNPs from recombinant regions<sup>31</sup>.

**Bayesian inference.** The presence of enough phylogenetic signal in the ST131 aligned sequences was investigated by likelihood mapping analysis using IQ-TREE (<http://www.iqtree.org/> accessed on 2019-03-05) and allowing the software to search for all possible quartets using the best-fitting nucleotide substitution model<sup>32</sup>. Temporal signal was assessed by plotting root-to-tip divergence versus sampling time with TempEst v1.5 (<http://tree.bio.ed.ac.uk/software/tempest/>) of the maximum likelihood (ML) tree including all aligned sequences<sup>33</sup>. ML tree was inferred with RaxML, using the best-fitting nucleotide substitution model bootstrapping (1,000 replicates) to assess statistical robustness for internal branching order in the phylogeny<sup>34</sup>. Bayesian phylogenetic inference was carried out with BEAST v1.10<sup>35</sup>. Path- and stepping stone sampling was used for the marginal likelihood estimation<sup>36,37</sup>. A Markov Chain Monte Carlo (MCMC) was run for 200 million generations, sampling every 20 million generations. Analysis files (xml) were generated in BEAUti using: collection year tip dating, site heterogeneity gamma model with 4 categories and ascertainment bias correction. Both strict and uncorrelated relaxed clock was tested with the following tree priors: constant size, exponential growth, GMRF Bayesian Skyride and Bayesian skyline with 3 groups. The substitution models HKY and GTR were also tested. Proper sampling of the Markov chain (Xml files are available upon request) was evaluated by calculating the effective sampling size (ESS) with Tracer 1.7; ESS values >200 for parameter estimates were considered acceptable<sup>38</sup>. Bayes factor (BF) was used in the estimation of best model fit for the dataset. TreeAnnotator was used to create Maximum Clade Credibility (MCC) trees with 10% burn in. Visualisation of trees and metadata was done using iTOL (<https://itol.embl.de/> accessed on 2019-03-05)<sup>39</sup>. Bayesian skyline plot was generated in Tracer v1.6 and visualized in R v3.5.1. The posterior distribution of trees was summarized into the maximum clade credibility (MCC) tree with TreeAnnotator v1.8.4 (BEAST package) after 10% burn-in and the final tree edited with FigTree (<http://tree.bio.ed.ac.uk/software/figtree/> accessed on 2019-03-05). All computations were performed on resources provided by the Swedish National Infrastructure for Computing (SNIC) at Uppsala Multidisciplinary Center for Advanced Computational Science (UPPMAX).

**Ethical approval and informed consent.** This study was conducted in accordance with the Declaration of Helsinki and national and institutional standards and was approved by the research ethics committee “Regionala Etikprövningsnämnden i Stockholm, EPN” in Stockholm, Sweden (Record: 2012/1204-31/4). Written informed consent was obtained from all participants contributing with samples.

### Data Availability

Additional data and analysis files produced during this study are available upon request. Sequence data were deposited at NCBI under BioSample accession SAMN10839615 to 43.

## References

- Zakour, B. N. L. *et al.* Sequential Acquisition of Virulence and Fluoroquinolone Resistance Has Shaped the Evolution of *Escherichia coli* ST131. *mBio* **7**, <https://doi.org/10.1128/mBio.00347-16> (2016).
- Price, L. B. *et al.* The epidemic of extended-spectrum-beta-lactamase-producing *Escherichia coli* ST131 is driven by a single highly pathogenic subclone, H30-Rx. *mBio* **4**, e00377–00313, <https://doi.org/10.1128/mBio.00377-13> (2013).
- Petty, N. K. *et al.* Global dissemination of a multidrug resistant *Escherichia coli* clone. *Proceedings of the National Academy of Sciences* **111**, 5694–5699, <https://doi.org/10.1073/pnas.1322678111> (2014).
- Nicolas-Chanoine, M. H., Bertrand, X. & Madec, J. Y. *Escherichia coli* ST131, an intriguing clonal group. *Clin Microbiol Rev* **27**, 543–574, <https://doi.org/10.1128/CMR.00125-13> (2014).
- Johnson, J. R. *et al.* Abrupt emergence of a single dominant multidrug-resistant strain of *Escherichia coli*. *The Journal of infectious diseases* **207**, 919–928, <https://doi.org/10.1093/infdis/jis933> (2013).
- Mathers, A. J., Peirano, G. & Pitout, J. D. D. *Escherichia coli* ST131: The Quintessential Example of an International Multiresistant High-Risk Clone. **90**, 109–154, <https://doi.org/10.1016/bs.aamb.2014.09.002> (2015).
- Stoesser, N. *et al.* Evolutionary History of the Global Emergence of the *Escherichia coli* Epidemic Clone ST131. *mBio* **7**, <https://doi.org/10.1128/mBio.02162-15> (2016).
- Downing, T. T. D. Resistant Infection Outbreaks of Global Pandemic *Escherichia coli* ST131 Using Evolutionary and Epidemiological Genomics. *Microorganisms* **3**, 236–267, <https://doi.org/10.3390/microorganisms3020236> (2015).
- McNally, A. *et al.* Combined Analysis of Variation in Core, Accessory and Regulatory Genome Regions Provides a Super-Resolution View into the Evolution of Bacterial Populations. *PLoS genetics* **12**, e1006280, <https://doi.org/10.1371/journal.pgen.1006280> (2016).
- Shaik, S. *et al.* Comparative Genomic Analysis of Globally Dominant ST131 Clone with Other Epidemiologically Successful Extraintestinal Pathogenic *Escherichia coli* (ExPEC) Lineages. *mBio* **8**, <https://doi.org/10.1128/mBio.01596-17> (2017).
- Pitout, J. D. & DeVinney, R. *Escherichia coli* ST131: a multidrug-resistant clone primed for global domination. *F1000Res* **6**, 1–7 (2017).
- Mathers, A. J., Peirano, G. & Pitout, J. D. The role of epidemic resistance plasmids and international high-risk clones in the spread of multidrug-resistant Enterobacteriaceae. *Clin Microbiol Rev* **28**, 565–591, <https://doi.org/10.1128/CMR.00116-14> (2015).
- Nicolas-Chanoine, M. H. *et al.* Intercontinental emergence of *Escherichia coli* clone O25:H4-ST131 producing CTX-M-15. *The Journal of antimicrobial chemotherapy* **61**, 273–281, <https://doi.org/10.1093/jac/dkm464> (2008).
- Matsumura, Y. *et al.* CTX-M-27- and CTX-M-14-producing, ciprofloxacin-resistant *Escherichia coli* of the H30 subclonal group within ST131 drive a Japanese regional ESBL epidemic. *The Journal of antimicrobial chemotherapy* **70**, 1639–1649, <https://doi.org/10.1093/jac/dkv017> (2015).
- Matsumura, Y. *et al.* Global *Escherichia coli* Sequence Type 131 Clade with blaCTX-M-27 Gene. *Emerging infectious diseases* **22**, 1900–1907, <https://doi.org/10.3201/eid2211.160519> (2016).
- Sweden, P. H. A. o. Swedish work on containment of antibiotic resistance Tools, methods and experiences. (2014).
- Tängdén, T., Cars, O., Melhus, Å. & Löwdin, E. Foreign Travel Is a Major Risk Factor for Colonization with *Escherichia coli* Producing CTX-M-Type Extended-Spectrum  $\beta$ -Lactamases: a Prospective Study with Swedish Volunteers. *Antimicrobial Agents and Chemotherapy* **54**, 3564–3568, <https://doi.org/10.1128/aac.00220-10> (2010).
- Vading, M. *et al.* Frequent acquisition of low-virulence strains of ESBL-producing *Escherichia coli* in travellers. *The Journal of antimicrobial chemotherapy* **71**, 3548–3555, <https://doi.org/10.1093/jac/dkw335> (2016).
- Ny, S. *et al.* Community carriage of ESBL-producing *Escherichia coli* is associated with strains of low pathogenicity: a Swedish nationwide study. *The Journal of antimicrobial chemotherapy* **72**, 582–588, <https://doi.org/10.1093/jac/dkw419> (2017).
- Falgenhauer, L. *et al.* Complete Genome Sequence of Phage-Like Plasmid pECOH89, Encoding CTX-M-15. *Genome announcements* **2**, <https://doi.org/10.1128/genomeA.00356-14> (2014).
- Hiren Ghosh, B. B. *et al.* Complete Genome Sequence of blaCTX-M-27-Encoding *Escherichia coli* Strain H105 of Sequence Type 131 Lineage C1/H30R. *Genome announcements* **5**(31), e00736–17 (2017).
- Karim A, P. L., Nagarajan, S. & Nordmann, P. Plasmid-mediated extended-spectrum beta-lactamase (CTX-M-3 like) from India and gene association with insertion sequence ISEcp1. *FEMS Microbiol Letters* **24**:201(2):237–41 (2001).
- Johnson, J. R., Johnston, B., Clabots, C., Kuskowski, M. A. & Castanheira, M. *Escherichia coli* Sequence Type ST131 as the Major Cause of Serious Multidrug-Resistant *E. coli* Infections in the United States. *Clinical Infectious Diseases* **51**, 286–294, <https://doi.org/10.1086/653932> (2010).
- Brolund, A. *et al.* Epidemiology of extended-spectrum  $\beta$ -lactamase-producing *Escherichia coli* in Sweden 2007–2011. *Clinical Microbiology and Infection* **20**, 1–9, <https://doi.org/10.1111/1469-0691.12413> (2013).
- Alikhan, N.-E., Zhou, Z., Sergeant, M. J. & Achtman, M. A genomic overview of the population structure of *Salmonella*. *PLoS genetics* **14**, e1007261, <https://doi.org/10.1371/journal.pgen.1007261> (2018).
- Ruan, Z. & Feng, Y. BacWGSTdb, a database for genotyping and source tracking bacterial pathogens. *Nucleic acids research* **44**, D682–D687, <https://doi.org/10.1093/nar/gkv1004> (2016).
- Kolmogorov, M., Yuan, J., Lin, Y. & Pevzner, P. A. Assembly of Long Error-Prone Reads Using Repeat Graphs. *bioRxiv*, 247148, <https://doi.org/10.1101/247148> (2018).
- Lin, Y. *et al.* Assembly of long error-prone reads using de Bruijn graphs. *Proceedings of the National Academy of Sciences* **113**, E8396, <https://doi.org/10.1073/pnas.1604560113> (2016).
- Zankari, E. *et al.* Identification of acquired antimicrobial resistance genes. *The Journal of antimicrobial chemotherapy* **67**, 2640–2644, <https://doi.org/10.1093/jac/dks261> (2012).
- Forde, B. M. *et al.* The complete genome sequence of *Escherichia coli* EC958: a high quality reference sequence for the globally disseminated multidrug resistant *E. coli* O25b:H4-ST131 clone. *PLoS One* **9**, e104400, <https://doi.org/10.1371/journal.pone.0104400> (2014).
- Croucher, N. J. *et al.* Rapid phylogenetic analysis of large samples of recombinant bacterial whole genome sequences using Gubbins. *Nucleic Acids Res* **43**, e15, <https://doi.org/10.1093/nar/gku1196> (2015).
- Schmidt, H. A., Martin Vingron, K. S. & Haeseler, A. V. TREE-PUZZLE: maximum likelihood phylogenetic analysis using quartets and parallel computing. *Bioinformatics* **18**, 502–504 (2002).
- Rambaut, A., Lam, T. T., Max Carvalho, L. & Pybus, O. G. Exploring the temporal structure of heterochronous sequences using TempEst (formerly Path-O-Gen). *Virus Evol* **2**, vew007, <https://doi.org/10.1093/ve/vew007> (2016).
- Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* **30**, 1312–1313, <https://doi.org/10.1093/bioinformatics/btu033> (2014).
- Suchard, M. A. *et al.* Bayesian phylogenetic and phylodynamic data integration using BEAST 1.10. *Virus Evol* **4**, vey016, <https://doi.org/10.1093/ve/vey016> (2018).
- Baele, G. *et al.* Improving the accuracy of demographic and molecular clock model comparison while accommodating phylogenetic uncertainty. *Mol Biol Evol* **29**, 2157–2167, <https://doi.org/10.1093/molbev/mss084> (2012).
- Baele, G., Li, W. L., Drummond, A. J., Suchard, M. A. & Lemey, P. Accurate model selection of relaxed molecular clocks in bayesian phylogenetics. *Mol Biol Evol* **30**, 239–243, <https://doi.org/10.1093/molbev/mss243> (2013).
- Rambaut, A., Drummond, A. J., Xie, D., Baele, G. & Suchard, M. A. Posterior Summarization in Bayesian Phylogenetics Using Tracer 1.7. *Syst Biol* **67**, 901–904, <https://doi.org/10.1093/sysbio/syy032> (2018).
- Letunic, I. & Bork, P. Interactive tree of life (iTOL) v3: an online tool for the display and annotation of phylogenetic and other trees. *Nucleic Acids Res* **44**, W242–245, <https://doi.org/10.1093/nar/gkw290> (2016).

## Acknowledgements

We would like to acknowledge the help and advice from a few colleagues. Lionel Guy for helpful advice and discussions during the data analysis. A big thank you to Carla Scarlatta and Taylor Paisie for their patient advice and very helpful guidance. Last but not least thank you to all volunteers and patients that provided data to this study.

## Author Contributions

Sofia Ny.: Study design, data management, data analysis, data interpretation, producing Figures and Tables, writing of manuscript. Linus Sandegren: Study design, data interpretation, writing of manuscript. Marco Salemi: Study design, data analysis, data interpretation, writing of manuscript. Christian Giske: Study design, data interpretation, writing of manuscript

## Additional Information

**Supplementary information** accompanies this paper at <https://doi.org/10.1038/s41598-019-46580-3>.

**Competing Interests:** The authors declare no competing interests.

**Publisher's note:** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2019