Automated extraction of typing information for bacterial pathogens from whole genome sequence data: Neisseria meningitidis as an exemplar

K A Jolley¹, M C Maiden (martin.maiden@zoo.ox.ac.uk)¹

1. Department of Zoology, University of Oxford, Oxford, United Kingdom

Citation style for this article: Jolley KA, Maiden MC. Automated extraction of typing information for bacterial pathogens from whole genome sequence data: Neisseria meningitidis as an exemplar. Euro Surveill. 2013;18(4):pii=20379. Ávailable online: http://www.eurosurveillance.org/ViewArticle.aspx?ArticleId=20379

Article submitted on 07 December 2012 / published on 24 January 2013

Whole genome sequence (WGS) data are increasingly used to characterise bacterial pathogens. These data provide detailed information on the genotypes and likely phenotypes of aetiological agents, enabling the relationships of samples from potential disease outbreaks to be established precisely. However, the generation of increasing quantities of sequence data does not, in itself, resolve the problems that many microbiological typing methods have addressed over the last 100 years or so; indeed, providing large volumes of unstructured data can confuse rather than resolve these issues. Here we review the nascent field of storage of WGS data for clinical application and show how curated sequence-based typing schemes on websites have generated an infrastructure that can exploit WGS for bacterial typing efficiently. We review the tools that have been implemented within the PubMLST website to extract clinically useful, straincharacterisation information that can be provided to physicians and public health professionals in a timely, concise and understandable way. These data can be used to inform medical decisions such as how to treat a patient, whether to instigate public health action, and what action might be appropriate. The information is compatible both with previous sequence-based typing data and also with data obtained in the absence of WGS, providing a flexible infrastructure for WGSbased clinical microbiology.

Introduction

The application of whole genome sequencing (WGS) technology to clinical microbiology has been described as revolutionary: the opportunities are certainly immense, but so too are the challenges of implementing this technology effectively [1]. Above all, clinical microbiology and epidemiology are pragmatic sciences, which require accurate and understandable information to be delivered to those who need to make medical judgements in real time. Often these judgements have to be made in the absence of complete information, and it is essential that widely understood, accepted and reproducible typing methods are employed to guide these decisions [2]. Just as the advent of molecular

techniques challenged phenotypic methodologies over a decade ago – replacing imperfect but at least widely accepted techniques with a plethora of non-standardised alternatives [3] - the high volumes of sequence data have to be carefully managed if they are to provide enlightenment rather than confusion.

The multilocus sequence typing (MLST) paradigm was established in 1998 [4], a time when molecular techniques were beginning to be widely used in the clinical laboratory, but when there was no universally agreed way forward [5]. It was intended as a standardised, reproducible and portable approach that could replace and enhance previous methods, particularly multilocus enzyme electrophoresis (MLEE) [6]. MLST was the first sequence-based approach to the genome-wide characterisation of bacterial isolates to be widely adopted and automated methods for performing the reactions and extracting the sequence information have subsequently been developed [7-9]. At the time MLST was introduced, it was impractical to sequence whole genomes on very large numbers of isolates and early analyses showed that in many cases this was not required. The first MLST scheme, for example, was designed to identify major clones within populations of Neisseria meningitidis, the meningococcus, and was able to do this reliably and reproducibly with just seven gene fragments, totalling only 3,284 bp or about 0.15% of the whole genome [10,11]. Similar numbers and sizes of loci have been successful for MLST schemes covering a wide range of organisms, which is an indication of the high degree of structuring present in many bacterial populations. For many bacteria, including the meningococcus, the extent of genetic diversity present even in this small number of genes under stabilising selection is extensive [12]: as of November 2012, each of the gene fragments used as meningococcal MLST loci had between 424 to 675 distinct alleles recorded on the PubMLST *Neisseria* website [13], with 54–94% (mean: 71%) sites that were polymorphic. Furthermore, in the representative *abcZ* locus, all four bases were present at a given site over the known population in 54/433 (12%) of the nucleotide positions (Figure 1). Much of this variation is at low frequency and transitory, but

the variants for which this is the case for cannot be known without exhaustive, or at least extensive, sampling over time.

The MLST approach catalogues this extreme diversity, which is seen in many microbial populations and which remains only partially explored, by the maintenance of curated libraries of allele sequences for each MLST locus. Each unique sequence (allele) is assigned a unique arbitrary number, effectively compressing 400–600 bp of information into a single integer. Further organisation and compression of genetic variation is attained by combining the data from all MLST loci into allelic profiles or sequence types (STs), which are also assigned arbitrary numeric designations, each of which defines a unique string of several thousand nucleotides [12]. This approach has proved to be both efficient and effective: as of November 2012, there were 9,927 STs in the *Neisseria* MLST database, for example, each precisely characterising a particular seven-locus *Neisseria* genotype. Similar levels of diversity have been observed in other bacteria hosted at PubMLST and on other MLST repositories [14]. The fact that nearly 10,000 distinct variants of only 3,284 bp of coding sequence under stabilising selection are known to exist in one human-associated bacterium with a genome of about 2.2 Mbp indicates the scale of the cataloguing problem facing us in the era of genomic microbiology.

Nevertheless, there are instances when even the very high levels of diversity routinely seen in MLST datasets do not provide sufficient information for clinical decision-making. This is because even populations of

FIGURE 1

Schematic of one of the *Neisseria meningitidis* MLST loci (*abcZ*) showing the number and positions of known polymorphic sites within the gene fragment (unmodified PubMLST.org screenshot)



MLST: multilocus sequence typing. Source: PubMLST *Neisseria* website [13].

diverse organisms, such as the meningococcus, are highly structured, with most isolates belonging to clonal complexes of related bacteria, many of which share identical STs [15]. This detection of population structuring is one of the strengths of the MLST approach, as these clusters are frequently associated with phenotypes of clinical interest such as virulence or expression of vaccine antigens [16]. This clustering, however, can mean that isolates with the same ST may not have the same point source, so ST alone is insufficient to unambiguously identify strains belonging to an outbreak. For this reason, additional highly variable antigenic loci are included in the recommended typing scheme for meningococci [17] and for other organisms such as Campylobacter [18] that are regularly typed by MLST. For meningococci, there are also curated sequence-based schemes for genes that encode antimicrobial resistance that provide additional clinically valuable information [19,20]. Other schemes, such as variable-number tandem repeat (VNTR), also allow high discrimination of isolates in outbreak situations [21,22]. Combining these high-resolution typing approaches with seven-locus MLST and spatial and temporal epidemiology techniques permits the proactive identification of outbreaks of infectious disease [23].

For a small number of bacteria, the so-called single clone pathogens, there is insufficient variation in seven-locus MLST to provide epidemiological resolution, usually because these pathogens have evolved recently from single clones, undergo little recombination and contain too little genetic variation [24]. These include organisms of great medical importance such Mycobacterium tuberculosis [25], Yersinia pestis [26], Bacillus anthracis [27] and Salmonella enterica var Typhi [28]. For these bacteria, data from the whole genome, often in the form of single nucleotide polymorphisms (SNPs) [29], but also including other types of variation such as VNTRs, is essential for epidemiological purposes. These data will also have to be stored and interpreted in an accessible way that produces data usable by clinical decision-makers and which is both forwards and backwards compatible.

One of the motivations that drove the development of MLST was future-proofing. Even at a time when the costs of sequencing were seen by some as prohibitive [30], nucleotide sequence data had major advantages: they might be added to, but they would never become obsolete - as they represented the fundamental level of genetic information – and they are readily understood, stored, compared and distributed [12]. Obtaining WGS data is now becoming so inexpensive that it is becoming the fastest and most economical way of obtaining information at multiple loci for determining MLST or other STs [31]. When used in this way, these data are directly comparable to the extensive sequence databases that have been established since the first use of MLST [32,33]. Here we describe how the suite of databases hosted at PubMLST [34] has been updated

to accommodate WGS data and describe the tools that are available to rapidly extract typing information from such data. We also describe how these tools can be exploited further to achieve very high resolution from such data when required.

Database structure

As of November 2012, the majority of the typing databases hosted at PubMLST [34] were using the Bacterial Isolate Genome Sequence Database (BIGSdb) platform to archive isolate and sequence diversity data [35]. This software was developed to facilitate the flexible storage and exploitation of the whole range of sequence data that might be available from a clinical specimen, from single Sanger sequencing reads through to whole genomes, which may be either complete or consisting of multiple contiguous sequences ('contigs'), as assembled from data from the current generation of sequencing instruments. The BIGSdb platform consists of two kinds of database: (i) a definition database that contains the sequences of known alleles of loci under study, as well as allelic profiles (combinations of alleles at specific loci) for schemes such as MLST; and (ii) an isolate database that contains isolate provenance and other metadata along with nucleotide sequences associated with that isolate. An isolate database can interact with any number of definition databases and vice versa, allowing networks of authoritative nomenclature servers and partitioning of isolate datasets and projects, with curator access controlled by specific permissions set by an administrator.

Reference databases

The definition databases are central to genome analysis using the gene-by-gene (MLST-like) analysis approach implemented in BIGSdb. By storing all known allelic diversity for any locus of interest, the definition databases provide a centralised queryable repository that provides a common language for expressing sequence differences, making it a trivial process to identify alleles that are different among isolates, and equally importantly, those that are identical. Because sequence differences are linked directly to a particular locus (which can be any definable sequence string, nucleotide or peptide) and with appropriate grouping of loci into 'schemes' (groups of related loci), the context of this locus is immediately apparent: identifying it, for example, as a member of a conventional MLST scheme, as responsible for antimicrobial resistance, as a participant of a biochemical pathway and so on. As of November 2012, the Neisseria PubMLST definition database had allelic sequences defined for 1,272 loci with 114,469 unique alleles.

Extracting typing information

Web-based and stand-alone tools have been developed that facilitate identification of STs directly from short-read data [36,37]. These methods are, of course, dependent on the sequence and profile definitions made available on PubMLST.org, which also has functionality to extract typing information directly from submitted assembled genomes that are routinely scanned for known alleles. As the locations of these loci are 'tagged' in the sequence data for future reference within BIGSdb, this means that the genome sequences are automatically annotated for those loci for which definition databases exist. The definition database can also be queried using genome data not uploaded to the isolate database to identify a strain directly from sequence data. The BIGSdb platform also has functionality that enables an administrator to define scanning rules and report formatting. This uses a built-in script interpreter so that analysis paths can be taken by following a decision tree defined by the rules. This has been implemented within the PubMLST Neisseria sequence definition database to automatically extract the strain typing information for the meningococcus (ST, clonal complex and antigen sequence type comprising PorA variable regions and FetA variable region) [17,33], along with antibiotic resistance information from sequence data that is pasted in to a web form (Figure 2, panel A). The script instructs the software to first scan the MLST alleles and, if these are all identified, to identify the ST and clonal complex by querying the reference data tables. It then scans the typing antigens and formats the results of these with the MLST results in to a standardised strain designation [17]. Following this, the sequences of the penA and rpoB genes are extracted and then compared with isolates with matching sequences within the PubMLST isolate database to determine the most likely penicillin and rifampicin sensitivity. All of this is displayed in a plain language report (Figure 2, panel B). The whole analysis is extremely rapid, taking about 40 seconds within the web interface.

Comparing genomes

Because genomic diversity is recorded within BIGSdb as allele numbers, WGS analysis is possible using the highly scalable techniques developed for seven-locus MLST. Once loci have been defined and alleles identified, they can be used essentially as a whole-genome MLST scheme, or any chosen subset of predefined loci combined to form a scheme. This is the principle behind the Genome Comparator analysis [38], which can use either the defined loci or extract coding sequences from an annotated reference genome to perform comparisons against genomes within the database. Using a reference genome, or set of predefined reference loci, each of the coding sequences are compared against the test genomes using BLAST. Allele sequences that are the same as the reference are designated allele 1, while each unique allele different from the reference is assigned a sequential number. Once each locus has been tested, a distance matrix is then generated based on allelic identities between each pair of isolates. This can then be visualised using standard algorithms – the PubMLST website incorporates the Neighbor-net algorithm [39] implemented in SplitsTree4 [40]. Because analysis relies only on using BLAST to compare each locus within a genome in turn, either against the single annotated reference sequence or against all known

alleles if using defined loci, the analysis is again very rapid, allowing multiple genomes to be compared within minutes, with the time taken to analyse only increasing linearly, not geometrically, with additional genomes.

The Genome Comparator approach is generic and any number of loci in any groups can be used for this type of analysis. Many loci have been defined for the meningococcus, including the 53 ribosomal (r) genes that are used as a basis of rMLST [41-44]. The full complement of ribosomal genes has a number of advantages for indexing variation. These genes are universally present in members of the domain, are protein encoding and therefore generally assemble well from short-read sequences and are distributed throughout the genome. They encode proteins that form part of a coherent, macromolecular structure and contain variation that is informative at a wide range of levels of discrimination. These data can be used within and among members of the same genus, for both species and strain definition [42].

Analysis of whole genome sequence data for meningococci

The *Neisseria* PubMLST database is continually expanding: as of November 2012, there were 221 isolate records with deposited genome sequence data linked to published studies [11,45-51]. Of these 221 genomes, 170 were meningococci, with the remainder belonging to other species within the genus [42]. The data consisted of a mixture of finished genomes, multiple contigs generated from de novo assembly, contigs generated by mapping to a reference sequence and sets of predicted coding sequences. These are treated identically by BIGSdb to identify and tag sequences of known loci, and where these loci are members of existing typing schemes, such as MLST or antigen typing, these genomes could be compared to legacy data (Table).

Neighbor-net visualisation of distance matrices generated with Genome Comparator from allelic rMLST data [44] provides a highly scalable, rapid and easily understood way of placing isolates within the known diversity of a bacterial species. For example, the interrelationships of 139 *N. meningitidis* isolates present in the PubMLST *Neisseria* database [13] can be efficiently represented by this method. Since rMLST alleles are automatically tagged within the database, this analysis is rapid and the Neighbor-net trees can be generated in a few minutes. The rMLST analysis differentiates clonal complexes; however, in addition it provides much higher resolution than conventional seven-locus MLST [38], robustly indicating both relationships among and diversity within clonal complexes (Figure 3).

The locations of isolates belonging to major clonal complexes identified by conventional MLST are indicated (cc1, etc.). The figure illustrates relationships not apparent from seven-locus MLST, including the

FIGURE 2

Extracting antigen and antibiotic resistance data from Neisseria meningitidis whole genome sequences



A whole genome sequence, which may consist of multiple contigs, can be pasted in to the Neisseria PubMLST website (panel A) with typing and antibiotic resistance data for penicillin and rifampicin rapidly extracted (panel B) (unmodified PubMLST.org screenshots).

Source: PubMLST Neisseria website [13].

TABLE

Meningococcal whole genome sequencing data linked to published studies, deposited in the PubMLST *Neisseria* database as of November 2012

| Clonal complex | Number of genome sequences | Number of STs | Serogroups | PorA variant combinations | FetA variants |
|----------------|-------------------------------|---------------|-------------------------------------|---------------------------|---------------|
| CC11 | 31 | 6 | C (22), W (4), B(2), NG (1), NA (2) | 8 | 8 |
| cc41/44 | 20 | 12 | B (14), NA (5), NG (1) | 10 | 5 |
| CC32 | 17 | 4 | B (14), C (1), NG (1), NA (1) | 10 | 5 |
| сс5 | 16 | 5 | A (16) | 3 | 5 |
| cc4 | 14 | 1 | A (14) | 4 | 1 |
| CC1 | 13 | 3 | A (13) | 4 | 5 |
| cc8 | 9 | 5 | B (5), C (3), NA (1) | 6 | 5 |
| cc18 | 5 | 4 | B (4), C (1) | 5 | 4 |
| CC23 | 5 | 2 | Y (5) | 3 | 2 |
| CC22 | 4 | 1 | W (4) | 1 | 2 |
| cc167 | 4 | 4 | Y (4) | 1 | 2 |
| cc269 | 4 | 3 | B (2), NA (2) | 4 | 3 |
| cc37 | 2 | 2 | B (2) | 1 | 2 |

NA: not available; NG: non-groupable; ST: sequence type.

The table shows the clonal complex and indicates the diversity of ST, serogroup and typing antigens. Only clonal complexes represented by two or more genomes are included.

diversity of some clonal complexes (e.g. cc1) and the interrelationships of others, e.g. cc8 and cc11 clonal complexes, and the relationships of the ET-15 and ET-37 variants within cc11.

Conclusions and future prospects

Nucleotide sequences are a universal language that can be interpreted in a number of ways. For clinical and epidemiological purposes, sequences from clinical specimens have to be rapidly and effectively translated into a meaningful term or set of terms that define those properties of the aetiological agents of disease that direct medical and public health action. One of the factors behind the success of seven-locus MLST was the introduction of standard sets of nomenclature that reflected the structure of microbial populations and their phenotypic properties. For organisms with well-established and accepted MLST and other typing schemes in place, the impact of the application of WGS data will be to rapidly identify properties such as strain type. In some cases, novel nomenclature may be required, but this is a process that has to be approached with care, if confusion in the wider clinical community is to be avoided.

The suite of database subsites on PubMLST, which now includes a site that catalogues the ribosomal diversity across the whole domain for the purposes of rMLST typing [44,52], provides an example of how WGS data can be used to efficiently designate specimens to current strain types. It can be also used to establish additional typing schemes which can coexist with each other side

by side, as there is no limit to the number of loci and schemes that can be defined. As the database stores the sequence information that is available for an isolate, be that a single read or a whole genome, it means that it is possible to seamlessly compare isolates for which different types of information are available, achieving backwards compatibility with previous typing schemes, as well as compatibility with diagnostic tests that may target only one or a few loci. The extent to which isolates can be compared depends only on the quality of the sequence data available for the locus in question, but given that clinical specimens are often imperfect, it is important for clinical and epidemiological purposes that incomplete or partial information can be used. While many studies place short-read data in a sequence read archive, this is not easily accessible or readily analysed. PubMLST curators do proactively assemble short-read data and incorporate the resultant contigs into the database where metadata are available. Links are made to the sequence read archive within PubMLST isolate records so that original data can be retrieved and analysed when required. While the Neisseria databases described are exemplars, databases for other species can be hosted on request and the open-source BIGSdb software is freely available for local installation.

The first analyses of WGS data on bacterial specimens relied on SNP analysis of closely related bacteria, with mapping of sequence reads to a predefined reference genome. These have required pre-analysis of the samples by an approach such as MLST to limit the extent

FIGURE 3

Relationships of 139 Neisseria meningitidis genomes in the PubMLST Neisseria database, generated with Genome Comparator and Neighbor-net from allelic profiles data for rMLST loci

r: ribosomal; MLST: multilocus sequence typing.

The locations of isolates belonging to major clonal complexes identified by conventional MLST are indicated (cc1, etc.). The figure illustrates relationships not apparent from seven-locus MLST, including the diversity of some clonal complexes (e.g. cc1) and the interrelationships of others, e.g. cc8 and cc11 clonal complexes, and the relationships of the ET-15 and ET-37 variants within cc11.

of variation being analysed [53-58]. This approach is also appropriate and can be very effective for 'single clone' pathogens [25-28]; however, it is not feasible for the general analysis for diagnosis or surveillance of bacteria such as the meningococcus that exhibit more typical levels of sequence diversity. Indeed, the use of the term SNP when discussing bacterial genome variation outside the examples described above, is unfortunate and can be misleading. The concept of the 'SNP' has been taken from human medical genomics to microbial genomics: in humans, it is in some cases appropriate to discuss SNPs, when they are associated with a particular genetic disease, but genetic variation in terms of sequence polymorphism is much more complex in bacteria. As seen here, the great majority of microbial populations contain tens of thousands of polymorphisms even within organisms that are closely related – not to mention large amounts of variation due to insertions, deletions and rearrangements, which cannot even remotely be described as 'SNPs'. The term sequence variation is more appropriate as individual polymorphisms, especially in bacteria, are invariably embedded with many other variants into alleles and it is these alleles – each often with many variable sites – that are associated with particular phenotypes.

Although the typing of bacterial specimens with existing schemes is a valuable contribution of WGS data to clinical microbiology and epidemiology, it is not, of course, the only use for these data. There are many other possible applications for both research and detailed investigation of outbreaks [38]; however, it is important that the analysis of these data is driven by the question that is being asked. If an outbreak can be resolved with a few loci, then there is no need to pursue the data further and certainly no need to report more detail than necessary to a hard-pressed frontline clinician or epidemiologist who, in general, will only require the information necessary to resolve the medical problem at hand. In other cases, resolution of a particular outbreak may require data from the whole genome [53]. For this reason, it will be increasingly necessary to store WGS data from clinical specimens in an understandable form, that is, as assembled sequences, within flexible structures, such as that offered by the PubMLST platform powered by BIGSdb, where WGS information can be hierarchically queried in real time by individuals with limited bioinformatics expertise to generate the data at the resolution required to address their problem. In this context these data will provide an exciting opportunity to extend our understanding of infectious disease caused by bacteria and will enhance our ability to combat it.

References

- 1. Pallen MJ, Loman NJ, Penn CW. High-throughput sequencing and clinical microbiology: progress, opportunities and challenges. Curr Opin Microbiol. 2010;13(5):625-31.
- van Belkum A, Tassios PT, Dijkshoorn L, Haeggman S, Cookson B, Fry NK, et al. Guidelines for the validation and application of typing methods for use in bacterial epidemiology. Clin Microbiol Infect. 2007;13 Suppl 3:1-46.
- 3. Achtman M. A surfeit of YATMs? J Clin Microbiol. 1996;34(7):1870.
- Maiden MC, Bygraves JA, Feil E, Morelli G, Russell JE, Urwin R, et al. Multilocus sequence typing: a portable approach to the identification of clones within populations of pathogenic microorganisms. Proc Natl Acad Sci USA. 1998;95(6):3140-5.
- van Belkum A, Struelens M, de Visser A, Verbrugh H, Tibayrenc M. Role of genomic typing in taxonomy, evolutionary genetics, and microbial epidemiology. Clin Microbiol Rev. 2001;14(3):547-60.
- Selander RK, Caugant DA, Ochman H, Musser JM, Gilmour MN, Whittam TS. Methods of multilocus enzyme electrophoresis for bacterial population genetics and systematics. Appl Environ Microbiol. 1986;51(5):837-84.
- Sullivan CB, Jefferies JM, Diggle MA, Clarke SC. Automation of MLST using third-generation liquid-handling technology. Mol Biotechnol. 2006;32(3):219-26.
- Platt S, Pichon B, George R, Green J. A bioinformatics pipeline for high-throughput microbial multilocus sequence typing (MLST) analyses. Clin Microbiol Infect. 2006;12(11):1144-6.
- 9. O'Farrell B, Haase JK, Velayudhan V, Murphy RA, Achtman M. Transforming microbial genotyping: a robotic pipeline for genotyping bacterial strains. PLoS One. 2012;7(10):e48022.
- Holmes EC, Urwin R, Maiden MC. The influence of recombination on the population structure and evolution of the human pathogen Neisseria meningitidis. Mol Biol Evol. 1999;16(6):741-9.
- Parkhill J, Achtman M, James KD, Bentley SD, Churcher C, Klee SR, et al. Complete DNA sequence of a serogroup A strain of Neisseria meningitidis Z2491. Nature. 2000;404(6777):502-6.
- 12. Maiden MC. Multilocus sequence typing of bacteria. Annu Rev Microbiol. 2006;60:561-88.
- Neisseria sequence typing home page. Oxford: University of Oxford. [Accessed 31 Nov 2012]. Available from: http:// pubmlst.org/neisseria/
- All species MLST databases and published schemes. [Accessed 31 Nov 2012]. Available from: http://pubmlst.org/databases. shtml
- Caugant DA, Maiden MC. Meningococcal carriage and diseasepopulation biology and evolution. Vaccine. 2009;27 Suppl 2:B64-70.
- 16. Yazdankhah SP, Kriz P, Tzanakaki G, Kremastinou J, Kalmusova J, Musilek M, et al. Distribution of serogroups and genotypes among disease-associated and carried isolates of Neisseria meningitidis from the Czech Republic, Greece, and Norway. J Clin Microbiol. 2004;42(11):5146-53.
- Jolley KA, Brehony C, Maiden MC. Molecular typing of meningococci: recommendations for target choice and nomenclature. FEMS Microbiol Rev. 2007;31(1):89-96.
- Dingle KE, McCarthy ND, Cody AJ, Peto TE, Maiden MC. Extended sequence typing of Campylobacter spp., United Kingdom. Emerg Infect Dis. 2008;14(10):1620-2.
- 19. Taha MK, Hedberg ST, Szatanik M, Hong E, Ruckly C, Abad R, et al. Multicenter study for defining the breakpoint for rifampin resistance in Neisseria meningitidis by rpoB sequencing. Antimicrob Agents Chemother. 2010;54(9):3651-8.
- 20. Taha MK, Vázquez JA, Hong E, Bennett DE, Bertrand S, Bukovski S, et al. Target gene sequencing to characterize the penicillin G susceptibility of Neisseria meningitidis. Antimicrob Agents Chemother. 2007;51(8):2784-92.
- Schouls LM, van der Ende A, Damen M, van de Pol I. Multiple-locus variable-number tandem repeat analysis of Neisseria meningitidis yields groupings similar to those obtained by multilocus sequence typing. J Clin Microbiol. 2006;44(4):1509-18.
- 22. Elias J, Schouls LM, van de Pol I, Keijzers WC, Martin DR, Glennie A, et al. Vaccine preventability of meningococcal clone, Greater Aachen Region, Germany. Emerg Infect Dis. 2010;16(3):464-472.
- 23. Elias J, Harmsen D, Claus H, Hellenbrand W, Frosch M, Vogel U. Spatiotemporal analysis of invasive meningococcal disease, Germany. Emerg Infect Dis. 2006;12(11):1689-95.
- 24. Achtman M. Evolution, population structure, and phylogeography of genetically monomorphic bacterial pathogens. Annu Rev Microbiol. 2008;62:53-70.

- 25. Wirth T, Hildebrand F, Allix-Béguec C, Wölbeling F, Kubica T, Kremer K, et al. Origin, spread and demography of the *Mycobacterium tuberculosis* complex. PLoS Pathogens. 2008;4(9):e1000160.
- 26. Haensch S, Bianucci R, Signoli M, Rajerison M, Schultz M, Kacki S, et al. Distinct clones of *Yersinia pestis* caused the Black Death. Plos Pathogens. 2010;6(10):e1001134.
- 27. Pearson T, Okinaka RT, Foster JT, Keim P. Phylogenetic understanding of clonal populations in an era of whole genome sequencing. Infect Genet Evol. 2009;9(5):1010-9.
- Holt KE, Parkhill J, Mazzoni CJ, Roumagnac P, Weill FX, Goodhead I, et al. High-throughput sequencing provides insights into genome variation and evolution in Salmonella Typhi. Nat Genet. 2008;40(8):987-93.
- Baker L, Brown T, Maiden MC, Drobniewski F. Silent nucleotide polymorphisms and a phylogeny for Mycobacterium tuberculosis. Emerg Infect Dis. 2004;10(9):1568-77.
- Olive DM, Bean P. Principles and applications of methods for DNA-based typing of microbial organisms. J Clin Microbiol. 1999;37(6):1661-9.
- Medini D, Serruto D, Parkhill J, Relman DA, Donati C, Moxon R, et al. Microbiology in the post-genomic era. Nat Rev Microbiol. 2008;6(6):419-30.
- 32. Mellmann A, Harmsen D, Cummings CA, Zentz EB, Leopold SR, Rico A, et al. Prospective genomic characterization of the German enterohemorrhagic Escherichia coli 0104:H4 outbreak by rapid next generation sequencing technology. PLoS One. 2011;6(7):e22751.
- 33. Vogel U, Szczepanowski R, Claus H, Jünemann S, Prior K, Harmsen D. Ion torrent personal genome machine sequencing for genomic typing of Neisseria meningitidis for rapid determination of multiple layers of typing information. J Clin Microbiol. 2012;50(6):1889-94.
- 34. PubMLST. Oxford: University of Oxford. [Accessed 31 Nov 2012]. Available from: http://pubmlst.org/
- 35. Jolley KA, Maiden MC. BIGSdb: Scalable analysis of bacterial genome variation at the population level. BMC Bioinformatics. 2010;11:595.
- Larsen MV, Cosentino S, Rasmussen S, Friis C, Hasman H, Marvig RL, et al. Multilocus sequence typing of total-genomesequenced bacteria. J Clin Microbiol. 2012;50(4):1355-61.
- Inouye M, Conway TC, Zobel J, Holt KE. Short read sequence typing (SRST): multi-locus sequence types from short reads. BMC Genomics. 2012;13:338.
- 38. Jolley KA, Hill DM, Bratcher HB, Harrison OB, Feavers IM, Parkhill J, et al. Resolution of a meningococcal disease outbreak from whole-genome sequence data with rapid webbased analysis methods. J Clin Microbiol. 2012;50(9):3046-53.
- 39. 39. Bryant D, Moulton V. Neighbor-net: an agglomerative method for the construction of phylogenetic networks. Mol Biol Evol. 2004;21(2):255-65.
- 40. Huson DH, Bryant D. Application of phylogenetic networks in evolutionary studies. Mol Biol Evol. 2006;23(2):254-67.
- 41. Read DS, Woodcock DJ, Strachan NJ, Forbes KJ, Colles FM, Maiden MC, et al. Evidence for phenotypic plasticity among multihost Campylobacter jejuni and C. coli lineages, obtained using ribosomal multilocus sequence typing and Raman spectroscopy. Appl Environ Microbiol. 2013;79(3):965-73.
- 42. Bennett JS, Jolley KA, Earle SG, Corton C, Bentley SD, Parkhill J, Maiden MC. A genomic approach to bacterial taxonomy: an examination and proposed reclassification of species within the genus Neisseria. Microbiology. 2012;158(Pt 6):1570-80.
- Ussery DW, Gordon SV. Two novel methods for using genome sequences to infer taxonomy. Microbiology. 2012;158(Pt 6):1414.
- 44. Jolley KA, Bliss CM, Bennett JS, Bratcher HB, Brehony CM, Colles FM, et al Ribosomal multilocus sequence typing: universal characterisation of bacteria from domain to strain. Microbiology. 2012;158(Pt 4):1005-15.
- 45. Hao W, Ma JH, Warren K, Tsang RS, Low DE, Jamieson FB, et al. Extensive genomic variation within clonal complexes of Neisseria meningitidis. Genome Biol Evol. 2011;3:1406-18.
- 46. Budroni S, Siena E, Hotopp JCD, Seib KL, Serruto D, Nofroni C, et al. Neisseria meningitidis is structured in clades associated with restriction modification systems that modulate homologous recombination. Proc Natl Acad Sci USA. 2011;108(11):4494-9.
- 47. Schoen C, Blom J, Claus H, Schramm-Glück A, Brandt P, Müller T, et al. Whole-genome comparison of disease and carriage strains provides insights into virulence evolution in Neisseria meningitidis. Proc Natl Acad Sci USA. 2008;105(9):3473-8.
- Katz LS, Humphrey JC, Conley AB, Nelakuditi V, Kislyuk AO, Agrawal S, et al. Neisseria Base: a comparative genomics

database for Neisseria meningitidis. Database (Oxford). 2011;2011:bar035.

- 49. Rusniok C, Vallenet D, Floquet S, Ewles H, Mouzé-Soulama C, Brown D, et al. NeMeSys: a biological resource for narrowing the gap between sequence and function in the human pathogen Neisseria meningitidis. Genome Biol. 2009;10(10):R110.
- 50. Bentley SD, Vernikos GS, Snyder LA, Churcher C, Arrowsmith C, Chillingworth T, et al. Meningococcal genetic variation mechanisms viewed through comparative analysis of serogroup C strain FAM18. PLoS Genet. 2007;3(2):e23.
- 51. Peng J, Yang L, Yang F, Yang J, Yan Y, Nie H, et al. Characterization of ST-4821 complex, a unique Neisseria meningitidis clone. Genomics. 2008;91(1):78-87.
- 52. Ribosomal multilocus sequence typing (rMLST). Oxford: University of Oxford. [Accessed 31 Nov 2012]. Available from: http://rmlst.org/
- 53. Köser CU, Holden MT, Ellington MJ, Cartwright EJ, Brown NM, Ogilvy-Stuart AL, et al. Rapid whole-genome sequencing for investigation of a neonatal MRSA outbreak. N Engl J Med. 2012;366(24):2267-75.
- 54. Croucher NJ, Harris SR, Fraser C, Quail MA, Burton J, van der Linden M, et al. Rapid pneumococcal evolution in response to clinical interventions. Science. 2011;331(6016):430-4.
- 55. Harris SR, Feil EJ, Holden MT, Quail MA, Nickerson EK, Chantratita N, et al. Evolution of MRSA during hospital transmission and intercontinental spread. Science. 2010;327(5964):469-74.
- 56. Grad YH, Lipsitch M, Feldgarden M, Arachchi HM, Cerqueira GC, Fitzgerald M, et al. Genomic epidemiology of the Escherichia coli 0104:H4 outbreaks in Europe, 2011. Proc Natl Acad Sci USA. 2012;109(8):3065-70.
- 57. Young BC, Golubchik T, Batty EM, Fung R, Larner-Svensson H, Votintseva AA, et al. Evolutionary dynamics of Staphylococcus aureus during progression from carriage to disease. Proc Natl Acad Sci USA. 2012;109(12):4550-5.
- 58. Eyre DW, Golubchik T, Gordon NC, Bowden R, Piazza P, Batty EM, et al. A pilot study of rapid benchtop sequencing of Staphylococcus aureus and Clostridium difficile for outbreak detection and surveillance. BMJ Open. 2012;2(3). pii: e001124.