# Microbial Genomic Data Analysis for Infectious Diseases

**Dr. Margaret Varga, Dr. Charlene Rodrigues, Dr. Keith Jolley, Dr. Holly Bratcher, Dr. Jenny MacLennan, Dr. Odile Harrison, Dr. Fran Colles, Dr. Alison Cody, Dr. James Bray and Prof. Martin Maiden**

Department of Zoology, University of Oxford, South Parks Road, Oxford, OX1 3SY,
UNITIED KINGDOM

margaret.varga@zoo.ox.ac.uk

## ABSTRACT

*Infectious diseases are caused by pathogenic micro-organisms which can be bacteria, viruses, parasites or fungi. The diseases can be spread through many different routes, either directly or indirectly. Military personnel are at high risk of contracting infections, in particular vector-borne and zoonotic infections, during overseas deployments, where they may be exposed to endemic or emerging infections to which they do not have immunity. Additionally, overcrowded settings with poor sanitation are high risks for disease.*

*Genomics is having a transformational impact on medicine. It is enabling advances in accurate diagnosis of infectious disease, development of effective and targeted treatment strategies and opportunities to assess pathogenicity. Further, it supports the detection, surveillance of infectious diseases, the development and assessment of vaccines, as well as the assessment and prediction of anti-microbial resistance. These capabilities are all key military needs to protect personnel in this inter-connected world.*

*The advances in sequencing technologies have resulted in an explosion of genomic data. However, making sense of genomic data requires advances in computational analysis technologies together with cross-disciplinary scientific approaches, skill sets and people. There are extensive reference databases of genomic data. One such open access database is PubMLST.org: it contains well curated genomes for more than 100 microbial species and genera integrated with provenance and phenotype information. All levels of sequence data, from single gene sequences up to and including complete, finished genomes can be accessed on this platform. This data is, however, both large and complex and intractable to analyse and understand using traditional analysis tools.*

*This paper will discuss the challenges of analysing such genomic data for bacterial infections and consider the application of bioinformatics tools and techniques to analyse and communicate microbial genomic data in healthcare.*

## 1.0 INTRODUCTION

There are three different categories of military casualties, namely, battle casualties, non-battle injuries and those suffering from non-battle related diseases or infections. Infectious diseases are caused by pathogenic micro-organisms which can be bacteria, viruses, parasites or fungi, with varying global epidemiology. Diseases can be spread through many different routes, including direct person-to-person spread, or indirectly through the air, food, water, and vectors such as arthropods or animals. Military personnel are at high risk of contracting infections, in particular vector-borne and zoonotic infections, especially during overseas deployments, as they travel to areas with different exposures to the home region against which they had developed immunity, Additionally, living in barracks or settings where sanitation may be poor and living condition crowded can lead to other public health risks [12].

The health status of military personnel has direct bearing on the success of a military campaign's outcome. The 1918 influenza pandemic (January 1918–December 1920), termed the Spanish flu, occurred during World War I and represented the most severe pandemic in recent history. The war and the pandemic intertwined. The infection was caused by an H1N1 virus with genes of avian origin which spread worldwide during 1918-1919. It was first identified in US military personnel in the spring of 1918. The compact living conditions, poor hygiene, malnutrition and the global troop movement contributed to the spread of the infection. The lack of vaccine and treatments resulted in an estimated 500 million people, or one-third of the world's population, becoming infected with this virus. The number of deaths was estimated to be at least 50 million worldwide, mostly young previously healthy adults, greater than the almost 20 million death toll of World War I [5, 20 and 32]. The US Army medical department concluded that secondary bacterial pneumonia, a common secondary infection complicating influenza, were the cause of nearly all the deaths rather than the primary infection [5]. Indeed, in World War I more deaths in the US armed forces (army and navy) engaged in the battlefield were due to pneumonia associated with the Spanish flu than from direct combat injuries [5]. The epidemic resulted in the diversion of personnel, resources, human attention as well as energy from the military campaign. Even during peace time, infectious diseases may spread easily in military bases and training camps resulting in epidemics with material effect on the military power of the units or the country in question. Vaccination and other preventative treatments play an important role in protecting the armed forces [3, 4, 5, 13, 18, 31 and 41].

## 2.0 MICROBIAL GENOMICS DATA

The sequencing of disease-causing microorganisms and study of genetic code is having a transformational impact on medicine, enabling advances in accurate diagnosis of infectious diseases, development of targeted treatment strategies and opportunities to assess pathogenicity of disease-causing organisms. Further, it supports the detection and surveillance of infectious diseases, vaccine development and post-implementation assessment and prediction of anti-microbial resistance. As the military, civilians and animals move with ease around the world, genomic technology can help in the areas of disease prevention, surveillance and management. These capabilities are all key military needs so as to protect personnel in this inter-connected world.

Multilocus Sequence Typing (MLST) was the most widely used portable, reproducible sequence-based approach for bacterial typing in the 21st century, it is used in identifying relationship among bacteria and cataloguing these using stable nomenclatures [25]. Seven MLST loci are indexed in most of the MLST schemes where an arbitrary and unique 'allele' number is assigned to each unique sequence for each locus. The designations for each of the loci are incorporated into a sequence type (ST), e.g. *Campylobacter jejuni* ST-21 or *Neisseria meningitidis* ST-11. Each ST summarizes thousands of base pairs of information; for some species, many hundreds of alleles at each locus and thousands of STs have been identified. Further STs can be grouped by similarity at more than four of the seven loci, into clonal complexes.

The Bacterial Isolate Genome Sequence Database (BIGSdb) platform is designed to store and analyse sequence data for bacterial isolates in which any number of sequences can be linked to isolate records [19]. It extends the principle of MLST to genomic data, where large numbers of loci can be defined, with alleles assigned by reference to sequence definition databases. Loci can also be grouped into schemes, most commonly of biological relevance, so that types can be defined by combinations of alleles (allelic profiles), a concept similar to MLST, Figure 1. This provides an effective mean to conduct gene-by-gene typing [26].
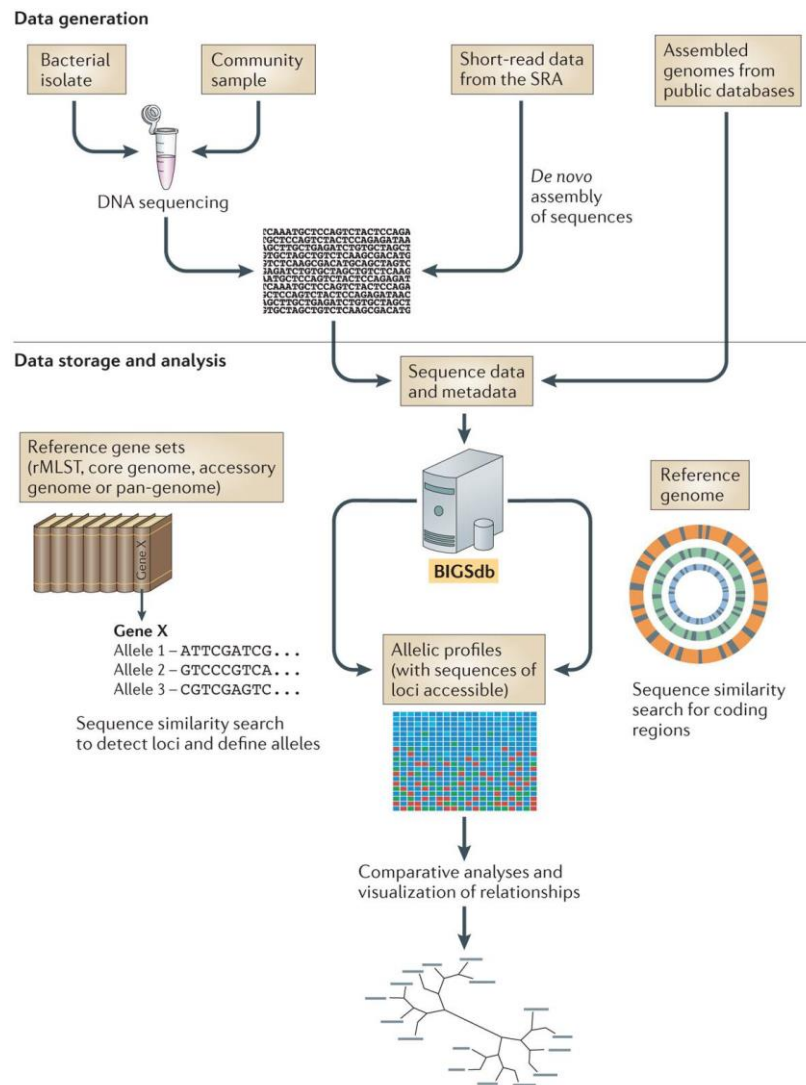
**Figure 1: Illustrative representation of the gene-by-gene approach to whole genome sequence analysis using the Bacterial Isolate Genome Sequence Database (BIGSdb) platform. Reproduced with permission (Maiden et al. 2013) [26].**

There are extensive reference databases of microbial genomic data. One such open access database is PubMLST.org which hosts databases using the Bacterial Isolate Genome Sequence Database platform (BIGSdb) [19 and 26]. Users can download genomic and associated metadata from those databases via PubMLST.org (or indeed from other public databases) and can analyse the data using the open source BIGSdb platform. PubMLST contains well curated genomes for more than 100 microbial species and genera integrated with provenance and phenotype information. It includes all levels of sequence data, from single gene sequences up to and including complete, reference genomes, see Figure 2. The allelic variants in the MLST data support sequence-based analysis. At the time of writing, PubMLST.org held more than 750,000 isolates from different countries and of different species.
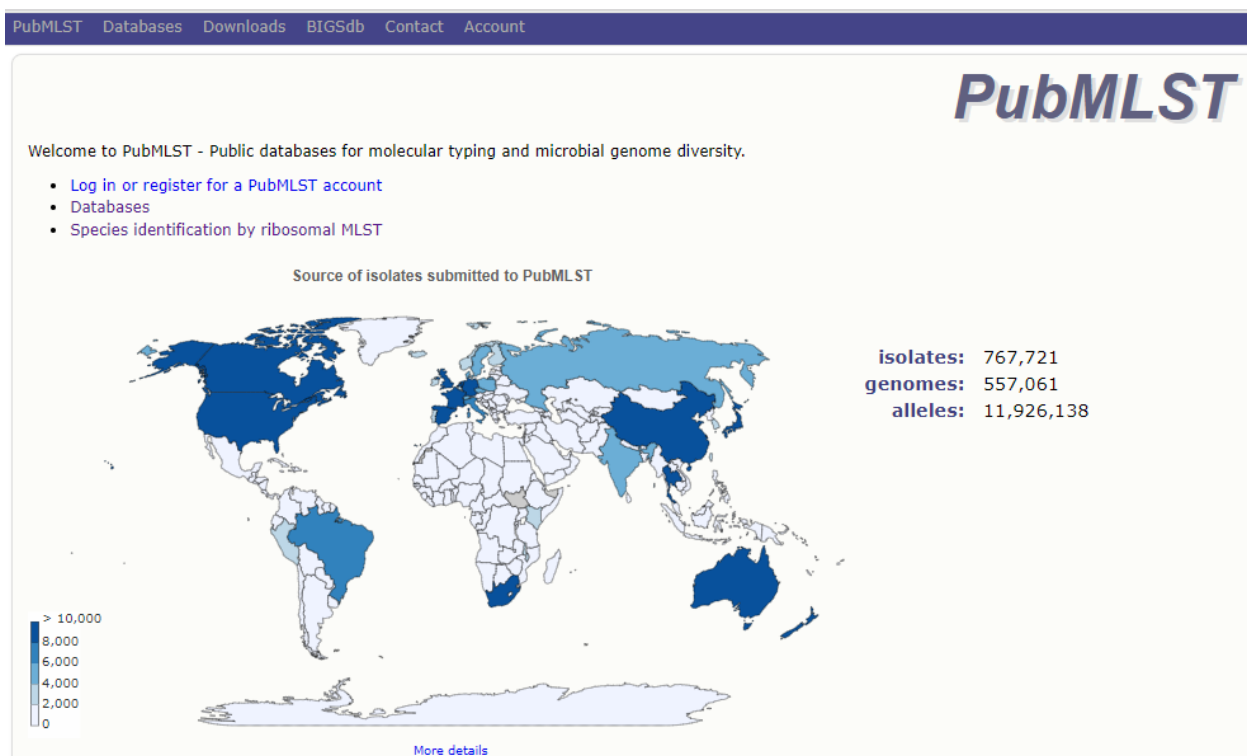
**Figure 2: PubMLST (https://pubmlst.org/) The colour coded geospatial map shows from where the isolates were submitted; the darker the blue the more isolates submitted.**

This data is both large and complex and is intractable to analyse and understand using traditional analysis tools.

## 3.0   MENINGOCCOCAL DISEASE AND EPIDEMIOLOGY

Significant progress has been made in developing effective vaccines against meningococcal disease and treatment in terms of antibiotic therapy and wider healthcare, but invasive disease remains a life-threatening disease. Meningococcal disease is considered to be a medical emergency; it requires immediate action for accurate diagnosis, rapid stabilisation and appropriate antimicrobial treatment [17 and 35]. A patient can present with non-specific symptoms and can progress to death in just a few hours. Permanent disabilities such as brain damage, hearing loss, learning disabilities and loss of limbs can result from meningococcal disease in up to one third of survivors [47]. The infection is spread from person-to-person; human carriers rarely become sick although they can be infectious. The epidemiology of meningococcal disease is highly variable, for reasons that are incompletely understood. It is a major cause of meningitis and sepsis in industrialised settings but the meningitis belt in Africa continues to have the highest burden of disease [40].

Military personnel and recruits are high risk groups for meningococcal disease. It has been reported that the incidence is between four and ten times higher than the in general population [30]. There have been many cases of meningococcal disease outbreaks in the armed forces across different countries, including those of the United Kingdom [29] and the USA [28 and 39]. The living conditions, in particular overcrowding on military bases and training camps, were implicated as far back as 1917 [10].

*Neisseria meningitis* has thirteen significant serogroups, among them, six serogroups, namely A, B, C, W, Y and X are responsible for the majority of cases of disease [14 and 15]. The distribution of each serogroup varies geographically and also changes over time [35 and 40]. Vaccines have been developed targeted at

specific serogroup(s). In view of the increasing use of these vaccines worldwide, high-throughput genomic methods are required to understand variations in capsular type and meningococcal protein diversity [36].

Vaccines have been developed to control the diseases, with development initiated at the Walter Reed Army Institute in the 1960's to protect military personnel [2 and 11]. The US military was involved in the development of the first licenced polysaccharide vaccine for serogroup A in 1970 followed by a combined serogroup A/C vaccine in 1978. Polysaccharide vaccines were later developed against serogroups W and Y, with X in development, and all have prevented countless cases of meningococcal disease globally. Meningococcal serogroup B vaccine was harder to develop as its polysaccharides are similar to the human neuronal proteins [9]. However, through the application of reverse vaccinology, a vaccine was designed utilising other expressed antigens predicted from the genome sequence of a meningococcal B strain (MenB) [8]. In the UK different types of vaccines are offered to different age groups to prevent meningitis and septicaemia; the 6-in-1, MenB, pneumococcal, Hib/Men C vaccines for babies, and meningococcal ACWY conjugate vaccine offered to teenagers and first-time university students. The UK armed forces personnel are offered meningococcal conjugate ACWY vaccine [46]. Vaccines can also be deployed in an outbreak setting, due to the high mortality and morbidity rates, and so the provision of accurate information is crucial in understanding, analysising and communicating an outbreak 'situation' to allow the infection prevention and control measures to be implemented in a timely manner [23].

## 4.0   VISUAL ANALYTICS AND VISUALIZATION

It is often said that a picture is worth a thousand words; in the case of genomic data a picture is worth millions of isolates. Humans have long used visual aids to help them define, understand, analyse and navigate their way through their problems, so as to understand the situation and thus enable informed decision-making. Visual analytics is the science of analytics reasoning supported by interactive visualizations [42]. The analytics approaches include statistical analysis, knowledge discovery, data management and knowledge representation [21]. Visual analytics and visualization are widely used to enable the exploration and analysis of the genomic data to gain insight of the situation. Visualization transforms inherently non-visual data into a form that supports efficient exploration, analysis, discovery, understanding and communication of large volumes of complex data. Visualization can be a static or interactive presentation of data that reinforces human cognition. Interactive visualization enables dynamic exploration and analysis of huge amounts of data that support users not only to detect expected patterns, trends, or correlations but also to discover unexpected associations more efficiently than using tradition approaches. When given a dataset, there are many different ways to represent the data and encode the various variables. However, certain displays are more effective than others for a given task, user, data and / or user, due to the different manner in which they exploit the underlying human perception and cognition [1, 21, 24, 37, 38, 4, 44 and 45]. In short, disease situational awareness encompasses the human perception and cognition skills as well as the processing of data by the machine. In such a complex and dynamic environment, acute situational awareness is likely to greatly enhance the rate and the quality of the decision-making process and infection management [6 and 7].

A wide range of visualization approaches have been developed and applied to support visual analytics of genomic data [33, 34 and 48]. Commonly, the large amount of data is aggregated to provide an overview of the situation and more detailed information can be 'drilled down' into for further analysis.

Figure 3 shows an exploration of *Neisseria menigitidis* for the period 2011-2017 in the United Kingdom using an integrated dashboard. In the integrated dashboard, each display represents a different aspect of the data: clicking on any element in any of the display will result in filtering of all the related elements in the other parts of the dashboard display. This provides an effective means of exploring different aspects of the data and their inter-relations and impacts. The top left hand area shows geographic data and the top right shows the temporal trend of the different serogroups and one non-group (NG), namely A, B, C, E,

W, W/ Y, X, Y and Z, these are colored coded so as to facilitate identification of the serogroup. The lower left stacked bar chart shows the temporal patterns of the clonal complexes. The lower right heatmap shows the nested relationship between clonal complex and serogroup. It can be seen that there is a noticeable decrease in serogroup B and a noticeable increase in serogroup W.
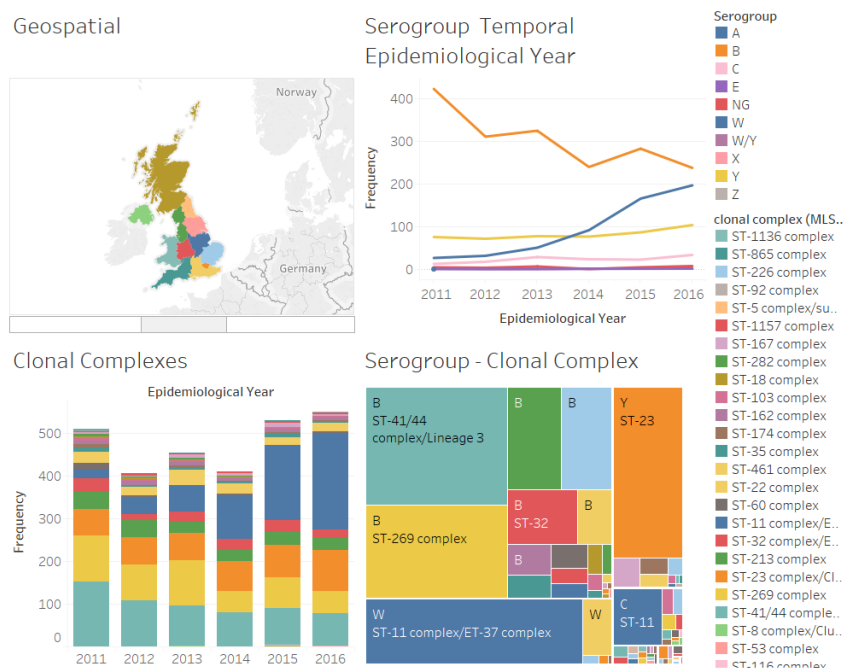


**Figure 3: Geo-temporal Analysis of the *Neisseria Menigitidis***

Figure 4 shows 3,506 invasive meningococcal disease isolates from 2010/11 to 2016/17 analysed using the 24 loci OMVT (outer membrane vesicle typing) scheme and visualised using GrapeTree software. The OMVT cluster by clonal complex (cc), represented by different colours, with unfilled nodes representing isolates from the other ccs. The size of the circles is proportional to the number of isolates. [36].
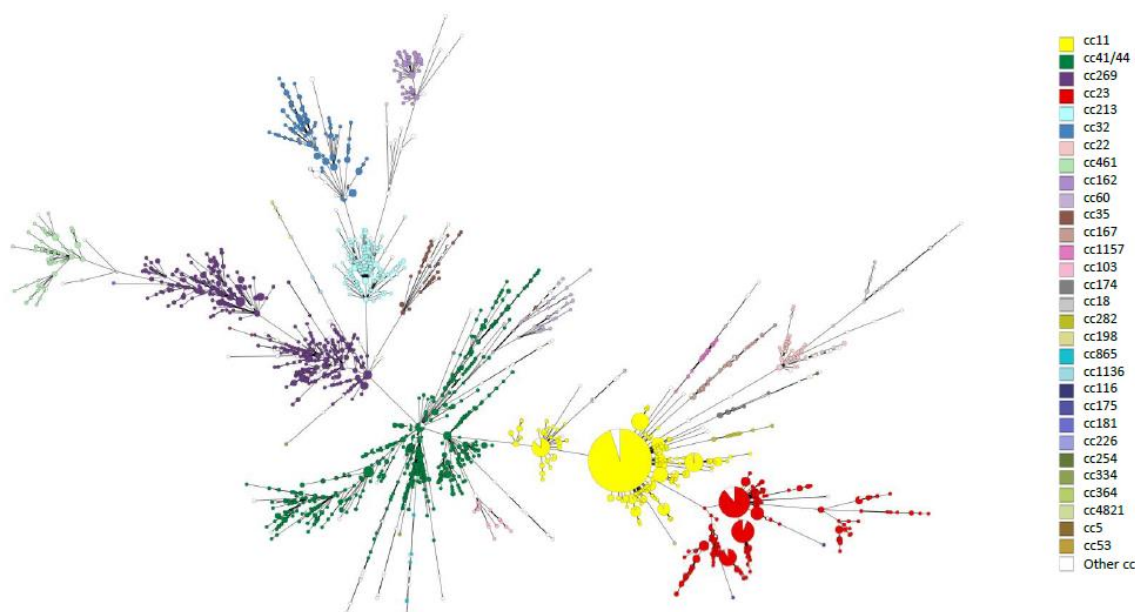


**Figure 4: UK meningococcal disease isolates clustered by OMVT [36]**

The above examples show that the transformation of the non-visual data into a visual form enables effective analysis of large and complex microbial genomic data in an intuitive and understandable manner.

## 5.0  SUMMARY

The advances in sequencing technologies have resulted in an explosion of genomic data.  Genomic research is vital in the fight against infections: for prevention, surveillance and treatment strategies.  The study of microbial genomes has allowed the discovery of anomalies and patterns that have public health implications. Genomic and visual analytics tools can support and enhance our understanding of the biology of infections and lead to new techniques or approaches to diagnosing and treating disease.  Visual analytics and visualization of genomic data are now widely used to support the analysis, presentation and understanding of these invaluable, big and complex data, as well as in the communication of the findings and discovery of insights.  New technologies and evolving perception and cognitive frameworks enhance the development of better and more contextual and user-centred visual analytics approaches.  Genomic and visual analytics combined with machine learning based AI now offer tremendous new opportunities.

## 6.0 REFERENCES

[1] Albers, M. J., Human–Information interaction with complex information for decision-making. Informatics, 2(2):4-19, 2015.

[2] Artenstein, M. S., Control of Meningococcal Meningitis with Meningococcal Vaccines. Yale Journal of Biology and Medicine 48(3): 197-200, 1975.

[3] Barry, J., M., The site of origin of the 1918 influenza pandemic and its public health implications, Journal of Translational Medicine, volume 2, Article number: 3, 2004.

[4] Brundage J. F. and Shanks, G. D., What really happened during the 1918 influenza pandemic? The importance of bacterial secondary infections. The Journal of Infectious Diseases. 196 (11): 1717–8, December 2007.  Author reply 1718–9. doi:10.1086/522355. PMID 18008258.

[5] Byerly, C. R., The U.S. Military and the Influenza Pandemic of 1918 – 1919, Public Health Reports, 2919 Supplement 3/ Volume 125, 2010.

[6] Endsley, M. R., Toward a theory of situation awareness in dynamic systems. Human Factors, 37(1): 32 - 64, March, 1995.

[7] Endsley, M. R.  and Jones, D. G., Designing for Situation Awareness: An Approach to User-Centered Design, Second Edition, CRC Press, 2004, ISBN 9781420063554.

[8] Findlow, J., Meningococcal group B vaccines. Human Vaccines & Immunotherapeutics, 9(6):1387–8. doi:10.4161/hv.24689 45, 2013.

[9] Finne, J., Leinonen, M. and Makela, P. H., Antigenic similarities between brain components and bacteria causing meningitis. Implications for vaccine development and pathogenesis. Lancet 2(8346): 355-357, 1983.

[10] Glover, J. A., Observations on the Meningococcus Carrier-Rate in relation to density of population in Sleeping Quarters. The Journal of Hygiene 17(4): 367-379, 1918.

[11] Gotschlich, E. C., Goldschneider, I. and Artenstein, M. S., Human immunity to the meningococcus IV. Immunogenicity of group A and group C meningococcal polysaccharides. Journal of Experimental Medicine 129: 1367-1384, 1969.

[12] Gray, G. C., Feighner, B., Trump, D. H, Berg S.W., Zajdowicz, M.J. and Zajdowicz, T. R., Disease spread by close personal contact. In: Lenhart MK, Lounsbury D. E, editors. Military preventive medicine: Mobilization and deployment, Volume 2. Washington DC: Borden Institute; 2005. p. 1127–8.

[13] Hall, M. W., Inflammatory diseases of the respiratory tract (bronchitis, influenza, bronchopneumonia, lobar pneumonia). In: Siler JF, ed. The Medical Department of the United States Army in the World War. Volume IX: communicable and other diseases. Washington, DC: US Government Printing Office, 1928: 138. http://www.ibiblio.org/hyperwar/AMH/XX/WWI/Army/Medical/IX/USA-Med-IX-2.html (accessed 3/10/2019).

[14] Harrison L. H, Trotter, C. L. and Ramsay, M. E., Global epidemiology of meningococcal disease.

Vaccine. 2009; 27 Supplement 2: B51-63.

[15] Harrison, O. B., Claus, H., Jiang, Y., Bennett, J. S., Bratcher, H. B., Jolley, K. A., Corton, C., Care, R., et al. (2013). Description and Nomenclature of Neisseria meningitidis Capsule Locus. Emerging Infectious Diseases 19(4): 566-573.

[16] Heer, J. and Shneiderman, B., Interactive Dynamics for Visual Analysis', ACM Queue 10(2), pp 1 - 30, February 2012.

[17] Jackson L. A., Schuchat, A., Reeves. M. W., Wenger J. D. Serogroup C meningococcal outbreaks in the United States. An emerging threat. JAMA. 1995 Feb 1;273(5):383-9.

[18] Jester, B., Uyeki, T. M., Jernigan, D. B. and Tumpey T. M., Historical and clinical aspects of the 1918 H1N1 pandemic in the United States. Virology. 2019 Jan 15;527:32-37. doi: 10.1016/j.virol.2018.10.019.

[19] Jolley, K. A. and Maiden, M. C., BIGSdb: Scalable analysis of bacterial genome variation at the population level. BMC Bioinformatics 11(1): 595, 2010.

[20] Johnson N P. A. S. and Mueller, Updating the accounts: global mortality of the 1918- 1920 "Spanish" influenza pandemic. Bulletin of the history of medicine, 2002; 6:105-15.

[21] Keim, D., Andrienko, G. and Fekete, J. D. & Görg, C., Kohlhammer, J. & Melançon, G., Visual Analytics: Definition, Process, and Challenges. 10.1007/978-3-540-70956-5_7., 2008.

[22] Koppes, G.M., Ellenbogen, C. and Gebhart, R. J., Group Y meningococcal disease in United States Air Force recruits. Am J Med. 1977; 62(5):661–6. [PubMed: 404877].

[23] Leca, M., Bornet, C., Montana, M., Curti, C., Vanelle, P. Meningococcal vaccines: current state and future outlook. Pathologie Biologie (Paris) (2015) 63(3):144–51. doi:10.1016/j.patbio.2015.04.003

[24] Lurie, N. H., Mason, C.H., Visual representation: Implications for decision making. Journal of Marketing 2007, 71(1):160- 177.

[25] Maiden, M. C. J., Bygraves, J. A., Feil, E., Morelli, G., Russell, J. E., Urwin, R., Zhang, Q., Zhou, J., Zurth, K., Caugant, D. A., Feavers, I. M., Achtman, M. and Spratt, B. G. (1998). Multilocus sequence typing: a portable approach to the identification of clones within populations of pathogenic microorganisms. Proceedings of the National Academy of Sciences of the United States of America 95(6): 3140-3145.

[26] Maiden, M. C., Jansen van Rensburg, M. J., Bray, J. E., Earle, S. G., Ford, S. A., Jolley, K. A. and McCarthy, N. D., MLST revisited: the gene-by-gene approach to bacterial genomics. Nature Reviews Microbiology 11(10): 728-736, 2013.

[27] Maiden, M. C., The impact of protein-conjugate polysaccharide vaccines: an endgame for meningitis? Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences 368(1623): 20120147, 2013.

[28] Makras, P., Alexiou-Daniel, S., Antoniadis, A. and Hatzigeorgiou, D., Outbreak of meningococcal disease after an influenza B epidemic at a Hellenic Air Force recruit training center. Clin Infect Dis. 2001;33(6):e48–50. doi: 10.1086/322609. [PubMed: 11512107].

[29] Masterton, R. G., Youngs, E. R., Wardle, J. C., Croft, K. F. and Jones, D. M. Control of an outbreak of group C meningococcal meningitis with a polysaccharide vaccine. J Infect. 1988;17(2):177–82. [PubMed: 3141518].

[30] Mehrdadi, S., Acute Bacterial Meningitis: Diagnosis, Treatment and Prevention, Journal of Archives in Military Medicine, 6(4):e84749, December 2018.

[31] Morens, D. M., Fauci, A. S., The 1918 influenza pandemic: insights for the 21st century. The Journal of Infectious Diseases. 195 (7): 1018–28. April 2007. doi:10.1086/511989. PMID 17330793.

[32] Mougel, N., CVCE, 2011, 2011, English translation, Gratz, J., Centre European Roobert Schuman, Reperes – module 1-0 – explanatory notes – World War I casualties – EN.

[33] Overmars, L., Kerkhoven R., Siezen R. J., Francke, C.: MGcV: the microbial genomic context viewer for comparative genome analysis. BMC Genomics 14 (Apr. 2013), 209. doi: 10.1186/1471-2164-14-209.

[34] Ramírez F, Ryan DP, Grüning B, Bhardwaj V, Kilpert F, Richter AS, Heyne S, Dündar F, Manke T. deepTools2: a next generation web server for deep-sequencing data analysis. Nucleic Acids Research. 2016 Apr 13.

[35] Rosenstein N. E., Perkins B. A., Stephens D. S., et al. Meningococcal disease. N Engl J Med. 2001; 344(18): 1378-88.

[36] Rodrigues, C. M. C., Genomic Epidemiology of Meningococcal Vaccine Antigens in the United Kingdom, PhD thesis, December 2018.

[37] Sacha, D., Stoffel, A., Stoffel, F., Kwon, B. C., Ellis, G. and Keim, D. A.: Knowledge generation model for visual analytics. Visualization and Computer Graphics, IEEE Transactions on 2014, 20(12):1604-1613.

[38] Sedig, K. and Parsons, P., Design of visualizations for human information interaction: A pattern-based framework. Synthesis Lectures on Visualization 2016, 4(1):1-185.

[39] Smilack, J. D., Group-y meningococcal disease. Twelve cases at an army training center. Ann Intern Med. 1974;81(6):740–5.[ doi: 10.7326/0003-4819-81-6-740] [PubMed: 4215351]

[40] Sridhar, S., Greenwood, B., Head, C., Plotkin, S. A., Safadi, M. A., Saha, S., Taha, M. K., Tomori, O. and Gessner, B. D. (2015). Global incidence of serogroup B invasive meningococcal disease: a systematic review. Lancet Infectious Diseases 15(11): 1334-1346.

[41] Taubenberger, J. K. and Morens, D. M., 1918 Influenza: the Mother of All Pandemics, Emerging Infectious Diseases, Volume. 12, No. 1, pp 15 -22, January 2006.

[42] Thomas, J. J. and Cook, K. A. (Eds.), Illuminating the Path: The Research and Development Agenda for Visual Analytics, National Visualization and Analytics Center, 2005.

[43] Tory, M., and Moller, T., Human Factors In Visualization Research, IEEE Transactions on visualization and Computer Graphics, Volume 10, No. 1, January/February 2004.

[44] Tufte, E. The Visual Display of Quantitative Information, Graphic Press, Second Edition, ISBN-13: 978- 0961392147, May 2001.

[45] Ware, C., Information Visualization: Perception for Design, 3rd Edition. Morgan Kaufmann Publishers, Inc., San Francisco, California, 2012.

[46] Vaccination in the armed forces – what you need to know, Annex B, TO AFC(H) JIs Dated SEP/OCT 17.

[47] Viner, R. M., Booy, R., Johnson, H., Edmunds, W. J., Hudson, L., Bedford, H., Kaczmarski, E., Rajput, K., et al. (2012). Outcomes of invasive meningococcal serogroup B disease in children and adolescents (MOSAIC): a case-control study. Lancet Neurology 11(9): 774-783.

[48] Yachdav, G., Wilzbach, S., Rauscher, B., Sheridan R., Sillitoe, I., Procter, J., Lewis, S. E., Rost, B., Goldberg, T.: MSAViewer: interactive JavaScript visualization of multiple sequence alignments. Bioinformatics 32, 22 (Nov. 2016), 3501–3503.