# varVAMP: automated pan-specific primer design for tiled full genome sequencing and qPCR of highly diverse viral pathogens.

Jonas Fuchs [1*], Johanna Kleine [1], Mathias Schemmerer [2], Julian Kreibich [3], Wolfgang Maier [4], Namuun Battur [5], Thomas Krannich [5], Somayyeh Sedaghatjoo [5], Lena Jaki [1], Anastasija Maks [1], Christina Boehm [2], Carina Wilhelm [2], Jessica Schulze [6], Christin Mache [6], Elischa Berger [4], Jessica Panajotov [7], Lisa Eidenschink [8], Björn Grüning [4], Markus Bauswein [8], Sindy Böttcher [3], Reimar Johne [7], Jürgen Wenzel [2], Martin Hölzer [5], Marcus Panning [1]

 * corresponding author

[1] Institute of Virology, Freiburg University Medical Center, Faculty of Medicine, University of Freiburg, Freiburg, Germany

[2] Institute of Clinical Microbiology and Hygiene, National Consultant Laboratory for HAV and HEV, University Medical Center Regensburg, Regensburg, Germany

[3] National Reference Center for Poliomyelitis and Enteroviruses, Robert Koch Institute, Berlin, Germany

[4] Bioinformatics Group, Department of Computer Science, Albert-Ludwigs-University Freiburg, Freiburg, Germany

[5] Genome Competence Center (MF1), Robert Koch Institute, Berlin, Germany

[6] Unit 17 "Influenza and Other Respiratory Viruses", Robert Koch-Institute, Berlin, Germany

[7] Department of Biological Safety, German Federal Institute for Risk Assessment (BfR), Berlin, Germany

[8] Institute of Clinical Microbiology and Hygiene, University Hospital Regensburg, Regensburg, Germany

## ABSTRACT

Time- and cost-saving surveillance of viral pathogens is achieved by tiled sequencing in which a viral genome is amplified in overlapping PCR amplicons and qPCR. However, designing pan-specific primers for viral pathogens that have high genomic variability represents a major challenge. Here, we present a bioinformatics command-line tool, called varVAMP (variable virus amplicons). It relies on multiple sequence alignments of highly variable virus sequences and enables automatic pan-specific primer design for qPCR or tiled amplicon whole genome sequencing.

The varVAMP software guarantees pan-specificity by two means: it designs primers in regions with minimal variability and introduces degenerate nucleotides into primer sequences to compensate for common sequence variations. We demonstrate varVAMP's utility by designing and evaluating novel pan-specific primer schemes suitable for sequencing the genomes of SARS-CoV-2, Hepatitis E virus, rat Hepatitis E virus, Hepatitis A virus, Borna-disease-virus-1, and Poliovirus. Moreover, we established highly sensitive and specific Poliovirus qPCR assays that could potentially simplify current Poliovirus surveillance. Importantly, wet-lab and bioinformatic techniques established for SARS-CoV-2 tiled amplicon sequencing were readily transferable to these new primer schemes and will allow sequencing laboratories to extend their established methodology to other human pathogens.

## INTRODUCTION

In recent years, next-generation full-genome sequencing of viruses has become an irreplaceable method to track the evolution of viral pathogens, study outbreaks in the human population and animal kingdom, and identify novel zoonotic threats[1–3]. While metagenomic analyses enable the broad analysis of viromes and potentially identify novel pathogens[4], a high genome coverage is required to sufficiently analyze the genomic makeup of a specific viral population in order to e.g. reconstruct viral intra-host evolution[5,6]. This can be achieved by prior virus cultivation or increased sequencing depth, which have drawbacks. Virus cultivation is not always successful and can lead to cell-culture adaptations[7]. Moreover, increased sequencing depth is costly and might still result in poor genome coverage[8]. Targeted sequencing approaches via PCR-tiling or DNA hybridization allow highly specific sequencing on smaller machines without prior pathogen cultivation[9,10]. Particularly, PCR-tiling, in which the viral genome is amplified in overlapping fragments, has gained popularity due to its cost-effectiveness, low input requirement, and simple library preparation. The most prominent viral amplicon schemes were developed for SARS-CoV-2 in early 2020 and have allowed the sequencing of millions of viral genomes during the pandemic[11,12]. However, such amplicon schemes often need to be updated to reflect evolutionary changes or they have not been developed at all for many viral pathogens. Therefore, quantitative real-time PCR (qPCR) remains the diagnostic gold standard for analyzing patient samples for the presence of a viral pathogen[13].

In an optimal setting, tiled-sequencing and qPCR primer designs for viral pathogens should be pan-specific. This can be challenging for viruses with a high genomic variability and common insertions and deletions (INDELs) sites. Thus, primers have to be designed in conserved regions with minimal genomic variation and should not span INDELs. As potential

67 primer target regions might still display sequence variation, degenerate nucleotides can be
68 introduced into primer sequences to further broaden their binding capacity. Optimal
69 pan-specific primers should target highly conserved regions while keeping degeneracy
70 minimal. This problem, coined maximum coverage degenerate primer design (MC-DGD), is
71 a trade-off between specificity and sensitivity[14]. Primer-specific parameters complicate
72 MC-DGD as not all potential regions are also potential primer binding sites[15]. Notably, qPCR
73 designs have even more constraints due to additional hydrolysis probe-specific parameters
74 and a low Gibbs free energy change ($\Delta G$) of the target region[16].

75 Various commercial and open-source primer design applications are available and often
76 utilize primer3 at their core to calculate various primer parameters[17]. However, many of these
77 tools were developed for a particular primer design problem and each of them only
78 addresses some of the previously mentioned problems[18]. Primalscheme is the gold standard
79 for designing tiled primer schemes for viral full genome sequencing[10]. However,
80 primalscheme only handles genomic variations up to a sequence divergence of 5%,
81 precluding its use for viral pathogens with significantly higher sequence divergence, such as
82 Hepatitis E virus (HEV) or Hepatitis A virus (HAV)[19]. Moreover, primalscheme does not
83 introduce degenerate nucleotides into primer sequences, limiting or even abolishing the
84 binding affinity if mutations are located in the primer binding site[20]. Degenerate primer design
85 has been elegantly addressed by software packages like easyPAC or DegePrime[21,22], but
86 they are not suited for the automatic design of tiled or qPCR schemes. For qPCR primer
87 design, there are only a few open-source projects like QuantPrime[23], but most software is
88 not open access and available through commercial companies. However, none of these
89 applications address pan-specific primer design, and not all calculate $\Delta G$, resulting in
90 time-intensive manual primer and amplicon evaluation.

91 Here, we present the command-line tool varVAMP (variable virus amplicons) that enables
92 fully automated pan-specific degenerate primer design for single amplicons, tiled amplicon
93 schemes and qPCR and was tailored to viral genomics. We show varVAMP's utility by
94 designing and testing pan-specific tiled and qPCR primer sets for SARS-CoV-2, HEV
95 (*Paslahepevirus balayani*), ratHEV (*Rocahepevirus ratti*), HAV (*Hepatovirus A*),
96 Borna-disease-virus 1 (BoDV-1, *Orthobornavirus bornaense*), and Poliovirus (*Enterovirus C,*
97 PV) 1-3, that represent different levels of sequence variability.

## MATERIAL AND METHODS

### Software

100 varVAMP is an easy-to-use cross-platform command line tool that is available via PyPI,
101 DOCKER, BIOCONDA, and the Galaxy platform. It enables primer design for a variety of

102 molecular techniques, including single amplicons, tiled full genome sequencing, and qPCR.
103 In the following, we describe the different algorithms and steps that have been implemented
104 for varVAMP. The primer design pipeline only requires an already computed multiple
105 sequence alignment as input. Importantly, all parameters are highly customizable via direct
106 arguments or a config file. varVAMP performs the following major steps (Fig. 1a): (1)
107 Automatic parameter selection based on the input alignment, (2) alignment masking, (3)
108 consensus sequence generation, (4) potential primer region search, (5) evaluation for primer
109 or qPCR probe suitability of digested kmers, and (6) amplicon scheme creation. The last
110 step differs depending on the three different modes available for varVAMP: Single, tiled, and
111 qPCR. To evaluate potential off-targets, varVAMP can use a BLAST database to predict
112 off-target effects and preferentially selects amplicons without off-targets in the final amplicon
113 scheme.

114 In a first step, varVAMP can estimate some of the user-parameters. For a minimal primer
115 length $l_{min}$, two main parameters influence the primer design: The number of ambiguous
116 nucleotides tolerated within a primer sequence $n_a$ with $n_a \in \mathbb{N}$, $0 \leq n_a \leq l_{min}$ and the identity
117 threshold for a nucleotide $t$ to be considered a consensus nucleotide with $t \in \mathbb{R}$, $0 \leq t \leq 1$.
118 Optimization is only performed for $n_a$ or $t$, the other parameter has to be set manually. If $n_a$
119 and $t$ are both not given, $t$ is optimized and $n_a = 2$. For each optimization iteration, $n_a$ or $t$
120 are incremented by $-1$ or $0.1$, respectively. To perform parameter selection, the highest
121 nucleotide frequency at each alignment position is determined. For each optimization step,
122 the lengths of nucleotide stretches that consist of nucleotides reaching the current threshold
123 $\{l_1, l_2 \dots l_m\}$ are calculated. The coverage $c$ of the given alignment that can be considered for
124 potential primer regions is then estimated by:

125
$$c \approx \sum_{l_i \in L} l_i, \qquad L = \{\, l_i \mid i = 1, \dots, m;\ l_i + n_a \geq l_{min} \,\}$$

126 We define that optimization is reached if less than 50% of the alignment can be considered
127 for potential primers. If varVAMP is used to design qPCR schemes, the number of
128 ambiguous characters for the qPCR probe $n_p$ is set to $n_p = n_a - 1$ to ensure a higher
129 specificity of the probe compared to the flanking primers.

130 In the following preprocessing step, gaps in the alignment are masked. Given $n_s$, the number
131 of sequences in the alignment, common gaps are defined as gaps present in more
132 sequences than $n_s \cdot (1 - t)$. Common gaps are then masked with 'N' (single nucleotide
133 deletions) or 'NN' (larger deletions) in the multiple sequence alignment. This ensures that

134 primers will not span regions that could be potential INDEL sites and that for the large
135 majority of sequences, the amplicon size is not overestimated. Moreover, qPCR amplicons
136 are not considered if they would span large deletions as small deviations in the amplicon
137 length are particularly problematic for smaller qPCR amplicons.

138 In the next step, two consensus sequences are deduced from the gap-masked alignment. At
139 each alignment position, the sorted list of observed nucleotide frequencies is calculated. If a
140 sequence in the alignment contains ambiguous nucleotides, all permutations of these
141 nucleotides are considered and added to the nucleotide frequencies proportionally to the
142 number of permutations. The first consensus sequence is generated simply from the most
143 frequent nucleotide at each site. This majority consensus sequence is the basis for the
144 primer search. For the second consensus sequence, the observed nucleotide frequencies at
145 each site are added up, starting from the highest frequency, until their sum reaches $t$. The
146 IUPAC symbol for the set of nucleotides that contributed to the frequency sum is then taken
147 as the consensus at the site. This symbol will be identical to the corresponding nucleotide in
148 the majority consensus if the frequency of that nucleotide alone reaches or exceeds $t$. This
149 second consensus sequence allows searching for regions that only have a certain amount of
150 ambiguous characters within the minimal primer length.

151 Next, the consensus sequence with ambiguous nucleotide characters is searched for
152 potential primer regions. The algorithm opens a region window at the start of the sequence.
153 The window is closed if $> n_a$ ambiguous nucleotides are found within a sequence of $l_{min}$
154 nucleotides or a gap is reached. We define the resulting window as $w = [w_{start}, w_{end}]$ with
155 ambiguous character positions $x_1, \ldots, x_m$. If the window was closed due to a gap, a new
156 window is opened at the subsequent nucleotide $w_{end} + 1$. If the window was closed due to
157 the number of ambiguous characters, the new window is opened at the position after the first
158 ambiguous character counting towards $n_a$ that led to closing the previous window, $x_1 + 1$.
159 Regions are only considered for the primer search if $w_{end} - w_{start} \geq l_{min}$.

160 In the identified primer search regions, the majority consensus is digested into all possible
161 unique $k$-mers for $l_{min} \leq k \leq l_{max}$. Afterwards, each kmer is evaluated for its suitability as a
162 primer. For this primer3[17] is used and some of the rationals and functions were adapted from
163 primalscheme[10]. First, each kmer is hard-filtered independent of its direction for
164 unacceptable temperature, size, GC content, homopolymer length, di-nucleotide repeats,
165 and homodimer formation. Moreover and similar to primalscheme, a base penalty $p_b$ for the
166 kmers' deviations from the optimal temperature, size and GC content is calculated and also

167 hard-filtered if it exceeds the base penalty threshold. Then primers are evaluated for their

168 suitability as forward or reverse primers by filtering out unacceptable hairpin formation

169 temperatures, the 3-prime presence of ambiguous characters, and the absence of a GC

170 clamp. For primers surviving all filtering steps, a permutation penalty $p_p$ and a 3' mismatch

171 penalty $p_m$ are calculated. $p_p$ is calculated as the number of primer permutations of the

172 primer version that has the ambiguous characters (deduced from the ambiguous consensus

173 sequence) multiplied by the permutation penalty. For a primer with $n_p$ characters, we define

174 position-specific penalties $\{c_1, c_2, \dots c_{n_p}\}$ and calculate the mismatch frequencies at each

175 position $\{f_1, f_2, \dots f_{n_p}\}$. We then calculate:

176 $$p_m = \sum_{i=1}^{n_p} c_i \cdot f_i$$

177 Note that only the last 5 positions receive non-zero multipliers in the standard settings. The

178 final primer penalty $p_{primer}$ is then calculated as:

179 $$p_{primer} = p_b + p_m + p_p$$

180 $p_{primer}$ reflects the primers' deviations from base parameters, its number of permutations

181 and number of mismatches at the 3'-prime end. The closer $p_{primer}$ is to zero, the better it

182 represents an optimal primer.

183 To reduce the number of potential primers, all primers are penalty-sorted from low to high.

184 From this sorted list, primers are retained if they do not overlap with the middle third of an

185 already retained lower scoring primer, improving a final selection of primers with minimal

186 overlap and minimal $p_{primer}$. Next, all potential non-dimer forming combinations of forward

187 and reverse primers within a given amplicon range length are computed. A resulting

188 amplicon $a_i$ is defined by primers $primer_{fw}$, $primer_{rv}$. Given the length of the amplicon $l_i$

189 and the user-defined optimal amplicon length $l_{opt}$, we define the amplicon penalty $p_i$:

190 $$p_i = (p_{primer_{fw}} + p_{primer_{rv}}) \cdot e^{\frac{l_i}{l_{opt}}}$$

191 This ensures that the amplicon selection is length dependent and that it favors amplicons

192 with a length closer to the optimal amplicon length. For a single-amplicon design, amplicons

193 are sorted by their penalties from low to high and only low-scoring non-overlapping

194 amplicons are retained.

195 For the tiled approach, a weighted directed graph $G = (V, E)$ with vertices $V$ and edges $E$ is

196 created. Each vertex $v_i \in V$ represents an amplicon $a_i$ and the set of vertices is given by

197 $V = \{v_1, v_2, \ldots v_m\}$. We define the vertex start $start_{v_i}$ as the position of the first nucleotide in

198 $primer_{fw}$ belonging to $a_i$ and the vertex stop $stop_{v_i}$ as the position of the last nucleotide in

199 $primer_{rv}$ belonging to $a_i$. An edge $e \in E$ is defined as a tupel $v_i, v_j, w_j$ of two distinct nodes.

200 The edge weight $w_j = (o_j, p_j)$ incorporates the information about whether or not amplicon

201 $a_j$ generated an off-target hit with the optional BLAST database and amplicon penalty $p_j$. If

202 an off-target is generated $o_j$ is 1 otherwise it is 0. The set of all edges $E$ is defined as

203 $E = \{(v_i, v_j, w_j) \mid v_i, v_j \in V \text{ and } i \neq j \text{ and } v_i, v_j \text{ overlap}\}$. We say that $v_i, v_j$ overlap if they

204 satisfy the user-defined reciprocal pairwise sequence overlap and $start_{v_j}$ is not located in the

205 first half of $v_i$. Next, varVAMP searches for the shortest path in $G$ from a source vertex $v_s$ with

206 Dijkstra's algorithm[24]. The stop position with the highest genomic index of all $v_i \in G$ is

207 denoted $stop_{max} = max\{stop_{v_i} \mid i = 1, \ldots, m\}$ and the lowest penalized amplicon with the

208 furthest stop position reached by Dijkstra's search is termed $v_{max}$. The amplicon coverage

209 over the consensus sequence is defined as $c_{con} = stop_{v_{max}} - start_{v_s}$. We store the current

210 highest coverage $c^*_{con}$ and the shortest path search is repeated for all $v_i$ until

211 $start_{v_i} + c^*_{con} > stop_{max}$. Therefore, the shortest path resulting in the highest coverage is

212 the path that resulted in $c^*_{con}$. varVAMP defines two amplicon pools containing non-adjacent

213 amplicons for the final scheme to allow primer multiplexing. In the last step, both pools are

214 analyzed for the presence of primer heterodimers. If heterodimers are found, varVAMP

215 considers the previously excluded primers overlapping with the middle third of the

216 heterodimer-forming pair and tries to find primers that do not form heterodimers within the

217 respective primer pools.

218 For the qPCR mode, the consensus sequence containing ambiguous nucleotides is

219 searched for regions that satisfy the qPCR probe specific length and $n_p$ constraints and is

220 again digested into all possible unique kmers within the probe size range. These kmers are

221 tested and evaluated for their suitability as qPCR probes in a manner analogous to primer

222 screening. However, here we apply additional constraints: (i) probes are not allowed to have

223 ambiguous bases at either end, (ii) probes cannot have a guanine at the 5' end as this might

224 result in quenching, and (iii) their direction is defined so that the qPCR probes have more
225 cytosines than guanines. Next, varVAMP searches for potential qPCR amplicons. This is
226 achieved by searching for primer subsets within the amplicon length constraint flanking a
227 qPCR probe. Potential amplicons are excluded if they violate the GC content constraint or
228 contain large deletions. The flanking primers must be within a narrow temperature range, the
229 probe has to have a higher temperature than the primers, they cannot form dimers with each
230 other, and the probe has to be within a certain distance to the primer on the same strand.
231 varVAMP also evaluates the presence of dimers in all probe-primer permutations and
232 excludes primer-probe combinations that overlap at their ends, as this might also lead to
233 unspecific probe hydrolysis. Lastly, amplicons are sorted by their amplicon penalty $p_i =$

234 $p_{primer_{fw}} + p_{primer_{rw}} + p_{probe}$ and tested for their ΔG at the lowest primer temperature using
235 seqfold (https://github.com/Lattice-Automation/seqfold). varVAMP reports amplicons that
236 pass the ΔG cutoff.

237 varVAMP has an optional BLAST feature that allows evaluation if amplicon primers could
238 result in off-target products using a custom BLAST database[25]. Here, we perform a relaxed
239 BLAST search with the BLAST settings published for primerBLAST[26]. Afterwards, the results
240 are filtered for matches with a user-definable minimal overlap of identical nucleotides
241 considering both query coverage and mismatches. We now check each amplicon for
242 potential off-target hits defined as matches for both primers that are sufficiently close
243 together on the same reference sequence, but on opposite strands. Amplicons that result in
244 off-target hits, are preferentially not considered in the final scheme. In the single and qPCR
245 mode, amplicons are first sorted for the absence of off-targets and then by their penalty. In
246 the tiled mode, the shortest path is first evaluated on the amount of off-target hits generated
247 by the path before considering the cumulative amplicon penalty, thereby avoiding, but not
248 excluding amplicons with off-target effects.

249 The final primers (independently from the varVAMP mode) are deduced from the consensus
250 sequence incorporating degenerate nucleotides.

## Primer design for the individual pathogens

252 Data selection for the multiple sequence alignments (MSAs) that were used as the inputs
253 for varVAMP were highly dependent on the individual pathogens.

254 For SARS-CoV-2, we obtained 920,323 full-length genome sequences, sampled between
255 2021-10-11 and 2023-09-26, and their lineage assignments from
256 https://github.com/robert-koch-institut/SARS-CoV-2-Sequenzdaten_aus_Deutschland
257 (accessed 2023-10-13). Covsonar (https://github.com/rki-mf1/covsonar, v1.1.9) was used to

258 calculate mutation profiles for all sequences, followed by a Python script 259 (https://github.com/rki-mf1/sc2-mutation-frequency-calculator, v0.0.2-alpha) to select 260 characteristic mutations per lineage (75% frequency). The script then uses these 261 characteristic mutations to construct a single, representative consensus sequence per 262 lineage. A representative consensus sequence was only calculated if at least ten genomes 263 were available for a particular lineage. This resulted in representative consensus genomes 264 for 865 SARS-CoV-2 lineages which were then used as input for varVAMP.

265 For BoDV-1 we downloaded all available full-length sequences that belong to the 266 *Orthobornavirus bornaense* species (BoDV-1: 54 sequences, BoDV-2: 1 sequence). For 267 HAV we downloaded all available full-length sequences that belong to the *Hepatovirus A* 268 species (326 sequences). Patent and artificial clone sequences were excluded resulting in 269 309 HAV sequences. PV sequences were filtered in a similar manner and we excluded, by 270 manual alignment inspection, highly divergent sequences that were likely the result of 271 recombination events with other Enteroviruses yielding 944 sequences. For qPCR designs, 272 we split this dataset, based on metadata, into the individual serotypes 1-3 resulting in 241, 273 494 and 209 sequences, respectively. For HEV data selection, we downloaded all available 274 full-length sequences of the *Hepeviridae* family (1377 sequences). Patent and artificial clone 275 sequences were excluded resulting in 1349 sequences. The remaining sequences were 276 compared to the HEV reference set by Smith et al. 2020[27], extended with the reference 277 sequences for rat, bird, bat, fish, frog and planthopper HEV (NC_038504.1, NC_023425.1, 278 NC_018382.1, NC_015521.1, NC_040835.1 and NC_040710.1, respectively), using the 279 ggsearch36 algorithm[28]. Classification resulted in 1222 HEV sequences and 71 ratHEV 280 sequences. Next, we used the greedy clustering algorithm of vsearch 2.22.1[29] with global 281 clustering thresholds of 0.82 and 0.71, respectively, to further split the HEV and ratHEV 282 datasets by similarity. Clustering results were manually inspected in phylogenetic trees 283 constructed with IQ-TREE 2 under the GTR+F+R10 substitution model and 1000 bootstrap 284 replicates[30]. For HEV we choose two clusters that reflect the most common European HEV-3 285 subgenotypes (HEV-3 f, e and HEV-3 c, h1, m, i, uc, l). For ratHEV we focused on the 286 largest cluster essentially excluding ratHEV from non-rat species and further excluded 287 sequences that were too short resulting in a total of 41 sequences.

288 Next, the pairwise sequence identity within each sequence batch was calculated with Identity 289 (https://github.com/BioinformaticsToolsmith/Identity)[31] and the sequences were aligned with 290 MAFFT[32] with default settings. These alignments were then used as the input for varVAMP. 291 Based on the sequence identity, we chose, for tiled sequencing, to fix the allowed max 292 number of ambiguous characters within the minimum primer length depending on mean 293 sequence identity within the batch ($n_a$ = 2 with 90% identity, $n_a$ = 4 between 70 and 80%, $n_a$

294 = 5 below 70%). Next, the identity threshold $t$ was maximized until varVAMP could not find a

295 tiled scheme that covered the whole genome (Table 1). For the qPCR designs we chose to

296 allow one less ambiguous base for the probe compared to the primers (Table 2). Here,

297 settings were selected based on if varVAMP was able to find a qPCR scheme under the ΔG

298 constraints rather than solely on sequence similarity. All input alignments and varVAMP

299 outputs are available at: https://github.com/jonas-fuchs/ViralPrimerSchemes.

## HEV qPCR and HEV sub-genotyping

301 Patient serum samples were tested for HEV using the HEV RT-PCR Kit 1.5 (AS0271543,

302 Altona Diagnostics, AltoStar®). For HEV sub-genotyping of HEV-positive samples, we used

303 an in-house nested RT-PCR (210212, Qiagen, Hilden, Germany) protocol based on

304 previously published primers and nested primers in a conserved region of ORF1[33]. RT-PCR

305 was performed at 42 °C for 60 min, 15 min at 95°C followed by 40 cycles at 94 °C (30 s),

306 56.5 °C (30 s) and 74 °C (45 s). Final elongation was at 74 °C for 5 min. Agarose gel

307 negative PCR reactions were subjected to a nested PCR reaction. Afterwards, PCR

308 products were Sanger sequenced.

## Production of virus stocks

310 The BoDV-1 strains were derived from native human brain sections[34–36] and were

311 propagated in Vero cells, which were grown in DMEM supplemented with 10%

312 heat-inactivated fetal calf serum (FCS), 90 U/ml streptomycin, 0.3 mg/ml glutamine, and 200

313 U/ml penicillin (all PAN Biotech, Aidenbach, Germany). Permanently infected cells were split

314 twice a week and monitored for *Mycoplasma* spp. contamination every 12 weeks. To obtain

315 a cell-free viral stock, cell culture supernatants were centrifuged at 1,000 rcf to remove cell

316 debris and filtered through Rotilab syringe filters with a pore size of 0.22 μm (Carl Roth,

317 Karlsruhe, Germany).

318 HEV-containing supernatant was harvested from persistently infected cell culture. The liver

319 carcinoma cell lines (PLC/PRF/5, ATCC: CRL-8024) persistently infected with HEV-3c strain

320 14-16753, HEV-3e strain 14-22707 or HEV-3f strain 15-22016 (provided by National

321 Consultant Laboratory for HAV and HEV, University Hospital Regensburg) were maintained

322 in modified Minimum Essential Medium at 37 °C and 5% $CO_2$[37].

323 HAV genotype IB strains MBB[38] and V18-35519 (derived from plasma of a patient with acute

324 hepatitis A) were both propagated in HuH-7 cells maintained in BMEM and incubated at 34.5

325 °C and 5% $CO_2$. MBB and V18-35519 strains were harvested at 665 and 378 days post

326 inoculation, respectively.

327 RD-A cells were infected with PV for virus propagation and cultivated in MEM Earls media
328 with L-glutamine, 1x Non-essential amino acids, 100 U/ml penicillin and streptomycin
329 (61100087, 11140050, 15140122, Thermo Fisher Scientific, Germany) and 7.3% heat
330 inactivated fetal calf serum (BioWest, South America). RD-A cells were split once a week
331 and an internal quality control performed at passage five to ensure cell sensitivity. The cell
332 culture was conducted at 37 °C and 5% $CO_2$. Cultures were checked daily for a cytopathic
333 effect (CPE) for a maximum of seven days. Up to 2 ml of cell cultures with an observed CPE
334 were centrifuged 4 min at 4000 rpm and the supernatant used as viral stock.

335 RatHEV strain R63[39], which was originally detected in a Norway rat from Germany, and
336 ratHEV strain pt2[40], which was identified in a human patient in Hong Kong, were generated
337 and propagated in HuH-7-Lunet BLR cells under conditions as described previously[41,42] . The
338 ratHEV positive culture supernatants were harvested after 66 days post infection

### 339 Tiled Illumina sequencing for HEV-3

340 Viral RNA was isolated using the QIAamp® Viral RNA kit (52904, Qiagen, Hilden, Germany)
341 following the manufacturer's protocol. Subsequently, a one-step RT-PCR using the
342 SuperScript™ IV One-Step RT-PCR System (12594025, Thermo Fisher) was performed for
343 each amplicon separately with a total primer concentration of 1 µM. To reduce non-specific
344 amplification, reverse transcription was performed at 55 °C for 60 min and we gradually
345 reduced the primer annealing temperature during the PCR in the first 10 cycles (10 sec at 98
346 °C, 10 sec at 63 °C (-0.5°C/cycle), 2 min at 72 °C) and then performed another 35 cycles at
347 a constant annealing temperature (10 sec at 98 °C, 10 sec at 58 °C, 2 min at 72 °C). Next,
348 amplicons were pooled and purified using AMPure XP beads (A63881, Beckman Coulter).
349 50-100 ng of DNA were prepared for Illumina sequencing using the NEBNext Ultra II FS
350 DNA Library Prep Kit (E6177, NEB, Frankfurt am Main, Germany). Normalized and pooled
351 sequencing libraries were denatured with 0.2 N NaOH and sequenced on an Illumina MiSeq
352 instrument using the 300-cycle MiSeq Reagent Kit v2 (MS-102-2002, Illumina).

### 353 Tiled Illumina sequencing for BoDV-1, HAV, ratHEV

354 Viral RNA was isolated using the QIAamp Viral RNA Mini Kit (Qiagen, Hilden, Germany) or
355 the EMAG Nucleic Acid Extraction System (Biomeriéux Deutschland GmbH, Nürtingen,
356 Germany). RNA was transcribed into cDNA with LunaScript RT SuperMix Kit (New England
357 Biolabs, Ipswich, MA, USA) for 2 min at 25 °C, 10 min at 55 °C and 1 min at 95 °C. The
358 cDNA was then amplified in single reactions with primer pairs, as well as multiplex PCRs
359 with primer pools using the Q5 Hot Start-Fidelity DNA Polymerase Kit (New England Biolabs,
360 Ipswich, MA, USA) with an initial step at 98 °C for 30 sec, followed by 35 cycles (15 sec at

361 98 °C, 5 min 65 °C) and a final extension step at 65 °C for 5 min. Illumina sequencing was
362 performed analogous to the HEV-3 sequencing protocol.

## Tiled Illumina sequencing for PV

364 PV vaccine strain 1-3 (Sabin 1-3) RNA was isolated using the QIAamp® Viral RNA kit
365 (52904, Qiagen, Hilden, Germany) following the manufacturers' protocol from infected RD-A
366 cells, with PV2 being archived RNA due to containment reasons. Subsequently, a one-step
367 RT-PCR using the Qiagen One-step RT-PCR kit (210212, Qiagen, Hilden, Germany) was
368 performed for each amplicon separately or for the respective pools with a total primer
369 concentration of 0.6 µM. Reverse transcription step was done for 30 min at 50 °C followed
370 by an initial PCR activation step at 95 °C for 15 min. Product amplification was done in 40
371 cycles with a stepwise reduction of the primer annealing temperature during the first 10
372 cycles (30 sec at 94 °C, 45 sec at 70 °C ($\Delta$T -1 °C/cycle), 90 sec 72 °C) and a constant
373 annealing temperature for the next 30 cycles (30 sec at 94 °C, 45 sec at 60 °C , 90 sec
374 72°C) and a final extension step for 10 min at 72 °C. Multiplex RT-PCR pools of each sample
375 were combined and purified using MagSi-NGSPREP-PLUS beads (MDKT00010075,
376 Steinbrenner, Germany) according to the manufacturer's manual and DNA concentration
377 measured using the Qubit™ 1X dsDNA Assay-Kit (Q33230, Thermo Fisher Scientific,
378 Germany).

379 For Illumina sequencing, library preparation was done using 70-400 ng DNA with the
380 Nextera XT DNA Library Preparation Kit (FC-131-1096, Illumina) and sequenced on an
381 Illumina MiSeq Instrument (2 x 300 bp read length).

## Tiled ONT sequencing for HAV

383 Nucleic acid was extracted from samples on an EZ1® Advanced XL workstation using the
384 EZ1 Virus Mini Kit v2.0 (Qiagen, Hilden, Germany) and transcribed into cDNA with
385 LunaScript RT SuperMix Kit (New England Biolabs, Ipswich, MA, USA) for 2 min at 25 °C,
386 10 min at 55 °C and 1 min at 95 °C. The cDNA was then amplified in multiplex PCRs with
387 HAV-specific primer pools using the Q5 Hot Start-Fidelity DNA Polymerase Kit (New England
388 Biolabs, Ipswich, MA, USA) with an initial step at 98 °C for 30 sec, followed by 35 cycles (15
389 sec at 98 °C, 5 min 65 °C) and a final extension step at 65 °C for 5 min. Barcoding was
390 performed with eight samples per run using the Rapid Barcoding Kit 96 V14 (Oxford
391 Nanopore Technologies, Oxford, UK). The library was sequenced on an Mk1C (MinKNOW
392 software version 23.07.12) for 72 hours using an R10.4.1 Flow Cell.

## Tiled ONT sequencing for SARS-CoV-2

A total of 14 SARS-CoV-2 positive extracts, tested at the Robert Koch Institute were selected. These samples originated from IMSSC2-lab networks that were received under the RKI integrated genomic surveillance program. All samples were nasopharyngeal or oropharyngeal swabs originating from patients in Germany during January-February 2024. Total nucleic acid extraction was done using MagNA Pure 96 DNA and viral NA Small Volume kit (Roche Life Science, Mannheim, Germany) on an automated extraction instrument (MagNA Pure 96 system, Roche Diagnostics) according to the manufacturer's manual. Reverse-transcription was performed on viral RNA extracts using LunaScript® RT SuperMix (New England Biolabs, as part of the NEBNext ARTIC SARS-CoV-2 Companion Kit (Oxford Nanopore Technologies) according to the manufacturer's protocol. Amplification was performed with 35 cycles of annealing temperature of 60 °C for 2 min and elongation at 72 °C for 3 min for both pools, without amplicons cleanup afterward. The barcoding was done using ONT Native Barcoding Expansion kit (EXP-NBD196). Fourteen samples together with 2 negative controls were multiplexed on a FLO-MIN 114 flow cell version R10 and sequenced on a GridION Mk1 device for 16 hours.

## RT-qPCR (PV 1-3)

The different varVAMP RT-qPCR assays were performed on serial dilution series of PV 1-3 (Sabin) RNA with a RNA concentration ranging from 1-3 ng/µl of the stock solution to evaluate the performance and sensitivity of the different assays. For proof of lack of cross-detection between the PV vaccine strains, all possible combinations were tested. Quantitative realtime RT-PCR was carried out using the 4X CAPITAL™ 1-Step qRT-PCR Probe Master Mix (BR0502002, Biotechrabbit, Germany) on a Roche instrument, following the manufacturer's instructions. The reaction contained 0.4 µM of each primer with a total reaction volume of 20 µl. qRT-PCR cycling program started with reverse transcription at 50 °C for 10 min followed by an activation step at 95 °C for 3 min and 45 cycles for target amplification (95 °C for 10 sec, 59 °C for 30 sec). Primer annealing temperature was chosen after preliminary tests with different annealing temperature settings ranging from 56 - 64 °C on a Biometra TAdvanced (analytik jena, Germany) and product observation on an 1.5% agarose gel. The fluorescence was measured during the extension step with three different channel setups respective to the probe fluorophore (FAM = 465-510 nm, JOE = 533-580 nm, CY5 = 618-660 nm).

## Sequencing data analysis

The de-multiplexed raw Illumina reads were subjected to a custom Galaxy pipeline which we had initially developed for tiled amplicon sequencing of SARS-CoV-2[6]. These reads were

428 pre-processed with fastp (v0.20.1) [43] and mapped to respective closely related genomes
429 using BWA-MEM[44] (v0.7.17). Importantly, the 3' and 5' regions of the viral reference genome
430 were masked prior mapping until the 5' and 3' end of the flanking primer binding regions,
431 respectively, as no novel information can be generated in these regions. Post mapping,
432 primer sequences were trimmed with ivar trim (v1.3.1). Variants (SNPs and INDELs) were
433 called with the ultrasensitive variant caller LoFreq[45] (v2.1.5), demanding a minimum base
434 quality of 30 and a coverage of at least 20-fold. Afterwards, the called variants were filtered
435 based on a minimum variant frequency of 10% and on strand bias support. Finally,
436 consensus sequences were constructed with bcftools (v1.15.1)[46]. Regions at both genome
437 ends that lie outside the amplicons, regions with low coverage (<20x) or variant frequencies
438 between 0.3 and 0.7 were masked with Ns.

439 For Oxford nanopore sequencing of SARS-CoV-2, poreCov (v1.9.3), a Nextflow workflow
440 specifically tailored for SARS-CoV-2 genome reconstruction from nanopore amplicon data,
441 was used to perform mapping (minimap2; v2.17)[47], primer clipping, variant calling (Medaka;
442 v1.8.0), and consensus genome reconstruction[48]. We ran poreCov to initially filter reads
443 below 400 bp (--minLength 400) and above 1 kbp (--maxLength 1000) while coverage
444 downsampling was disabled (--artic_normalize 0). The r1041_e82_400bps_sup_v4.2.0
445 model was used for variant calling with Medaka and detected variant calls filtered by a
446 minimum base quality of 20 and a coverage of at least 20-fold. We set the allelic frequency
447 of called mutations to 1 to ensure compatibility with the Illumina data for the in silico analysis.

448 For Oxford nanopore sequencing of HAV samples, pod5 raw data was duplex basecalled
449 with Dorado version 0.5.3 and demultiplexed. Subsequent, fastq files were again processed
450 using a custom Galaxy pipeline. First, raw data was pre-processed with fastp[43] excluding
451 reads <50 bp and >2000 bp. Afterwards, reads were mapped to the HAV reference genome
452 NC_001607 using minimap2 (v2.17)[47] and trimmed with ivar trim (v1.3.1). Variants were
453 called with medaka (v1.3.2) and consensus sequences were constructed with bcftools
454 (v.1.15.1)[46]. Regions at both genome ends that lie outside the amplicons and regions with
455 low coverage (<20x) were masked with Ns.

## Data analysis and visualization

457 qPCR amplification curves were analyzed using the Roche LightCycler 480 II device
458 software version LCS480 1.5.1.62 (Roche Applied Science). Data was analyzed using the
459 2nd derivative. Mapped bam files were analyzed and visualized with BAMdash v.0.2.4
460 (https://github.com/jonas-fuchs/BAMdash). We used GraphPad Prism 8 (genome recovery
461 and qPCR), R 4.3.2 (phylogenetic tree) or python 3.11 (remaining figures) for data analysis
462 and visualization. The data and code to reproduce the figures is available at:

463 https://github.com/jonas-fuchs/varVAMP_in_silico_analysis. The schematic varVAMP
464 workflow and data preparation workflow were created with biorender
465 (https://www.biorender.com).

## RESULTS

### Software and output

468 The command-line tool varVAMP was written in python3 and requires only a pre-computed
469 MSA as input. Notably, varVAMP is cross-platform (Windows 11, MacOS and Linux) with
470 python 3.9 or higher being the single requirement prior to installation. varVAMP can design
471 primers for single amplicons, tiled amplicon schemes and qPCR. The pipeline consists of
472 multiple steps that are common to all different modes (alignment preprocessing, consensus
473 generation and primer evaluation), mode-specific or optional (automatic parameter search
474 and BLAST evaluation) (Fig. 1a). At its core, varVAMP wraps Primer3[17] and uses a
475 kmer-based approach to find all potential primers in a consensus sequence calculated from
476 the input MSA. varVAMP addresses the MC-DGD problem by first calculating two consensus
477 sequences that consist either of the majority nucleotides at each position or integrate
478 degenerate nucleotides. The latter is used to find potential primer regions that are regions
479 with a user-defined maximum amount of degenerate nucleotides within the minimal primer
480 length. Afterwards, kmers of the majority consensus sequence that lie within these potential
481 primer regions are tested for all relevant primer parameters. varVAMP evaluates these
482 primers via a penalty system that incorporates information about primer parameters, 3'
483 mismatches, and degeneracy. In its tiled sequencing mode, varVAMP finds overlapping
484 amplicons spanning the alignment while minimizing primer penalties by using Dijkstra's
485 algorithm for finding the shortest paths between nodes in a weighted graph[24]. For qPCRs,
486 varVAMP evaluates probe and primer parameters independently and tests the ΔG of
487 potential qPCR amplicons. The final primers are then deduced from the consensus
488 sequence incorporating degenerate nucleotides. For some of the more computationally
489 intensive tasks, the program is capable of using multicore processing. Although we have not
490 extensively evaluated the running times, varVAMP typically finishes within seconds to
491 minutes. This is highly dependent on the alignment's size and number of sequences, as well
492 as the alignment's sequence variability that directly influences the amount of found primers.
493 varVAMP produces multiple outputs in standardized formats and a plot displaying the
494 alignment's normalized Shannon's entropy, all potential target regions, all primers that
495 passed the initial filtering steps and the final amplicon design with low penalty primers (Fig.
496 1b).

## Design and evaluation of HEV pan-specific primers

HEV of the genus *Paslahepevirus* is the most common cause of acute viral hepatitis worldwide and is phylogenetically separated into four distinct genotypes (genotypes 1-4). In risk groups such as immunocompromised patients, the zoonotic HEV genotype 3 (HEV-3) can cause acute or chronic hepatitis[49]. HEV-3 has a high prevalence in industrialized countries and is further classified into subgenotypes with varying prevalence depending on the geographic region[50]. Most genome sequences show exceptional variability[51] and have to be generated from the initial patient material as virus isolations require optimized cell culture systems[37]. To provide a simple sequencing procedure from patient material and to test varVAMP's real-world applicability, we set out to design primers for HEV-3 tiled sequencing as a proof-of-principle. We initially downloaded all available full-genome HEV sequences from NCBI's genbank and classified the (sub-) genotypes using fasta36 as previously described (Fig. 2a)[2]. Our aim was to design primers that would be specific for multiple HEV-3 sub-genotypes. Therefore, sequences were clustered based on their similarity using vsearch[29] and the clustering result evaluated by constructing a maximum-likelihood phylogenetic tree with IQ-TREE 2[30]. Clustering resulted in seven clusters with more than 6 sequences (Fig. 2b). Four large clusters belonged to HEV-3 each comprising multiple subgenotypes. We decided to design primers for cluster 2 (HEV-3 f, e) and cluster 4 (HEV-3 c, h1, m, i, uc, l) to reflect the most common European HEV-3 subgenotypes[2]. Therefore, sequences of cluster 2 and 4 were separately aligned with MAFFT[32] and the alignments used as the input for varVAMP yielding seven and six 1-1.5 kb amplicons, respectively (Fig. 2a and Table 1). Next, we evaluated these primer schemes on persistently HEV-3 f (strain: 15-22016) and c (strain: 14-16753) infected cell cultures using a one-step RT-PCR protocol. Agarose gel electrophoresis showed consistent and strong amplification for all amplicons with only a few unspecific bands for amplicon 2 and 3 of cluster 4 (Fig. 2c). Next generation Illumina sequencing of the pooled PCR products resulted in an even and high coverage for both samples (Fig. 2d). To further evaluate the primer schemes, we applied our protocol to a third HEV-3 e (strain: 14-22707) persistently infected cell culture for cluster 2 and to HEV-3 positive patient material for both clusters. In order to select the proper amplicon scheme, we first subclassified HEV-3 positive blood samples. Next, we evaluated the cluster 2 and 4 primer schemes on HEV-3 e (n=2) or HEV-3 c (n=4) samples, respectively. Next-generation sequencing results were comparable to the prior results and allowed HEV-3 genome reconstruction (Fig. 2e and S1). However, for patient 2 of cluster 2 and patient 4 of cluster 4 we observed a single amplicon dropout (S1). Interestingly, both dropouts were caused by amplicons with a forward primer close to the HEV-3 hypervariable region[51] suggesting the presence of potential INDELs or variations that might have restricted primer binding.

533 In summary, we used varVAMP to design two tiled primer schemes each specific for multiple 534 HEV-3 sub-genotypes and used Illumina sequencing to recover near-to-complete viral 535 genomes for both infected cell cultures and patient material.

## *In silico* design and evaluation of primer schemes for multiple viral pathogens with diverse sequencing variability

538 We designed primer schemes for tiled full-genome sequencing for SARS-CoV-2, BoDV-1, 539 HAV, PV and ratHEV that display different degrees of sequence conservation over the whole 540 genome with SARS-CoV-2 having the lowest (99 % pairwise identity) and ratHEV the highest 541 overall sequence variability (57 % pairwise identity) (Table 1 and Fig. 3a). Similar to HEV-3, 542 pan-specific amplicon sequencing protocols would massively simplify diagnostics and 543 surveillance. The initial data selection and pre-processing was highly dependent on the 544 individual data sets and inspired by our experiences with HEV-3. Only for SARS-CoV-2, we 545 did not directly align sequences from public databases, but created representative 546 consensus sequences of circulating lineages in Germany between October 2021 and 547 September 2023 (920k samples) to represent the most prevalent variations for each lineage 548 within the alignment. Depending on the mean pairwise identity, we chose to tolerate one to 549 five ambiguous bases within the primer sequences and optimized the identity threshold 550 (Table 1). With the exception of BoDV-1 and SARS-CoV-2, we aimed for an amplicon size of 551 over 1000 bp so amplicons could span regions with an overall higher variability in which 552 potential primers are scarce (Fig. 3a). Next, we evaluated the designed primer schemes *in* 553 *silico* prior to wet-lab evaluation.

554 First, we analyzed the degeneracy per primer as this is highly penalized by varVAMP to keep 555 the number of permutations minimal. Two, four and five tolerated ambiguous bases within a 556 primer sequence can lead to a maximum degeneracy of 4, 256 and 1024, respectively. 557 However, for the primer schemes with four and five tolerated ambiguous bases the mean 558 number of permutations was over 10-fold lower than theoretically possible, indicating a 559 preferential selection of primers with a low degeneracy (Fig. 3b). With the integration of 560 ambiguous bases, varVAMP aims to minimize mismatches with the input MSA. Therefore, 561 we analyzed for each scheme the number of mismatches against sequences in the MSA 562 (Fig. 3c). Indeed, the large majority of build consensus sequences did not have sequence 563 variations not covered by any primer permutations with ratHEV having higher number of 564 mismatches likely due to the sequence variability of the MSA (Fig. 3 a). varVAMP penalizes 565 mismatches in the last five bases of a primer's 3' end to ensure stable target binding. By 566 analyzing the position-dependent mismatches of all primers in a scheme, we indeed 567 observed that most sequences of the input MSA displayed a low frequency of mismatches at

568 the 3' end of the primers (Fig. 3d). Starting at the 3' end, the number of mismatches 569 increased in consistent waves of three nucleotides for all primer schemes except 570 SARS-CoV-2. We hypothesized that this might be due to synonymous codon usage caused 571 by variations in the second and third codon position[52,53]. Manual inspection of the primers' 572 locations in the MSAs indeed confirmed that our 3' penalty system preferentially selected 573 primers if their 3' ends are located at the first and not at the second or third position of a 574 codon. Lastly, we explored the hypothesis that the mean primer parameters of all primer 575 permutations would lie within our target range even if they were initially calculated on the 576 basis of the primer sequence including the most common nucleotides. We therefore 577 calculated melting temperature, hairpin temperature, homo-dimer temperature and GC 578 content (Fig. 3 e-h). In most cases, the mean of the primer permutations were within the 579 target range or below the cutoff but showed a higher deviation from the optimum compared 580 to the primer that was initially used for parameter calculation. The GC content is the least 581 penalized parameter by varVAMP and other parameters should have a more pronounced 582 effect on primer selection with the current settings. Indeed, the GC content was also within 583 the target range but overall more dependent on the MSA's GC content (Fig. 3 h).

584 In conclusion, the newly designed primers should recognize the majority of sequences in the 585 initial MSA while minimizing degeneracy, overall mismatches, and 3' mismatches.

## 586 Full genome tiled amplicon sequencing of SARS-CoV-2, BoDV-1, HAV, PV and 587 ratHEV

588 In a multi-center study with specialists for the respective pathogens, we evaluated if the 589 newly designed primers for SARS-CoV-2, BoDV-1, HAV, PV and ratHEV were suitable for 590 whole-genome sequencing and genome reconstruction. Similar to the HEV-3 primer 591 schemes, we performed amplicon-based Illumina and, in the case of SARS-CoV-2 and some 592 of the HAV samples, ONT sequencing on various samples in either singleplex or multiplex 593 PCR reactions. Sequencing protocols and selection of samples differed due to 594 center-specific preferences. For SARS-CoV-2, we tested the novel primer scheme in 595 multiplex PCR reactions on a random set of respiratory patient specimens from currently 596 circulating variants with different viral loads. Although some amplicons had a lower 597 coverage, we were able to construct complete genomes in the majority of cases (Fig. 4a and 598 S1). We evaluated the BoDV-1 primers in multiplex reactions on three different virus stocks 599 that had been isolated from brains of diseased patients in 2019, 2020 and 2022[34–36] and 600 were cultivated on Vero cells. For all isolates we were able to recover highly covered 601 genome sequences (Fig. 4b and S1). Only for the 2022 isolate, the last three amplicons 602 were poorly amplified, leading to a slightly lower genome recovery (S1). For HAV, we tested

603 the HAV-specific primers on the cell culture derived lab strain V18-35519. Illumina
604 sequencing yielded consistent and high coverage over all amplicons independent of multi- or
605 singleplex reactions (Fig. 4c). Next, we transferred the protocol to three different
606 HAV-positive patient samples: genotype IB-positive feces (patient 1) and sera (patient 3) as
607 well as genotype IA positive feces (patient 2). Finally, we sequenced another four patients
608 sera via Oxford Nanopore: IA-positive (patient 4 and 5), IB-positive (patient 6) and
609 IIIA-positive (patient 7). Full-genome recovery was achieved with all samples. However, for
610 some of the amplicons we observed a lower overall coverage compared to the virus isolate
611 (S1). Next, the PV primer scheme was tested on the Sabin 1-3 vaccine strains. Similar to our
612 prior results, sequencing resulted in high coverage and full-genome recovery. However, we
613 observed that the third amplicon overall under-performed in multiplex but not in singleplex
614 reactions (Fig. 4d and S1). Lastly, we evaluated the ratHEV primers that we designed to test
615 the limits of varVAMP given the highest sequence variability and low number of MSA
616 sequences (Table 1). We tested either single- or multiplex PCR reactions for the two
617 previously described isolates R63 and pt2[39,40] (Fig. 4e). While we were able to achieve high
618 coverage and genome recovery for the R63 isolate, we observed one and two amplicon
619 dropouts for the pt2 single- and multiplex reactions, respectively (S1).

620 We systematically evaluated the coverage and amplicon recovery for all primer schemes
621 and samples (Fig. 5a). Most amplicons performed in a sample-dependent manner but in
622 some cases multiplex performance was intrinsic to specific amplicons as exemplified by the
623 third amplicon of the PV scheme (Fig. 4c). As all multiplex reactions across tested schemes
624 were performed with equimolar primer concentrations, we hypothesized that the
625 performance could be improved by balancing primer concentrations. Therefore, we adjusted
626 the molarity of the PV primers in two consecutive rounds as a proof-of-concept for further
627 wet-lab optimization. For the final iteration, we achieved a coverage for all Sabin strains that
628 was comparable to the respective singleplex reactions (S2). Analogous to the mismatch
629 analysis with the input MSA (Fig. 3c), we also examined how many nucleotide mismatches
630 between the primers and their target regions are present in our sequencing results (Fig. 5b).
631 The primer target regions of the new sequences showed up to two mismatches to the
632 degenerated primer sequences with the majority having no mismatches. Similar to the
633 amplicon performance, the number of mismatches were mostly sample-dependent. As the
634 samples used for evaluation were selected in the respective centers based on availability
635 and not sequence diversity, we tested if this could have produced an unintended selection
636 bias towards specific viral strains. Therefore, we evaluated if the novel consensus
637 sequences of each scheme have a variability that is comparable to that of the respective
638 input alignments. Pairwise sequence identities of these small datasets were highly similar or

639 lower to that of the alignment with only the newly produced sequences for HEV cluster 4
640 showing a significantly higher mean pairwise sequence identity of 7 %, respectively (Fig. 5c),
641 indicating a slight bias for higher conserved sequences.

642 In summary, all primer schemes were suitable for tiled amplicon Illumina or ONT sequencing
643 and resulted in highly covered full-genome sequences by applying center-specific
644 sequencing protocols and bioinformatic pipelines initially developed for tiled amplicon
645 sequencing of SARS-CoV-2.

## Design and wet-lab evaluation of PV qPCR primers designed with varVAMP

647 The WHO gold standard for PV detection is based on time- and resource-consuming virus
648 cultivation. Molecular detection by qPCR is available but was designed for virus isolates
649 propagated in cell culture[54]. Moreover, primers and probes display a high level of
650 degeneracy, decreasing the assay's sensitivity and increasing the risk of unspecific non-viral
651 amplification products for other specimen sources like wastewater samples. Therefore, we
652 used varVAMP to design PV serotype-specific assays. Optimal primer annealing
653 temperature of each assay was tested with a gradient PCR ranging from 56 - 64 °C. All
654 annealing temperatures resulted in the expected product size with 59 °C showing the lowest
655 abundance of unspecific products. PV serotype specific RT-qPCR assays were performed in
656 a serial dilution experiment between $10^{-1}$ and $10^{-8}$ with RNA extracted from viral supernatant.
657 For all three types, PV detection was achieved up to a dilution of $10^{-7}$ (Fig. 6 a-c). Absolute
658 quantification analysis based on the dilution series calculated an efficiency value (PV1 =
659 1.88, PV2 = 1.87, PV 3 = 1.90) close to 2 corresponding to a perfect amplification reaction.
660 Cross-specificity testing showed no detection between the Sabin strains indicating that the
661 designed primers and probes are highly specific for their respective PV serotype, despite
662 primer and probe degeneration.

## DISCUSSION

664 Here, we describe varVAMP, a command-line software tailored to pan-specific primer design
665 for highly variable MSAs. Importantly, varVAMP is available through various bioinformatic
666 repositories and has also been deployed to Galaxy Europe (usegalaxy.eu), a web-based
667 platform for bioinformatic data analysis[55]. On the basis of the input MSA, varVAMP
668 generates consensus sequences and analyzes them for the presence of potential primer
669 sequences. From the subsequent pool of found primers, varVAMP chooses optimal
670 amplicons for specific molecular techniques such as tiled sequencing or qPCR and
671 introduces degenerated nucleotides into primer sequences to compensate for sequence
672 variations.

673 To demonstrate varVAMP's applicability, we used the software to design and evaluate
674 primers for tiled sequencing schemes of SARS-CoV-2, BoDV-1, HAV, HEV, PV and ratHEV
675 and for qPCR of PV. With the exception of SARS-CoV-2, these pathogens have a high
676 genomic diversity restricting conventional primer design. While varVAMP provides a fully
677 automated solution, selecting appropriate reference data prior to primer design can pose a
678 challenge. The core concept of varVAMP is based on the assumption that the sequences in
679 the input MSA reflect the majority of variations within the viral genome. However, in
680 sequence data repositories, particular in the case of viral sequences, there can be a lack of
681 associated metadata, the presence of recombinant sequences, bias towards sequencing
682 labs, geographic biases towards circulating strains or an underrepresented amount of
683 recently discovered or understudied viral pathogens[56,57]. Therefore, a prior careful data
684 selection is a key requirement for a successful primer design and might require the
685 additional use of clustering algorithms and phylogenetic assessment tools[29,30,58,59]. This was
686 most prominently observed for our ratHEV tiled sequencing scheme. Here, we could only
687 successfully generate full length sequences for the R63 but not for the pt2 isolate. The
688 ratHEV primer design was challenging as ratHEV full-length sequences are highly variable
689 and scarce in NCBI's genbank. In such cases, PCR reactions might not be successful for
690 viral genomes that are only distantly related to the majority of sequences in the input MSA.
691 The pathogenicity of ratHEV for humans has been only recently identified[60]. We expect an
692 increase in reported sequence data over the coming years, which will help to overcome
693 current limitations. The second challenge for varVAMP users is the proper selection of input
694 parameters. While varVAMP offers an automated solution, it often requires multiple rounds of
695 manual optimizations that should, however, be computationally inexpensive.

696 It is worth mentioning that we did not compare our primer design suite to other primer design
697 software, as we are not aware of other tools that handle similar sequence variability while
698 being able to design degenerate primers for tiled sequencing or qPCR. Primalscheme was
699 highly successful during the SARS-CoV-2 pandemic and has been also used to design tiled
700 amplicon schemes for viruses such as West-Nile or Monkeypox virus[61–63] but limits the
701 variability of the input alignment at 5%[10] . Despite the fact that varVAMP adapted some of
702 primalscheme's primer assessment functionality, a head-to-head comparison is not possible
703 as varVAMP was not designed as an alternative but as a solution for highly variable MSAs.

704 We did observe sample-dependent and independent amplicon performance. The molecular
705 reasons for sample-dependent PCR performance can be diverse. We show that primers lead
706 to amplification even if there are one or two mismatches between primer and target
707 sequence. However, if more mismatches are present, poorer amplification or complete
708 dropouts could be the result. For patient samples, we used archived material with varying

709 storage times and viral loads that both could have impacted PCR amplification. Another
710 problem can be the presence of sample-specific non-viral nucleic acid that could provide
711 additional primer binding sites and result in unspecific amplification. To address this issue,
712 varVAMP provides a BLAST module that allows users to check for potential off-target effects
713 with a custom database. Sample-independent amplicon performance is likely a PCR
714 optimization issue. Compared to popular SARS-CoV-2 schemes[64], we have not optimized
715 the PCR conditions, nor the primer pooling ratios. That was reflected by varying amplicon
716 performance in multiplex reactions. To give an example of how to further optimize the
717 multiplex PCRs, we adjusted the primer concentrations of the PV scheme and show this can
718 lead to a more balanced overall coverage.

719 All primer schemes described here have been developed because of a methodology gap.
720 HEV phylogeny and sub-typing is mostly restricted to Sanger sequencing of ORF2
721 fragments thereby neglecting viral evolution in the remaining genome[65,66]. Current
722 methodology to generate highly covered HEV full-genome sequences via Illumina or
723 Nanopore sequencing requires costly RNA-Seq protocols using hybridization probes[67,68] or
724 powerful sequencing machines. Our newly developed primer schemes could simplify
725 sequencing and aid, for example, analyses of Ribavirin resistance-associated mutations that
726 can develop in immunocompromised HEV patients[69]. For PV, molecular assays used for
727 detection and intratypic differentiation of serotypes are well suited for cell-culture isolates[70].
728 However, these techniques are laborious and time-consuming. qPCR and amplicon-based
729 NGS protocols for the rapid analysis of PV from patient and environmental samples are
730 imperative for fast public health decisions. There are protocols for high-throughput
731 sequencing of PV but similar to HEV they are restricted to a conserved part of the
732 enterovirus genome[71]. Our novel methods for PV detection by qPCR and whole-genome
733 sequencing could not only benefit existing surveillance programs but might also lay the
734 foundation for wastewater surveillance strategies within the global PV eradication
735 program[72,73].

736 In conclusion, the varVAMP pipeline was developed because primer design on the basis of
737 highly variable MSAs is difficult, time consuming and there are no automated solutions for
738 qPCR and tiled sequencing. The designed and validated primer schemes for the different
739 viruses are not only a proof-of-concept for varVAMP's applicability but have been developed
740 because they could directly benefit viral diagnostics and epidemiology. Laboratories that
741 have already established SARS-CoV-2 sequencing pipelines or in-house qPCR protocols
742 should be able to adapt their methodologies to these new primer schemes with only minor
743 modifications.

## DATA AVAILABILITY

varVAMP v.1.2.0 and BAMdash v.0.2.4 are open source and available at https://github.com/jonas-fuchs/varVAMP (DOI: 10.5281/zenodo.11125498) and https://github.com/jonas-fuchs/BAMdash (DOI: 10.5281/zenodo.10804160). The Galaxy version of varVAMP is available at https://usegalaxy.eu/root?tool_id=toolshed.g2.bx.psu.edu/repos/iuc/varvamp/varvamp/. The code and data to reproduce the figures is available at https://github.com/jonas-fuchs/varVAMP_in_silico_analysis (DOI: 10.5281/zenodo.10942525). All input multiple sequence alignments and varVAMP outputs for primers that have been evaluated in this study are available at: https://github.com/jonas-fuchs/ViralPrimerSchemes (DOI: 10.5281/zenodo.10562882). Raw sequencing data has been deposited at ENA under the study accession number: PRJEB74744.

## AUTHOR CONTRIBUTIONS

JF, MH, MS, WM, MP conceptualized the project. JF wrote the manuscript. JF developed varVAMP with help from WM. WM and BG deployed varVAMP to CONDA, DOCKER and Galaxy. TK critically evaluated the algorithms employed by varVAMP. MS and JF selected the input data for primer design. JF and JKl designed primers with varVAMP. JKl, MS, JKr, CW, LE, LJ, AM, CB, MB, JP, RJ, JW, JS, CM, SB, SS generated and evaluated the experimental data. JF, NB, JKl and MH performed the bioinformatic analyses. JF performed phylogenetic analyses and the *in silico* primer evaluation. All authors reviewed and edited the manuscript.

## ACKNOWLEDGEMENTS

## FUNDING

## CONFLICT OF INTEREST

The authors declare no conflict of interest.

## REFERENCES

1. Gire, S. K. *et al.* Genomic surveillance elucidates Ebola virus origin and transmission during the 2014 outbreak. *Science* **345**, 1369–1372 (2014).

2. Schemmerer, M., Wenzel, J. J., Stark, K. & Faber, M. Molecular epidemiology and genotype-specific disease severity of hepatitis E virus infections in Germany, 2010–2019. *Emerg. Microbes Infect.* **11**, 1754–1763 (2022).

3. Metsky, H. C. *et al.* Zika virus evolution and spread in the Americas. *Nature* **546**, 411–415 (2017).

4. Simner, P. J., Miller, S. & Carroll, K. C. Understanding the Promises and Hurdles of Metagenomic Next-Generation Sequencing as a Diagnostic Tool for Infectious Diseases. *Clin. Infect. Dis.* **66**, 778–788 (2018).

5. Jaki, L. *et al.* Total escape of SARS-CoV-2 from dual monoclonal antibody therapy in an immunocompromised patient. *Nat. Commun.* **14**, 1999 (2023).

6. Weigang, S. *et al.* Within-host evolution of SARS-CoV-2 in an immunosuppressed COVID-19 patient as a source of immune escape variants. *Nat. Commun.* **12**, 6405 (2021).

7. Lamers, M. M. *et al.* Human airway cells prevent SARS-CoV-2 multibasic cleavage site cell culture adaptation. *eLife* **10**, e66815 (2021).

8. Rose, R., Constantinides, B., Tapinos, A., Robertson, D. L. & Prosperi, M. Challenges in the analysis of viral metagenomes. *Virus Evol.* **2**, vew022 (2016).

9. Viral Hemorrhagic Fever Consortium *et al.* Capturing sequence diversity in metagenomes

811 with comprehensive and scalable probe design. *Nat. Biotechnol.* **37**, 160–168 (2019).

812 10. Quick, J. *et al.* Multiplex PCR method for MinION and Illumina sequencing of Zika

813 and other virus genomes directly from clinical samples. *Nat. Protoc.* **12**, 1261–1276

814 (2017).

815 11. Gohl, D. M. *et al.* A rapid, cost-effective tailed amplicon method for sequencing

816 SARS-CoV-2. *BMC Genomics* **21**, 863 (2020).

817 12. Chen, Z. *et al.* Global landscape of SARS-CoV-2 genomic surveillance and data

818 sharing. *Nat. Genet.* **54**, 499–507 (2022).

819 13. Kralik, P. & Ricchi, M. A Basic Guide to Real Time PCR in Microbial Diagnostics:

820 Definitions, Parameters, and Everything. *Front. Microbiol.* **8**, (2017).

821 14. Linhart, C. & Shamir, R. The degenerate primer design problem. in *ISMB* 172–181

822 (2002).

823 15. Dieffenbach, C. W., Lowe, T. M. & Dveksler, G. S. General concepts for PCR primer

824 design. *Genome Res.* **3**, S30–S37 (1993).

825 16. Bustin, S. & Huggett, J. qPCR primer design revisited. *Biomol. Detect. Quantif.* **14**,

826 19–28 (2017).

827 17. Untergasser, A. *et al.* Primer3—new capabilities and interfaces. *Nucleic Acids Res.*

828 **40**, e115–e115 (2012).

829 18. Guo, J., Starr, D. & Guo, H. Classification and review of free PCR primer design

830 software. *Bioinformatics* **36**, 5263–5268 (2021).

831 19. Smith, D. B. & Simmonds, P. Classification and Genomic Diversity of Enterically

832 Transmitted Hepatitis Viruses. *Cold Spring Harb. Perspect. Med.* **8**, a031880 (2018).

833 20. Borcard, L. *et al.* Investigating the Extent of Primer Dropout in SARS-CoV-2 Genome

834 Sequences During the Early Circulation of Delta Variants. *Front. Virol.* **2**, 840952 (2022).

835 21. Hugerth, L. W. *et al.* DegePrime, a Program for Degenerate Primer Design for

836 Broad-Taxonomic-Range PCR in Microbial Ecology Studies. *Appl. Environ. Microbiol.* **80**,

837 5116–5123 (2014).

838 22. Rosenkranz, D. easyPAC: A Tool for Fast Prediction, Testing and Reference Mapping

839 of Degenerate PCR Primers from Alignments or Consensus Sequences. *Evol. Bioinforma.*

840 **8**, EBO.S8870 (2012).

841 23. Arvidsson, S., Kwasniewski, M., Riaño-Pachón, D. M. & Mueller-Roeber, B.

842 QuantPrime – a flexible tool for reliable high-throughput primer design for quantitative

843    PCR. *BMC Bioinformatics* **9**, 465 (2008).

844 24.    Dijkstra, E. W. A Note on Two Problems in Connexion with Graphs. in *Edsger Wybe*
845    *Dijkstra* (eds. Apt, K. R. & Hoare, T.) 287–290 (ACM, New York, NY, USA, 2022).
846    doi:10.1145/3544585.3544600.

847 25.    Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local
848    alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990).

849 26.    Ye, J. *et al.* Primer-BLAST: A tool to design target-specific primers for polymerase
850    chain reaction. *BMC Bioinformatics* **13**, 134 (2012).

851 27.    Smith, D. B. *et al.* Update: proposed reference sequences for subtypes of hepatitis E
852    virus (species Orthohepevirus A). *J. Gen. Virol.* **101**, 692–698 (2020).

853 28.    Pearson, W. R. Finding Protein and Nucleotide Similarities with FASTA. *Curr. Protoc.*
854    *Bioinforma.* **53**, (2016).

855 29.    Rognes, T., Flouri, T., Nichols, B., Quince, C. & Mahé, F. VSEARCH: a versatile open
856    source tool for metagenomics. *PeerJ* **4**, e2584 (2016).

857 30.    Minh, B. Q. *et al.* IQ-TREE 2: New Models and Efficient Methods for Phylogenetic
858    Inference in the Genomic Era. *Mol. Biol. Evol.* **37**, 1530–1534 (2020).

859 31.    Girgis, H. Z. MeShClust v3.0: high-quality clustering of DNA sequences using the
860    mean shift algorithm and alignment-free identity scores. *BMC Genomics* **23**, 423 (2022).

861 32.    Katoh, K. & Toh, H. Recent developments in the MAFFT multiple sequence alignment
862    program. *Brief. Bioinform.* **9**, 286–298 (2008).

863 33.    Johne, R. *et al.* Detection of a novel hepatitis E-like virus in faeces of wild rats using
864    a nested broad-spectrum RT-PCR. *J. Gen. Virol.* **91**, 750–758 (2010).

865 34.    Niller, H. H. *et al.* Zoonotic spillover infections with Borna disease virus 1 leading to
866    fatal human encephalitis, 1999–2019: an epidemiological investigation. *Lancet Infect. Dis.*
867    **20**, 467–477 (2020).

868 35.    Neumann, B. *et al.* Antibodies against viral nucleo-, phospho-, and X protein
869    contribute to serological diagnosis of fatal Borna disease virus 1 infections. *Cell Rep.*
870    *Med.* **3**, 100499 (2022).

871 36.    Bauswein, M. *et al.* Human Infections with Borna Disease Virus 1 (BoDV-1) Primarily
872    Lead to Severe Encephalitis: Further Evidence from the Seroepidemiological BoSOT
873    Study in an Endemic Region in Southern Germany. *Viruses* **15**, 188 (2023).

874 37.    Schemmerer, M., Johne, R., Erl, M., Jilg, W. & Wenzel, J. J. Isolation of Subtype 3c,

875  3e and 3f-Like Hepatitis E Virus Strains Stably Replicating to High Viral Loads in an
876  Optimized Cell Culture System. *Viruses* **11**, 483 (2019).

877  38.  Paul, A. V. *et al.* The entire nucleotide sequence of the genome of human hepatitis A
878  virus (isolate MBB). *Virus Res.* **8**, 153–171 (1987).

879  39.  Johne, R. *et al.* Novel Hepatitis E Virus Genotype in Norway Rats, Germany. *Emerg.*
880  *Infect. Dis.* **16**, 1452–1455 (2010).

881  40.  Sridhar, S. *et al.* Transmission of Rat Hepatitis E Virus Infection to Humans in Hong
882  Kong: A Clinical and Epidemiological Analysis. *Hepatology* **73**, 10–22 (2021).

883  41.  Schemmerer, M., Erl, M. & Wenzel, J. J. HuH-7-Lunet BLR Cells Propagate Rat
884  Hepatitis E Virus (HEV) in a Cell Culture System Optimized for HEV. *Viruses* **14**, 1116
885  (2022).

886  42.  Panajotov, J., Falkenhagen, A., Gadicherla, A. K. & Johne, R. Molecularly generated
887  rat hepatitis E virus strains from human and rat show efficient replication in a human
888  hepatoma cell line. *Virus Res.* **344**, 199364 (2024).

889  43.  Chen, S., Zhou, Y., Chen, Y. & Gu, J. fastp: an ultra-fast all-in-one FASTQ
890  preprocessor. *Bioinformatics* **34**, i884–i890 (2018).

891  44.  Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows–Wheeler
892  transform. *Bioinformatics* **26**, 589–595 (2010).

893  45.  Wilm, A. *et al.* LoFreq: a sequence-quality aware, ultra-sensitive variant caller for
894  uncovering cell-population heterogeneity from high-throughput sequencing datasets.
895  *Nucleic Acids Res.* **40**, 11189–11201 (2012).

896  46.  Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,
897  2078–2079 (2009).

898  47.  Li, H. New strategies to improve minimap2 alignment accuracy. *Bioinformatics* **37**,
899  4572–4574 (2021).

900  48.  Brandt, C. *et al.* poreCov-An Easy to Use, Fast, and Robust Workflow for
901  SARS-CoV-2 Genome Reconstruction via Nanopore Sequencing. *Front. Genet.* **12**,
902  711437 (2021).

903  49.  Aslan, A. T. & Balaban, H. Y. Hepatitis E virus: Epidemiology, diagnosis, clinical
904  manifestations, and treatment. *World J. Gastroenterol.* **26**, 5543–5560 (2020).

905  50.  Nicot, F. *et al.* Classification of the Zoonotic Hepatitis E Virus Genotype 3 Into Distinct
906  Subgenotypes. *Front. Microbiol.* **11**, 634430 (2021).

907 51.    Lhomme, S. *et al.* Insertions and Duplications in the Polyproline Region of the
908    Hepatitis E Virus. *Front. Microbiol.* **11**, 1 (2020).

909 52.    Belalov, I. S. & Lukashev, A. N. Causes and Implications of Codon Usage Bias in
910    RNA Viruses. *PLoS ONE* **8**, e56642 (2013).

911 53.    Gu, H., Fan, R. L. Y., Wang, D. & Poon, L. L. M. Dinucleotide evolutionary dynamics
912    in influenza A virus. *Virus Evol.* **5**, vez038 (2019).

913 54.    Gerloff, N. *et al.* Diagnostic Assay Development for Poliovirus Eradication. *J. Clin.*
914    *Microbiol.* **56**, e01624-17 (2018).

915 55.    Jalili, V. *et al.* The Galaxy platform for accessible, reproducible and collaborative
916    biomedical analyses: 2020 update. *Nucleic Acids Res.* **48**, W395–W402 (2020).

917 56.    Liu, P., Song, Y., Colijn, C. & MacPherson, A. The impact of sampling bias on viral
918    phylogeographic reconstruction. *PLOS Glob. Public Health* **2**, e0000577 (2022).

919 57.    Kieft, K. & Anantharaman, K. Virus genomics: what is being overlooked? *Curr. Opin.*
920    *Virol.* **53**, 101200 (2022).

921 58.    Balaban, M., Moshiri, N., Mai, U., Jia, X. & Mirarab, S. TreeCluster: Clustering
922    biological sequences using phylogenetic trees. *PLOS ONE* **14**, e0221068 (2019).

923 59.    Zou, Q., Lin, G., Jiang, X., Liu, X. & Zeng, X. Sequence clustering in bioinformatics:
924    an empirical study. *Brief. Bioinform.* (2018) doi:10.1093/bib/bby090.

925 60.    Sridhar, S. *et al.* Rat Hepatitis E Virus as Cause of Persistent Hepatitis after Liver
926    Transplant. *Emerg. Infect. Dis.* **24**, 2241–2250 (2018).

927 61.    Isabel, S. *et al.* Targeted amplification-based whole genome sequencing of
928    *Monkeypox virus* in clinical specimens. *Microbiol. Spectr.* **12**, e02979-23 (2024).

929 62.    Tešović, B. *et al.* Development of multiplex PCR based NGS protocol for whole
930    genome sequencing of West Nile virus lineage 2 directly from biological samples using
931    Oxford Nanopore platform. *Diagn. Microbiol. Infect. Dis.* **105**, 115852 (2023).

932 63.    Hourdel, V. *et al.* Rapid Genomic Characterization of SARS-CoV-2 by Direct
933    Amplicon-Based Sequencing Through Comparison of MinION and Illumina iSeq100TM
934    System. *Front. Microbiol.* **11**, 571328 (2020).

935 64.    Lambisia, A. W. *et al.* Optimization of the SARS-CoV-2 ARTIC Network V4 Primers
936    and Whole Genome Sequencing Protocol. *Front. Med.* **9**, 836728 (2022).

937 65.    Bruni, R. *et al.* Hepatitis E virus genotypes and subgenotypes causing acute
938    hepatitis, Bulgaria, 2013–2015. *PLOS ONE* **13**, e0198045 (2018).

66. Inoue, J., Takahashi, M., Yazaki, Y., Tsuda, F. & Okamoto, H. Development and validation of an improved RT-PCR assay with nested universal primers for detection of hepatitis E virus strains with significant sequence divergence. *J. Virol. Methods* **137**, 325–333 (2006).

67. Davis, C. A. *et al.* Hepatitis E virus: Whole genome sequencing as a new tool for understanding HEV epidemiology and phenotypes. *J. Clin. Virol.* **139**, 104738 (2021).

68. Schilling-Loeffler, K. *et al.* Cell Culture Isolation and Whole Genome Characterization of Hepatitis E Virus Strains from Wild Boars in Germany. *Microorganisms* **9**, 2302 (2021).

69. Kamar, N. *et al.* Ribavirin for Hepatitis E Virus Infection After Organ Transplantation: A Large European Retrospective Multicenter Study. *Clin. Infect. Dis.* **71**, 1204–1211 (2020).

70. Kilpatrick, D. R. *et al.* Poliovirus serotype-specific VP1 sequencing primers. *J. Virol. Methods* **174**, 128–130 (2011).

71. Shaw, A. G. *et al.* Rapid and Sensitive Direct Detection and Identification of Poliovirus from Stool and Environmental Surveillance Samples by Use of Nanopore Sequencing. *J. Clin. Microbiol.* **58**, e00920-20 (2020).

72. Chumakov, K., Ehrenfeld, E., Agol, V. I. & Wimmer, E. Polio eradication at the crossroads. *Lancet Glob. Health* **9**, e1172–e1175 (2021).

73. Cooper, L. V. *et al.* Risk factors for the spread of vaccine-derived type 2 polioviruses after global withdrawal of trivalent oral poliovirus vaccine and the effects of outbreak responses with monovalent vaccine: a retrospective analysis of surveillance data for 51 countries in Africa. *Lancet Infect. Dis.* **22**, 284–294 (2022).

## TABLE AND FIGURE LEGENDS

**Table 1. Summary of varVAMP designs for tiled sequencing.** The table lists important alignment statistics, varVAMP input parameters and output information including the varVAMP version. Pairwise sequence identity was calculated with Identity (https://github.com/BioinformaticsToolsmith/Identity). All primers and respective varVAMP outputs are accessible at: https://github.com/jonas-fuchs/ViralPrimerSchemes.

**Table 2**. **Summary of varVAMP designs for qPCR.** Important alignment statistics are listed together with varVAMP input parameters and output information including the varVAMP version. Pairwise sequence identity was calculated with Identity (https://github.com/BioinformaticsToolsmith/Identity). All primers and respective varVAMP outputs are accessible at: https://github.com/jonas-fuchs/ViralPrimerSchemes.

**972** **Figure 1. Schematic varVAMP overview and example output. (a)** Overview of the
**973** varVAMP workflow. White boxes represent steps of the pipeline that are common to all
**974** modes. Consecutive steps are connected by arrows and optional steps are indicated with a
**975** dotted border. Colored boxes mark unique steps for each varVAMP mode (blue - single,
**976** orange - tiled, green - qPCR). Steps that produce outputs end in schematic folder icons for
**977** the main output and the additional data subfolder. (n - number, nt - nucleotide). **(b)** Example
**978** overview plot that is produced when running varVAMP. This plot was generated with
**979** varVAMP's tiled mode on the example MSA of HEV-3 sequences provided as example data
**980** within the varVAMP github repository. Shown is the normalized Shannon's entropy for each
**981** alignment position (gray) and its rolling average over 10 nucleotides (black curve). The
**982** orange boxes below the plot mark the start and stop MSA positions of potential primer
**983** regions (regions that have, in this case, a maximum of 4 ambiguous bases within the
**984** minimal primer length of 19) and the gray and light gray boxes mark all considered forward
**985** and reverse primers, respectively. The final scheme that was selected by the graph search
**986** for overlapping amplicons (blue) with low-penalty primers (red) is depicted at the bottom.

**987** **Figure 2. Primer design and tiled sequencing of HEV-3. (a)** Schematic overview of the
**988** data preparation steps preceding primer design. All full-length sequences of HEV were
**989** downloaded from NCBI, sub-genotyped with fasta36 and clustered by similarity with vsearch.
**990** The clustering result was evaluated by phylogenetic tree construction. Afterwards, clusters
**991** comprising multiple subgenotypes were aligned with MAFFT and the MSA used as the input
**992** for varVAMP **(b)** Phylogenetic tree of full-length HEV sequences constructed with IQ-TREE 2
**993** (GTR+F+R10, 1000 bootstrap replicates). The vsearch clustering results for each sequence
**994** is displayed in colors and the HEV genotypes and subgenotypes are indicated at the
**995** respective branches (n = number of sequences). **(c)** Agarose electrophoresis pictures of the
**996** individual PCR products for the cluster 2 (upper plot) and cluster 4 (lower plot) primer
**997** schemes tested with supernatant of HEV-3 f or HEV-3 c stably infected PLC/PRF/5 cells,
**998** respectively. Triangles indicate bands at the expected molecular weight of the PCR products
**999** (kb - kilobases) **(d)** Coverage plots of the Illumina sequencing results of the in (c) amplified
**1000** PCR products for cluster 2 (upper plot) and cluster 4 (lower plot) mapped to their respective
**1001** NCBI reference sequences MK089847 and MK089849. Below each coverage plot the
**1002** genomic start and stop positions of each amplicon are displayed as gray boxes with their
**1003** respective amplicon number. Dotted lines indicate mean coverages. Coverage plots were
**1004** created with BAMdash (individual coverage plots are given in S1). **(e)** Genome recovery of
**1005** HEV-3 persistently infected cell cultures and sub-genotyped HEV-3 positive blood samples
**1006** that were subjected to their respective tiled amplicon workflow for cluster 2 (upper plot) or

1007 cluster 4 (lower plot). Genome recovery was calculated as % of reference nucleotides 1008 covered at least 20 fold. All PCR reactions were performed in the singleplex setting.

1009 **Figure 3. *In silico* evaluation of novel tiled primer schemes for SARS-CoV-2, BoDV-1,** 1010 **HAV, HEV, PV and ratHEV. (a)** Normalized Shannon's entropy (1% rolling average) of the 1011 MSA used as the varVAMP input. **(b)** Number of permutations (degeneracy) of each primer 1012 in the tiled sequencing scheme for the respective viruses. Each dot shows the degeneracy of 1013 a single primer. Horizontal lines indicate the mean. (n - number) **(c)** Cumulative counts of 1014 mismatches between primers and sequences in the varVAMP input MSA. For each primer 1015 the number of mismatches with each sequence of the MSA was counted if it was not 1016 covered by any primer permutation. Shown are the cumulative mismatches between primers 1017 and MSA sequences in the tiled primer schemes for the respective viruses. Dot area size is 1018 proportionate to the percentage. **(d)** Analogous to (c) the mismatches with the MSA 1019 sequences were counted per primer nucleotide position and averaged over all primers in a 1020 scheme. As primers vary in their length, the % mismatches are displayed starting at the 1021 primer's 3' end (position 0 is the most 3' nucleotide position). The gray triangle schematically 1022 indicates the primer positions that varVAMP penalizes and the position-specific penalty 1023 multipliers (32, 16, 8, 4, 2). **(e-h)** Primer melting temperatures **(e)**, hairpin temperatures **(f)**, 1024 homo-dimer temperatures **(g)** or the GC content **(h)** were calculated either for the primer 1025 sequence including the most common nucleotides or averaged over all permutations of the 1026 final primer sequences that include degenerate nucleotides. (e-f) were calculated with 1027 primer3. Each dot represents a single primer of the respective tiled primer schemes. Dotted 1028 lines indicate the upper target cut-offs or target ranges employed by varVAMP (nt - 1029 nucleotide).

1030 **Figure 4. Whole genome sequencing utilizing the SARS-CoV-2, BoDV-1, HAV, PV and** 1031 **ratHEV primer schemes.** Representative coverage plots (left) and % genome recovery 1032 (right) of the different **(a)** SARS-CoV-2, **(b)** BoDV-1, **(c)** HAV, **(d)** PV and **(e)** ratHEV samples 1033 subjected to their respective tiled amplicon whole genome sequencing workflow. Coverage 1034 plots were created with BAMdash. Dotted lines indicate mean coverages. Reference 1035 genomes used for mapping are indicated in the header of the coverage plots (individual 1036 coverage plots are given in S1). Genome recovery was calculated as % of reference 1037 nucleotides covered at least 20 fold (sp - single plex, mp - multiplex). Dark grey bars - ONT 1038 generated data, light grey bars - Illumina generated data.

1039 **Figure 5. Amplicon performance and mismatch analysis. (a)** For each sequencing result 1040 using the virus specific primer scheme the amplicon recovery (upper panel) and normalized 1041 coverage (lower panel) was calculated. Each color represents an individual amplicon tracked 1042 over different samples. Amplicon recovery was calculated as % of reference nucleotides

covered at least 20 fold between the genomic start and stop position of the individual amplicons. For the normalized amplicon coverage, the mean coverage was calculated for each amplicon and normalized to the highest covered amplicon of each scheme (set to 100). **(b)** For each sequencing result, each primer binding region was analyzed for the number of mismatches not covered by any permutation of the corresponding primer sequence. Mutations were only considered if their variant frequency was >= 0.7. Primers were excluded from the analysis if any primer binding position was not covered at least 20-fold. **(c)** Dumbbell plot showing the pairwise identities of the newly generated fasta consensus sequences (blue dot) or the sequences of the varVAMP input MSA (dark gray dot) of each respective primer scheme. Light gray and red lines indicate the percent pairwise identity increase or decrease, respectively. Significance was calculated with a Welch's *t*-test (n.d. - not determined as n < 3, n.s. - not significant, *: $p \leq 0.05$, **: $p \leq 0.05$).

**Figure 6. Specificity and sensitivity of the novel PV qPCR schemes.** qPCR primers specific for **(a)** PV1, **(b)** PV2 and **(c)** PV3 were tested on serial RNA dilutions extracted from viral supernatants of Sabin 1, Sabin 2 and Sabin 3 infected cell cultures (n=3). The fluorescence was measured during the extension step with three different channel setups respective to the probe fluorophore (FAM = 465-510 nm, JOE = 533-580 nm, CY5 = 618-660 nm). Amplification curves were analyzed using the Roche LightCycler 480 II device software.

**S1. Coverage plots for all sequencing results.** Coverage plots of the different SARS-CoV-2, BoDV-1, HAV, PV and ratHEV samples subjected to their respective tiled amplicon whole genome Illumina sequencing workflow. Dotted lines indicate mean coverages. NCBI accession numbers of the reference sequences used for mappings are indicated in the headers. Coverage plots were created with BAMdash.
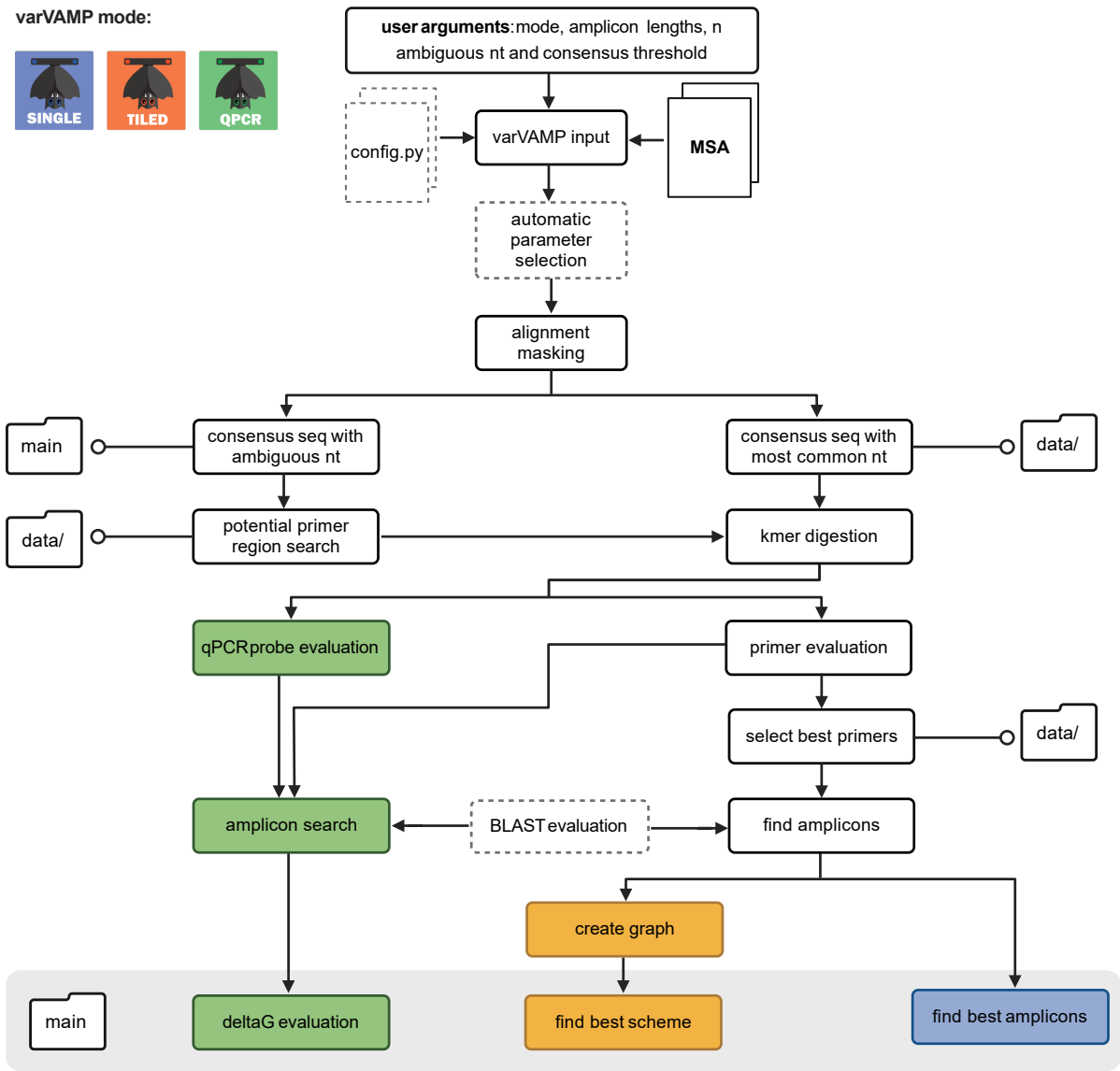
**S2. Primer balancing for PV whole genome sequencing.** The PV tiled primers initially used in equimolar concentrations for multiplex reactions were balanced in two consecutive rounds based on prior results and then the balanced primers were used in multiplex PCR reactions for Sabin 1-3 and the amplicons subjected to Illumina sequencing. The respective concentrations for each iteration are given above the coverage plots (blue arrow - increase in molarity, gray arrow - no change in molarity, red arrow - decrease in molarity). Dotted lines indicate mean coverages. NCBI accession numbers of the reference sequences used for mappings are indicated in the headers. Coverage plots were created with BAMdash.

**Table 1**

| alignment statistics | | | | varVAMP parameters | | | | varVAMP output | | |
|---|---|---|---|---|---|---|---|---|---|---|
| virus | subtypes | n sequences | % mean sequence identiy | max ambig bases | threshold | optimal amplicon size | maximum amplicon size | expected recovery | n amplicons | varVAMP version |
| SARS-CoV-2 | B.1 - XBB | 865 | 99 ± 1 | 1 | 0.99875 | 700 | 800 | 99.72 % | 55 | v.0.9.4 |
| BoDV-1 | all | 55 | 89 ± 8 | 2 | 0.94 | 400 | 550 | 98.6 % | 27 | v.0.6 |
| HAV | all | 309 | 81 ± 10 | 4 | 0.93 | 1000 | 1600 | 95.65 % | 7 | v.0.8.3 |
| HEV genotype 3 | f, e | 376 | 76 ± 6 | 4 | 0.91 | 1000 | 1500 | 99.02 % | 7 | v.0.8.2 |
| HEV genotype 3 | c, h1, m, i, uc, l | 201 | 75 ± 9 | 4 | 0.90 | 1000 | 1500 | 99.28 % | 6 | v.0.8.2 |
| PV | 1-3 | 944 | 71 ± 13 | 4 | 0.91 | 1000 | 1400 | 99.63 % | 7 | v.0.8 |
| ratHEV | all | 41 | 57 ± 10 | 5 | 0.82 | 1200 | 1700 | 97.41 % | 6 | v.0.8.3 |

**Table 2**

| alignment statistics | | | | varVAMP parameters | | | | varVAMP output | |
|---|---|---|---|---|---|---|---|---|---|
| virus | subtypes | n sequences | % mean sequence identiy | primer max ambig bases | probe max ambig bases | threshold | ΔG cutoff (kcal/mol) | n found schemes | varVAMP version |
| PV | 1 | 241 | 86 ± 12 | 2 | 1 | 0.93 | -3 | 3 | 0.7 |
| PV | 2 | 494 | 88 ± 8 | 1 | 0 | 0.98 | -3 | 4 | 0.7 |
| PV | 3 | 209 | 91 ± 9 | 2 | 1 | 0.93 | -3 | 3 | 0.7 |

a



b



**Figure 1**

**Figure 2**

**Figure 3**

**Figure 4**

**Figure 5**

**Figure 6**