

# Genetic variation regulates the activation and specificity of Restriction-Modification systems in *Neisseria gonorrhoeae*

Leonor Sánchez-Busó<sup>1,\*</sup>, Daniel Golparian<sup>2</sup>, Julian Parkhill<sup>1</sup>, Magnus Unemo<sup>2</sup> and Simon R. Harris<sup>1,3,\*</sup>

<sup>1</sup> Pathogen Genomics, Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridgeshire, CB10 1SA, United Kingdom.

<sup>2</sup> WHO Collaborating Centre for Gonorrhoea and other Sexually Transmitted Infections, National Reference Laboratory for Sexually Transmitted Infections, Department of Laboratory Medicine, Clinical Microbiology, Faculty of Medicine and Health, Örebro University, Örebro, SE-701 85, Sweden.

<sup>3</sup> Microbiotica Ltd, Biodata Innovation Centre, Wellcome Genome Campus, Hinxton, Cambridgeshire, CB10 1DR, UK

\* To whom correspondence should be addressed. Leonor Sánchez-Busó. Tel: +44 1223 494761; Email: [lsb@sanger.ac.uk](mailto:lsb@sanger.ac.uk). Correspondence may also be addressed to Simon R. Harris ([sharris@microbiotica.com](mailto:sharris@microbiotica.com)).

## ABSTRACT

Restriction-Modification systems (RMS) are one of the main mechanisms of defence against foreign DNA invasion and can have an important role in the regulation of gene expression. The obligate human pathogen *Neisseria gonorrhoeae* carries one of the highest loads of RMS in its genome; between 13 to 15 of the three main types. Previous work has described their organization in the reference genome FA1090 and has experimentally inferred the associated methylated motifs. Here, we studied the structure of RMS and target methylated motifs in 25 gonococcal strains sequenced with Single Molecule Real-Time (SMRT) technology, which provides data on DNA modification. The results showed a variable picture of active RMS in different strains, with phase variation switching the activity of Type III RMS, and both the activity and specificity of a Type I RMS. Interestingly, the Dam methylase was found in place of the NgoAXI endonuclease in two of the strains, despite being previously thought to be absent in the gonococcus. We also identified the real methylation target of NgoAX as 5'-GCAGA-3', different from that previously described. Results from this work give further insights into the diversity and dynamics of RMS and methylation patterns in *N. gonorrhoeae*.

# Introduction

*Neisseria gonorrhoeae* is a sexually-transmitted pathogen that causes gonorrhoea. Antimicrobial resistance in this pathogen is of great public health concern and there have recently been an increasing number of publications aimed at analysing its transmission and the evolution of genetic determinants of antimicrobial resistance using high throughput short-read sequencing<sup>1-4</sup>. However, Single Molecule Real-Time (SMRT) PacBio sequencing provides deeper information about microbial genomes, including data on DNA modification, mainly methylation in the form of 6mA, 5mC or 4mC. Several studies have focused on the reference genome *N. gonorrhoeae* FA1090 to characterize the population of Restriction-Modification systems (RMS) in the gonococcus<sup>5-8</sup> and its methylation landscape<sup>9</sup>. However, RMS and methylation status are known to vary significantly between strains of the same species, so this provides only limited knowledge about RMS and methylation in the gonococcus as a whole.

*N. gonorrhoeae* is among the bacterial species with the highest numbers of RMS despite being naturally competent for DNA uptake<sup>10</sup>. Indeed, highly transformable bacteria have been found to contain a higher number of RMS compared with those which are less competent, acting as a defence system against the invasion of foreign DNA<sup>11</sup>. RMS are formed by a restriction endonuclease (REase) and a DNA methyltransferase (MTase) that recognize a specific pattern in the genome. If the MTase is active, it will methylate the associated motifs and thus protect them from cleavage by the cognate REase. Bacteria can contain up to four types of RMS, although only Types I, II and III have been found in *Neisseria*<sup>10,12</sup>. Briefly, Type I RMS consist of three subunits named *hsdM* (MTase), *hsdS* (specificity, S, unit) and *hsdR* (REase)<sup>13</sup> that form a multi-enzyme complex that recognizes specific asymmetric sequences, methylating both strands. Type II RMS are generally formed by individual MTase and REase enzymes that recognize the same short palindromic motifs and perform a two-strand methylation or cleavage, respectively. Finally, Type III RMS forms a complex of Mod (MTase) and Res (REase) subunits<sup>14</sup> that recognize specific short asymmetrical motifs. The specificity of Type I RMS is encoded by an independent gene (*hsdS*) that must be active for the other subunits to be functional, while Type III RMS harbour a DNA recognition domain (DRD) within the MTase<sup>12</sup>. RMS have also been associated with other roles apart from defending the genome against foreign DNA invasion, such as being involved in the epigenetic control of gene expression<sup>11</sup>. They have been described as selfish elements, tending to propagate on mobile genetic elements and promoting their own survival<sup>15</sup>. They are often associated with mobility genes, such as integrases or transposases and are sometimes flanked by repeats<sup>16</sup>. Indeed, they have been shown to

have an important role in genetic flux among bacteria, as cleavage by REases provides fragments of double-stranded DNA that could be incorporated into the host genome<sup>11,17</sup>.

The action of some type III RMS is known to be regulated by phase variation<sup>10,12</sup>. The variation in the number of repeat copies in homopolymers or short tandem repeats, mainly due to slipped-strand mispairing during DNA replication, can cause the switch between a functional, non-functional or a different version of the gene<sup>18</sup>. It has been hypothesized that phase variation can be used by RMS to regulate genome flux<sup>11</sup> and it has been shown to be associated with a random switching in the expression of multiple genes (the so-called 'phasevarion', for phase-variable regulon<sup>12,19</sup>) in *Haemophilus influenzae*<sup>20,21</sup>, *N. meningitidis*, *N. gonorrhoeae*<sup>5</sup>, *Helicobacter pylori*<sup>22</sup>, and *Moraxella catarrhalis*<sup>23</sup>. Apart from the type III RMS, the type II 'orphan' Dam MTase has also been shown to be involved in gene regulation in several Gram-negative bacteria<sup>18</sup> and is even required for virulence in some pathogenic organisms such as *Escherichia coli*, *Salmonella*, *Yersinia* or *Vibrio* species<sup>24</sup>. However, the Dam MTase is also associated with the regulation of phase variation and participates in the methyl-directed mismatch repair system. Thus, Dam-defective bacteria have more flexibility to undergo phase variation<sup>25</sup>. The Dam MTase<sup>26</sup> has been found in some *N. lactamica* and *N. meningitidis* strains, however, there have been no reports of *N. gonorrhoeae* carrying this enzyme until now. Instead, the gonococcus and other Dam-defective *Neisseria* contain a *dam replacing gene (drg)*, characterized as an endonuclease<sup>26</sup> that has been reported to be important for adhesion and biofilm formation<sup>25</sup>.

In this study, we characterised the population of RMS and their associated methylation specificities in 25 *N. gonorrhoeae* strains using SMRT PacBio data. Results from this study give a more comprehensive insight into the link between genomics and epigenomics in the gonococcus.

## Results

### Detection of active RMS and associated DNA methylation

From 13 to 15 complete RMS were found in each of the 25 *N. gonorrhoeae* strains included in the study (Supplementary Table 1): 2 of Type I, 11 of Type II and 2 of Type III, all of which are present in the REBASE database (Figure 1). A detailed scan of the genomes for genes with a Pfam annotation associated with an REase or an MTase did not detect any new RMS. Most of the systems were found surrounded by core genes coding for essential functions, such as tRNA aminoacylation during protein translation, ATP binding or magnesium ion binding, but also DNA transposition and isomerase activity (Supplementary Table 2). This could potentially be related to the fact that Type I REases have an important requirement of

ATP for cleavage, and Type II REases require magnesium ion as a cofactor<sup>13</sup>. The amino acid translation of the repetitive region in both Type III RMS is rich in prolines and, interestingly, the *ppiA* gene upstream the methylase in NgoAX is an isomerase that catalyzes the isomerization of peptide bonds in prolyl residues, assisting protein folding<sup>27</sup>.

Analysis of patterns of DNA modification using PacBio sequencing information for the 25 gonococcal strains produced a curated list of 13 different methylated motifs (Table 1, raw motif summaries in Supplementary Table 3). Manual curation and visualization of per-base IPD ratios for each motif was essential to filter out low-frequency motifs, as several of them were found to overlap with those methylated in a higher frequency or were predicted because of higher IPD values in guanines. Indeed, O-6-methylguanine modification has been associated with DNA damage by alkylating agents<sup>28</sup>, and such modifications could be detected by the SMRT technology.

A general investigation of the average IPD ratios for methylated and unmethylated bases was conducted. 6mA modifications extracted from our final motif list showed the highest IPD values, with a mean above 5 (quantile (Q)25 = 4.3, Q75 = 6.7; Supplementary Figure 1). Cytosine methylations (4mC and 5mC) are difficult to detect using native DNA, specially 5mC, because of the lower effect on delaying the polymerase during PacBio sequencing. Still, they showed an average IPD ratio of around 3 (Q25 = 2.3, Q75 = 3.6 for 4mC; Q25 = 2.2, Q75 = 2.9 for 5mC), higher than the mean values for any unmethylated base (Q25 = 0.9, Q75 = 1.3; Supplementary Figure 1).

### **Type I RMS: NgoAV is multispecific**

Two Type I RMS were present in the *N. gonorrhoeae* strains (Figures 1 and 2). Type I NgoAV RMS showed at least three different sources of genetic variability within the specificity unit (*hsdS*). An insertion of a T at position 186 of the gene creates a frameshift that causes a premature stop codon in position 62 of the protein that inactivates *hsdS*, and thus, the whole RMS (Figure 3). From 0 to 3 ATLE repeats were observed in the amino acid sequence in different strains. Some strains also contained a downstream frameshift caused by the deletion of two Gs in a homopolymer. This second frameshift does not cause the inactivation of the gene, but in combination with the number of ATLE repeats, was found to be associated with a change in the recognition pattern (Figure 3 and Supplementary Table 4). Three different methylation targets were associated with this RMS. Strains NCTC10931 and WHO M with 1 ATLE repeat and no frameshift showed two-strand 6mA methylation in the 5'-GACN[6]TGC-3' motif (Supplementary Figure 2). However, four strains with 2 ATLE repeats instead showed the same starting and ending triplet nucleotides as the recognition pattern but the spacer was 1 bp longer, 5'-GACN[7]TGC-3' (Supplementary Figure 3). Finally,

two strains (FA1090 and WHO W) showed a slightly different motif which is a further 1 bp longer, 5'-GCAN[8]TGC-3' (Supplementary Figure 4). In this case, they also contained 2 ATLE repeats but with the frameshift downstream that causes a premature end of the protein, which therefore does not inactivate it but changes its specificity (Figure 3). The only exception was strain NCTC12700, which contained an apparently active methylase (*hsdM*) and complete specificity unit but no associated methylated motif. This could be caused by *hsdM* or *hsdS* being not functional due to unrecognised regulatory or mutational changes.

The second Type I RMS, NgoAXVII, did not show specific regions of variability, only the inactivation of the methylase and/or the specificity unit causes the inactivation of the whole RMS. All strains with an active *hsdM* and *hsdS* in this RMS showed two-strand 6mA methylation in the 5'-GAGN[5]TAC-3' motif (Supplementary Figure 5 and Supplementary Table 5). Two additional motifs were detected in three strains (FA1090, NCTC12700 and WHO O) that we hypothesized are off-target methylations by this enzyme because of the similarity of the motifs: double-stranded methylation in 5'-GAGN[4]GTYAC-3'/3'-CTCN[4]CARTG-5' and single-stranded methylation in 5'-YGTABYN[3]CCC-3', which is the reverse complement of 5'-GGGN[5]TAC-3', very similar to the main motif methylated by NgoAXVII, 5'-GAGN[5]TAC-3' (Table 1). The specificity unit of this enzyme was found to be very conserved, with only two non-synonymous mutations that had no association with the observed off-target methylations.

## Cytosine methylation by Type II RMS

*N. gonorrhoeae* contains up to 11 Type II RMS (Figures 1 and 2), 7 of them strictly associated with cytosine methylation (mC). In spite of apparently showing functional cytosine methylases, only three motifs containing 5mC or 4mC methylation were detected by the PacBio pipeline (Table 1 and Supplementary Table 3). However, as we show above, both types of cytosine modification can be distinguished from an unmethylated base (Supplementary Figure 1). In order to perform a deeper analysis of the motifs associated with mC in Type II RMS, we specifically identified those reported in REBASE and evaluated their per-base IPD ratios (Supplementary Table 6). A Mann-Whitney test was performed between the distribution of IPD values for the predicted mCs and a random set of unmethylated cytosines for each strain to assess whether there was a statistically significant difference. Significant cytosine modification was found for the motifs associated with the 7 RMS expressing an active methylase (Supplementary Figure 6 and Supplementary Table 6). Four strains showed a gene annotated as "virulence-associated" instead of the NgoAIII RMS (WHO F, WHO L, NCTC10928 and NCTC12700) and three of those did not show significant signals of methylation (Supplementary Table 6). The motif associated with NgoAIII in

REBASE is 5'-GGCGCC-3' (5mC) and this was detected as a 4mC signal (reported as the degenerate motif 5'-GGCSCCND-3') in WHO V (Supplementary Table 3), although we showed the signal was in fact present in all the strains with a functional RMS (Supplementary Figure 6). The WHO G data produced a non-significant test for mC in the forward strand of 5'-GCGGC-3', associated with NgoAVII, which could be due to hemi-methylation or poor signal in the data.

RMS NgoAXIII was found as a VSR (Very Short patch Repair) protein followed by a MTase and there is no report of its target motif in REBASE. A series of BLAST searches revealed that only *N. polysaccharea* M18661 (Genbank accession number CP031325.1) contained a homologous region to NgoAXIII. Instead, genomes of *N. lactamica* (i.e. Y92-1009, Genbank accession number CP019894.1) showed a longer version of the MTase (362 amino acids compared to 112 in *N. gonorrhoeae* FA1090 and 133 in the rest). This finding could mean that *N. gonorrhoeae* does not have a complete NgoAXIII system, but only remnants of a non-active MTase. Interestingly, several *N. meningitidis* genomes contain the VSR protein followed by a different RMS instead (NmeDI).

The 5'-AAANCGGTTNNC-3' motif was detected by the PacBio SMRT pipeline containing single-strand m4C methylation in the underlined base. The study of the distribution of per-base IPD ratios in the WHO strains revealed that the complementary sequence was probably also methylated at 5'-AAANCGGTTNNC-3'/3'-TTTNGCC<sup>u</sup>AANNG-5' (Supplementary Figure 7). A Mann-Whitney test revealed the distribution of IPD values for these two bases to be significantly different than that of an unmethylated cytosine for most of the strains (Supplementary Figure 8). FA1090 and, especially, the NCTC strains, consistently showed very noisy results, but the test was still statistically significant for them (Supplementary Figure 8). We hypothesize that the methylation of this motif could result from a secondary target of a Type II RMS or might be the target of NgoAXIII, for which the methylated motif is currently unknown, although this is unlikely as the fraction of methylated 5'-AAANCGGTTNNC-3' is very low (Table 1).

### 6mA methylation by Type II RMS

Only three of the Type II RMS contain 6mA methylases. The NgoAVIII system has a Bcgl-like structure<sup>29</sup> as it is formed by an enzyme with restriction and 6mA methylation properties (RM) and a specificity unit (S). A BLASTn search of the genomic region spanning both units against the NCBI database revealed that this RMS has not been described in any other organism. A transposase immediately upstream of the RM unit indicates that this system may be mobile and acquired from an unknown source. The target motif in REBASE for this RMS is that of Bcgl, 5'-GACN[5]TGA-3', although we do not detect this motif as methylated



or any other motif with 6mA methylation that could be associated with NgoAVIII (Supplementary Table 3), despite the system appearing intact. Interestingly, BcgI-like MTases are known to prefer hemi-methylated DNA, although their recognition sequences are very similar to those of Type I MTases, which methylate both strands<sup>30,31</sup>. Thus, NgoAVIII may have a maintenance role for the motifs targeted by Type I MTases, methylating the growing strand immediately after chromosomal replication, a point when these motifs will be hemi-methylated.

The Dam methylase was found in place of the restriction enzyme in NgoAXI in two strains (NCTC10931 and NCTC12700), which showed double-strand 6mA 5'-GATC-3' methylation (Supplementary Figure 9). This MTase has been reported to be substituted by the Dam-replacing endonuclease (Drg) in many *Neisseria* species, but has not been found in *N. gonorrhoeae* until now<sup>25,26</sup>. Finally, the NgoAXVI RMS is formed by two methylases (5mC and 6mA, respectively) and a restriction enzyme. A clear 6mA signal is found in the associated 5'-GGTGA-3' motif in all the strains (Supplementary Figure 10). A more detailed comparative analysis of the distribution of IPD values of the three cytosines in the reverse complement of this motif against unmethylated Cs for each strain identified the middle cytosine as the most plausible candidate for 5mC methylation (3'-CCACT-5' Supplementary Figure 11).

### **Type III RMS: the real target of NgoAXII is 5'-GCAGA-3'**

Two Type III RMS were present in the *N. gonorrhoeae* strains (Figures 1 and 2). These contain a phase-variable methylase (*Mod*) that controls the activation or inactivation of the system and includes a DNA-recognition domain (DRD) that controls its specificity. NgoAX contained a variable number of 5'-CCCAA-3' repeats that switch the system ON/OFF, and the *modB1* DRD allele (Supplementary Table 7). In contrast, the NgoAXII contained a variable number of 5'-AGCC-3' repeats and the *modA13* DRD allele (Supplementary Table 8). Both *modA* and *modB* DRD alleles are those typically found in *N. gonorrhoeae*<sup>8,32</sup>.

Active NgoAX RMS showed 6mA 5'-CCACC-3' methylation in 6 strains (Supplementary Figure 12 and Supplementary Table 7), and active NgoAXII showed 6mA 5'-GCAGA-3' methylation in 3 strains (Supplementary Figure 13 and Supplementary Table 8). Interestingly, a previous report suggested the methylated motif for the *modA13* allele of this enzyme was 5'-AGAAA-3' by evaluating the digestion pattern of Apol and other enzymes<sup>5</sup>, and this is the recognition motif associated with NgoAXII in REBASE. However, our analysis of the distribution of the per-base IPD ratios for both motifs in the three strains with an active *modA* (NCTC13798, NCTC13799 and WHO N; Figure 4) revealed that less than 0.1% of the AGAAA motifs contained an IPD ratio above 3, while most of the GCAGA (>79%) showed

high IPDs in the underlined base (Supplementary Table 9). We only observed between 4 and 8 motifs in the genomes that included both motifs overlapping with each other and that expanded to include a complete A<sub>po</sub>I recognition sequence (5'-RAATTY-3'), 5'-GCAGAAATTY-3'. Therefore, the digestion in these cases could have been impeded by the methylation of the second A instead (in bold). A low fraction of motifs with high IPDs in the first A of 5'-GCAGA-3' was found in the strains with an active Type I NgoAV, and this may be due to overlapping signals between both RMS (Supplementary Figures 1-3 and Supplementary Figure 13). Strangely, results from the PacBio motif analysis in FA1090 for 5'-GCAGA-3' did not show signals of methylation, although Mod was apparently complete in the reference genome. An assembly of the PacBio reads associated with the methylation data for this strain<sup>9</sup> revealed that the methylase in the NgoAXII RM system in the subclone of FA1090 used for sequencing in the cited work contained 35 AGCC repeats instead of 37 in the reference genome, causing a protein truncation and thus a lack of 5'-GCAGA-3' methylation.

### **Orphan methylases in *N. gonorrhoeae* include the Dam methylase**

Disrupted methylases can cause the inactivation of the whole RMS, due to mechanisms that avoid self-degradation of the chromosomal DNA<sup>12</sup> in this situation. However, methylases can be functional with a disrupted or absent endonuclease. These methylases have been associated with a regulatory rather than defensive role and have been described as 'orphan'. In the strains under study we observed multiple examples of this in the Type I NgoAXVII (Supplementary Table 5) and several Type II RMS (Supplementary Table 6). The Dam methylase in NgoAXI acts as an 'orphan' methylase, as it completely replaces the endonuclease of the system, and the downstream MTase is disrupted. A BLASTp search of the NCBI non-redundant protein database revealed that the gonococcal Dam sequence has 97-98% amino acid identity to *N. polysaccharea*, 97% to *N. flavescens*, 96-97% to *N. lactamica*, and 92-96% to *N. meningitidis*. A broader screening of the Dam protein in the *Neisseria* genus showed a reasonable amount of diversity at the amino acid level, with recombination swapping sequence variants among *Neisseria* species (Supplementary Figure 14).

A scan of the Pfam domains of the proteome of the 25 strains under study did not reveal further DNA methylases without an accompanying restriction enzyme in the main chromosomes. However, the strains carrying the conjugative plasmid (WHO G, L, M, N, O, W, and NCTC10931) showed an orphan MTase in this element that has a 100% amino acid and nucleotide identity to the Type II M.Ngo8107ORF11P and M.Ngo5289ORFAP enzymes in REBASE. Visualization of the annotation of the WHO strains revealed that it is present in



the conjugative plasmid (Supplementary Figure 15). No further matches were found to any other bacteria in a BLASTn search against a non-redundant nucleotide database.

## Discussion

Restriction-modification systems have been shown to perform several other roles apart from being a defence system against foreign DNA<sup>11</sup>. For example, some of them, especially those from type III, are involved in the control of gene expression through phase variable repeats in the Mod unit<sup>19</sup>. Previous papers have characterised different RMS in the *N. gonorrhoeae* FA1090 reference strain<sup>6-8,33</sup> and also data on its methylation status has been described<sup>9</sup>. In this manuscript, we provide a detailed view of the structure and variability of DNA MTases and specificity units in RMS in 25 strains of *N. gonorrhoeae* and link them to the detected methylated motifs.

The analysed gonococcal genomes showed from 13 to 15 complete RMS, all of them described in the reference strain FA1090 in REBASE<sup>34</sup>. However, a detailed screening of the genome assemblies did not reveal any new systems or orphan MTases, apart from an enzyme homologous to M.Ngo8107ORF11P on the conjugative plasmid in the strains carrying this element. Thirteen curated methylated motifs were detected in the 25 strains, most of them containing 6mA methylation (Table 1). A detailed analysis of the genetic variability of the *hsdS* units of Type I MTases and the phase variable repeats in the DRD of Type III MTases allowed us to link them to particular methylated motifs. In the case of the Type I NgoAV, earlier work confirmed the functionality of the naturally truncated *hsdS* in FA1090, which was predicted to recognize the 5'-GCAN[8]TGC-3' motif<sup>35</sup>. Later, it was proven that the cause of this truncation was a phase-variable G homopolymer and that the restitution of a complete form of this unit changed the specificity of the MTase to 5'-GCAN[7]STGC-3'<sup>7</sup>. Our study reveals further insight into the diversity of specificities of the NgoAV RMS. Three different motifs were detected in strains with a functional NgoAV. The change in specificity was found to be related to a combination of two factors: the number of ATLE repeats in the protein sequence of *hsdS* and the presence of the frameshift that causes the known truncation downstream of these repeats (Figure 3 and Supplementary Table 4). Strains with 1 ATLE repeat and a complete *hsdS* showed 5'-GACN[6]TGC-3' double-stranded 6mA methylation, strains with 2 ATLE repeats and a complete *hsdS* showed 5'-GACN[7]TGC-3' and strains with 2 ATLE repeats and a downstream frameshift causing truncation, which is the case in FA1090, showed 5'-GCAN[8]TGC-3' (Supplementary Table 4) instead. The presence of amino acid repeats in the centre region of *hsdS* has already been described<sup>35,36</sup>, but not its role in sequence specificity in combination

with the downstream frameshift. Off-target methylated motifs were associated with the second Type I RMS NgoAXVII, although no associated variability was found in *hsdS* in this system.

Seven out of 11 Type II RMS in *N. gonorrhoeae* are associated with cytosine methylation. However, this type of DNA modification is difficult to detect by current PacBio sequencing, and high coverage in combination with Tet-conversion is preferable<sup>37</sup>. In this study, we relied on available data from other projects, which performed sequencing on untreated native DNA, thus, most of the analysis were performed on 6mA methylation. Only three motifs associated with 5mC or 4mC methylation were detected in a subset of the strains (Table 1). Nonetheless, a statistical comparison between the IPD values of target cytosines in the motifs annotated in REBASE for these RMS to those from unmethylated cytosines revealed significantly higher values for those in the strains carrying active 5mC MTases (Supplementary Table 5). 5mC can be converted into T by deamination, generating mismatches that can produce mutations. Bacteria can modulate this effect by methylating 4mC instead<sup>11</sup>. In fact, double-stranded methylation of 5'-CCGCGG-3' by NgoAIII was detected by the PacBio system as 4mC, apart from the potential off-target 5'-AAANCGGTTNNC-3' motif (Table 1). Interestingly, secondary methylation by these types of enzymes has been observed in *H. influenzae*<sup>38</sup>.

*N. gonorrhoeae* harbours two VSR endonucleases to correct T:G mismatches, V.NgoAXIII and V.NgoAXIV. The specificity of NgoAXIII is still unknown and, in this study, we hypothesize that the MTase may be truncated, as a longer version is found in *N. lactamica* (Genbank accession number CP019894.1). However, the VSR may still be functional, as it has been found to recognize T:G mismatches in every nucleotide context<sup>33</sup>. V.NgoAXIV has also been described to recognize mismatches in sequences other than 5'-CCGG-3'<sup>33</sup>. No methylated motif has been found for the Bcg-like RMS NgoAVIII, however, previous reports describe this type of MTase as having a strict preference for hemi-methylated sites<sup>30,31</sup> in Type I-like motifs. Thus, we propose that NgoAVIII may have a maintenance role, methylating the replicated strand in sites targeted by Type I RMS during DNA replication.

The Dam methylase is present in some strains of *N. lactamica* and *N. meningitidis*<sup>26</sup>. However, it has not been described in *N. gonorrhoeae* until now. Instead, *Neisseria* species or strains lacking this enzyme are known to carry an endonuclease encoded by *drg* (*dam* replacing gene) recognizing the 5'-GATC-3' motif<sup>26,39</sup>. Here, we observed two *N. gonorrhoeae* strains (NCTC10931 and NCTC12700) carrying the Dam MTase, instead of the *drg* gene, in the NgoAXI RMS locus, followed by a truncated MTase. Thus, Dam is acting as an orphan enzyme, a form that has been related to gene expression regulation<sup>40</sup>. Interestingly, previous publications have suggested that strains carrying the *drg* gene have

an advantage over those carrying the Dam MTase because they have more flexibility for phase-variation, as Dam participates in DNA mismatch repair during replication<sup>25,39</sup>. Also, the same publications showed experimentally that strains carrying *drg* form more stable biofilms and have better adhesion to human cells during infection.

Several works have studied the so-called ‘phasevarion’ in *N. gonorrhoeae* and *N. meningitidis*, referred to the set of genes with an altered expression due to phase variation in an MTase<sup>5,12</sup>. Type I NgoAV has been considered to regulate a phasevarion as the phase variation of the G homopolymer is involved in a change in specificity<sup>19</sup>. In our work, we complement this result, as we show that the change in specificity is produced as a combination of the phase-variable poly-G and the number of ATLE repeats in the amino acid sequence (Supplementary Table 4). However, the best example of this is the phase-variable Type III RMS<sup>12</sup>, in which for example *modD1* has been associated with hypervirulent *N. meningitidis* clonal complexes and *modA13* with enhanced biofilm formation and intracellular survival<sup>19</sup>. In a recent publication, the NgoAX Mod was experimentally inactivated in *N. gonorrhoeae* and a deregulation of the expression of 121 genes was observed, along with an effect on the adherence to and invasion of host epithelial cells, which are essential for infection<sup>8</sup>. In contrast, inactivation of the NgoAXII Mod enzyme showed a deregulation of 54 genes under iron-limiting conditions<sup>5</sup>. Studies have reported that *N. gonorrhoeae* carries a very conserved *modB1* allele in NgoAX compared to *modA13* in NgoAXII. These results have led to the suggestion that NgoAX Mod is the main regulator of the epigenome, while NgoAXII regulates in specific growth conditions<sup>8</sup>. Here, we observe that indeed all the strains under study carry the *modA13* and *modB1* alleles, with all possible combinations of activity apparent: we detected strains with only the NgoAX Mod active (i.e. FA1090 or WHO F), some with only the NgoAXII Mod active (WHO N and NCTC13799), one with both active (NCTC13798), and the rest with both systems inactive (Supplementary Table 7 and Supplementary Table 8).

In summary, *N. gonorrhoeae* contains several RMS to protect the genome against foreign DNA invasion. Genetic diversity created through phase variation or hypervariable domains controls their activation or inactivation, along with a variation of specificity in one case. This work gives further insight into the RMS in *N. gonorrhoeae* and the consequent methylation landscape, which is able to change to modulate gene expression. We also show the importance of detailed PacBio SMRT analysis to enhance and complete the methylation information contained in public databases.

## Methods

### Genomes included in the study

We analysed the genomes of a total of 25 *N. gonorrhoeae* strains that had been sequenced using PacBio. Of those, 14 complete genomes were obtained from the 2016 WHO reference panel<sup>41</sup> and 10 additional genomes were selected from the Public Health England NCTC 3000 collection (<http://www.sanger.ac.uk/resources/downloads/bacteria/nctc/>) avoiding overlap with the WHO panel, to further improve the representativeness of the *N. gonorrhoeae* genome diversity (Supplementary Table 1). Additionally, PacBio raw data and predicted motifs were retrieved for the reference genome *N. gonorrhoeae* FA1090 from a recent study<sup>9</sup>. Sequencing was run using native DNA in all cases. Assembly and annotation were performed as reported in the publications cited above, using an automatic and improved pipeline at the Wellcome Sanger Institute<sup>41</sup>. FA1090 was reassembled using the PacBio data with Canu v1.6<sup>42</sup> to compare the number of tandem repeats in the type III MTases with those from the original assembly available in the public databases (GenBank accession AE004969).

### Analysis of DNA methylation

Genome-wide base modifications and predicted modified motifs were called using the RS\_Modification\_and\_Motif\_Analysis protocol from the PacBio SMRT Analysis software v2.3.0. Coordinates of the predicted motifs were localized in all the genomes and plasmids using the EMBOSS application *fuzznuc*<sup>43</sup>. Per-base IPD ratios for the predicted motifs were extracted from the raw data and visualized in the 25 strains using R<sup>44</sup>.

The distribution of IPD values for each of the four bases outside any of the predicted motifs was used as the distribution of values for the four unmethylated bases. A random subsample of 10,000 unmethylated sites of each base was used. Cytosine methylation from type II RMS was inferred by evaluating the distribution of IPD ratios in the target base in the associated motif in REBASE<sup>34</sup> compared to that of unmethylated cytosines. A Mann-Whitney test was performed to evaluate statistical significance and p-values were corrected using Bonferroni correction<sup>45</sup>.

### Detection and specificity of RMS

An extended and manually curated list of Pfam domains associated with REases and MTases from a previous work<sup>46</sup> was used to detect these genes in the annotations of the genomes under study. HMMER<sup>47</sup> *hmmsearch* was run independently on the proteome of each genome (.faa files) against Pfam(A) v30<sup>48</sup> to complement the annotation information. R language<sup>44</sup> was used to parse the annotation files and the results from *hmmsearch* and extract genes with an inferred Pfam domain in the list of target REases and MTases. Genes predicted to be an REase or a MTase that were less than 10 genes distant were considered

to be part of the same RMS. The rest were tagged as 'orphan'. Results were compared to the RMS annotated for FA1090 in REBASE<sup>34</sup>. The presence or absence of every RMS in each strain was visualized using phandango<sup>49</sup>.

The motif annotated in REBASE<sup>34</sup> for each MTase was compared to the list of predicted motifs obtained by the PacBio SMRT Analysis software for all the strains. Nucleotide and protein sequences of every unit of the RMS were extracted for each strain and compared using SeaView v4.6.1<sup>50</sup> to look for sources of variability within the MTases or the S units. These and the flanking genes were visualized using Artemis v16.0.0<sup>51</sup>. Disrupted proteins were confirmed by a protein-protein BLAST against a non-redundant protein database<sup>52</sup>.

Representative protein sequences of the Dam methylase in the *Neisseria* genus were downloaded from the Identical Protein Groups (IPG) tool of the NCBI database (accessed on 21/01/2019). These sequences were aligned using SeaView v4.6.1<sup>50</sup> and the resulting alignment trimmed using Gblocks v0.91b<sup>53</sup> considering only positions in which at least half of the sequences do not have a gap. PhyML<sup>54</sup> was used to build a maximum likelihood tree using the LG model and performing 100 bootstrap replicates. Final tree and metadata were visualised using iTOL<sup>55</sup>.

## GO enrichment analysis

A Gene Ontology (GO) enrichment analysis was performed for the genes flanking all the RMS using the *topGO* R package<sup>56</sup>. The three ontologies were scanned ('BP', biological process; 'MF', molecular function; and 'CC', cellular component). The *classic* and *weight01* algorithms were used which do not, and do, use hierarchy information for scoring a particular GO term, respectively<sup>57</sup>. The statistical significance of the enrichment was calculated using a Fisher's exact test between the observed and expected number of genes in each term for each algorithm. Terms with a p-value<0.05 in both algorithms were considered as significant in the results. P-values were not corrected for multiple testing to avoid excluding significant GO terms near the cut-off. Besides, the tests are not independent when using the *weight01* algorithm as it is conditioned on neighbouring terms<sup>56</sup>.

## Data availability

Scripts used to perform the analyses and plots in this work are available in the GitHub repository <https://github.com/leosanbu/MethylationProject>. The 14 *N. gonorrhoeae* PacBio raw data and complete genomes from the 2016 WHO panel are available under the ENA Bioproject PRJEB14020 (Sample accessions SAMEA2448460-SAMEA2448470 and SAMEA2796326-SAMEA2796328)<sup>41</sup>. Accession numbers for the PacBio data from the ten strains downloaded from the NCTC3000 project are available in the following link: <https://www.sanger.ac.uk/resources/downloads/bacteria/nctc/> (Sample accessions

SAMEA3174297-SAMEA3174299, SAMEA4076737, SAMEA4076741, SAMEA4076765, SAMEA4076768-SAMEA4076770, SAMEA4076773). GFF files for the WHO and NCTC sequence data are available in the GitHub repository <https://github.com/leosanbu/MethylationProject>. The *N. gonorrhoeae* FA1090 reference genome was obtained from NCBI accession number AE004969 and the PacBio raw data from the study by Blow *et al* (2016)<sup>9</sup>. Supplementary Table 1 contains the detailed information for each strain.



## References

- 1 Grad, Y. H. *et al.* Genomic epidemiology of gonococcal resistance to extended spectrum cephalosporins, macrolides, and fluoroquinolones in the US, 2000-2013. *J Infect Dis* **214**, 1579-1587, doi:10.1093/infdis/jiw420 (2016).
- 2 De Silva, D. *et al.* Whole-genome sequencing to determine transmission of *Neisseria gonorrhoeae*: an observational study. *Lancet Infect Dis* **16**, 1295-1303, doi:10.1016/S1473-3099(16)30157-8 (2016).
- 3 Sánchez-Busó, L. *et al.* Antimicrobial exposure in sexual networks drives divergent evolution in modern gonococci. (*under final review*), *bioRxiv* 334847, doi:<https://doi.org/10.1101/334847> (2019).
- 4 Harris, S. R. *et al.* Public health surveillance of multidrug-resistant clones of *Neisseria gonorrhoeae* in Europe: a genomic survey. *Lancet Infect Dis* **18**, 758-768, doi:10.1016/S1473-3099(18)30225-1 (2018).
- 5 Srikhanta, Y. N. *et al.* Phasevarions mediate random switching of gene expression in pathogenic *Neisseria*. *PLoS Pathog* **5**, e1000400, doi:10.1371/journal.ppat.1000400 (2009).
- 6 Adamczyk-Poplawska, M., Lower, M. & Piekarowicz, A. Characterization of the NgoAXP: phase-variable type III restriction-modification system in *Neisseria gonorrhoeae*. *FEMS Microbiol Lett* **300**, 25-35, doi:10.1111/j.1574-6968.2009.01760.x (2009).
- 7 Adamczyk-Poplawska, M., Lower, M. & Piekarowicz, A. Deletion of one nucleotide within the homonucleotide tract present in the *hsdS* gene alters the DNA sequence specificity of type I restriction-modification system NgoAV. *J Bacteriol* **193**, 6750-6759, doi:10.1128/JB.05672-11 (2011).
- 8 Kwiatek, A., Mrozek, A., Bacal, P., Piekarowicz, A. & Adamczyk-Poplawska, M. Type III Methyltransferase M.NgoAX from *Neisseria gonorrhoeae* FA1090 Regulates Biofilm Formation and Interactions with Human Cells. *Front Microbiol* **6**, 1426, doi:10.3389/fmicb.2015.01426 (2015).
- 9 Blow, M. J. *et al.* The Epigenomic Landscape of Prokaryotes. *PLoS Genet* **12**, e1005854, doi:10.1371/journal.pgen.1005854 (2016).
- 10 Rotman, E. & Seifert, H. S. The genetics of *Neisseria* species. *Annu Rev Genet* **48**, 405-431, doi:10.1146/annurev-genet-120213-092007 (2014).
- 11 Vasu, K. & Nagaraja, V. Diverse functions of restriction-modification systems in addition to cellular defense. *Microbiol Mol Biol Rev* **77**, 53-72, doi:10.1128/MMBR.00044-12 (2013).
- 12 Srikhanta, Y. N., Fox, K. L. & Jennings, M. P. The phasevarion: phase variation of type III DNA methyltransferases controls coordinated switching in multiple genes. *Nat Rev Microbiol* **8**, 196-206, doi:10.1038/nrmicro2283 (2010).
- 13 Roberts, R. J. *et al.* A nomenclature for restriction enzymes, DNA methyltransferases, homing endonucleases and their genes. *Nucleic Acids Res* **31**, 1805-1812 (2003).
- 14 Roberts, R. J., Vincze, T., Posfai, J. & Macelis, D. REBASE: restriction enzymes and methyltransferases. *Nucleic Acids Res* **31**, 418-420 (2003).

- 15 Kobayashi, I. Behavior of restriction-modification systems as selfish mobile elements and their impact on genome evolution. *Nucleic Acids Res* **29**, 3742-3756 (2001).
- 16 Furuta, Y., Abe, K. & Kobayashi, I. Genome comparison and context analysis reveals putative mobile forms of restriction-modification systems and related rearrangements. *Nucleic Acids Res* **38**, 2428-2443, doi:10.1093/nar/gkp1226 (2010).
- 17 Oliveira, P. H., Touchon, M. & Rocha, E. P. Regulation of genetic flux between bacteria by restriction-modification systems. *Proc Natl Acad Sci U S A* **113**, 5658-5663, doi:10.1073/pnas.1603257113 (2016).
- 18 Casadesus, J. & Low, D. Epigenetic gene regulation in the bacterial world. *Microbiol Mol Biol Rev* **70**, 830-856, doi:10.1128/MMBR.00016-06 (2006).
- 19 Seib, K. L., Jen, F. E., Scott, A. L., Tan, A. & Jennings, M. P. Phase variation of DNA methyltransferases and the regulation of virulence and immune evasion in the pathogenic *Neisseria*. *Pathog Dis* **75**, ftx080, doi:10.1093/femspd/ftx080 (2017).
- 20 Srikhanta, Y. N., Maguire, T. L., Stacey, K. J., Grimmond, S. M. & Jennings, M. P. The phasevarion: a genetic system controlling coordinated, random switching of expression of multiple genes. *Proc Natl Acad Sci U S A* **102**, 5547-5551, doi:10.1073/pnas.0501169102 (2005).
- 21 Fox, K. L. *et al.* Haemophilus influenzae phasevarions have evolved from type III DNA restriction systems into epigenetic regulators of gene expression. *Nucleic Acids Res* **35**, 5242-5252, doi:10.1093/nar/gkm571 (2007).
- 22 Srikhanta, Y. N. *et al.* Phasevarion mediated epigenetic gene regulation in *Helicobacter pylori*. *PLoS One* **6**, e27569, doi:10.1371/journal.pone.0027569 (2011).
- 23 Blakeway, L. V. *et al.* ModM DNA methyltransferase methylome analysis reveals a potential role for *Moraxella catarrhalis* phasevarions in otitis media. *FASEB J* **28**, 5197-5207, doi:10.1096/fj.14-256578 (2014).
- 24 Wion, D. & Casadesus, J. N6-methyl-adenine: an epigenetic signal for DNA-protein interactions. *Nat Rev Microbiol* **4**, 183-192, doi:10.1038/nrmicro1350 (2006).
- 25 Kwiatek, A., Bacal, P., Wasiluk, A., Trybunko, A. & Adamczyk-Poplawska, M. The dam replacing gene product enhances *Neisseria gonorrhoeae* FA1090 viability and biofilm formation. *Front Microbiol* **5**, 712, doi:10.3389/fmicb.2014.00712 (2014).
- 26 Cantalupo, G. *et al.* Evolution and function of the neisserial dam-replacing gene. *FEBS Lett* **495**, 178-183 (2001).
- 27 Unal, C. M. & Steinert, M. Microbial peptidyl-prolyl cis/trans isomerases (PPIases): virulence factors and potential alternative drug targets. *Microbiol Mol Biol Rev* **78**, 544-571, doi:10.1128/MMBR.00015-14 (2014).
- 28 Sedgwick, B. Repairing DNA-methylation damage. *Nat Rev Mol Cell Biol* **5**, 148-157, doi:10.1038/nrm1312 (2004).
- 29 Loenen, W. A., Dryden, D. T., Raleigh, E. A. & Wilson, G. G. Type I restriction enzymes and their relatives. *Nucleic Acids Res* **42**, 20-44, doi:10.1093/nar/gkt847 (2014).

- 30 Kong, H. Analyzing the functional organization of a novel restriction modification system, the BcgI system. *J Mol Biol* **279**, 823-832, doi:10.1006/jmbi.1998.1821 (1998).
- 31 Smith, R. M., Jacklin, A. J., Marshall, J. J., Sobott, F. & Halford, S. E. Organization of the BcgI restriction-modification protein for the transfer of one methyl group to DNA. *Nucleic Acids Res* **41**, 405-417, doi:10.1093/nar/gks1000 (2013).
- 32 Tan, A. *et al.* Distribution of the type III DNA methyltransferases modA, modB and modD among *Neisseria meningitidis* genotypes: implications for gene regulation and virulence. *Sci Rep* **6**, 21015, doi:10.1038/srep21015 (2016).
- 33 Kwiatek, A., Luczkiewicz, M., Bandyra, K., Stein, D. C. & Piekawicz, A. *Neisseria gonorrhoeae* FA1090 carries genes encoding two classes of Vsr endonucleases. *J Bacteriol* **192**, 3951-3960, doi:10.1128/JB.00098-10 (2010).
- 34 Roberts, R. J., Vincze, T., Posfai, J. & Macelis, D. REBASE--a database for DNA restriction and modification: enzymes, genes and genomes. *Nucleic Acids Res* **43**, D298-299, doi:10.1093/nar/gku1046 (2015).
- 35 Piekawicz, A., Klyz, A., Kwiatek, A. & Stein, D. C. Analysis of type I restriction modification systems in the *Neisseriaceae*: genetic organization and properties of the gene products. *Mol Microbiol* **41**, 1199-1210 (2001).
- 36 MacWilliams, M. P. & Bickle, T. A. Generation of new DNA binding specificity by truncation of the type IC EcoDXXI hsdS gene. *EMBO J* **15**, 4775-4783 (1996).
- 37 Clark, T. A. *et al.* Enhanced 5-methylcytosine detection in single-molecule, real-time sequencing via Tet1 oxidation. *BMC Biol* **11**, 4, doi:10.1186/1741-7007-11-4 (2013).
- 38 Cohen, H. M., Tawfik, D. S. & Griffiths, A. D. Promiscuous methylation of non-canonical DNA sites by HaeIII methyltransferase. *Nucleic Acids Res* **30**, 3880-3885 (2002).
- 39 Bucci, C. *et al.* Hypermutation in pathogenic bacteria: frequent phase variation in meningococci is a phenotypic trait of a specialized mutator biotype. *Mol Cell* **3**, 435-445 (1999).
- 40 Low, D. A., Weyand, N. J. & Mahan, M. J. Roles of DNA adenine methylation in regulating bacterial gene expression and virulence. *Infect Immun* **69**, 7197-7204, doi:10.1128/IAI.69.12.7197-7204.2001 (2001).
- 41 Unemo, M. *et al.* The novel 2016 WHO *Neisseria gonorrhoeae* reference strains for global quality assurance of laboratory investigations: phenotypic, genetic and reference genome characterization. *J Antimicrob Chemother* **71**, 3096-3108, doi:10.1093/jac/dkw288 (2016).
- 42 Koren, S. *et al.* Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res* **27**, 722-736, doi:10.1101/gr.215087.116 (2017).
- 43 Rice, P., Longden, I. & Bleasby, A. EMBOS: the European Molecular Biology Open Software Suite. *Trends Genet* **16**, 276-277 (2000).
- 44 R Development Core Team, R. I. R: A Language and Environment for Statistical Computing. (2008).
- 45 Miller, R. G. (ed New York: Springer-Verlag) (1981).

- 46 Croucher, N. J. *et al.* Diversification of bacterial genome content through distinct mechanisms over different timescales. *Nat Commun* **5**, 5471, doi:10.1038/ncomms6471 (2014).
- 47 Eddy, S. R. Accelerated Profile HMM Searches. *PLoS Comput Biol* **7**, e1002195, doi:10.1371/journal.pcbi.1002195 (2011).
- 48 Finn, R. D. *et al.* The Pfam protein families database: towards a more sustainable future. *Nucleic Acids Res* **44**, D279-285, doi:10.1093/nar/gkv1344 (2016).
- 49 Hadfield, J. *et al.* Phandango: an interactive viewer for bacterial population genomics. *Bioinformatics*, 1-2, doi: 10.1093/bioinformatics/btx610 (2017).
- 50 Gouy, M., Guindon, S. & Gascuel, O. SeaView version 4: A multiplatform graphical user interface for sequence alignment and phylogenetic tree building. *Mol Biol Evol* **27**, 221-224, doi:10.1093/molbev/msp259 (2010).
- 51 Carver, T., Harris, S. R., Berriman, M., Parkhill, J. & McQuillan, J. A. Artemis: an integrated platform for visualization and analysis of high-throughput sequence-based experimental data. *Bioinformatics* **28**, 464-469, doi:10.1093/bioinformatics/btr703 (2012).
- 52 Johnson, M. *et al.* NCBI BLAST: a better web interface. *Nucleic Acids Res* **36**, W5-9, doi:10.1093/nar/gkn201 (2008).
- 53 Castresana, J. Selection of conserved blocks from multiple alignments for their use in phylogenetic analysis. *Mol Biol Evol* **17**, 540-552 (2000).
- 54 Guindon, S. *et al.* New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0. *Syst Biol* **59**, 307-321, doi:10.1093/sysbio/syq010 (2010).
- 55 Letunic, I. & Bork, P. Interactive tree of life (iTOL) v3: an online tool for the display and annotation of phylogenetic and other trees. *Nucleic Acids Res* **44**, W242-W245, doi:10.1093/nar/gkw290 (2016).
- 56 Alexa, A. & Rahnenfuhrer, J. topGO: Enrichment Analysis for Gene Ontology. R package version 2.26.0. (2016).
- 57 Alexa, A., Rahnenfuhrer, J. & Lengauer, T. Improved scoring of functional groups from gene expression data by decorrelating GO graph structure. *Bioinformatics* **22**, 1600-1607, doi:10.1093/bioinformatics/btl140 (2006).

## Table and figure legends

Table 1. Restriction-Modifications systems and target motifs in *N. gonorrhoeae*. Those detected as methylated by the PacBio SMRT pipeline are marked as such and the fraction of modified motifs indicated. For those not directly detected by the pipeline, the associated motifs in REBASE are shown and marked as detected if a significant difference in IPD ratios is found between the target base and the unmethylated equivalent by a Mann-Whitney test.

Supplementary Table 1. List of *N. gonorrhoeae* strains used in the analysis. The ENA accession information of the raw and complete genomic data are indicated along with the average per-base sequencing coverage. Sample accession numbers are used for raw data as some of the strains include more than one sequencing run.

Supplementary Table 2. Results from a GO enrichment analysis of the flanking genes of the 15 RMs in *N. gonorrhoeae*. The number of genes annotated in each GO term and the expected number of genes to fall in each category are shown together with the number of significant hits in the flanking genes. GO terms significantly enriched by both the *classic* and *weight01* algorithms (p-value < 0.05) are shown. Results are shown for the three sub-ontologies (BP = Biological Process, MF = Molecular Function, CC = Cellular Component).

Supplementary Table 3. Compilation of the motif\_summary.csv output files of the PacBio methylation pipeline run on the 25 *N. gonorrhoeae* included in the study. Fields are explained in the following technical note: <https://github.com/PacificBiosciences/Bioinformatics-Training/wiki/Methylome-Analysis-Technical-Note>.

Supplementary Table 4. Type I NgoAV restriction-modifications system (RMS). Active (A) or disrupted (D) components are indicated. Red highlight is used to mark a cause of disruption. The different sources of variability in the pattern recognition domain of the specificity unit, which cause different methylated motifs (in different colours) are shown. An active methyltransferase requires both the methylase and the specificity unit to be active. In that case, if the restriction enzyme is not functional, the methylase is tagged as 'orphan'. *hsdR*: Restriction endonuclease; *hsdM*: Methyltransferase; *hsdS*: Specificity unit.

Supplementary Table 5. Type I NgoAXVII restriction-modifications system (RMS). Active (A) or disrupted (D, in red) components are indicated. An active methyltransferase requires both the methylase and the specificity unit to be active. In that case, if the restriction enzyme is not functional, the methylase is tagged as 'orphan' (blue). *hsdR*: Restriction endonuclease; *hsdM*: Methyltransferase; *hsdS*: Specificity unit.

Supplementary Table 6. Type II restriction-modifications systems (RMSs). Active (A) or disrupted (D, in red) components are indicated. An active methyltransferase requires both the methylase and the specificity unit to be active. In that case, if the restriction enzyme is not functional, the methylase is

tagged as 'orphan' (blue). R: Restriction enzyme; M: Methyltransferase. The Bonferroni-corrected p-values of a Mann-Whitney test between the IPD values of methylated cytosines of each motif (as indicated in REBASE) and the distribution of unmethylated cytosines for each strain are shown. \*\*\*\*p-value < 0.001; \*\*\*p-value<0.001; \*\*p-value<0.01; \*p-value<0.05; ns: non-significant.

Supplementary Table 7. Type III NgoAX RM system. Active (A) or disrupted (D, in red) components are indicated. The number of CCCAA repeats and the DNA recognition domain (DRD) allele are shown, which are the source of variability in the Mod unit. An active methyltransferase requires both the methylase and the specificity unit to be active. In that case, if the restriction enzyme is not functional, the methylase is tagged as 'orphan'. Res = Restriction enzyme, Mod = Methyltransferase.

Supplementary Table 8. Type III NgoAXII RM system. Active (A) or disrupted (D, in red) components are indicated. The number of AGCC repeats and the DNA-recognition domain (DRD) allele are shown, which are the source of variability in the Mod unit. An active methyltransferase requires both the methylase and the specificity unit to be active. In that case, if the restriction enzyme is not functional, the methylase is tagged as 'orphan'. Res = Restriction enzyme, Mod = Methyltransferase.

Supplementary Table 9. Number of GCAGA and AGAAA and GCAGAAATTY motifs in each strain and number and proportion of those with an IPD>3 in the following underlined bases. GCAGA is the motif associated to NgoAXII detected by the PacBio SMRT analysis pipeline (NCTC13798, NCTC13799 and WHO N), AGAAA is the predicted motif for this methylase in REBASE and GCAGAAATTY is the motif resulting from overlapping GCAGA, AGAAA and the recognition sequence of the Apol enzyme (5'-RAATTY-3').

Figure 1. Genome organization of restriction-modification systems in *N. gonorrhoeae*. VSR: Very short patch repair endonuclease.

Figure 2. Active and inactive restriction-modification systems (RMS) in the 25 *N. gonorrhoeae* strains under study. Phandango visualization of active and inactive subunits of each RMS. The Dam methylase is marked in a different colour (see legend). An extra column ('?') in each RMS is plotted showing if methylation is expected for each strain. The tree on the left is a maximum likelihood reconstruction of the core genome SNPs of the 25 strains.

Figure 3. Protein alignment of the *hsdS* specificity unit of the Type I NgoAV RMS. The three main sources of variation within the unit are shown as a premature stop codon in position 62 of the protein that causes inactivation of the system, a variable number of ATLE repeats and a downstream stop codon in position 201 which, in combination, are related to a different recognition motif. The strains with a complete *hsdS* are labelled in green (left), and the detected methylated motifs using PacBio are shown in different colours on the right. '-' represents gaps and '\*' stop codons.

Figure 4. Comparison of the per-base distribution of IPD ratios for the 5'-GCAGA-3' motif detected in this study as associated with NgoAXII, the previously inferred 5'-AGAAA-3' and the motif overlapping among the previous two and the recognition pattern of the Apol restriction enzyme, 5'-RAATTY-3' (5)



in the three strains with a functional NgoAXII system. Two boxplots are shown per base, corresponding to the IPD ratios of the forward (red) and reverse (blue) strands.

Supplementary Figure 1. Distribution of the IPD ratios for methylated and unmethylated bases. 6mA, 4mC and 5mC IPD values were extracted from the target bases in the final list of curated motifs. Those below an IPD ratio of 2 were excluded to minimize the inclusion of unmethylated motifs. Values for each unmethylated base correspond to a random sample of 10,000 sites per strain outside all the non-redundant predicted motifs detected by the PacBio SMRT pipeline. Dashed vertical lines and numbers indicate the mean value of each distribution.

Supplementary Figure 2. IPD ratio values for each base in all instances of the 5'-GACN[6]TGC-3' motif, target of the Type I NgoAV RMS, in each of the 25 *N. gonorrhoeae* strains included in the study. Per-base distribution of IPD ratio is shown as two separate boxplots containing the values in the forward (red) and reverse (blue) strand.

Supplementary Figure 3. IPD ratio values for each base in all instances of the 5'-GACN[7]TGC-3' motif, target of the Type I NgoAV RMS, in each of the 25 *N. gonorrhoeae* strains included in the study. Per-base distribution of IPD ratio is shown as two separate boxplots containing the values in the forward (red) and reverse (blue) strand.

Supplementary Figure 4. IPD ratio values for each base in all instances of the 5'-GCAN[8]TGC-3' motif, target of the Type I NgoAV RMS, in each of the 25 *N. gonorrhoeae* strains included in the study. Per-base distribution of IPD ratio is shown as two separate boxplots containing the values in the forward (red) and reverse (blue) strand.

Supplementary Figure 5. IPD ratio values for each base in all instances of the 5'-GAGN[5]TAC-3' motif, target of the Type I NgoAXVII RMS, in each of the 25 *N. gonorrhoeae* strains included in the study. Per-base distribution of IPD ratio is shown as two separate boxplots containing the values in the forward (red) and reverse (blue) strand.

Supplementary Figure 6. IPD ratios for methylated cytosines (mCs) in Type II RMS. The distribution of IPD ratios for the forward (red) and reverse (green) mCs as stated in REBASE are shown as boxplots together with a random 10k subsampling of unmethylated cytosines for each strain (blue). Statistical significance of the comparison among the distribution of the forward and reverse methylations with that from a random 10k subsampling of unmethylated cytosines is indicated next to the boxplots. \*\*\*\*p-value < 0.0001; \*\*\*p-value<0.001; \*\*p-value<0.01; \*p-value<0.05; ns: non-significant.

Supplementary Figure 7. IPD ratio values for each base in all instances of the 5'-AAANCGGTTNNC-3' motif in each of the 25 *N. gonorrhoeae* strains included in the study. Per-base distribution of IPD ratio is shown as two separate boxplots containing the values in the forward (red) and reverse (blue) strand.

Supplementary Figure 8. IPD ratios for the underlined forward and reverse cytosines in the 5'-AAANCGGTTNNC-3' motif. These values were compared to the distribution of a random 10k

subsampling of unmethylated cytosines for each strain. Statistical significance is shown above the boxplots.

\*\*\*\*p-value < 0.0001; \*\*\*p-value<0.001; \*\*p-value<0.01;\*p-value<0.05; ns: non-significant.

Supplementary Figure 9. IPD ratio values for each base in all instances of the 5'-GATC-3' motif, target of the Type II NgoAXI RMS, in each of the 25 *N. gonorrhoeae* strains included in the study. Per-base distribution of IPD ratio is shown as two separate boxplots containing the values in the forward (red) and reverse (blue) strand.

Supplementary Figure 10. IPD ratio values for each base in all instances of the 5'-GGTGA-3' motif, target of the Type II NgoAXVI RMS, in each of the 25 *N. gonorrhoeae* strains included in the study. Per-base distribution of IPD ratio is shown as two separate boxplots containing the values in the forward (red) and reverse (blue) strand.

Supplementary Figure 11. Inference of the methylated cytosine in the reverse strand of 5'-GGTGA-3' (NgoAXVI). Bonferroni-corrected p-values are shown in logarithmic scale for the comparison between the distribution of IPD values for the underlined cytosine and a random 10k subsampling of unmethylated cytosines for each strain. The cytosine with a higher significance and thus, candidate to be methylated, is the second one.

Supplementary Figure 12. IPD ratio values for each base in all instances of the 5'-CCACC-3' motif, target of the Type III NgoAX RMS, in each of the 25 *N. gonorrhoeae* strains included in the study. Per-base distribution of IPD ratio is shown as two separate boxplots containing the values in the forward (red) and reverse (blue) strand.

Supplementary Figure 13. IPD ratio values for each base in all instances of the 5'-GCAGA-3' motif, target of the Type III NgoAXII RMS, in each of the 25 *N. gonorrhoeae* strains included in the study. Per-base distribution of IPD ratio is shown as two separate boxplots containing the values in the forward (red) and reverse (blue) strand.

Supplementary Figure 14. Maximum likelihood phylogenetic reconstruction of representative Dam protein sequences available in genomes of the *Neisseria* genus in the RefSeq database (accessed on 21st January 2019). These representative sequences (WP\_\*) were obtained using the Identical Protein Groups (IPG) tool in NCBI. Coloured strips represent the species in which these sequences are found. Those appearing only once have been coloured in grey as 'Others' but can be accessed using the ID on the tip labels in <https://www.ncbi.nlm.nih.gov/ipg>.

Supplementary Figure 15. Circular plot (58) representing the annotation of the WHO G pConjugative plasmid, which shows the presence of an orphan DNA methylase that has a 100% protein identity with two entries in REBASE.

## Acknowledgements

This work was supported by Wellcome grant number 098051 (LSB, SRH and JP) and the Foundation for Medical Research at Örebro University Hospital, Örebro, Sweden (DG and MU). We thank them for the financial support.

## Authors contributions statement

All co-authors planned and designed the study. LSB and SRH performed the bioinformatic analyses, with support from DG, MU and JP. LSB and SRH drafted the manuscript and all co-authors commented and approved the final version.

## Additional information

**Supplementary information** is available at XXX.

**Competing interests.** The authors declare no competing interests.

**Accession codes.** Stated in the Data Availability section.

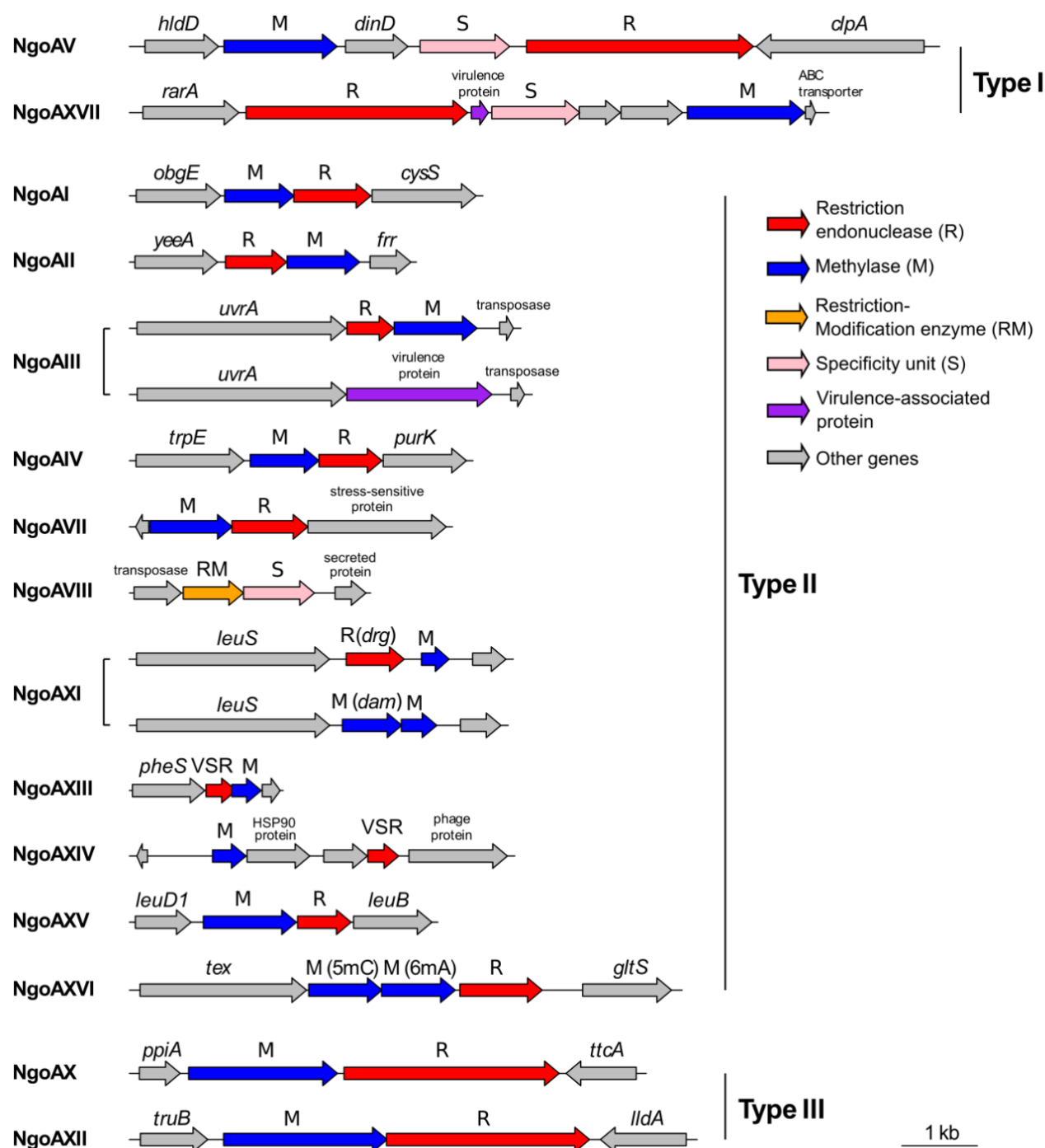


Figure 1. Genome organization of restriction-modification systems in 25 *N. gonorrhoeae* strains. VSR: Very short patch repair endonuclease.

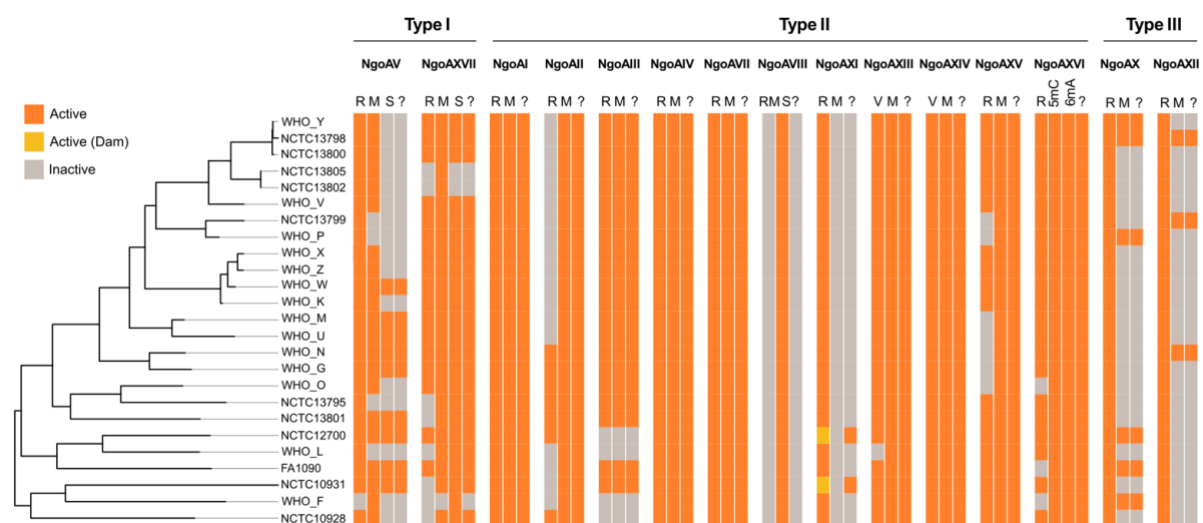


Figure 2. Active and inactive restriction-modification systems (RMS) in the 25 *N. gonorrhoeae* strains under study. Phandango visualization of active and inactive subunits of each RMS. The Dam methylase is marked in a different colour (see legend). An extra column ('?') in each RMS is plotted showing if methylation is expected for each strain. The tree on the left is a maximum likelihood reconstruction of the core genome SNPs of the 25 strains.

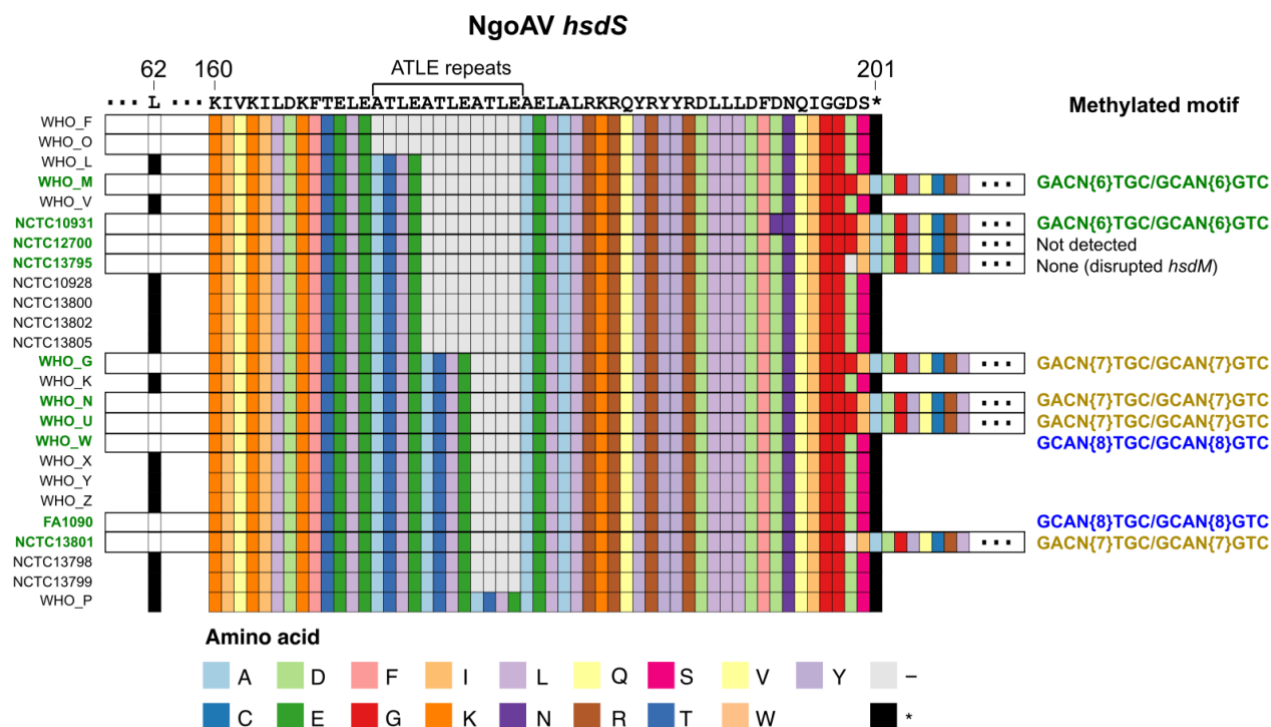


Figure 3. Protein alignment of the *hsdS* specificity unit of the Type I NgoAV restriction-modification systems (RMS). The three main sources of variation within the unit are shown as a premature stop codon in position 62 of the protein that causes inactivation of the RMS, a variable number of ATLE repeats and a downstream stop codon in position 201 which, in combination, are related to a different recognition motif. The strains with a complete *hsdS* are labelled in green (left), and the detected methylated motifs using PacBio are shown in different colours on the right. '-' represents gaps and '\*' stop codons.



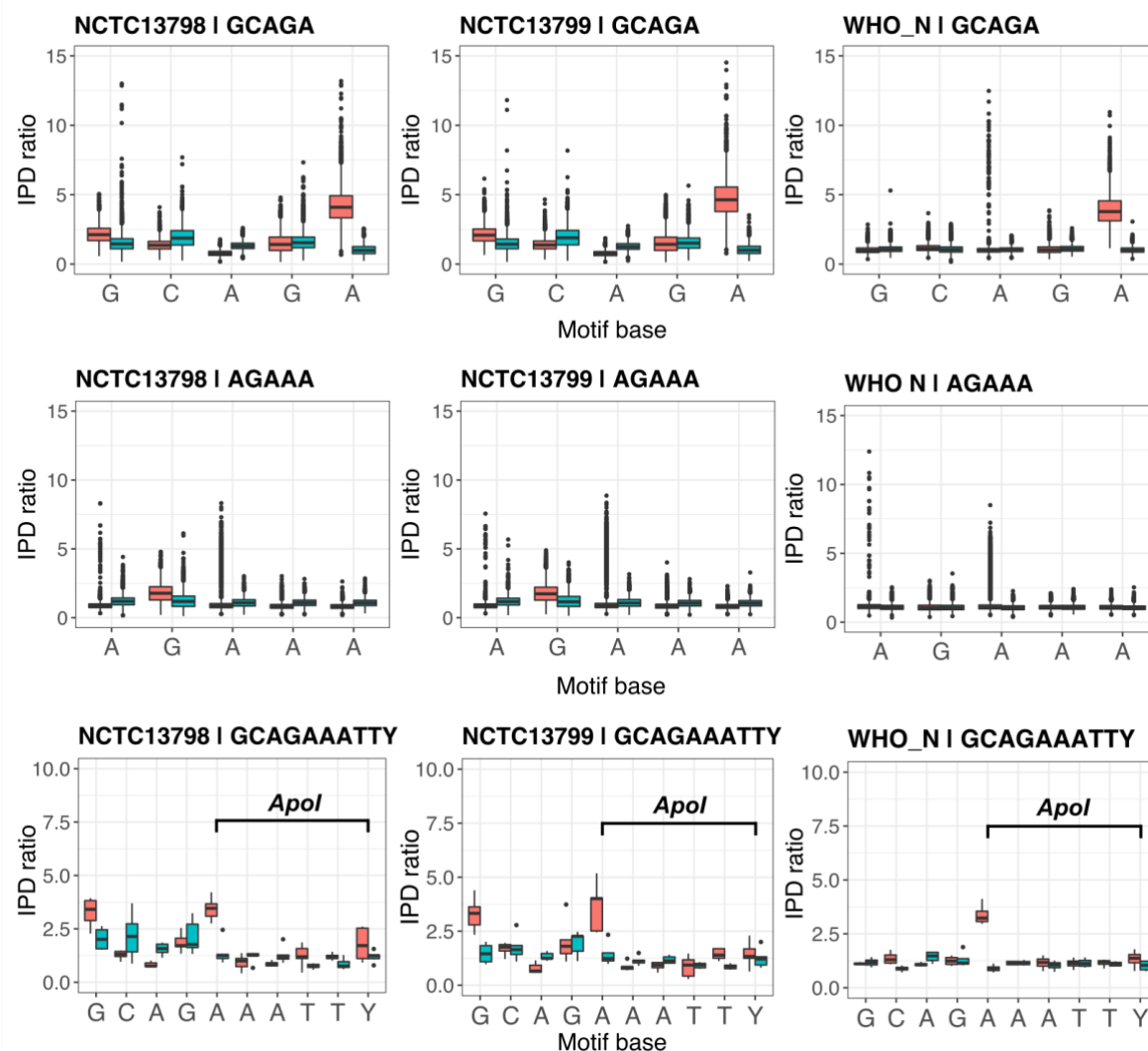


Figure 4. Comparison of the per-base distribution of IPD ratios for the 5'-GCAGA-3' motif detected in this study as associated with NgoAXII, the previously inferred 5'-AGAAA-3' and the motif overlapping among the previous two and the recognition pattern of the Apol restriction enzyme, 5'-RAATTY-3' <sup>5</sup> in the three strains with a functional NgoAXII system. Two boxplots are shown per base, corresponding to the IPD ratios of the forward (red) and reverse (blue) strands.

Table 1. Restriction-Modifications systems (RMS) and target motifs in *N. gonorrhoeae*. Those detected as methylated by the PacBio SMRT pipeline are marked as such and the fraction of modified motifs indicated. For those not directly detected by the pipeline, the associated motifs in REBASE are shown and marked as detected if a significant difference in IPD ratios is found between the target base and the unmethylated equivalent by a Mann-Whitney test.

RMS type	RMS	RMS subtype (REBASE)	Detected motif	Methylation type	Detected ?	Detected by PacBio or REBASE	Fraction (%)
I	NgoAV	gamma	5'-GACN[6]TGC-3' 3'-CTGN[6]ACG-5'	Two-strand 6mA	Yes	PacBio	100
			5'-GACN[7]TGC-3' 3'-CTGN[7]ACG-5'	Two-strand 6mA	Yes	PacBio	100
			5'-GCAN[8]TGC-3' 3'-CGTN[8]ACG-5'	Two-strand 6mA	Yes	PacBio	99.9
	NgoAXVII	gamma	5'-GAGN[5]TAC-3' 3'-CTCN[5]ATG-5'	Two-strand 6mA	Yes	PacBio	100
			5'-YGTABYNNNCCC-3'	One-strand 6mA	Yes	PacBio	59.2
			5'-GAGNNNNNGTYAC-3' 3'-CTCNNNNNCARAG-5'	Two-strand 6mA	Yes	PacBio	83.3
II	NgoAI	P	5'-RGC <sup>C</sup> GCY-3' 3'-YCG <sup>C</sup> GGR-5'	Two-strand 5mC	Yes	REBASE	-
	NgoAII	P	5'-GG <sup>C</sup> CC-3' 3'-CC <sup>C</sup> GG-5'	Two-strand 5mC	Yes	REBASE	-
	NgoAIII	P	5'-CC <sup>C</sup> GCGG-3' 3'-GGC <sup>C</sup> GC-5'	Two-strand 4mC	Yes	PacBio	50.4
	NgoAIV	P	5'-G <sup>C</sup> CGGC-3' 3'-CGGC <sup>C</sup> G-5'	Two-strand 5mC	Yes	REBASE	-
	NgoAVII	S	5'-G <sup>C</sup> CGGC-3' 3'-CGC <sup>C</sup> G-5'	Two-strand 5mC	Yes	PacBio	36.3-39.7
	NgoAVIII	G,S	5'-GACN[5]TGA-3' 3'-CTGN[4]ACT-5'	Two-strand 6mA	No	REBASE (base undetermined)	-
	NgoAXI (Dam)	Beta	5'-GATC-3' 3'-CTAG-5'	Two-strand 6mA	Yes	PacBio	98-100
	NgoAXIII	-	Unknown	Unknown	No	REBASE	-
	NgoAXIV	P	5'-CC <sup>C</sup> GG-3' 3'-GG <sup>C</sup> C-5'	Two-strand 5mC	Yes	REBASE	-
	NgoAXV	P	5'-GGNN <sup>C</sup> C-3' 3'-CC <sup>C</sup> NNGG-5'	Two-strand 5mC	Yes	REBASE	-

	NgoAXVI	S	5'-GGTGA <u>A</u> -3' 3'-CC <u>C</u> ACT-5'	Forward 6mA Reverse 5mC	Yes	PacBio	93.1- 99.9
II *	Unknown	Unknown	5'-AAAN <u>C</u> GGTTNNC-3' 3'-TTTNGC <u>C</u> AANNG-5'	Two-strand m4C	Yes	PacBio	17.5- 20.6
III	NgoAX	-	5'-CC <u>A</u> CC-3'	One-strand 6mA	Yes	PacBio	97.3- 100
	NgoAXII	-	5'-GCAG <u>A</u> -3'	One-strand 6mA	Yes	PacBio	62.6- 98.2