**Title:** *NGMASTER: in silico* Multi-Antigen Sequence Typing for *Neisseria gonorrhoeae*

**Running title:** *NGMASTER: in silico* NG-MAST for *Neisseria gonorrhoeae*

**Authors:** Jason C Kwong[1,2,3], Anders Gonçalves da Silva[1,4], Kristin Dyet[5], Deborah A Williamson[1,4], Timothy P Stinear[1,2], Benjamin P Howden[1,2,3,4], Torsten Seemann[1,6]

**Affiliations:**
1. Doherty Applied Microbial Genomics, Peter Doherty Institute for Infection & Immunity, Melbourne, Australia
2. Department of Microbiology & Immunology, University of Melbourne, Parkville, Australia
3. Department of Infectious Diseases, Austin Health, Heidelberg, Australia
4. Microbiological Diagnostic Unit Public Health Laboratory, Peter Doherty Institute for Infection & Immunity, Melbourne, Australia
5. Institute of Environmental Science and Research, Wellington, New Zealand
6. Victorian Life Sciences Computation Initiative, Carlton, Australia

**Keywords:**
Neisseria gonorrhoeae
Multi-antigen sequence typing
NG-MAST
Whole-genome sequencing
In silico typing

**Correspondence:**
Dr Jason Kwong
Doherty Applied Microbial Genomics
Peter Doherty Institute for Infection & Immunity
792 Elizabeth Street
Melbourne, Victoria
Australia 3000
Tel: +61 3 8344 5701
Email: jason.kwong (AT) unimelb.edu.au

# NGMASTER – *in silico* Multi-Antigen Sequence Typing for *Neisseria gonorrhoeae*

## ABSTRACT

Whole-genome sequencing (WGS) provides the highest resolution analysis for comparison of bacterial isolates in public health microbiology. However, although increasingly being used routinely for some pathogens such as *Listeria monocytogenes* and *Salmonella enterica*, the use of WGS is still limited for other organisms, such as *Neisseria gonorrhoeae*. Multi-antigen sequence typing (NG-MAST) is the most widely performed typing method for epidemiologic surveillance of gonorrhoea. Here, we present *NGMASTER* – a command-line software tool for performing *in silico* NG-MAST on assembled genome data. *NGMASTER* rapidly and accurately determined the NG-MAST of 630 assembled genomes, facilitating comparisons between WGS and previously published gonorrhoea epidemiological studies. The source code and user documentation are available at https://github.com/MDU-PHL/ngmaster.

## DATA SUMMARY

1. The Python source code for *NGMASTER* is available from GitHub under GNU GPL v2. (URL: https://github.com/MDU-PHL/ngmaster)
2. The software is installable via the Python "pip" package management system. Install using "pip install --user git+https://github.com/MDU-PHL/ngmaster.git"
3. Sequencing data used are available for download from the EBI European Nucleotide Archive under BioProject accessions PRJEB2999, PRJNA29335, PRJNA266539, PRJNA298332, and PRJEB14168.

**We confirm all supporting data, code and protocols have been provided within the article or through supplementary data files. ☑**

34

35 ——————————————————————————————

## IMPACT STATEMENT

37

Whole-genome sequencing (WGS) offers the potential for high-resolution comparative analyses of microbial pathogens. However, there remains a need for backward compatibility with previous molecular typing methods to place genomic studies in context. NG-MAST is currently the most widely used method for epidemiologic surveillance of *Neisseria gonorrhoeae*. We present *NGMASTER*, a command-line software tool for performing Multi-Antigen Sequence Typing (NG-MAST) of *Neisseria gonorrhoeae* from WGS data. This tool is targeted at clinical and research microbiology laboratories that have performed WGS of *N. gonorrhoeae* isolates and wish to understand the molecular context of their data in comparison to previously published epidemiological studies. As WGS becomes more routinely performed, *NGMASTER* was developed to completely replace PCR-based NG-MAST, reducing time and labour costs.

49

50

51 ——————————————————————————————

## INTRODUCTION

53

*Neisseria gonorrhoeae* is one of the most common sexually transmitted bacterial infections worldwide. There is growing concern about the global spread of resistant epidemic clones, with extensively drug-resistant gonorrhoea being listed as an urgent antimicrobial resistance threat (CDC, 2013; WHO, 2014).

58

Multi-Antigen Sequence Typing of *N. gonorrhoeae* (NG-MAST) has been important in tracking these resistant clones, such as the NG-MAST 1407 clone associated with decreased susceptibility to third-generation cephalosporins (Unemo & Dillon, 2011). It involves sequence-based typing using established PCR primers of two highly variable and polymorphic outer membrane protein genes, *porB* and *tbpB* by comparing the sequences to an open-access database (http://www.ng-mast.net/) (Martin *et al.*, 2004). Although NG-MAST is the most frequently performed molecular typing method for *N. gonorrhoeae*, it requires multiple PCR amplification and sequencing reactions, making it more laborious than other typing methods (Heymans *et al.*, 2012).

68

69   Whole-genome sequencing (WGS) is increasingly being used for molecular typing and
70   epidemiologic investigation of microbial pathogens as it provides considerably higher
71   resolution. A number of studies using genomic data to understand the epidemiology of
72   *N. gonorrhoeae* have already been published (Grad *et al.*, 2014) (Demczuk *et al.*, 2015)
73   (Ezewudo *et al.*, 2015) (Demczuk *et al.*, 2016). However, the ability to perform retrospective
74   comparisons with previous epidemiological studies is reliant on conducting both traditional
75   typing (such as NG-MAST) as well as more modern WGS analyses on the same isolates.
76
77   *NGMASTER* is a command-line software tool for rapidly determining NG-MAST types *in*
78   *silico* from genome assemblies of *N. gonorrhoeae*.
79
80
81   ─────────────────────────────────────────────────────────

## DESCRIPTION

83
84   *NGMASTER* is an open source tool written in Python and released under a GPLv2 Licence.
85   The source code can be downloaded from Github (https://github.com/MDU-PHL/ngmaster).
86   It has two software dependencies: *isPcr* (http://hgwdev.cse.ucsc.edu/~kent/src/) and
87   BioPython (Cock *et al.*, 2009), and uses the allele databases publicly available at
88   http://www.ng-mast.net/, which *NGMASTER* can automatically download and update locally
89   for running.
90
91   *NGMASTER* is based on the laboratory method published by Martin *et al.* (Martin *et al.*,
92   2004), and uses *isPcr* to retrieve allele sequences from a user-specified genome assembly
93   in FASTA format by locating the flanking primers. These allele sequences are trimmed to a
94   set length from starting key motifs in conserved gene regions, and then checked against the
95   allele databases. Results are printed in machine readable tab- or comma-separated format.
96
97
98   ─────────────────────────────────────────────────────────

## METHODS AND RESULTS

100
101   *NGMASTER* was validated against 630 publicly available *N. gonorrhoeae* genome
102   sequences derived from published studies (Table 1). This included 8 well characterised
103   WHO reference genomes with published data and 50 local isolates that had undergone
104   "traditional" NG-MAST by PCR and Sanger sequencing (Martin *et al.*, 2004). A further 572

105    isolates that had undergone manual *in silico* NG-MAST from WGS data (Demczuk *et al.*,
106    2015; Demczuk *et al.*, 2016; Grad *et al.*, 2014), including the fully assembled reference
107    genome NCCP11945, were also tested. Raw WGS data for these sequences were retrieved
108    from the European Nucleotide Archive (ENA). Average sequencing depth was >30x for all
109    ENA sequences, with a combination of 100 bp, 250 bp and 300 bp paired-end Illumina
110    reads. Local isolates also underwent WGS on the Illumina MiSeq/NextSeq using Nextera
111    libraries and manufacturer protocols, with an average sequencing depth >50x. The raw
112    sequencing reads for these local isolates have been uploaded to the ENA (BioProject
113    accession PRJEB14168).

114

115    Sequencing reads were trimmed to clip Illumina adapters and low-quality sequence
116    (minimum Q20) using *Trimmomatic* v0.35 (Bolger *et al.*, 2014). Draft genomes were
117    assembled *de novo* with *MEGAHIT* v1.0.3 and *SPAdes* v3.7.1 (Li *et al.*, 2015) (Bankevich *et*
118    *al.*, 2012) to investigate whether the faster, but approximate genome assembler, *MEGAHIT*,
119    would be sufficient for *NGMASTER*. A list of the commands and parameters used is
120    included in the Appendix 1.

121

122    The *de novo* assembled draft genomes and the fully assembled NCCP11945 reference
123    genome in FASTA format were used as input to *NGMASTER* with the overall results shown
124    in Table 1. Complete *NGMASTER* results with sequencing and assembly metrics are
125    included in Appendix 2. Running *NGMASTER* on 630 genome assemblies using a single
126    Intel(R) Xeon(R) 2.3GHz CPU core was completed in less than two minutes.

127

128    Overall, *NGMASTER* assigned NG-MAST types that were concordant with published results
129    for 93-97% of the tested *N. gonorrhoeae* genomes using *MEGAHIT* or *SPAdes* assemblies.
130    Notably, comparisons with results from traditional NG-MAST were 100% concordant (57/57).
131    Reasons for discordant results are shown in Table 2. In general, running *NGMASTER* using
132    *SPAdes* assemblies resolved more NG-MAST types than when using *MEGAHIT*
133    assemblies. However, 10 genomes assembled with *SPAdes* were found to have assembly
134    errors in either *por* or *tbpB* introduced in the repeat resolution stage, resulting in discordant
135    NG-MAST types for those isolates (major errors). Running *NGMASTER* on preliminary
136    contigs prior to this process ("before_rr.fasta") alleviated these major errors, and were
137    concordant with *MEGAHIT* results and the published results (Appendix 2). In contrast, minor
138    errors (due to incomplete NG-MAST types or multiple alleles detected) were more frequent
139    using *MEGAHIT* assemblies, particularly those with poor assembly metrics (e.g. >500
140    contigs, N50 <10 kbp). When *MEGAHIT* assemblies successfully produced complete

141 *NGMASTER* results, these NG-MAST types were highly concordant with the published

142 results.

143

144 To overcome this issue, a two-stage assembly approach was also tested, where a draft

145 genome was first assembled using *MEGAHIT* for initial testing. If a complete NG-MAST

146 result was obtained, this was recorded as the final result for that isolate. If the result was

147 incomplete or suggested multiple alleles were present, the genome was also assembled

148 using *SPAdes*. Using this combined approach, 620/630 (98%) NG-MAST types derived from

149 *NGMASTER* were concordant with the published results, with only 42 genomes requiring

150 additional assembly with the slower, but more thorough *SPAdes* assembler.

151

152 For the remaining 10 discordant results, seven of these were likely due to errors in the

153 published data, including for NCCP11945. A further two isolates were found to have multiple

154 *tbpB* alleles in both *SPAdes* and *MEGAHIT* assemblies, with the dominant allele (indicated

155 by higher read coverage and better flanking assembly) matching the published result. The

156 *tbpB* allele for the final isolate was not able to be determined by *NGMASTER* due to a

157 mutation in the conserved starting key motif required for sequence trimming to a standard

158 size.

159

160

161 _____

## ISSUES WITH IMPLEMENTATION

163

164 The NG-MAST procedure involves sequencing the internal regions of *por* and *tbpB* that

165 encode two variable outer membrane proteins. The sequences are trimmed to a standard

166 length from a starting key motif in conserved regions of each gene. However, despite being

167 relatively conserved, a number of variations of this starting motif appear in the NG-MAST

168 database (Fig. 1), causing one discordant result (Table 2). Some sequences appeared to

169 lack a *tbpB* gene due to the presence of non-typeable *tbpB* genes acquired from

170 *N. meningitidis*, though this was also noted in the published data. Another source of

171 discordant results were genomes that appeared to have multiple alleles, suggesting isolate

172 contamination or polyclonal infection.

173

174 A number of isolates were found to have novel alleles or allele combinations that were not in

175 the most recent version of the database available at http://www.ng-mast.net. For

176 convenience, *NGMASTER* includes an option to save these allele sequences in FASTA

177 format for manual submission to the database and allele type assignment.

178

179 Notably, results were dependent on the accuracy and quality of the *de novo* draft genome

180 assembly. It should be noted that for this study, draft genomes were assembled *de novo*

181 using relatively standard parameters for *MEGAHIT* and *SPAdes* without post-assembly error

182 checking (see Appendix 1). We were alerted to the presence of *SPAdes* assembly errors

183 after finding the corresponding *MEGAHIT* assemblies produced different *NGMASTER*

184 results. Concordant results were able to be obtained for each of these genomes after

185 identifying and correcting assembly errors through re-mapping each isolate's reads back to

186 the respective draft *SPAdes* assembly. Results from running *NGMASTER* on the *SPAdes*

187 interim "before_rr.fasta" contigs also produced concordant results. Assuming accurate

188 closed genome assemblies are used with an accurate and well curated database, based on

189 our testing, we anticipate that *NGMASTER* would produce NG-MAST results that were

190 >99% if not 100% accurate.

191

192

193 _____

194 # CONCLUSION

195

196 *NGMASTER* rapidly and accurately performs *in silico* NG-MAST typing of *N. gonorrhoeae*

197 from assembled WGS data, and may be a useful command-line tool to help contextualise

198 genomic epidemiological studies of *N. gonorrhoeae*.

199

200

201

202

# REFERENCES

**Bankevich, A., Nurk, S., Antipov, D., Gurevich, A. A., Dvorkin, M., Kulikov, A. S., Lesin, V. M., Nikolenko, S. I., Pham, S., Prjibelski, A. D., Pyshkin, A. V., Sirotkin, A. V., Vyahhi, N., Tesler, G., Alekseyev, M. A. & Pevzner, P. A. (2012).** SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *J Comput Biol* **19**, 455-477.

**Bolger, A. M., Lohse, M. & Usadel, B. (2014).** Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* **30**, 2114-2120.

**CDC (2013).** Antibiotic Resistance Threats in the United States, 2013: Centers for Disease Control and Prevention, US Department of Health and Human Services.

**Cock, P. J., Antao, T., Chang, J. T., Chapman, B. A., Cox, C. J., Dalke, A., Friedberg, I., Hamelryck, T., Kauff, F., Wilczynski, B. & de Hoon, M. J. (2009).** Biopython: freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics* **25**, 1422-1423.

**Demczuk, W., Lynch, T., Martin, I., Van Domselaar, G., Graham, M., Bharat, A., Allen, V., Hoang, L., Lefebvre, B., Tyrrell, G., Horsman, G., Haldane, D., Garceau, R., Wylie, J., Wong, T. & Mulvey, M. R. (2015).** Whole-genome phylogenomic heterogeneity of Neisseria gonorrhoeae isolates with decreased cephalosporin susceptibility collected in Canada between 1989 and 2013. *J Clin Microbiol* **53**, 191-200.

**Demczuk, W., Martin, I., Peterson, S., Bharat, A., Van Domselaar, G., Graham, M., Lefebvre, B., Allen, V., Hoang, L., Tyrrell, G., Horsman, G., Wylie, J., Haldane, D., Archibald, C., Wong, T., Unemo, M. & Mulvey, M. R. (2016).** Genomic Epidemiology and Molecular Resistance Mechanisms of Azithromycin-Resistant Neisseria gonorrhoeae in Canada from 1997 to 2014. *J Clin Microbiol* **54**, 1304-1313.

**Ezewudo, M. N., Joseph, S. J., Castillo-Ramirez, S., Dean, D., Del Rio, C., Didelot, X., Dillon, J. A., Selden, R. F., Shafer, W. M., Turingan, R. S., Unemo, M. & Read, T. D. (2015).** Population structure of Neisseria gonorrhoeae based on whole genome data and its relationship with antibiotic resistance. *PeerJ* **3**, e806.

**Grad, Y. H., Kirkcaldy, R. D., Trees, D., Dordel, J., Harris, S. R., Goldstein, E., Weinstock, H., Parkhill, J., Hanage, W. P., Bentley, S. & Lipsitch, M. (2014).** Genomic epidemiology of Neisseria gonorrhoeae with reduced susceptibility to cefixime in the USA: a retrospective observational study. *Lancet Infect Dis* **14**, 220-226.

239   **Heymans, R., Golparian, D., Bruisten, S. M., Schouls, L. M. & Unemo, M. (2012).**

240       Evaluation of Neisseria gonorrhoeae multiple-locus variable-number tandem-repeat

241       analysis, N. gonorrhoeae Multiantigen sequence typing, and full-length porB gene

242       sequence analysis for molecular epidemiological typing. *J Clin Microbiol* **50**, 180-183.

243   **Li, D., Liu, C. M., Luo, R., Sadakane, K. & Lam, T. W. (2015).** MEGAHIT: an ultra-fast

244       single-node solution for large and complex metagenomics assembly via succinct de

245       Bruijn graph. *Bioinformatics* **31**, 1674-1676.

246   **Martin, I. M., Ison, C. A., Aanensen, D. M., Fenton, K. A. & Spratt, B. G. (2004).** Rapid

247       sequence-based identification of gonococcal transmission clusters in a large

248       metropolitan area. *J Infect Dis* **189**, 1497-1505.

249   **Unemo, M. & Dillon, J. A. (2011).** Review and international recommendation of methods for

250       typing neisseria gonorrhoeae isolates and their implications for improved knowledge

251       of gonococcal epidemiology, treatment, and biology. *Clin Microbiol Rev* **24**, 447-458.

252   **WHO (2014).** Antimicrobial resistance: global report on surveillance 2014: World Health

253       Organization.

254

255

# FIGURES AND TABLES

257

258    **Table 1:** Concordance between *NGMASTER* results from draft genome assemblies using

259    *MEGAHIT* and *SPAdes*, and previously published NG-MAST results.

260

|  | *MEGAHIT* | *SPAdes* | 2-stage [$] | TOTAL |
|---|---|---|---|---|
| PRJEB2999 [†] | 176 (95%) | 184 (99%) | 184 (99%) | 186 |
| PRJNA29335 [#][*] | - | - | - | 1 |
| PRJNA266539 [#] | 162 (91%) | 169 (94%) | 178 (99%) | 179 |
| PRJNA298332 [§] | 199 (93%) | 207 (97%) | 208 (97%) | 214 |
| PRJEB14168 [‡] | 50 (100%) | 50 (100%) | 50 (100%) | 50 |
| TOTAL | 587 (93%) | 610 (97%) | 620 (98%) | 630 |

† Grad *et al.*, Lancet Infect Dis 2014

# Demczuk *et al.*, J Clin Microbiol 2015

§ Demczuk *et al.*, J Clin Microbiol 2016

* Closed reference genome NCCP11945 (Genbank accession CP001050.1) - *in silico* NG-MAST results reported by Demczuk et al. (Demczuk *et al.*, 2015)

‡ Local isolates with NG-MAST performed by PCR/Sanger sequencing

$ 2-stage assembly: 1. NGMASTER run using rapid assembly with *MEGAHIT*; 2. *NGMASTER* also run using *SPAdes* if no result or mixed result using *MEGAHIT* assembly

261

262

263

264    **Table 2:** Reasons for discordant results between *NGMASTER* and published data using

265    *SPAdes* assemblies

266

| Reason for discordant result | MEGAHIT | SPAdes |
|---|---|---|
| *Major errors (incorrect result)* | | |
| Assembly error | 0 | 10 |
| *Minor errors (incomplete/missing result)* | | |
| Alternate conserved key motif | 1 | 1 |
| Multiple alleles detected | 6 | 2 |
| Allele not detected | 29 | 0 |
| *Errors in published data* | | |
| Possible sequence mix-up in published data | 4 | 4 |
| Probable transcription error in published data | 1 | 1 |
| Error in published data | 1 | 1 |

267

268

269

270

271

272 **Figure 1:** Number and frequency of alternate starting key motifs within "conserved" gene

273 regions for trimming allele sequences.

274



**por conserved start motif**

| Motif | Count |
|---|---|
| TTGAA | ▬ 7789 |
| TTTGA | 8 |
| GATTT | 8 |
| TTGGA | 5 |
| No sequence | 3 |
| TTGAG | 2 |
| TGAAG | 2 |
| CTTGA | 2 |
| TTTTA | 1 |
| TTTAA | 1 |
| TTAAA | 1 |
| GTGCC | 1 |
| GAATT | 1 |
| CTTAA | 1 |
| CCTTA | 1 |
| ATGAA | 1 |

**tbpB conserved start motif**

| Motif | Count |
|---|---|
| CGTCTGAA | ▬ 2139 |
| CGTCTGGA | ▪ 64 |
| CGTCGTCT | 13 |
| CCGTCTGA | 10 |
| TTGAATTG | 9 |
| CGTCTGCA | 8 |
| CGTCTAAA | 7 |
| TGCCGTCT | 4 |
| CTGAAAAC | 4 |
| TCTGAAAC | 3 |
| CTGGAAAC | 3 |
| CGTTTGAA | 3 |
| CGTTGTCG | 3 |
| CGTGTCAG | 3 |
| CGGTTGAA | 3 |
| TGTCTGAA | 2 |
| TCTGAAAA | 2 |
| TCGTCTGA | 2 |
| CGGTTGTC | 2 |
| No sequence | 2 |
| TTTCCGCA | 1 |
| TTGAAGGG | 1 |
| TGCGCCCT | 1 |
| TCGTCTGC | 1 |
| TCGCCTTG | 1 |
| GTCTGGAA | 1 |
| GTCTGAAG | 1 |
| GGTTGTCG | 1 |
| GAGCTGAA | 1 |
| CTGTCTGA | 1 |
| CTGAAGAC | 1 |
| CGTGTCAA | 1 |
| CGTCTGTC | 1 |
| CGTCTGAT | 1 |
| CGTCTGAC | 1 |
| CGTCTCAA | 1 |
| CGTCAGAA | 1 |
| CGTATCAA | 1 |
| CGGCTTGA | 1 |
| CGCGTCAG | 1 |
| CGCGTCAA | 1 |
| CGCGGGAA | 1 |
| CGCCGTCT | 1 |
| CGAAAACC | 1 |
| CCTCTGAA | 1 |
| CCAGAGAC | 1 |
| AGGTTGGA | 1 |
| AGGAAGCG | 1 |
| ACCGAAAA | 1 |
| ACCAGAGA | 1 |
| AAAACCAA | 1 |
| AAAACACT | 1 |

275