

1 Deep sequencing of norovirus genomes defines evolutionary patterns in an urban tropical
2 setting

3

4 Matthew Cotten^a, Velislava Petrova^a, My Vu Tra Phan^b, Maia A. Rabaa^{b,c}, Simon J. Watson^a,
5 Swee Hoe Ong^a, Paul Kellam^{a,d,#} and Stephen Baker^{b,e,f}

6

7 The Wellcome Trust Sanger Institute, Hinxton, CB10 1SA, UK^a; Hospital for Tropical
8 Diseases, Wellcome Trust Major Overseas Programme, Oxford University Clinical Research
9 Unit, Ho Chi Minh City, Vietnam^b; Centre for Immunity, Infection and Evolution, University of
10 Edinburgh, Edinburgh, UK^c; Division of Infection & Immunity, University College London,
11 London, UK^d; The London School of Hygiene and Tropical Medicine, London, UK^e; Centre for
12 Tropical Medicine, University of Oxford, Oxford, UK^f;

13

14 #Address correspondence to: Professor Paul Kellam, pk5@sanger.ac.uk

15 **Running head:** Norovirus genomics in Vietnam

16 Abstract word count: 165

17 Text word count: 4809

18

19 **Key words:**

20 Genomics, deep sequencing, norovirus, diarrhea, evolution

21

22 **ABSTRACT**

23 Norovirus is a highly transmissible infectious agent that causes epidemic gastroenteritis
24 in susceptible children and adults. Norovirus infections can be severe and can be initiated from
25 an exceptionally small number of viral particles. Detailed genome sequence data are useful for
26 tracking norovirus transmission and evolution. To address this need we have developed a
27 whole-genome deep sequencing method that generates entire genome sequences from small

28 amounts of clinical specimens. This novel approach employs an algorithm for reverse
29 transcription and PCR amplification primer design using all publically-available norovirus
30 sequence data. Deep sequencing and *de novo* assembly were used to generate norovirus
31 genomes from a large set of diarrheal patients attending 3 hospitals in Ho Chi Minh City,
32 Vietnam, over a 2.5-year period. Positive selection analysis and direct examination of protein
33 changes in the virus over time identified codons in the regions encoding proteins VP1, p48
34 (NS1-2) and p22 (NS4) under positive selection and expands the known targets of norovirus
35 evolutionary pressure.

36

37 **Importance**

38 The high transmissibility and rapid evolution rate of norovirus, combined with short-lived
39 host immune responses, are thought to be responsible for the virus causing a majority of
40 pediatric viral diarrhea cases. The evolutionary patterns of this RNA virus have only been
41 described in detail for a portion of the virus genome and never from a detailed urban tropical
42 setting. We have developed robust deep sequencing methods for generating complete
43 genome sequences directly from small amounts of patient fecal material. We use this method
44 to provide a detailed sequence description of the noroviruses circulating in three Ho Chi Minh
45 City hospitals over a 2.5-year period. The study identified patterns of virus change in known
46 sites of host immune response and identified three additional regions of the virus genome
47 under selection that were not previously recognized. In addition, the methods described here
48 provide a robust full-genome sequencing platform for community-based virus surveillance.

49

50 **INTRODUCTION**

51 Norovirus is a non-enveloped positive-sense single-stranded RNA virus of
52 approximately 7.5-7.7 kb in length (reviewed in reference (1)). The viral genome is organized
53 into three (or four in the case of murine norovirus (2)) open reading frames (ORFs), encoding
54 several structural and non-structural proteins. The ORF1 encodes for a large polyprotein,

55 which is proteolytically cleaved into 6 non-structural proteins, including N-terminal p48 protein
56 (NS1-2), an NTPase protein (NS3), a 3A-like p22 protein (NS4), a viral protein-genome linked
57 VpG protein (NS5), a 3C-like protease 3CLpro protein (NS6) and an RNA-dependent RNA
58 polymerase RdRp (NS7). Note that the nomenclature for the NS proteins is currently in flux
59 and both existing names have been included (3). The ORF2 overlaps ORF1 by a short region
60 and encodes the major capsid protein VP1, comprising an S (shell) domain connecting the two
61 P (protruding) sub-domains, P1 and P2 with the P2 domain binding to histo-blood group
62 antigens (HBGAs) on target host cells. ORF3, located at the 3' end of the genome, encodes
63 the minor capsid protein VP2.

64 Norovirus comprises one of the genera in the Caliciviridae family of viruses, and can be
65 further classified into different genogroups (reviewed in (1)). Noroviruses are known to cause
66 diseases in humans (genogroup GI, GII and GIV) and a number of other mammals including
67 porcine (GII), ovine/bovine (GIII), canine (GIV) and murine (MNV, forming a distinct GV)
68 viruses(4-12).

69 In humans, norovirus is a highly infectious pathogen that causes a severe
70 gastrointestinal disease in susceptible individual after the ingestion of an exceptionally small
71 number of viral particles. The virus is so infectious that the probability of symptomatic disease
72 from a single norovirus virion has been estimated to be as high as 0.5 (13). A determination of
73 the infectious dose required to infect 50% of subjects (ID_{50}) estimates this as 1,000-3,000 virus
74 genome equivalents(14). A typical norovirus infection can result in profuse volumes of feces
75 and vomitus containing 10^6 - 10^9 stable, non-enveloped virions per ml of excreta, creating
76 almost infinite opportunities for onward transmission and additional infections. An inability to
77 culture human noroviruses in a laboratory prevents the testing of inactivation and disinfection
78 methods and further complicates control efforts. These issues highlight some of the difficulties
79 in eliminating infectious norovirus from food supplies and the environment, and indicate the
80 need for the development of intelligent approaches to prevent norovirus transmission and
81 infection.

82 An effective approach to controlling norovirus may be to understand how norovirus
83 evades the human immune system and use this information to develop novel therapeutic
84 options. Norovirus infection in a “healthy” individual is typically short and self-limiting, which
85 results in transient or long-lived immunity (15, 16). No approved drugs exist to block virus
86 replication. Accordingly, public health measures to identify and eliminate sources of infection
87 or behavior leading to virus spread are warranted (17, 18). The utility of viral sequencing to
88 track norovirus in transmission studies has been explored with fragments of the viral genome
89 (19-22). As a consequence of the speed of disease onset and high transmissibility, the number
90 nucleotide and amino acid sequence changes within a local outbreak may be rare so the
91 sequencing of larger genomic fragments should provide greater resolution for defining
92 transmission patterns.

93 The natural duration and specificity of immune responses to norovirus are difficult to
94 measure due to the lack of a cell culture system for norovirus neutralization studies and the
95 inability to grow a defined virus for such trials (reviewed in (16, 23)). The duration of norovirus
96 immunity may be limited due to the short period of a typical infection and a corresponding
97 short exposure to viral antigens. Periodic population-level replacement of norovirus lineages
98 with viruses with surface residues under positive selection is evidence of immune response-
99 driven antigenic change and suggests that these immune responses are of sufficient strength
100 to drive viral evolution(24-26). Immune studies have identified blockade epitopes in VP1, the
101 major capsid protein. These epitopes are important for interacting with histo-blood group
102 antigens (HBGAs) on target host cells; high titers of antibodies that block virus-like particles
103 binding to HBGAs correlate with protection to norovirus challenge (27-29).

104 Diarrheal diseases are a serious health problem, especially in developing countries
105 when combined with nutritional problems, co-infection with other pathogens, crowding and
106 limited access to health care. It is clear that norovirus and rotavirus are frequently associated
107 with diarrhea in this setting (30) and it is essential to closely follow the local evolution of
108 norovirus. We describe here a method for deep sequencing the approximately 7,500

109 nucleotide norovirus RNA genome directly from patient material and used this method to
110 provide a detailed description of genome and community-wide norovirus evolution.

111

112

113 **MATERIALS AND METHODS**

114 **Primer design.** Primers were designed using Python algorithms to identify highly
115 conserved primer targets in the appropriate genome locations. Briefly the algorithm takes as
116 input all complete genome sequence of human norovirus available in GenBank (January 2012,
117 260 GII.4 entries, 5 GI entries, total sequence 1.9×10^6 nt). A counting method was employed
118 to identify all highly conserved primer-like sequences with a G+C percentage between 30 and
119 75 %, a calculated melting temperature between 55 and 59 °C, and with no single nucleotide
120 comprising greater than 40 % of the sequence. The norovirus genome was divided into three
121 overlapping 2.5-3 kb amplicons and the highest frequency primer site in the first and last 800
122 nt of each amplicon were selected. Finally the primers were used in a virtual PCR reaction to
123 determining the binding behavior of the primer set all available full norovirus genomes (see
124 Figure 1). Primer details are summarized in Table 1.

125 **Sample collection.** Stool samples were obtained as part of a larger study examining
126 causes of pediatric diarrhea in subjects presenting to Children's Hospital 1, Children's Hospital
127 2 and The Hospital for Tropical Diseases, Ho Chi Minh City (HCMC), Vietnam (30, 31).
128 Additional samples came from an *ad hoc* enrollment of children admitted to Children's Hospital
129 2 with potentially hospital-acquired norovirus diarrhea or prolonged norovirus incubation. In the
130 *ad hoc* collection, pediatric patients were admitted to the hospital due to diseases other than
131 diarrheal diseases and had no diarrhea when presenting at the hospital. These patients
132 eventually developed diarrhea after 48 hours of hospital admission and diarrhea lasted for at
133 least 3 days after the onset. Ethical approval was granted from the institutional ethical review
134 boards and the University of Oxford Tropical Research Ethics Committee (OxTREC No. 0109).

135 **Generation of amplified cDNA for deep sequencing.** For RNA extraction, 140 μ l of
136 each stool specimen was subjected to automated extraction into a final 50 μ l elution using the
137 MagNA Pure 96 automated extraction machine according to manufacturer's instruction
138 (Roche). Reverse transcription was performed as previously described(32). Briefly, a primer
139 mix was prepared separately for each amplicon: the reverse primers for the amplicon were

140 pooled in an equimolar ratio and water up to 7 μ l primer mix (7.6 pmol of each primer;
141 0.38pmol/ μ l per reaction). Extracted norovirus RNA was diluted in 1:10 in water; 5 μ l of this
142 dilution was added to the primer mix and heated for 5 min at 65 °C and immediately
143 transferred to an ice block for 1 min. An enzyme mix was then added to each reaction and
144 mixed by pipetting. Each 20 μ l reaction contained 4 μ l 5x First Strand buffer (250 mM Tris-HCl
145 (pH 8.3), 375mM KCl, 15 mM MgCl₂), 1 μ l 0.1 M DTT, 1 μ l 10mM dNTPs, 1 μ l RNase Inhibitor
146 (Promega, 40U/ μ l), 1 μ l SuperScript III RT (Life Technologies, 200U/ μ l). Reverse transcription
147 was performed at 50°C for 60 min, followed by 70 °C for 15 min.

148 **PCR amplification.** Amplification was performed with primer mix solutions prepared for
149 each amplicon. Primer mix (per 25 μ l reaction): the forward and reverse primers from each
150 amplicon were pooled together in 1.5:1 ratio (1.9 pmol of each forward primer and 1.26 pmol
151 of each reverse primer; 0.08 pmol/ μ l and 0.05 pmol/ μ l respectively). A 5 μ l aliquot of the
152 reverse transcription reaction for the each amplicon was used as a template for the PCR step.
153 The following thermal cycling conditions were used: enzyme activation: 98°C for 30 sec;
154 cycling (35 cycles): 98°C for 10 sec, 53°C for 30 sec, 72°C for 3.0 min; final extension: 72°C
155 for 10 min; hold: 4°C.

156 **Sequencing and genome assembly.** Pooled amplicons for each sample
157 (approximately 1.2 μ g) were individually indexed and subjected to sequencing with Illumina
158 MiSeq (33, 34) to generate approximately 300,000 reads of 149 nt per sample (median value
159 302904 reads). All reads were processed using QUASR (35) to remove sequencing adapters
160 and index sequences and to trim primer sequences present within a fixed distance from the 5'
161 or 3' end of a read. Reads were then trimmed from the 3' end to reach a minimum median
162 Phred quality score of 35 and reads below 125 nt in length were removed. After primer
163 trimming and quality control for each sample, *de novo* assembly using SPAdes (36) was used
164 to generate full norovirus genomes. Intact open reading frames (ORFs) were checked using
165 Python scripts as a measure of correct genome assembly.

166 **Recombination detection.** The 119 complete genome of all GII noroviruses from this
167 study and from global data (retrieved from GenBank) were manually aligned using Se-AL v2.0
168 (<http://tree.bio.ed.ac.uk/software/seal/>). Only full-length sequences with information on sample
169 date and location were included in this analysis. Potential presence of recombination in these
170 complete sequences was screened using the Recombination Detection Program version 4
171 (RDP4) software (37). RDP, GENECONV, 3SEQ and MAXCHI methods were employed for
172 primary screening and the BOOTSCAN and SISCAN methods for automatic checking of the
173 recombination signals, as described previously (38). The automask xfunction in RDP4 was
174 selected for optimal recombination detection, *i.e.* one representative strain within each group
175 of similar sequences was examined during the primary/exploratory search for recombination
176 signals while the remaining sequences within groups of sequences with high similarity were
177 automatically masked. Using this method, masked sequences were examined for the
178 presence of recombination if the program detected a recombination signal in the
179 representative unmasked sequence. Each test of recombination used a 400 nt sliding window,
180 and any recombination signals with significant p-values for ≥ 3 test parameters were
181 considered as potential recombination events. A further analysis on these potential
182 recombinants, comparing tree topologies with likelihood (Shimodaira-Hidetoshi, SH test) was
183 employed to determine which of the testing strains were likely to be true recombinants and
184 which were not. All intra-ORF recombinant strains (GenBank accession numbers EU921388,
185 AB541275, GU991355 and AB541254) were excluded from the estimation of selection rates.

186 **Phylogenetic analysis.** An alignment of non-recombinant sequences, including all full
187 genomes determined in this analysis and global background sequences obtained from
188 GenBank, was utilized to reconstruct evolutionary relationships among norovirus sequences.
189 A phylogenetic tree was inferred using aligned nucleotide sequences, employing a maximum
190 likelihood (ML) method in RaxML (39), under the GTR+ Γ model of substitution, determined to
191 be the best fit model to our data using jModelTest version 2.1.1 (40). Tree topology was

192 assessed through bootstrapping with 1,000 pseudo-replicates. The resulting phylogenetic tree
193 was visualized and edited in FigTree v1.4.0 (<http://tree.bio.ed.ac.uk/software/figtree/>).

194 **Evolutionary rate estimations.** Evolutionary rates were estimated using a Bayesian
195 Markov chain Monte Carlo (BMCMC) method implemented in BEAST version 1.7.2 (41). A
196 relaxed uncorrelated lognormal molecular clock was employed to account for lineage-specific
197 rates, and a GMRF Bayesian skyride coalescent (42) was used to model the population
198 dynamics. The relevant substitution models for each alignment were selected using
199 jModelTest version 2.1.1 (40). The mean evolutionary rate and the 95% upper and lower
200 highest posterior density (HPD) intervals were inferred from the posterior tree distribution
201 generated from the BMCMC runs using Tracer version 1.6
202 (<http://tree.bio.ed.ac.uk/software/tracer/>).

203 **Positive selection analysis.** To determine evolutionary patterns of norovirus, selection
204 analyses were performed on the regions encoding VP1, VP2, and ORF1-encoded p48 (NS1-2)
205 and p22 (NS4) proteins. Norovirus codons under selective pressure were first determined
206 using the Mixed Effects Model of Evolution (MEME; p -value < 0.05) (43) and Fast
207 Unconstrained Bayesian AppRoximation (FUBAR; posterior probability >0.9) (44) implemented
208 through the Datamonkey web browser (45). Codons that were found to be under positive
209 selection by either method were inspected at the sequence alignment, and those with no
210 evidence of polymorphisms were considered a false positive and discarded.

211 Ancestral sequences were reconstructed from the sequence alignment and inferred
212 phylogeny using the joint likelihood method implemented in Hyphy (46) under a GTR+ Γ model
213 of evolution.

214 The GenBank accession numbers for all new norovirus sequences reported here are
215 reported in Table 2. Also listed are the samples collection dates, the genetic cluster (see
216 Figure 2 and Table 4) and the European Nucleotide Archive accession numbers for the raw
217 sequence data. In addition, 89 GII.4 genomes from the same HCMC study were also publically
218 available in GenBank, with the following accession numbers: Cluster 1: KC409244,

219 KC409245, KC409246, KC409257, KC409258, KC409259, KC409260, KC409261, KC409262,
220 KC409264, KC409265, KC409266, KC409267, KC409268, KC409269, KC409270,
221 KC409271, KC409272, KC409273, KC409274, KC409275, KC409276, KC409277,
222 KC409279, KC409280, KC409281, KC409282, KC409283, KC409284, KC409285,
223 KC409286, KC409287, KC409288, KC409289, KC409290, KC409291, KC409293,
224 KC409294, KC409295, KC409296, KC409297, KC409298, KC409304, KC409305,
225 KC409306, KC409307, KC409308, KC409309, KC409310, KC409312, KC409313,
226 KC409314, KC409315, KC409318, KC175360, KC175365, KC175366, KC175371,
227 KC175373, KC175381, KC175388, KC175389, KC175390, KC175391, KC175392,
228 KC175393, KC175394, KC175395, KC175396, KC175406, KC175407, KC175408,
229 KC175409, KC175410
230 Cluster 3: KC409256, KC409263, KC409278
231 Cluster 4: KC409240, KC409241, KC409242, KC409243, KC409299, KC409301, KC409302,
232 KC409303, KC175384, KC175385, KC175386, KC175387.
233
234
235

236 **RESULTS**

237 **Norovirus sequencing strategy.** A novel general strategy for designing PCR primers
238 was developed that would permit the production of complete norovirus genome sequences.
239 Deep sequencing of RNA virus genomes requires reverse transcription (RT) of viral RNA and
240 the amplification of the resulting cDNA that encompass the entire viral genome. Python
241 algorithms were used to process all available norovirus full genome data (265 full genomes,,
242 January 2012) and to select primer target sequences suitable for whole genome amplification.
243 Briefly the algorithm processes the norovirus sequence data into primer-sized sequences
244 trimmed to a calculated melting temperature. The frequency of each sequence in the entire set
245 is calculated, with high frequency sequences correlating with conserved sites across the viral
246 genome. The norovirus genome was divided into three overlapping amplicons, potential
247 primers were mapped to a reference genome and the highest frequency sequences mapping
248 within the terminal 800 nt of each amplicon were identified. Reverse complements of the
249 primers mapping to the 3' end were prepared. A virtual PCR was performed to examine the
250 potential function of the primers across all known full norovirus genomes, the output of such an
251 analysis is shown in Figure 1 (left panel) with blue markers indicating the position of the each
252 primer and grey bars indicating the expected PCR product. The actual function of the primer
253 set is demonstrated in Figure 1 (right panel) with each lane showing the PCR products from 14
254 samples, present by amplicon. Each reverse transcription reaction contained the 2 (or 3 for
255 amplicon 3) reverse primers for each amplicon 1, 2 or 3, and each PCR reaction contained the
256 2 (or 3 for amplicon 3) forward and reverse primers for amplicon 1, 2, or 3. Within these
257 samples, sample 7 failed, however the remaining 13 samples provided sufficient material for
258 deep sequencing.

259 A summary of the predicted performance of the norovirus primer set on all available
260 norovirus genomes is included in Table 1. All full length Norovirus GII genomes (TaxonID
261 142786, length 7000-8000 nt, 517 entries) or all Norovirus genomes (TaxonID 122929, length
262 7000-8000 nt, 753 entries) were retrieved from GenBank. The target sequence for each

263 primer was examined in these genome sets and the percentage of the genomes with a perfect
264 match to the target sequence or with a functional match (0-3 mismatches) to the target
265 sequence was reported. For the Norovirus GII genomes, the primers have a perfect match to
266 79% of the genomes, and a functional match (up to 3 mismatches) for 97% of the genomes.
267 For the complete set of Norovirus genomes (this includes all GI, GII and all animal
268 noroviruses), the primers have a perfect match to 65% of the genomes, and a functional match
269 (0-3 mismatches) for 82% of the genomes. These values and the details of the analysis as
270 well as the GC content and calculated T_m for all primers are listed in Table 1.

271 A summary of the performance of the norovirus primer set for amplifying and
272 sequencing 188 fecal-derived RNA samples is presented in Table 3. PCR success was
273 defined as obtaining the three amplicon-specific RT-PCR products of the predicted size with
274 sufficient yield for sequencing library preparation. The overall RT-PCR success rate was
275 78.2% (147 out of the 188 tested clinical samples). The most common genotype globally
276 (GII.4), had the highest PCR success rate (93.7%, 74 of 79 samples), followed by GII.6 (88%,
277 7 of 8 samples), GII.13 (83%, 5 of 6 samples), and GII.3 (77%, 26 of 34 samples). Much lower
278 amplification efficiency was observed for genogroup I (GI) strains, with successful PCR
279 genome amplification in only two of 10 tested samples. The high success with GII with respect
280 to GI strains (especially GII.4) was predictable given that GII.4 genomes dominate sequences
281 in public databases. Future primer sets could be reiteratively designed using targeted and
282 revised genome data sets.

283 **Norovirus diversity in HCMC.** Using the developed whole genome sequencing
284 technique, 112 novel GII norovirus genomic sequences were generated. In addition, 89 GII.4
285 genomes from the same HCMC study were also publically available in GenBank; these were
286 included in the following analysis for a total of 201 complete genomes, with collection dates
287 between April 2009 and December 2011. A phylogenetic analysis of the 201 genomes defined
288 eight genotypes of GII norovirus using maximum likelihood methods (Figure 2). Consistent
289 with previous characterization of norovirus infections in HCMC (30), and global norovirus

290 patterns, the most prevalent genotype GII.4 found in this study belonged to the GII.4 Den
291 Haag genotype lineage (Figure 2, Clusters 1, 2 and 3), which is most genetically similar to the
292 GII.4 Minerva_2006b partial sequence and Taiwan_2006 (GenBank JN400601).
293 Phylogenetically, the GII.4 strains in Cluster 4 (Figure 2A) were most closely related to GII.4
294 New_Orleans_2010 (GenBank JN595867), while GII strains in Cluster 5 were classified as
295 GII.P21_GII.3 and most closely related to the strain NV_Pune_2007 (GenBank EU921389). A
296 small number of strains belonged to genotype GII.Pg_GII.12 (Cluster 6), while viruses of the
297 GII.P7_GII.6 genotype fell into the 2 distinct lineages, Clusters 7 and 8 (Figure 2). Our
298 genotype assignment using phylogenetic reconstruction was consistent with the genotype
299 designation generated by the RIVM algorithm (47) (Table 4). Additionally, the relative
300 frequency of each genotype observed in the full genome set was similar to the frequencies
301 determined by My *et al.* (31) from a larger set of HCMC samples using ORF1 and 2 fragments
302 (Table 4), indicating that the generation of full genome sequences was not strongly influenced
303 by genotype-based selection biases. Viruses of the GII.4 Den Haag and GII.P21/GII.3
304 genotypes, (Cluster 1 and 5), were identified over two sampling periods, from 2009-2010 and
305 later in 2011, while other virus genotypes were detected only in the first sampling period.

306 The temporal occurrence of sampled noroviruses is shown in Figure 3, with samples
307 stratified by genotype cluster. The three GII.P4/GII.4(2006) genotypes (Clusters 1,2, and 3)
308 were present in the first half of 2009, with the GII.P4/GII.4(2010) genotype (Cluster 4, grey)
309 first appearing at the end of 2009. There was a pause in sampling in the first half of 2011,
310 followed by sampling in the second half of 2011. Reduced diversity was observed in 2011, with
311 only Clusters 3 and 5 sequenced from these samples. Changes in sampling protocols between
312 2010 and 2011 preclude inference on how this reduced diversity may relate to norovirus
313 epidemiology and evolution. However, the identification of clusters of phylogenetically-related
314 viruses undergoing *in situ* evolution in this region over the observation period allowed an
315 examination of evolutionary processes that may allow the continued transmission and
316 maintenance of viral lineages in the presence of population immune responses. Characterizing

317 such changes in the norovirus population may provide important clues about how the virus
318 evades host immunity.

319 **Evolutionary rates within each cluster.** A sufficient number of genomes were
320 available from Clusters 1, 4 and 5 for well-supported evolutionary rate estimations (Table 5).
321 Mean evolutionary rates of 6.15, 5.73 and 5.34 ($\times 10^{-3}$ substitutions per site per year) were
322 estimated from the full genomes for Cluster 1, 4 and 5. Figure 4 plots the rates for the GII.4
323 Cluster 1 viruses by the region of the genome used for each calculation.

324 ORF-specific rates estimated for the three genetic clusters show that the ORF1 regions
325 exhibited a reduced rate compared to those of the ORF2 (VP1) regions. For all three clusters
326 the ORF1 and ORF2 (VP1) regions showed modestly reduced rates relative to the full
327 genome, while the ORF3 (VP2) substitution rates for both Cluster 1 (8.99×10^{-3} substitutions
328 per site per year) and Cluster 5 viruses (7.38×10^{-3} substitutions per site per year) were higher
329 than those of the whole genome. The overlapping confidence intervals for these estimations
330 make these conclusions less secure. The amount of signal available for the Cluster 4 ORF3
331 was not sufficient to yield a reliable rate estimate.

332 The norovirus ORF1 encodes a large polyprotein containing the viral polymerase,
333 protease and several essential replicase components. Evolutionary rates were estimated
334 separately for these individual coding regions of the Cluster 1 ORF1 (Table 5, Figure 4). The
335 region encoding p22 (NS4) showed the highest levels of change (6.60 and 8.21×10^{-3}
336 substitutions per site per year, Figure 4), greater than the whole genome rates for Cluster 1
337 (6.15×10^{-3} substitutions per site per year). The enzymes (NTPase (NS3), protease and RdRp
338 (NS7)) and VP1 show modestly lower substitution rates than observed across the whole
339 genome.

340 **Amino acid changes in norovirus proteins.** The evolutionary patterns of four
341 norovirus-encoded proteins with the higher evolutionary rates were examined (VP1, VP2 and
342 p48 (NS1-2) and p22 (NS4)). An alignment of protein sequences ordered by time was used to
343 detect sustained versus sporadic changes in the protein relative to a reconstructed ancestral

344 sequence. Information about the biochemical properties of the protein was gathered from
345 published literature. Positive selection analysis was performed with the mixed effects model of
346 evolution (MEME) (43) or fast unconstrained Bayesian approximation (FUBAR) (44).

347 Cluster 1 VP1 showed changes in multiple patients relative to the ancestral sequence:
348 Q106R, S174P, N298D (in blockade epitope A), G340E, G393S (in blockade epitope D)
349 (Figure 5). Additional substitutions were seen at lower frequency suggesting evolution during
350 the course of transmission through HCMC. Position 298 in the blockade epitope A was found
351 to be positively selected using FUBAR, while both FUBAR and MEME identified position 106
352 within the shell domain (Figure 5) to be under positive selection (Table 6).

353 An alignment of VP2 protein sequences ordered by time was used to detect sustained
354 versus sporadic changes in the protein relative to the ancestral sequence. Several changes,
355 including T139M/A, I144V/T, and Y169H occurred in multiple HCMC Cluster 1 viruses with a
356 much higher frequency of changes in the internal region of the protein (Figure 6). It was
357 previously noted that changes in this region of VP2 (VP1 interacting domain (VP_ID) were
358 associated with changes in VP1(48). Both MEME and FUBAR identified VP2 codon 144
359 (marked with red asterisk in Figure 6) as being under positive selection.

360 The region encoding p22 (NS4) from the Cluster 1 viruses showed elevated
361 evolutionary rates relative to the full genome (Table 5, Figure 4). Analysis of all encoded p22
362 (NS4) molecules from Cluster 1 (Figure 7) showed amino acid changes relative to the
363 ancestral sequence. Substitutions were observed in multiple isolates suggesting neutral or
364 positive selection (I29V, E46D, N77S, R82K, T86S and D174V). Analysis of all encoded p48
365 (NS1-2) molecules from Cluster 1 (Figure 8) showed amino acid changes appearing in multiple
366 isolates suggesting neutral consequences with no constraints to limit change or positive
367 selection (D7V, N15D, R55K, V79T or A, and S184P. Both MEME and FUBAR identified p48
368 (NS1-2) codon 79 to be under positive selection (Table 6).

369

370

371 **DISCUSSION**

372 Our work outlines a strategy for full genome deep sequencing of norovirus directly from
373 fecal specimens and we have applied the strategy to characterize norovirus samples collected
374 across a clinical spectrum of pediatric norovirus infections in HCMC, Vietnam. An essential
375 component of the methods is a primer design algorithm that takes as input all available
376 sequence data for a virus and quickly provides a set of functional primers. The flexible design
377 of the primer design algorithm avoids a cumbersome alignment step in the process and
378 facilitates regular updates with new sequence data. This is essential to avoid perpetuating a
379 bias in the sequence data whereby sequences are only obtained if primers have functioned,
380 and primers are designed on antiquated data sets. The method showed a high success rate of
381 full genome sequencing of GII noroviruses, especially GII.4, which was predictable given that
382 GII.4 genomes dominated the sequence data set used to design the primers. Future primer
383 sets will be designed using more targeted and updated genomes sets and including more
384 sequence data from other genogroups.

385 Results using this method have provided a large set of norovirus genome sequences
386 derived from longitudinal samples from one location. At the start of this study, 265 full norovirus
387 genomes were available in GenBank and the current study added an additional 112 genomes.
388 The data allowed estimation of evolutionary rates for several genotypes, for full genomes as
389 well as for sub-genomic regions. The evolutionary pressures and the constraints to avoid
390 change are not expected to be uniform across the virus genome. Selection pressures are likely
391 to vary greatly depending on the function of the encoded proteins, with enzymatic and
392 structural regions more constrained than surface and immune-exposed or spacer regions with
393 less defined functions. ORF-specific substitution rates estimated for the three phylogenetic
394 clusters show that for the ORF1 regions exhibited a reduced evolutionary rate compared to
395 those of the ORF2 (VP1) regions rates.

396 Previous studies have estimated norovirus GII.4 and GII.3 VP1 capsid regions to evolve
397 at 5.1 to 5.8 x 10⁻³ substitutions per site per year (49) (50) (51), while the GII.4 polymerase

398 region was estimated to evolve between $4.33 - 8.98 \times 10^{-3}$ substitutions per site per year
399 depending on the data set used (49). Our estimates using the HCMC data are consistent with
400 these published figures. The evolutionary rate determined for GII.4 Cluster 1 was higher than
401 the estimated rates for GII.4 Cluster 4 and the GII.3 Cluster 5 viruses, which is perhaps due to
402 a greater number of Cluster 1 infections/unit time and thus a greater number of replication
403 events. Alternatively, the three virus genotypes might have intrinsically different replication
404 properties, polymerase fidelity, or immune selection pressure that result in the differing rates.

405 The norovirus sequence data obtained from this study allowed an analysis of the
406 evolutionary patterns of the second viral capsid protein VP2. The high evolutionary rates
407 reported here (Cluster 5: 7.38×10^{-3} substitutions per site per year, Cluster 1: 8.99×10^{-3}
408 substitutions per site per year) have not been observed previously, as this region is seldom
409 included in previous sequencing projects. The structure of VP2 is not defined although there is
410 evidence that the protein is interior to the VP1 shell and may be important for assembly of the
411 VP1 structure (52). The protein is moderately basic and the C-terminal half of the protein is
412 rich in serine and threonine residues (providing possible phosphorylation sites) and proline
413 residues (perhaps accounting for the inability to define a structure for this protein). Evidence
414 that changes in VP2 accompany changes in VP1 has been presented (48). Recently the
415 murine norovirus (MNV) VP2 has been shown to influence the host immune response to the
416 virus, with MNV1 VP2 interfering with antigen presenting cell (APC) function and MNV3 VP2
417 promoting the response (53). These observations identify a possible site of virus/host
418 interaction that could be a source of selective pressure. The evolutionary rates for the VP2-
419 encoding regions were found here to be much higher than the well-studied norovirus VP1
420 region and the higher rates are consistent with less-constrained protein product, stronger
421 selection pressures or both. Positive selection analysis across the VP2 region identified
422 position 144 to be under selection, this position region of the protein was previously found to
423 be involved in interactions of VP2 with VP1(52). A high evolutionary rate in a virus capsid
424 protein suggests a region of the virion experiencing immune selection. Vaccine development

425 efforts should take this accelerated rate of change into consideration when selecting
426 components for a vaccine.

427 Humoral immunity to norovirus (at least GII.4) may involve blockade antibodies that
428 bind and block the VP1 residues required for binding to HBGAs (16, 54, 55). The correlation of
429 high titer blockade antibodies with protection from gastroenteritis in challenge studies (29) and
430 the frequent evolution of these sites (blockade epitopes A, D and E) suggest that these amino
431 acid residues may be frequent targets of immune selection(55). Blockade epitope D may be
432 directly involved in HBGA binding (16, 54). Our observation of changes in VP1 position 298
433 epitope A, position 393 epitope D, and position 412 epitope E support these previous
434 conclusions. Several additional changes were located outside of blockade epitopes (S78G,
435 S174P, G340E, T502N, Figure 5). Further studies should investigate whether these are
436 founder effect changes of neutral consequence, or if they provide an advantage for the virus.

437 Similar mean evolutionary rates for full genome were found for Cluster 1, 4 and 5, with
438 95% confidence ranges largely overlapping. One might expect a higher evolutionary rate for
439 GII.4 viruses compared to GII.3 viruses if the 10-fold higher detection frequency than GII.3,
440 viruses directly reflects the community prevalence of these two infections. The similar full
441 genome rates suggests that either the number of active infections is not a large factor in rate,
442 or that the less frequently diagnosed GII.3 infection may be as frequent in the population as
443 GII.4, but does not appear as frequently in clinics.

444 The ORF1-encoded p22 (NS4) regions showed a higher evolutionary rate compared to
445 the full genome and p48 (NS1-2) codon 79 was found to be under positive selection. The
446 function of p22 (NS4) is not known but the protein has been observed to localize to the
447 Golgi/ER and influence the host secretory pathway with a centrally located MERES motif
448 required for localization (56, 57). The function of p48 (NS1-2) in norovirus infection is also
449 largely unexplored, although the protein is reported to localize to vesicles and has been
450 proposed to influence protein trafficking (58). The evidence that these viral proteins interact
451 with host proteins, combined with higher evolutionary rate or positive selection described here,

452 suggest that these proteins may be interacting with host restriction factors. Alternatively these
453 regions with higher rates of change could encode proteins with no constraint. Further studies
454 are needed to clarify this.

455 Extensive work has been done using the Feline Calicivirus (FCV) and MNV models to
456 elucidate the roles and interactions of the non-structural (NS1-7) and structural proteins (VP1
457 and VP2) in regulating virus replication and infectivity, as comprehensively reviewed (1).
458 However a functional profiling of human norovirus is not yet possible due to the lack of tissue
459 culture and animal models for human norovirus replication. The full genome sequences of
460 human norovirus available from this study provides valuable data on the spectrum of changes
461 in the viral proteins allowed by the virus while awaiting alternative models for functional
462 experiments.

463 This study has provided a description of norovirus evolution rates across HCMC over a
464 2.5-year period for full genomes as well as for subgenomic regions of the virus. We reveal for
465 the first time a higher evolutionary rate for three regions of the genome (VP2, p22 (NS4) and
466 p48 (NS1-2)) and provide evidence of positive selection in two coding regions (VP2 and p48
467 (NS1-2)). We suggest that these regions should be monitored for interactions with the host
468 that might be a source of selective pressure. Finally, we believe this study and the methods
469 we have described will provide a useful template for community-wide studies of full genome
470 evolution for many RNA virus pathogens.

471

472 **Acknowledgements.**

473 We thank the Sanger Institute Illumina C team for their sequencing support. We are
474 grateful to members of the HCMC study teams for their support and efforts for enrolling
475 patients and sample collections at three collaborative hospitals: Children's Hospital 1,
476 Children's Hospital 2 and Hospital for Tropical Diseases (HCMC, Vietnam). This work was
477 supported by the Wellcome Trust Strategic Award (VIZIONS – Vietnam Initiative on Zoonotic
478 Infections) program and the European Community's Seventh Framework Programme

479 (FP7/2007_2013) under the project EMPERIE, European Community grant agreement number
480 223498.
481

482 **TABLES**

483 **Table 1. Primers used.**

Primer	Sequence	Strand	Position ¹	Tm ²	GC_ fraction	Norovirus GII Genomes (517) ³		Norovirus Genomes (753) ⁶	
						0 MM ⁴	0-3 MM ⁵	0 MM ⁷	0-3 MM ⁸
UNP_47	GTGAATGAAGATGGCGTCTAAC	Plus	1	55.52	0.45	98	100	83	84
UNP_45	TCTAACGACGCTTCCGCTG	Plus	17	58.30	0.58	75	96	62	80
UNP_201R	GCAATGGCCACCTCCTCAT	Minus	2808	57.95	0.58	97	100	80	85
UNP_226R	TTGCCTCCTCCTTCACA	Minus	2850	58.21	0.55	92	99	76	82
UNP_339	GGCAAGAAGCACACAGCC	Plus	2660	57.48	0.61	88	92	73	78
UNP_1316	TGGTCCAAGCCACAAGTGG	Plus	2519	58.05	0.58	11	100	13	89
UNP_82	GACCTCTGGGACGAGGTTG	Minus	5150	57.41	0.63	87	96	70	80
UNP_135	CTCCACCAGGGCTTGAC	Minus	5271	57.63	0.63	89	94	73	78
UNP_2	GGGAGGGCGATCGCAAT	Plus	5049	57.57	0.65	88	96	72	79
UNP_23	TTGTGAATGAAGATGGCGTCGA	Plus	5079	58.53	0.45	56	100	42	84
UNP_100	GCCAGTCCAGGAGTCCAA	Minus	7447	56.43	0.61	74	97	61	83
UNP_44	GCACGGTTGAGACTGTGC	Minus	7418	57.27	0.61	84	98	69	82
UNP_20	CGAGGGGAGTCACGGGT	Minus	7493	58.34	0.71	86	97	70	83

484

485 **Footnotes.**

- 486 1. Primer mapping position in norovirus GII4 norovirus genome JQ613552
- 487 3.Tm (melting temperature) calculated using a Python script that approximates the Breslauer
- 488 method(59).
- 489 3. All GenBank entries (July 2014) for Norovirus GII (TaxonID 142786, length 7000-8000 nt)=
- 490 517 entries.
- 491 4. Percentage of Norovirus GII genomes (n=517) showing perfect homology to primer.
- 492 5. Percentage of Norovirus GII genomes (n=517) showing the target sequence for the primer
- 493 with up to 3 mismatches.
- 494 6. All GenBank entries (July 2014) for Norovirus (TaxonID 122929, length 7000-8000 nt) = 753
- 495 entries.
- 496 7. Percentage of Norovirus genomes (n=753) showing perfect homology to primer.
- 497 5. Percentage of Norovirus GII genomes (n=753) showing the target sequence for the primer
- 498 with up to 3 mismatches.

499

500

501 Table 2. GenBank and ENA accession numbers, genetic cluster and samples dates.

502

Virus	GenBank_Acc_No ¹	ENA_Acc_No ²	Cluster ³	Sample_Date ⁴
Hu_GII_10116_2009_VNM	KM198480	ERR212491	1	09-Jul-2009
Hu_GII_10054_2009_VNM	KM198481	ERR225641	1	21-May-2009
Hu_GII_10114_2009_VNM	KM198482	ERR212490	1	09-Jul-2009
Hu_GII_10313_2010_VNM	KM198483	ERR217285	4	22-Feb-2010
Hu_GII_30212_2009_VNM	KM198484	ERR223539	5	06-Oct-2009
Hu_GII_10148_2009_VNM	KM198485	ERR212498	2	11-Aug-2009
Hu_GII_C2H-18_2011_VNM	KM198486	ERR225628	3	30-Aug-2011
Hu_GII_10110_2009_VNM	KM198487	ERR212489	1	06-Jul-2009
Hu_GII_10325_2010_VNM	KM198488	ERR217290	4	26-Feb-2010
Hu_GII_10101_2009_VNM	KM198489	ERR212487	1	25-Jun-2009
Hu_GII_10002_2009_VNM	KM198490	ERR225635	1	04-May-2009
Hu_GII_30351_2009_VNM	KM198491	ERR138007	4	17-Dec-2009
Hu_GII_30448_2010_VNM	KM198492	ERR223547	6	29-Jan-2010
Hu_GII_30468_2010_VNM	KM198493	ERR223549	5	24-Feb-2010
Hu_GII_10247_2009_VNM	KM198494	ERR217278	1	10-Dec-2009
Hu_GII_10193_2009_VNM	KM198495	ERR138002	1	05-Oct-2009
Hu_GII_20419_2010_VNM	KM198496	ERR223554	5	01-Feb-2010
Hu_GII_10236_2009_VNM	KM198497	ERR217280	1	19-Nov-2009
Hu_GII_20088_2009_VNM	KM198498	ERR223553	8	28-Jul-2009
Hu_GII_20118_2009_VNM	KM198499	ERR212481	1	28-Aug-2009
Hu_GII_C2H-20_2011_VNM	KM198500	ERR225629	5	05-Sep-2011
Hu_GII_10173_2009_VNM	KM198501	ERR212503	1	11-Sep-2009
Hu_GII_10136_2009_VNM	KM198502	ERR212495	1	03-Aug-2009
Hu_GII_C2033_2010_VNM	KM198503	ERR212484	6	28-Jun-2010
Hu_GII_20151_2009_VNM	KM198504	ERR212467	1	16-Sep-2009
Hu_GII_20460_2010_VNM	KM198505	ERR223530	5	04-Mar-2010
Hu_GII_10199_2009_VNM	KM198506	ERR217283	1	20-Oct-2009
Hu_GII_C2007_2010_VNM	KM198507	ERR138011	4	02-Apr-2010
Hu_GII_20066_2009_VNM	KM198508	ERR212470	1	14-Jul-2009
Hu_GII_20479_2010_VNM	KM198509	ERR223531	5	16-Mar-2010
Hu_GII_10012_2009_VNM	KM198510	ERR225637	1	07-May-2009
Hu_GII_C2H-24_2011_VNM	KM198511	ERR225631	5	16-Sep-2011

Hu_GII_10062_2009_VNM	KM198512	ERR225642	1	29-May-2009
Hu_GII_20494_2010_VNM	KM198513	ERR223534	1	19-Mar-2010
Hu_GII_10079_2009_VNM	KM198514	ERR212486	1	11-Jun-2009
Hu_GII_C2H-31_2011_VNM	KM198515	ERR217269	3	28-Sep-2011
Hu_GII_10285_2010_VNM	KM198516	ERR217288	4	18-Jan-2010
Hu_GII_30399_2010_VNM	KM198517	ERR138008	1	11-Jan-2010
Hu_GII_10182_2009_VNM	KM198518	ERR212507	1	17-Sep-2009
Hu_GII_30082_2009_VNM	KM198519	ERR223537	8	23-Jun-2009
Hu_GII_10158_2009_VNM	KM198520	ERR212499	1	21-Aug-2009
Hu_GII_10176_2009_VNM	KM198521	ERR212504	1	14-Sep-2009
Hu_GII_20150_2009_VNM	KM198522	ERR212466	1	15-Sep-2009
Hu_GII_10204_2009_VNM	KM198523	ERR217287	1	29-Oct-2009
Hu_GII_10034_2009_VNM	KM198524	ERR225638	1	15-May-2009
Hu_GII_10163_2009_VNM	KM198525	ERR212501	1	28-Aug-2009
Hu_GII_10075_2009_VNM	KM198526	ERR225644	1	09-Jun-2009
Hu_GII_10074_2009_VNM	KM198527	ERR225643	1	09-Jun-2009
Hu_GII_C2H-25_2011_VNM	KM198528	ERR225632	5	20-Sep-2011
Hu_GII_C2H-27_2011_VNM	KM198529	ERR225633	5	21-Sep-2011
Hu_GII_20486_2010_VNM	KM198530	ERR223532	7	18-Mar-2010
Hu_GII_30116_2009_VNM	KM198531	ERR223538	7	09-Jul-2009
Hu_GII_10108_2009_VNM	KM198532	ERR212488	1	02-Jul-2009
Hu_GII_30241_2009_VNM	KM198533	ERR223540	1	26-Oct-2009
Hu_GII_30443_2010_VNM	KM198534	ERR223546	8	28-Jan-2010
Hu_GII_20092_2009_VNM	KM198535	ERR212474	1	31-Jul-2009
Hu_GII_20079_2009_VNM	KM198536	ERR212473	1	24-Jul-2009
Hu_GII_10137_2009_VNM	KM198537	ERR212496	1	04-Aug-2009
Hu_GII_10051_2009_VNM	KM198538	ERR225640	1	21-May-2009
Hu_GII_C2H-36_2011_VNM	KM198539	ERR217266	3	25-Oct-2011
Hu_GII_20145_2009_VNM	KM198540	ERR212464	1	11-Sep-2009
Hu_GII_20188_2009_VNM	KM198541	ERR138004	1	07-Oct-2009
Hu_GII_20094_2009_VNM	KM198542	ERR212476	2	03-Aug-2009
Hu_GII_20357_2009_VNM	KM198543	ERR138005	4	30-Dec-2009
Hu_GII_C2418_2010_VNM	KM198544	ERR138012	4	01-Nov-2010
Hu_GII_10195_2009_VNM	KM198545	ERR217284	1	13-Oct-2009
Hu_GII_20067_2009_VNM	KM198546	ERR212471	1	15-Jul-2009

Hu_GII_C2H-47_2011_VNM	KM198547	ERR217275	5	03-Nov-2011
Hu_GII_20107_2009_VNM	KM198548	ERR212477	1	20-Aug-2009
Hu_GII_30473_2010_VNM	KM198549	ERR223550	7	01-Mar-2010
Hu_GII_10078_2009_VNM	KM198550	ERR225645	1	10-Jun-2009
Hu_GII_20108_2009_VNM	KM198551	ERR212478	1	20-Aug-2009
Hu_GII_20154_2009_VNM	KM198552	ERR212469	1	16-Sep-2009
Hu_GII_30381_2010_VNM	KM198553	ERR223544	5	04-Jan-2010
Hu_GII_C2H-48_2011_VNM	KM198554	ERR217271	5	04-Nov-2011
Hu_GII_C2035_2010_VNM	KM198555	ERR138006	4	28-Jun-2010
Hu_GII_10127_2009_VNM	KM198556	ERR212492	1	17-Jul-2009
Hu_GII_10194_2009_VNM	KM198557	ERR217286	1	13-Oct-2009
Hu_GII_30257_2009_VNM	KM198558	ERR223541	1	30-Oct-2009
Hu_GII_10129_2009_VNM	KM198559	ERR212493	1	20-Jul-2009
Hu_GII_C2H-50_2011_VNM	KM198560	ERR217268	3	22-Nov-2011
Hu_GII_30303_2009_VNM	KM198561	ERR223542	5	23-Nov-2009
Hu_GII_C2H-55_2011_VNM	KM198562	ERR212509	3	25-Nov-2011
Hu_GII_C2365_2010_VNM	KM198563	ERR212485	5	15-Sep-2010
Hu_GII_C2H-44_2011_VNM	KM198564	ERR217276	3	31-Oct-2011
Hu_GII_10169_2009_VNM	KM198565	ERR212502	1	04-Sep-2009
Hu_GII_20093_2009_VNM	KM198566	ERR212475	1	03-Aug-2009
Hu_GII_10255_2009_VNM	KM198567	ERR217273	1	15-Dec-2009
Hu_GII_C2H-62_2011_VNM	KM198568	ERR217267	3	14-Dec-2011
Hu_GII_10235_2009_VNM	KM198569	ERR217274	1	19-Nov-2009
Hu_GII_20146_2009_VNM	KM198570	ERR212465	1	11-Sep-2009
Hu_GII_20123_2009_VNM	KM198571	ERR212479	1	01-Sep-2009
Hu_GII_20370_2010_VNM	KM198572	ERR212461	5	06-Jan-2010
Hu_GII_C2H-45_2011_VNM	KM198573	ERR217277	5	02-Nov-2011
Hu_GII_10183_2009_VNM	KM198574	ERR217289	1	22-Sep-2009
Hu_GII_20069_2009_VNM	KM198575	ERR212472	1	16-Jul-2009
Hu_GII_C2H-43_2011_VNM	KM198576	ERR217281	3	14-Oct-2011
Hu_GII_10145_2009_VNM	KM198577	ERR212497	1	07-Aug-2009
Hu_GII_C2H-52_2011_VNM	KM198578	ERR212508	3	24-Nov-2011
Hu_GII_10160_2009_VNM	KM198579	ERR212500	1	26-Aug-2009
Hu_GII_10223_2009_VNM	KM198580	ERR217272	1	06-Nov-2009
Hu_GII_10003_2009_VNM	KM198581	ERR225636	1	05-May-2009

Hu_GII_30192_2010_VNM	KM198582	ERR138003	1	21-Sep-2009
Hu_GII_20493_2010_VNM	KM198583	ERR223533	5	19-Mar-2010
Hu_GII_10131_2009_VNM	KM198584	ERR212494	1	22-Jul-2009
Hu_GII_10238_2009_VNM	KM198585	ERR217282	1	19-Nov-2009
Hu_GII_30400_2010_VNM	KM198586	ERR223545	5	11-Jan-2010
Hu_GII_20153_2009_VNM	KM198587	ERR212468	1	16-Sep-2009
Hu_GII_10037_2009_VNM	KM198588	ERR225639	1	15-May-2009
Hu_GII_20144_2009_VNM	KM198589	ERR212463	1	10-Sep-2009
Hu_GII_C2H-39_2011_VNM	KM198590	ERR217270	5	24-Oct-2011
Hu_GII_20122_2009_VNM	KM198591	ERR212480	1	31-Aug-2009

503 Footnotes.

504 1. GenBank accession number, accessible at <http://www.ncbi.nlm.nih.gov/nuccore/>

505 2. European Nucleotide Archive accession number, accessible at <http://www.ebi.ac.uk/ena/>

506 3. Genetic cluster as defined in Figure 2.

507 4. Date of sample collection.

508

509 **Table 3. PCR and genome sequencing success by norovirus genotypes.**

Genotype ¹	Samples	Amplicon1 ²	Amplicon2	Amplicon3	Genomes ³	(%) Successful
GII.4	60	55	55	57	55	92
GII (non GII.4)	58	48	45	53	45	74
GI	10	7	4	4	2	20
GII.2	5	4	1	5	2	40
GII.3	34	26	27	27	26	77
GII.6	8	8	7	8	7	88
GII.7	2	2	0	2	0	0
GII.9	1	1	0	1	1	100
GII.12	2	1	2	2	1	50
GII.13	6	5	6	6	5	83

510 Footnotes.

- 511 1. Sample genotype previously determined (30, 31).
 512 2. Successful reverse transcription/ PCR amplification of sufficient DNA (ca. 0.4 µg) for
 513 Illumina library preparation.
 514 3. Yield of >95% full genome.

515

516 **Table 4. Summary of phylogenetic clusters identified in this study.**

Phylogenetic cluster ¹	Closest genome ²	Genotype by RIVM algorithm ³	Number of genomes	Frequency in 201 genomes
1	NV_GII_VNM_2009_KC175360 NV_GII_VNM_2009_KC175395	GII.P4.DH06b_GII.4.DH06b	140	69.65 (67.6) ⁴
2	NV_GII_VNM_2009_KC175402	GII.P4.DH06b_GII.4.DH06b	2	
3	NV_GII4_TW_2007_JN400615 NV_GII4_Ehime_2007_AB541241	GII.P4.DH06b_GII.4.DH06b	12	
4	NV_GII_VNM_2010_KC175383	GII.P4.NO09_GII.4.NO09	20	9.95 (9.5%)
5	NV_Pune_2007_EU921389	GII.P21_GII.3	19	9.45 (10.2%)
6	NV_Pune_2007_EU921389 NV_GII2_12_Wahroonga_2009_JQ613568	GII.Pg_GII.12	2	0.99 (0.6%)
7	NV_GII_Gifu_1999_AB084071 (<50%)	GII.P7_GII.6	3	1.49 (2.5%)
8	NV_GII_Gifu_1999_AB084071 (<50%)	GII.P7_GII.6	3	1.49 (2.5%)

517 Footnotes

- 518 1. Phylogenetic classification (see Figure 2A).
 519 2. Based on number of reads mapped.
 520 3. Based on the RIVM algorithm(47).
 521 4. Values in parentheses are the genotype frequency values determined by My *et al.* (31).

522

523 **Table 5. Evolutionary rates**

Sequence sets	Genomic region	Mean rate ¹	95% HPD ¹	Substitution Model
Cluster 1: GII.P4 Den Haag 2006b_GII.4 Den Haag 2006b	Whole genome	6.15	5.39 - 6.86	SRD06
	ORF1	5.94	5.04 - 6.94	SRD06
	ORF2	5.69	4.54 - 6.90	SRD06
	ORF3	8.99	6.59 - 11.6	SRD06
	p48 (NS1-2)	6.60	4.83 - 8.47	GTR+G
	NTPase (NS3)	5.41	4.04 - 6.89	GTR+G
	p22 (NS4)	8.21	5.48 - 11.11	GTR+G
	VPg (NS5)	5.95	3.27 - 8.94	HKY+G
	3CLpro (NS6)	6.03	3.57 - 8.61	GTR+G
	RdRp (NS7)	4.74	3.50 - 6.10	GTR+G
Cluster 4: GII.P4 New Orleans 2009_GII.4 New Orleans 2009	Whole genome	5.73	3.74 - 7.81	GTR+G
	ORF1	4.03	1.77 - 6.33	HKY+G
	ORF2	5.60	0.68 - 9.82	HKY+G
	ORF3 ¹			
Cluster 5: GII.P21_GII.3	Whole genome	5.34	4.06 - 6.82	SRD06
	ORF1	4.81	3.45 - 6.17	SRD06
	ORF2	5.99	3.75 - 8.39	SRD06
	ORF3	7.38	2.06 - 13.9	SRD06

524 Footnotes.

525 1. Evolutionary rates were measured as $\times 10^{-3}$ substitutions per site per year. The mean
 526 evolutionary rate ($\times 10^{-3}$ substitutions per site per year) and the 95% upper and lower highest
 527 posterior density (HPD) were determined as described in the Materials and Methods section.

528 2. There was insufficient signal for the algorithms to return a reliable evolutionary rate for
 529 ORF3 region of sequences from GII.4 Cluster 4.

530

531

532 **Table 6. Positive selection analysis**

Codon position ¹	FUBAR ²	MEME ³
ORF1 (p48 (NS1-2)) 79	0.991	0.037
ORF2 (VP1) 106	0.993	0.014
ORF2 (VP1) 298	0.984	>0.05
ORF3 (VP2) 144 ⁴	0.983	0.043

533 Footnotes.

534 1. Codons under positive selection in 140 Cluster 1 norovirus genomes sequenced in this
 535 study as detected using FUBAR(44) and MEME(43).

536 2.The values of the posterior probability (FUBAR) or the *p*-value (MEME) are indicated.

537 3. Analysis for ORF3 covered the first 247 of the protein's 268 codons.

538

539

540

541

542

543 **FIGURE LEGENDS**

544 **Figure 1.** Primer design and function for full genome deep sequencing amplification. **Left**
 545 **panel:** Virtual PCR showing the mapping of the designed primers to a norovirus GII.4 genome
 546 (GenBank JQ613552). Colored circles indicate the position of the each primer and number of
 547 mismatches, grey bars indicate the predicted size of the PCR products. A schematic of the
 548 open reading frame organization of the virus is shown at the top of the panel. **Right panel:**
 549 the PCR products from 14 samples for individual primer pairs for Amplicon 1, 2 and 3. Sample
 550 7 failed to amplify, **c:** water control, **m:** size marker. The size of relevant marker bands (in
 551 kilobasepairs) is indicated to the right of the image.

552
 553 **Figure 2.** Maximum likelihood phylogenetic tree of all HCMC GII genomes from this study
 554 (112) and 89 GII.4 genomes from the same HCMC cohort that were sequenced separately
 555 and made publically available in GenBank, plus selected global reference genomes. The eight
 556 phylogenetic clusters of norovirus identified in this study are marked with colored bars.
 557 Bootstrap support ≥ 0.85 at key nodes is indicated with asterisks. The tree is midpoint rooted
 558 for purposes of clarity, and all horizontal branch lengths are drawn to a scale of nucleotide
 559 substitutions per site.

560
 561 **Figure 3.** The temporal appearance of the HCMC norovirus GII genotypes during the study
 562 period. Genomes were stratified by genotype (from Figure 2), color-coded and plotted by date
 563 of sample isolation.

564
 565 **Figure 4.** Summary of evolutionary rates inferred for the genomic regions GII.4 Cluster 1.
 566 Evolutionary rates were estimated as described in the Methods section, mean values are
 567 indicated by colored circles with the 95% credible interval (CI) indicated. The region of the
 568 norovirus genome used for calculation is labelled, and the two regions with rates elevated
 569 relative to full genome are marked in red.

570
 571 **Figure 5.** Changes in the GII.4 Cluster 1 VP1 protein. The protein sequences were aligned,
 572 and amino acid differences from the reconstructed ancestral sequence of Cluster 1 were
 573 determined and marked with vertical colored bars, with the new amino acid residue color-
 574 coded as shown in the bottom panel; gray bars indicate a gap in the query sequence. The
 575 sequences were ordered by sample date with earliest samples at the bottom of the graph.
 576 Functional domains of the VP1 protein are indicated at the top of the graph and include the
 577 shell domain the protruding 1 (P1) and protruding 2 (P2) domains. The locations of the
 578 blockade epitopes A-E are also indicated (B_A - B_E). A histogram in the second panel
 579 indicates total changes at each position. The protein changes occurring in more than four

580 samples are annotated with the parental amino acid, the position and the new amino acid.
 581 Codons found to be under positive selection by MEME or FUBAR are indicated with a red
 582 asterisk.

583

584 **Figure 6.** Changes in the Cluster 1 minor capsid protein VP2 over time. Protein changes were
 585 analyzed and depicted as described for Figure 5. Functional domains of the VP2 protein
 586 including the VP1 interacting region (VP1-ID) are marked at the top of the graph. A histogram
 587 in the second panel indicates total changes at each position. Codons found to be under
 588 positive selection by MEME or FUBAR are indicated with a red asterisk.

589

590 **Figure 7.** Changes in the Cluster 1 p22 proteins over time. Protein changes were analyzed
 591 and depicted as described for Figure 5. Functional domains of the p22 protein including the
 592 MERES domain are marked at the top of the graph. A histogram in the second panel indicates
 593 total changes at each position.

594

595 **Figure 8.** Changes in the Cluster 1 p48 proteins over time. Protein changes were analyzed
 596 and depicted as described for Figure 5. Functional domains of the p48 protein including the
 597 transmembrane domain are marked at the top of the graph. A histogram in the second panel
 598 indicates total changes at each position. Codons found to be under positive selection by
 599 MEME or FUBAR are indicated with a red asterisk.

600

601 REFERENCES

602

- 603 1. **Thorne LG, Goodfellow IG.** 2014. Norovirus gene expression and replication. *J Gen*
 604 *Viro* **95**:278-291.
- 605 2. **McFadden N, Bailey D, Carrara G, Benson A, Chaudhry Y, Shortland A, Heeney J,**
 606 **Yarovinsky F, Simmonds P, Macdonald A, Goodfellow I.** 2011. Norovirus regulation of the
 607 innate immune response and apoptosis occurs via the product of the alternative open reading
 608 frame 4. *PLoS pathogens* **7**:e1002413.
- 609 3. **Sosnovtsev SV, Belliot G, Chang KO, Prikhodko VG, Thackray LB, Wobus CE,**
 610 **Karst SM, Virgin HW, Green KY.** 2006. Cleavage map and proteolytic processing of the
 611 murine norovirus nonstructural polyprotein in infected cells. *Journal of virology* **80**:7816-7831.
- 612 4. **Wolf S, Williamson W, Hewitt J, Lin S, Rivera-Aban M, Ball A, Scholes P, Savill M,**
 613 **Greening GE.** 2009. Molecular detection of norovirus in sheep and pigs in New Zealand
 614 farms. *Veterinary microbiology* **133**:184-189.

- 615 5. **Wang QH, Han MG, Cheetham S, Souza M, Funk JA, Saif LJ.** 2005. Porcine
616 noroviruses related to human noroviruses. *Emerging infectious diseases* **11**:1874-1881.
- 617 6. **Wang QH, Chang KO, Han MG, Sreevatsan S, Saif LJ.** 2006. Development of a new
618 microwell hybridization assay and an internal control RNA for the detection of porcine
619 noroviruses and sapoviruses by reverse transcription-PCR. *J Virol Methods* **132**:135-145.
- 620 7. **Sugieda M, Nakajima S.** 2002. Viruses detected in the caecum contents of healthy
621 pigs representing a new genetic cluster in genogroup II of the genus "Norwalk-like viruses".
622 *Virus research* **87**:165-172.
- 623 8. **Oliver SL, Dastjerdi AM, Wong S, El-Attar L, Gallimore C, Brown DW, Green J,
624 Bridger JC.** 2003. Molecular characterization of bovine enteric caliciviruses: a distinct third
625 genogroup of noroviruses (Norwalk-like viruses) unlikely to be of risk to humans. *Journal of
626 virology* **77**:2789-2798.
- 627 9. **Martella V, Lorusso E, Decaro N, Elia G, Radogna A, D'Abramo M, Desario C,
628 Cavalli A, Corrente M, Camero M, Germinario CA, Banyai K, Di Martino B, Marsilio F,
629 Carmichael LE, Buonavoglia C.** 2008. Detection and molecular characterization of a canine
630 norovirus. *Emerging infectious diseases* **14**:1306-1308.
- 631 10. **Martella V, Campolo M, Lorusso E, Cavicchio P, Camero M, Bellacicco AL, Decaro
632 N, Elia G, Greco G, Corrente M, Desario C, Arista S, Banyai K, Koopmans M,
633 Buonavoglia C.** 2007. Norovirus in captive lion cub (*Panthera leo*). *Emerging infectious
634 diseases* **13**:1071-1073.
- 635 11. **Liu BL, Lambden PR, Gunther H, Otto P, Elschner M, Clarke IN.** 1999. Molecular
636 characterization of a bovine enteric calicivirus: relationship to the Norwalk-like viruses. *Journal
637 of virology* **73**:819-825.
- 638 12. **Hsu CC, Riley LK, Livingston RS.** 2007. Molecular characterization of three novel
639 murine noroviruses. *Virus Genes* **34**:147-155.
- 640 13. **Teunis PF, Moe CL, Liu P, Miller SE, Lindesmith L, Baric RS, Le Pendu J,
641 Calderon RL.** 2008. Norwalk virus: how infectious is it? *J Med Virol* **80**:1468-1476.
- 642 14. **Atmar RL, Opekun AR, Gilger MA, Estes MK, Crawford SE, Neill FH, Ramani S, Hill
643 H, Ferreira J, Graham DY.** 2013. Determination of the 50% Human Infectious Dose for
644 Norwalk Virus. *The Journal of infectious diseases*.
- 645 15. **Simmons K, Gambhir M, Leon J, Lopman B.** 2013. Duration of immunity to norovirus
646 gastroenteritis. *Emerging infectious diseases* **19**:1260-1267.
- 647 16. **Debbink K, Lindesmith LC, Donaldson EF, Baric RS.** 2012. Norovirus immunity and
648 the great escape. *PLoS pathogens* **8**:e1002921.
- 649 17. **Sukhrie FH, Teunis P, Vennema H, Copra C, Thijs Beersma MF, Bogerman J,
650 Koopmans M.** 2012. Nosocomial transmission of norovirus is mainly caused by symptomatic

- 651 cases. *Clinical infectious diseases : an official publication of the Infectious Diseases Society of*
 652 *America* **54**:931-937.
- 653 18. **Siebenga JJ, Beersma MF, Vennema H, van Biezen P, Hartwig NJ, Koopmans M.**
 654 2008. High prevalence of prolonged norovirus shedding and illness among hospitalized
 655 patients: a model for in vivo molecular evolution. *The Journal of infectious diseases* **198**:994-
 656 1001.
- 657 19. **Verhoef L, Williams KP, Kroneman A, Sobral B, van Pelt W, Koopmans M.** 2012.
 658 Selection of a phylogenetically informative region of the norovirus genome for outbreak
 659 linkage. *Virus Genes* **44**:8-18.
- 660 20. **Sukhrie FH, Teunis P, Vennema H, Bogerman J, van Marm S, Beersma MF,**
 661 **Koopmans M.** 2013. P2 domain profiles and shedding dynamics in prospectively monitored
 662 norovirus outbreaks. *Journal of clinical virology : the official publication of the Pan American*
 663 *Society for Clinical Virology* **56**:286-292.
- 664 21. **Bull RA, Eden JS, Luciani F, McElroy K, Rawlinson WD, White PA.** 2012.
 665 Contribution of intra- and interhost dynamics to norovirus evolution. *Journal of virology*
 666 **86**:3219-3229.
- 667 22. **Kundu S, Lockwood J, Depledge DP, Chaudhry Y, Aston A, Rao K, Hartley JC,**
 668 **Goodfellow I, Breuer J.** 2013. Next-generation whole genome sequencing identifies the
 669 direction of norovirus transmission in linked patients. *Clinical infectious diseases : an official*
 670 *publication of the Infectious Diseases Society of America* **57**:407-414.
- 671 23. **Lindesmith L, Moe C, Lependu J, Frelinger JA, Treanor J, Baric RS.** 2005. Cellular
 672 and humoral immunity following Snow Mountain virus challenge. *Journal of virology* **79**:2900-
 673 2909.
- 674 24. **Cannon JL, Lindesmith LC, Donaldson EF, Saxe L, Baric RS, Vinje J.** 2009. Herd
 675 immunity to GII.4 noroviruses is supported by outbreak patient sera. *Journal of virology*
 676 **83**:5363-5374.
- 677 25. **Donaldson EF, Lindesmith LC, Lobue AD, Baric RS.** 2008. Norovirus pathogenesis:
 678 mechanisms of persistence and immune evasion in human populations. *Immunol Rev*
 679 **225**:190-211.
- 680 26. **Lindesmith LC, Beltramello M, Donaldson EF, Corti D, Swanstrom J, Debbink K,**
 681 **Lanzavecchia A, Baric RS.** 2012. Immunogenetic mechanisms driving norovirus GII.4
 682 antigenic variation. *PLoS pathogens* **8**:e1002705.
- 683 27. **LoBue AD, Lindesmith L, Yount B, Harrington PR, Thompson JM, Johnston RE,**
 684 **Moe CL, Baric RS.** 2006. Multivalent norovirus vaccines induce strong mucosal and systemic
 685 blocking antibodies against multiple strains. *Vaccine* **24**:5220-5234.

- 686 28. **Lindesmith LC, Donaldson E, Leon J, Moe CL, Frelinger JA, Johnston RE, Weber**
687 **DJ, Baric RS.** 2010. Heterotypic humoral and cellular immune responses following Norwalk
688 virus infection. *Journal of virology* **84**:1800-1815.
- 689 29. **Reeck A, Kavanagh O, Estes MK, Opekun AR, Gilger MA, Graham DY, Atmar RL.**
690 2010. Serological correlate of protection against norovirus-induced gastroenteritis. *The Journal*
691 *of infectious diseases* **202**:1212-1218.
- 692 30. **My PV, Thompson C, Phuc HL, Tuyet PT, Vinh H, Hoang NV, Minh PV, Vinh NT,**
693 **Thuy CT, Nga TT, Hau NT, Campbell J, Chinh NT, Thuong TC, Tuan HM, Farrar J, Baker**
694 **S.** 2013. Endemic norovirus infections in children, Ho Chi Minh City, Vietnam, 2009-2010.
695 *Emerging infectious diseases* **19**:977-980.
- 696 31. **Tra My PV, Lam HM, Thompson CN, Phuc HL, Tuyet PT, Vinh H, Hoang NV, Minh**
697 **P, Vinh NT, Thuy CT, Nga TT, Hau NT, Chinh NT, Thuong TC, Tuan HM, Campbell JI,**
698 **Clements AC, Farrar J, Boni MF, Baker S.** 2013. The dynamics of GII.4 Norovirus in Ho Chi
699 Minh City, Vietnam. *Infect Genet Evol* **18**:335-343.
- 700 32. **Cotten M, Lam TT, Watson SJ, Palser AL, Petrova V, Grant P, Pybus OG, Rambaut**
701 **A, Guan Y, Pillay D, Kellam P, Nastouli E.** 2013. Full-genome deep sequencing and
702 phylogenetic analysis of novel human betacoronavirus. *Emerging infectious diseases* **19**:736-
703 742B.
- 704 33. **Quail MA, Otto TD, Gu Y, Harris SR, Skelly TF, McQuillan JA, Swerdlow HP, Oyola**
705 **SO.** 2012. Optimal enzymes for amplifying sequencing libraries. *Nature methods* **9**:10-11.
- 706 34. **Quail MA, Kozarewa I, Smith F, Scally A, Stephens PJ, Durbin R, Swerdlow H,**
707 **Turner DJ.** 2008. A large genome center's improvements to the Illumina sequencing system.
708 *Nature methods* **5**:1005-1010.
- 709 35. **Watson SJ, Welkers MR, Depledge DP, Coulter E, Breuer JM, de Jong MD, Kellam**
710 **P.** 2013. Viral population analysis and minority-variant detection using short read next-
711 generation sequencing. *Philos Trans R Soc Lond B Biol Sci* **368**:20120205.
- 712 36. **Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM,**
713 **Nikolenko SI, Pham S, Prjibelski AD, Pyshkin AV, Sirotkin AV, Vyahhi N, Tesler G,**
714 **Alekseyev MA, Pevzner PA.** 2012. SPAdes: a new genome assembly algorithm and its
715 applications to single-cell sequencing. *J Comput Biol* **19**:455-477.
- 716 37. **Martin D, Rybicki E.** 2000. RDP: detection of recombination amongst aligned
717 sequences. *Bioinformatics* **16**:562-563.
- 718 38. **Dubot-Peres A, Vongphrachanh P, Denny J, Phetsouvanh R, Linthavong S,**
719 **Sengkeopraseuth B, Khasing A, Xaythideth V, Moore CE, Vongsouvath M, Castonguay-**
720 **Vanier J, Sibounheuang B, Taojaikong T, Chanthongthip A, de Lamballerie X, Newton**

- 721 **PN**. 2013. An epidemic of dengue-1 in a remote village in rural Laos. *PLoS Negl Trop Dis*
722 **7:e2360**.
- 723 39. **Stamatakis A**. 2006. RAxML-VI-HPC: maximum likelihood-based phylogenetic
724 analyses with thousands of taxa and mixed models. *Bioinformatics* **22**:2688-2690.
- 725 40. **Darriba D, Taboada GL, Doallo R, Posada D**. 2012. jModelTest 2: more models, new
726 heuristics and parallel computing. *Nature methods* **9**:772.
- 727 41. **Drummond AJ, Suchard MA, Xie D, Rambaut A**. 2012. Bayesian phylogenetics with
728 BEAUti and the BEAST 1.7. *Mol Biol Evol* **29**:1969-1973.
- 729 42. **Minin VN, Bloomquist EW, Suchard MA**. 2008. Smooth skyride through a rough
730 skyline: Bayesian coalescent-based inference of population dynamics. *Mol Biol Evol* **25**:1459-
731 1471.
- 732 43. **Murrell B, Wertheim JO, Moola S, Weighill T, Scheffler K, Kosakovsky Pond SL**.
733 2012. Detecting individual sites subject to episodic diversifying selection. *PLoS Genet*
734 **8:e1002764**.
- 735 44. **Murrell B, Moola S, Mabona A, Weighill T, Sheward D, Kosakovsky Pond SL,**
736 **Scheffler K**. 2013. FUBAR: a fast, unconstrained bayesian approximation for inferring
737 selection. *Mol Biol Evol* **30**:1196-1205.
- 738 45. **Delport W, Poon AF, Frost SD, Kosakovsky Pond SL**. 2010. Datamonkey 2010: a
739 suite of phylogenetic analysis tools for evolutionary biology. *Bioinformatics* **26**:2455-2457.
- 740 46. **Pond SLK, Frost SD, Muse SV**. 2005. HyPhy: hypothesis testing using phylogenies.
741 *BIOINFORMATICS* **21**:676-679.
- 742 47. **Kroneman A, Vennema H, Deforche K, v d Avoort H, Penaranda S, Oberste MS,**
743 **Vinje J, Koopmans M**. 2011. An automated genotyping tool for enteroviruses and
744 noroviruses. *Journal of clinical virology : the official publication of the Pan American Society for*
745 *Clinical Virology* **51**:121-125.
- 746 48. **Chan MC, Lee N, Ho WS, Law CO, Lau TC, Tsui SK, Sung JJ**. 2012. Covariation of
747 major and minor viral capsid proteins in norovirus genogroup II genotype 4 strains. *Journal of*
748 *virology* **86**:1227-1232.
- 749 49. **Siebenga JJ, Lemey P, Kosakovsky Pond SL, Rambaut A, Vennema H, Koopmans**
750 **M**. 2010. Phylodynamic reconstruction reveals norovirus GII.4 epidemic expansions and their
751 molecular determinants. *PLoS pathogens* **6**:e1000884.
- 752 50. **Bok K, Abente EJ, Realpe-Quintero M, Mitra T, Sosnovtsev SV, Kapikian AZ,**
753 **Green KY**. 2009. Evolutionary dynamics of GII.4 noroviruses over a 34-year period. *Journal of*
754 *virology* **83**:11890-11901.

- 755 51. **Boon D, Mahar JE, Abente EJ, Kirkwood CD, Purcell RH, Kapikian AZ, Green KY,**
 756 **Bok K.** 2011. Comparative evolution of GII.3 and GII.4 norovirus over a 31-year period.
 757 *Journal of virology* **85**:8656-8666.
- 758 52. **Vongpunsawad S, Venkataram Prasad BV, Estes MK.** 2013. Norwalk Virus Minor
 759 Capsid Protein VP2 Associates within the VP1 Shell Domain. *Journal of virology* **87**:4818-
 760 4825.
- 761 53. **Zhu S, Regev D, Watanabe M, Hickman D, Moussatche N, Jesus DM, Kahan SM,**
 762 **Naphtine S, Brierley I, Hunter RN, 3rd, Devabhaktuni D, Jones MK, Karst SM.** 2013.
 763 Identification of immune and viral correlates of norovirus protective immunity through
 764 comparative study of intra-cluster norovirus strains. *PLoS pathogens* **9**:e1003592.
- 765 54. **de Rougemont A, Ruvoen-Clouet N, Simon B, Estienney M, Elie-Caille C, Aho S,**
 766 **Pothier P, Le Pendu J, Boireau W, Belliot G.** 2011. Qualitative and quantitative analysis of
 767 the binding of GII.4 norovirus variants onto human blood group antigens. *Journal of virology*
 768 **85**:4057-4070.
- 769 55. **Debbink K, Donaldson EF, Lindesmith LC, Baric RS.** 2012. Genetic mapping of a
 770 highly variable norovirus GII.4 blockade epitope: potential role in escape from human herd
 771 immunity. *Journal of virology* **86**:1214-1226.
- 772 56. **Sharp TM, Crawford SE, Ajami NJ, Neill FH, Atmar RL, Katayama K, Utama B,**
 773 **Estes MK.** 2012. Secretory pathway antagonism by calicivirus homologues of Norwalk virus
 774 nonstructural protein p22 is restricted to noroviruses. *Virology* **9**:181.
- 775 57. **Sharp TM, Guix S, Katayama K, Crawford SE, Estes MK.** 2010. Inhibition of cellular
 776 protein secretion by norwalk virus nonstructural protein p22 requires a mimic of an
 777 endoplasmic reticulum export signal. *PloS one* **5**:e13130.
- 778 58. **Ettayebi K, Hardy ME.** 2003. Norwalk virus nonstructural protein p48 forms a complex
 779 with the SNARE regulator VAP-A and prevents cell surface expression of vesicular stomatitis
 780 virus G protein. *Journal of virology* **77**:11790-11797.
- 781 59. **Breslauer KJ, Frank R, Blocker H, Marky LA.** 1986. Predicting DNA duplex stability
 782 from the base sequence. *Proc Natl Acad Sci U S A* **83**:3746-3750.
- 783
 784















