

## SRST2: Rapid genomic surveillance for public health and hospital microbiology labs

Michael Inouye<sup>1,2</sup>, Harriet Dashnow<sup>3,4</sup>, Lesley Raven<sup>1</sup>, Mark B. Schultz<sup>3</sup>, Bernard J. Pope<sup>4,5</sup>, Takehiro Tomita<sup>2,6</sup>, Justin Zobel<sup>5</sup>, Kathryn E. Holt<sup>3,#</sup>

<sup>1</sup> Medical Systems Biology, Department of Pathology, The University of Melbourne, Parkville, Victoria, Australia

<sup>2</sup> Department of Microbiology and Immunology, The University of Melbourne, Parkville, Victoria, Australia

<sup>3</sup> Department of Biochemistry and Molecular Biology, Bio21 Molecular Science and Biotechnology Institute, University of Melbourne, Parkville, Victoria 3010, Australia

<sup>4</sup> Victorian Life Sciences Computation Initiative, The University of Melbourne, 187 Grattan Street Carlton, Melbourne, Victoria, Australia

<sup>5</sup> Department of Computing and Information Systems, The University of Melbourne, Parkville, Victoria, Australia

<sup>6</sup> Microbiological Diagnostic Unit, The University of Melbourne, Parkville, Victoria, Australia

# To whom correspondence should be addressed: [kholt@unimelb.edu.au](mailto:kholt@unimelb.edu.au)

### Availability

SRST2 Python code is freely available (<https://github.com/katholt/srst2>) and utilises *bowtie2*<sup>1</sup> for read mapping and *SAMtools*<sup>2</sup> for alignment processing.

## Abstract

Rapid molecular typing of bacterial pathogens is critical for public health epidemiology, surveillance and infection control, yet routine use of whole genome sequencing (WGS) for these purposes poses significant challenges. Here we present SRST2, a tool for fast and accurate detection of genes, alleles and multi-locus sequence types from WGS data, which outperforms assembly-based methods. Using >900 genomes from common pathogens, we demonstrate SRST2's utility for rapid genome surveillance in public health laboratory and hospital infection control settings.

## Text

Rapid molecular typing of bacterial pathogens is critical for public health epidemiology, surveillance and infection control<sup>3,4</sup>. Whole genome sequencing (WGS) has revolutionised pathogen research and promises to revolutionise public health and hospital microbiology<sup>4-6</sup>, especially in the management of antimicrobial resistance. Yet, despite several demonstrative studies<sup>5,7,8</sup>, the lack of suitable analytic WGS tools for public health and diagnostic laboratories poses significant challenges<sup>3,9</sup>.

Two key goals for routine genomic surveillance of pathogens are: (i) to detect the presence of genes linked to clinically relevant phenotypes - including virulence genes, antimicrobial resistance genes or serotype determinants; and (ii) to classify isolates into clonal groups via multi-locus sequence typing (MLST<sup>10</sup>), detection of clone-specific or other epidemiological markers. Current methods rely on interrogating genome assemblies using BLAST or other search algorithms to identify genes or alleles<sup>10-13</sup>, yet these methods have limited sensitivity, efficiency and reproducibility. Production of quality assemblies can be highly variable due to dependence on data quality and to the requirements for data pre-processing and optimisation of parameters such as kmer length (**Methods**). This makes assembly-based analyses difficult to standardise and quality-control, features that are critical for routine use in hospital and public health settings.

Here we present SRST2, a reconstructed and substantially enhanced version of SRST<sup>14</sup>, for fast and accurate detection of genes, alleles, and/or MLST directly from WGS short reads using a mapping-based approach (**Methods, Supplementary Fig. 1**). To assess the accuracy of allele identification with SRST2, we analysed publicly available Illumina data from 543 bacterial genomes of five different species for which independent MLST data was available (**Supplementary Table 1**). With seven loci in each MLST scheme, this yielded 3,801 allele calls across 35 loci to assess call rate and false positive rate. The read sets represented a wide range of average read depths, with 90% in the range 12x - 130x and 50% between 20x - 60x (**Supplementary Table 1**). For each species, we used SRST2 to download the latest MLST database from [pubmlst.org](http://pubmlst.org) and subsequently ran SRST2 using default parameters. Median run time was 6 minutes per sample (interquartile range, 4-10 minutes) and increased linearly with number of reads (**Supplementary Fig. 2**). Efficiency can be easily improved or standardised, without data pre-processing, by instructing SRST2 to map the first N reads only.

SRST2 call rates and false positive rates increased with average read depth, stabilizing with depths  $\geq 15x$  (**Fig. 1a, Supplementary Fig. 3**). For comparison, we also assembled each read set using *Velvet*<sup>15</sup> and used nucleotide BLAST to identify MLST alleles (assembly+BLAST method; **Methods**). At read depths  $\geq 15x$ , SRST2 made significantly more allele calls than assembly+BLAST (call rates 99.9% vs 95.9%, respectively;  $p < 1 \times 10^{-15}$ ), with significantly greater accuracy (false positive rates 0.46% vs 0.90%;  $p = 0.05$ ). The heuristic information provided by SRST2 (that is, confident mismatches, insertions, deletions or truncations reported from read mapping) was a strong indicator of accuracy in the result: where an exact match was reported (98% of calls with depth  $\geq 15x$ ), the false positive rate was 0.2%; where an inexact match was reported, the false positive rate was 11.7%. For assembly, false positive rates were 0.2% for exact matches (95% of calls) and 76% for inexact matches. Hence, the key difference between the two methods was the ability of SRST2 to make correct calls where assembly+BLAST could not: for read depths  $\geq 15x$ , SRST2 made a call with the correct allele 99.4% of the time, compared to only 95% for assembly analysis ( $p < 1 \times 10^{-15}$  for difference in rates). At sequence type (ST) level, the difference was even greater: SRST2 achieved accurate ST assignment for 96% of isolates with average depth  $\geq 15x$ , whereas assembly+BLAST correctly identified only 76%.

To assess performance at low read depths ( $\leq 15x$ ), ten *S. aureus* read sets were subsampled to low depths (**Methods**). This confirmed that an average depth of only 10x was required for SRST2 to achieve  $>90\%$  call rate and  $<0.5\%$  false positives (**Fig. 1a, Supplementary Fig. 4**). MLST databases can be expected to grow indefinitely due to increasing diversity and broader sampling. However simulations (**Methods**) indicated that doubling the size of the *S. aureus* MLST database had no impact on SRST2 accuracy (**Fig. 1a, Supplementary Fig. 4**).

In addition to reliably distinguishing alleles of a given gene, SRST2 can also accurately determine the presence or absence of genes of interest, such as those encoding antimicrobial resistance or virulence. To evaluate this, we used 43 *E. faecium* genomes (**Supplementary Table 1**), previously screened for vancomycin susceptibility and presence of the VanB vancomycin resistance operon *vanABHSXY*<sup>16,17</sup>. Seventeen isolates were vancomycin resistant (VRE), and all were PCR positive for the *vanA-B* gene. These genomes were sequenced to  $\sim 1,000x$  depth and SRST2 correctly detected *vanA-B* in 17/17 VRE. In five vancomycin sensitive (VSE) isolates PCR negative for *vanA-B*, SRST2 detected *VanA-B* sequences at very low depths ( $<0.2\%$  of average depth), probably caused by minor but easily identifiable contamination during VRE-VSE multiplexed sequencing. SRST2 also confirmed the presence of the entire VanB operon, which is strongly predictive of the VRE phenotype. For comparison, assembly+BLAST identified full-length *vanA-B* sequences in just 7/17 VRE genomes, with multiple smaller hits spanning the full-length gene in five VRE and  $<50\%$  coverage of the gene identified in the remaining five VRE. To investigate the effect of sequencing depth on gene detection, we randomly selected five VRE and five VSE read sets for subsampling at  $<10x$  average read depth. *VanA-B* was only ever detected in confirmed VRE genomes, and sensitivity of detection with SRST2 reached 100% for read sets with  $\geq 5x$  average read depth (**Fig. 1c**).

To further explore the relative sensitivity of gene detection with SRST2, we screened all the read sets used for MLST validation (**Supplementary Table 1**) for antimicrobial resistance genes in the ARG-Annot database of acquired resistance genes<sup>13</sup> (**Methods**). SRST2's detection of whole genes was more sensitive than detection of whole or partial gene sequences by assembly+BLAST (**Supplementary Fig. 5**): 6.8% of genes detected at

$\geq 90\%$  coverage by SRST2 at depths  $\geq 15x$  were not found at  $\geq 90\%$  coverage in assemblies. For most of these genes, smaller fragments were detected by BLAST (**Supplementary Fig. 5**); however, SRST2 has the advantage of sensitive detection and confident allele-calling across the full length of genes, even at low depths (**Fig. 1c, Supplementary Fig. 5**).

To validate SRST2 in a public health laboratory setting, we analysed 231 clinical isolates of *Listeria monocytogenes* and compared MLST data obtained from gold-standard PCR and amplicon sequencing with those obtained from SRST2 or assembly+BLAST analysis of Illumina MiSeq data (**Fig. 1b**). Sequencing and analysis was performed by the Microbiological Diagnostic Unit Public Health Laboratory in Melbourne, Australia, the national reference laboratory for *L. monocytogenes*. For average read depths  $\geq 15x$ , SRST2 had a substantially higher call rate than assembly-based analysis (99.6% vs. 95.7%;  $p < 1 \times 10^{-12}$ ), with similar low false positive rates (0.7% vs. 0.6%;  $p = 0.9$ ). Hence, for samples with  $\geq 15x$  data, a total of 99% of all alleles were called correctly by SRST2, a significantly higher proportion than the 95% achieved by assembly+BLAST ( $p < 1 \times 10^{-12}$ ). At  $< 15x$  read depths, SRST2 also performed better than assembly-based analysis (87% vs 72% of alleles correctly called, respectively,  $p < 1 \times 10^{-3}$ ; **Fig. 1b**).

Further, SRST2 is already being assessed for routine MLST analysis of *Streptococcus pneumoniae* at Public Health England (Anthony Underwood, personal communication), and the open-source SRST2 code has been adapted by Public Health Ontario, Canada to perform specialist *emm* typing of Group A *Streptococcus*<sup>18</sup>.

In a hospital setting, the combination of MLST and gene detection can provide rapid and powerful insights for infection control without specialist bioinformatics knowledge. SRST2 analysis of 69 *K. pneumoniae* and 74 *E. coli* genomes from a UK hospital<sup>8</sup> revealed that each was dominated by a single ST with a high rate of the extended-spectrum beta-lactamase (ESBL) gene CTX-M-15 (*K. pneumoniae* ST490 comprising 25% of total, 71% of ESBL; *E. coli* ST131 comprising 40% of total, 77% of ESBL; **Supplementary Fig. 6**). Routine SRST2 surveillance of ESBL infections could be indicative of hospital outbreaks and used to identify which isolates should be investigated via transmission analysis.

Using the *E. faecium* genome data, collected as part of a 12-year hospital study of vancomycin resistance<sup>25</sup>, SRST2 took  $\sim 30$  minutes to generate the results in **Figure 2a-c**, showing (i) increasing vancomycin resistance over time; (ii) a shift in dominant ST during the same period; and importantly (iii) that this was not attributable to the introduction nor transmission of a new resistant clone, as the resistance rates were steady (approximately 50%) across all dominant STs. Similar conclusions typically require many days of labour and specialised assays in the diagnostic laboratory<sup>19</sup> and have been confirmed by detailed WGS analysis showing frequent acquisition of VanB transposons by diverse circulating strains<sup>16</sup>.

We next applied SRST2 to analyse data from real-world small-scale infection control investigations<sup>7</sup>. SRST2 took 5 minutes to generate results for suspected outbreaks of VRE and *E. cloacae* (**Fig. 3**), in which suspected outbreak isolates were readily distinguishable from epidemiologically unrelated isolates, consistent with WGS phylogenies and manual analysis of antimicrobial resistance markers<sup>7</sup>. SRST2 typing of 18 plasmid replicons<sup>20</sup> also indicated specific plasmid replicons (InCHI2, IncA/C) associated with two of the resistance profiles. The authors also reported use of a complex hybrid of assembly, mapping and manual inspection to determine carbapenem resistance mechanisms in five Gram-negative

bacteria isolated in close proximity<sup>7</sup>. SRST2 analysis of these five read sets identified the acquired beta-lactamases OXA-23 in AB223; IMP, SHV-12 and TEM-1 in EC1a; CTX-M-15 and TEM-1 in Eco216; CTX-M-15 and SHV-133 in KP652; and no acquired carbapenemase genes in EC302. These results are consistent with those reported from manual analysis<sup>7</sup>.

Here we have demonstrated the use of SRST2 for microbial genome surveillance in a variety of public health and hospital settings. In the face of rising threats of antimicrobial resistance and emerging virulence amongst bacterial pathogens, SRST2 represents a powerful tool for rapidly extracting clinically useful information from raw WGS data.

### **Availability**

SRST2 Python code is freely available (<https://github.com/katholt/srst2>) and utilises *bowtie2*<sup>1</sup> for read mapping and *SAMtools*<sup>2</sup> for alignment processing.

### **Acknowledgments**

This work was supported by the NHMRC of Australia (grant #1043830; fellowships #1061409 (KEH) and #1061435 (MI, co-funded with the Australian Heart Foundation)) and the Victorian Life Sciences Computation Initiative (VLSCI) (grant #VR0082).

### **Author Contributions**

Wrote code: MI, HD, BJP, KEH. Designed the study and algorithm: MI, BJP, JZ, KEH. Performed DNA extraction and sequencing: TT. Analyzed data: MI, HD, LR, MBS, KEH. All authors read and approved the final manuscript.

## Figure Legends

### Figure 1. Overall accuracy of SRST2 allele calling and gene detection.

(a) MLST analysis of public data from 5 species (N=543 genomes, 3801 loci, details Supplementary Table 1). Tests were grouped by read depth and accuracy rates (left y-axis, correct allele calls as a proportion of tests), calculated at each depth (x-axis, red slashes indicate scale change). Grey bars, number of tests at each depth (right y-axis); Lines, accuracy of allele calling. (b) MLST analysis of *Listeria monocytogenes* data (N=231 genomes, 1671 loci) conducted in a public health laboratory; colours and axes as in a. (c) Accuracy of *vanB* resistance gene detection for *E. faecium* read sets subsampled to low depth; y-axis shows proportion of correct (presence vs. absence) calls as a proportion of 100 tests at each depth; colours and axes as in a. A call of “present” implies detection of  $\geq 90\%$  of the length of the gene at  $\geq 90\%$  nucleotide identity.

### Figure 2. SRST2 analysis of *E. faecium* hospital data and hospital outbreak investigation.

Temporal distribution of isolates is shown in (a) coloured by vancomycin resistance as inferred from *vanA-B* detection with SRST2, and in (b) by coloured by sequence type inferred by SRST2. (c) Summary of all SRST2 results by sequence type (ST), in order from left to right: single linkage clustering of STs by number of shared alleles; MLST allele profiles; heatmap indicating the proportion of isolates that carries each resistance gene (scale as indicated), frequency of the ST (axis as indicated, coloured as in b).

### Figure 3. SRST2 analysis of hospital outbreak investigation.

(a) Isolate genetic profiles obtained from SRST2 analysis, indicating that case EF4 was distinct in both sequence type and resistance gene profile from the outbreak cases EF2 and EF3. Full WGS analysis showed a similar result<sup>7</sup>. (b) Isolate genetic profiles obtained from SRST2 analysis, including plasmid replicons detected (pink). The profiles indicate that case EC3 shared the same sequence type as the linked cases EC1 and EC2 (ST94), but lacked the IncA/C plasmid and had a distinct resistance gene profile. Full WGS analysis showed that EC1 and EC2 isolates were much closer to each other ( $\leq 22$  SNPs) than to EC3 ( $>150$  SNPs)<sup>7</sup>.

## Methods

### Approach and implementation

Given a read set and database of reference allele sequences, SRST2 is designed to perform two key tasks: (i) detect the presence of a gene or locus, and (ii) determine the precise or closest matching allele for that locus, amongst a set of possible reference allele sequences. The approach is illustrated in **Supplementary Figure 1**. A database of reference sequences must be provided in fasta format, in which the fasta headers indicate both the locus (so that alleles of the same locus can be compared) and a unique name for each allele. In the case of MLST data an additional database of ST profiles is provided as tab-delimited text, which assigns STs to unique combinations of alleles. Current MLST data (allele sequences and profile definitions), suitable for use with SRST2, can be downloaded from [pubmlst.org](http://pubmlst.org) automatically using the *getmlst.py* script supplied with SRST2. Other sequence databases can be easily formatted for use with SRST2 using the scripts supplied with the program. Any number of sequence databases can be analysed in a single run, allowing for simultaneous typing of MLST, resistance genes and virulence genes.

For each input database, reads are aligned using *bowtie2* v2.1.0 or above with the ‘--very-sensitive-local’ and ‘-a’ settings, and all alignments are reported to a file in SAM format. Mapping sensitivity can be fine-tuned by specifying to SRST2 any of the parameters available within the *bowtie2-align* command or a maximum number of mismatches per read (default 10 mismatches allowed). Flags in the resulting SAM file are modified so that each read is included in the pileup for every allele to which it is aligned. Pileups are generated using *SAMtools* v0.1.18 *mpileup* and parsed by SRST2 to determine percent coverage, divergence, and mismatches, and to calculate a score for each possible allele.

### Allele scoring

An overview of the scoring approach is given in **Supplementary Figure 1**. We begin with an alignment of reads from sample  $s$  to a reference sequence  $r$ . At each position  $i$  in the reference sequence  $r$  ( $r_i$ ), let  $s_i$  be the set of reads in sample  $s$  that align to  $r_i$ . Let  $a_i$  be the total number of reads in  $s_i$ , and let  $b_i$  be the number of reads in  $s_i$  in which the aligned base does not match the reference base at  $r_i$ . If sample  $s$  contains the precise sequence  $r$ , then the probability of a mismatched base at any position in an aligned read is equal to the per-base error rate of the sequencing technology  $e_i$ , which for Illumina is taken to be 0.01, although this can vary depending on what pre-processing steps are implemented<sup>21,22</sup>.

To quantify the evidence against the presence of the reference sequence  $r$  in  $s$ , we perform a Binomial test at each position  $r_i$ , to generate a 1-sided P-value  $P_i$  to assess the probability of observing  $a_i - b_i$  successes in  $a_i$  trials, with a probability of success of  $1 - e_i$ . Any change at position  $r_i$  - including a base substitution, an insertion of any size or a deleted base - is treated as a mismatch, incrementing  $b_i$  by 1. For large deletions that result in an absence of any aligned reads (including truncations of the end of the sequence),  $a_i = 0$  and no Binomial test is possible. In this case, the evidence for the deletion is provided by the reads which align adjacent to the deletion but do not align across the deletion. Hence we calculate the average number of reads aligned to the two bases preceding the deletion,  $d_i$ , and conduct the Binomial test with  $a_i = b_i = d_i$ .

We then utilize a non-parametric approach to score each allele by considering the set of all P-values calculated for reference sequence  $r$ . First, to minimise artefacts associated with

fluctuation in read depths, we (a) set  $P_i=1$  where  $b_i=0$ , and weight  $P_i$  by the relative read depth (i.e. weight of evidence) at position  $r_i$  compared to those of other positions in  $r$ :

$$\text{weighted } P_i (P_{i,w}) = P_i * (a_i / r_{\text{max depth}})$$

We then compare the sorted  $-\log_{10}(P_i)$  values versus those of the theoretical distribution of  $-\log_{10}(x_j/n)$  where  $n=\text{length}(r)$  and  $x_j = 1, 2, \dots, n$ , analogous to inspecting a quantile-quantile (QQ) plot (**Supplementary Fig. 1**). A linear model is fitted to the two probability distributions and the resulting slope is taken as the score for reference sequence  $r$ ,  $score_r$ . Here we leverage a common criticism of linear models to our advantage: the susceptibility to outliers at the tails of the distribution. In this case, outliers are typically SNPs or indels relative to the sequence  $r$  which, because they result in low P-values in the Binomial test and thus very high values of  $-\log_{10}(P)$ , are at the end of the observed distribution (**Supplementary Fig. 1**). Thus when a linear model is fitted, its slope increases with the number of well-supported SNPs and indels compared to the reference. As a result, among reference alleles of the same locus, the sequence  $r$  with the lowest  $score_r$  (flattest slope in the QQ plot) is the most likely match for  $s$ .

### Reporting outputs

SRST2 output tables report, for each sample  $s$  and each locus or gene cluster, the lowest scoring allele sequence  $r$ , the average read depth of  $s$  across  $r$  and indicators of any evidence against a precise match with  $r$  (including mismatches supported by >50% of aligned reads, or read depth falling below a cutoff). Only matches passing the user-set coverage and divergence cut-offs (by default, >90% coverage and <10% divergence) are reported. For MLST data, STs are calculated according to the MLST profiles database provided, based on the closest matching alleles at each locus.

Normally, an exact match between  $r$  and  $s$  would be assigned if (a)  $r$  has the lowest  $score_r$  amongst the set of alleles of the same locus or gene cluster, and (b) there are no SNPs or indels between  $r$  and  $s$ . If (a) holds but (b) does not, this is indicative of a novel allele and SRST2 will flag the result in output tables. In such cases, we recommend that users who are interested in defining novel alleles should inspect the raw sequence data (which may be assisted by the alignments, pileups and consensus fastq files generated by SRST2).

Optionally, SRST2 can report the full details of scoring  $s$  against all reference sequences  $r$ , to enable users to parse and interpret the results to suit specific needs. These include average depth of  $s$  across  $r$ , average depth across the first and last two bases of  $r$ , the number of positions in  $r$  in which the majority of aligned reads in  $s$  show a mismatch against  $r$  (with SNPs, insertion/deletions and truncations reported separately), the depth of bases neighbouring truncations and, for the position with the greatest proportion of mismatching reads, the total aligned reads, total mismatching, proportion mismatching, and Binomial p-value.

### Bacterial isolates and sequencing

A total of 231 *Listeria monocytogenes* isolates were analysed in this study, at the Microbiological Diagnostic Unit (MDU) Public Health Laboratory in Victoria, Australia. MDU is the national reference laboratory for *L. monocytogenes* and the isolates analysed include several from recent outbreaks as well as from the laboratory's reference collection. Cultures of *L. monocytogenes* isolated from food, environmental or clinical specimens were purified by two successive single colony selections after streaking onto horse blood agar



(HBA) incubated for 18-24 h at 37°C. Resultant bacterial growth on the surface of HBA medium was aseptically collected and resuspended in a cryotube (Nalgene) containing 1 mL of sterile glycerol storage broth (1.6% w/v Tryptone, Oxoid Pty Ltd, LP0042 containing 20% v/v glycerol) prior to storage at -70°C. Cultures were retrieved from storage as required and freshly grown (HBA, 18-24h at 37°C) in preparation for DNA extraction. DNA was extracted from each isolate using QIAmp DNA Mini Kit (Qiagen) and eluted in EB buffer (Qiagen) (Tris buffer, no EDTA).

DNA samples were subjected to traditional *L. monocytogenes* MLST analysis<sup>23</sup> (<http://www.pasteur.fr/recherche/genopole/PF8/mlst/Lmono.html>), with a minor modification to the annealing temperature for the *bglA* PCR (52°C not 45°C). The PCR products were purified with FastAP Thermosensitive Alkaline Phosphatase (Thermo Scientific) and Exonuclease I (Thermo Scientific). The purified PCR products were sequenced using BigDye Terminator v3 chemistry followed by capillary sequencing using a 3130xL Genetic Analyzer (Applied Biosystems). Trace analysis was conducted using BioNumerics version 6.6 with MLST Online plugin version 2.13 and Batch Sequence Assembly plugin version 1.34.

DNA was subjected to multiplex library preparation using Nextera XT followed by sequencing using an Illumina MiSeq. DNA was quantified by Qubit dsDNA HS Assay Kit (Invitrogen) and normalized to 0.2ng/μl. Total 1 ng of DNA was used for Nextera XT DNA Sample Preparation Kit (Illumina). Tagmentation of genomic DNA, PCR amplification with dual index primers, PCR clean-up using Agencourt AMPure XP (Beckman Coulter), DNA libraries normalization, library pooling and MiSeq sample loading were performed according to the manufacturer's instruction with minor modifications. For longer than 2×250 bp runs on the MiSeq, 25 μl of AMPure XP beads was added to each PCR-amplified product during the PCR purification step otherwise 30 μl of AMPure XP beads was added. For some samples, after PCR purification, DNA fragment size and library concentration was analysed by 2100 Bioanalyzer (Agilent Technologies) and Qubit dsDNA HS Assay Kit (Invitrogen). DNA libraries were normalized manually to 4 nM and libraries with unique indexes were pooled in equal volumes. Each resulting pooled library was denatured and diluted with 0.2N NaOH and pre-chilled HT1 (Illumina) to produce a 20 pM denatured library in 1 mM NaOH. Prior to the MiSeq run, the denatured library was further diluted with pre-chilled HT1 to approximately 12-13.5 pM. 600μl of library including 2% (v/v) 20 pM denatured PhiX library (Illumina) was loaded together with MiSeq reagent kit v3 (Illumina) according to the manufacturer's instructions.

### **Publicly available short read data used in this study**

Details of Illumina read sets used in this study are provided in **Supplementary Table 1**. Data tables specifying the expected STs of each read set, summarised from published papers, are available on the SRST2 website (<https://github.com/katholt/srst2>).

### **Subsampling of read sets**

To explore accuracy at low read depths, ten genomes each of *S. aureus* and *E. faecium* were selected for random subsampling of reads to simulate genomes sequenced to low read depth. To do this, we used the mean read depth across MLST loci to calculate the sampling fraction required to achieve approximately 1x, 2x, ... 10x mean read depth. We randomly sampled reads from the forward reads file at the required sampling fraction, and extracted the corresponding reverse reads, using Perl scripts. Ten random samples were generated from each read set at each depth level, generating a total of 1,000 read sets for each species.

### Sequence databases used in this study

MLST databases for *S. aureus*, *S. pneumoniae*, *S. enterica*, *E. coli*, *E. faecium*, *L. monocytogenes* and *E. cloacae* were downloaded from [pubmlst.org](http://pubmlst.org) using the *getmlst.py* script included with SRST2 (June 2014).

Antimicrobial resistance gene detection was performed using the ARG-Annot database of acquired resistance genes<sup>13</sup>. Allele sequences (DNA) were downloaded in fasta format from <http://www.mediterranee-infection.com/article.php?laref=282&titer=arg-annot> (May, 2014). Sequences were clustered into gene groups with  $\geq 80\%$  identity using CD-hit<sup>24</sup> and the headers formatted for use with SRST2 using the scripts provided (*cdhit\_to\_csv.py*, *csv\_to\_gene\_db.py*). A copy of the formatted sequence database used in this study is available on the SRST2 website (<https://github.com/katholt/srst2>).

Representative sequences for 18 plasmid replicons were extracted from GenBank using the accessions and primer sequences specified by Carattoli *et al*<sup>20</sup>. A copy of the formatted sequence database used in this study is available on the SRST2 website (<https://github.com/katholt/srst2>).

### Simulation of expanded *S. aureus* MLST database

As more genomes are sequenced and as bacteria continue to evolve, novel alleles will continue to be discovered and thus the size of allele databases will increase. To explore the impact of database size on accuracy of allele detection with SRST2, we simulated expansion of the current *S. aureus* MLST database from 2,161 alleles (mean 309 per locus) to 5,578 alleles (mean 797 per locus). The additional  $\sim 500$  alleles per locus were generated using *netrecodon* v6.0.0<sup>25</sup>. Sequences derived from the true MLST database were used to seed the simulation at each locus as follows. Existing alleles were translation-aligned between start (alignment start) and stop (alignment end) codons, those containing a frameshift or stop codon were removed, and the modal consensus sequence was exported. The best-fit DNA substitution model of each true alignment was determined using the AIC in *MrModeltest* v2.3, as implemented in *PAUP\** v4.0b. In *netrecodon*, the modal sequences were forward evolved under the coalescent, using the parameters of the best-fitting model for each locus, mutation rate  $1E-7$  and recombination rate  $1E-7/15$  (based on reported  $r/m$  of  $1/15$ <sup>26</sup>). A total of 100 independent replicates of forward evolution were performed per locus, retaining 2,000 sequences per replicate ( $N = 200,000$  simulated sequences per locus). The first 500 unique simulated sequences at each locus were added to the MLST database, and duplicate sequences were removed.

### Analysis runs and time calculations

All SRST2, assembly and BLAST analysis was run on a Linux cluster (iDataplex x86 system, “Barcoo” cluster at VLSCI – <http://vlsci.org.au>). SRST2 was run with default parameters. Details of *Velvet* assembly and BLAST analysis are given below. Run times were calculated from time stamps extracted from log files for SRST2 and *Velvet Optimiser* assembly runs.

### Assembly-based analysis

Assemblies were generated using the de novo assembler *Velvet* v1.2.10<sup>15</sup>, with optimal kmer choice for each readset refined through iterative calls to *Velvet Optimiser* v2.2.5 (<http://bioinformatics.net.au/>). Briefly, each read set was assembled using a call to *Velvet Optimiser* with kmers from 29 up to 89, in steps of 12. The optimal kmer,  $k_1$ , was

extracted and a second call to *VelvetOptimiser* was made using kmers from  $k_1-12$  up to  $k_1+12$ , in steps of 4. A final call to *VelvetOptimiser* was run using kmers from  $k_2-4$  up to  $k_2+4$ , in steps of 2. The final assembly was that output from the third and final call to *VelvetOptimiser*.

For MLST analysis from assemblies, a nucleotide BLAST+ (v2.2.25) search was performed for each locus and each contig set. In this BLAST search, the contig set was used to query the database containing all known allele sequences for a given locus, and the top BLAST hit was reported. If this hit had  $\geq 90\%$  nucleotide identity across  $\geq 90\%$  of the length of the reference allele sequence, an allele call was recorded. If the hit was an exact match to a known allele (i.e. 100% nucleotide identity across 100% of the length of the allele sequence), this was considered a precise allele call. The Python code used is available within the SRST2 distribution. For gene detection analysis from assemblies, a nucleotide BLAST search was performed in which the set of reference sequences (sequence database, i.e. antimicrobial resistance gene database) was used to query the database of all contigs for that assembly.

### Statistical analysis

All statistical analysis and data plotting was performed in *R*. Allele calling performance of SRST2 and assembly+BLAST was assessed via three metrics: (i) call rate = total number of allele calls made, for SRST2 this was a call with  $\geq 90\%$  coverage and no uncertainty recorded (i.e. with  $\geq 2x$  read depth at both ends and also neighbouring any truncations or deleted bases), for BLAST this was a call with  $\geq 90\%$  coverage and  $\geq 90\%$  nucleotide identity; (ii) false positive rate = total number of correct allele calls as a proportion of all calls; (iii) proportion of all tests resulting in a call with a correct allele, equal to (call rate) \* [1 - (false positive rate)]. As these metrics are proportions, the significance of differences in performance metrics was calculated using a two-sided test for equality of proportions (*prop.test* function in *R*). Resistance gene detection was assessed using a cut-off of  $\geq 90\%$  coverage and  $\geq 90\%$  identity to define the presence of a gene.

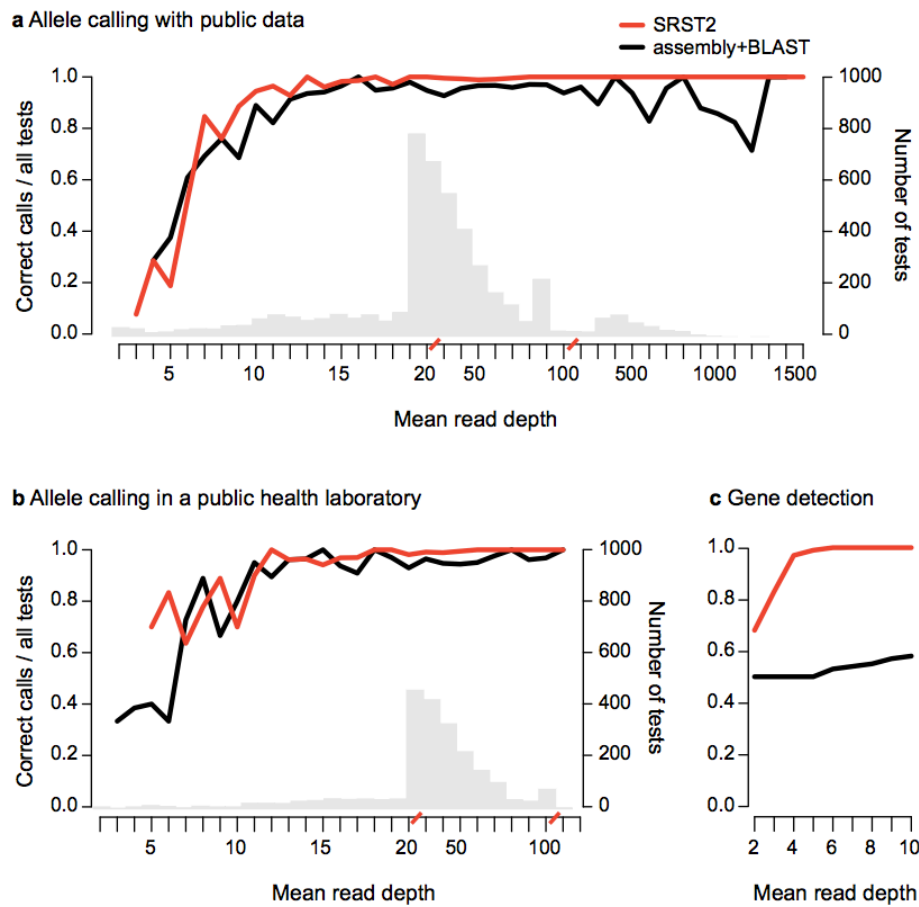
## References

1. Langmead, B. & Salzberg, S.L. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**, 357-9 (2012).
2. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-9 (2009).
3. Sabat, A.J. *et al.* Overview of molecular typing methods for outbreak detection and epidemiological surveillance. *Euro Surveill* **18**, 20380 (2013).
4. Bertelli, C. & Greub, G. Rapid bacterial genome sequencing: methods and applications in clinical microbiology. *Clin Microbiol Infect* **19**, 803-13 (2013).
5. Didelot, X., Bowden, R., Wilson, D.J., Peto, T.E. & Crook, D.W. Transforming clinical microbiology with bacterial genome sequencing. *Nat Rev Genet* **13**, 601-12 (2012).
6. Koser, C.U. *et al.* Routine use of microbial whole genome sequencing in diagnostic and public health microbiology. *PLoS Pathog* **8**, e1002824 (2012).
7. Reuter, S. *et al.* Rapid bacterial whole-genome sequencing to enhance diagnostic and public health microbiology. *JAMA Intern Med* **173**, 1397-404 (2013).
8. Stoesser, N. *et al.* Predicting antimicrobial susceptibilities for Escherichia coli and Klebsiella pneumoniae isolates using whole genomic sequence data. *J Antimicrob Chemother* **68**, 2234-44 (2013).
9. D'Auria, G., Schneider, M.V. & Moya, A. Live Genomics for Pathogen Monitoring in Public Health. *Pathogens* **3**, 93-108 (2014).
10. Jolley, K.A. & Maiden, M.C. Automated extraction of typing information for bacterial pathogens from whole genome sequence data: Neisseria meningitidis as an exemplar. *Euro Surveill* **18**, 20379 (2013).
11. Zankari, E. *et al.* Identification of acquired antimicrobial resistance genes. *The Journal of antimicrobial chemotherapy* **67**, 2640-4 (2012).
12. Larsen, M.V. *et al.* Multilocus Sequence Typing of Total Genome Sequenced Bacteria. *J Clin Microbiol* (2012).
13. Gupta, S.K. *et al.* ARG-ANNOT, a new bioinformatic tool to discover antibiotic resistance genes in bacterial genomes. *Antimicrob Agents Chemother* **58**, 212-20 (2014).
14. Inouye, M., Conway, T.C., Zobel, J. & Holt, K.E. Short read sequence typing (SRST): multi-locus sequence types from short reads. *BMC Genomics* **13**, 338 (2012).
15. Zerbino, D.R. & Birney, E. Velvet: algorithms for de novo short read assembly using de Bruijn graphs. *Genome Res* **18**, 821-9 (2008).
16. Howden, B.P. *et al.* Genomic insights to control the emergence of vancomycin-resistant enterococci. *MBio* **4**(2013).
17. Stinear, T.P., Olden, D.C., Johnson, P.D., Davies, J.K. & Grayson, M.L. Enterococcal vanB resistance locus in anaerobic bacteria in human faeces. *Lancet* **357**, 855-6 (2001).
18. Athey, T.B. *et al.* Deriving Group A Streptococcus Typing Information from Short-Read Whole Genome Sequencing Data. *J Clin Microbiol* (2014).
19. Johnson, P.D. *et al.* A sustained hospital outbreak of vancomycin-resistant Enterococcus faecium bacteremia due to emergence of vanB E. faecium sequence type 203. *J Infect Dis* **202**, 1278-86 (2010).

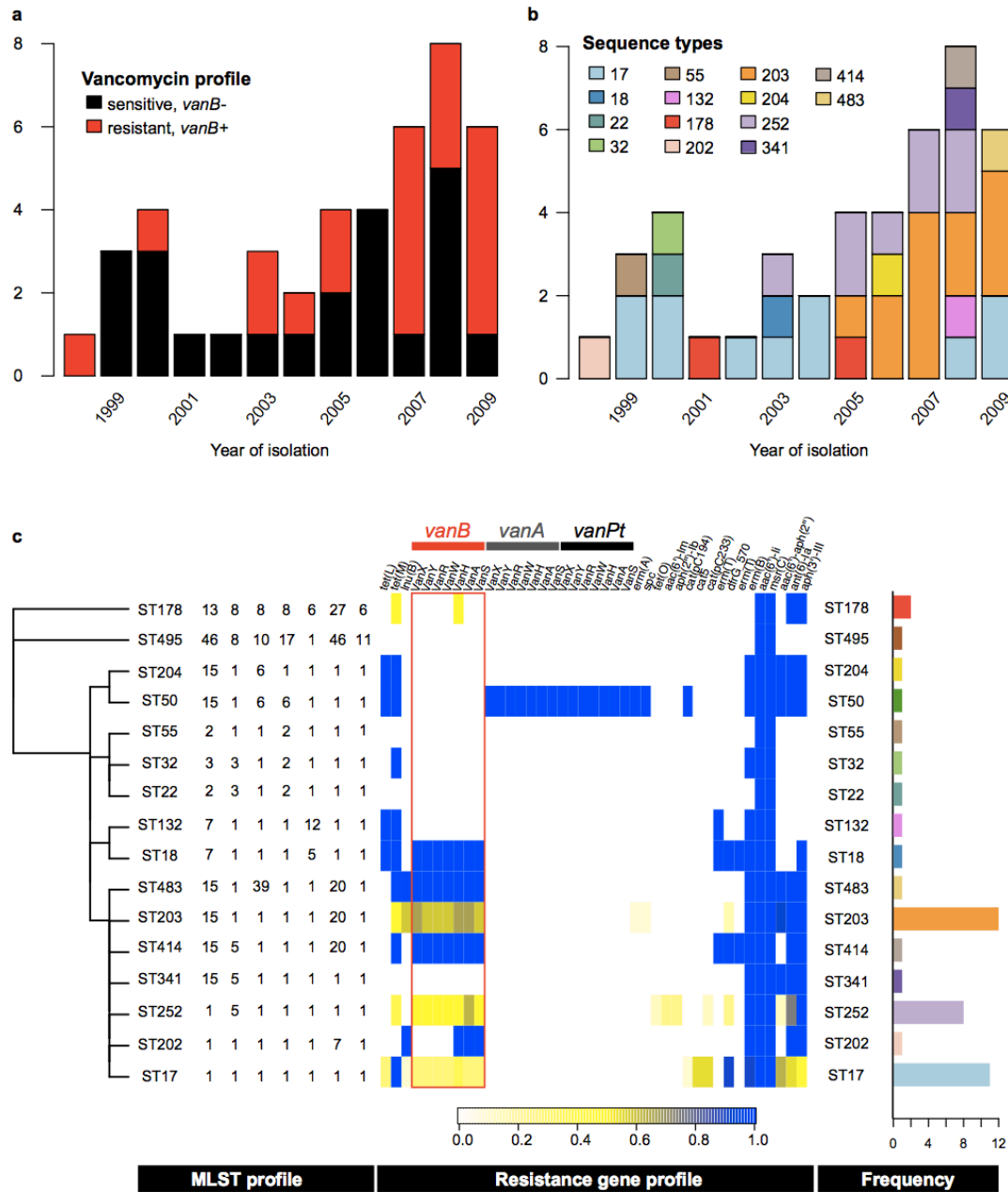
20. Carattoli, A. *et al.* Identification of plasmids by PCR-based replicon typing. *J Microbiol Methods* **63**, 219-28 (2005).
21. Loman, N.J. *et al.* High-throughput bacterial genome sequencing: an embarrassment of choice, a world of opportunity. *Nat Rev Microbiol* **10**, 599-606 (2012).
22. Meacham, F. *et al.* Identification and correction of systematic error in high-throughput sequence data. *BMC Bioinformatics* **12**, 451 (2011).
23. Ragon, M. *et al.* A new perspective on *Listeria monocytogenes* evolution. *PLoS Pathog* **4**, e1000146 (2008).
24. Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics* **22**, 1658-9 (2006).
25. Arenas, M. & Posada, D. Coalescent simulation of intracodon recombination. *Genetics* **184**, 429-37 (2010).
26. Feil, E.J. *et al.* How clonal is *Staphylococcus aureus*? *J Bacteriol* **185**, 3307-16 (2003).

## Figures

**Figure 1: Overall accuracy of SRST2 allele calling and gene detection**



**Figure 2: SRST2 analysis of *E. faecium* hospital data and hospital outbreak investigation**

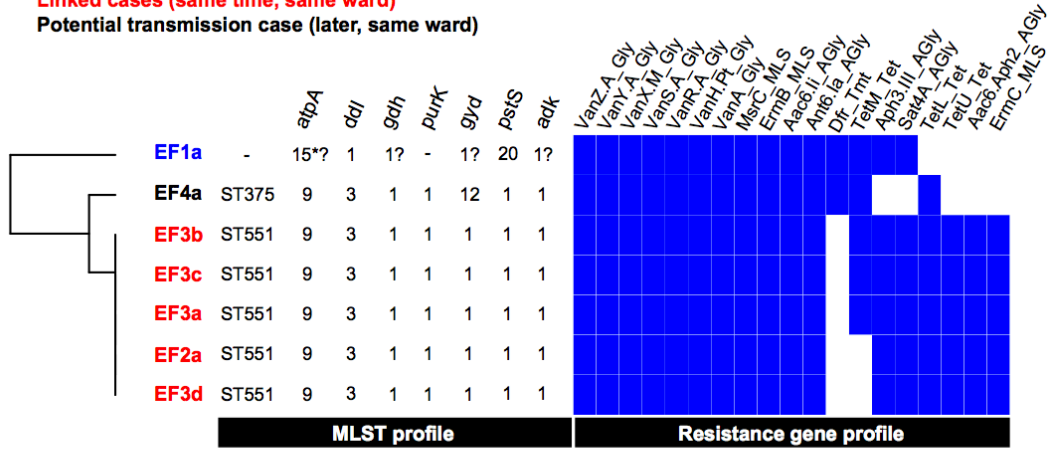


### Figure 3: SRST2 analysis of hospital outbreak investigation

a **Unrelated case (earlier time, same ward)**

**Linked cases (same time, same ward)**

**Potential transmission case (later, same ward)**



b **Unrelated case (same time, different wards)**

**Linked case (same time, same ward)**

**Potential transmission case? (later, same ward)**

