

RESEARCH ARTICLE

Radical genome remodelling accompanied the emergence of a novel host-restricted bacterial pathogen

Gonzalo Yebra¹, Andreas F. Haag², Maan M. Neamah^{2,3,4}, Bryan A. Wee¹, Emily J. Richardson¹, Pilar Horcajo⁵, Sander Granneman⁶, María Ángeles Tormo-Más^{7,8}, Ricardo de la Fuente⁵, J. Ross Fitzgerald^{1*}, José R. Penadés^{2,7,9*}

1 The Roslin Institute, University of Edinburgh, Edinburgh, United Kingdom, **2** Institute of Infection, Immunity & Inflammation, University of Glasgow, Glasgow, United Kingdom, **3** Faculty of Veterinary Medicine, University of Kufa, Kufa, Iraq, **4** Middle Euphrates Centre for Cancer and Genetic Research, University of Kufa, Kufa, Iraq, **5** Facultad de Veterinaria, Universidad Complutense de Madrid, Madrid, Spain, **6** Centre for Synthetic and Systems Biology, University of Edinburgh, Edinburgh, United Kingdom, **7** Departamento de Ciencias Biomédicas, Facultad de Ciencias de la Salud, Universidad CEU Cardenal Herrera, Valencia, Spain, **8** Severe Infection Group, Health Research Institute Hospital La Fe, Valencia, Spain, **9** MRC Centre for Molecular Bacteriology and Infection, Imperial College London, United Kingdom

* These authors contributed equally to this work.

* ross.fitzgerald@ed.ac.uk (JRF); j.penades@imperial.ac.uk (JRP)



OPEN ACCESS

Citation: Yebra G, Haag AF, Neamah MM, Wee BA, Richardson EJ, Horcajo P, et al. (2021) Radical genome remodelling accompanied the emergence of a novel host-restricted bacterial pathogen. *PLoS Pathog* 17(5): e1009606. <https://doi.org/10.1371/journal.ppat.1009606>

Editor: Andreas Peschel, University of Tubingen, GERMANY

Received: February 24, 2021

Accepted: May 3, 2021

Published: May 20, 2021

Copyright: © 2021 Yebra et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: Illumina read data is available in the European Nucleotide Archive under the study accession number PRJEB30965. The MVF7 full genome assembly was deposited at NCBI with GenBank accession number CP062279.

Funding: This work was supported by the Biotechnology and Biological Sciences Research Council (<https://bbsrc.ukri.org/>) (project grant BB/K00638X/1 and institute strategic grant funding ISP2 BB/P013740/1 to J.R.F.); the Medical Research Council (<https://mrc.ukri.org/>) (grant

Abstract

The emergence of new pathogens is a major threat to public and veterinary health. Changes in bacterial habitat such as a switch in host or disease tropism are typically accompanied by genetic diversification. *Staphylococcus aureus* is a multi-host bacterial species associated with human and livestock infections. A microaerophilic subspecies, *Staphylococcus aureus* subsp. *anaerobius*, is responsible for Morel's disease, a lymphadenitis restricted to sheep and goats. However, the evolutionary history of *S. aureus* subsp. *anaerobius* and its relatedness to *S. aureus* are unknown. Population genomic analyses of clinical *S. aureus* subsp. *anaerobius* isolates revealed a highly conserved clone that descended from a *S. aureus* progenitor about 1000 years ago before differentiating into distinct lineages that contain African and European isolates. *S. aureus* subsp. *anaerobius* has undergone limited clonal expansion, with a restricted population size, and an evolutionary rate 10-fold slower than *S. aureus*. The transition to its current restricted ecological niche involved acquisition of a pathogenicity island encoding a ruminant host-specific effector of abscess formation, large chromosomal re-arrangements, and the accumulation of at least 205 pseudogenes, resulting in a highly fastidious metabolism. Importantly, expansion of ~87 insertion sequences (IS) located largely in intergenic regions provided distinct mechanisms for the control of expression of flanking genes, including a novel mechanism associated with IS-mediated anti-anti-sense decoupling of ancestral gene repression. Our findings reveal the remarkable evolutionary trajectory of a host-restricted bacterial pathogen that resulted from extensive remodelling of the *S. aureus* genome through an array of diverse mechanisms in parallel.

MR/N02995X/1 to J.R.F.); and the Wellcome Trust (<https://wellcome.org/>) (collaborative award 201531/Z/16/Z to J.R.F. and J.R.P.). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

Author summary

The emergence of new pathogens is a major threat to public and veterinary health. Some bacteria such as *Staphylococcus aureus*, have the capacity to infect many different host species including humans and livestock while others such as the closely-related *S. aureus* subsp. *anaerobius*, associated with a single type of pathology called Morel's disease in small ruminants, are highly niche-restricted. However, our understanding of the genetic basis for such differences in bacterial host-tropism is very limited. Here, we discovered that *S. aureus* subsp. *anaerobius* evolved from an *S. aureus* ancestor and underwent an array of extensive changes to its genome that accompanied the transition to its current restricted lifestyle. We observed genome decay involving loss of function of hundreds of genes, large intra-chromosomal rearrangements affecting most of the genome, acquisition of a pathogenicity island, and expansion of large numbers of insertion sequences that are inserted at intergenic sites around the genome. Importantly, we found that IS elements affect the expression of neighbouring genes in different ways including a novel mechanism of IS-enabled disruption of ancestral gene repression. Taken together, we provide a remarkable example of radical genomic changes associated with evolutionary transition from a multi-host to highly restricted host ecology.

Introduction

Bacteria have a remarkable capacity to adapt to new environmental niches, a characteristic that underpins their potential to become successful pathogens. Some pathogens can infect multiple host-species [1], while others become specialized and restricted to a single host. The evolutionary process of host restriction, observed across distantly-related bacterial groups, typically involves genomic events such as gene loss, gene acquisition, or chromosomal rearrangements [2–5]. The impact of such events may be most apparent after long-term associations between bacterium and host, with obligate intracellular pathogens the most extreme examples, exhibiting extremely small and compact genomes [6].

Staphylococcus aureus is a highly versatile bacterial pathogen, associated with an array of diseases in humans and livestock representing a threat to public and livestock health [7]. *S. aureus* has undergone extensive host-switching events during its evolutionary history leading to the emergence of new endemic and epidemic clones in livestock and humans [8]. In particular, the highest number of host-jump events appeared to have been between humans and ruminants in either direction [8]. A combination of gene acquisition, loss of gene function, and allelic diversification have been central to the capacity of *S. aureus* to undergo successful host-adaptation [8].

A subspecies of *S. aureus*, *Staphylococcus aureus* subsp. *anaerobius*, responsible for Morel's disease [9], is highly restricted to small ruminants and its clinical manifestation is confined to abscesses located at major superficial lymph nodes. Outbreaks of this disease are associated with significant economic losses and have been reported in Europe, Africa and the Middle East [10]. In contrast to *S. aureus* subsp. *aureus*, *S. aureus* subsp. *anaerobius* is microaerophilic and catalase-negative [11], but their genetic and evolutionary relatedness is poorly understood.

Here, we investigate the evolutionary history of *S. aureus* subsp. *anaerobius*, revealing that it descended from an ancestor resembling *S. aureus* and adapted to its current ecological niche by extensive genome diversification involving massive genome decay, gene acquisition, and genome rearrangements. Furthermore, acquisition and expansion of ISSau8-like insertion sequences associated with attenuation of transcription of flanking genes through multiple

mechanisms has led to rewiring of the transcriptome. Overall, our results reveal a remarkable example of a bacterial pathogen that has undergone extensive genome remodelling in transition to a highly niche-restricted ecology.

Results

S. aureus subsp. *anaerobius* has a highly conserved genome

Previously, it has been reported that *S. aureus* subsp. *anaerobius* isolates belong to a single clonal complex as determined by multi-locus sequence typing (MLST) [10,12]. To date only one fragmented draft genome of *S. aureus* subsp. *anaerobius* isolated in Sudan has been publicly available, and therefore the genome architecture has been unclear [13]. Here, we applied PacBio sequencing to the type strain MVF7, isolated in Spain in 1981–1982 [11], to provide the first complete genome sequence for a strain of *S. aureus* subsp. *anaerobius*. The genome consisted of a circular chromosome of 2,755,024 bp with 2,888 coding DNA sequences (CDS), 56 tRNAs and 5 complete rRNA copies, with a GC content of 32.74% (Fig 1A).

In order to further explore the genome content of *S. aureus* subsp. *anaerobius*, we obtained genomic DNA for 39 additional strains and performed Illumina whole genome sequencing to examine genome diversity and population genetic structure. Of these strains, 30 were isolated from different outbreaks in Spain between 1981 and 2012, with the others from Sudan ($n = 3$), Italy ($n = 3$), Poland ($n = 2$) and Denmark ($n = 1$) [10]. Among the isolates, we identified 4 closely-related novel sequence types (ST) (see [Supplementary S1 Text](#) for more detail).

The combined dataset of 41 genomes had an average length of 2,685,366 bp (range: 2,604,446 to 2,758,945 bp) that contained 2,479 genes (92.7%) shared among all isolates (core-genome) of the total 2,675 genes identified (pan-genome). On average, each strain has a highly compact accessory genome of only 130 genes (range 107–143).

S. aureus subsp. *anaerobius* evolved from an ancestor that resembled *S. aureus* subsp. *aureus* over a millennium ago

The evolutionary origin of *S. aureus* subsp. *anaerobius* and its phylogenetic relatedness to *S. aureus* subsp. *aureus* is unknown. In order to investigate its evolutionary history, we constructed a core genome alignment of *S. aureus* subsp. *anaerobius* and 807 *S. aureus* isolates that were representative of the species diversity [8] and reconstructed a maximum-likelihood phylogeny (Fig 2). *S. aureus* subsp. *anaerobius* formed a distinct clade on a long branch within the *S. aureus* subsp. *aureus* diversity but was not genetically allied to any particular clonal complex that has been previously assigned. These data indicate that *S. aureus* subsp. *anaerobius* emerged from an unknown *S. aureus* lineage likely after a human to ruminant host-switch event that occurred over a millennium ago.

A core genome sequence alignment of the 41 isolates was produced with 5 recombinant regions (size range: 54 bp to 2,950 bp, see [S1 Fig](#)) excluded leaving an alignment of 2,324,010 sites with 3,443 variable sites. Phylogenetic analyses by the maximum likelihood approach revealed two distinct clades among the *S. aureus* subsp. *anaerobius* examined representing the isolates sampled in Sudan, and all others (sampled in Europe), respectively (Fig 3). The average number of SNPs between pairs of *S. aureus* subsp. *anaerobius* genomes was 494 SNPs (range: 0–1,746), and the number of SNPs between the two clades was 1,697 SNPs (1,663–1,746).

Linear regression analysis provided evidence of a molecular clock-like evolution in the ML phylogeny (S2 Fig) and time-scaled Bayesian phylogenetic analysis indicated that the model combination that best fitted the data was a strict molecular clock paired with a constant population size, though other models performed similarly (S1 Table). This analysis dated the most

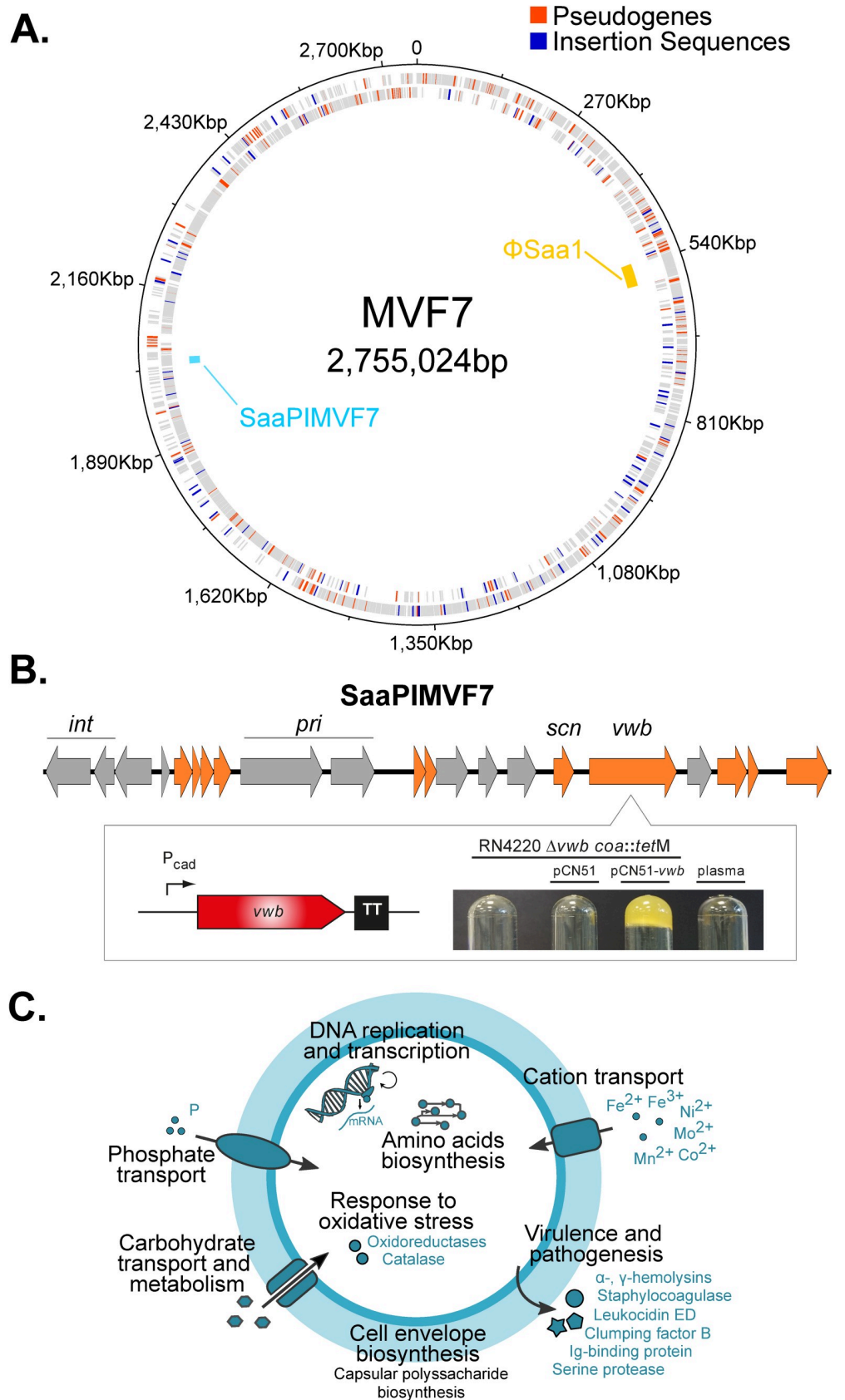


Fig 1. Pseudogenes, insertion sequences and mobile genetic elements in the *Staphylococcus aureus* subsp. *anaerobius* isolate MVF7. (A) Circular map of the chromosome. Rings show, from outside to inside: annotated genes in the positive strand, those in the negative strand (with the 205 pseudogenes and 87 insertion sequences shown in red and blue, respectively), and the mobile genetic elements found across all isolates (gold: prophage Φ Saa1; cyan: pathogenicity island SaaPIMVF7). (B) Gene map of SaaPIMVF7. Genes in grey are pseudogenes and genes in orange are intact by comparison other SAPI relatives. *int*: integrase; *pri*: primase; *scn*: Staphylococcal complement inhibitor (SCIN); *vwf*: von Willebrand factor-binding protein (vWBP). The box shows the expression of the SaaPIMVF7-encoded *vwf*. The SaaPIMVF7 *vwf* gene was cloned into the expression vector pCN51 under the control of a cadmium-inducible promoter, transformed into a coagulase and vWbp-deficient derivative of strain RN4220 (RN4220 *coa::tetM* Δ *vwf*) and the ability of SaaPIMVF7 vWbp to coagulate ruminant plasma was assessed. (C) Graphical summary of the main biological functions potentially disrupted by the presence of pseudogenes.

<https://doi.org/10.1371/journal.ppat.1009606.g001>

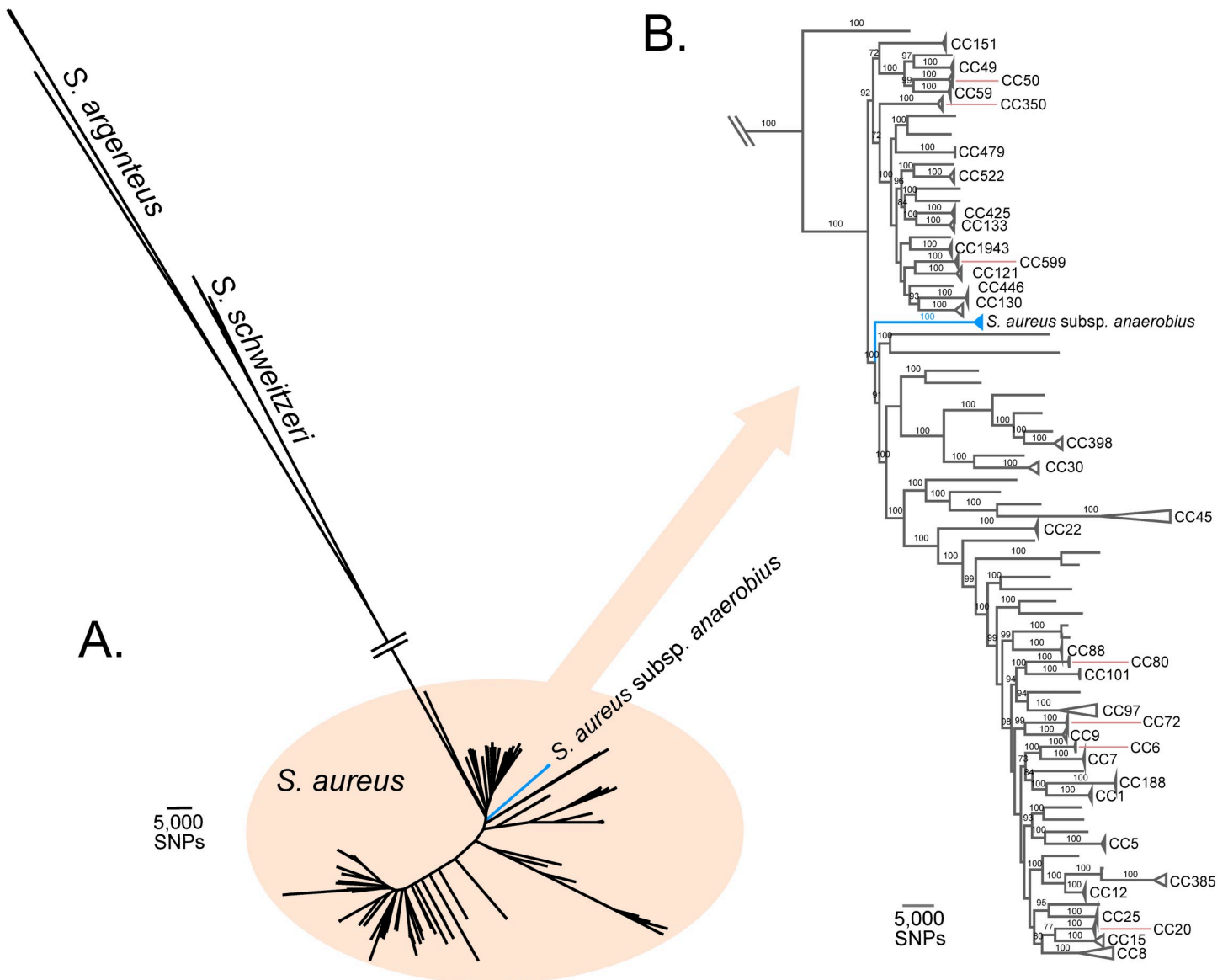


Fig 2. *Staphylococcus aureus* subsp. *anaerobius* represents a single clade within the *S. aureus* phylogenetic tree. Maximum Likelihood tree constructed from a SNP alignment of the studied *S. aureus* subsp. *anaerobius* sequences (clades in blue) and 787 sequences of different *S. aureus* subsp. *aureus* (in black). (A) Unrooted tree showing the divergence of 17 sequences of other *Staphylococcus* (*S. schweitzeri* and *S. argenteus*) whereas *S. aureus* subsp. *aureus* is embedded in the *S. aureus* diversity. (B) Subtree showing the position of the *S. aureus* subsp. *anaerobius* clade (collapsed in blue) with respect to the other *S. aureus* subsp. *aureus* clonal complexes (CC).

<https://doi.org/10.1371/journal.ppat.1009606.g002>

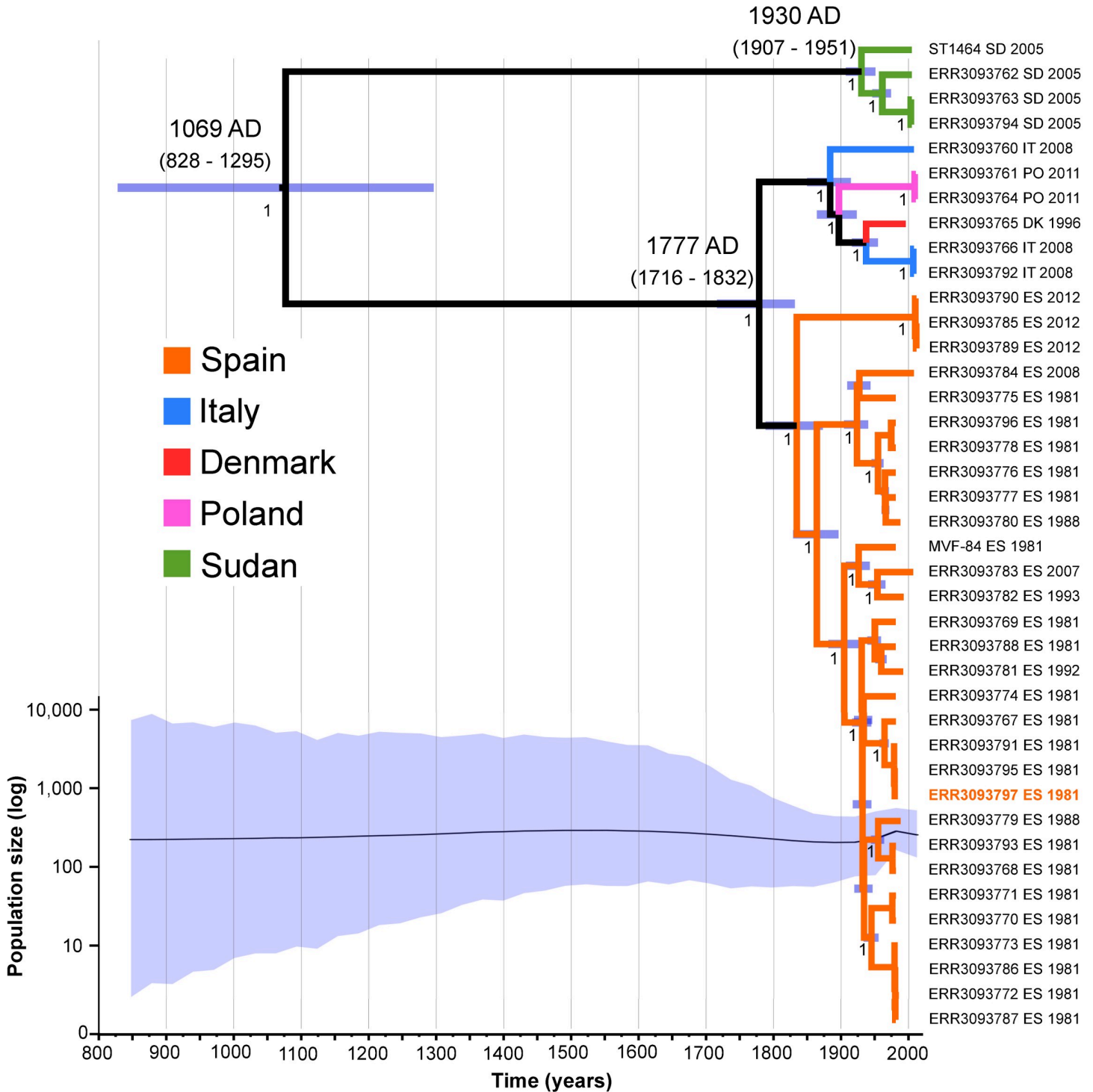


Fig 3. *S. aureus* subsp. *anaerobius* evolved over the last 1000 years with a stable population size. Bayesian maximum clade credibility tree of the *Staphylococcus aureus* subsp. *anaerobius* isolates. The x-axis is expressed in calendar years, branch colours denote sampling country. Purple bars at each node show its most recent common ancestor (MRCA) confidence interval. Numbers in nodes indicate the posterior probability. The estimated dates of the MRCAs of the main lineages are indicated. The isolate that was whole genome sequenced using PacBio technology (MV7) is highlighted in orange. The bottom plot represents the changes in the effective population size over time, with the shadowed area representing the 95% credible interval, following the same timescale as the tree.

<https://doi.org/10.1371/journal.ppat.1009606.g003>

recent common ancestor (MRCAs) of *S. aureus* subsp. *anaerobius* in 943 years before present (YBP) (95% Bayesian confidence interval (BCI): 1,184–716) (Fig 3). The clades containing European and Sudanese isolates diverged between 235 YBP (95% BCI: 296–180) and 82 YBP (95% BCI: 105–61), respectively. The estimated substitution rate was 4.4 (95% BCI: 3.3–5.5) $\times 10^{-7}$ substitutions/site/year (s/s/y), approximating to 1.2 SNPs per year which is ~one order of magnitude slower than *S. aureus* subsp. *aureus* and may reflect the relatively slow growth observed for *S. aureus* subsp. *anaerobius* [9,14].

***S. aureus* subsp. *anaerobius* has undergone large intra-chromosomal rearrangements**

During the evolutionary genomic analysis, we discovered that *S. aureus* subsp. *anaerobius* has undergone six large genomic translocations that ranged in size between 70 kb and 346 kb since separation from its *S. aureus* subsp. *aureus* progenitor (Figs 4 and S2). As the edges of each translocated region were flanked by identical insertion sequence elements, we suggest that each event was the result of transposase-mediated homologous recombination between the leading and lagging strands [15]. For 2 pairs of translocations (I, VI, and II, V, respectively) the original genomic locations were exchanged whereas translocations III and IV were the result of excision and subsequent insertion at a distinct genomic location, or by sequential single inversion events that are reciprocal (or near-reciprocal) around the origin or terminus (Fig 4). Because all translocations occurred via an inversion event, genes in the translocated areas conserved their original orientation and the GC skew remained unaffected (Fig 4).

PCR-based analysis of the distribution of the 6 rearrangements identified in strain MVF7 among all *S. aureus* subsp. *anaerobius* strains revealed that all 6 were conserved in the majority of European isolates (clade II) whereas the isolates from Sudan (clade I) lacked translocations numbers II and V, suggesting they occurred in the clade I lineage since separation from a common ancestor (S2 Table). Although the impact of the re-arrangements on global gene expression is unclear and it is feasible that these events have become fixed due to absence of any negative selective consequence, large chromosomal rearrangements in *S. aureus* subsp. *aureus* have been demonstrated to mediate transition to less virulent strains associated with persistent infection [16,17].

***S. aureus* subsp. *anaerobius* exhibits extensive genome decay**

Analysis of the complete genome of *S. aureus* subsp. *anaerobius* strain MVF7 revealed the existence of 205 pseudogenes, with a range of 201 to 210 pseudogenes per genome across the 41 *S. aureus* subsp. *anaerobius* strains examined with similar numbers in each major clade (see **Supplementary S1 Text**). Pseudogenes originated through point mutations that caused frame-shifts, premature stop codons or alternative downstream start codons and were evenly distributed around the genome (Fig 1A). As a comparison, the same pseudogene detection pipeline found an average of 14 (range: 2–30) pseudogenes per genome among the closed *S. aureus* subsp. *aureus* genomes ($n = 167$) in the reference dataset [8]. The high frequency of pseudogenes in *S. aureus* subsp. *anaerobius* represented an average of 9.5% (range: 9.2–10.1%) of the total genome length representing a remarkable example of extensive gene loss that is likely to have a major impact on bacterial phenotype. A total of 164 pseudogenes (approximately 80% of those found in any given isolate) were shared among all *S. aureus* subsp. *anaerobius* genomes, of which 92.1% were caused by the same mutation in all isolates consistent with early events in the evolutionary history of *S. aureus* subsp. *anaerobius*. Of note, 149 (90.8%) of these pseudogenes were part of the *S. aureus* “core” pangenome, as calculated from the full genomes within the *S. aureus* reference set [8].

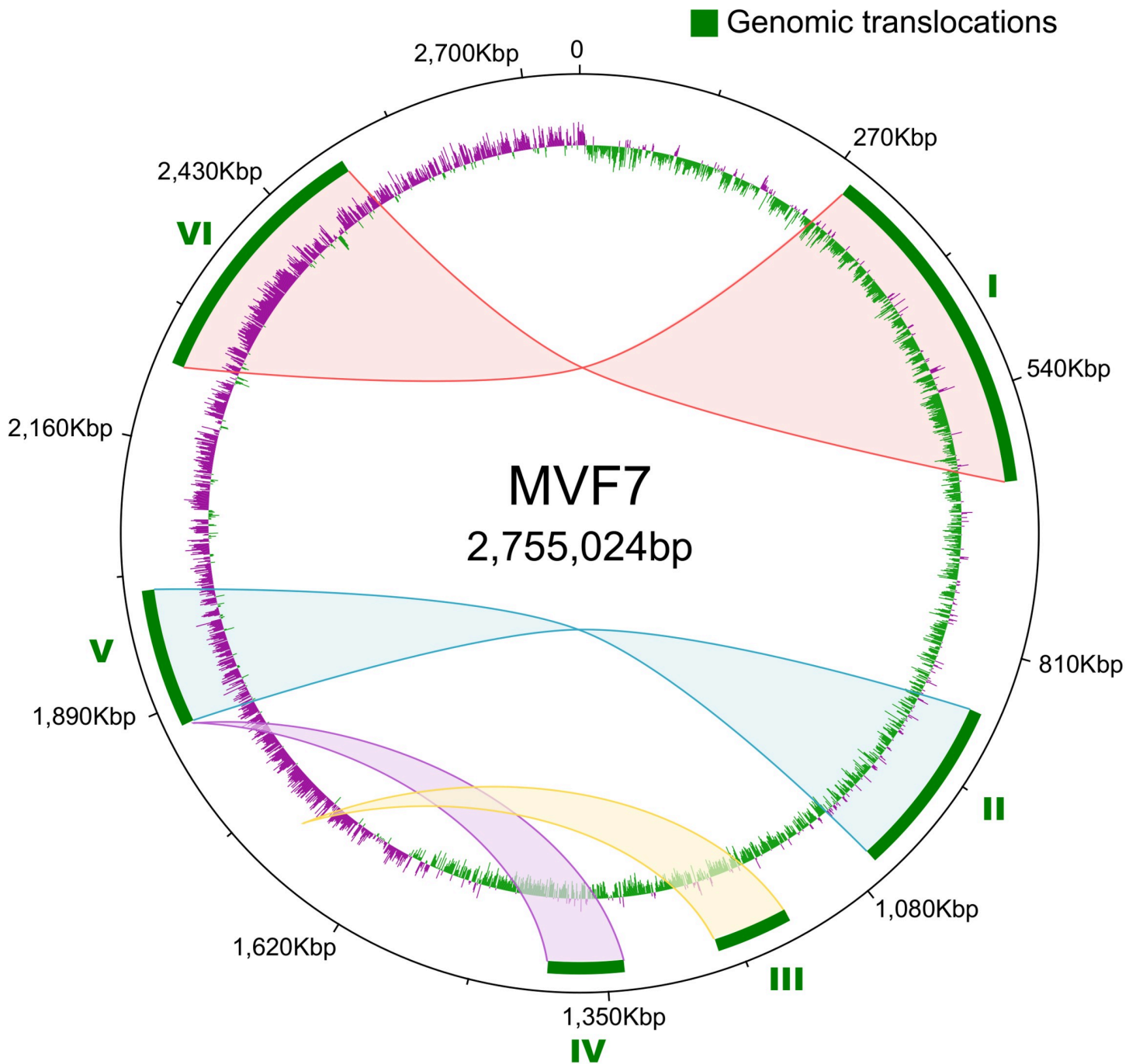


Fig 4. *S. aureus* subsp. *anaerobius* has undergone 6 large intra-chromosomal rearrangements. Map of the genomic translocations across the *Staphylococcus aureus* subsp. *anaerobius* isolate MVF7 chromosome. The green blocks represent the 6 large translocations detected in MVF7 when compared to the *Staphylococcus aureus* subsp. *aureus* RF122 strain (see also S3 Fig). The shaded areas indicate the original location of the translocated portions in the putative ancestral genome. The inner ring shows the GC skew (green for positive skew, purple for negative), which is unaffected by the translocations.

<https://doi.org/10.1371/journal.ppat.1009606.g004>

The original functions of the 205 pseudogenes present in MVF7 were classified into clusters of orthologous groups (COGs) of proteins (see S3 Table). Of these, 82 (50%) were associated with metabolic functions, 29 (17.7%) to cellular processes and signalling, and 10 (6.1%) to information storage and processing, and 43 (26.2%) had unknown functions. Enrichment

analyses of both KEGG pathways and GO terms (see [S4 Table](#) for a full list) revealed a statistically significant enrichment of pseudogenes involved in biosynthesis of amino acids and metabolic pathways (see [Supplementary S1 Text](#) and [S3 Table](#)). Of note, loss of function of genes mediating resistance to oxidative killing mechanisms such as catalase and other oxidoreductases may underpin the micro-aerophilic growth phenotype. In addition, an array of genes with known roles in pathogenicity including capsule biosynthesis, adherence and secreted proteins mediating interactions with the host, are predicted to be non-functional. Taken together, the extensive loss of gene function impacting on bacterial metabolism and pathogenicity is consistent with the highly fastidious growth requirements and defined disease tropism of *S. aureus* subsp. *anaerobius* ([Fig 1C](#)).

***S. aureus* subsp. *anaerobius* contains a novel pathogenicity island encoding a ruminant specific effector of abscess formation**

Our genomic analysis revealed a small accessory genome with limited strain-dependent variation in gene content. However, 2 putative mobile genetic elements were identified in all strains examined including a 43.2 kb, novel prophage (Φ Saa1) belonging to the *Siphoviridae* family, and a novel 13 kb *S. aureus* subsp. *anaerobius* pathogenicity island (SaaPIMVF7) ([Fig 1B](#)). While Φ Saa1 contained no putative determinants of virulence, SaaPIMVF7 encodes novel variants of known virulence factors Staphylococcal complement inhibitor (SCIN) and the von Willebrand factor-binding protein (vWbp) ([Fig 1B](#)) [[18,19](#)]. A phylogenetic network of all members of the SaPI family constructed with SplitsTree v4.14.6 and the NeighborNet algorithm [[20](#)] ([S4 Fig](#)) revealed that SaaPIMVF7 is genetically closest (93.1% nucleotide identity) to SaPIov2, found among isolates of *S. aureus* subsp. *aureus* CC133 which has a tropism for small ruminants ([Fig 1B](#)). Of note, SaaPIMVF7 contained 10 pseudogenes out of a total of 21 genes identified, including the integrase and primase genes, suggesting that it can no longer be mobilized but is a stable feature of the chromosome. Of note, both the *vwb* and *scn* genes were intact and transcribed ([S5 Fig](#)), consistent with functionality ([Fig 1C](#)). The vWBP has previously been demonstrated to promote coagulation of plasma, and is involved with abscess formation during invasive infection [[21,22](#)]. In order to test the hypothesis that SaaPIMVF7-encoded vWBP was functional, we cloned the *vwb* variant into the expression vector pCN51, and introduced this plasmid into a coagulase and vWbp-deficient derivative of strain RN4220 (RN4220 *coa::tetM* Δ *vwb*). Importantly, expression of the variant vWBP protein conferred the capacity for coagulation of ruminant plasma ([Fig 1B](#)). In summary, all strains of *S. aureus* subsp. *anaerobius* contain a novel SaPI indicating an acquisition event that occurred over 1000 years ago prior to divergence of the 2 major subclades. Strikingly, the accumulation of loss-of-function mutations in genes required for mobilization indicate that SaaPIMVF7 is a stable element of the genome and functional expression of vWBP with ruminant coagulation activity suggests a potential role in abscess formation, a defining characteristic of Morel's disease.

Expansion and intergenic location of insertion sequences in the genome of *S. aureus* subsp. *anaerobius*

An additional striking feature of *S. aureus* subsp. *anaerobius* was the presence of 87 insertion sequences (IS) distributed around the genome in strain MVF7 ([Fig 1A](#)). Each IS had 99.3% nucleotide identity with each other, and 97% identity with the previously described ISSau8 from bovine *S. aureus* strain RF122. The great majority (77 of 87) had premature stop codons that truncated the transposase gene whereas only 10 IS contained full-length intact transposase genes. The distribution of the IS elements identified in MVF7 among the rest of *S. aureus*

subsp. *anaerobius* isolates was examined by mapping Illumina reads for each strain against the border between each IS and its flanking gene in the closed genome of strain MVF7 to identify reads overlapping this border. This analysis revealed that 68 out of the 87 MVF7 ISs (78.2%) were inserted at an identical location in all 41 *S. aureus* subsp. *anaerobius* isolates consistent with an ancient acquisition that occurred during the emergence of *S. aureus* subsp. *anaerobius*. An additional 7 IS elements were shared among all the European isolates, but were absent from isolates from Sudan.

Previously, IS elements have been shown to contribute to loss of gene function or phase variation via direct insertion into CDS sequences [23]. However, all except one IS element were inserted in intergenic regions (the single exception was inserted into the *msbA* CDS), and we speculated that the non-random distribution of IS elements may influence the expression of the chromosomal genes located adjacent to the IS insertion site.

IS affects *neighbouring* gene regulation in *S. aureus* subsp. *anaerobius* by disruption of promoter region and operon structure

Our *in silico* analysis of promoter and transcriptional terminator (TT) sequences indicated that each IS had a functional promoter but lacked a TT between the *tnp* gene and the downstream gene, supporting the idea that the *tnp* transcript could influence expression of the 3'-located gene. Furthermore, 2 putative TT were identified at different sites in the IS: TT1 is located in the 5' region of the IS genes, and is bidirectional, while TT2 is located at the 3' end region of the IS and is unidirectional in the antisense orientation (see Fig 5). The putative role for the bidirectional TT1 would be to prevent the interference between transcripts initiated either from the IS or from the gene located upstream to the IS. By contrast, TT2 would prevent the impact that an antisense transcript originated from a gene located downstream of the IS could have on the expression of the IS transcript (Fig 5). This organisation suggests a strategy that allows IS transcripts to interfere with the genes located 3' to the IS, while blocking any interference that transcripts from flanking genes could have on the IS-derived regulatory transcript.

To validate the functionality and directionality of the TTs, we employed plasmid pCN42 [24], in which the *blaZ* reporter is under the control of the P_{cad} promoter. The conserved 5' and 3' region of the ISs containing the bi- and unidirectional TTs, respectively, were inserted, in both orientations, between the P_{cad} promoter and the *blaZ* reporter (see Fig 5), and the expression of the reporter gene measured with and without induction of the P_{cad} promoter. The results confirm the existence, functionality and directionality of both TTs, consistent with the hypothesis that ISs can control the expression of 3'-located chromosomal genes.

It is likely that some of the IS insertions have been fixed in the population due to the effects of genetic drift after a bottleneck. However, it is also possible that intergenic IS insertions are fixed by selection because they beneficially influence the expression of neighbouring genes. In order to test this hypothesis, we examined the impact of IS on the expression of selected downstream genes located in the same orientation as the IS. Since *S. aureus* subsp. *anaerobius* is not genetically tractable, we designed reporter constructs reconstructing the genetic organisation found in *S. aureus* subsp. *anaerobius* (S6 Fig) and tested gene expression levels in *S. aureus* subsp. *aureus*. Specifically, we selected four different genes (*rpsP*, *dnaD*, *metC* and *metN*) that were identified in *S. aureus* subsp. *anaerobius* to contain an IS at distinct distances upstream of the gene start codon, and constructed *blaZ* transcriptional reporter fusions using pCN41 [24] (Fig 6A–6D). The expression of these genes in the presence or absence of upstream ISs was then tested. Importantly, and in support of the hypothesis, absence of an IS impacted on the expression of the neighbouring genes, although this effect was different depending on the gene

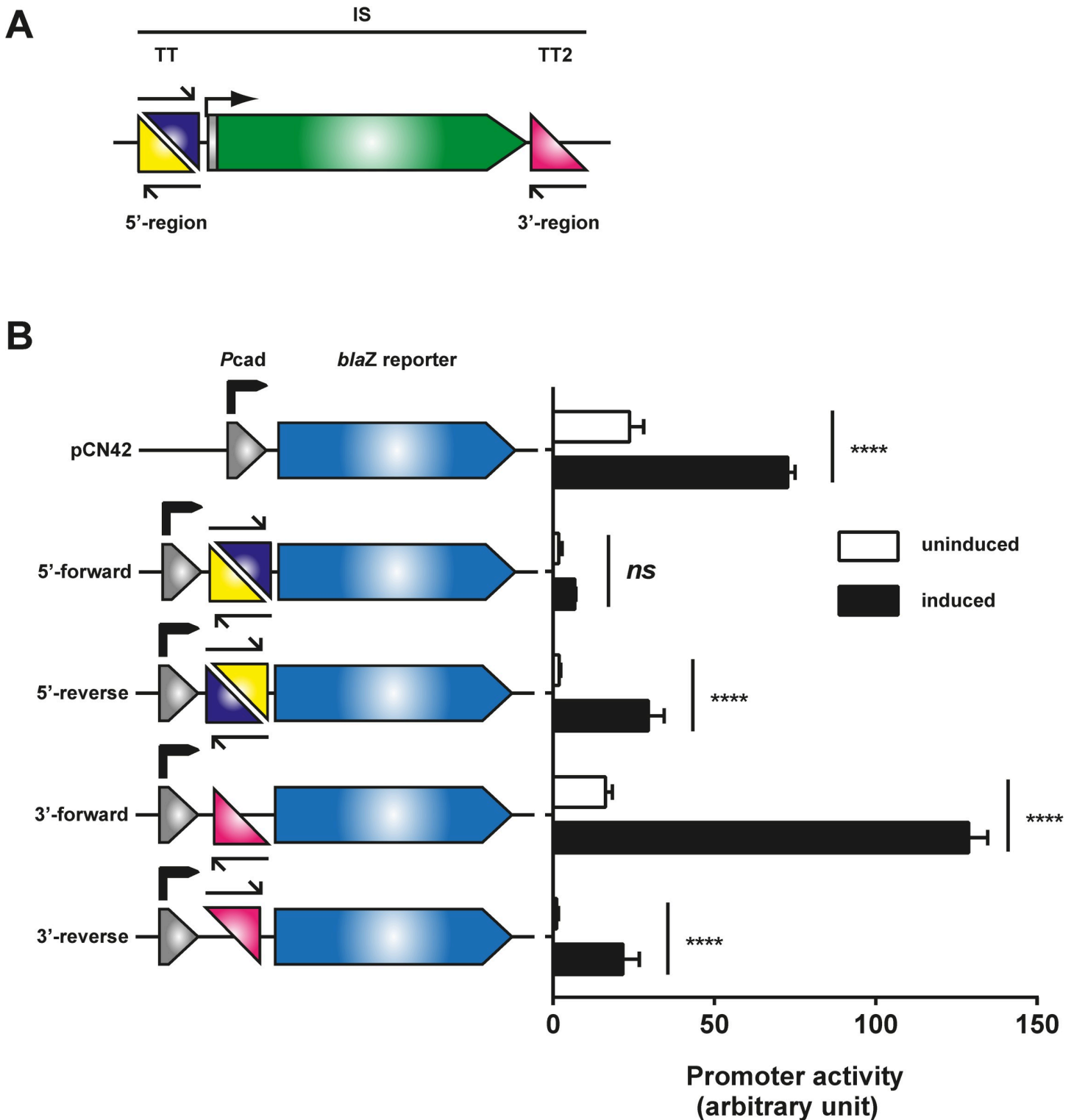


Fig 5. Characterisation of the transcriptional terminators (TTs) present in the ISs. (A) Schematic representation of IS-associated TTs. (B) Schematic representation of and functional assessment of the bi-directional TT1 localised in the 5' region of the IS. The bi-directional TT1 is composed of two sub-TTs (represented by the blue and yellow triangles, respectively) that share the same sequence but are inverted repeats, thus generating a master bi-directional TT1. The unidirectional TT2 is represented by a pink triangle. The expression of the *blaZ* reporter in plasmid pCN42 is controlled by a cadmium-inducible promoter. The identified transcriptional terminators were cloned between the promoter and the *blaZ* reporter gene. β -lactamase expression was monitored either in uninduced (black bars) or Cd-induced (open bars) cultures 180 min after induction. Data are shown as the means of three independent biological replicates and error bars show the standard deviation from the mean. Statistical analysis was

performed using two-way ANOVA followed by Dunnett's post-test. Multiplicity adjusted p-values are shown for uninduced to induced comparison and all terminator containing reporters showed significant differences when compared to their respective empty plasmid control. **** $p < 0.0001$, ns not significant.

<https://doi.org/10.1371/journal.ppat.1009606.g005>

including increased transcription of *dnaD*, *metC* or *metN*, or decreased *rpsP* transcription (Fig 6A–6D). The impact on downstream gene expression likely depends on the position of the IS relative to the gene's start codon and the nature of its integration (conservative or deleterious; see S6 Fig). In some cases, the IS will reduce expression by altering either the promoter or the DNA binding domain of a putative inducer of the flanking gene, while in other cases it may promote enhanced transcription by eliminating the binding site of a repressor. These hypotheses are currently under investigation.

ISs affect *S. aureus* subsp. *anaerobius* gene expression by a novel anti-antisense RNA interference mechanism

Intriguingly, the majority of ISs (63 of 87) were located in an antisense orientation to putative target genes. In such an orientation, the previously described effects on gene expression would not be feasible, and we reasoned that distinct mechanisms of IS-mediated control of neighbouring gene expression may exist. To test this hypothesis, we selected 2 examples of *S. aureus* subsp. *anaerobius* genes that were located in an anti-sense orientation to the flanking IS, including *mgrA*, encoding a global regulator of virulence [25–27], and *adhA* encoding a zinc-dependent alcohol dehydrogenase, induced under low oxygen conditions [28] (S4 and S5 Figs). Since the antisense orientation of these genes relative to IS precludes monitoring of gene expression by transcriptional fusions, they were cloned to express 3xFLAG tagged versions of the encoded proteins, which facilitated measurement of expression levels by western immunoblot. To examine the impact of IS on gene expression, expression plasmid constructs with and without IS in the antisense orientation were introduced to *S. aureus* subsp. *aureus* strain RN4220 Δspa and expression levels compared (see scheme in S7 Fig). Surprisingly, no differences in the expression levels of MgrA or AdhA were observed in any of the constructs suggesting that the presence of ISs did not actively alter expression levels of the downstream genes by a classical anti-sense mechanism (S7A and S7B Fig, respectively).

However, the analysis of the genetic context of *mgrA* revealed another possible mechanism of IS-mediated control of gene expression. In the ancestral *S. aureus* subsp. *aureus*, the gene downstream of *mgrA* is in an antisense orientation (S6 Fig). Furthermore, overlapping antisense transcripts for *mgrA* and its downstream gene in *S. aureus* have been identified in a previous study [29], suggesting the possibility of antisense interference of *mgrA* gene expression. Accordingly, we hypothesized that the presence of an IS between the 2 genes in *S. aureus* subsp. *anaerobius* could disrupt the putative interference (antisense) mechanism. To examine this possibility, we constructed another set of reporters (see scheme in Fig 6E) in which we used the inducible expression plasmid pCN51 to mimic expression from the gene adjacent to *mgrA*. Since the gene located downstream of *mgrA* belongs to a large operon that was impossible to clone into pCN51, we used the cadmium-inducible promoter present in plasmid pCN51 as the origin of the anti-sense mRNA for *mgrA*, mimicking the natural genetic context (Fig 6E). The different constructs were introduced into the *S. aureus* subsp. *aureus* strain RN4220 Δspa and expression from the cadmium-inducible pCN51 promoter induced. Absence of the IS resulted in much reduced MgrA expression levels suggesting that presence of a strong transcriptional terminator protects *mgrA* mRNA from antisense mRNA interference (Fig 6E). In support of this idea, insertion of a TT between *mgrA* and the plasmid-encoded cadmium-inducible promoter produced the same effect as the presence of the complete IS (Fig 6E). A similar effect was

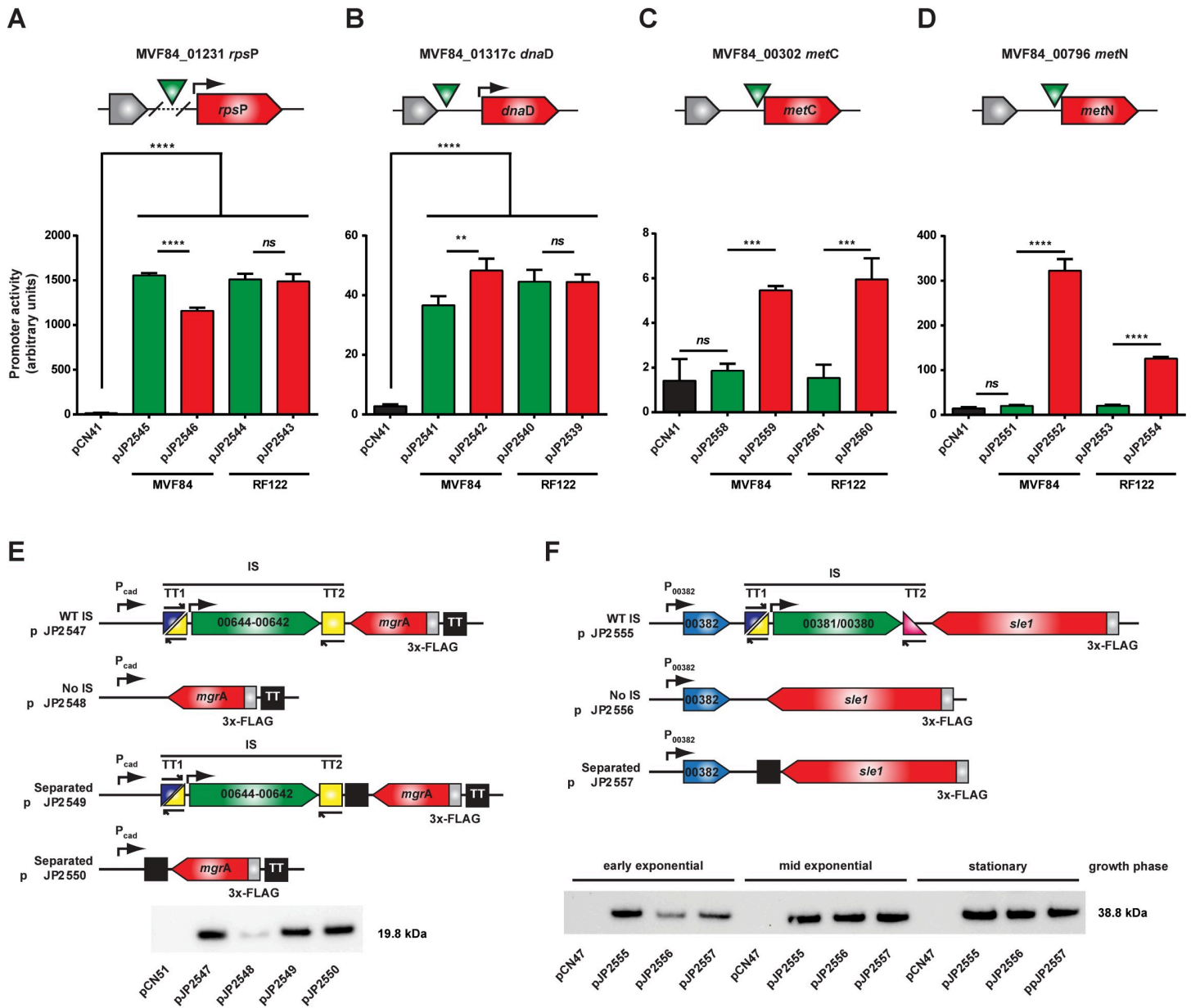


Fig 6. Presence of IS influences the expression of downstream genes and genes product through multiple mechanisms. (A-D) Reporter constructs for presence of the IS upstream and in sense orientation of the target gene were designed in plasmid pCN41 containing a β -lactamase reporter gene. The green triangle indicates the IS position relative to the target gene's start. All reporters were constructed either by removing the IS gene from *S. aureus* subsp. *anaerobius* strain MVF84 or by introducing the IS into *S. aureus* subsp. *aureus* strain RF122 to confirm the IS role on gene expression variation. In (A) the IS removal results in a 207bp deletion to the promoter region not present in RF122. Reporter plasmids were introduced into RN4220 as described in Methods. Data show the mean of three biological replicates, error bars represent the mean's standard deviation. One-Way ANOVA was performed followed by Tukey's multiple comparisons test. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$, **** $p < 0.0001$, *ns* not significant. (E&F) Western blot analysis of (E) *mgrA* and (F) *sle1* expression constructs to address whether the IS impacts on their expression levels by interrupting expression of an antisense transcript. (E) *mgrA* encoding a 3xFLAG tag was cloned in antisense to a cadmium-inducible promoter into plasmid pCN51 to mimic the expression of an antisense transcript to *mgrA*. Constructs presented different combinations of presence/absence of IS and/or transcriptional terminator the IS transcript or expression from the plasmid-encoded cadmium-inducible promoter from the *mgrA* transcript. Expression of the cadmium-inducible promoter was induced with 5 μ M CdCl₂. (F) To assess the IS impact in the context of a natural antisense transcript, 3xFLAG-encoding expression constructs of *sle1* including the gene downstream of *sle1* in the ancestral genome were cloned into the plasmid pCN47. Samples were taken during different growth phases. Plasmids were introduced into RN4220 *Aspa* as described in Methods. Western blots shown are representative of at least two independent biological replicates. TT, transcriptional terminator.

<https://doi.org/10.1371/journal.ppat.1009606.g006>

identified on *sle1* encoding an N-acetylmuramyl-L-alanine amidase involved in staphylococcal cell separation and β -lactam resistance [30,31], whereby IS insertion has disrupted the anti-sense regulation of *Sle1* expression by a neighbouring gene in a growth phase-dependent manner (Fig 6F). Thus, IS in *S. aureus* subsp. *anaerobius* can uncouple genes from ancestral anti-sense regulation by acting as anti-antisense elements, a hitherto unknown regulatory role.

Overall, our data suggest that the non-random fixation of IS into *S. aureus* subsp. *anaerobius* intergenic regions has influenced the expression of neighbouring genes through multiple distinct mechanisms. The putative attenuated expression of up to 87 different chromosomal genes involved in an array of different functions (S5 Table) has led to the re-wiring of the transcriptome.

Discussion

While *S. aureus* is an opportunistic pathogen responsible for an array of different pathologies in different anatomical sites in humans and animals, *S. aureus* subsp. *anaerobius* is restricted to a specific infection of superficial lymph nodes in sheep and goats. Here, we provide a remarkable example of the evolutionary transition of a versatile multi-host bacterium to a fastidious highly niche-restricted endemic pathogen of small ruminants. The transition was marked by multiple distinct evolutionary processes mediating drastic changes to the genome that resulted in a complete re-modelling of bacterium-niche interactions. Previous studies have indicated a human ancestral host for *S. aureus* and endemic livestock clones are the result of host jump events that have occurred during the evolutionary history of *S. aureus* [8,32]. Accordingly, the co-segregation of all isolates of the *S. aureus* subsp. *anaerobius* within a single monophyletic clade in the *S. aureus* phylogeny suggests a likely human to ruminant host-switch event that preceded the evolutionary transition to a highly niche-specific ecology.

Such switches to a host-restricted lifestyle [33], have occurred across the bacterial kingdom in diverse lineages such as *Yersinia* [34], *Mycobacterium* [35], *Shigella* [36], *Salmonella* [37], and *Burkholderia* [38], and extreme examples are represented by endosymbionts that have evolved from free-living organisms to become dependent on a single host species [2,6]. While most known examples of host-restrictive evolution have time-scales of hundreds of thousands up to many millions of years [6,33], our analysis indicates that *S. aureus* subsp. *anaerobius* evolved about 1000 or more years ago offering a unique insight into the relatively early stages of evolution towards niche-restriction.

Host shifts are associated with a radical change in habitat, typically with a genetic bottleneck that diminishes effective population size and a corresponding reduction in purifying selection activity [2,33]. We discovered that approximately 10% of the *S. aureus* subsp. *anaerobius* genome is made up of pseudogenes affecting an array of metabolic and pathogenic pathways associated with the fastidious nutritional requirements and limited virulence of *S. aureus* subsp. *anaerobius* (Fig 1C). In particular, the microaerophilic metabolism, is likely to be due to the loss of function of catalase and other oxidoreductase genes leading to increased sensitivity to oxidative free radicals. These data suggest adaptation to a nutrient-rich, oxygen-limited niche, such as that provided by the lymphatic system [11]. The repair of mildly deleterious mutations such as those resulting in loss of gene function is likely compromised by lower levels of purifying selection leading to fixation by genetic drift, and strikingly, there is no evidence for deletion of genes that are not functional in the new habitat, perhaps because sufficient time has not elapsed to facilitate this process, possibly also impacted by the observed lower rate of mutation.

We identified a large number of closely-related IS elements distributed around the genome of *S. aureus* subsp. *anaerobius* (Fig 1A; S3 Table), a phenomenon that has been observed among other bacteria evolving towards host-restriction [2,4,6,38]. Ineffective purifying selection after a population bottleneck also allows insertion sequences (IS), normally present in

bacterial genomes in low numbers (<10), to expand and disseminate around the genome in the early stages of host-restriction. Over time the IS elements may be eventually purged by deletion of regions of the genome via IS-mediated recombination, such that recently host-restricted bacteria will generally contain more IS than ancient symbionts [2,4]. It is well established that IS elements often insert into coding regions resulting in gene inactivation and can also mediate chromosomal rearrangements of genetic segments flanked by ISs via homologous recombination events [15]. We identified five large rearrangements that occurred in the common ancestor of all isolates and an additional event exclusive to the clade comprising European isolates. Although the effects on bacterial phenotype resulting from the identified rearrangements is unknown, large rearrangements can change the distance of a gene from the origin of chromosome replication leading to altered gene copy number and expression [39,40] and similar rearrangements in *S. aureus* have been reported to mediate transition to reduced virulence phenotypes such as single colony variants associated with persistent infections [16,17].

While IS elements often inactivate genes, it is striking that only one of 87 IS insertions in *S. aureus* subsp. *anaerobius* results in gene disruption. It is established that IS elements can also activate expression of neighbouring genes, either by an extended transcription from an internal promoter or by the generation of a hybrid promoter [23,41]. Uniquely, the IS-mediated gene regulation of *S. aureus* subsp. *anaerobius* involves one established and one novel mechanism of control. Upstream IS elements oriented in the sense orientation relative to the flanking gene (~25% of IS) act mainly through modulation of promoter and operon structure as previously reported [42]. However, in *S. aureus* subsp. *anaerobius*, the majority of ISs (~75%) are located downstream in the antisense orientation. In the ancestral *S. aureus* subsp. *aureus* strain, the expression of some genes is controlled by transcripts of a downstream gene in the antisense orientation suggesting expression of the corresponding proteins is mutually exclusive [43]. Here, we have shown that in *S. aureus* subsp. *anaerobius*, IS inserted in an antisense orientation to flanking genes can effectively uncouple the targeted genes from their interdependent expression control and act as anti-antisense regulatory elements. This novel regulatory mechanism might provide a selective advantage to *S. aureus* subsp. *anaerobius* in facilitating the simultaneous expression of proteins required concurrently in the new niche.

While some bacterial pathogens such as *Bordetella pertussis* utilise IS elements to control the expression of their flanking genes in a strain-specific manner [42], the genomic localisation of the majority of IS elements in *S. aureus* subsp. *anaerobius* is conserved across the phylogeny, indicating that insertion events happened early in the evolution of *S. aureus* subsp. *anaerobius*, and suggesting that the complement of genes affected may be important for the ecology of *S. aureus* subsp. *anaerobius*.

In conclusion, using a combination of phylodynamic, comparative genomic and molecular biology approaches, we have dissected the relatively recent evolution of a host-restricted bacterial pathogen underpinned by drastic changes to the genome via numerous distinct processes. In particular, IS elements have been domesticated by *S. aureus*, contributing to the chromosomal architecture via intra-chromosomal rearrangements, and control of gene expression via multiple mechanisms. Taken together, our findings provide a unique and remarkable example of the capacity of bacterial pathogens to expand into new host niches.

Methods

Bacterial strains and growth conditions

The bacterial strains employed are detailed in [S6 Table](#). *S. aureus* was grown in Tryptic soy broth (TSB) or on Tryptic soy agar and *Escherichia coli* was grown in Luria-Bertani broth (LB)

or on LB agar. Antibiotic selection was used where appropriate (erythromycin 10 $\mu\text{g ml}^{-1}$ for *S. aureus* and ampicillin 100 $\mu\text{g ml}^{-1}$ for *E. coli*).

Whole genome sequencing and genomic analysis

Forty *S. aureus* subsp. *anaerobius* ovine isolates previously reported [10] were selected, of which 31 were sampled in Spain during different outbreaks spanning 3 decades (1981–2012). The rest of the isolates, sampled between 1996 and 2011, were from Sudan ($n = 3$), Italy ($n = 3$), Poland ($n = 2$) and Denmark ($n = 1$). Isolates were sequenced using a MiSeq machine (Illumina, San Diego, CA, USA); the paired-end short reads were trimmed using Trimmomatic v0.36 [44] and *de novo* assembled into contigs using SPAdes v3.10.0 [45] and Velvet v1.2.10 [46]. Illumina read data is available in the European Nucleotide Archive under the study accession number PRJEB30965. We provide read quality statistics for the Illumina samples in [S9 Table](#).

One of the Spanish isolates (MVF7, the type strain) was sequenced on the RSII platform (Pacific Biosciences, Menlo Park, CA, USA) using SMRT technology in the Centre for Genomic Research (University of Liverpool, UK). The long reads were *de novo* assembled into a single contig using the Hierarchical Genome Assembly Process (HGAP) method. The MVF7 full genome assembly was deposited at NCBI with GenBank accession number GCA_014876765.1. The average read coverage depth was 479X (range: 120–790X) and the average read length was 18,464 bp.

Also included in the analysis was the only *S. aureus* subsp. *anaerobius* draft genome publicly available, isolate ST1464 (assembly GCA_000588835.1) [13]. Genes of all genomes and draft assemblies were annotated using Prokka v1.12 [47], and the pan-genome was determined using Roary v3.8.2 [48] applying a 95% identity cut-off. *In silico* MLST of the assemblies was performed using the mlst tool v2.8 (<https://github.com/tseemann/mlst>). Phage sequences and genomic islands were identified using PHASTER [49] and Island Viewer v4 [50], respectively, and manually inspected for reliability. Antimicrobial resistance genes were detected using ResFinder v3.0 [51].

Pseudogene detection

Here, pseudogenes were defined as protein-coding sequences with homologous genes in an *S. aureus* subsp. *aureus* reference that were split or truncated (i.e. <80% of the reference gene length). For this purpose, the isolate RF122 (assembly GCA_000009005.1) was employed as a reference and a custom python script (<https://github.com/GonzaloYebra/anaerobius>) was developed to identify pseudogenes from the output of a Roary analysis of *S. aureus* subsp. *anaerobius* isolates compared to the *S. aureus* RF122 genome. The analyses were repeated using multiple *S. aureus* subsp. *aureus* genome sequence references, but this did not alter significantly the pseudogenes detected. Ancestral functions of pseudogenes were predicted by assigning them to clusters of orthologous groups (COGs) using eggNOG v4.5.1 [52]. Enrichment analyses of GO (Gene Ontology) terms and KEGG (Kyoto Encyclopedia of Genes and Genomes) pathways assigned by eggNOG were performed using the R packages topGO v2.34.0 and ClusterProfiler v3.10.0 [53]. Additionally, the Integrated Microbial Genomes (IGM) annotation tool [54] was used to infer phenotypic metabolic characteristics from the presence/absence of protein pathways.

Insertion sequence detection

Transposase coding sequences were identified in the MVF7 complete genome using the Prokka output and the ISSaga web tool hosted in the ISfinder platform [55]. Determining the

location of IS elements and other repetitive sequences by assembly of short-read sequences is often not feasible [56]. Therefore, in order to examine the distribution of IS identified in strain MVF7 among other isolates, Illumina reads were mapped to fragments of the MVF7 genome representing each transposase edge together with their flanking genes using BWA-MEM [57]. In this manner, identification of reads that spanned the border between the IS and the flanking gene indicated the presence of that IS at the same genomic location relative to MVF7.

Phylogenetic analyses

A core genome alignment was created by aligning the Illumina short reads and the ST1464 assembly to the MVF7 whole genome using Snippy v3.1 (<https://github.com/tseemann/snippy>). Sites containing any gap character ('-') or unknown nucleotide ('N') were discarded. Gubbins v2.2.0 [58] was used to detect recombinant regions which were then discarded. A maximum likelihood (ML) tree was constructed using IQ-TREE [59], applying the GTR nucleotide substitution model together with a gamma-distributed rate heterogeneity across sites and 1,000 ultrafast bootstrap replicates. Temporal signal was investigated using TempEst v1.5 [60] by means of the correlation between root-to-tip distances and sampling dates (S2 Fig). Bayesian phylogenetic analysis was performed with BEAST v1.9.0 [61] using the HKY model for nucleotide substitution. Different models were tested for the molecular clock (strict and uncorrelated lognormal relaxed) and demographic (constant, exponential and Bayesian sky-grid) models. Each of these model combinations were run for 100 million generations, with sampling every 10,000 and discarding the initial 10% as burn-in. Runs were compared via a marginal likelihood estimation (MLE) using path sampling and stepping stone sampling methods implemented in BEAST. The posterior distribution of trees was summarised into a maximum clade credibility tree.

In order to examine the relatedness of *S. aureus* subsp. *anaerobius* and *S. aureus*, another core genome SNP tree was created, also using Snippy v3.1 and IQ-TREE (same settings than above). It included a sample of 790 *S. aureus* genomes (corresponding to 43 different host species and 77 clonal complexes (CCs), isolated in 50 different countries) along with 17 isolates of the most closely-related staphylococcal species (*S. schweitzeri* and *S. argenteus*) [8]. All bioinformatic analyses were carried out using the Cloud Infrastructure for Microbial Bioinformatics (CLIMB) facility [62].

Gene cloning

General DNA manipulations were performed using standard procedures. PCR fragments were amplified from genomic DNA using Kapa Hifi DNA polymerase (Kapa Biosystems). PCR fragments were either digested using restriction endonucleases and ligated into the respective plasmid backbone or assembled directly into a linearised plasmid backbone using Gibson assembly (NEBuilder HiFi DNA Assembly, NEB) according to the manufacturer's instructions. The plasmids and oligonucleotides used in this study are listed in S7 and S8 Tables, respectively.

Western blot analysis

S. aureus cultures harbouring the respective pCN47 or pCN51 derivative plasmids were diluted 1:50 from overnight cultures and grown in TSB at 37°C and 120 rpm until sample collection. For pCN51 constructs induced with CdCl₂, cultures were split at early exponential phase (OD₅₄₀ ~0.15), 5 µM CdCl₂ was added to half the cultures and incubation continued for another 2 h. *S. aureus* strains carrying pCN47 reporter constructs were grown to early exponential (OD₅₄₀ ~0.15), mid-exponential (OD₅₄₀ ~0.8) or stationary phase (OD₅₄₀ ~2 after

overnight culture). A sample amount corresponding to and OD₅₄₀ of 0.5 in 1 ml was pelleted and stored at -20°C. The sample pellets were re-suspended in 100 µl digestion/lysis buffer (50 mM Tris-HCl, 20 mM MgCl₂, 30% (w/v) raffinose) plus 1 µl of lysostaphin (12.5 µg ml⁻¹) and incubated at 37°C for 1 h. Samples for SDS-PAGE were prepared by adding 4X NuPAGE LDS Sample Buffer (Invitrogen) and proteins denatured at 95°C for 10 min. Samples were separated by SDS-PAGE (NuPAGE 4–12% Bis-Tris Protein Gels, Invitrogen) and then transferred to a PVDF transfer membrane (Thermo Scientific, 0.2 µM) following standard procedures [63,64]. FLAG-tagged proteins were detected using mouse anti-FLAG-HRP antibody (Monoclonal ANTI-FLAG M2-Peroxidase (HRP), Sigma-Aldrich) following the manufacturer's protocol.

Enzyme assay for the quantification of β-lactamase activity in transcriptional fusion plasmids

Samples for pCN41 based reporters were grown to either early exponential phase (OD₅₄₀~0.15) or stationary phase for promoters with low activity (*metC*) (OD₅₄₀ ~2 after overnight culture) and 1 ml of culture snap frozen. β-Lactamase assays, using nitrocefin (BD Diagnostic systems) as substrate, were performed as previously described [65] using an ELx808 microplate reader (BioTek) measuring absorbance at 490 nm. Promoter activity was calculated using the following equation:

$$\text{Promoter activity} = \frac{dA_{490}}{dt(h)} \frac{1}{OD_{540} * DF * V}$$

Supporting information

S1 Fig. Recombinant regions within the *S. aureus* subsp. *anaerobius* core genome alignment. Colour blocks (blue if affecting a single isolate, red if affecting more) represent recombinant regions detected with Gubbins. The tree at the left was built from the core genome alignment. As shown here, these recombination events affect small genomic regions and are clade-specific.

(TIF)

S2 Fig. Root-to-tip regression analysis. Root-to-tip genetic distance against sampling time estimated from a maximum-likelihood phylogenetic tree built from a core genome alignment of *S. aureus* subsp. *anaerobius* sequences.

(TIF)

S3 Fig. Pairwise genome alignment of *S. aureus* subsp. *anaerobius* versus *S. aureus* subsp. *aureus*. Artemis Comparison Tool (ACT) was used to compare both genomes (MVF7 and RF122, respectively). Red and blue bars indicate regions of similarity in the same and inverted orientation, respectively. The main 6 inverted chromosomal regions are highlighted in green and numbered.

(TIF)

S4 Fig. Phylogenetic network of representative examples of *S. aureus* Pathogenicity Islands. In bold and blue the SaPI found in *S. aureus* subsp. *anaerobius* (SaaPIMVF7). The red and blue circles indicate those SaPIs that harbour the genes *vwb* and *scn*. Reference sequences are labelled indicating SaPI name, isolate and accession number (of the SaPI sequence when available, of the original genome otherwise).

(TIF)

S5 Fig. Transcription profile of the Pathogenicity Island SaaPIMVF7. Genes in grey are pseudogenes and genes in orange are intact, according to homology against genes present in

previously described *S. aureus* pathogenicity islands (SaPIs). The histogram in blue represent the genes' transcription levels (inferred from RNA-seq read coverage).

(TIF)

S6 Fig. Schematic representation of IS loci selected for analysis. IS loci are shown for *S. aureus* subsp. *anaerobius* strain MVF84 and *S. aureus* strain RF122 representing the ancestral genomic context. (A-D) IS inserted at various distances from the downstream gene start codon. Note that in (A) IS insertion results in a 207 bp deletion in the intergenic region in strain MVF84 relative to strain RF122. (E-G) IS inserted downstream and in antisense orientation of target gene. (E&G) Locus in RF122 shows antisense orientation of downstream gene while in (F) downstream gene is in the same orientation as target gene for IS.

(TIF)

S7 Fig. Presence of the insertion sequence (IS) does not affect expression of the downstream gene product through active transcription. Western blot analysis of the depicted expression constructs for assessing the impact of IS on the expression of (A) MgrA or (B) AdhA from the IS encoded promoter. 3x-FLAG-tagged protein-encoding genes containing or missing the IS were cloned into pCN47 and plasmids introduced into the *S. aureus* subsp. *aureus* strain RN4220 Δspa for analysis. For a schematic of the locus in either *S. aureus* subsp. *anaerobius* MVF84 or *S. aureus* subsp. *aureus* RF122 refer to [S4 Fig.](#)

(TIF)

S8 Fig. Phylogenetic tree of the integrase gene from representative examples of *S. aureus* prophages. Top: genome map of the phage found in *S. aureus* subsp. *anaerobius* ($\Phi Saa1$). Genes in grey are pseudogenes and genes in orange are intact by comparison to $\Phi 2958PVL$. *Int*: integrase; *pol*: polymerase; *virE*: virulence protein E; *hel*: helicase; *ter*: terminase; *mtp*: measure tape protein; *hol*: holin; *ami*: amidase. Bottom: integrase tree. In bold and blue the integrase of $\Phi Saa1$. Reference sequences are labelled indicating integrase major group, phage and accession number.

(TIF)

S1 Table. Results from BEAST runs testing different model combinations.

(XLSX)

S2 Table. Results from PCRs design to test the presence of rearrangements.

(XLSX)

S3 Table. Features of the pseudogenes found in MVF7.

(XLSX)

S4 Table. GO terms and KEGG pathways enriched in pseudogenes.

(XLSX)

S5 Table. Features of the genes located downstream from Insertion Sequences.

(XLSX)

S6 Table. Strains used in this study.

(XLSX)

S7 Table. Plasmids used in this study.

(XLSX)

S8 Table. Oligonucleotides used in this study.

(XLSX)

S9 Table. Statistics for Illumina sequencing performed in this study.
(XLSX)

S1 Text. Supplementary results.
(PDF)

Acknowledgments

We wish to thank P. Vohra and C. Chintoan-Uta (The Roslin Institute, University of Edinburgh) for advice regarding culture of the *S. aureus* subsp. *anaerobius* isolates under microaerophilic conditions.

Author Contributions

Conceptualization: J. Ross Fitzgerald, José R. Penadés.

Formal analysis: Gonzalo Yebra, Andreas F. Haag, Maan M. Neamah, Bryan A. Wee, Emily J. Richardson.

Funding acquisition: J. Ross Fitzgerald, José R. Penadés.

Methodology: Sander Granneman.

Resources: Pilar Horcajo, María Ángeles Tormo-Más, Ricardo de la Fuente.

Writing – original draft: Gonzalo Yebra, Andreas F. Haag, J. Ross Fitzgerald, José R. Penadés.

Writing – review & editing: Gonzalo Yebra, Andreas F. Haag, J. Ross Fitzgerald, José R. Penadés.

References

1. Woolhouse ME, Taylor LH, Haydon DT. Population biology of multihost pathogens. *Science*. 2001; 292(5519):1109–12. <https://doi.org/10.1126/science.1059026> PMID: 11352066
2. Moran NA, Plague GR. Genomic changes following host restriction in bacteria. *Curr Opin Genet Dev*. 2004; 14(6):627–33. <https://doi.org/10.1016/j.gde.2004.09.003> PMID: 15531157
3. Baumler A, Fang FC. Host specificity of bacterial pathogens. *Cold Spring Harb Perspect Med*. 2013; 3(12):a010041. <https://doi.org/10.1101/cshperspect.a010041> PMID: 24296346
4. Bobay LM, Ochman H. The evolution of bacterial genome architecture. *Front Genet*. 2017; 8:72. <https://doi.org/10.3389/fgene.2017.00072> PMID: 28611826
5. Sheppard SK, Guttman DS, Fitzgerald JR. Population genomics of bacterial host adaptation. *Nat Rev Genet*. 2018; 19(9):549–65. <https://doi.org/10.1038/s41576-018-0032-z> PMID: 29973680
6. Weinert LA, Welch JJ. Why might bacterial pathogens have small genomes? *Trends Ecol Evol*. 2017; 32(12):936–47. <https://doi.org/10.1016/j.tree.2017.09.006> PMID: 29054300
7. Fitzgerald JR. Livestock-associated *Staphylococcus aureus*: origin, evolution and public health threat. *Trends in Microbiology*. 2012; 20(4):192–8. <https://doi.org/10.1016/j.tim.2012.01.006> PMID: 22386364
8. Richardson EJ, Bacigalupe R, Harrison EM, Weinert LA, Lycett S, Vrieling M, et al. Gene exchange drives the ecological success of a multi-host bacterial pathogen. *Nat Ecol Evol*. 2018. <https://doi.org/10.1038/s41559-018-0617-0> PMID: 30038246
9. de la Fuente R, Suarez G. Respiratory deficient *Staphylococcus aureus* as the aetiological agent of "abscess disease". *Zentralbl Veterinarmed B*. 1985; 32(6):397–406. <https://doi.org/10.1111/j.1439-0450.1985.tb01977.x> PMID: 4050206
10. de la Fuente R, Ballesteros C, Bautista V, Medina A, Orden JA, Domínguez-Bernal G, et al. *Staphylococcus aureus* subsp. *anaerobius* isolates from different countries are clonal in nature. *Vet Microbiol*. 2011; 150(1–2):198–202. <https://doi.org/10.1016/j.vetmic.2010.12.022> PMID: 21236606
11. De la Fuente R, Suarez G, Schleifer KH. *Staphylococcus aureus* subsp. *anaerobius* subsp. nov., the causal agent of abscess disease of sheep. *Int J Syst Bacteriol*. 1985; 35(1):99–102.

12. Szaluś-Jordanow O, Krysztopa-Grzybowska K, Czopowicz M, Moroz A, Mickiewicz M, Lutyńska A, et al. MLST and RAPD molecular analysis of *Staphylococcus aureus* subsp. *anaerobius* isolated from goats in Poland. *Archives of Microbiology*. 2018; 200(9):1407–10. <https://doi.org/10.1007/s00203-018-1568-1> PMID: 30182255
13. Elbir H, Robert C, Nguyen TT, Gimenez G, El Sanousi SM, Flock JI, et al. *Staphylococcus aureus* subsp. *anaerobius* strain ST1464 genome sequence. *Stand Genomic Sci*. 2013; 9(1):1–13. <https://doi.org/10.4056/signs.3748294> PMID: 24501641
14. Gibson B, Wilson DJ, Feil E, Eyre-Walker A. The distribution of bacterial doubling times in the wild. *Proc Biol Sci*. 2018; 285(1880). <https://doi.org/10.1098/rspb.2018.0789> PMID: 29899074
15. Parkhill J, Sebahia M, Preston A, Murphy LD, Thomson N, Harris DE, et al. Comparative analysis of the genome sequences of *Bordetella pertussis*, *Bordetella parapertussis* and *Bordetella bronchiseptica*. *Nat Genet*. 2003; 35(1):32–40. <https://doi.org/10.1038/ng1227> PMID: 12910271
16. Cui L, Neoh H-m, Iwamoto A, Hiramatsu K. Coordinated phenotype switching with large-scale chromosome flip-flop inversion observed in bacteria. *Proceedings of the National Academy of Sciences*. 2012; 109(25):E1647–E56.
17. Guérillot R, Kostoulias X, Donovan L, Li L, Carter GP, Hachani A, et al. Unstable chromosome rearrangements in *Staphylococcus aureus* cause phenotype switching associated with persistent infections. *Proceedings of the National Academy of Sciences*. 2019; 116(40):20135–40. <https://doi.org/10.1073/pnas.1904861116> PMID: 31527262
18. Guinane CM, Ben Zakour NL, Tormo-Mas MA, Weinert LA, Lowder BV, Cartwright RA, et al. Evolutionary genomics of *Staphylococcus aureus* reveals insights into the origin and molecular basis of ruminant host adaptation. *Genome Biol Evol*. 2010; 2:454–66. <https://doi.org/10.1093/gbe/evq031> PMID: 20624747
19. Viana D, Blanco J, Tormo-Mas MA, Selva L, Guinane CM, Baselga R, et al. Adaptation of *Staphylococcus aureus* to ruminant and equine hosts involves SaPI-carried variants of von Willebrand factor-binding protein. *Mol Microbiol*. 2010; 77(6):1583–94. <https://doi.org/10.1111/j.1365-2958.2010.07312.x> PMID: 20860091
20. Huson DH, Bryant D. Application of phylogenetic networks in evolutionary studies. *Mol Biol Evol*. 2006; 23(2):254–67. <https://doi.org/10.1093/molbev/msj030> PMID: 16221896
21. Malachowa N, Kobayashi SD, Porter AR, Braughton KR, Scott DP, Gardner DJ, et al. Contribution of *Staphylococcus aureus* coagulases and clumping factor A to abscess formation in a rabbit model of skin and soft tissue infection. *PLoS One*. 2016; 11(6):e0158293. <https://doi.org/10.1371/journal.pone.0158293> PMID: 27336691
22. Cheng AG, DeDent AC, Schneewind O, Missiakas D. A play in four acts: *Staphylococcus aureus* abscess formation. *Trends in Microbiology*. 2011; 19(5):225–32. <https://doi.org/10.1016/j.tim.2011.01.007> PMID: 21353779
23. Siguier P, Gourbeyre E, Chandler M. Bacterial insertion sequences: their genomic impact and diversity. *FEMS Microbiol Rev*. 2014; 38(5):865–91. <https://doi.org/10.1111/1574-6976.12067> PMID: 24499397
24. Charpentier E, Anton AI, Barry P, Alfonso B, Fang Y, Novick RP. Novel cassette-based shuttle vector system for gram-positive bacteria. *Appl Environ Microbiol*. 2004; 70(10):6076–85. <https://doi.org/10.1128/AEM.70.10.6076-6085.2004> PMID: 15466553
25. Luong TT, Lee CY. The *arl* locus positively regulates *Staphylococcus aureus* type 5 capsule via an *mgrA*-dependent pathway. *Microbiology (Reading)*. 2006; 152(Pt 10):3123–31. <https://doi.org/10.1099/mic.0.29177-0> PMID: 17005991
26. Crosby HA, Schlievert PM, Merriman JA, King JM, Salgado-Pabon W, Horswill AR. The *Staphylococcus aureus* global regulator MgrA modulates clumping and virulence by controlling surface protein expression. *PLoS Pathog*. 2016; 12(5):e1005604. <https://doi.org/10.1371/journal.ppat.1005604> PMID: 27144398
27. Crosby HA, Tiwari N, Kwiecinski JM, Xu Z, Dykstra A, Jenul C, et al. The *Staphylococcus aureus* ArlRS two-component system regulates virulence factor expression through MgrA. *Mol Microbiol*. 2020; 113(1):103–22. <https://doi.org/10.1111/mmi.14404> PMID: 31618469
28. Pagels M, Fuchs S, Pane-Farre J, Kohler C, Menschner L, Hecker M, et al. Redox sensing by a Rex-family repressor is involved in the regulation of anaerobic gene expression in *Staphylococcus aureus*. *Mol Microbiol*. 2010; 76(5):1142–61. <https://doi.org/10.1111/j.1365-2958.2010.07105.x> PMID: 20374494
29. Lasa I, Toledo-Arana A, Dobin A, Villanueva M, de los Mozos IR, Vergara-Irigaray M, et al. Genome-wide antisense transcription drives mRNA processing in bacteria. *Proc Natl Acad Sci U S A*. 2011; 108(50):20172–7. <https://doi.org/10.1073/pnas.1113521108> PMID: 22123973
30. Kajimura J, Fujiwara T, Yamada S, Suzawa Y, Nishida T, Oyamada Y, et al. Identification and molecular characterization of an N-acetylmuramyl-L-alanine amidase Sle1 involved in cell separation of

- Staphylococcus aureus*. *Mol Microbiol*. 2005; 58(4):1087–101. <https://doi.org/10.1111/j.1365-2958.2005.04881.x> PMID: 16262792
31. Thalsø-Madsen I, Torrubia FR, Xu L, Petersen A, Jensen C, Frees D. The Sle1 cell wall amidase is essential for β -lactam resistance in community-acquired methicillin-resistant *Staphylococcus aureus* USA300. *Antimicrob Agents Chemother*. 2019; 64(1). <https://doi.org/10.1128/AAC.01931-19> PMID: 31685469
 32. Weinert LA, Welch JJ, Suchard MA, Lemey P, Rambaut A, Fitzgerald JR. Molecular dating of human-to-bovid host jumps by *Staphylococcus aureus* reveals an association with the spread of domestication. *Biol Lett*. 2012; 8(5):829–32. <https://doi.org/10.1098/rsbl.2012.0290> PMID: 22628096
 33. Toft C, Andersson SG. Evolutionary microbial genomics: insights into bacterial host adaptation. *Nat Rev Genet*. 2010; 11(7):465–75. <https://doi.org/10.1038/nrg2798> PMID: 20517341
 34. Parkhill J, Wren BW, Thomson NR, Titball RW, Holden MTG, Prentice MB, et al. Genome sequence of *Yersinia pestis*, the causative agent of plague. *Nature*. 2001; 413(6855):523–7. <https://doi.org/10.1038/35097083> PMID: 11586360
 35. Cole ST, Eiglmeier K, Parkhill J, James KD, Thomson NR, Wheeler PR, et al. Massive gene decay in the leprosy bacillus. *Nature*. 2001; 409(6823):1007–11. <https://doi.org/10.1038/35059006> PMID: 11234002
 36. Feng Y, Chen Z, Liu SL. Gene decay in *Shigella* as an incipient stage of host-adaptation. *PLoS One*. 2011; 6(11):e27754. <https://doi.org/10.1371/journal.pone.0027754> PMID: 22110755
 37. Langridge GC, Fookes M, Connor TR, Feltwell T, Feasey N, Parsons BN, et al. Patterns of genome evolution that have accompanied host adaptation in *Salmonella*. *Proc Natl Acad Sci USA*. 2015; 112(3):863–8. <https://doi.org/10.1073/pnas.1416707112> PMID: 25535353
 38. Song H, Hwang J, Yi H, Ulrich RL, Yu Y, Nierman WC, et al. The early stage of bacterial genome-reductive evolution in the host. *PLoS Pathog*. 2010; 6(5):e1000922. <https://doi.org/10.1371/journal.ppat.1000922> PMID: 20523904
 39. Block DH, Hussein R, Liang LW, Lim HN. Regulatory consequences of gene translocation in bacteria. *Nucleic Acids Res*. 2012; 40(18):8979–92. <https://doi.org/10.1093/nar/gks694> PMID: 22833608
 40. Periwal V, Scaria V. Insights into structural variations and genome rearrangements in prokaryotic genomes. *Bioinformatics*. 2015; 31(1):1–9. <https://doi.org/10.1093/bioinformatics/btu600> PMID: 25189783
 41. Prentki P, Teter B, Chandler M, Galas DJ. Functional promoters created by the insertion of transposable element IS1. *J Mol Biol*. 1986; 191(3):383–93. [https://doi.org/10.1016/0022-2836\(86\)90134-8](https://doi.org/10.1016/0022-2836(86)90134-8) PMID: 3029382
 42. Amman F, D'Halluin A, Antoine R, Huot L, Bibova I, Keidel K, et al. Primary transcriptome analysis reveals importance of IS elements for the shaping of the transcriptional landscape of *Bordetella pertussis*. *RNA Biol*. 2018; 15(7):967–75. <https://doi.org/10.1080/15476286.2018.1462655> PMID: 29683387
 43. Toledo-Arana A, Lasa I. Advances in bacterial transcriptome understanding: From overlapping transcription to the excludon concept. *Mol Microbiol*. 2020; 113(3):593–602. <https://doi.org/10.1111/mmi.14456> PMID: 32185833
 44. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics*. 2014; 30(15):2114–20. <https://doi.org/10.1093/bioinformatics/btu170> PMID: 24695404
 45. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *J Comput Biol*. 2012; 19(5):455–77. <https://doi.org/10.1089/cmb.2012.0021> PMID: 22506599
 46. Zerbino DR, Birney E. Velvet: algorithms for de novo short read assembly using de Bruijn graphs. *Genome Res*. 2008; 18(5):821–9. <https://doi.org/10.1101/gr.074492.107> PMID: 18349386
 47. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics*. 2014; 30(14):2068–9. <https://doi.org/10.1093/bioinformatics/btu153> PMID: 24642063
 48. Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MT, et al. Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics*. 2015; 31(22):3691–3. <https://doi.org/10.1093/bioinformatics/btv421> PMID: 26198102
 49. Arndt D, Grant JR, Marcu A, Sajed T, Pon A, Liang Y, et al. PHASTER: a better, faster version of the PHAST phage search tool. *Nucleic Acids Res*. 2016; 44(W1):W16–21. <https://doi.org/10.1093/nar/gkw387> PMID: 27141966
 50. Bertelli C, Laird MR, Williams KP, Simon Fraser University Research Computing G, Lau BY, Hoard G, et al. IslandViewer 4: expanded prediction of genomic islands for larger-scale datasets. *Nucleic Acids Res*. 2017. <https://doi.org/10.1093/nar/gkx343> PMID: 28472413

51. Zankari E, Hasman H, Cosentino S, Vestergaard M, Rasmussen S, Lund O, et al. Identification of acquired antimicrobial resistance genes. *J Antimicrob Chemother.* 2012; 67(11):2640–4. <https://doi.org/10.1093/jac/dks261> PMID: 22782487
52. Huerta-Cepas J, Szklarczyk D, Forslund K, Cook H, Heller D, Walter MC, et al. eggNOG 4.5: a hierarchical orthology framework with improved functional annotations for eukaryotic, prokaryotic and viral sequences. *Nucleic Acids Res.* 2016; 44(D1):D286–93. <https://doi.org/10.1093/nar/gkv1248> PMID: 26582926
53. Yu G, Wang LG, Han Y, He QY. clusterProfiler: an R package for comparing biological themes among gene clusters. *OMICS.* 2012; 16(5):284–7. <https://doi.org/10.1089/omi.2011.0118> PMID: 22455463
54. Markowitz VM, Chen IM, Palaniappan K, Chu K, Szeto E, Grechkin Y, et al. IMG: the Integrated Microbial Genomes database and comparative analysis system. *Nucleic Acids Res.* 2012; 40(Database issue):D115–22. <https://doi.org/10.1093/nar/gkr1044> PMID: 22194640
55. Varani AM, Siguier P, Gourbeyre E, Charneau V, Chandler M. ISsaga is an ensemble of web-based methods for high throughput identification and semi-automatic annotation of insertion sequences in prokaryotic genomes. *Genome Biol.* 2011; 12(3):R30. <https://doi.org/10.1186/gb-2011-12-3-r30> PMID: 21443786
56. Treangen TJ, Salzberg SL. Repetitive DNA and next-generation sequencing: computational challenges and solutions. *Nat Rev Genet.* 2011; 13(1):36–46. <https://doi.org/10.1038/nrg3117> PMID: 22124482
57. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics.* 2009; 25(14):1754–60. <https://doi.org/10.1093/bioinformatics/btp324> PMID: 19451168
58. Croucher NJ, Page AJ, Connor TR, Delaney AJ, Keane JA, Bentley SD, et al. Rapid phylogenetic analysis of large samples of recombinant bacterial whole genome sequences using Gubbins. *Nucleic Acids Res.* 2015; 43(3):e15. <https://doi.org/10.1093/nar/gku1196> PMID: 25414349
59. Nguyen LT, Schmidt HA, von Haeseler A, Minh BQ. IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol Biol Evol.* 2015; 32(1):268–74. <https://doi.org/10.1093/molbev/msu300> PMID: 25371430
60. Rambaut A, Lam TT, Max Carvalho L, Pybus OG. Exploring the temporal structure of heterochronous sequences using TempEst (formerly Path-O-Gen). *Virus Evol.* 2016; 2(1):vew007. <https://doi.org/10.1093/ve/vew007> PMID: 27774300
61. Drummond AJ, Suchard MA, Xie D, Rambaut A. Bayesian phylogenetics with BEAUti and the BEAST 1.7. *Mol Biol Evol.* 2012; 29(8):1969–73. <https://doi.org/10.1093/molbev/mss075> PMID: 22367748
62. Connor TR, Loman NJ, Thompson S, Smith A, Southgate J, Poplawski R, et al. CLIMB (the Cloud Infrastructure for Microbial Bioinformatics): an online resource for the medical microbiology community. *Microb Genom.* 2016; 2(9):e000086. <https://doi.org/10.1099/mgen.0.000086> PMID: 28785418
63. Ausubel F, Brent R, Kingston R, Moore D, Seidman J, Smith J, et al. *Current Protocols in Molecular Biology.* New York, NY: John Wiley & Sons; 1990.
64. Sambrook J, Fritsch EF, Maniatis T. *Molecular cloning: a laboratory manual.* Cold Spring Harbor, NY: Cold Spring Harbor Laboratory Press; 1989.
65. Ji G, Beavis R, Novick RP. Bacterial interference caused by autoinducing peptide variants. *Science.* 1997; 276(5321):2027–30. <https://doi.org/10.1126/science.276.5321.2027> PMID: 9197262