# Refined analyses suggest recombination is a minor source of genomic diversity in *Pseudomonas aeruginosa* chronic cystic fibrosis infections.

## METHODS

### Short Read Preparation

Preprocessing of short read data was implemented in the PrepareReads module of BAGA. Large FASTQ read files were subsampled to provide a 80x read coverage for mapping to a 6.6 Mb genome. Reads were trimmed using position-specific quality scores with Sickle ver. 1.33 (Joshi & Fass 2011) . Adaptor sequences from library preparation were removed using Cutadapt ver. 1.9.dev0 (Martin 2011) .

### Short Read Alignment to reference genome

Short read alignment was implemented in the AlignReads module of BAGA. Paired end short reads were aligned to the LESB58 reference genome sequences with the MEM algorithm of BWA  (Li & Durbin 2009) (github.com repository revision 9812521) estimating insert size from the data for designation of "proper pair" reads. SAMTools (github.com repository revision b0525f3) (Li et al. 2009) was used for conversion of SAM to BAM files and BAM file sorting. Duplicate reads were removed from alignments using Picard version1.135. The IndelRealigner module of the Genome Analysis Toolkit (GATK) (McKenna et al. 2010; DePristo et al. 2011) was used to simultaneously re-align reads that were likely to be near insertions or deletions.

### Reference genome repetitive region variant filter

The algorithm for identifying reference genome repeat regions in which variants were deemed unreliable and omitted was implemented in the Repeats module of BAGA. First, each open reading frame nucleotide sequence in the published annotation was aligned back to the full LESB58 chromosome sequence using BWA MEM lowering alignment stringency by setting mismatch penalty to 2 and gap open penalty to 3. These initial alignments using BWA are used to seed an alignment procedure that uses optimal global (Needleman-Wunsch) alignment for better accuracy. The lower alignment stringency produces alignments more divergent than the target 98% nucleotide identity. After omitting alignments between the same sequences, the resulting collection of ORFs within

37  repeat regions were organised into contiguous blocks of homologous sequence. Each pair of
38  homologous sequences, which may now contain more than one pair of ORFs, were aligned using the
39  Needleman-Wunsch optimal global alignment algorithm implemented in the seq-align package
40  (github.com/noporpoise/seq-align revision 0589383). The percent identity between the two aligned
41  regions is calculated within each 100 bp window at 20 bp increments along a pairwise alignment and
42  regions with more than 98% nucleotide sequence identity were deemed similar enough to yield
43  ambiguous read alignments. Regions greater than the mean insert size are retained for filtering
44  because smaller regions are expected to be resolved by the relative positions of paired-end reads.

## 45 Genome rearrangements variant filter

46  This algorithm was implemented in the Structure module of BAGA and is applied on a per sample
47  basis. The threshold for the ratio of reads not in a proper pair, to those in a proper pair, which when
48  exceeded defines regions in which variant calls should be omitted, was set to 0.15. At this threshold,
49  false positive variants in the simulations where large deletions were addedwere successfully filtered
50  out. This ratio, collected at each position in the reference genome, is "smoothed" by averaging over a
51  500 bp moving window. Proper pair status was determined by the BWA MEM aligner and was
52  acquired from BAM files using PySAM version 0.8.3. Each region for variant omission was extended in
53  each direction if a moving window equal to the mean insert size had at least 15% without any aligned
54  reads. Plotting is implemented in the Structure module.
55

## 56 Variant calling

57  Variant calling was implemented in the CallVariants module of BAGA which follows the "Best Practices
58  workflow" provided by the developers of GATK (Van der Auwera et al. 2013)
59  (https://www.broadinstitute.org/gatk/guide/best-practices accessed 15/04/2015). SNPs and indels
60  were called simultaneously via local re-assembly of haplotypes using the HaplotypeCaller module of
61  the GATK setting –sample_ploidy to 1,  –heterozygosity to 0.0001 and –indel_heterozygosity to
62  0.00001 to approximate the 100-1400 SNPs in the 6.6Mb genomes reported by Williams et al. (2015)
63  and Darch et al. (2015) and the ten-fold few indels reported by Williams et al. (2015). The two
64  datasets were called in separate joint analyses. The standard GATK hard filters for SNPs and indels
65  were applied. Base quality scores in BAM files were recalibrated.

## 66 Population structure analysis

67  The following procedures were implemented in the ComparativeAnalysis module of the BAGA. A
68  multiple sequence alignment including the reference genome and all isolates was generated from the
69  called SNPs with all filters applied. Invariant positions were retained and columns with missing data
70  were included except if missing in all samples. For example, all 22 Nottingham samples are missing
71  LESB58 Genomic Island 5 which is 50 kb long. A pair of polymorphisms 50 bases either side of the
72  insertion site of Genomic Island 5 should be considered 100 bases apart in the 22 sampled genomes,
73  not 50,100 bases apart, as it is only in the reference genome. Maximum likelihood phylogenetic
74  reconstruction was performed using PhyML (github.com/stephaneguindon/phyml revision e71c553)
75  (Guindon et al. 2010) using the default settings. Recombination was inferred using ClonalFrameML

76  (github.com/xavierdidelot/ClonalFrameML revision 2d793a3) (Didelot & Wilson 2015) using the
77  PhyML phylogeny and -ignore_incomplete_sites true. Plotting including mapping the alignment
78  positions reported by ClonalFrameML back to the reference genome is implemented in the
79  ComparativeAnalysis module of BAGA. Phylogeny manipulations were performed using DendroPy
80  4.0.3 (Sukumaran & Holder 2010). Further detection of recombination events, not implemented in
81  BAGA, was performed using BRAT NextGen version 4/18/2011(Marttinen et al. 2012) with 100
82  iterations to achieve negligible changes in model parameter estimates over the last 50 iterations. A
83  threshold of 5% applied to 100 permutations was used for significance testing of each event.

84 ## In silico read generation

85  Ten sets of paired reads were generated from the LESB58 reference sequence using GemSim version
86  1.6 (McElroy et al. 2012) and modelling the error profile of the Illumina GA IIx with Illumina
87  Sequencing Kit v4 chemistry. Read length was 100 bp, with insert size 350 bp and standard deviation
88  20. To provide 60x coverage depth, 1,980,527 read pairs were generated. Each set of reads was
89  generated from a reference sequence on which 50 SNPs (30 shared with ~20 additionally selected at
90  random to approximate a tree-like distribution) and 15 insertions and deletions (five shared, 10
91  genome specific) had been applied *in silico*. These simulated reads can be reproduced using the
92  SimulateReads module of BAGA and the scripts provided in the figshare
93  repository(http://dx.doi.org/10.6084/m9.figshare.2056365). The read data was then analysed with
94  the BAGA pipeline described above.
95

96 ## Region-specific de novo assembly of reads for SNP validation

97  For each isolate read set, at each SNP position, reads mapping to 5000 bp each side were extracted
98  from BAM files using pySAM. These were combined with the unmapped and poorly mapped reads for
99  de novo assembly using SPAdes. By reducing the total reads per assembly we aimed to decrease
100 assembly complexity and increase accuracy. A set of contigs were also assembled from just the
101 unmapped and poorly mapped reads and subtracted from the contigs of each SNP-associated
102 assembly to leave those contigs relevant to that region. These remaining contigs were pairwise
103 aligned to the region using the optimal global Needleman-Wunsch algorithm implemented in
104 Seqalign. The pairwise alignment was then check for the presence or absence of the SNP.
105

106 ## RESULTS AND DISCUSSION
107

108 ## Impact of reference genome repeats on variant calling
109

110 In the present analysis, two novel filters were applied to mitigate false positive errors. For each filter,
111 variants were omitted if called in regions where assumptions of the method were likely to be violated:
112 either within long repeats in the reference sequence, or near break-points at rearrangements
113 between a sample and the reference. Calling variants by mapping whole-genome shotgun reads to a
114 reference genome assumes each gene sequence in the sampled individual must have a unique

115 equivalent in the reference. This 1-to-1 relationship means that each read will align unambiguously to
116 the position in the reference genome that is positionally orthologous to the read origin in the sampled
117 genome. If a read is aligned to the wrong (non-orthologous) part of a reference genome, false positive
118 variants calls are likely.
119
120 The first novel variant filter in this analysis omitted variants in long repeats in the reference sequence.
121 Sequence repeats can be caused by duplications (paralogy) or horizontally acquired DNA incorporated
122 into a chromosome whilst also being homologous to another part of the chromosome. The Burrow-
123 Wheeler Aligner used in this analysis reports an alignment quality of zero where a read maps perfectly
124 to two parts of a genome. In these cases, the risk of false positive variants is easily mitigated by
125 excluding reads with an alignment quality of zero. However, these zero alignment qualities will not be
126 reported for slightly divergent repeats, which would be the case when one or more substitutions have
127 occurred at one repeat unit but not another.
128
129 The algorithm we developed to identify repeat regions was implemented in the Repeats module of
130 the Bacterial and Archaeal Genome Analyser (BAGA). Of the 6.6 Mb LESB58 reference genome,
131 75.1 kb of chromosome in 35 regions were deemed repetitive. These regions required at least 98%
132 aligned nucleotide identity and a length greater than the mean sequencing library insert size because
133 smaller regions are expected to be resolved by the relative positions of paired-end reads. Only
134 contiguous repeat regions, tandem or otherwise, longer than the insert size used in library
135 preparation were included in the filter because shorter regions should be resolvable as unique. Of the
136 42.7 kb in previously described long duplications in LESB58 (Winstanley et al. 2009), 23.0 kb were
137 included in the high identity regions reported here. Among the Nottingham data, eight SNPs within
138 these regions were omitted (Table S2). In the Liverpool data, nine SNPs were omitted. No indels were
139 found in these regions.
140

141

## Impact of genome sequence rearrangements on variant calling

143 The second novel variant filter used in this analysis omitted variants in regions around break-points of
144 chromosomal rearrangements. The assumed 1-to-1 orthology required for variant calls from aligned
145 reads can be violated at the regions bordering rearrangements between a sample and the reference
146 genome. These structural changes in bacteria include large deletions, inversions, duplications and
147 integrations, the last typically caused by acquisition of prophage or genomic islands. Some phage are
148 known to carry a fragments of tRNA genes with homology to those in the host organism and
149 integration involves replacement of parts the host version of a tRNA gene with the phage version
150 (Campbell 1992) . Thus, regions flanking the resulting prophage may have sequence divergence
151 caused by a more complex history than the simple orthology assumed by the read mapping method
152 with an increased risk of artefactual, false-positive variants.
153

154  The algorithm we developed for omitting putatively unreliable variants near regions of structural
155  rearrangements was implemented in the Structure module of the BAGA software package. The
156  algorithm exploits the fact that sequence reads in this analysis were "paired-end": from either end of
157  many template DNA fragments of similar lengths. The Burrows-Wheeler Aligner used in this analysis
158  assigns a read "proper pair" status if the aligned distance and orientation to its paired read is close to
159  that expected given the mean template DNA fragment length. Paired-end reads on either side of a
160  rearrangement break-point typically align to the reference at distances far greater than the fragment
161  lengths or in different orientations. For example, in many samples, no reads mapped to the LESB58
162  chromosome region annotated as "LES prophage 5" indicating absence of the prophage in those
163  samples. Paired reads would align kilobases apart on either side of the prophage if they originated
164  from a DNA fragment spanning the orthologous site of prophage integration, in a sample lacking the
165  prophage. These reads would not be assigned "proper pair" status being so much further apart than
166  the mean fragment size.
167
168  Fig. 3 illustrates how our algorithm determines regions near genome rearrangements in which variant
169  calls may be unreliable. In this example the rearrangement is a duplication, but other common
170  rearrangements include the loss or gain of a prophage or other genomic island. These regions are
171  defined by a deviation in the ratio of proper-pair-assigned reads to non-proper-pair-assigned reads:
172  the more non-proper-pair-assigned reads, the less reliable the region. The initially selected region is
173  extended while regions of the reference have zero aligned read coverage and variants called within
174  the combined region are omitted. In the Nottingham data, 33 SNPs and 10 indels within these regions
175  were omitted (Table S2). In the Liverpool data, 39 SNPs and eight indels were omitted. This analytic
176  approach, by accounting for potential false positive variant calls caused by missing regions of
177  chromosome, is of value because turnover of prophage and genomic islands among microbial
178  genomes is often rapid(Wilmes et al. 2009). In many cases, including the analyses of the Liverpool and
179  Nottingham datasets, guaranteeing a similar genomic composition to permit a contiguous alignment
180  between reference and sampled genomes is difficult. Prophage and genomic islands and are
181  abundant throughout Bacteria and Archaea(Zhou et al. 2011)so few read mapping analyses will be
182  unaffected by these challenges.
183

## TABLES

Table S1: Small insertions and deletions predicted to cause frame-shift mutations among the Nottingham isolates

| Chromosome | Position (bp) | Reference | Variant | Gene ID | Gene name | Frequency | Frequency (filtered) |
|---|---|---|---|---|---|---|---|
| NC_011770.1 | 14914 | T | TGC | PLES_RS00075 | | 10 | 10 |
| NC_011770.1 | 335247 | C | CG | PLES_RS01530 | | 2 | 2 |
| NC_011770.1 | 467972 | GATCCTCCTCGT | G | PLES_RS02170 | mexB | 22 | 22 |
| NC_011770.1 | 1037925 | A | AC | PLES_RS05000 | mpl | 3 | 3 |
| NC_011770.1 | 1267074 | CCCGCTGGAGCT | C | PLES_RS06070 | | 3 | 3 |
| NC_011770.1 | 2532081 | G | GT | PLES_RS12255 | | 1 | 0 |
| NC_011770.1 | 2639073 | AG | A | PLES_RS12750 | | 1 | 1 |
| NC_011770.1 | 3280614 | A | AAGCGCGCGACC CGGAAGCCGTTG CCCAGCGGGGCG CCCTGT | PLES_RS15460 | | 1 | 0 |
| NC_011770.1 | 3378367 | AG | A | PLES_RS15875 | pslJ | 21 | 21 |
| NC_011770.1 | 3388596 | GTACTTGCCGCT GTCCAGGTAGGA CTGGGCGGTGGC CTGGTCGGGTTT CT | G | PLES_RS15915 | pslB | 1 | 0 |
| NC_011770.1 | 3484676 | CGCCCGAGCAA | C | PLES_RS16380 | | 22 | 22 |
| NC_011770.1 | 3898444 | C | CACGATTCGTTG TCAAAAATAGCC AAGGACCCGGAC ACACGCCTGATG CG | PLES_RS18160 | mltD | 1 | 1 |
| NC_011770.1 | 4044433 | C | CGG | PLES_RS18915 | stk1 | 1 | 1 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| NC_011770.1 | 4931526 | G | GGGGATGTCGAC ATGCAA | PLES_RS23105 | | 2 | 0 |
| NC_011770.1 | 5398560 | TCG | T | PLES_RS25200 | ampD | 22 | 22 |
| NC_011770.1 | 5903286 | GCT | G | PLES_RS27535 | | 22 | 22 |
| NC_011770.1 | 6060546 | C | CAT | PLES_RS28165 | | 22 | 22 |
| NC_011770.1 | 6557920 | GGGGCGGCAGCA GTGGCGTTCGGC GCAGT | G | PLES_RS30445 | | 22 | 22 |

Table S2: Total polymorphisms called for the Liverpool and Nottingham datasets as cumulative filters are applied, divided by type

| | Liverpool | | Nottingham | | All data | |
|---|---|---|---|---|---|---|
| Cumulative filters applied | SNPs[*] | Indels[†] | SNPs | Indels | SNPs | Indels |
| None | 318 | 68 | 170 | 47 | 440 | 89 |
| GATK[‡] standard 'hard' filter | 318 | 68 | 170 | 47 | 440 | 89 |
| BAGA[§] reference genome repeats | 309 | 68 | 162 | 47 | 424 | 89 |
| BAGA genome rearrangements | 270 | 60 | 129 | 37 | 364 | 78 |

*Single nucleotide polymorphisms; †Small insertion or deletion polymorphisms; ‡Genome Analysis Toolkit; §Bacterial and Archaeal Genome Analyser

Table S3: Polymorphisms called as fixed in samples and differing with reference for the Liverpool and Nottingham datasets as cumulative filters are applied, divided by type

| | Liverpool | | Nottingham | | All data | |
|---|---|---|---|---|---|---|
| Cumulative filters applied | SNPs[*] | Indels[†] | SNPs | Indels | SNPs | Indels |
| None | 43 | 5 | 101 | 17 | 113 | 21 |
| GATK[‡] standard 'hard' filter | 43 | 5 | 101 | 17 | 113 | 21 |
| BAGA[§] reference genome repeats | 43 | 5 | 99 | 17 | 111 | 21 |

| | | | | | | |
|---|---|---|---|---|---|---|
| BAGA genome rearrangements | 43 | 5 | 92 | 15 | 104 | 19 |

194 *Single nucleotide polymorphisms; †Small insertion or deletion polymorphisms; ‡Genome Analysis Toolkit; §Bacterial and Archaeal

195 Genome Analyser

196

197 Table S4:  Polymorphisms called among the Liverpool and Nottingham datasets as cumulative filters are applied, divided by type

| | Liverpool | | Nottingham | | All data | |
|---|---|---|---|---|---|---|
| Cumulative filters applied | SNPs* | Indels† | SNPs | Indels | SNPs | Indels |
| None | 275 | 63 | 69 | 30 | 337 | 71 |
| GATK‡ standard 'hard' filter | 275 | 63 | 69 | 30 | 337 | 71 |
| BAGA§ reference genome repeats | 266 | 63 | 63 | 30 | 322 | 71 |
| BAGA genome rearrangements | 227 | 55 | 37 | 22 | 264 | 60 |

198 *Single nucleotide polymorphisms; †Small insertion or deletion polymorphisms; ‡Genome Analysis Toolkit; §Bacterial and Archaeal
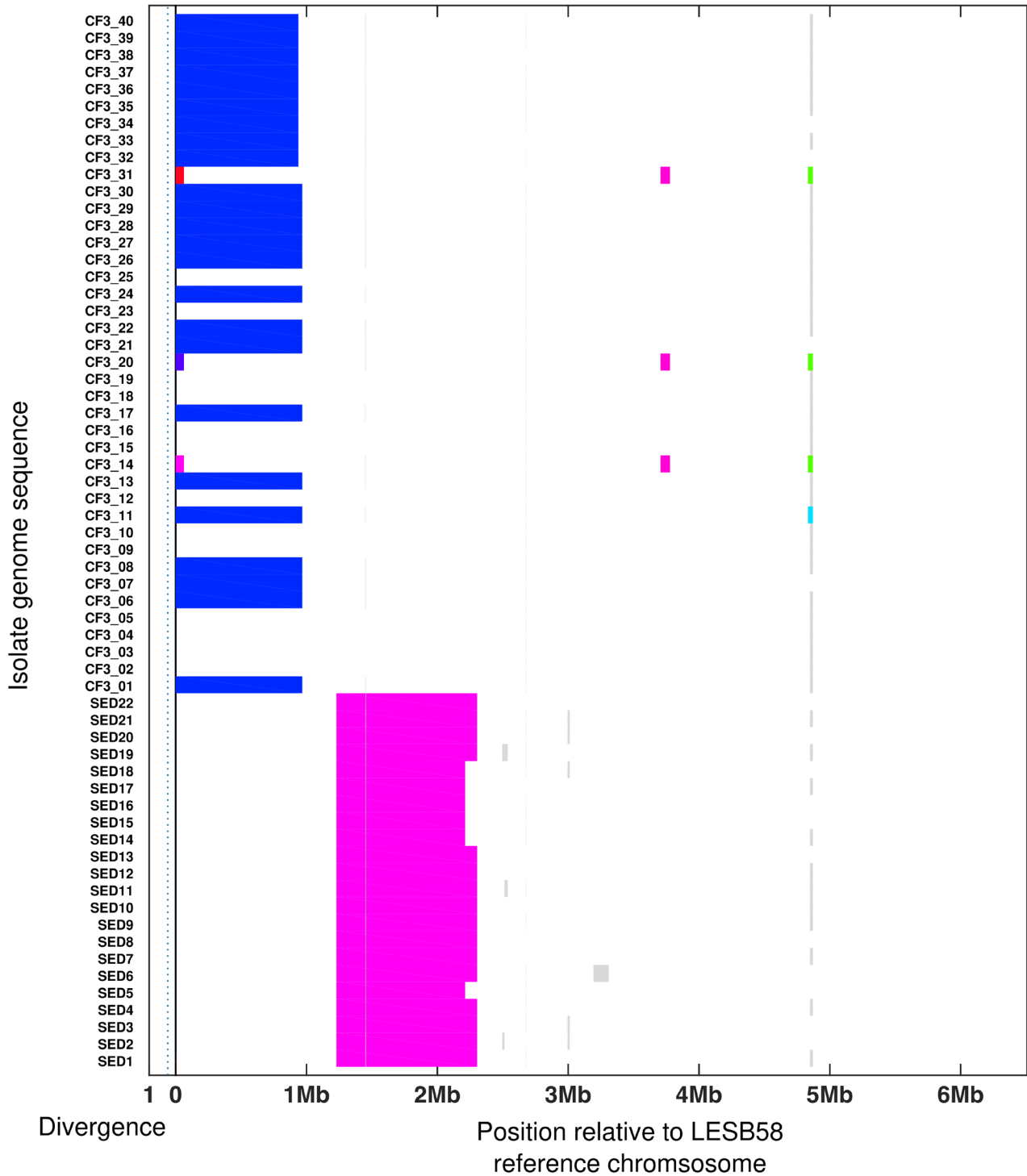
199 Genome Analyser

200

201 Table S5: Proportion of single nucleotide polymorphisms called by the GATK Haplotype Caller recovered in de novo assembled contigs.

202

| Sample | Proportion all variants corroborated | Proportion filtered variants corroborated* |
|---|---|---|
| SED01 | 96.83% | 100.00% |
| SED02 | 92.56% | 96.00% |
| SED03 | 97.54% | 100.00% |
| SED04 | 95.97% | 100.00% |
| SED05 | 92.48% | 98.04% |
| SED06 | 94.92% | 100.00% |
| SED07 | 91.41% | 97.94% |
| SED08 | 93.75% | 97.98% |
| SED09 | 92.48% | 99.02% |

| | | |
|---|---|---|
| SED10 | 94.26% | 97.96% |
| SED11 | 89.84% | 97.96% |
| SED12 | 93.10% | 96.94% |
| SED13 | 88.10% | 100.00% |
| SED14 | 94.40% | 100.00% |
| SED15 | 91.41% | 97.96% |
| SED16 | 93.75% | 99.00% |
| SED17 | 94.49% | 100.00% |
| SED18 | 95.31% | 100.00% |
| SED19 | 92.91% | 97.98% |
| SED20 | 92.11% | 96.94% |
| SED21 | 95.97% | 98.98% |
| SED22 | 96.09% | 100.00% |

203   *Most of the variants removed by the novel filters were not present in the *de novo* assemblies

# FIGURES
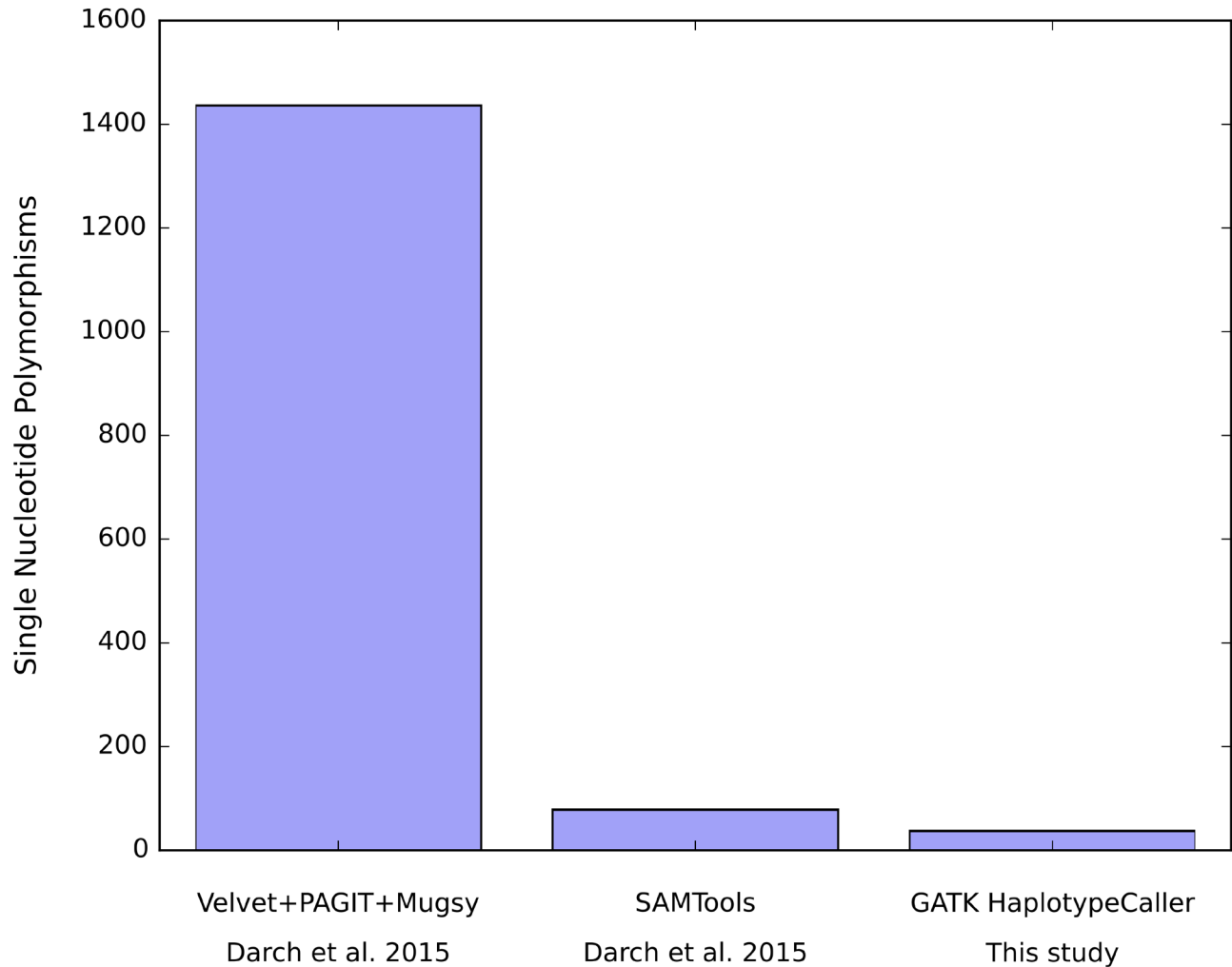


Figure S1. Segments of chromosome indicated by coloured bars, among the Liverpool and Nottingham
isolates, inferred to have undergone homologous recombination. Regions that are the same colour

209    share a common originwhereas grey regions were not deemed significant in a permutation test at
210    alpha = 0.05.
211



212
213    Figure S2. Total single nucleotide polymorphisms called by different methods from the same dataset.
214    In the first method from a previously published analysis, short reads from 22 *P. aeruginosa* isolates
215    were each assembled into contigs *de novo*, improved using an automated pipeline and finally
216    combined into a multiple sequence alignment within which polymorphismic positions were counted.
217    In the second two methods one from the previous study and the third from this report, variants were
218    called among short reads aligned against a complete reference sequence from a closely related isolate
219    (*P. aeruginosa* LESB58).
220
221

222    **REFERENCES**

223     Sickle: A sliding-window, adaptive, quality-based trimming tool for FastQ files (Version 1.29)

224     [https://github.com/najoshi/sickle]

225     **Martin M. (2011)**. Cutadapt removes adapter sequences from high-throughput sequencing reads.

226     *EMBnetjournal***17** 1.

227     **Li H., Durbin R. (2009)**. Fast and accurate short read alignment with Burrows–Wheeler transform.

228     *Bioinformatics***25** 1754-1760.

229     **Li H., Handsaker B., Wysoker A., Fennell T., Ruan J., Homer N., Marth G., Abecasis G., Durbin R. (2009)**. The

230     Sequence Alignment/Map format and SAMtools. *Bioinformatics***25** 2078-2079.

231     Picard Sequence Alignment/Map file manipulation library [http://picard.sourceforge.net/]

232     **McKenna A., Hanna M., Banks E., Sivachenko A., Cibulskis K., Kernytsky A., Garimella K., Altshuler D., Gabriel**

233     **S., other authors (2010)**. The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation

234     DNA sequencing data. *Genome Res***20** 1297-1303.

235     **DePristo M.A., Banks E., Poplin R., Garimella K.V., Maguire J.R., Hartl C., Philippakis A.A., del Angel G., Rivas**

236     **M.A., other authors (2011)**. A framework for variation discovery and genotyping using next-generation DNA

237     sequencing data. *Nat Genet***43** 491-498.

238     **Van der Auwera G.A., Carneiro M.O., Hartl C., Poplin R., del Angel G., Levy-Moonshine A., Jordan T., Shakir**

239     **K., Roazen D., other authors (2013)**. From FastQ Data to High-Confidence Variant Calls: The Genome Analysis

240     Toolkit Best Practices Pipeline. *Curr Protoc Bioinform***43** 11.10.1-11.10.33.

241     **Guindon S., Dufayard J., Lefort V., Anisimova M., Hordijk W., O. G. (2010)**. New algorithms and methods to

242     estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0. *Syst Biol***59** 307-21.

243     **Didelot X., Wilson D.J. (2015)**. ClonalFrameML: Efficient Inference of Recombination in Whole Bacterial

244     Genomes. *PLoS Comput Biol***11** e1004041-.

245     **Sukumaran J., Holder M.T. (2010)**. DendroPy: a Python library for phylogenetic computing. *Bioinformatics***26**

246     1569-1571.

247     **McElroy K.E., Luciani F., Thomas T. (2012)**. GemSIM: general, error-model based simulator of next-generation

248     sequencing data. *BMC Genomics***13** 74.

249     **Winstanley C., Langille M.G., Fothergill J.L., Kukavica-Ibrulj I., Paradis-Bleau C., Sanschagrin F., Thomson N.R.,**

250     **Winsor G.L., Quail M.A., other authors (2009)**. Newly introduced genomic prophage islands are critical

251     determinants of *in vivo* competitiveness in the Liverpool Epidemic Strain of *Pseudomonas aeruginosa*. *Genome*

252     *Res***19** 12-23.

253     **Campbell A.M. (1992)**. Chromosomal insertion sites for phages and plasmids. *J Bacteriol***174** 7495-7499.

254     **Wilmes P., Simmons S.L., Denef V.J., Banfield J.F. (2009)**. The dynamic genetic repertoire of microbial

255     communities. *FEMS Microbiol Rev***33** 109-132.

256     **Zhou Y., Liang Y., Lynch K.H., Dennis J.J., Wishart D.S. (2011)**. PHAST: A Fast Phage Search Tool. *Nucleic Acids*

257     *Res***39** W347-W352.