



The emergence and diversification of a zoonotic pathogen from within the microbiota of intensively farmed pigs

Murray, Gemma G.R.; Hossain, A. S.Md Mukarram; Miller, Eric L.; Bruchmann, Sebastian; Balmer, Andrew J.; Matuszewska, Marta; Herbert, Josephine; Hadjirin, Nazreen F.; Mugabi, Robert; Li, Ganwu

Total number of authors:
31

Published in:
Proceedings of the National Academy of Sciences of the United States of America

Link to article, DOI:
[10.1073/pnas.2307773120](https://doi.org/10.1073/pnas.2307773120)

Publication date:
2023

Document Version
Publisher's PDF, also known as Version of record

[Link back to DTU Orbit](#)

Citation (APA):
Murray, G. G. R., Hossain, A. S. M. M., Miller, E. L., Bruchmann, S., Balmer, A. J., Matuszewska, M., Herbert, J., Hadjirin, N. F., Mugabi, R., Li, G., Ferrando, M. L., de Oliveira, I. M. F., Nguyen, T., Yen, P. L. K., Phuc, H. D., Moe, A. Z., Wai, T. S., Gottschalk, M., Aragon, V., ... Weinert, L. A. (2023). The emergence and diversification of a zoonotic pathogen from within the microbiota of intensively farmed pigs. *Proceedings of the National Academy of Sciences of the United States of America*, 120(47), Article e2307773120.
<https://doi.org/10.1073/pnas.2307773120>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.



The emergence and diversification of a zoonotic pathogen from within the microbiota of intensively farmed pigs

Gemma G. R. Murray^{a,b,1} , A. S. Md. Mukarram Hossain^c , Eric L. Miller^d, Sebastian Bruchmann^b , Andrew J. Balmer^b, Marta Matuszewska^{b,e} , Josephine Herbert^f , Nazreen F. Hadjirin^g, Robert Mugabi^h , Ganwu Li^h , Maria Laura Ferrandoⁱ , Isabela Maria Fernandes de Oliveira^j, Thanh Nguyenⁱ , Phung L. K. Yen^j, Ho D. Phuc^j, Aung Zaw Moe^k, Thiri Su Wai^k, Marcelo Gottschalk^l, Virginia Aragon^{m,n} , Peter Valentin-Weigand^o, Peter M. H. Heegaard^p , Manouk Vrieling^q, Min Thein Maw^k, Hnin Thidar Myint^k, Ye Tun Win^k, Ngo Thi Hoa^{i,r,s}, Stephen D. Bentley^t , Maria J. Clavijo^h, Jerry M. Wells^{b,j}, Alexander W. Tucker^b, and Lucy A. Weinert^b

Edited by Michael S. Gilmore, Harvard Medical School, Boston, MA; received May 22, 2023; accepted October 2, 2023 by Editorial Board Member Caroline S. Harwood

The expansion and intensification of livestock production is predicted to promote the emergence of pathogens. As pathogens sometimes jump between species, this can affect the health of humans as well as livestock. Here, we investigate how livestock microbiota can act as a source of these emerging pathogens through analysis of *Streptococcus suis*, a ubiquitous component of the respiratory microbiota of pigs that is also a major cause of disease on pig farms and an important zoonotic pathogen. Combining molecular dating, phylogeography, and comparative genomic analyses of a large collection of isolates, we find that several pathogenic lineages of *S. suis* emerged in the 19th and 20th centuries, during an early period of growth in pig farming. These lineages have since spread between countries and continents, mirroring trade in live pigs. They are distinguished by the presence of three genomic islands with putative roles in metabolism and cell adhesion, and an ongoing reduction in genome size, which may reflect their recent shift to a more pathogenic ecology. Reconstructions of the evolutionary histories of these islands reveal constraints on pathogen emergence that could inform control strategies, with pathogenic lineages consistently emerging from one subpopulation of *S. suis* and acquiring genes through horizontal transfer from other pathogenic lineages. These results shed light on the capacity of the microbiota to rapidly evolve to exploit changes in their host population and suggest that the impact of changes in farming on the pathogenicity and zoonotic potential of *S. suis* is yet to be fully realized.

Streptococcus suis | pathogen emergence | bacterial pathogens | comparative genomics | livestock pathogens

Global livestock populations have grown rapidly over the past few centuries, with the global biomass of livestock now exceeding that of humans and wild mammals combined (1, 2). This has been facilitated by intensive farming systems that have also led to increased livestock population density, lower genetic diversity, and the long-distance movement of live animals. These changes are predicted to promote the emergence of pathogens (3, 4). While pathogen emergence typically arises through a pathogen jumping into a new host, pathogens can also emerge from within the microbiota already associated with a host population (5, 6). This route to pathogen emergence may be particularly important in intensive farming systems, where large population size and high population density may select for traits associated with pathogenicity, while biosecurity reduces the risk of novel pathogens entering the population (7, 8).

Streptococcus suis was first reported as a cause of disease in farmed pigs in 1954 (9) and is now a major cause of bacterial disease in piglets and an emerging human zoonotic pathogen (10, 11). As well as being an important pathogen, *S. suis* is a ubiquitous component of the microbiota of the upper respiratory tract of all pigs. It is one of the most common bacterial species on the surface of the palatine tonsil, which is considered its main niche (12). *S. suis* disease in pigs takes the form of septicemia with sudden death, meningitis, arthritis, and endocarditis and most often affects piglets. It is also associated with respiratory disease, although these infections tend to be polymicrobial (13). Humans can be infected by *S. suis* either through contact with pigs or consumption of raw pork or other pig products. These infections result in similar pathologies to pigs and have high fatality rates. The first reported human case was in 1968 (14), and since then, *S. suis* has led to large outbreaks in China and has become a major cause of adult meningitis and septicemia in South-East Asia (15–18).

While *S. suis* is a diverse species, only a small number of strains, typically characterised by multilocus sequence type (ST) or serotype, are responsible for most cases of disease

Significance

There is growing concern that rapid growth in livestock production and major changes in farming practices are driving the emergence of pathogens capable of causing disease in both livestock and humans. However, most studies neglect livestock microbiota as a potential source of emerging pathogens. Here, we show how the global transport of live animals has facilitated the emergence of an important livestock and human zoonotic pathogen from a common member of the pig respiratory microbiota. Our results indicate that pathogenic lineages are likely to continue to emerge and diversify and recommend ways of controlling this.

Author contributions: G.G.R.M., J.M.W., A.W.T., and L.A.W. designed research; G.G.R.M., A.S.M.M.H., E.L.M., S.B., A.J.B., M.M., J.H., N.F.H., R.M., G.L., M.L.F., I.M.F.d.O., T.N., P.L.K.Y., H.D.P., A.Z.M., T.S.W., M.G., V.A., P.V.-W., P.M.H.H., M.V., M.T.M., and H.T.M. performed research; G.G.R.M., A.S.M.M.H., S.B., and R.M. analyzed data; G.G.R.M., E.L.M., Y.T.W., N.T.H., S.D.B., M.J.C., J.M.W., A.W.T., and L.A.W. supervised and managed project; and G.G.R.M. and L.A.W. wrote the paper.

The authors declare no competing interest.

This article is a PNAS Direct Submission. M.S.G. is a guest editor invited by the Editorial Board.

Copyright © 2023 the Author(s). Published by PNAS. This open access article is distributed under Creative Commons Attribution License 4.0 (CC BY).

¹To whom correspondence may be addressed. Email: ggrmurray@gmail.com.

This article contains supporting information online at <https://www.pnas.org/lookup/suppl/doi:10.1073/pnas.2307773120/-DCSupplemental>.

Published November 14, 2023.

(19). What determines the pathogenicity of these strains remains poorly understood despite the identification of more than 100 putative virulence genes or factors (20). Difficulties in identifying the determinants of pathogenicity in *S. suis* have been attributed to its complex pathogenesis and high level of genetic diversity (20). Few studies have considered virulence factors in strains other than ST 1, which is responsible for most cases of *S. suis* disease in both pigs and humans worldwide (19).

In this study, we carried out a population-genomic analysis of 3,070 bacterial isolates sampled from tonsil and nasal swabs of pigs and wild boar, and blood and sites of infection in pigs and humans with *S. suis* disease, from Europe, North America, Asia, and Australia, dating from 1960 to 2020, to investigate the emergence, diversification, and geographic spread of pathogenic lineages of *S. suis*. Through development of a whole-genome typing schema, we identified 10 pathogenic lineages with broad geographic distributions, dated their origins, and mapped their movements between countries. We identified genomic changes associated with the emergence of these lineages and investigated their origins. We also considered the impact of pathogenicity on broader evolutionary dynamics, the ongoing diversification of pathogenic lineages, and how farming practices may have contributed to these processes.

Results

Pathogenic Lineages Emerged from a Subpopulation of a Diverse and Largely Commensal Species. We sequenced the genomes of isolates identified in laboratory assays as *S. suis* from pigs from farms in Denmark (n = 173), Germany (n = 166), the Netherlands (n = 168), Spain (n = 200), the United Kingdom (n = 49), Myanmar (n = 701), and the United States (n = 293) (SI Appendix, Table S1). These isolates date from 1960 to 2020. Those from Denmark, Germany, the Netherlands, and Spain included similar numbers of isolates from the tonsils or noses of pigs without

S. suis-associated disease (carriage isolates; n = 188) and isolates from blood or sites of infection in pigs with respiratory (n = 196) or systemic (n = 196) forms of *S. suis*-associated disease (disease isolates). Those from Myanmar only included carriage and environmental isolates and those from the United States only disease isolates. The collection from Spain included isolates from a wild boar population (n = 34). We combined these with previously published genome sequence data from isolates from Australia (n = 143), Canada (n = 200), China (n = 217), Denmark (n = 1), Vietnam (n = 191), the United Kingdom (n = 441), the Netherlands (n = 101), Spain (n = 10), and the United States (n = 16). This led to a collection of 3,070 high-quality genome assemblies, including 48 complete genomes.

To investigate the population structure of *S. suis*, we considered genetic variation estimated from both single nucleotide polymorphisms (SNPs) in genes that are present across all isolates (core genes) and from the presence/absence of accessory genes (Fig. 1A and SI Appendix, Figs. S1–S4). Both revealed high levels of diversity and distinguished a cluster of closely related isolates that includes most of our collection (2,424/3,070 isolates). Isolates in this cluster have a maximum pairwise distance of 0.08 differences per nucleotide in a core genome alignment, while the remaining 647 isolates diverge from members of this cluster by between 0.13 and 0.32 differences per nucleotide site. We refer to this cluster as the “central population” of *S. suis* [a previous study referred to it as “normal” *S. suis* (21)]. The isolates that fall outside of the central population form several distinct clades in a core genome phylogeny (SI Appendix, Fig. S1). While it has been suggested that these may represent distinct species, a previous study concluded that phenotypic similarities and extensive gene exchange between some of these lineages and the central *S. suis* population meant that there was insufficient evidence to assign them to a new species or subspecies (21). In support of this conclusion, we found that divergent isolates are less clearly distinguished from the central population of *S. suis*

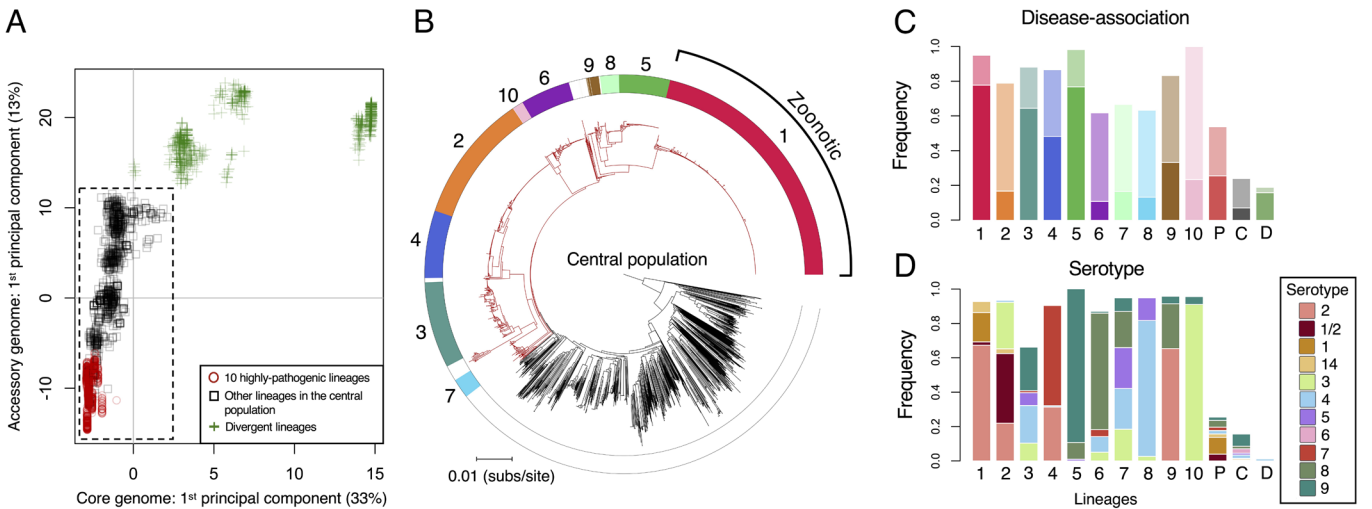


Fig. 1. The relationship between pathogenicity and genetic structure within *S. suis*. (A) The first components of principal component analyses of the presence/absence of accessory genes and SNPs in core genes plotted against one another for all 3,070 isolates in our collection (other components and their eigenvalues are shown in SI Appendix, Figs. S3 and S4). Points represent individual isolates and different shapes/colours represent categories of lineages. Isolates in the dashed box fall within the central population of *S. suis*. (B) A core genome phylogeny of the central population of *S. suis* (excluding divergent isolates). Branches in the pathogenic clade are coloured red and branches outside of this clade are coloured black. Colours in the outer ring indicate the 10 most common lineages in our collection, which we also find to be the most pathogenic. Lineage 1 largely corresponds to ST1, which is associated with most cases of zoonotic disease in humans. (C) The frequency of disease-associated isolates in each of the 10 pathogenic lineages (1 to 10), in other lineages within the pathogenic clade (P, coloured red in B), in the central population outside of the pathogenic clade (C, coloured black in B), and in divergent lineages (D, green addition signs in A). Bars 1 to 10 are coloured to match the outer ring in B with the lower part of the bar (deeper colour) representing the frequency of systemic disease isolates and the upper part of the bar (paler colour) representing the frequency of respiratory disease isolates. (D) The frequency of disease-associated serotypes in the same groups as shown in C.

Table 1. STs and serotypes associated with the 10 pathogenic lineages

Lineage	STs	Serotypes	Number of isolates in our collection	Frequency of disease association (systemic disease association)
1	1	1, 2, 14	530	0.95 (0.78)
2	27, 28	1/2, 2, 3	260	0.79 (0.17)
3	94, 108, 123	3, 4, 5, 9, 23	174	0.88 (0.64)
4	25, 29	2, 7	136	0.87 (0.48)
5	16, 136, 198	8, 9	193	0.98 (0.77)
6	23, 87, 89	3, 4, 8	100	0.62 (0.11)
7	54	3, 4, 5, 8, 9	39	0.67 (0.17)
8	17	4, 5	38	0.63 (0.13)
9	14, 20	2, 8	23	0.83 (0.33)
10	15	3	22	1.00 (0.24)

Each of the 10 pathogenic lineages includes multiple STs and serotypes. STs and serotypes present in >5% of the isolates from our collection of each lineage are shown here, along with the number of isolates from these lineages in our collection, and the frequency of disease association (and systemic disease association) in the subset of our collection we used to characterize pathogenicity. A description of the STs and serotypes of all isolates in our collection is provided in [SI Appendix, Table S1](#).

by differences in gene content than by SNPs in core genes (Fig. 1A and [SI Appendix, Fig. S1](#)).

To mitigate any geographic bias, we characterised variation in disease association across *S. suis* using a subset of our collection that includes isolates that a) have a well-characterised association with disease or carriage states and b) are from a country from which we have large samples of both disease and carriage isolates in our collection (>40 of each). This subset includes isolates from Denmark, Germany, the Netherlands, Spain, Canada, and the United Kingdom (n = 1,193). In agreement with previous studies, our analysis revealed that while disease isolates are present across the entire genetic diversity of *S. suis* (including in divergent clades), they are concentrated in a subpopulation of the central population (Fig. 1B and C) (22).

Using variation in both core and accessory genes, we partitioned the central population of *S. suis* into clusters of closely related isolates (lineages) ([SI Appendix, Table S1](#)). A reference database describing these lineages, which will allow for classification of further isolates using the same schema, is available online (www.bacpop.org/poppunk). Combining this with data on disease association, allowed us to identify 10 “pathogenic” lineages. In each of these lineages, at least 60% of the isolates from the subset of our collection described above are disease associated, compared to 26% of isolates from other lineages in the central population, and 19% of isolates from divergent lineages (Fig. 1C). Combined, the 10 lineages account for 80% of disease-associated isolates in our collection (Fig. 1C and Table 1) (19). These lineages also have much higher frequencies of serotypes that are commonly associated with disease: 88% of isolates from these 10 lineages have a disease-associated serotype (1, 1/2, 2, 3, 4, 5, 6, 7, 8, 9, or 14) compared to 16% from other lineages (Table 1 and Fig. 1D).

The 10 pathogenic lineages of *S. suis* fall within a subpopulation of the central *S. suis* population. Core nucleotide distances between these 10 lineages are, on average, lower than between lineages across the rest of the central population ([SI Appendix, Figs. S1 and S2](#)). Nevertheless, they are not clearly distinguished from other lineages in the central population; there are an additional 34 lineages (98 isolates) that fall within the clade that includes the 10 pathogenic lineages (Fig. 1B). Combined, these 34 lineages show an elevated rate of disease association: 54% of isolates are disease-associated compared to 26% from the central *S. suis* population outside of this clade (Fig. 1C). They are also enriched for pathogenic serotypes; 41% of isolates have a pathogenic serotype compared to 8% in other lineages of the central population

(Fig. 1D). This suggests that some of these lineages may also be pathogenic, but we cannot confidently characterise them as such due to small sample sizes (all have <15 isolates in the subset of our collection used for characterising disease association).

Consistent Genomic Changes Reveal Evolutionary Constraints on Pathogen Emergence. Previous studies that have investigated genes associated with pathogenicity in *S. suis* have tended to focus on isolates from the zoonotic lineage ST1, which is both highly pathogenic in pigs and responsible for most cases of human disease. Here, we aimed to identify genes associated with pathogenicity more broadly across *S. suis* (i.e., across all pathogenic lineages of *S. suis*). To do this, we identified genes that are present at >70% higher frequency in isolates from any of the 10 pathogenic lineages relative to both lineages from the central *S. suis* population outside of the pathogenic clade and divergent lineages. While no genes are uniquely present in pathogenic lineages, a small number (n = 37) are present at >70% higher frequencies in isolates from pathogenic lineages than both isolates from other lineages from the central population and from divergent lineages ([SI Appendix, Fig. S5 and Table S2](#)). Most of these genes (26/37) fall within just three genomic islands: Island 1 (SSU_RS05400-SSU_RS05325 in the published annotation of P1/7), Island 2 (SSU_RS02325-SSU_RS02355), and Island 3 (SSU_RS01130-SSU_RS01185) (Fig. 2 and [SI Appendix, Figs. S6–S8 and Table S2](#)). Five of the 11 genes outside of the three genomic islands have previously been described as associated with pathogenicity in *S. suis*, these are *mrrp*, *sspA*, *sspep*, SSU0587, and *sstgase* (20).

Island 1 is present in >95% of isolates from 9/10 pathogenic lineages, only 17% of isolates from the central population outside of the pathogenic clade (and these tend to be from lineages closely related to the pathogenic clade; Fig. 2A), and <1% of isolates from divergent lineages. The island contains genes predicted to encode proteins involved in the breakdown of host glycans present in the extracellular matrices of mammalian tissues. It includes a putative heparan sulfatase (SSU_RS05330 in P1/7) and a putative hyaluronate lyase (SSU_RS05335-SSU_RS05350 in P1/7) which cleave hyaluronic acid found in all connective body tissue. However, in zoonotic lineage 1, the gene encoding the hyaluronate lyase protein is always truncated, likely leading to loss of function. The phosphotransferase system (PTS) and other enzymes present in this island may be associated with sugar transport and metabolism of host aminoglycans which are degraded extracellularly by *S. suis*. The degradation of extracellular matrix proteins may both

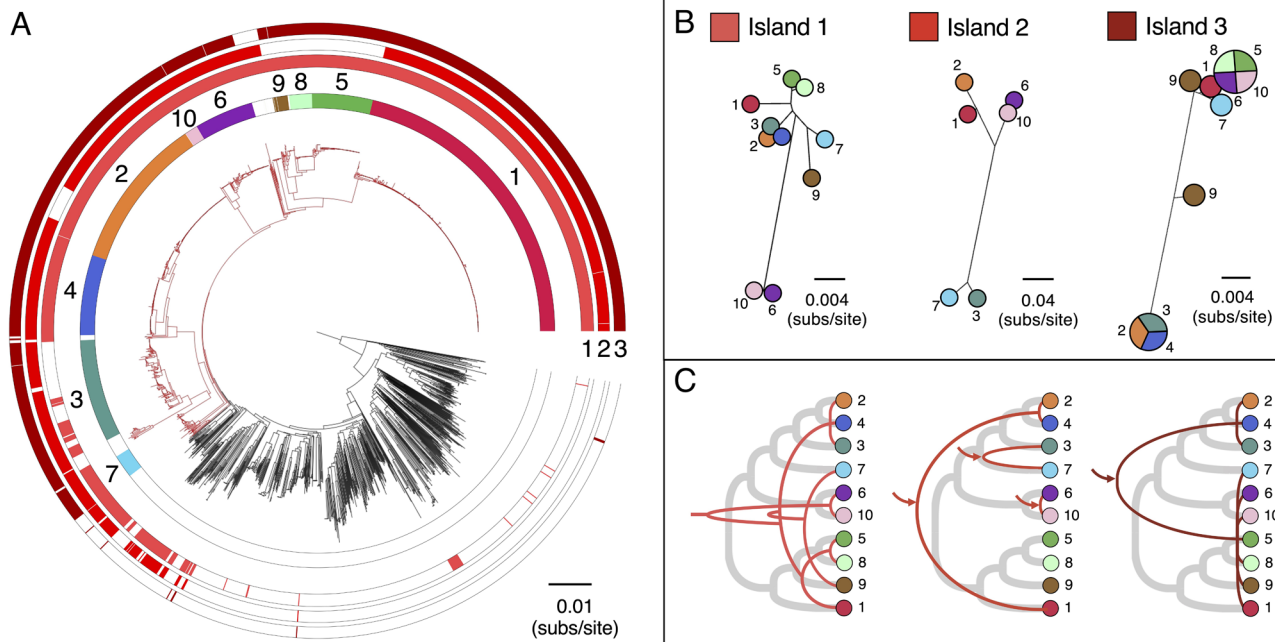


Fig. 2. Divergent evolutionary histories of the three pathogenicity-associated genomic islands. (A) The same core genome phylogeny as shown in Fig. 1B with three additional outer rings describing the presence of the three pathogenicity-associated genomic islands (1, 2, and 3). (B) Trees representing median genetic distances between pathogenic lineages based on coding regions of each of the three genomic islands (as described in *SI Appendix*, Figs. S6–S8). For Island 3, lineages that have average distances of zero are represented by the same circle split into segments and lineage 9 is represented by two circles due to the presence of two divergent versions of the island. (C) Cladograms representing inferred patterns of acquisition and inheritance of these three islands (red lines) compared to the core genome (thicker gray lines). Arrows indicate inferred acquisitions from outside of *S. suis*. A single inferred recombination event is omitted from the description of the history of Island 3 shown in C.

contribute to the spread of *S. suis* through the connective tissue and provide a source of energy at epithelial surfaces.

Nucleotide diversity within Island 1 is similar to and broadly correlated with diversity in core genes both within and between pathogenic lineages (Fig. 2 and *SI Appendix*, Figs. S9 and S10). The island is also found next to the same core gene (*nrdF* or SSU_RS05320 in P1/7) in 97% of isolates. Together, this suggests that the island was acquired once by a common ancestor of all pathogenic lineages. As homologous recombination is common in *S. suis* (23), the inheritance of this island is unlikely to have been entirely vertical, but our observations suggest that it has followed similar patterns to core genes (Fig. 2C).

Island 2 is present in 84% of isolates from the 10 pathogenic lineages, and only 11% from the central population outside of the pathogenic clade, where it tends to be carried by isolates that are both closely related to the pathogenic clade and also carry Island 1 (Fig. 2A). It is not found in any isolates from divergent lineages. Its presence is variable across pathogenic lineages: It is present in >80% of isolates from lineages 1, 2, 3, 4, 6, 7, and 10, but absent from lineages 5, 8, and 9. Island 2 encodes a major pilin subunit (SSU_RS02345), a minor pilin subunit (SSU_RS02335 and SSU_RS02340), and a pilin-specific sortase (SSU_RS02350) that is involved in maturation. The minor pilin subunit (SSU_RS02335 and SSU_RS02340) was previously described as a pseudogene in several pathogenic strains of *S. suis* (including P1/7) (24). We find evidence that this is true of all isolates in lineage 1 and 2 and several from other lineages due to premature stop codons. Nevertheless, Fittipaldi et al. (24) showed that despite this truncation, *S. suis* P1/7 still produces a major pilin subunit. While they found that this did not play a role in adherence to porcine brain microvascular endothelial cells or virulence in a mouse model of sepsis, Faulds-Pain et al. (25) later showed that the pilin-specific sortase in this island was essential for causing

disease in pigs via the intranasal route of infection. We find that in a small number of isolates (largely from lineages 3 and 7), the major pilin subunit is truncated (see *SI Appendix*, Fig. S6 for an example). Further work is required to understand the impact of this on the function of the island.

Diversity within Island 2 suggests multiple acquisitions from outside of *S. suis*: Three divergent versions are carried by pathogenic lineages (lineages 1 and 2 carry a similar version that is distinct from a version carried by lineages 6 and 10, and a version carried by lineages 3 and 7; Fig. 2B and C and *SI Appendix*, Fig. S9). The island is found next to the same core gene (*murD* or SSU_RS02355 in P1/7) in 99% of isolates, and divergence within each of the pathogenic lineages is similar to and correlated with diversity in both core genes and Island 1 (*SI Appendix*, Figs. S9 and S11). Divergence between the pairs of lineages that carry each of the three versions of the island is variable, and similar to distances between these lineages based on core genes and Island 1. This suggests that the presence of each version of Island 2 in a pair of pathogenic lineages is a consequence of a single acquisition by a common ancestor. As the most recent common ancestor of lineages 1 and 2 is also that of the entire pathogenic clade, this suggests that the version of this island carried by lineages 1 and 2 was maintained since the common ancestor of the pathogenic clade, while the version carried by lineages 6 and 10 has been acquired much more recently.

Island 3 shows the strongest association with pathogenic lineages. It is present in >95% of isolates in the 10 of the pathogenic lineages and only 1% of isolates from the central population outside of the pathogenic clade. It is entirely absent from divergent lineages. It contains genes that code for an ABC transporter, a ROK (repressor, open reading frame, kinase) family gene (26), and a large predicted surface protein with similarity to a PTSII subunit. The ROK is similar to an *E. coli* repressor (Mlc), which

represses several genes including two PTS genes (27). In *Salmonella Typhimurium*, Mlc positively regulates expression of the *Salmonella Typhimurium* pathogenicity island 1 genes by reducing the expression of the negative regulator HilE (28).

There are two divergent versions of Island 3 in our collection: one carried by lineages 1, 5, 6, 7, 8, 9, and 10 and the other by lineages 2, 3, and 4. While some isolates from lineage 9 carry a version that differs from these two, this appears to be the result of recombination between the two versions. In these isolates, there are extended tracts of SNPs that are shared with either of the two main versions, and there is only one unique SNP. Island 3 is found next to the same core gene in 98% of isolates in our collection (SSU_RS01115 in P1/7) and diversity in the island is positively correlated with diversity in core genes within lineages (SI Appendix, Fig. S12). Within-lineage diversity is generally lower in this island than in the other two islands (SI Appendix, Fig. S9), but this may reflect stronger selective constraint rather than a more recent origin. This is supported by lower divergence at 1st and 2nd codon positions relative to 3rd codon positions (SI Appendix, Fig. S13). Together, this suggests that the presence of Island 3 in each of the pathogenic lineages tends to be the result of a single acquisition by a common ancestor of the lineage. We find evidence of only one case of a second acquisition event by a pathogenic lineage: One isolate from lineage 2 carries the version of Island 3 associated with lineages 1, 5, 6, 7, 8, 9, and 10 (PH2016-139). While very low diversity within the two versions of the island means that the exact relationships between those carried by different lineages cannot be reconstructed, it indicates that for most pathogenic lineages Island 3 was obtained from another pathogenic lineage and that the transfer did not occur long before the most recent common ancestor of each lineage.

Each of the three pathogenicity-associated islands are sometimes present in lineages other than the 10 pathogenic lineages (Fig. 2A and SI Appendix, Table S1). All three are present together in only six additional lineages (15 isolates): lineages 18, 107, 144, 164, 466, and 511. Of these, three fall inside the pathogenic clade and three within the central population but outside of the

pathogenic clade. While the three in the pathogenic clade are all represented by multiple isolates, the three outside the pathogenic clade are all only represented by a single isolate. Combined, the six lineages have a broad geographic distribution that spans Europe, North America, Australia, and Asia. Lineage 18 alone includes isolates from Canada, the United States, and Myanmar, that date from 1987 to 2019. Some of these lineages may represent additional pathogenic lineages that are not well represented by our collection, while others may represent the chance acquisition of these islands by isolates that are otherwise ill-suited to a pathogenic ecology.

The Emergence and Spread of Pathogenic Lineages Are Linked to Growth in Pig Farming and International Trade. We used the temporal structure in our collection of isolates from the six most common pathogenic lineages to construct dated phylogenies (Fig. 3 and SI Appendix, Fig. S14). These reveal that the dates of the most recent common ancestors of these lineages range from 1827 (lineage 2; 95% HPD: 1798-1854) to 1951 (lineage 5; 95% HPD: 1944 to 1958) (SI Appendix, Table S3). These origins largely predate a rapid period of growth in pig numbers in several European countries that accompanied the wide-scale shift to larger farms and indoor rearing in the early 20th century (29, 30) and instead accompany a period of human population growth in Europe and North America that followed the Industrial Revolution (31) (SI Appendix, Fig. S15).

The 10 pathogenic lineages all have a broad geographic spread: Each is found in at least 6/11 countries in our collection. Similarly, outside of these lineages, we find little evidence of geographic structure in *S. suis*: Isolates from each well-sampled country span the diversity of *S. suis* (SI Appendix, Fig. S16). This is also true of isolates in our collection from a wild boar population in Spain (SI Appendix, Fig. S17). Using our dated phylogenies and a discrete asymmetric model, we estimated the number of between-country movements for each of the six most common pathogenic lineages and their relative rates. We estimated between 16 and 45 between-country movements for each lineage, and 182 across all

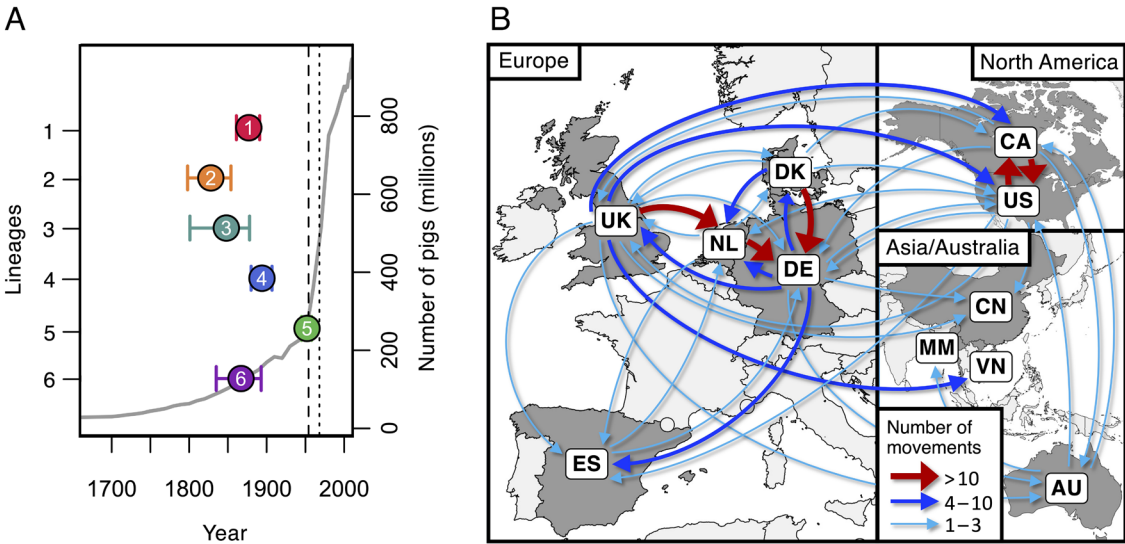


Fig. 3. Dates of emergence and paths of between-country transmission for the six most common pathogenic lineages. (A) Estimates of the dates of the most recent common ancestors of the six most common pathogenic lineages (coloured points) against an estimate of the global number of pigs (gray line). Country-specific estimates of pig numbers are shown in SI Appendix, Fig. S15. The vertical dashed line shows the date of the first reported case of *S. suis* disease in pigs (1954), and the dotted line shows the first reported human case (1968). (B) Map showing inferred routes of transmission of these six pathogenic lineages between the countries in our collection. Arrows represent routes with at least one inferred transmission event. Routes with more than ten inferred transmission events are shown in red, those with more than three in blue, and those with one to three in turquoise. Further details of the numbers and rates of movements between countries across our six lineages are shown in SI Appendix, Figs. S18-S20 and Table S4.

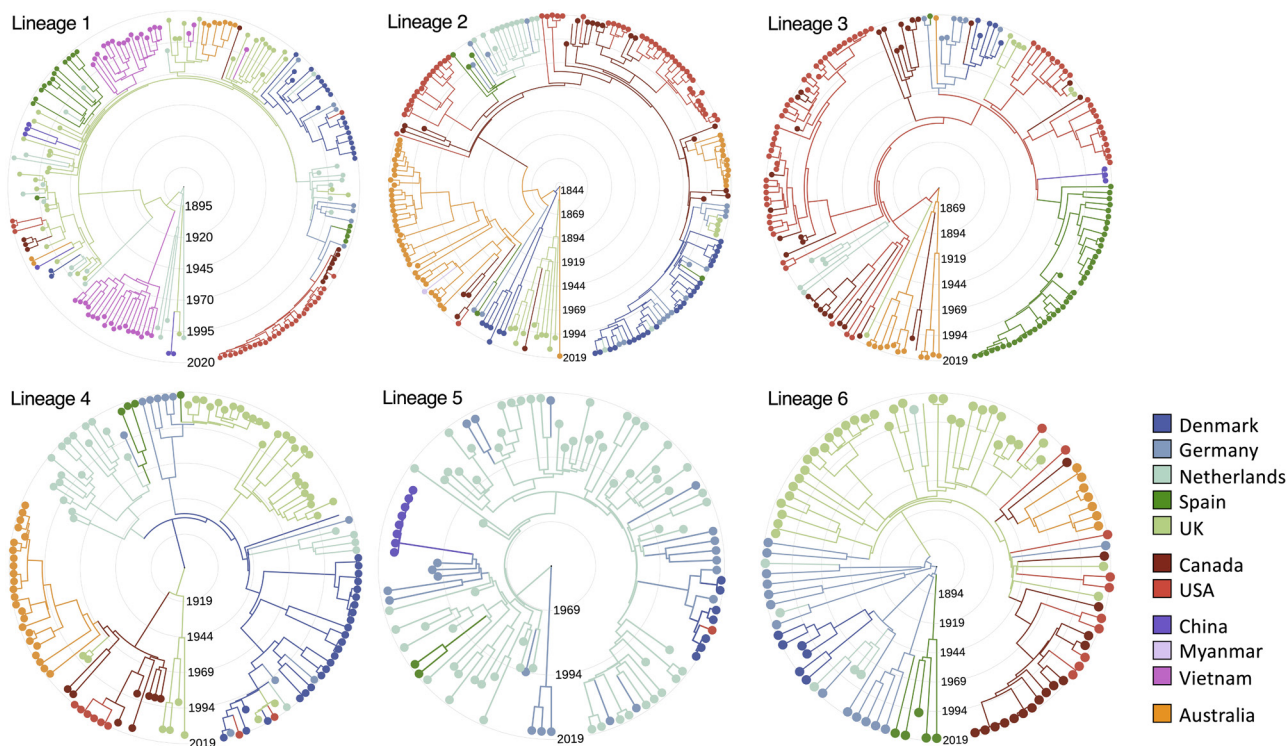


Fig. 4. Reconstructions of the emergence and spread of the six most common pathogenic lineages. Time-scaled phylogenies for the six most common pathogenic lineages of *S. suis* coloured by country of origin. Colours of branches represent the most likely ancestral location inferred from ancestral state reconstructions with an asymmetric discrete transition model in BEAST.

six lineages (*SI Appendix, Table S4*). This is likely to be a substantial underestimate of the true number of movements due to incomplete sampling.

Across all lineages, our highest inferred rate of between-country transmission was from Denmark to Germany, followed closely by Canada to the United States, the Netherlands to Germany, and the United States to Canada (Figs. 3 and 4 and *SI Appendix, Fig. S18 and Table S4*). This is consistent with these movements being driven by international trade in live pigs; over the last 80 y, the top three global exporters of live pigs have been Denmark, the Netherlands, and Canada, and top two importers have been Germany and the United States (www.fao.org/faostat). Additionally, we find no evidence of the transmission of pathogenic lineages into Australia since its ban on live pig imports in the mid-1980s (32). While we infer several transmission events from Europe and Canada to Australia, the most recent is a movement of zoonotic lineage 1 from the United Kingdom that we estimate to have occurred between 1981 (95% CI: 1976 to 1987) and 1986 (95% CI: 1981 to 1991).

We also observe an unusually low frequency of pathogenic lineages in our collection from Myanmar. Pig farming in Myanmar is typically small-scale and most of our samples from this country are from backyard or semi-intensive farms that typically have 10 to 30 pigs. While all of our isolates from Myanmar are nonclinical and so we would expect a low frequency of pathogenic lineages, they are less likely to be from pathogenic lineages than non-clinical isolates from other countries: <1% of isolates from Myanmar compared to a mean of 18% across other countries (*SI Appendix, Table S5*). In fact, none of our isolates from backyard and semi-intensive farms in Myanmar are from pathogenic lineages (0/109 isolates from the central *S. suis* population); the only pathogenic lineage isolate in our collection is from a larger commercial farm (1/39) with around 2,000 pigs.

As well as the spread of pathogenic lineages between pig farms in different countries, we find evidence of transmission between farmed pigs and wild boar. The two isolates from our collection of Spanish wild boar that are from a pathogenic lineage are both from lineage 1. They form a single clade that falls within a clade that appears to have been circulating in Spain since the 1970s. This indicates recent transmission between farmed pigs and local wild boar.

Individual lineages show different patterns of between-country transmission (Fig. 4 and *SI Appendix, Figs. S19 and S20*). This may reflect variation in both the prevalence of these lineages and their dates of introduction to particular countries. In particular, we find evidence that zoonotic lineage 1 has been circulating in the United Kingdom and the Netherlands for at least 50 y, with frequent transmission between these two countries, while the earliest evidence we find of its presence in North America is 2009 (95% CI: 2007 to 2011). This is consistent with this lineage both being less commonly associated with disease in North America and reports of an increase in its rate of detection over the last few years. In contrast, lineage 2 shows evidence of circulation in Canada for around 80 y and repeated transmission from Canada to the USA. Unfortunately, uncertainty in the inferred locations of the common ancestors of each of the lineages means that we cannot confidently infer the country of origin for any pathogenic lineage.

Pathogenic Lineages Have Variable Associations with Systemic Disease and Carry Unique Genes. Genetic and phenotypic diversity can make it more difficult to control the spread of a pathogen. For example, it can lead to difficulties in developing cross-protective vaccines. It can also make a pathogen more capable of evolving to evade control measures. In *S. suis*, we find evidence of both phenotypic and genetic variation between

pathogenic lineages. The six most common pathogenic lineages in our collection vary in their frequency of both association with disease relative to carriage (chi-squared test, $P = 3.6 \times 10^{-8}$) and systemic disease relative to respiratory disease (chi-squared test, $P = 2.0 \times 10^{-20}$). Pairwise comparisons reveal two distinct groups: lineages 1, 3, 4, and 5 have higher frequencies of association with systemic disease, while lineages 2 and 6 have higher frequencies of association with respiratory disease (SI Appendix, Fig. S21). This is consistent with previous studies that found that ST28, which is part of lineage 2, is less virulent in mouse models than both ST1 (lineage 1) and ST25 (lineage 4) (33).

Pathogenic lineages also show evidence of variation in their evolutionary dynamics and genome content. Average rates of nucleotide substitutions per site were highest for lineage 5, while the ratio of transitions to transversions is lower for lineages 2 and 4 (SI Appendix, Fig. S22). The six most common pathogenic lineages can be distinguished by their carriage of 2 to 32 lineage-specific genes (that are present in >95% of isolates in that lineage and <5% of isolates from other lineages; SI Appendix, Table S6). For instance, 17 genes distinguish zoonotic lineage 1. Far more genes are associated with multiple pathogenic lineages; 309 genes are present in >95% of isolates in at least one of the six most common pathogenic lineages and in <5% of isolates outside of the pathogenic clade (SI Appendix, Table S6). Apart from lineages 2 and 4 that share more genes than other pathogenic lineages, there is little evidence of clustering of lineages by shared gene content (SI Appendix, Fig. S23).

Pathogenic Lineages Continue to Diverge from Commensal Lineages and Diversify. Pathogenic lineages of *S. suis* tend to have smaller genomes than those from outside of the pathogenic clade, and pathogenic lineages with higher frequencies of association with disease tend to have smaller genomes than those with higher frequencies of association with carriage (Fig. 5A). This pattern was previously described in a study based on a smaller sample of *S.*

suis isolates (22). Using the temporal structure in our collection, we additionally found that older isolates of pathogenic lineages tend to have larger genomes than more recent isolates. This suggests a gradual and ongoing process of genome reduction in pathogenic lineages (Fig. 5B). While all pathogenic lineages show this pattern, we observe variation in the rate of genome reduction across lineages, with zoonotic lineage 1 showing the slowest rate and lineage 5 the fastest.

Against this backdrop of genome reduction, we found evidence of adaptive gene acquisitions. In particular, we found evidence of multiple capsular switches in all of the six most common pathogenic lineages (Fig. 1 and SI Appendix, Fig. S24). Serotype 2 is often linked with invasive disease in pigs and zoonotic disease in humans, particularly when carried by ST1 (zoonotic lineage 1). We found that the serotype 2 capsule locus is present in the genomes of the majority of isolates from lineage 1 and also in isolates from 3/9 of the other pathogenic lineages (lineages 2, 4, and 9) and a few other lineages within the pathogenic clade (SI Appendix, Table S1). While there is little diversity within the genes in the serotype 2 capsular locus either within or between lineages, there is greater diversity within zoonotic lineage 1 than within other lineages (SI Appendix, Fig. S25). Combined with an older estimate of the most recent common ancestor of lineage 1 (1876, 95% CI: 1860-1891) compared to the clades that carry serotype 2 in lineages 2 (1935, 95% CI: 1923 to 1948) and 4 (1914, 95% CI: 1902 to 1926), this suggests that this serotype has been horizontally transferred from lineage 1 to these other pathogenic lineages. We also observe evidence of repeated transitions between serotypes 2 and 1/2 [whose capsular loci are nearly identical (34)]. These transitions appear to be particularly common in lineage 2 (SI Appendix, Figs. S24 and S25).

Serotypes 1 and 14 are also common in zoonotic lineage 1. The genes encoding these capsules are largely shared with serotypes 2 and 1/2 (34), but divergence within the capsular genes is large enough to indicate independent origins of these two pairs of

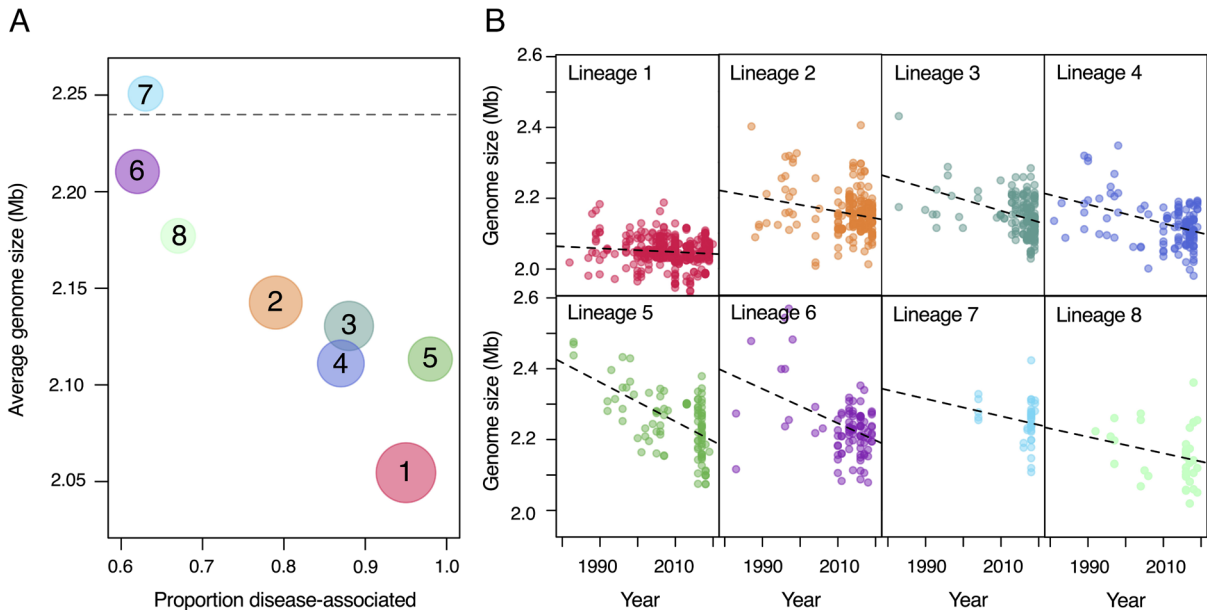


Fig. 5. Gradual and ongoing genome reduction in pathogenic lineages. (A) Average genome sizes for recently sampled isolates (2018 to 2020) are shown against the proportion of disease-associated isolates for pathogenic lineages 1 to 8. Pathogenic lineages 9 and 10 are not shown as we have no or very few recent isolates from these lineages in our collection. The size of the points reflects the number of isolates in our collection on a log-scale. The dashed line represents the average genome size of isolates of *S. suis* from outside of the pathogenic clade, where disease isolates are present at a frequency of approximately 25%. (B) Genome sizes of individual isolates against sampling year for lineages 1 to 8. Dashed lines represent best fits for linear models of genome size against year of isolation; all have a negative slope, with a gradient ranging from a loss of 516 bases per year (lineage 1) to a loss of 5,550 bases per year (lineage 5).

serotypes (*SI Appendix, Fig. S24*). Serotype 14 is also present in lineage 2 and in two other lineages (lineages 11 and 78). Diversity within the capsular genes of serotypes 1 and 14 is again consistent with them having been horizontally transmitted from zoonotic lineage 1 to other lineages and repeated transitions between serotypes 1 and 14.

Discussion

Studies of the impact of agricultural intensification on the risk of pathogen emergence generally focus on the most typical route: a pathogen jumping from one host species to another (4, 7, 35). In this study, we instead described the emergence of an important zoonotic pathogen from within a largely commensal member of the respiratory microbiota of pigs. Our results suggest that while this form of pathogen emergence is gradual, it can lead to a diverse pathogen whose impact on its host population is difficult to control, particularly alongside efforts to reduce the use of antibiotics in farming.

Over the last two hundred years global pig numbers have increased more than 10-fold (31). While the global rate of increase was highest in the second half of the 20th century, in some regions, such as the United States, growth was faster in the 19th century (*SI Appendix, Fig. S15*). Population growth has been accompanied by increased population density. For example, between 1921 and 2011 the number of pigs in Canada increased from 3.3 to 14.6 million, while the number of farms declined from 453 to just over 7 thousand (36). Modern farming practices also involve frequent movement of pigs between farms, with a trend toward production specialisation meaning that piglets often move to a different farm after weaning, and breeding for genetic improvement leading to the transport of breeding pigs around the world. These kinds of changes are predicted to promote pathogen emergence by facilitating the transmission of pathogens between hosts and therefore reducing the selective cost associated with increased host morbidity (8, 37).

Our analyses reveal high levels of genetic diversity in *S. suis* in pigs, that is mirrored in isolates sampled from a wild boar population in Spain. This diversity might reflect a long-standing association between *S. suis* and pigs (both domestic and wild). Our dating of the most recent common ancestors of the six most common pathogenic lineages in our collection indicates that they all emerged in the 19th and 20th centuries. The conclusion that these dates reflect an ecological shift toward pathogenicity in at least some of these lineages is supported by evidence that they coincided with the acquisition of a pathogenicity-associated genomic island (Island 3). It is further supported by patterns of genome reduction in each of the pathogenic lineages. In comparisons across bacterial species, it has been shown that bacterial pathogenicity is broadly associated with smaller genomes and fewer genes (22). While the drivers of this pattern are not well understood, and may be multiple, our observation of a gradual decline in the genome sizes of isolates from pathogenic lineages over our sampling period suggests a recent transition whose effect has yet to stabilise. Further evidence of an impact of intensive farming practices on pathogenic lineages of *S. suis* is found in our estimates of frequent long-distance movements, which are most likely to be a consequence of international trade in live animals.

While our results suggest that pathogenic lineages have emerged and spread globally over the last two centuries, they are also consistent with pathogenicity in *S. suis* predating this. Genetic diversity in all three of the pathogenicity-associated genomic islands we identified is consistent with their presence in *S. suis* long before the origins of individual pathogenic lineages. Genetic distances between the versions of Islands 1, 2, and 3 that are carried by lineages 1 and 2 are similar to distances between these lineages

estimated from core genes. This may reflect the presence of these islands in a common ancestor of these two lineages. The prior existence of pathogenic lineages could therefore have aided the recent emergence of several new pathogenic lineages of *S. suis*.

Traits that promote virulence are thought to be selected for because they increase within-host growth or between-host transmission (38). The pathogenicity-associated genomic islands we identified have putative functions that may influence patterns of within-host growth. Islands 1 and 3 both have putative functions linked to metabolism, either the capacity to exploit particular sources of sugar within a host, or to regulate their metabolism. Metabolic capacity and growth rate have been linked to virulence in several bacterial species (39). The maintenance of both commensal and pathogenic lineages of *S. suis* could therefore be a consequence of a partitioning of the within-host niche. Pathogenic lineages may be better able to exploit particular regions of the tonsil than commensal lineages and vice versa, thereby reducing within-host competition. This could lead to segregation of these populations and reduced gene flow between them, which could in turn lead to the genome reduction in more pathogenic lineages due to fewer opportunities for gene acquisition from more diverse commensal lineages. This partitioning of the within-host niche may also be aided by Island 2. Pili are often involved in adhesion and evasion of cells, and therefore, this island might aid in the colonization of a particular region of the tonsil.

As we observe a high rate of spread of pathogenic lineages between countries, and previous studies have found that pathogenic strains of *S. suis* are capable of spreading rapidly between pigs within a farm (40), it is likely that pathogenic lineages have faster rates of between-host transmission than commensal lineages. As none of the pathogenicity-associated genes we identify have putative functions that are likely to be directly associated with between-host transmission, faster transmission rates may instead be driven by a trait with a diverse molecular basis, such as the capsular polysaccharide. In the human pathogen *Streptococcus pneumoniae*, the capsular structure is known to both protect against immune recognition and immune clearance within a host (41), and promote between-host transmission through aiding survival outside a host (42), and increasing shedding (43).

The pathogenic lineages we have identified carry a mosaic of genes shared with other pathogenic and commensal lineages of *S. suis*. This is accompanied by diversity at key loci, such as the capsular locus, and in their association with systemic vs. respiratory disease. Our results reveal that the emergence of novel pathogenic lineages therefore has not only led to more pathogenic lineages but also to their diversification. This diversification of pathogenic lineages has led to difficulties in developing a cross-protective *S. suis* vaccine (44). We also find evidence of horizontal gene transfer between lineages leading to further diversification. In particular, we find evidence that the serotype 2 and serotype 1/14 capsules have been transferred from zoonotic lineage 1 (ST1) to other lineages. While we have been unable to explore this with our collection, which includes only a small sample of isolates from human infections, previous studies have suggested that strains of *S. suis* vary in their capacity to infect humans, and most human infections are caused by strains with a serotype 2 or 14 capsule (19). These acquisitions of the serotype 2 and 14 capsular loci may therefore have increased the zoonotic potential of these lineages.

Our results provide a framework for understanding the genomic diversity in *S. suis* and its association with pathogenicity. This is likely to be of widespread use in *S. suis* research and in informing strategies for controlling the burden of this disease on pig farming and human health. As our collection spans only a small proportion of the countries that farm pigs globally, further sampling from a broader range

countries and more extensive sampling within countries, particularly those with large and growing pig populations, is needed to investigate the existence of additional pathogenic lineages that are geographically restricted or have recently emerged. While our analyses suggest that the global movement of *S. suis* is driven by trade, they also suggest a possible role for wild boar. Further sampling of wild boar populations is needed to determine the role of wild boar in the transmission of pathogenic lineages between pig farms and into human populations. Further research is also required to experimentally characterise the functions of the pathogenicity-associated genomic islands we have identified and establish their relationship with pathogenicity. Our results suggest that they represent evolutionary constraints on the emergence of pathogenic lineages: The emergence of new pathogenic lineages is contingent on the horizontal transfer of genes from an existing pathogenic lineage to another susceptible lineage. The conditions that allow for the spread of pathogenic lineages may therefore also promote the emergence of new pathogenic lineages and the diversification of existing pathogenic lineages through generating opportunities for the transfer of pathogenicity-associated genes. While this process has so far been gradual, we might predict that the pace will increase. The intensification of pig farming is ongoing in some regions, and the common pathogenic lineages we have identified have only recently spread to some parts of the world. This could lead to both an expanding niche for pathogenic lineages and more opportunities for the emergence of new pathogenic lineages. Controlling the spread of pathogenic lineages of *S. suis* through pig populations should therefore be a priority to limit the potential impact of this pathogen on our future food security and public health.

Methods

Sampling of Isolates and Characterisation of Disease Association. We generated a collection of genome assemblies of 3,070 isolates of *S. suis* and its close relatives including both isolates that we sequenced and previously published data (*SI Appendix, Table S1*). The read data generated by this study is publicly available from the SRA; the BioProject IDs are provided in *SI Appendix, Table S1*. Our collection includes isolates from Australia, Canada, China, Denmark, Germany, Myanmar, the Netherlands, Spain, the United Kingdom, the United States, and Vietnam. We also included 29 published reference genomes.

Isolates from Denmark, Germany, the Netherlands, and Spain are largely from two newly sequenced collections. The first aimed at collecting similar numbers of systemic disease, respiratory disease, and carriage isolates from pigs from each country and were sampled from 2014 to 2018 ($n = 593$). The second was a sample of historic isolates that aimed at capturing the breadth of disease-associated strains from each country from sample archives; they date from 1960 to 2007 ($n = 99$). Isolates from the Netherlands also included a published collection of systemic disease isolates from humans and pigs that date from 1982 to 2008 ($n = 97$) (45). Isolates from Spain also included newly sequenced isolates from a herd of wild boars sampled in 2015 as part of a published study ($n = 41$) (46).

Isolates from the United States are from a newly sequenced collection of isolates from pigs from 2017 to 2020. Isolates were obtained from clinical cases submitted to the Iowa State University Veterinary Diagnostic Laboratory for routine diagnostics. Isolates from Myanmar are from another newly sequenced collection of isolates from pigs from farms and slaughterhouses in Yangon. They were sampled from pig farms that included small backyard farms, small-scale traditional farms, and modern industrial farms between 2016 and 2019. They were predominantly from throat swabs from pigs without *S. suis* disease but were also sampled from farm and slaughterhouse drainage systems.

Isolates from the United Kingdom came from four different collections. The first sampled nonclinical and clinical isolates from pigs across England and Wales in 2009–2011 (described in ref. 47) ($n = 193$). The second in 2013 to 2014 sampled nonclinical isolates from five farms (described in ref. 48) ($n = 117$). The third in 2013 to 2014 targeted clinical isolates from pigs across England and Wales (described in ref. 49) ($n = 129$). The fourth is a newly sequenced collection of archived clinical isolates from pigs that date from 1987 to 2000 ($n = 49$).

Isolates from Vietnam are from a collection that aimed at sampling closely related populations from humans and pigs (described in ref. 47). These included systemic disease isolates ($n = 153$) from human clinical cases of meningitis from provinces in southern and central Vietnam, and systemic disease ($n = 6$) or nonclinical isolates ($n = 32$) from pigs, collected between 2000 and 2010. These isolates were exclusively serotype 2 or 14. Isolates from Canada are from a previously published collection that aimed to target similar numbers of clinical and nonclinical isolates and dates from 1983 to 2016 (described in ref. 50).

Modern isolates from pigs from Denmark, Germany, the Netherlands, Canada, the United Kingdom, and Vietnam were characterised as associated with systemic or respiratory disease, or nonclinical carriage based on clinical symptoms and the isolation site. In pigs that showed clinical symptoms consistent with *S. suis* infections (e.g., meningitis, septicemia, and arthritis), the site of isolation was classified as "systemic" if recovered from systemic sites (i.e., brain, liver, blood, joints). The site of recovery was classified as "respiratory" if derived from lungs with gross lesions of pneumonia. *S. suis* isolates from the tonsils, nose, or tracheo-bronchus of healthy pigs or dead pigs without any typical signs of *S. suis* infections were defined as "nonclinical". Isolates that could not confidently be assigned to these categories (e.g., a tonsil isolate from a pig with systemic signs) were classified as unknown.

Whole Genome Sequencing and Assembly. Illumina whole genome sequencing was undertaken for all newly sequenced isolates. For isolates from Europe and the United Kingdom, DNA extraction, library preparation, and sequencing were undertaken using a HiSeq 2500 instrument (Illumina, San Diego, CA, USA) by MicrobesNG (Birmingham, UK). For isolates from the United States, multiplex genome libraries were prepared using the Nextera XT DNA library preparation kit (Illumina, San Diego, CA, USA). The genomic library was quantified using a Qubit fluorometer dsDNA HS kit (Life Technologies Carlsbad, CA, USA) and normalized to the recommended amplification concentrations. The pooled libraries were sequenced on an Illumina Miseq sequencer using Miseq Reagent V3 for 600 cycles (Illumina, San Diego, CA, USA). Raw reads were demultiplexed automatically on the Miseq.

Raw sequence reads were preprocessed using Trimmomatic V0.36 to remove adaptors, trim poor-quality ends, and delete short sequences (<36 nt) (51). Raw and preprocessed reads were assessed for quality to ensure cleaning efficiency using FastQC (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc>). We generated de novo assemblies with Spades v.3.12.0 (52). All assemblies were evaluated using QUAST v.5.0.1 (53) and we mapped reads back to de novo assemblies to investigate polymorphism (indicative of mixed cultures) using Bowtie2 v1.2.2 (54). Low-quality assemblies were excluded from the collection.

We assembled high-quality reference genomes for 19 isolates. For 12/19 isolates, long-read sequencing library preparation was performed using Genomic-tips and a Genomic Blood and Cell Culture DNA Midi kit (Qiagen, Hilden, Germany). Sequencing was performed on the Sequel instrument from Pacific Biosciences using v2.1 chemistry and a multiplexed sample preparation. Reads were demultiplexed using Lima in the SMRT link software (<https://github.com/PacificBiosciences/barcoding>). Reads shorter than 2,500 bases were removed using prinseq-lite.pl (<https://sourceforge.net/projects/prinseq>). Hybrid assemblies, using filtered PacBio and Illumina reads, and preliminary assemblies of long-read data assembled with Canu v1.9 (55) were generated with Unicycler v0.4.7 using the normal mode and default settings (56). Assembly graphs were visualised and, if necessary, manually corrected with Bandage (57). For the remaining 7/19 isolates, library preparation, short-read and long-read sequencing, and hybrid assembly with Unicycler were undertaken by MicrobesNG as part of their Enhanced Genome Service, which uses both Illumina and Oxford Nanopore Technologies. Four of these complete assemblies were published in a previous study (58).

Serotyping and Sequence-Typing. Serotypes and STs were determined in silico using the Athey et al. (59) pipeline.

Identification of Homologous Genes and Analysis of Pathogenicity-Associated Genomic Islands. All genomes were annotated using Prokka v1.14.5 (60). Panaroo v1.2.2 (61) was used to identify orthologous genes (using recommended parameter settings) and create alignments of core genes. Pathogenicity-associated genomic islands were identified through comparison of frequencies of homologous genes identified by Panaroo and analysis of their relative genomic locations. Concatenated alignments of genes within each genomic island were generated and checked by eye. Regions too divergent to

align were excluded from further analysis. Distance matrices were estimated and a neighbour-joining tree was created using the *ape* package in R (62). Trees were visualised using iTOL and Grapetree (63, 64).

Analysis of Population Structure, Phylogeny, and Geographic Spread. We used PopPunk to cluster our genomes (65). We first used the software to identify divergent genomes and then used the pipeline to cluster genomes in the central population into lineages. We used the standard pipeline followed by refinement using only core distances to determine lineages.

We generated reference-mapped assemblies of the six most common pathogenic lineages with Bowtie2 using reference genomes within the lineages. We identified recombination using Gubbins v2.3.1 (66) and masked all identified recombinant sites in our alignments. We tested for temporal signal in each lineage using a regression of root-to-tip distances against sampling year using the trees output from Gubbins and a published R script (67). Roots were chosen so as to minimise the residual mean squares of a linear regression. For lineage 1, we downsampled our data from 530 to 200 isolates that were randomly selected under the constraint of maintaining the temporal and geographic breadth of the collection. For each lineage, we tested for temporal signal using 1,000 random permutations of dates over clades sampled from the same year to account for any confounding of temporal and genetic structure. This yielded significant evidence ($P < 0.05$) of temporal signal in all lineages except for in lineage 2. Nevertheless, as our estimate of evolutionary rate for lineage 2 was similar to lineages 1 and 3, we considered it likely that there was sufficient temporal signal to inform our estimates of dates.

We constructed dated phylogenies using BEAST v1.10 with a HKY+ Γ model, a strict molecular clock, and an exponential population size coalescent model (68). We also constructed dated phylogenies using a relaxed molecular clock, a constant population size model, and a skyline population model, and found that our results were robust to different model choices. We undertook ancestral state reconstruction in BEAST to infer the geographic spread of these lineages by fitting an asymmetric discrete traits model to the posterior distributions of trees, with each state representing a country.

Data, Materials, and Software Availability. DNA sequence reads, genome assemblies and annotations have been deposited in the NCBI Sequence Read Archive and are available under BioProjects [PRJNA1009400](#), [PRJNA1012585](#), [PRJNA1010137](#) and [PRJNA972671](#) (supporting information). Previously published data were used for this work (supporting information).

ACKNOWLEDGMENTS. This work was primarily funded by an EU Horizon 2020 grant "PIGSs" (727966) and a ZELS BBSRC award "Myanmar Pigs Partnership (MPP)" (BB/L018934/1). G.G.R.M., E.L.M., and L.A.W. were supported by a Sir Henry Dale Fellowship to L.A.W. jointly funded by the Wellcome Trust and the Royal Society (109385/Z/15/Z). N.H. was supported by a Challenge grant from the Royal Society (CH16011) and an Isaac Newton Trust Research Grant [17.24(u)]. G.G.R.M. was also supported by a Research Fellowship at Newnham College. S.B. is supported by the Medical Research Council (MR/V032836/1). PIC North America provided part of the funds for the sequencing of the isolates from the USA. A.J.B. and M.M. were funded by Medical Research Council and Biotechnology and Biological Sciences Research Council studentships respectively, and M.M. was co-funded by the Raymond and Beverly Sackler Fund. We would like to acknowledge Susanna Williamson at the APHA for providing samples, Oscar Cabezon for sampling of the wild boar population in Spain, Mark O'Dea for access to sequence data from Australian isolates, the PIGSs and MPP consortiums for providing samples and helpful discussions, Julian Parkhill and John Welch for helpful discussions, and two anonymous reviewers for their valuable suggestions for improving the manuscript. This research was funded in whole or in part by the Wellcome Trust. For the purpose of Open Access, the author has applied a CC BY public copyright license to any Author Accepted Manuscript (AAM) version arising from this submission.

Author affiliations: ^aDepartment of Genetics, Evolution and Environment, University College London, London WC1E 6BT, United Kingdom; ^bDepartment of Veterinary Medicine, University of Cambridge, Cambridge CB3 0ES, United Kingdom; ^cCancer Research UK Manchester Institute, University of Manchester, Manchester M20 4BX, United Kingdom; ^dDepartment of Biology, Haverford College, Haverford, PA 19041; ^eDepartment of Medicine, University of Cambridge, Cambridge CB2 2ZQ, United Kingdom; ^fCentre for Enzyme Innovation, University of Portsmouth, Portsmouth PO1 2DD, United Kingdom; ^gNuffield Department of Population Health, University of Oxford, Oxford OX3 7LF, United Kingdom; ^hCollege of Veterinary Medicine, Iowa State University, Ames, IA 50011; ⁱAnimal Sciences Department, Wageningen University, 6700 AH Wageningen, The Netherlands; ^jOxford University Clinical Research Unit, Ho Chi Minh City, Vietnam; ^kLivestock Breeding and Veterinary Department, Yangon, Myanmar; ^lDépartement de Pathologie et Microbiologie, Université de Montréal, Québec J2S 2M2, Canada; ^mUnitat Mixta d'Investigació IRTA-UAB en Sanitat Animal, Centre de Recerca en Sanitat Animal, Campus de la Universitat Autònoma de Barcelona, Barcelona 08193, Spain; ⁿOIE Collaborating Centre for the Research and Control of Emerging and Re-Emerging Swine Diseases in Europe (IRTA-CReSA), Barcelona 08193, Spain; ^oInstitute for Microbiology, University of Veterinary Medicine Hannover, Hannover 30559, Germany; ^pDepartment of Health Technology, Technical University of Denmark, Kgs. Lyngby 2800, Denmark; ^qWageningen Bioveterinary Research, 8221 RA Lelystad, The Netherlands; ^rCentre for Tropical Medicine, Nuffield Department of Medicine, University of Oxford, Oxford OX3 7LG, United Kingdom; ^sMicrobiology Department and Center for Tropical Medicine Research, Ngoc Thach University of Medicine, Ho Chi Minh City, Vietnam; and ^tParasites and Microbes Programme, Wellcome Sanger Institute, Cambridge CB10 1RQ, United Kingdom

1. Y. M. Bar-On, R. Phillips, R. Milo, The biomass distribution on Earth. *Proc. Natl. Acad. Sci. U.S.A.* **115**, 6506–6511 (2018).
2. Food and Agriculture Organization of the United Nations Statistics Division (FAOSTAT), Food and agriculture data. <https://www.fao.org>. 2 March 2023.
3. W. Gilbert, L. F. Thomas, L. Coyne, J. Rushton, Review: Mitigating the risks posed by intensification in livestock production: The examples of antimicrobial resistance and zoonoses. *Animal* **15**, 100123 (2021).
4. B. A. Jones *et al.*, Zoonosis emergence linked to agricultural intensification and environmental change. *Proc. Natl. Acad. Sci. U.S.A.* **110**, 8399–8404 (2013).
5. A. Engering, L. Hogewerf, J. Slingenbergh, Pathogen-host-environment interplay and disease emergence. *Emerg. Microbes Infect.* **2**, e5 (2013).
6. K. E. Jones *et al.*, Global trends in emerging infectious diseases. *Nature* **451**, 990–993 (2008).
7. H. Bartlett *et al.*, Understanding the relative risks of zoonosis emergence under contrasting approaches to meeting livestock product demand. *R. Soc. Open Sci.* **9**, 211573 (2022).
8. A. Mennerat, F. Nilsen, D. Ebert, A. Skorpington, Intensive farming: Evolutionary implications for parasites and pathogens. *Evol. Biol.* **37**, 59–67 (2010).
9. H. Field, Studies on piglet mortality. Streptococcal meningitis and arthritis. *Vet Rec.* **66**, 453–455 (1954).
10. D. Vötsch, M. Willenborg, Y. B. Weldearegay, P. Valentin-Weigand, *Streptococcus suis*—The "Two Faces" of a pathobiont in the porcine respiratory tract. *Front. Microbiol.* **9**, 480 (2018).
11. Z.-R. Lun, Q.-P. Wang, X.-G. Chen, A.-X. Li, X.-Q. Zhu, *Streptococcus suis*: An emerging zoonotic pathogen. *Lancet Infect. Dis.* **7**, 201–209 (2007).
12. S. Fredriksen *et al.*, Environmental and maternal factors shaping tonsillar microbiota development in piglets. *BMC Microbiol.* **22**, 224 (2022).
13. T. Opriessnig, L. G. Giménez-Lirola, P. G. Halbur, Polymicrobial respiratory disease in pigs. *Anim. Health Res. Rev.* **12**, 133–148 (2011).
14. B. Perch, P. Kristjansen, K. Skadhauge, Group R streptococci pathogenic for man. Two cases of meningitis and one fatal case of sepsis. *Acta Pathol. Microbiol. Scand.* **74**, 69–76 (1968).
15. J. Tang *et al.*, Streptococcal toxic shock syndrome caused by *Streptococcus suis* serotype 2. *PLoS Med.* **3**, e151 (2006).
16. H. Yu *et al.*, Human *Streptococcus suis* outbreak, Sichuan, China. *Emerg. Infect. Dis.* **12**, 914–920 (2006).
17. P. Praphasiri *et al.*, *Streptococcus suis* infection in hospitalized patients, Nakhon Phanom Province, Thailand. *Emerg. Infect. Dis.* **21**, 345–348 (2015).
18. V. T. L. Huong *et al.*, Epidemiology, clinical manifestations, and outcomes of *Streptococcus suis* infection in humans. *Emerg. Infect. Dis.* **20**, 1105–1114 (2014).
19. G. Goyette-Desjardins, J.-P. Auger, J. Xu, M. Segura, M. Gottschalk, *Streptococcus suis*, an important pig pathogen and emerging zoonotic agent—an update on the worldwide distribution based on serotyping and sequence typing. *Emerg. Microbes Infect.* **3**, e45 (2014).
20. M. Segura, N. Fittipaldi, C. Calzas, M. Gottschalk, Critical *Streptococcus suis* virulence factors: Are they all really critical? *Trends Microbiol.* **25**, 585–599 (2017).
21. A. Baig *et al.*, Whole genome investigation of a divergent clade of the pathogen *Streptococcus suis*. *Front. Microbiol.* **6**, 1191 (2015).
22. G. G. R. Murray *et al.*, Genome reduction is associated with bacterial pathogenicity across different scales of temporal and ecological divergence. *Mol. Biol. Evol.* **38**, 1570–1579 (2021).
23. I. P. Lee, C. P. Andam, Frequencies and characteristics of genome-wide recombination in *Streptococcus agalactiae*, *Streptococcus pyogenes*, and *Streptococcus suis*. *Sci. Rep.* **12**, 1515 (2022).
24. N. Fittipaldi *et al.*, Mutations in the gene encoding the ancillary pilin subunit of the *Streptococcus suis* srtF cluster result in pili formed by the major subunit only. *PLoS One* **5**, e8426 (2010).
25. A. Faulds-Pain *et al.*, The *Streptococcus suis* sortases SrtB and SrtF are essential for disease in pigs. *Microbiology (Reading)* **165**, 163–173 (2019).
26. M. D. Kazanov, X. Li, M. S. Gelfand, A. L. Osterman, D. A. Rodionov, Functional diversification of ROK-family transcriptional regulators of sugar catabolism in the Thermotogae phylum. *Nucleic Acids Res.* **41**, 790–803 (2013).
27. J. Plumbridge, Regulation of gene expression in the PTS in *Escherichia coli*: The role and interactions of Mlc. *Curr. Opin. Microbiol.* **5**, 187–193 (2002).
28. S. Lim *et al.*, Mlc regulation of Salmonella pathogenicity island I gene expression via hIle repression. *Nucleic Acids Res.* **35**, 1822–1832 (2007).
29. A. Woods, Rethinking the history of modern agriculture: British pig production, c.1910–65. *20th Century Br. Hist.* **23**, 165–191 (2012).
30. M. R. Finlay, "Hogs, antibiotics, and the industrial environments of postwar agriculture" in *Industrializing Organisms* (Routledge, 2003).
31. K. K. Goldewijk, Total Number of Pigs | Clío Infra | Reconstructing global inequality. (2013). (15 March 2023).
32. Productivity Commission Government of Australia, Pig and pigmeat industries: Safeguard action against imports. Australian Productivity Commission Working Paper. No 1557. <https://doi.org/10.2139/ssrn.278061> (1998).
33. N. Fittipaldi *et al.*, Lineage and virulence of *Streptococcus suis* serotype 2 isolates from North America. *Emerg. Infect. Dis.* **17**, 2239–2244 (2011).

34. M. Okura *et al.*, Genetic analysis of capsular polysaccharide synthesis gene clusters from all serotypes of *Streptococcus suis*: Potential mechanisms for generation of capsular variation. *Appl. Environ. Microbiol.* **79**, 2796–2806 (2013).
35. M. E. J. Woolhouse, D. T. Haydon, R. Antia, Emerging pathogens: The epidemiology and evolution of species jumps. *Trends Ecol. Evol.* **20**, 238–244 (2005).
36. Statistics Canada: Canada's National Statistical Agency, The changing face of the Canadian hog industry, <https://www150.statcan.gc.ca/n1/pub/96-325-x/2014001/article/14027-eng.htm> (22 March 2023).
37. S. Alizon, A. Hurford, N. Mideo, M. Van Baalen, Virulence evolution and the trade-off hypothesis: History, current state of affairs and the future. *J. Evol. Biol.* **22**, 245–259 (2009).
38. B. R. Levin, J. J. Bull, Short-sighted evolution and the virulence of pathogenic microorganisms. *Trends Microbiol.* **2**, 76–81 (1994).
39. L. Rohmer, D. Hocquet, S. I. Miller, Are pathogenic bacteria just looking for food? Metabolism and microbial pathogenesis *Trends Microbiol.* **19**, 341–348 (2011).
40. N. Dekker *et al.*, Effect of spatial separation of pigs on spread of *Streptococcus suis* Serotype 9. *PLOS One* **8**, e61339 (2013).
41. N. Fittipaldi, M. Segura, D. Grenier, M. Gottschalk, Virulence factors involved in the pathogenesis of the infection caused by the swine pathogen and zoonotic agent *Streptococcus suis*. *Future Microbiol.* **7**, 259–279 (2012).
42. S. Hamaguchi, M. A. Zafar, M. Cammer, J. N. Weiser, Capsule prolongs survival of *Streptococcus pneumoniae* during starvation. *Infect. Immun.* **86**, e00802–17 (2018).
43. M. A. Zafar, S. Hamaguchi, T. Zangari, M. Cammer, J. N. Weiser, Capsule type and amount affect shedding and transmission of *Streptococcus pneumoniae*. *mBio* **8**, e00989–17 (2017).
44. M. Segura *et al.*, Update on *Streptococcus suis* research and prevention in the era of antimicrobial restriction: 4th international workshop on *S. suis*. *Pathogens* **9**, 374 (2020).
45. N. Willemsse *et al.*, An emerging zoonotic clone in the Netherlands provides clues to virulence and zoonotic potential of *Streptococcus suis*. *Sci. Rep.* **6**, 28984 (2016).
46. X. Fernández-Aguilar *et al.*, Urban wild boars and risk for zoonotic *Streptococcus suis* Spain. *Emerg. Infect. Dis.* **24**, 1083–1086 (2018).
47. L. A. Weinert *et al.*, Genomic signatures of human and animal disease in the zoonotic pathogen *Streptococcus suis*. *Nat. Commun.* **6**, 6740 (2015).
48. G. Zou *et al.*, Effects of environmental and management-associated factors on prevalence and diversity of *Streptococcus suis* in clinically healthy pig herds in China and the United Kingdom. *Appl. Environ. Microbiol.* **84**, e02590–17 (2018).
49. T. M. Wileman *et al.*, Pathotyping the zoonotic pathogen *Streptococcus suis*: Novel genetic markers to differentiate invasive disease-associated isolates from non-disease-associated isolates from England and Wales. *J. Clin. Microbiol.* **57**, e01712–18 (2019).
50. N. F. Hadjirin *et al.*, Large-scale genomic analysis of antimicrobial resistance in the zoonotic pathogen *Streptococcus suis*. *BMC Biol.* **19**, 191 (2021).
51. A. Bankevich *et al.*, SPAdes: A new genome assembly algorithm and its applications to single-cell sequencing. *J. Comput. Biol.* **19**, 455–477 (2012).
52. A. M. Bolger, M. Lohse, B. Usadel, Trimmomatic: A flexible trimmer for Illumina sequence data. *Bioinformatics (Oxford, England)* **30**, 2114–2120 (2014).
53. A. Gurevich, V. Saveliev, N. Vyahhi, G. Tesler, QUAST: Quality assessment tool for genome assemblies. *Bioinformatics* **29**, 1072–1075 (2013).
54. B. Langmead, S. L. Salzberg, Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–359 (2012).
55. S. Koren *et al.*, Canu: Scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res.* **27**, 722–736 (2017).
56. R. R. Wick, L. M. Judd, C. L. Gorrie, K. E. Holt, Unicycler: Resolving bacterial genome assemblies from short and long sequencing reads. *PLOS Comput. Biol.* **13**, e1005595 (2017).
57. R. R. Wick, M. B. Schultz, J. Zobel, K. E. Holt, Bandage: Interactive visualization of de novo genome assemblies. *Bioinformatics* **31**, 3350–3352 (2015).
58. G. G. R. Murray *et al.*, Mutation rate dynamics reflect ecological change in an emerging zoonotic pathogen. *PLOS Genet.* **17**, e1009864 (2021).
59. T. B. T. Athey *et al.*, Determining *Streptococcus suis* serotype from short-read whole-genome sequencing data. *BMC Microbiol.* **16**, 162 (2016).
60. T. Seemann, Prokka: Rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069 (2014).
61. G. Tonkin-Hill *et al.*, Producing polished prokaryotic pangenomes with the Panaroo pipeline. *Genome Biol.* **21**, 180 (2020).
62. E. Paradis, J. Claude, K. Strimmer, APE: Analyses of phylogenetics and evolution in R language. *Bioinformatics* **20**, 289–290 (2004).
63. Z. Zhou *et al.*, GrapeTree: Visualization of core genomic relationships among 100,000 bacterial pathogens. *Genome Res.* **28**, 1395–1404 (2018).
64. I. Letunic, P. Bork, Interactive Tree Of Life (iTOL) v5: An online tool for phylogenetic tree display and annotation. *Nucleic Acids Res.* **49**, W293–W296 (2021).
65. J. A. Lees *et al.*, Fast and flexible bacterial genomic epidemiology with PopPUNK. *Genome Res.* **29**, 304–316 (2019).
66. N. J. Croucher *et al.*, Rapid phylogenetic analysis of large samples of recombinant bacterial whole genome sequences using Gubbins. *Nucleic Acids Res.* **43**, e15 (2015).
67. G. G. R. Murray *et al.*, The effect of genetic structure on molecular dating and tests for temporal signal. *Methods Ecol. Evol.* **7**, 80–89 (2016).
68. A. J. Drummond, M. A. Suchard, D. Xie, A. Rambaut, Bayesian phylogenetics with BEAUti and the BEAST 1.7. *Mol. Biol. Evol.* **29**, 1969–1973 (2012).