

INTRODUCING

# Genpax

An annotated introduction  
including Publications  
and Poster Presentations, and  
selected additional materials

**A new era of connected pathogen genomics**



# Genpax 2024

## Introducing Genpax

- Genpax is a research and development-based company building novel solutions for clinical pathogen genomics. We are focused on **the needs of infection prevention and control (IPC)**, providing optimal information to recognize and respond to the transmission of strains in the healthcare system and beyond.
- Since 2021, we have established a large team of specialist bioinformaticians focused upon bacterial pathogen genomics, building upon more than a century of prior collective experience to develop species-specific toolkits with new analytical capabilities. These are being made available through our **IDEM** platform, addressing over 30, healthcare-associated, public health, and food-associated pathogens.
- Considering pathogens can change at a rate of 0 to 5 SNPs per genome per year, the actionable information needed to proactively detect (and exclude) outbreaks, infer transmission, and effectively direct IPC is beyond what other analysis pipelines can reliably deliver.
- Compromises, such as Sequence Typing (e.g. cgMLST) achieve scalability and error tolerance at the expense of sensitivity and specificity. In contrast, SNP-solutions (e.g. wgSNP) cannot be scaled and have reference-dependent accuracy (and good references do not exist for many strains and species).
- Genpax exists to **eliminate these constraints**; to deliver a new generation of pathogen analysis capabilities which address the IPC challenges of emergent pathogens and AMR.
- This brochure highlights a selection of **key differentiating capabilities** in our quest to make the best possible genomic pathogen analysis accessible to everyone.

## Genpax offers an automated cloud-based solution with the following features:

- **A SNP-throughout analysis** that makes maximal use of the safely interpretable genome sequencing information.
- An analysis solution that **does not require typing or other steps** to select a reference genome.
- **Performance equivalent to wgSNP under its most optimal conditions**, delivering equally optimal and comparable results for all species and strains.
- Entirely **consistent findings from clinical replicates**, identical strains from common sources, and samples in clinical ring-trials.
- A **near-zero error rate** meaning results from the multiple labs can be reliably combined and compared.
- **Unprecedented accuracy** combined with **addressing more of the genome sequence** information than either Sequence Typing or previous whole genome SNP comparisons.
- Unmatched capabilities to identify stains and **infer their likely membership of outbreaks (or not) and order of transmission**, even with only two isolates within a transmission cluster.
- Effective analysis that can be used as **part of a clinical solution** to help optimize infection prevention and control workflows for improved patient care and safety at the same time as reducing the costs of healthcare.
- **Open scalability**, so that each newly analyzed strain can be compared with all strains previously processed, within and between sites that choose to openly display their results.
- A **user-friendly platform** with interactive and continuously updated communication of findings.
- A **rapid turnaround time**, whether analyzing one or hundreds of isolates from a sequencing run, while comparing them with hundreds or tens of thousands (or more) previously analyzed strains.

## A novel multi-scale outbreak detection and strain identification capability for genome sequence-based infection control: an MRSA example

J.C. Littlefair, B.J. Uttley, D. Frampton, J.F. Peden, and Nigel J. Saunders

Materials from oral presentation at ASM Microbe 22

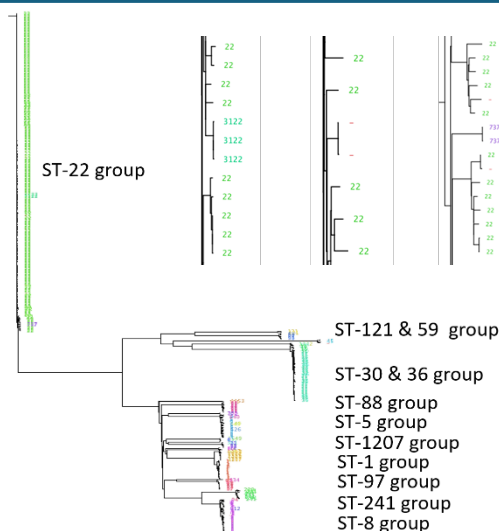
**Abstract:** The exclusion (as well as inclusion) of strains is vital for outbreak investigation and infection prevention and control of healthcare-associated infection, for both health resource management and patient-care. Current genome-sequence based diagnostics strategies have limitations in their analytical capabilities and scalability which restrict their utility for proactive diagnostic surveillance and to direct real-time infection control.

MLST is of limited value in the context of dominant epidemic or more virulent clones; such as MRSA where up to 90% of isolates fall into a small number of common STs. High-resolution SNP-based analysis can be pursued when there is an appropriate reference genome, but the greater the diversity between the sample and reference genome the lower the coverage, resolution, and accuracy. Further, the similarity of members of clonal clusters within hospital and surrounding community and health-care environments can be high. Thus, the detection and differentiation of hospital-associated acquisition and transmission from strains entering the hospital is non-trivial and requires improved sequence analysis to resolve.

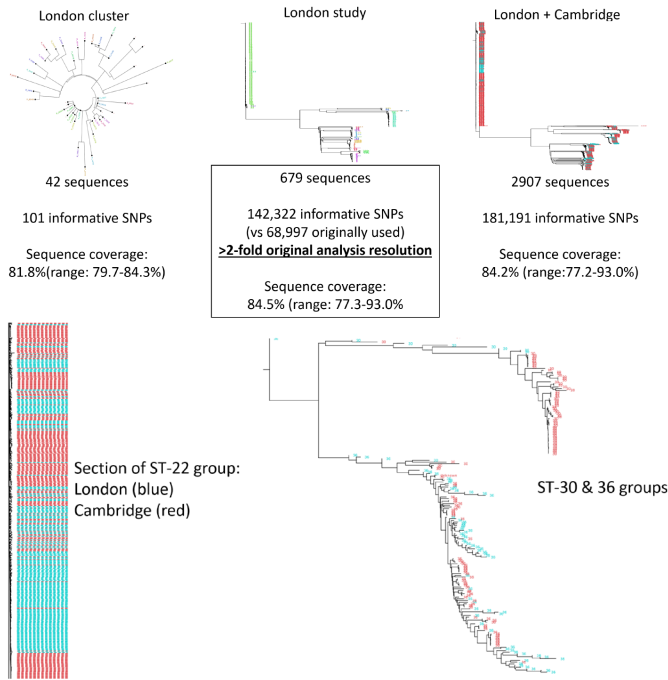
GenPax has developed a novel analysis pipeline that can perform multi-scale analysis for both outbreak detection and high-resolution determination of strain identity and relatedness (and implied transmission connections). This can be run without prior knowledge of potential outbreaks or associated strains and can detect single instances of known high-virulence and highly-resistant clones.

Using a well-established dataset of over 600 London-based MRSA we achieved clearly higher resolution and improved determination of connected isolates than using traditional methods: the number of sites used for identification within the study was more than doubled, and the dynamic range of differences and resolution of associated strains was greatly increased. While highly similar strains some with less than 10 nt difference were still identified, other previously linked strains became clearly separated, indicating that the number and size of hospital and health-care associated cross-infections was substantially lower than previously thought. There was also a better match between the size of the largest cluster and the number of epidemiologically determined interactions.

Thus, we have shown that this next generation of bacterial genomics resources can substantially increase the future diagnostic utility of genome sequencing for hospital infection control.



The first implementation of a novel multi-scale solution can analyze, connect, and compare diverse strains with a natural reference-free solution. A solution that works across a whole population with preserved resolution, while clearly linking isolates within outbreaks with divergent or no Sequence Type designations.

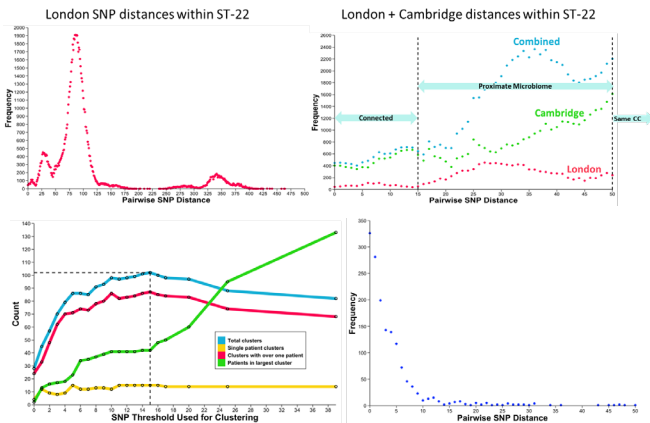


A unique multiscale detailed analysis that does not lose resolution when addressing larger numbers of strains. More than doubling the strain defining information at around 700 strains, and maintaining sequence addressed; illustrated by extending analysis to nearly 3000 strains.

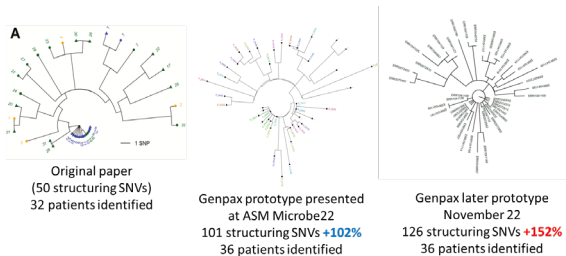
Large scale analysis led to the discovery of local healthcare microbiomes that mean that thresholds for the detection and discrimination must be re-defined in this context. Because outbreaks are occurring in a background of highly similar unlinked locally circulating strains.

Thresholds for detection and separation of outbreaks need to be lower than previously reported (with this analysis, around 15 SNVs), meaning that discrimination of outbreaks is compromised by the noise/error rates using other methods. The London study reported more outbreaks in their publication than had actually occurred.

Reanalysis of the largest outbreak identified four previously missed cases, and more than doubled the resolution for the identification of transmission chains (middle). A later 2022 prototype that was not ready for presentation but was referred to with further resolution (right) is also shown here.



**Conclusion:**  
There is an optimal window for connected strain detection between 10 and 15 SNPs, which for inclusivity we choose to use 15 SNP threshold for this population



- A minimally reducing SNP-level multiscale 'direct to identify' analysis is possible, addressing more genome content than cg/wgMLST or wgSNP
- In this analysis of 679 strains with an *S. aureus* prototype, a >2-fold increase in informative SNPs was achieved, coupled with noise reduction
- Broader comparisons identified local healthcare microbiomes, and these and the diversity of same-patient isolates informed the selection of thresholds for identifying probably connected strains (15 SNPs)
- Reanalysis resulted in BOTH recognizing that some previously suggested clusters were probably not connected strains AND that strains had been wrongly omitted from some true clusters
- The greater information accessible can be used to structure clusters to assist and improve outbreak investigation by infection control
- The 'each against all' capability of this analysis at scale enables this strategy to be applied to proactive surveillance and outbreak detection, in addition to confirmation and investigation of suspected clusters

As a commercial entity, Genpax cannot share its code and solutions. However, we are a team, comprised largely of former academics with hundreds of publications in the field between us, who want to share our system and its capabilities as openly and clearly as possible. To do this, we have performed a set of validation and demonstration analyses using information from the highest quality publications and studies that we could identify, selecting those with the best evidence for ‘ground truth’ against which to be measured.

- In tests of **reproducibility** and near-zero error, genuine biological replicates from clinical ring trials (ECCMID, *Staph. aureus*) and large well-documented outbreaks and re-isolation studies have been addressed (*E. coli* and *K. pneumoniae*).
- In tests of **accuracy** (*E. coli*, *Campylobacter jejuni*, and *Ps. aeruginosa*), exceptional situations in which published or specifically generated almost identical high-quality reference genomes were used in the original studies were selected and thus represent the most stringent findings to measure performance against that we could identify.
- In tests of **reference-free** performance and transmission-chain re-structuring, species were selected that represent extremes of highly recombining panmictic (*Campylobacter jejuni*) and deeply rooted, highly-clonally diverse (*Ps. aeruginosa*) population structures.
- In the test of **scalable** comparisons, using *Listeria monocytogenes*, we processed data generated in large studies from different European laboratories.
- The test of MRSA **gene finding** (ASM, *Staph. aureus*) used sequencing and MRSA/MSSA data from two published studies from an EU reference laboratory.
- Likewise, our **economic modeling** adopts conservative assumptions, taking a cautious approach towards outbreak sizes and containment speed, in contrast to the assumptions of the published models it is built upon. Additionally, our analysis includes up-to-date costs of sequencing and analysis ensuring accurate financial impacts.

Each poster primarily addresses one or two aspects of our platform’s performance for IPC applications: accuracy, low noise, high resolution, comparability, species applicability, and scalability. In combination they represent a real enhancement in what sequencing analysis can offer infection prevention and control.

If you have questions, please get in touch via [research@genpax.co](mailto:research@genpax.co).

# Calling Zero: A new foundation for diagnostic bacterial genomics

Presented at ECCMID 2023

James C. Littlefair, Benedict J. Uttley, Dan G. Frampton, Gareth M. Linsmith, John F. Peden, & Nigel J. Saunders: Genpax, London, United Kingdom

[jlittlefair@genpax.co](mailto:jlittlefair@genpax.co) +44 7473 871 175

## Introduction

- Bacterial strains typically diversify at rates below 10 SNPs per year and thresholds to recognize source-linked and transmission-associated strains typically range from 10 to 20 SNPs. Therefore, even low noise levels impact outbreak cluster detection and analysis for source attribution and transmission inference.
- Reproducibility, low noise, and the ability to call true zeros from resequencing of the same DNA, culture, patient, or closely-linked isolates, is essential for cluster recognition and accurate branching structures of within-cluster dendrograms.
- The most effective outbreak surveillance requires reproducibility within and between laboratories as it facilitates multi-site surveillance and comparability.
- A cloud-based solution which works directly from FASTQ, delivers SNP-resolution information, and requires no clonal reference can facilitate this enhanced infection control and prevention.

## Results

- The Genpax analysis pipeline consistently obtained 0 SNP-distances within all 20 replicate groups across all five participating laboratories, producing 17 clusters from 110 samples.
- The replicate groups that clustered together at 0 SNP-distance, in concordance with the original study, were:
  - NGSRT01& NGSRT02
  - NGSRT03 & NGSRT05
  - NGSRT14 & NGSRT15
- However, replicate groups NGSRT18 & NGSRT19, identical in the original cgMLST analysis, were separated by a single intergenic SNP not addressed by the sequences used in cgMLST.
- Replicates within 0 SNP clusters typically shared >80% of their genome length from which variants could be called (based on a 2.8Mb genome), over 500kb more sequence than cgMLST.
- This compares to reported whole genome SNP resolutions of 72% within Sequence Type and 57% across the species [2].

## Results

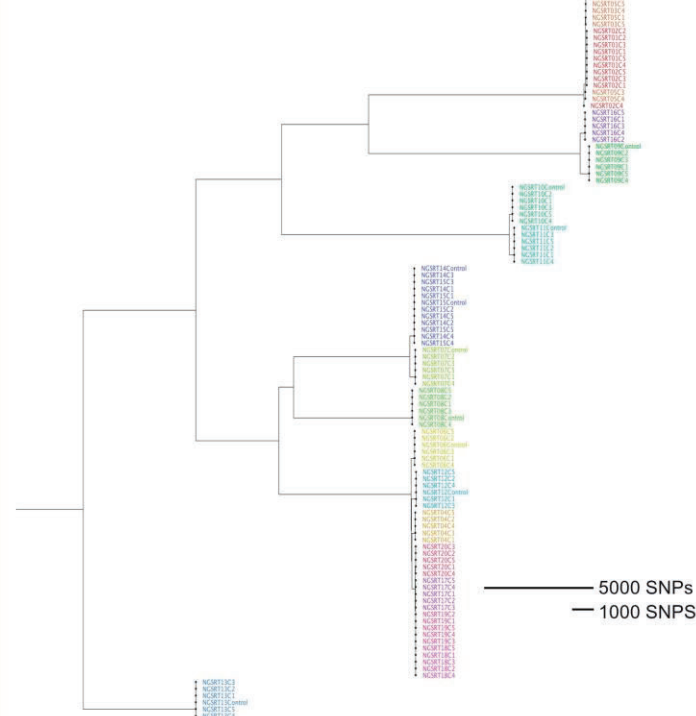


FIGURE 1 – NJ dendrogram showing all 110 Mellman ring trial *S. aureus* replicates as analyzed with the Genpax pipeline. Colours represent 0 SNP clusters.

## Methods

- 110 readsets from the Mellmann ring trial [1] published in 2017 were processed through the Genpax analysis pipeline in the version as of Q1 2023.
- The ring trial included *S. aureus* sequences from diverse sources sequenced on Illumina MiSeq platforms by five different laboratories across three European countries (Denmark, Germany and the Netherlands) following the same protocol (Nextera XT Library Prep and 250-bp paired end).

TABLE 1 Characteristics of the 20 human *S. aureus* isolates that were sent as DNA samples to the five participating laboratories in a blinded fashion and used as controls

Sample ID	Original strain	Spa type (based on Sanger sequencing)	Comment/reference	Genpax 0 SNP cluster
NGSRT01	469	1011	Livestock-associated MRSA	1
NGSRT02	551	1011	Livestock-associated MRSA, identical cgMLST genotype as NGSRT01	1
NGSRT03	1346	1011	Livestock-associated MRSA	2
NGSRT04	1364	1010	Classical hospital-acquired MRSA	3
NGSRT05	1360	1011	Livestock-associated MRSA, identical cgMLST genotype as NGSRT03	2
NGSRT06*	2180	1002	Central European community-acquired PVL <sup>+</sup> -positive MRSA	4
NGSRT07*	2482	1008	US typical community-acquired PVL <sup>+</sup> -positive MRSA	5
NGSRT08*	2560	1044	Central European community-acquired PVL <sup>+</sup> -positive MRSA	6
NGSRT09*	2638	1012	Classical hospital-acquired MRSA	7
NGSRT10*	2786	1843	mecC-positive MRSA	8
NGSRT11*	2949	1843	mecC-positive MRSA	9
NGSRT12*	2994	1003	Classical hospital-acquired MRSA	10
NGSRT13*	3039	1032	Classical hospital-acquired MRSA	11
NGSRT14*	COL	1008	MRSA strain COL	12
NGSRT15*	COL	1008	Duplicate of MRSA reference strain COL	12
NGSRT16	ATCC 25923	1021	MSSA quality control strain ATCC 25923	13
NGSRT17	P1	1001	Isolate P1 from reference 23	14
NGSRT18	P3	1001	Isolate P3 from reference 23	15
NGSRT19	P4	1001	Isolate P4 from reference 23, identical cgMLST genotype as NGSRT18	16
NGSRT20	P12	1001	Isolate P12 from reference 23	17

\*These samples were separately cultivated, and DNA was extracted and sequenced as controls.  
\*PVL<sup>+</sup>, Pantom-Valentine leukocidin.  
Genpax 0 SNP clusters in bold font comprised more than one ring trial replicate group.

## Results

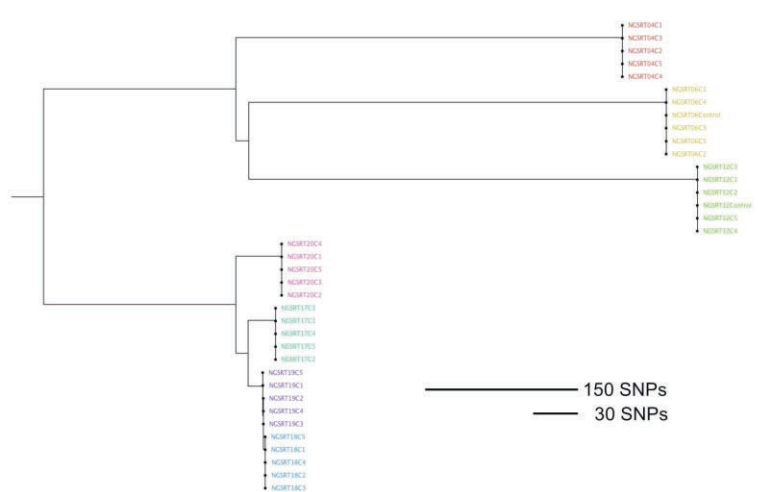


FIGURE 2 – NJ dendrogram showing 37 Mellman ring trial *S. aureus* replicates within the same clonal complex (CC5). Colours represent 0 SNP clusters.

## Conclusions

- Because of limitations in comparability and scalability, it has previously been necessary to use typing methods such as cgMLST, which do not make full use of the available WGS information when performing large-scale and multi-site public health surveillance.
- The Genpax analysis pipeline, despite the multi-site nature of the study with variability in sequence coverage and quality, demonstrates increased resolution whilst simultaneously reducing noise, giving accurate and comprehensive SNP-resolution information.
- By using a strategy that does not depend upon a clonal reference, which is necessary in the absence of prior knowledge and typing, the pipeline is applicable to all strains within the species diversity, including MSSA and MRSA, which is essential for optimal infection prevention and control.

## References

- Mellmann A, Andersen PS, Bletz S, Friedrich AW, Kohl TA, Lijie B, et al. High Interlaboratory Reproducibility and Accuracy of Next-Generation-Sequencing-Based Bacterial Genotyping in a Ring Trial. *J Clin Microbiol*. 2017;55(3):908-13
- Gorrie CL, Da Silva AG, Ingle DJ, Higgs C, Seemann T, Stinear TP, et al. Key parameters for genomics-based real-time detection and tracking of multidrug-resistant bacteria: a systematic analysis. *Lancet Microbe*. 2021;2(11):e575-e83

## Declaration

This research was entirely funded by Genpax.  
Genpax is a bioinformatics company founded in 2020 seeking to develop novel solutions that overcome the limitations of established analysis strategies to maximize the usefulness of bacterial genome sequences in infection control and prevention.

Check out our website:





# A novel genome comparison tool producing near-zero error for same-patient isolates of *E. coli* ST131

Presented at ASM Microbe 23

Poster No. 265  
Date: 6/17/2023

Rebecca J. Bengtsson, John F. Peden, Dan Frampton, Gareth Linsmith, Benedict J. Uttley, Arthur Poivet, Luis Montemayor & Nigel J. Saunders  
Genpax, London, United Kingdom

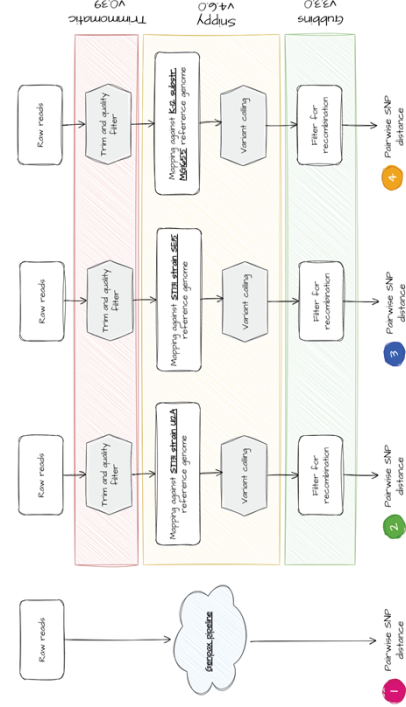
research@genpax.co  
+44 203 603 6889

## Introduction

- Escherichia coli ST131 is a globally disseminated clone and a significant contributor to the global burden of urinary tract infections (UTIs).
- Whole genome sequencing (WGS) based analysis provides higher resolution than traditional typing methods and its use in public health settings can significantly help improve surveillance and outbreak investigation of pathogens.
- Currently, WGS is primarily applied retrospectively and is limited by its speed and scalability.
- To achieve its potential impact, it needs to be used proactively to direct patient care and healthcare resource management.
- WGS analysis also requires expert knowledge, including selecting appropriate reference genomes and tools for analysis of the sequence data as well as interpreting the resulting data to clinically actionable results.
- To tackle this limitation, Genpax has developed a novel species-centric, automated genome comparison tool that does not require the selection of an appropriate reference genome to perform SNP-resolution analysis.
- In this study, the performance and accuracy of this method is compared to traditional core genome SNP analysis with UTI *E. coli* isolates from a published study [1].

## Methods

- The published study consisted of 65 *E. coli* ST131 faecal and urinary isolates sampled from a single patient with a long-term UTI were analyzed.
- This study was selected because it contained similar biological replicates, and had generated a reference genome (U12A, using PacBio) from one of the clonal isolates.
- The use of a nearly identical reference provides a best case for analysis that minimizes the introduction of systematic errors during mapping and maximizes the accuracy and resolution of pairwise SNP (pairwise-SNP) distance determination [2].
- The dataset was processed using two methods.
- The first is the Genpax analysis pipeline which does not require the selection of a reference genome and is referred to as analysis 1.
- For the second method, we used the current versions of the industry standard tools used for bacterial genomic analysis. This included Trimmomatic for short reads adapter trimming and quality filter, Snippy for reference mapping and variant calling, and Gubbins for recombination filtering. This is similar to the strategy used in the original publication [1].
- Three analyses using the second method were carried out with different reference genomes for mapping. Analysis 2 used the intra-clonal U12A study-generated reference genome. Analysis 3 used a standard reference genome for this clonal complex: *E. coli* ST131 strain SE15. Analysis 4 used standard *E. coli* strain K-12 substr. MG1655.
- To assess the accuracy of the developed pipeline, the pairwise-SNP distances from the four analyses were compared to the results from the original publication.



## References

- Forde, S.M., Roberts, L.W., Phan, M.D., Peters, K.M., Fleming, B.A., Russett, C.W., Leuther, S.M., Myers, J.B., Barker, A.P., Fisher, M.A. and Chong, T.M. 2019. Population dynamics of an *Escherichia coli* ST131 lineage during recurrent urinary tract infection. *Nature communications*, 10(7), 1-9.
- Coxe, C.L., De Silva, A.G., Iggle, D.J., Hogg, C., Steinhilber, T.J., Shearer, T.P., Williamson, D.A., Kwong, J.C., Grayson, M.L., Sherry, N.L. and Hooper, B.P. 2021. Comments on genomic-based outbreak detection and tracing of *Escherichia coli* ST131. *Microbes*, 16(11), 1949-1952.

Check out our website and other supporting information



## Declaration

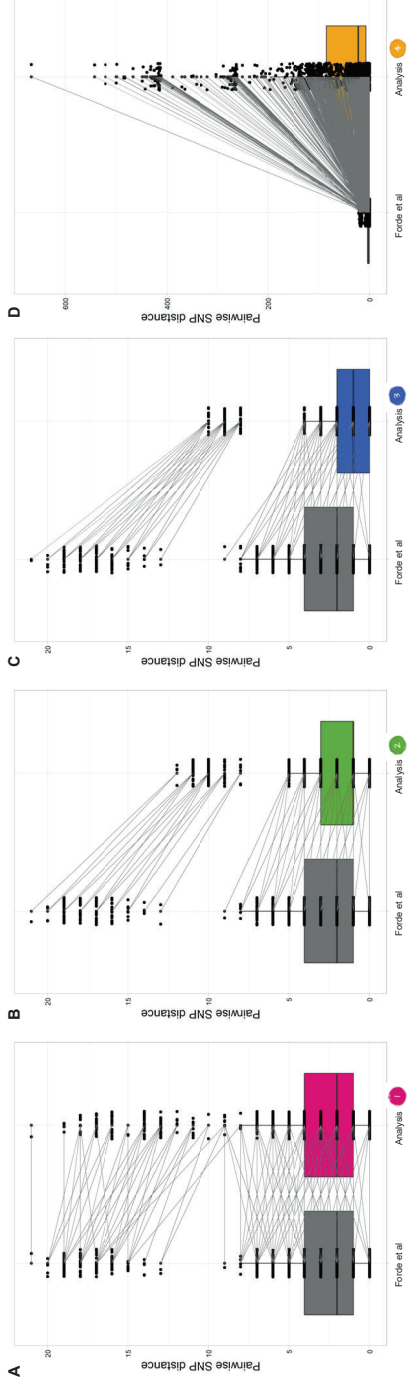
This research was entirely funded by Genpax. Genpax is a bioinformatics company developing novel solutions that overcome the limitations of established methods for genomic analysis of bacterial genomes. Genpax is a bioinformatics company developing novel solutions that overcome the limitations of established methods for genomic analysis of bacterial genomes. Genpax is a bioinformatics company developing novel solutions that overcome the limitations of established methods for genomic analysis of bacterial genomes.

## Results

Table 1. Number of samples processed and variants genotyped

	Genpax 1	U12A 2	SE15 3	K-12 4
Samples processed	65	63	63	61
Variants genotyped	59	24	22	1,640

All 65 samples were analyzed by the Genpax pipeline. However, using the Snippy/Gubbins pipeline, 2 samples were rejected using the U12A and SE15 reference genome, and 4 samples were rejected using the K-12 reference genome. Therefore, we present the 61 samples that are common to all analyses. Table 1 describes the total number of variants genotyped across the 61 samples analyzed.



Comparison of pairwise SNP distance distribution of 61 *E. coli* ST131 isolates from the four analyses under investigation. Each box plot demonstrates the distribution of pairwise-SNP distances between all isolate pairs compared to pairwise-SNP distances from the original publication. (A) Comparison of pairwise-SNP distances between isolate pairs from Forde et al relative to isolates pairs from analysis 1, (B) isolate pairs from analysis 2, (C) isolate pairs from analysis 3 and (D) isolate pairs from analysis 4.

- There are two main populations from the original publication, isolate pairs with distance between 0 and 9, and pairs with distance between 13 and 21.
- The boxplots revealed that analysis 1 processed using the Genpax pipeline produced pairwise-SNP distances that shared the most similarities to distances from the published results, for both populations.
- The current version of Snippy and Gubbins with both the intra-clonal U12A (analysis 2) and intra-ST SE15 (analysis 3) reference genomes detected lower pairwise-SNP distances relative to the original publication, with analysis 3 detecting even lower comparative pairwise-SNP distances than analysis 2.
- Analysis 4 using the general *E. coli*/K-12 reference genome resulted in a dramatic decrease in SNP detection precision, by which all isolate pairs displayed a significant increase in pairwise-SNP distance compared to the published results, reflecting the widely recognized impact of selecting distantly related reference genomes for mapping-based SNP analyses.

## Conclusions

- The Genpax reference-independent method produced isolate pairs with similar pairwise-SNP distances to those reported using a gold standard method with an optimal intra-clonal reference genome.
- Using the standard reference-based method, the intra-clonal strain reference U12A improves SNP detection sensitivity compared to using a reference of the same ST; but neither matched the performance of the Genpax pipeline.
- With a standard reference-based method, use of a species-specific reference genome outside of the clonal complex resulted in a dramatic drop in SNP precision and over-estimation of SNPs identified.
- The Genpax reference-independent method demonstrated substantially better performance when compared to using an industry standard methodology.

# Near zero error using large-scale hospital outbreak whole genome sequence data for *Klebsiella pneumoniae*

Presented at ASMI Microbe 23 **Ramiro Morales-Hojas, James C. Littlefair, Daniel Frampton, John F. Peden & Nigel J. Saunders: Genpax, London, United Kingdom**

research@genpax.co.uk  
+44 203 603 6669

## Introduction

- Klebsiella pneumoniae* is a common nosocomial pathogen responsible for a range of severe and life-threatening infections [1].
- K. pneumoniae* is a leading cause of extended-spectrum beta-lactamase (ESBL)-producing, and carbapenem-resistant healthcare-associated infections and is considered a critical public health threat by the WHO [2].
- Although the overall prevalence is lower than *E. coli*, *K. pneumoniae* is notably linked with higher rates of hospital transmission [3].
- Several genomic studies indicate that up to 1/3 of hospital infections are linked to within-hospital transmission events [4, 5], and proactive genomic surveillance is becoming the gold standard infection control practice to detect nosocomial infections and to prevent the transmission of highly resistant strains.
- Prospective surveillance requires accurate and high-resolution genomic analysis solutions to detect transmission events when bacterial species typically evolve at a rate of between 1 and 10 SNPs per year, which is challenging to accomplish using existing methods.

**Objective:** To test the resolving power of Genpax's WGS analysis, through the reanalysis of clinical outbreaks of *K. pneumoniae* using data from published studies with established epidemiological contexts.

## Methods

- The performance of the Genpax pipeline was evaluated using a recently published dataset obtained from two intensive care units (ICUs) in Vietnam [6]. This study was chosen due to the number of isolates associated with transmission events and low genomic distances, providing an ideal opportunity to assess the pipeline's performance.
- The novel pipeline was used to call genome-wide SNPs for 1314 *K. pneumoniae* isolates from the study.
- Cluster analysis was conducted based on informative SNPs. This analysis assigned samples to zero-SNP transmission clusters if they exhibited a distance of 0 SNPs to any other sample within the cluster.

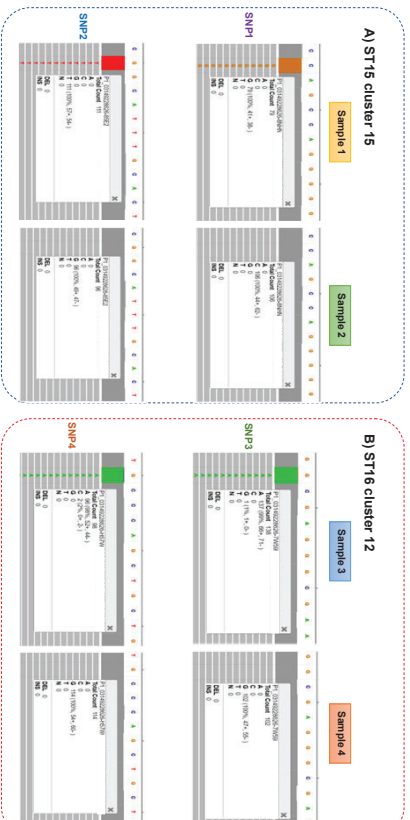
## References

- Parishos, N., Gounis, S. P. & Gounis, B. F. (2015). Characterisation of the *Escherichia* Pathogen. *Expert Rev Anti Infect Ther*, 11, 297-308.
- World Health Organization. (2017). Global Priority List of Antibiotic-Resistant Bacteria to Guide Research, Discovery, and Development of New Drugs. <https://www.who.int/publications/m/item/global-priority-list-of-antibiotic-resistant-bacteria-to-guide-research-discovery-and-development-of-new-drugs>
- Hilly et al. (2012). Transmission dynamics of extended-spectrum beta-lactamase-producing *Escherichia coli* in the tertiary care hospital and community. *PLoS One*, 7(12), e50729.
- Morales-Hojas, R. M., Littlefair, J. C., Frampton, D., Peden, J. F., & Saunders, N. J. (2023). Proactive genomic surveillance of nosocomial infections in the tertiary care hospital. *PLoS One*, 18(3), e0242082.
- Commun, L. S. (2015). *Escherichia coli* O157:H7 outbreak associated with consumption of unpasteurized apple juice. *Emerging Infectious Diseases*, 21(12), e2015-1234.
- Roberts, T. W. et al. (2022). Genomic characterization of multidrug-resistant *Escherichia coli* ST15 *Klebsiella pneumoniae* and 4 related clusters. *Emerging Infectious Diseases*, 28(12), e2022-1234.

## Results

- The isolates had an average genome length of 5.2 Mbp, of which 89% (4.6 Mbp) was suitable for SNP calling using the Genpax pipeline. This represents a 24% increase in genome length available for variation analysis compared to the 3.76 Mbp core genome described in the original publication.
- The reanalysis revealed a higher number of 0-SNP clusters for all *K. pneumoniae* samples compared to the original publication. This difference was primarily due to the increased number of clusters identified in the reanalysis of the ST15 isolates (Table 1).
- ST15 Cluster 15: The cluster revealed a revised 0-SNP cluster with 78 members. Among these, 67 (85%) were shared with the previously reported cluster. The reanalysis identified an additional 11 isolates, while 12 originally included isolates had 1 to 4 verified SNPs (Figure 1A). Notably, four of the samples originally reported to be part of the ST15 0-SNP cluster were assigned to a different 0-SNP group from the original publication (Table 2).
- ST16 Cluster 12: 13 out of the 14 ST16 isolates reported to be in a 0-SNP cluster were included in a redefined 0-SNP cluster. One isolate originally assigned to cluster 12 differed by 2 verified SNPs (Figure 1B). Six additional isolates were placed in this redefined cluster (Table 2), making a final group of 19 rather than 14.
- Cluster reanalysis of ST15 isolates in 0-SNP clusters utilizing a 5-SNP threshold was performed, as in the original publication. The largest cluster obtained in the reanalysis comprised 130 isolates of the 138 in the largest 5-SNP cluster in the publication (Table 3).
- Cluster reanalysis of ST16 isolates in 0-SNP clusters using a 5-SNP threshold identified 3 clusters instead of the single 5-SNP cluster reported previously (Table 3). Isolates excluded from the redefined 5-SNP clusters differed by 6 to 12 identified (and verified) SNPs.

**Figure 1:** IGV screenshots of SNP positions identified in the present reanalysis, highlighting the differentiation among isolates originally assigned to the 0-SNP clusters 12 and 15. **A)** Reference and alternative SNP alleles in two positions separating two samples of the ST15 0-SNP cluster (cluster 15). **B)** Reference and alternative SNP alleles in two positions identified in the reanalysis, differentiating ERR3585327 from the other samples in the ST16 0-SNP cluster 12.



**Table 1:** Number of 0-SNP clusters reported by Roberts et al. for each of the main sequence types (STs) and those obtained in this reanalysis with the Genpax pipeline.

	Roberts et al.	Genpax
<i>K. pneumoniae</i>	71	84
ST15	21	28
ST16	17	16
ST11	11	14
ST656	13	11

**Table 2:** The number of isolates in two 0-SNP clusters reported by Roberts et al. for ST15 and ST16 and the corresponding clusters obtained with the Genpax pipeline.

ST	Cluster	N (6)	N Genpax	Common	Unique (6)	Unique Genpax
15	15	79	78	67	12	11*
16	12	14	19	13	1	6**

\* 1 of these samples were reported in a different 0-SNP cluster by Roberts et al. (6); 7 were not in any cluster

\*\* 1 sample reported on a different 0-SNP cluster by Roberts et al. (6); 5 were not in any cluster

**Table 3:** The number of 5-SNP clusters reported by Roberts et al. [6] for ST15 and ST16 and redefined with the Genpax pipeline; and number of shared isolates between the largest original and redefined 5-SNP clusters.

ST	N (6)	N Genpax	Shared isolates in largest cluster
15	5	8	130/138
16	1	3	114/117

## Conclusions

- This reanalysis found that the Genpax pipeline accurately identified and differentiated transmission-linked samples. With lower noise, higher accuracy, higher resolution, and an improved determination of strain relationships compared to the standard tools used in the original study.
- The high coverage and resolution, reflected by the substantial (24%) increase in addressed genome space, was achieved while maintaining high accuracy and low noise (near-zero error).
- The resolution exhibited by this novel pipeline enables precise evaluation of transmission networks, offering more accurate insights into hospital outbreaks of *K. pneumoniae*.

Check out our website and other supporting information



**Declaration**  
This research was entirely funded by Genpax. Genpax is a bioinformatics company developing and providing genomic data analysis services. We have established analysis strategies to maximize the usefulness of bacterial genome sequences in infection prevention and control.

# Reference-free whole genome SNP analysis of *Pseudomonas aeruginosa*, with the restructuring of outbreaks analyzed with established methods

Presented at  
ASM Microbe 23

research@genpax.co  
+44 203 603 6869

Arthur Poivet, Rebecca J. Bengtsson, Dan Frampton, John F. Peden & Nigel J. Saunders: Genpax, London, United Kingdom

Poster No. 227  
Date: 6/18/2023

## Introduction

- Pseudomonas aeruginosa* is one of the main nosocomial pathogens with a high prevalence in burns units, Intensive Care Units (ICUs), and patients with cystic fibrosis [1].
- Ps. aeruginosa* has been classified as one of three critical priority pathogens, and is considered a major threat by the WHO and CDC due to the emergence of multidrug- and extended-drug-resistant clinical isolates [2,3].
- Ps. aeruginosa* is an ancient species with a diverse clonal population structure, for which reference genome solutions are unsuitable.
- The lack of suitable reference genomes typically restricts analysis to Sequence Typing (MLST, cgMLST, wgMLST).
- A general solution that does not depend upon local, high-quality references is needed to deliver clinical genomics and proactive sequencing for infection prevention and control for this AMR priority species.

**Objective:** To test a reference-free whole genome SNP analysis for *Ps. aeruginosa*, and to compare its performance against published studies.

## Methods

- The performance of the Genpax analysis pipeline was evaluated using datasets from two published studies, and compared to the original findings:
  - > 156 isolates from Magalhães et al [4].
  - Clinical and environmental isolates were originally typed and separated in 3 groups mainly corresponding to ST-1076, ST253, and ST17.
  - For each ST group, a complete reference genome was created by sequencing a clinical isolate with both PacBio and Illumina HiSeq technologies, and whole genome SNP distances were obtained using these references.
- > 38 isolates from Cunningham et al., which were originally analyzed using two cgMLST methods. One of those was an in-house method which addressed 4,041 alleles based on the PAO1 reference genome [5].

## References

- Shaw C, et al. "Whole-genome Sequencing of Pseudomonas aeruginosa PAO1, an Opportunistic Pathogen." *Nature* 406, no. 6799 (August 2000): 908-14. <https://doi.org/10.1038/35020200>
- Reyes J, et al. "Global Epidemiology and Clinical Outcome of Carbapenem-Resistant Pseudomonas aeruginosa and Associated Risk Factors in a Tertiary Care Hospital." *The Lancet Infectious Diseases* 4, no. 3 (March 1, 2004): e10-19. [https://doi.org/10.1016/S1473-3099\(03\)00262-9](https://doi.org/10.1016/S1473-3099(03)00262-9)
- Kerrin A, et al. "Evidence-Based Treatment of Pseudomonas aeruginosa Infections: A Critical Appraisal." *Antibiotics* 12, no. 2 (February 2023): 262-75. <https://doi.org/10.3390/ant12020262>
- Magalhães B, et al. "Combining Strand-Seq, PacBio and HiSeq Genome Sequencing to Investigate Pseudomonas aeruginosa Epidemiology in Intensive Care Units." *Frontiers in Public Health* 8 (January 2020): 1-12. <https://doi.org/10.3389/fpubh.2020.00013>
- Cunningham M, et al. "A Reference-Free Approach to Investigate Pseudomonas aeruginosa Outbreaks." *Microbiology Spectrum* 10, no. 6 (November 6, 2022): e02920-22. <https://doi.org/10.1128/microsp.02920-22>

## Results – SNP analysis

- The mean sequence length analyzed using the Genpax pipeline represents around 85% of the average genome length.
- The Genpax analysis replicated the results obtained from whole genome SNP analysis (which used optimal bespoke reference genomes), both for clonal outbreaks (figure 1) and more diverse isolates (figure 2).
- Previously identified clusters of closely related strains were identified by the Genpax pipeline, and the performance was assessed by comparing within-cluster SNP distances (Table 1).

- In the ST-1076 outbreak (Figure 1), all samples fall within 14 SNPs of each other, and local transmission chains can be inferred (examples in colored boxes).
- In some cases, the environmental source of the infection could be identified (e.g. Figure 2 – blue cluster).

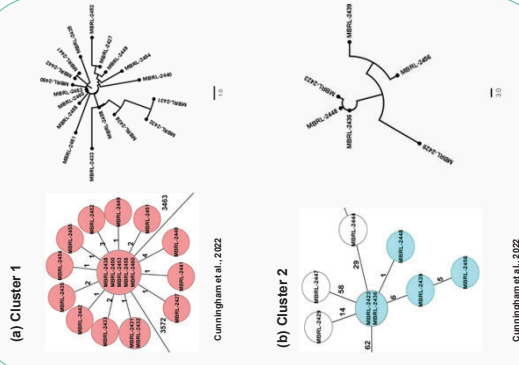
ST	target	Magalhães et al. (SNPs)	Genpax (SNPs)
1076	All iso.	≤ 14	≤ 14
1076	All same-patient iso.	≤ 10	≤ 11
1076	Patient 4 iso.	≤ 7	≤ 5
1076	Patient 24 iso.	≤ 2	≤ 4
253	All same-patient iso.	≤ 4	≤ 3
253	Burns unit	≤ 11	≤ 11
253	ICU 5 cluster	11 ≤ 14	10 ≤ 14
253	ICU 2 cluster (suspected outbreak)	≤ 1	≤ 0
17	All same-patient iso.	≤ 6	≤ 12
17	Suspected outbreak (figure 2 - green)	≤ 13	≤ 12
17	Patient 11 cluster (figure 2 - blue)	≤ 7	≤ 2

**Table 1:** Comparing the performance of the Genpax pipeline by comparing whole genome SNP distances of closely related isolates.

## Results – SNP vs cgMLST

- This new SNP-based approach successfully identified outbreak clusters from clustered strains of diverse clonal complexes.
- The analysis demonstrates superior resolution and more informative sub-structuring of the outbreak than the minimum spanning trees derived from either of the two cgMLST schema under comparison.

- In Figure 3a the SNP-level analysis identified differences in 3 strains previously found to be identical by cgMLST.
- In Figure 3b the analysis of cluster 2 confirmed that 3 isolates near the root of the outbreak were within 3 SNPs and yielded different transmission inferences:
  - MBRL-2439 and MBRL-2456 are no longer sequential in the transmission chain.
  - MBRL-2429 is not tangential to the outbreak structure.



**Figure 3:** Comparison between cgMLST minimum spanning trees (left) and Genpax dendrograms (right). The distances on the minimum spanning tree refer to allelic differences [5].

## Conclusions

- The novel Genpax methodology accurately determined strain identity without the need for a closely related reference genome, prior knowledge of strain types, or clonal clusters.
- The achieved resolution in this study matched that of an expensive and unscalable customized high-quality approach and surpassed the resolution achieved with cgMLST approaches.
- The utilization of this innovative analysis tool can enable real-time phylogenetic analysis in clinical settings.

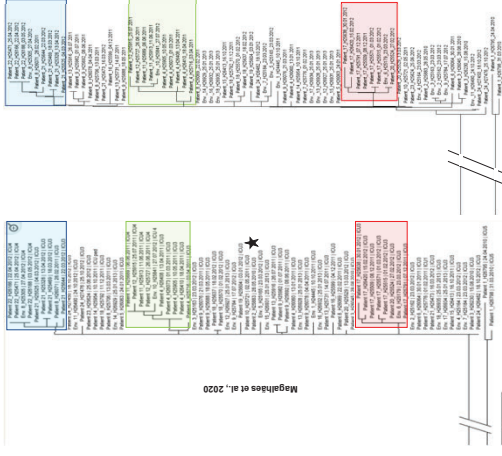
Check out our website and other supporting information



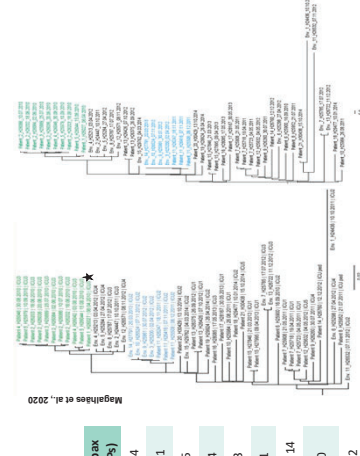
## Declaration

This research was entirely funded by Genpax. Genpax is a bioinformatics company developing novel solutions that overcome the limitations of established bacterial genome sequences in infection prevention and control.

**Figure 1:** Comparison of whole genome SNP analysis of all ST1076 isolates, between a tree from Magalhães et al. (left), and the reference-free Genpax solution (right) (\* = reference).



**Figure 2:** Comparison of whole genome SNP analysis of all ST17 isolates, between a tree from Magalhães et al. (left), and the reference-free Genpax solution (right).



**Table 1:** Comparing the performance of the Genpax pipeline by comparing whole genome SNP distances of closely related isolates.

# Reference-free WGS SNP-resolution analysis of *Campylobacter jejuni*

## Introduction

- Analyzing highly recombining species with panmictic populations poses challenges for whole genome sequence analysis due to the absence of reference genomes for accurate strain analysis and generation of universally applicable and comparable results.
- This means that analysis of these species is often limited to Sequence Typing methods that have sub-optimal resolution and provide poor source-attribution and inferred transmission information.
- Campylobacter jejuni* exemplifies this issue, as outbreaks can be spread through food distribution networks, resulting in transmission-linked isolates being potentially distributed across different laboratories.
- The objective was to evaluate the performance of Genpax's WGS analysis pipeline by applying it to sequence data from Pascoe's study [1] and compare the results.
- This study was selected because of its clearly defined transfer history of strains, low SNP differences, but counter-intuitive findings with respect to history and rates of evolutionary change - even when using a close to optimal WGS analysis strategy with a well established highly similar ancestral reference genome.
- It also presents a well controlled and documented simulation of the processes that occur in the context of a clinical outbreak, in which mutations are tracked to determine strain identity, and to infer transmission. It also represents a stringent test against which to compare a reference-independent analysis.

## Methods

- Genome-wide SNPs for the 22 publicly available *C. jejuni* strain 11168 derivative isolates from various UK laboratories were called using the Genpax pipeline.
- All metadata was taken from the Pascoe study [1].
- Dendrograms were created using a neighbour-joining method [2].
- Missing data was excluded from further analysis.

## References

1. Pascoe B, Williams JK, Cleland JK, McInnes G, Hutcheon S, Dyer M, Fisher J, Shaw S, Lopez RB, Chennamanchar C, Alan E, Vidal A, Fanning C, Esmail P, Parajuli A, Cogart A, Shewen M, Hartney T, Robinson T, Cook AJ, Collins M, Jolley KA, Maiden MCJ, Stevenson P, Munro G, Gow A, 2019. A reference-free approach to whole genome sequencing of *Campylobacter jejuni*. *Microb Drug Resist* 32:1169-1182. doi:10.1007/s12228-019-0947-4

2. Nei M, Li W-H, 1989. The neighbour-joining method: a new method for reconstructing phylogenetic trees. *Molecular Biology and Evolution*, Volume 6, Issue 4, July 1989. Pages 406-425. <https://doi.org/10.1093/molbev/6.4.406>

## Results

- Results were obtained addressing an average of 1.38 Mbp of *C. jejuni*, equivalent to 85.3% of the published genome.
- The AL11168.1 sequence which is the closest to the original isolate is clearly ancestral, as it should be, unlike in the original analysis.
- The strains strongly correlate with the sources and laboratories between which they were circulated, unlike in the original analysis.
- Genome-wide SNP analysis identified an average of 9 SNP differences to the reference (AL11168.1) per sample. In contrast to 29 SNP differences in the Pascoe et al. (2019) study.
- Two 0-SNP distance clusters were identified, totaling 12 of the 22 samples, of 5 and 7 members that are separated by 3 SNPs.

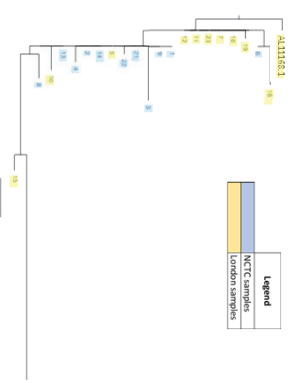
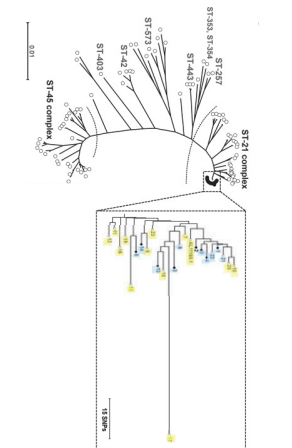


Figure 1: Phylogenetic tree from Pascoe et al. (2019), created using maximum-likelihood. Phylogeny was reconstructed in FastTree2 with the generalized time reversible substitution model.

Figure 2: Phylogenetic tree created using reanalyzed data from the Genpax pipeline. Two near 0-distance clusters are visualized.

**Table 1:** A) SNP distance matrix showing the pairwise comparisons between samples from Pascoe et al. (2019). B) SNP distance matrix showing pairwise comparisons between samples using the Genpax pipeline.

ID	Sample IDs (1-22)																								
	AL11168.1	5000	5001	5002	5003	5004	5005	5006	5007	5008	5009	5010	5011	5012	5013	5014	5015	5016	5017	5018	5019				
1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1			
2	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
3	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	
4	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
5	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
6	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
7	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
8	1	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
9	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
10	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
11	1	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1
12	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1
13	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1
14	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1
15	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1
16	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	1
17	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1
18	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1	1
19	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1
20	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1	1	1
21	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1	1
22	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	1
Reference AL11168	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1

**Table 2:** Metadata table from the Pascoe study [1]. London sourced samples cluster together and NCTC sourced samples cluster together (with the exception of isolate 5). Full table available in the Pascoe study.

Isolate	ID	Source	Original source	Archived
1	5000	Aberystwyth	NCTC	2005
2	5001	Aberystwyth	NCTC	2002
3	5002	Bristol	NCTC	2002
4	5003	London	NCTC	2002
5	5005	Glasgow	NCTC	2000
6	5006	Glasgow	NCTC (via Martin Stuvval)	2000
7	5007	Glasgow	London	2000
8	5008	London	NCTC	2004
9	5009	London	NCTC	2004
10	5010	London	NCTC	2000
11	5011	London	London	2000
12	5012	London	NCTC	2006
13	5013	London	NCTC	2006
14	5014	London	NCTC	2014
15	5015	Sheffield	London	2015
16	5016	Sheffield	NCTC	2010
17	5017	Sheffield	London (via Birmingham)	2010
18	5018	London	London	2002
19	5019	London	London	2002
20	5020	London	NCTC	2004
21	5042	Surrey	NCTC (via Cambridge)	2000
22	5043	Surrey	NCTC (via Cambridge)	2000
23	5044	Leicestershire	London (via Sheffield)	2013
Reference AL11168	NCTC	London	NCTC	2000

Phylogenetic trees reflect and are consistent with the relationship between the 23 isolates in terms of their known origin and distribution histories and can be used to trace the transfer of samples across different locations (See Figure 1 and Table 2).

Matrices show a significant reduction in pairwise SNP differences. The Genpax method exhibits reduced noise, as well as enhanced resolution and sensitivity. NCTC sourced samples tend to cluster together. This is also the case for London sourced samples.

## Conclusions

- The re-analysis generates an inferred set of relationships that is more parsimonious with respect to both the ancestral strain, and the pattern of distribution of strains between laboratories.
- The re-analysis does not confirm the original conclusions as to the substantial diversity and non-comparability of studies conducted in different laboratories using derivatives of *C. jejuni* strain 11168.
- The Genpax method shows increased sensitivity and decreased noise compared to previous analyses, and more accurate SNP-calling enables better identification and near-distance determination between isolates and the relationships between them, thereby enabling detection and definition of transmission and outbreaks in future applications.
- These results strongly validate the accuracy, resolution, and performance of a reference-free WGS analysis that is applicable to any strain, when measured against the performance of established methods with an ideal reference genome.

Check out our website and other supporting information



### Declaration

This research was funded by Genpax. Genpax is a bioinformatics company developing novel solutions that establish analysis strategies to maximize the usefulness of bacterial genome sequences in infection prevention and control.

# Openly comparable and scalable SNP-resolution analysis for *Listeria monocytogenes* using a novel genome comparison tool

Poster No. 269  
Date: 6/17/2023

Presented at ASM Microbe 23

Georgina R. Russell, Benedict J. Uttley, Dan Frampton, John F. Peden, Nigel J. Saunders: Genpax, London, United Kingdom

research@genpax.co  
+44 203 603 6669

## Introduction

- Listeria monocytogenes* is a food-borne pathogen with symptomatic infections resulting in a high hospitalization and mortality rate.
- L. monocytogenes* persists in food processing environments for extended periods and can be widely distributed through food transportation networks.
- Public health and food safety laboratories need to accurately and iteratively compare strains.
- It has a diverse clonal population structure and lacks the comprehensive set of reference genomes necessary to underpin traditional whole genome SNP analysis with optimum accuracy and resolution.
- It is therefore an ideal candidate for novel reference-free and scalable solutions that work at SNP-level resolution.

**Objective:** To assess the ability of the Genpax pipeline to integrate and compare findings from three previously published studies [1-3] that use a range of cgMLST and reference-genome SNP analyses from three different laboratories in Germany and Austria.

## Methods

- Sequences (n=587) from three different studies [1-3] spanning isolates from multiple European countries, mostly over the last 15 years were downloaded.
- These were analyzed with the latest Genpax developed pipeline and cluster analysis of SNP pairwise distances was conducted.

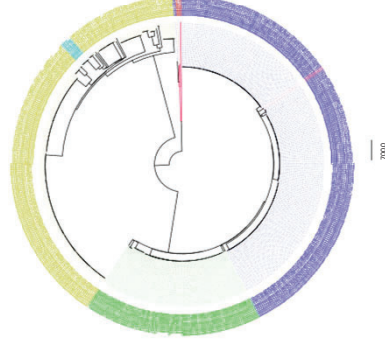
## Results

- The median of each input genome that was analyzed was 91.5% of the average genome length and 389,933 variant positions were called. In contrast the 1,701 cgMLST scheme for *L. monocytogenes* analyzes 53.5% of the genome (~1.5Mb).
- At a 20-SNP threshold, 54 clusters were identified ranging in size from 2 to 51 isolates: ten contained isolates from multiple studies.
- At a 2-SNP threshold, 41 clusters were identified ranging in size from 2 to 42 isolates: four contained isolates from multiple studies.
- We found previously unrecognized relationships spanning laboratories and countries of isolation.
- This analysis also clustered human and food isolates together (3-SNP threshold, not shown) providing links between source and patient.

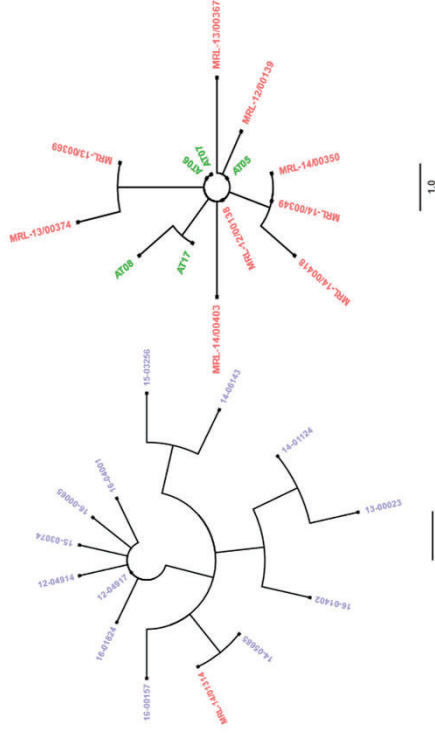
	Halbedel 2018	Hyden 2016	Schmid 2014	Source	Serogroup	Country	Lineage
Cluster 1	13	1	0	Human	IIa	Germany	Lineage II
Cluster 2	0	9	5	Food	IIb	Austria	Lineage I
Cluster 3	1	1	0	Human	IIa	Germany/Austria	Lineage II
Cluster 4	1	0	1	Human	IIb	Germany	Lineage I

**Table 1:** Clusters of isolates within 2SNPs containing isolates from multiple studies with selected metadata

## Results - continued



**Figure 1:**  
Concordance between serogroups and clustering  
Lineage I: I**lb**, I**vb**, I**b**-v1  
Lineage II: I**IIa**, I**IIc**  
Lineage III: I**IIla**, I**IIb**



**Figure 2:**  
Dendrogram of Cluster 1  
isolates, colour indicates source (see Table 1; isolate from the Hyden study is in red).

**Figure 3:**  
Dendrogram of Cluster 2  
isolates, colour indicates source (see Table 1; isolates from the Schmid study are in green. Isolates from the Hyden study are in red).

## Conclusions

- Clusters were found with isolates separated by time, source, and location.
- Our reference-free SNP-level resolution provided additional population structuring and transmission inference to cgMLST, and traditional whole genome SNP.
- These findings show the value and importance of being able to meaningfully compare strains over temporal and geographical space with SNP resolution at scale.

## Declaration

This research was entirely funded by Genpax.  
Genpax is a bioinformatics company developing novel solutions that overcome the limitations of established analysis strategies to maximize the usefulness of bacterial genome sequences in infection prevention and control.

Check out our website and other supporting information



## References

- Halbedel, S., Prager, R., Fuchs, S., Trost, E., Werner, G. and Fieger, A., 2018. Whole-genome sequencing of recent *Listeria monocytogenes* isolates from Germany reveals population structure and disease clusters. *Journal of clinical microbiology*, 56(6), pp. e00119-18.
- Hyden, P., Pietzka, A., Lemkh, A., Murer, A., Springer, B., Blaschitz, M., Indra, A., Huhulescu, S., Allerberger, F., Ruppitsch, W. and Samsen, C.W., 2016. Whole genome sequence-based serogrouping of *Listeria monocytogenes* isolates. *Journal of biotechnology*, 235, pp. 181-186.
- Schmid, D., Allerberger, F., Huhulescu, S., Pietzka, A., Anar, C., Kleta, S., Prager, R., Preussel, K., Aichinger, E. and Mellmann, A., 2014. Whole genome sequencing as a tool to investigate a cluster of seven cases of listeriosis in Austria and Germany, 2011–2013. *Clinical Microbiology and Infection*, 20(5), pp.431-436.

## Comparison of Reference-independent SNV-resolution genome comparison tool with existing retrospective methods.

Poster No. 6396  
Date: 16-June-2024

Georgina R. Russell, Luis F. Montemayor, John F. Peden, Nigel J. Saunders: Genpax, London, United Kingdom

research@genpax.co  
+44 20 3603 6699

### 1. Introduction

*Listeria monocytogenes* is a foodborne pathogen and the causative agent of listeriosis: an infection that is often self-limiting but when clinically significant is associated with high hospitalization and mortality rates. As such, it is a priority pathogen for public health surveillance and food safety.

Listeriosis cases within nationally agreed wgMLST thresholds are considered as potentially representing a single outbreak and prompt investigations by public health authorities. These thresholds vary between country e.g. 25 alleles in the USA and 10 in Canada.

Between 2016 and 2020 an outbreak of *Listeria* spanned multiple countries (US, Canada, France and Australia) and the source was determined to be *L. monocytogenes* from enoki mushrooms from a manufacturer in South Korea.

*L. monocytogenes* has a diverse clonal population structure with lineage specific differences but is typically analyzed with a single common reference strain. This can result in varying quality of analysis of outbreaks depending upon their similarity in sequence content, depending upon lineage.

**Objectives:** To assess the results of using the IDEM that uses a reference-independent solution that is applicable to all strains, compared to the published findings of the expert group at CDC using gold standard methods in a *L. monocytogenes* outbreak. To demonstrate the improvement in resolution possible through adding supplementary variant events (indels) with prototype methods for release in an upcoming version of the platform.

### 2. Methods

Read sets data (n=72) were gathered for analysis from the original study by the CDC. These were analyzed using the IDEM platform: a reference-independent SNV-level resolution tool currently available from Genpax. Analysis using IDEM generates continuously updated results, with the addition of new to existing isolate relationships with full SNV plus recombination event resolution within two hours of upload of a FASTQ.

Dendrograms were generated outside of the platform to generate readily comparable illustrations for this presentation. Using in-house tools, SNVs, complex events including recombination, and indels were independently verified to address events and differences identified between the new and previous analyses.

### 3. Results

Using IDEM, the median amount of each input genome analyzed was 2,420,186 bp. There are 119 variant sites identified in this outbreak, consisting of 117 SNVs and 2 recombination events, and the 72 samples in this study were all within 19 SNVs of each other.

**Congruence:** Comparison of the dendrograms shows broadly congruent relationship structures, as determined by IDEM and the original CDC analyses. There is consistent clustering down to small pairs and groups (demonstrated by pairs 1-4 and highlighted boxes in Fig 1 & Fig 2). As a result, we infer similar outcomes in terms of transmission chain, corroborating the original CDC results.

**0 distance pairs:** The original CDC analysis reported four 0-distance pairs. Pairs 1,2 and 3 (Fig 2) were corroborated by our analysis. Two SNVs (separately verified with high confidence) were identified between the samples in pair 4 (PNUUSA003927, CHAFRB2000090).

**Indels:** The robust determination of small indels is recognized to be analytically challenging. However, we have found (and verified) that indel analysis provides additional resolution in this outbreak dataset (results not shown).

### 4. Results

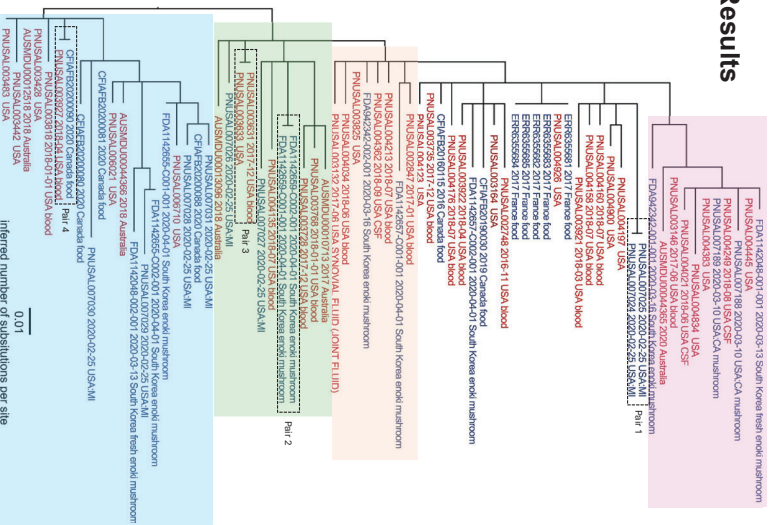


Fig 1  
Maximum-likelihood phylogenetic tree from Pereira et al., 2023. Boxes indicate clusters of interest. Text color indicates source, blue: food sample, red: clinical sample.

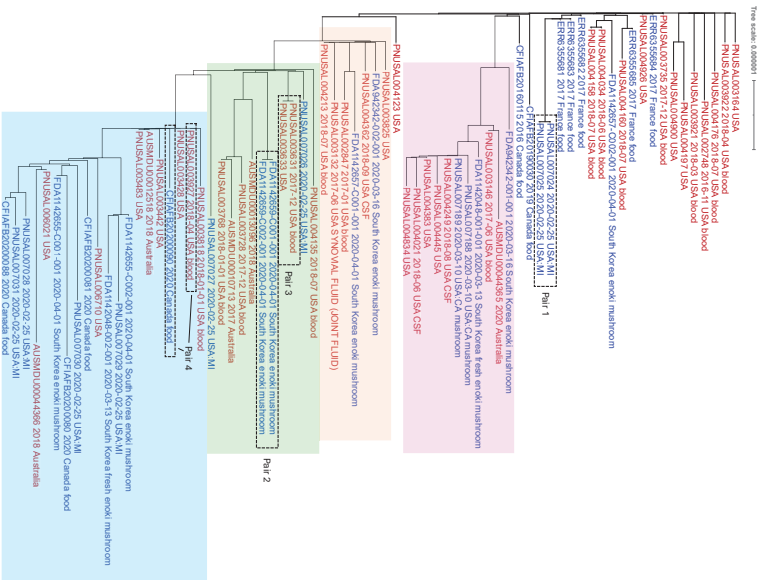


Fig 2  
Dendrogram produced using the current Genpax method available in IDEM (SNVs only). Boxes indicate cluster of interest. Metadata data and tip colors from Pereira et al., 2023 to match Fig 1.

### 5. Conclusions

The IDEM platform methodologies demonstrated at least equivalent performance to gold standard expert analysis. This platform enables prospective detection and continuously updated analysis of ongoing outbreaks, with reduced analysis time and high resolution. This methodology matched reference-based performance while being applicable to all strains in the species (separately demonstrated). Based on comparable branch length we can infer that the Genpax method has identified the SNVs previously reported, while performing filtering, and additionally identifying two recombination events.

We have found that the inclusion of indels provides additional resolution and structural information useful for transmission inference.

### Declaration

This research was entirely funded by Genpax. Genpax is a bioinformatics company developing novel solutions that overcome the limitations of established analysis strategies to maximize the usefulness of bacterial genome sequences in infection prevention and control.

### References

1. Pereira, E., Conrad, A., Tesfai, A., Palacios, A., Kandar, R., Kearney, A., Lucas, A., Jamieson, F., Elliot, E., Otto, M., and Kurndlla, K., 2023. Multinational Outbreak of *Listeria monocytogenes* Infections Linked to Enoki Mushrooms Imported from the Republic of Korea 2016–2020. *Journal of Food Protection*, 86(7), p.100101.

Check out our website and other resources



# Novel genome comparison tool reveals both false-positive and false-negative MRSA and MSSA strain identification and a failure to detect transmission-linked strains using phenotypic, PCR, and previous genomic strategies

Poster No. 266  
Date: 6/17/2023

James C. Littlefair, Benedict J. Uttley, Gareth Linsmith, Dan Frampton, John F. Peden, & Nigel J. Saunders: Genpax, London, United Kingdom

Presented at ASM Microbe 23

research@genpax.co  
+44 203 603 6869

## Introduction

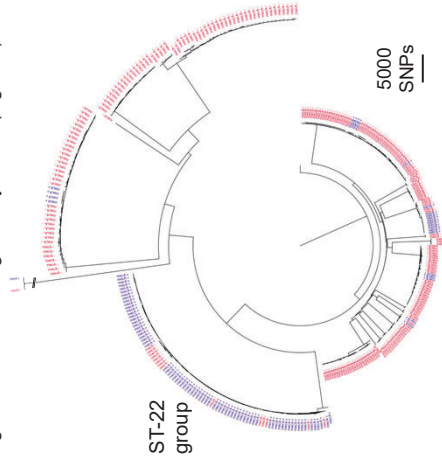
- *Staphylococcus aureus* is associated with >1 million deaths a year globally and is an AMR priority pathogen [1].
- Methicillin resistance in *S. aureus* is typically conferred by the presence of the *mecA* gene carried on the mobile SCCmec cassette, which is often spontaneously lost during culture, leading to discordant phenotypes [2].
- Transmission links may be overlooked due to the common approach of partitioning strains into MRSA and MSSA, especially where the determination of methicillin resistance can be unreliable using current methodologies such as indicator media, susceptibility testing, and PCR.
- Thus, it is necessary to use a genome comparison tool coupling accurate SNP-resolution strain identification with optimised gene detection.

## Methods

- All publicly-deposited readsets (n = 369, as of April 2023) from two published studies [3, 4] which collected MRSA and MSSA isolates in tandem within the same acute hospital between May 2017 and March 2019 were processed through the Genpax analysis pipeline. This included SNP-resolution strain identification and gene detection (which assigns gene detection confidence and identifies putative degeneracy).
- Isolates which met inclusion criteria (107/111 previously identified as MRSA and 243/258 previously identified as MSSA using screening methods including EUCAST susceptibility testing and PCR) were screened for *mecA/mecC* and *mupA/mupB*.
- The genes *lukS-PV / lukF-PV* (PVL), *tsst-1*, *eta*, *etb*, *etd* and *ete* were also screened for as they are important for the clinical management of *S. aureus*.

## Results – mecA

- 15/107 (~14%) of isolates previously identified as MRSA were both *mecA*- and *mecC*-.
- 8/243 (~3%) of isolates previously identified as MSSA were *mecA*+, all high confidence.
- Segregating strains into MRSA and MSSA based on susceptibility testing and PCR prior to genomic distance determination in the original analyses, led to missing transmission links in 9/48 (~19%) of transmission clusters (maximum pairwise distance of 15 SNPs).
- In the largest transmission cluster (n = 24), the original analyses excluded 4 isolates, of which 3 were identified as MSSA despite being *mecA*+, and 1 was genuinely *mecA*- (Figure 2).



**Figure 1:** NJ dendrogram showing the distribution of *mecA* among all processed isolates. *mecA*+ isolates are colored blue and *mecA*- isolates are colored red. *mecA*+ isolates are found mostly within ST22-group.

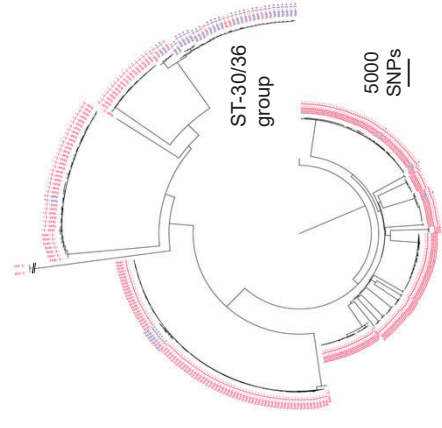
## Results – other genes

- 2 out of 3 of the previously identified mupirocin-resistant isolates were confirmed to be *mupA*+, both of which were high confidence (Table 1).
- 6 isolates were found to be PVL+ (5 MSSA; 1 MRSA). Both constituent genes of this complex were found with high confidence in all isolates.
- 41 isolates were found to be *tsst-1*+ (38 MSSA; 3 MRSA). All were found with high confidence bar one which was found with medium confidence. These were overwhelmingly found in ST30/36-group isolates, even outside of transmission clusters (Figure 3).
- 12 isolates were found to be *eta*+ (12 MSSA), all of which were high confidence except one which was medium confidence.
- 2 isolates were found to be *etd*+ (2 MSSA), both of which were found with high confidence.

Table 1: Isolates identified as mupirocin resistant in the original analysis or *mupA*+ in the Genpax reanalysis

Sample	<i>mecA</i> status	Available?	EUCAST mupirocin phenotype	<i>mupA</i> status	<i>mupA</i> + confidence
C-1310091	<i>mecA</i> +	Yes	Resistant	<i>mupA</i> +	High
C-1653983	<i>mecA</i> +	Yes	Susceptible	<i>mupA</i> + (frameshifted)	High
C-76440	<i>mecA</i> -	Yes	Resistant	<i>mupA</i> -	N/A
HND048-3	<i>mecA</i> -	Yes	Resistant	<i>mupA</i> -	N/A
P-00329	N/A	No	Resistant	N/A	N/A

**Figure 3:** NJ dendrogram showing the distribution of *tsst-1* among all processed isolates. *tsst-1*+ isolates are colored blue and *tsst-1*- isolates are colored red.



## Conclusions

- The reanalysis reveals putative false-positive and false-negative MRSA and MSSA determination by methods such as susceptibility testing and PCR, potential spontaneous gain/loss of *mecA* and *mupA*, and potential cryptic resistance.
- Partitioning strains into MRSA and MSSA leads to missed transmission-links, particularly when done unreliably, and indicates that transmission inference requires WGS of all clinically relevant *S. aureus* isolates supported by accurate, high-resolution, scalable genome analysis. To address the full diversity of isolates, a clonal reference is not required to get accurate and comprehensive results.
- The gene detection provided information on both resistance determinants and markers that indicate the requirements for different clinical management that were highly concordant with the strains underlying relationships, enabling confident interpretation of findings.

## References

1. WHO. *Staphylococcus aureus*. Weekly Epidemiol Rec. 2018;43(48):1059-1063.
2. Littlefair JC, Uttley BJ, Linsmith G, Frampton D, Peden JF, Saunders NJ. A multi-centre genomic investigation of MRSA and MSSA in a large acute hospital. *Antonie van Leeuwenhoek*. 2023;107(1):1-10. doi:10.1007/s10485-022-09810-3
3. Littlefair JC, Uttley BJ, Linsmith G, Frampton D, Peden JF, Saunders NJ. A multi-centre genomic investigation of MRSA and MSSA in a large acute hospital. *Antonie van Leeuwenhoek*. 2023;107(1):1-10. doi:10.1007/s10485-022-09810-3
4. Littlefair JC, Uttley BJ, Linsmith G, Frampton D, Peden JF, Saunders NJ. A multi-centre genomic investigation of MRSA and MSSA in a large acute hospital. *Antonie van Leeuwenhoek*. 2023;107(1):1-10. doi:10.1007/s10485-022-09810-3

## Declaration

This research was entirely funded by Genpax, a bioinformatics company developing novel solutions that overcome the limitations of existing genomic analysis tools to maximize the usefulness of bacterial genome sequences in infection prevention and control.

Check out our website and supporting information



# Precision Outbreak Surveillance of *Clostridioides difficile* Through Reference-free WGS SNP-resolution Analysis

Anastasia Pivnyuk, Dan Frampton, John F. Peden, & Nigel J. Saunders: Genpax, London, United Kingdom

## Introduction

*C. difficile* is an important cause of healthcare associated infections, with high prevalence in the US, and rising prevalence in the UK and elsewhere [1]. Partly related to demographic change, but also due to currently undefined factors, some clinically important clones also exhibit resistance to established treatments.

It has a highly diverse clonal population structure, that in some parts are more different than its generally considered to be a single species, and as a species it also exhibits substantial and relatively frequent recombination. This challenges traditional reference genome solutions. Sequence Typing (MLST, cgMLST, wgMLST) lacks the resolution needed for precise outbreak detection, transmission-inference and strain differentiation, especially in a species in which rates of diversification are currently considered to be lower than with others. Recombination in this species can also lead to conflicting findings when using different sequencing-based methodologies, depending upon the extent of Sequence Typing components that lie within recombined regions.

These species-specific issues may be one contributory reason for the persistence of Ce-htyping for this species, when WGS is becoming more widely used for others, even though this typing system uses the size of only 7 amplified product peaks and does not differentiate detailed isolate relationships within recurrent clinically important haplotypes.

**Objective:** To test a reference-free whole genome SNP analysis for *C. difficile* and compare its performance against published studies.

## Methods

The performance of the Genpax analysis pipeline was evaluated using datasets from two published studies that used WGS to determine transmission or relatedness. The results were compared to those presented in the original publications.

➢ 94 isolates from Knight et al. [2]

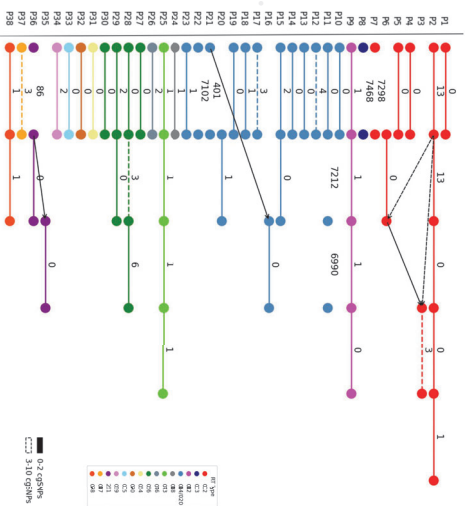
- Clinical isolates from 38 patients with recurrent *C. difficile* infection (CDI) are taken for WGS over 2 years of surveillance (94 isolates in total)
- Sequence Typing followed by cgSNP using ST-matched references, using Shiny v4.6.0 followed by snp-dists v0.8.2 with further cgSNP analysis. Study uses a mixture of ST types and RT types

➢ 40 isolates from Knight et al. [3]

- 40 RT10/14 isolates with 16 being from pigs (P) and 24 from human patients (H)
- SNV analysis using Small v0.7.6 compared with a high quality RT10/14 reference genome, mhap/SAWtools, and a combination of public (VCFtools, SnpEFF) and private tools to stringently filter results.

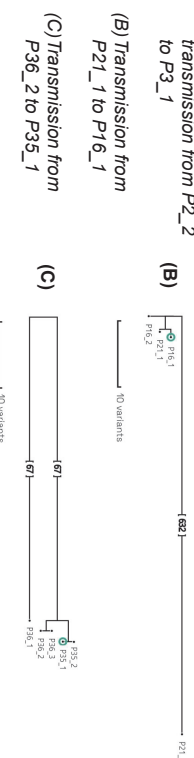
## Comparison with cgSNP analysis for epidemiological surveillance

- Thresholds used in the original study:  $\leq 2$ cgSNP for relapses;  $\geq 3-10$  for reinfection.
- In 54 out of 56 distances that could be compared the results using the Genpax pipeline are concordant with the original publication to the extent of the resolutions published. 100% of 0-2 SNP pairwise distances were resolved similarly. 3 out of 3 transmission events were identified correctly (black arrows in Fig 1)
- In Patient 2 (P2) there is a 13 SNP distance between the second (P2\_2) isolate, and the first (P2\_1) and third (P2\_3) isolates. However, P2\_1 and P2\_3 are identical, and P2\_2 is identical to both Patient 6 isolates (P6\_1, P6\_2). This suggests either there was another transmission event that was not previously recognized and reported in the original publication, or a sample swap during data submission.
- In Patient 28 event 3 for Patient 28 (P28) the Genpax IDEM pipeline identified more SNVs in event 3 than the original analysis (6 vs 0-2). By the studies' criteria this would alter the status from relapse to reinfection.



**Figure 2:** Trees generated by IDEM reporting interface, representing the transmission events.

- (A) Transmission from P6\_2 to P3\_1 and possible transmission from P2\_2 to P3\_1
- (B) Transmission from P2\_1 to P16\_1
- (C) Transmission from P36\_2 to P35\_1

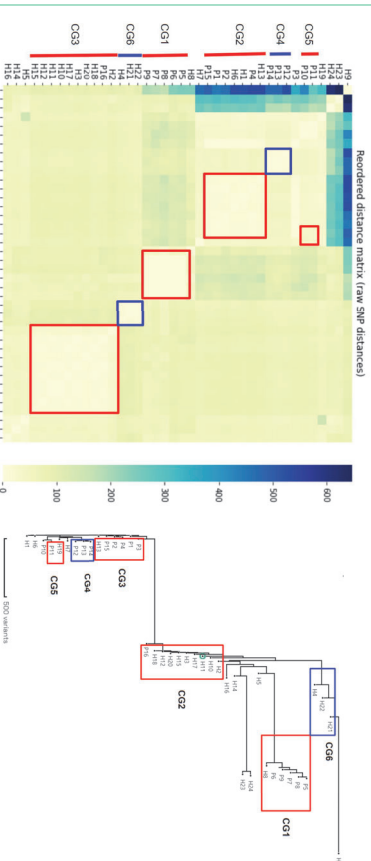


**Figure 1:**

Comparison between cgSNP distances and Genpax IDEM distances (prototype of patient report). Dashed lines and solid lines represent distances as per Knight et al. (2023) [2]. Each subplot represents a patient (P), each dot is an isolate take from that patient. Dots are shown in chronological order. Numbers above the lines represent SNP distances by Genpax IDEM pipeline. RT for isolates are shown on the right.

## Comparison with reference-based WGS

- 40 CDI RT10/14 isolates were assembled and mapped on RT10/14 reference strain CD630
- 6 clonal groups (CG1-6) with  $\leq 2$ SNV in the whole comparable genome (wgsNP) with CG1 corresponding to ST149, CG2 and CG6 to ST2, and the remainder being ST13
- The Genpax IDEM pipeline correctly identified all 6 clonal groups generating concordant results and detail to that generated in the original study by an expert group using a stringent analysis and a specifically selected local reference, while using a generally applicable natural reference-free solution.



**Figure 3: Distance matrix (output from Genpax IDEM and reorganized for visualization) between samples from Knight et al. (2017) [3]. Isolates were taken from human patients (H) or pigs (P). Left: Distance matrix generated by Genpax IDEM pipeline. Right: Phylogenetic tree generated by IDEM with CG1-6 highlighted**

## Conclusions

- The novel Genpax methodology accurately determined strain identity without the need for a closely related reference genome, prior knowledge of strain types, or clonal clusters.
- The achieved resolution matched, and in some cases surpassed, the resolution achieved and reported with core genome and whole genome SNP approaches performed by a specialist expert group, while being deliverable through a system that is not dependent upon sequence typing, is openly scalable, and with a 2-hour turnaround that can support both expert and non-expert users.
- The utilization of this innovative analysis pipeline can provide high resolution analysis that has the performance necessary to enable real-time outbreak detection and phylogenetic analysis for transmission inference in clinical settings in ways that are optimal, consistent, and not strain dependent.

## References

1. Kozak, C. et al. (2023). Burden of *Clostridioides difficile* infections in the United States. *New England Journal of Medicine*, 373(9), 825-834.  
 2. Knight, D. et al. (2023). Genomic epidemiology and transmission of a novel recurrent *Clostridioides difficile* infection in a Western Australian pig-breeding community. *Emerging Infectious Diseases*, 42, 607-615.  
 3. Knight, D. et al. (2017). Genomic epidemiology and transmission of a novel recurrent *Clostridioides difficile* infection in a Western Australian pig-breeding community. *Emerging Infectious Diseases*, 23, 100-108.  
 4. Saunders, N. J. et al. (2023). Disease Genetic Dependence and Spread of *Staphylococcus aureus* Infections in Microbiome. *Frontiers in Microbiology*, 7, 2138.

Check out our website and other supporting information



## Declaration

This research was entirely funded by Genpax. Genpax is a bioinformatics company developing novel solutions that overcome the limitations of established methods for the analysis of genomic data. The analysis pipeline was developed to maximize the usefulness of the whole genome sequences in infection prevention and control.



# Economic and health impact modeling of a whole genome sequencing-led intervention strategy for bacterial healthcare-associated infections for the USA

Poster No. 290  
Date: 6/17/2023

John M. Fox, Nigel J. Saunders & Susie H. Jerwood: Genpax, London, United Kingdom

Presented at ASM Microbe 23

research@genpax.co  
+44 203 603 6869

## Introduction

- Bacterial HAIs are a substantial source of global morbidity and mortality, resulting in increased length of hospital stay and high healthcare costs.
- Costs associated with HAI ranges from \$35 to \$45 billion in the USA [1].
- WGS has been promoted as a new gold standard for outbreak detection, but widespread adoption to date is limited.
- The upfront costs of WGS implementation have been identified as obstacles to adoption.
- Previous models addressing the impact of WGS on bacterial HAI have predicted a wide range of clinical and financial impacts from various methodologies and scope [2,3].
- It is timely to determine the economic viability and impact of routine diagnostic bacterial genomics.
- The aim of building this model was to evaluate the clinical and economic impact of a prospective WGS-led track and trace system of eleven common healthcare associated and AMR priority bacterial pathogens in England and the USA compared to the current standard of care, without WGS.

## Methods

- Using a synthesis of published models [2,3], inputs from national statistics, and peer-reviewed articles the clinical and financial impact models were created to address the most common nosocomial infections found in England and the USA.
- These are caused by *Staphylococcus aureus*, *Escherichia coli*, *Enterococcus* species, *Klebsiella* species, *Enterobacter* species, *Acinetobacter* species, *Stenotrophomonas maltophilia*, *Clostridioides difficile*, *Pseudomonas* species (mainly *Ps. aeruginosa*), *Citrobacter* species and *Serratia* species.
- All models were constructed, and analyses performed in Excel.

- Sensitivity analyses were conducted for each variable by varying the upper and lower limits within a wide range of available evidence.

## References

- Steer, J.A., Vayns, W., Zilka, N., Damm, C., Kabanek, J., Kelly, L., Price, M., L., Grayson, E., Kelly, and S. Alagarasu. 2017. Core components for effective infection prevention and control programmes: new WHO evidence based recommendations. *Antonie van Leeuwenhoek* 111: 641-650.
- Gordon, Louise G., Thomas M. Elliot, Brian Franks, Brent Michael, Philip L. Russo, David L. Paterson, and Patrick N. A. Harris. 2021. Budget impact analysis of routine using whole-genome sequencing of six multidrug-resistant bacterial pathogens in Queensland, Australia. *BMJ Open*. 11: e018986.
- Kumar, P. A., J. Sandermann, E. M. Math, G. M. Snyder, J. W. Marsh, L. H. Harrison, and M. S. Roberts. 2021. Method for Economic Evaluation of Bacterial Whole Genome Sequencing Surveillance Compared to Standard of Care in Detecting Hospital Outbreaks. *Clin Infect Dis*. 72: e94-95.

## Results

- The model shows bacterial HAI currently cost the NHS in England around £3 billion annually.
- WGS-based surveillance is predicted to cost £61.1 million associated with the prevention of 74,408 HAI and 1,257 deaths based on a cluster detection and intervention turnaround time of seven days.
- The net cost saving was £476.3 million, of which £65.8 million were from directly incurred savings (antibiotics, consumables etc.) and £412.5 million from opportunity cost savings due to re-allocation of hospital beds and healthcare professionals.
- The USA model indicates that the bacterial HAI care baseline costs are around \$18.3 billion.
- WGS surveillance cost \$169.2 million and resulted in a net saving of ca.\$3.2 billion, while preventing 169,260 HAIs and 4,862 deaths, also based on a cluster detection and intervention turnaround time of seven days.
- The average reduction of total infections was 18% from using WGS.
- This clinical impact model estimated *S. aureus* to be the most common bacterial HAI in both England and the USA, and the cause of most deaths in the USA, with 17,176 deaths annually. In comparison, *E. coli* was responsible for the most nosocomial deaths in England, with 1,456 deaths.
- The model predicts a return to the hospitals of £7.83 per £1 invested in diagnostic WGS in the UK, and US\$18.74 per \$1 in the USA.

## Results – England sensitivity analysis

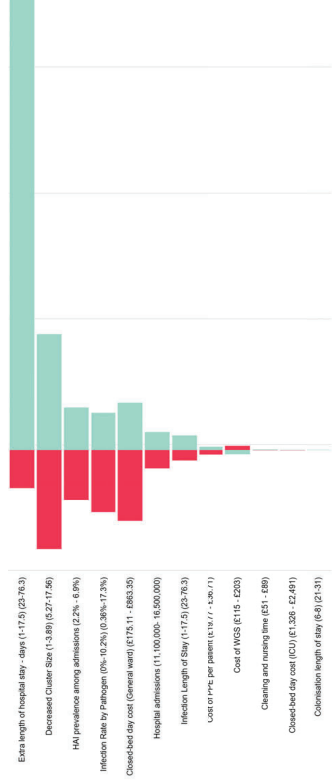


Figure 1: All savings are in millions of pounds relative to a base saving of £476.3 million. Low value (EM) represents the lower input for each variable and equal to the High value (EM) represents the higher model input. The higher value for Cost of WGS was the only variable to reduce savings

Check out our website and supporting information



## Declaration

This research was entirely funded by Genpax. Genpax is a for-profit company that has established analysis strategies to maximize the usefulness of bacterial genome sequences in infection prevention and control.

## Results - continued

Table 3: Estimated number of HAI and death for each pathogen group using current practice and estimated number of avoided infections and death with a WGS-based system

Organism	Current practice		WGS		Death avoided
	HAI	Death	HAI avoided	Death avoided	
<i>Staphylococcus aureus</i>	100,044	1,009	10,765	109	
<i>Stenotrophomonas maltophilia</i>	1,842	14	174	2	
<i>Enterococcus</i> spp.	1,426	14	208	3	
<i>Enterobacter</i> species	61,684	1,292	25,572	519	
<i>Klebsiella</i> species	58,559	986	3,793	64	
<i>Enterobacter</i> species	16,074	294	6,105	112	
<i>Acinetobacter</i> species	9,171	128	2,201	31	
<i>Clostridioides difficile</i>	41,685	375	7,816	70	
<i>Pseudomonas</i> species (mainly <i>Ps. aeruginosa</i> )	8,337	797	1,338	6	
<i>Citrobacter</i> species	14,173	88	1,338	6	
<i>Serratia</i> species	35,227	92	74,408	1,227	
<b>Total</b>	<b>302,073</b>	<b>3,074</b>	<b>217,43</b>	<b>1,848</b>	

Table 4: Estimated differences in costs in England and the USA with current practice and WGS surveillance

Direct hospital resources	England		USA	
	Current practice	WGS	Current practice	WGS
WGS	£ 61,168,073	£ -	\$ 169,245,276	\$ -
PPE	£ 266,750,089	£ 212,408,937	\$ 740,286,599	\$ 585,475,953
Consumables	£ 63,113,318	£ 51,710,819	\$ 177,685,114	\$ 144,713,545
Antibiotic treatment	£ 319,095,715	£ 299,090,469	\$ 1,069,939,350	\$ 1,303,622,919
Total cost of hospital resources	£ 840,558,122	£ 563,798,288	\$ 2,593,811,088	\$ 2,281,257,084
Allocation of hospital beds and healthcare professionals		WGS		
Extra length of stay - General Ward	£ 2,220,892,705	£ 1,820,127,623	\$ 14,278,211,667	\$ 12,501,524,165
Extra length of stay - ICU	£ 4,234,642	£ 3,479,079	\$ 168,859,625	\$ 146,123,367
Cleaning and nursing time	£ 24,289,160	£ 19,656,643	\$ 68,180,372	\$ 55,534,925
Cost of infection prevention and control team	£ 32,200,330	£ 26,918,080	\$ 203,127,450	\$ 165,544,527
Total cost of allocation of hospital beds and healthcare professionals	£ 2,282,716,838	£ 1,870,227,405	\$ 15,783,378,114	\$ 12,888,616,985
Overall total	£ 2,932,275,960	£ 2,454,019,702	\$ 18,301,988,182	\$ 15,129,873,979
Overall cost savings with WGS surveillance	£ 476,256,256	£ -	\$ 3,172,115,203	\$ -

## Conclusions

- This economic analysis indicates that substantial savings and improvements in clinical outcomes are generated using proactive clinical genomics of bacterial pathogens.
- Sensitivity analyses demonstrate that calculations involved in length of hospital stay affect the financial outcome the most.
- The largest savings are associated with improved use of healthcare resources, due to avoidance of prolonged patient stays and the ability to use facilities more effectively.
- Savings were retained when tested by sensitivity analyses, which showed small impacts from the costs of WGS.

## Limitations of cgMLST in current practice

Previously, scalability challenges of high-resolution analyses using SNVs, or SNV with other differences (e.g. indels, recombination), meant that detailed analysis for outbreak detection and investigation either had to address only small numbers of isolates, be periodically consolidated in major centres, or had to be preceded by an initial low-resolution but practically more deliverable step. Commonly this is a form of Sequence Typing (MLST, cgMLST/cgST, or wgMLST). After this, strains of the same type or with a certain range of differences are selected for more detailed analysis of some kind. cgST has recognized limitations, as do the minimum spanning trees generated from it. But it has higher resolution than MLST, and can group strains into smaller and more analyzable groups for subsequent investigation. Although this can still present a scaling challenge to compare strains beyond limited time-frames or geographies, especially for the most clinically important and common clones.

Because cgST correlates well with other phylogenomic information in population-scale studies when compared to findings using more detailed analyses, it has been assumed that it will perform reliably and sensitively in selecting the strains that are potentially parts of outbreaks and transmission-linked clusters. However, the nature and scale of differences spanning populations does not necessarily reflect performance in distinguishing more closely related strains, with a method that is reported to generate slightly different results when using different assemblers, less than 40x coverage, and varied addressed alleles, even when reanalyzing the same sequencing files (e.g. Abdel-Glil *et al.* J. Clin Micro, 2022). However, because cgST when used in clinical and epidemiology studies is used as a pre-selection step, this assumption is not normally tested because ungrouped strains are not compared in detail.

The resources within IDEM perform analysis in detail throughout. Further, because it is a natural-reference free solution, the results remain fully comparable because they are not divided into groups that have been analyzed using different references, or references with different degrees of divergence from analyzed isolates. This enables the performance of cgST in the identification of putative outbreak members and transmission-linked isolates to be assessed.

Such an analysis has been performed using and comparing findings to the results of two large published studies of *Campylobacter* species. One of *C. jejuni*, the other of *C. coli*. This genus illustrates several challenges for WGS analysis, because it undergoes relatively frequent recombination resulting in a panmictic population structure. This means that defining genes and alleles may be difficult using the search methods used in cgST (BLAST), and that it is difficult (perhaps impossible) to establish a consistent set of good quality reference genomes. The analysis methods used within IDEM are natural reference free, and overcome these barriers. Thereby enabling a test of the performance of cgST in these species. The findings have potential implications not only for analyses performed solely using cgST, but also for all pipelines that use it as a preliminary step.

# Limitations of sequence typing for isolate inclusion in outbreak investigations

Presented at ESCMID 2024

Poster No. P2151  
Date: 27/04/2024

S. Jackowska <sup>1</sup>, D. Frampton <sup>1</sup>, J.F. Peden <sup>1</sup>, N.J. Saunders <sup>1</sup>, <sup>1</sup>Genpax - London (United Kingdom)

Email via: [research@genpax.co](mailto:research@genpax.co)  
Tel: +41 44 586 86 06

## Introduction

- Campylobacteriosis, the leading cause of gastrointestinal disease in the EU, is primarily caused by two species: *Campylobacter jejuni* and *Campylobacter coli*.
- Correct detection of outbreaks is key to infection control.
- MLST and cgMLST are commonly used to identify clusters for further analysis.
- Here, we assess their performance against an openly scalable solution comparing strains at SNV resolution to determine their sensitivity and specificity for likely transmission-linked strains, using data from two recent publications. (Hsu *et al.*, Harrison *et al.*)

## Results - Overview

- Of the pairs which were linked by IDEM at a distance of  $\leq 10$  SNV, 80/318 (25%) of *C. jejuni* and 1374/2920 (47%) of *C. coli* isolate pairs were not clustered by cgMLST at an allelic difference of 10 (AD10); representing thresholds used for outbreak detection.
- Using a higher threshold of AD25, a portion (21/318 (7%)) of the  $\leq 10$  SNP *C. jejuni* isolates remained ungrouped.
- Strikingly, 4/39 (10%) *C. jejuni* and 47/117 (40%) *C. coli* 0-SNV isolates are also not grouped at AD10.
- Examination of reported sequence types (ST) showed that 11/318 (3%) of *C. jejuni* and 200/2920 (7%) *C. coli* sample pairs were identified as different sequence types (STs) while differing by  $\leq 10$  SNVs.
- In contrast, samples belonging to the same ST had a median pairwise SNV distance of 568 (IQR: 75-1961) for *C. jejuni* and 95 (IQR: 66-153) for *C. coli* respectively.

## Campylobacter jejuni results

Figure 1. Distribution of cgST/AD/MLST pairs at each SNV distance threshold for *C. jejuni*

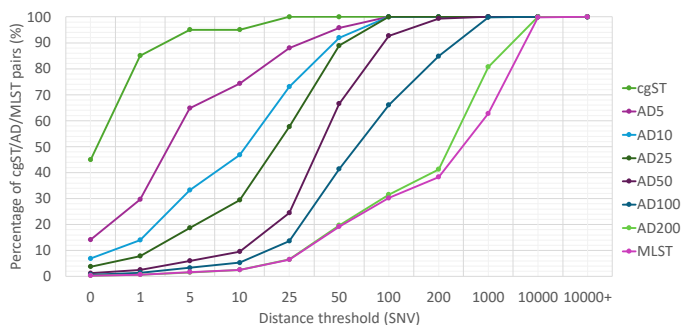
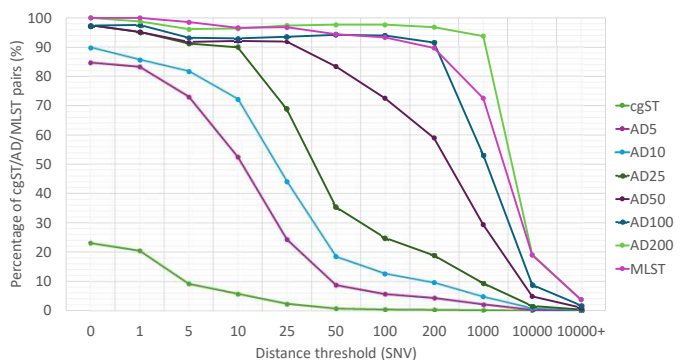


Figure 2. Percentage of pairs at each SNV distance threshold that have the same cgST/AD/MLST classification for *C. jejuni*



## Declaration

This research was entirely funded by Genpax. Genpax is a bioinformatics company founded in 2020 seeking to develop novel solutions that overcome the limitations of established analysis strategies to maximize the usefulness of bacterial genome sequences in infection prevention and control.

Copyright and IP owned by Genpax

## Methods

- In total, 3762 readsets were processed through the Genpax IDEM platform.
- All pairwise SNP distances for 844 *C. jejuni* and 2918 *C. coli* isolates were determined using the platform.
- Allele difference groups, as well as cgST and MLST were defined and derived for each sample in the two publications. (Hsu *et al.*, Harrison *et al.*)
- The results were combined and compared.

## Key

<b>AD0</b>	Allele difference of 0
<b>AD5</b>	Allele difference of 5
<b>AD10</b>	Allele difference of 10
<b>AD25</b>	Allele difference of 25
<b>AD100</b>	Allele difference of 100
<b>AD200</b>	Allele difference of 200
<b>cgST</b>	Core genome sequence type
<b>MLST</b>	Multi-locus sequence typing

## Campylobacter coli results

Figure 3. Distribution of AD/MLST pairs at each SNV distance threshold for *C. coli*

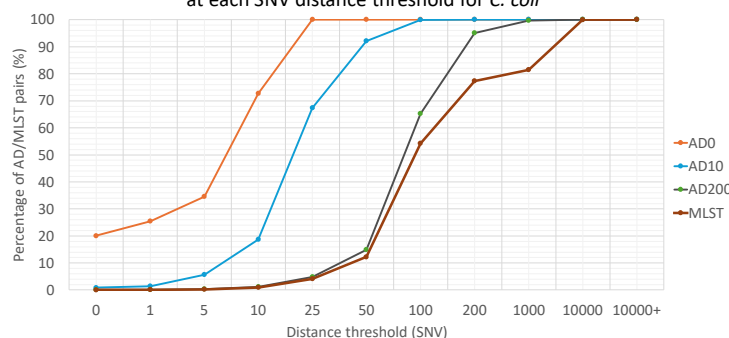
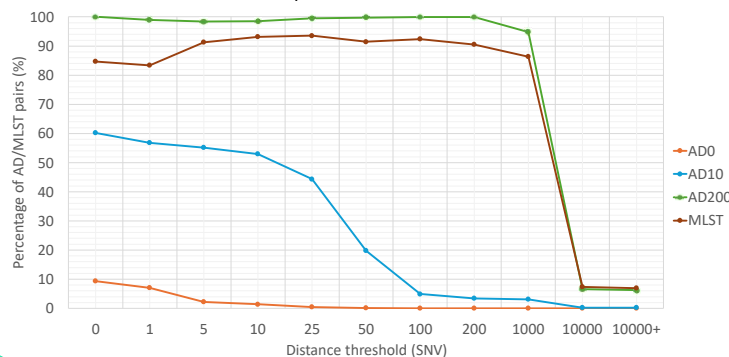


Figure 4. Percentage of pairs at each SNV distance threshold that have the same AD/MLST classification for *C. coli*



## Conclusions

- Both MLST and cgMLST under- and over-predict sample linkage in both published studies (generate both false negative and false positive results).
- Across the two studies, 211 pairs of isolates  $\leq 10$  SNV distance identified by IDEM do not share the same sequence type (ST) and 1466 pairs do not share the same cgMLST cluster at AD10, which would normally exclude them from being identified as potential members of outbreaks and subsequent more detailed comparisons and analysis.
- In addition, of the 6% of isolate pairs identified as the same ST (300959/4514920), only 1% (3039) fall within a distance of 10 SNVs. The majority (99%) of same ST pairs have a greater SNV distance than would normally be considered to indicate outbreak / transmission connection.
- This performance indicates that Sequence Typing (MLST or cgMLST/cgST) is not an optimal first stage analysis for the detection and investigation of transmission-linked strains and outbreaks in these species.
- This may reflect the combined effects of issues inherent to the underlying methodology in the context of the highly recombining and panmictic nature of the genus. Similar analyses seem warranted in other species.

## References

1. Hsu C-H, Harrison L, Mukherjee S, Strain E, McDermott P, Zhang Q, Zhao S. Core Genome Multilocus Sequence Typing for Food Animal Source Attribution of Human *Campylobacter jejuni* Infections. *Pathogens*. 2020; 9(7):532. <https://doi.org/10.3390/pathogens9070532>
2. Harrison L, Mukherjee S, Hsu C-H, et al. Core genome MLST for source attribution of *Campylobacter coli*. *Frontiers*. June 21, 2021. Accessed March 19, 2024. <https://www.frontiersin.org/journals/microbiology/articles/10.3389/fmicb.2021.703690/full>.

Check out our website:



## INTRODUCTION TO THE FOLLOWING HEALTH ECONOMICS PAPER

A clinical genomics solution must provide an increase in patient and public safety with improved patient care and improved management of healthcare resources, while also being economically viable. Several previous publications indicate that this is the case for proactive clinical genomics for IPC. (These are cited in the following paper.)

However, these papers do not always include all attributable costs of sequencing, analysis, staffing and other infrastructure costs, and others are dated or don't work with the most current epidemiological information or a more limited set of species. This new analysis draws upon the best of this previous work, updates it, addresses a core set of healthcare-associated hospital-acquired AMR priority pathogens, using inclusive current real world costs, and a best-available set of epidemiological information. The publication also makes its model available in an Excel format to enable local hospitals or others to modify it to generate locally informed versions and ongoing updates. Key findings of this modelling include:

- Addressing NHS England as an example:
  - it would be possible to **save over 70,000 bed days per year**, which is the equivalent of building and fully equipping and staffing a new 200-bed hospital with full occupancy
  - it is possible to **prevent more than 1,200 avoidable hospital care associated deaths**; representing 10-20% of estimated avoidable hospital deaths per year
  - it is possible to **save at least £480 million per year** in avoidable costs
- There is **no economic obstacle to adoption**, because the savings to hospitals and healthcare delivery systems are considerably greater than the costs of adopting and delivering proactive bacterial genomics surveillance for IPC.
- **Improved patient safety and actions to contain and prevent the spread of AMR** within the hospital can be achieved at negative costs.
- The **hospital-level costs savings** are dominated by improved use of healthcare resources, such that large savings remain with wide variations in the costs of sequencing and analysis.
- Larger savings and proportionate returns on investment are available in the US than the UK

The only remaining obstacles to adoption are sequencing, which is now available in-house to any modern laboratory capable of typical microbiology and pathology services or through external services, and the expertise and resources that are required for analysis and interpretation that are now openly available from Genpax.

Finally, it should be noted that **these models are intentionally conservative**. They do not include savings from other activities such as combined environmental, healthcare worker, and pre-admission screening; the additional benefits of addressing non-AMR/antibiotic sensitive strains with similar transmission mechanisms and clinical consequences (e.g. MSSA which has a 20 to 30% mortality); nor additional species. They also do not include costs associated with exceptional responses such as ward closure, rebuild and refits, equipment replacement, insurance company non-payment or claw-backs, or legal liabilities for hospital transmitted infections. Nor the savings from avoidable responses to 'non-outbreaks' that suspected on epidemiological grounds are caused by strains are unrelated and not connected, or being able to demonstrate that infections were not caused by hospital-associated strains.

# Economic and health impact modelling of a whole genome sequencing-led intervention strategy for bacterial healthcare-associated infections for England and for the USA

John M. Fox, Nigel J. Saunders and Susie H. Jerwood\*

## Abstract

Bacterial healthcare-associated infections (HAIs) are a substantial source of global morbidity and mortality. The estimated cost associated with HAIs ranges from \$35 to \$45 billion in the USA alone. The costs and accessibility of whole genome sequencing (WGS) of bacteria and the lack of sufficiently accurate, high-resolution, scalable and accessible analysis for strain identification are being addressed. Thus, it is timely to determine the economic viability and impact of routine diagnostic bacterial genomics. The aim of this study was to model the economic impact of a WGS surveillance system that proactively detects and directs interventions for nosocomial infections and outbreaks compared to the current standard of care, without WGS. Using a synthesis of published models, inputs from national statistics, and peer-reviewed articles, the economic impacts of conducting a WGS-led surveillance system addressing the 11 most common nosocomial pathogen groups in England and the USA were modelled. This was followed by a series of sensitivity analyses. England was used to establish the baseline model because of the greater availability of underpinning data, and this was then modified using USA-specific parameters where available. The model for the NHS in England shows bacterial HAIs currently cost the NHS around £3 billion. WGS-based surveillance delivery is predicted to cost £61.1 million associated with the prevention of 74 408 HAIs and 1257 deaths. The net cost saving was £478.3 million, of which £65.8 million were from directly incurred savings (antibiotics, consumables, etc.) and £412.5 million from opportunity cost savings due to re-allocation of hospital beds and healthcare professionals. The USA model indicates that the bacterial HAI care baseline costs are around \$18.3 billion. WGS surveillance costs \$169.2 million, and resulted in a net saving of ca.\$3.2 billion, while preventing 169 260 HAIs and 4862 deaths. From a 'return on investment' perspective, the model predicts a return to the hospitals of £7.83 per £1 invested in diagnostic WGS in the UK, and US\$18.74 per \$1 in the USA. Sensitivity analyses show that substantial savings are retained when inputs to the model are varied within a wide range of upper and lower limits. Modelling a proactive WGS system addressing HAI pathogens shows significant improvement in morbidity and mortality while simultaneously achieving substantial savings to healthcare facilities that more than offset the cost of implementing diagnostic genomics surveillance.

## Impact Statement

This article estimates the impact of effective whole genome sequencing-based surveillance for tracking and intervening in bacterial nosocomial outbreaks of the 11 most common healthcare-associated infection (HAI) species in both England and the USA. The projected outcome would be to reduce the bacterial morbidity and mortality of HAI in hospitals while simultaneously reducing the cost of patient care and increasing the wider cost savings of England and the USA by £478.3 million and \$3.2 billion respectively, with more efficient use of hospital resources.



## Transformative differences to practice enabled by Genpax IDEM

### From reactive to proactive WGS for hospital IPC

Current use of WGS-based genomics is largely limited to the investigation of otherwise suspected outbreak-associated isolates. It has a reactive and remediation function, but it is not primarily a detection and prevention resource. In the rare settings where it is currently used proactively, it is typically limited to low resolution analysis (e.g. cgMLST), followed by later detailed analyses of selected strains; or to limited sampling space and time-windows to limit the number of strains compared (normally to not more than 50-100). This is because of multiple factors, and requires local expert teams, all of which are overcome by the Genpax IDEM solution.

Using IDEM information from proactive sequencing of targeted pathogens can **identify transmission-connected infections, outbreak clusters**, and for tracking and source identification, that can be augmented by environmental surveillance. Enabling rapid responses, without depending upon other indicators, to contain high-risk strains within hospital and other healthcare environments. It can also **identify sites with more transmissible, virulent, and resistant strains** for targeted containment, pre-admission screening, and follow-up interventions.

By detecting otherwise unrecognized connections between infections early, from the second isolate of a strain transmitted or acquired in the hospital, or even the first isolate of a strain identified as high risk, IPC responses can be targeted **to prevent onward transmission and reduce the size of outbreaks** and the number of healthcare-associated infections that occur. This will result in smaller outbreaks, fewer outbreaks, more rapid detection and remediation of routes and sources of transmission, and **greater protection of both patients and staff** from emergent pathogens. Thus, protecting patient safety, the biosecurity of hospital environments, and reducing direct and indirect costs of care, that conservative models show cover the costs of WGS-led rapid surveillance many times over.

Due to its **open scalability**, the more surveillance data collected within IDEM, from patients of the environment, **the more informed the IPC team becomes**, and with it their ability to rapidly deliver effective interventions and identify sources. Thereby, **protecting both patients and the hospital** from the ever-increasing biosecurity challenges of more resistant and virulent healthcare adapted strains. The new paradigm being maximal immediate patient prevention, coupled with the creation and maintenance of the safest possible healthcare environment, though new IPC capabilities enabled by proactive pathogen sequencing, analyzed and connected through IDEM.

### So what?

- Smaller outbreaks
- Fewer outbreaks
- Better infection prevention
- Improved patient safety
- Improved hospital reputation and IPC practices
- Saves money while saving lives

## The clinical and monetary value of knowing that you don't have an outbreak

IPC resources are limited and must be used with maximum efficiency to protect the patients and hospital environments from highly transmissible, virulent, and resistant strains. Being able to distinguish outbreaks from non-outbreaks and to know which infection and colonizations are transmission-linked or not is fundamental to infection prevention. The pursuit of connections between patients with infections that are not linked, and the institution of control measures for coincidental but not connected similar infections, consumes limited resources and distracts IPC teams from investigating genuinely linked infections, and confuses those investigations. Meanwhile professional medical practice requires precautions and actions to address possible risks to patients, so recognizing unconnected cases is important.

**Analysis in IDEM provides pre-emptive information.** A team can see whether a *C. difficile* isolate is part of the hospital associated strains, or one that was imported by the patient; and costly cohorting and outbreak response meetings and actions avoided. When increased incidence of infections are noted, the proactively sequenced isolated information can be used to determine whether there is an outbreak, and which isolates are members of which part of coincident outbreaks. This latter issue is greatest relevance to early stage adopters in which multiple long-standing hospital-associated outbreaks often coexist.

Examples in early users of IDEM have already illustrated situations in which non-outbreaks have been rapidly recognized, allowing IPC teams **not to spend resources inappropriately** and to focus upon other impactful activities. Substantial IPC resources can be wasted investigating non-outbreak strains, where IDEM would enable more correctly targeted and effective responses. And, IDEM consistently identifies multiple transmission-linked clusters of only 2 or 3 isolates that would otherwise not have been recognized at all due to not being alert organisms, and therefore not being on the 'radar'.

### So what?

- Cost savings from avoiding unnecessary IPC meetings, interventions pending investigations, and precautionary ward closures for non-outbreaks.
- More effective IPC from better and more detailed information on connected strains and transmission chains
- More effective IPC from reducing time investigating non-outbreaks, freeing time to focus upon real outbreaks and transmission events and other prevention-focused activities.
- Defensible positions with patients demonstrably infected by patient-linked, rather than hospital associated strains

## Detailed and timely information communicated directly to the IPC team

IDEM is not genomics only for the bioinformaticians, report-focused epidemiologists, and academics. IDEM is about optimized genomics presented directly to those who can act on it to improve patient and public safety in a way that can be easily understood; which means the Infection Prevention and Control teams, and the patient-facing clinical staff. Batch and QC reports are generated for the sequencing laboratory to ensure that the data generated is of consistently high quality and for any issues to be quickly addressed. An individual sample report is also generated for record keeping. But, the most important report is an interactive, continuously updated, easy to access, interpret, and interrogate resource for front-line clinicians.

The IPC reporting system enables the IPC team, from the nurses and infectious disease physicians, to the Director of Infection Prevention and Control (and equivalents) to access and act upon the information as soon as the sequencing data has been analyzed and integrated. (Typically within 2 hours of a sequencing run being completed, or first thing in the morning if a run has completed during the night.) There are no delays or intermediate interpretations between the frontline patient-facing teams and the results of the sequencing analysis. With minimal induction training, any user with a professional understanding of infection prevention and control has direct access to the usable information, and is empowered to act upon it. No long per-sample multiple page pdf reports from multiple samples to work through, no large tables of information that doesn't impact clinical decision making, no periodic data consolidation. Just simple, focused, clearly communicated information in a format tailored for the people who need it that highlights connected strains and potential outbreak clusters. It also provides access to other findings of relevance to IPC and containment of more hazardous or resistant isolates such as resistance and virulence genes.

### So what?

- The information gets directly to the people who need it, not delayed or stuck in the lab
- Patient-focused responses are enabled more quickly
- Fewer and smaller outbreaks
- Improved patient care and safety



## Making immediate connections in Healthcare Surveillance and Public Health

**No more waiting for periodic consolidation of data.** All strains with evidence of outbreak- and/or transmission-connections are immediately detected and accessible through IDEM. Within a target turnaround time of under 2 hours from receipt of the FASTQ (the file from the DNA sequencer).

IDEM does not use a Sequence Typing step, thereby avoiding associated errors, avoiding missing some outbreak members, and putting strains into groups that are too large for optimal detailed scalable analyses. The IDEM pipeline is also **natural reference-free**, meaning that all strains are comparable at the same high-resolution regardless of reference sequence availability or quality; and that all strains are comparable and connectable with an optimum high-resolution analysis. Thereby **avoiding missing outbreak members** by separating them prior to detailed comparisons or using difference references. Thus, all sequenced isolates are directly compared in a system that is continuously updated to detect transmission connections and determine relationships in real clinical time.

Once weekly, or other periodic consolidation (as is typically practiced in reference laboratories, and other settings) **is no longer necessary to identify connections.** Previously a slow and computationally intensive costly process that could not reasonably be performed on a rolling basis for every newly sequenced isolate, this is now integrated into the core analysis process. Thus a delay associated with the final stage of analysis when addressing larger numbers of isolates in detail of 1 to 7 days, or more, is avoided in recognition and response times; and the actionable information can be obtained in close to the time it takes to isolate and sequence the DNA.

### So what?

- Surveillance in healthcare can operate over longer time periods, necessary to detect some outbreaks linked to the environment, healthcare workers, or patient re-admissions
- Critical information for public health responses available more quickly
- No consolidation delays with public health teams able to respond on a rolling basis, rather than following periodic updates
- Can be used to connect data between labs, enabling wider connected surveillance and closer to isolation laboratory sequencing
- Faster responses in pathogen eradication programs, such as for *M. tuberculosis*
- Connection possible between all historic, surveillance, and clinical samples with ability to look over multiple years and origins

Species	IDEM status
<i>Acinetobacter baumannii</i>	Release June 24
<b><i>Campylobacter coli</i></b>	Available
<b><i>Campylobacter jejuni</i></b>	Available
<b><i>Campylobacter lari</i></b>	Available
<i>Citrobacter freundii</i>	Release July 24
<b><i>Clostridioides difficile</i></b>	Available
<i>Corynebacterium diphtheriae</i> complex	Release July 24
<i>Cronobacter sakazakii</i>	Release July 24
<i>Enterobacter asburiae</i>	Release June 24
<i>Enterobacter cloacae</i> (and subspecies)	Release June 24
<i>Enterobacter hormaechei</i> (and subspecies)	Release June 24
<i>Enterobacter kobei</i>	Release June 24
<i>Enterobacter ludwigii</i>	Release June 24
<i>Enterobacter rogenkampii</i>	Release June 24
<b><i>Enterococcus faecalis</i></b>	Available
<b><i>Enterococcus faecium</i></b>	Available
<b><i>Escherichia coli</i></b>	Available
<i>Haemophilus influenzae</i>	Release July 24
<i>Klebsiella aerogenes</i>	Release July 24
<i>Klebsiella oxytoca</i>	Release June 24
<b><i>Klebsiella pneumoniae</i></b>	Available
<b><i>Klebsiella quasipneumoniae</i></b>	Available
<b><i>Klebsiella varicola</i></b>	Available
<i>Legionella pneumophila</i>	Release July 24
<b><i>Listeria monocytogenes</i></b>	Available
<b><i>Mycobacterium tuberculosis</i></b>	Available
<i>Mycobacteroides abscessus</i>	Release July 24
<i>Neisseria gonorrhoeae</i>	Release July 24
<i>Neisseria lactamica</i>	Release July 24
<i>Neisseria meningitidis</i>	Release July 24
<b><i>Pseudomonas aeruginosa</i></b>	Available
<i>Salmonella enterica</i>	Release June 24
<i>Serratia marcescens</i>	Release June 24
<i>Shigella boydii</i>	Release June 24
<i>Shigella dysenteriae</i>	Release June 24
<i>Shigella flexneri</i>	Release June 24
<i>Shigella sonnei</i>	Release June 24
<b><i>Staphylococcus aureus</i></b>	Available
<i>Staphylococcus epidermidis</i>	Release July 24
<i>Stenotrophomonas maltophilia</i>	Release July 24
<i>Streptococcus pneumoniae</i>	Release July 24
<i>Vibrio cholerae</i>	Release July 24
Species currently available	<b>14</b>
Species available by end of July	<b>42</b>
Species scheduled by end of 2024	<b>&gt;100</b>



# Genpax

[www.genpax.co](http://www.genpax.co)  
+44 20 3603 6869  
[research@genpax.co](mailto:research@genpax.co)