

1 Exploring the limit of using a deep neural 2 network on pileup data for germline 3 variant calling

4
5 Ruibang Luo*, Chak-Lim Wong, Yat-Sing Wong, Chi-Ian Tang, Chi-Man Liu, Chi-Ming Leung,
6 Tak-Wah Lam*

7
8 Department of Computer Science, The University of Hong Kong, Hong Kong, China

9
10 * Correspondence and requests for materials should be addressed to R. L. (email:
11 rbluo@cs.hku.hk) and T. L. (twlam@cs.hku.hk)

13 Abstract

14 Single-molecule sequencing technologies have emerged in recent years and revolutionized
15 structural variant calling, complex genome assembly, and epigenetic mark detection.

16 However, the lack of a highly accurate small variant caller has limited the new technologies
17 from being more widely used. In this study, we present Clair, the successor to Clairvoyante,
18 a program for fast and accurate germline small variant calling, using single molecule
19 sequencing data. For ONT data, Clair achieves the best precision, recall and speed as
20 compared to several competing programs, including Clairvoyante, Longshot and Medaka.

21 Through studying the missed variants and benchmarking intentionally overfitted models, we
22 found that Clair may be approaching the limit of possible accuracy for germline small variant
23 calling using pileup data and deep neural networks. Clair requires only a conventional CPU
24 for variant calling and is an open source project available at [https://github.com/HKU-](https://github.com/HKU-BAL/Clair)
25 [BAL/Clair](https://github.com/HKU-BAL/Clair).

26 Introduction

27 Fast and accurate variant calling is essential for both research and clinical applications of
28 human genome sequencing^{1,2}. Algorithms, best practices and benchmarking guidelines have
29 been established for how to use Illumina sequencing to call germline small variants,
30 including single-nucleotide polymorphisms (SNPs) and insertions/deletions (indels)³⁻⁶. In
31 recent years, single-molecule sequencing (SMS) technologies have emerged for a variety of
32 important applications⁷. These technologies, which are also known as the third-generation
33 sequencing technologies, generate sequencing reads two to three orders of magnitude
34 longer than Illumina reads (10–100kbp versus 100–250bp). The long read length has made
35 the new SMS technologies, including Pacific Biosciences (PacBio) and Oxford Nanopore
36 Technology (ONT), unprecedentedly powerful for resolving complex genome assembly
37 problems and for detecting large structural variants⁸. However, currently available SMS
38 technologies also have a significantly higher base error rate of 3–15%⁹, making the variant
39 calling methods previously designed for Illumina sequencing inapplicable to SMS
40 technologies. The lack of accurate tools for efficient variant calling has limited SMS
41 technologies from being applied to the many problems that require SNPs and small indels.
42

43 In our previous work, we developed Clairvoyante¹⁰, a germline small variant caller for single
44 molecule sequencing data. Clairvoyante does not require sequence assembly and calls
45 variants directly from read alignments. Clairvoyante adopts a deep convolutional neural
46 network, so that by using the truth variants called and orthogonally verified in seven human
47 individuals by the Genome In A Bottle (GIAB) consortium¹¹⁻¹³, Clairvoyante can be trained
48 for variant calling on any new type of sequencing data without the need to look into its

49 error profile and build a hand-crafted model. Clairvoyante takes pileup data as input and
50 runs quickly. However, Clairvoyante's design is unable to call multiallelic variants or indels
51 longer than four bases. These defects remain to be solved. Meanwhile, the limit of using
52 pileup data and deep neural networks for variant calling remains to be explored.

53

54 In this study, we present Clair, a fast and accurate system for germline small variant calling
55 using single molecule sequencing data. With an entirely different network architecture and
56 learning tasks (i.e. output components), Clair resolves the multiallelic and long indel variant
57 calling problems that have prevented Clairvoyante from calling all types of small variants.
58 We describe in detail the methods we tried that either worked or did not work for
59 improving Clair's performance. For ONT datasets¹⁴, our experiments on whole-genome
60 variant calling in GIAB samples show that Clair outperforms Clairvoyante and other variant
61 callers, including Longshot¹⁵ and Medaka¹⁶, in terms of precision, recall and speed. For high
62 accuracy reads, including both PacBio CCS (Circular Consensus Sequencing)¹⁷ and Illumina
63 datasets¹³, DeepVariant¹⁸ had modestly improved F1-scores over Clair by .11% to .13%,
64 although Clair was seven times faster. Looking into the false positive (FP) and false negative
65 (FN) variants of the three sequencing technologies showed that except for variants with
66 insufficient coverage by chance, most of the others could be resolved using complete read
67 alignments instead of pileup data or else could not be resolved at all, even with a manual
68 inspection.

69 Results

70 Overview of Clair

71 Clair is a four-task, five-layer recurrent neural network with two bi-directional LSTM layers
72 followed by three feedforward layers (**Figure 1**). Clair takes a BAM file as input to find
73 candidate variants with any minor allele frequencies larger than a threshold (typically
74 between 0.1 and 0.2), and then computes a pileup of the candidates and converts the
75 summaries into a tensor. In a tensor, the allelic counts of bases and gaps on both strands of
76 a candidate variant and its 16 flanking bases are encoded into 1,056 integer values. More
77 details and pseudo code are available in the Methods section. As discussed in the
78 Clairvoyante paper, one major unsolved problem was how to support the calling of multi-
79 allelic variants (i.e., variants with two alternative alleles). In Clair, the problem is solved by
80 using four new (deep learning) tasks that are entirely different from Clairvoyante. These are:
81 1) a 21-genotype probabilistic model with 21 probability outputs; 2) the use of three
82 probabilities for the input, including a homozygous reference (0/0 genotype), a
83 heterozygous variant (0/1) or a homozygous variant (1/1); 3) the length of the first indel
84 allele, with 33 probabilities representing a length of '<-15bp', '-15bp', '-14bp', ..., '-1bp',
85 '0bp', '1bp', ..., '15bp', '>15bp'; and 4) the length of the second indel allele. The 21-genotype
86 probabilistic model can represent all possible genotypes of a diploid sample at the genome
87 position. The length of indels longer than 15bp cannot be directly inferred from the third
88 and fourth tasks, so Clair includes an additional step that re-scans the alignments. More
89 details on each of these steps can be found in the Methods section. The four tasks make
90 their own decisions and are designed to cross-validate each other. For example, task two is
91 a coarse-grained version of task one and can veto the decision made by task one. Tasks

92 three and four should indicate 0bp indel length if an SNP variant is decided by task one.
93 More details on how the four tasks make a joint decision are available in the Methods
94 section. We used the ‘focal loss’ deep-learning technique to solve the problem of
95 unbalanced variant types in training data. We used the ‘cyclical learning rate’ deep learning
96 technique to achieve the maximum possible variant calling performance and speed up the
97 training process to be able to handle larger training datasets. To improve Clair’s
98 performance at lower sequencing coverages, we augmented the training data with 10
99 subsampled coverages of each dataset. The parameters of these three new techniques are
100 in the Methods section.

101

102 Clair has 2,377,818 parameters, which is 45.7% more than Clairvoyante (1,631,496
103 parameters) but only one tenth as many as DeepVariant (23,885,392 parameters). In terms
104 of variant calling speed, Clair takes about 30 minutes, 1.5 hour, and 5 hours for a 50-fold
105 coverage WGS sample using Illumina, PacBio CCS and ONT data, respectively, using 24 CPU
106 cores. In our experiments, Clair was 10–20% slower than Clairvoyante, but significantly
107 faster than DeepVariant, Longshot and Medaka.

108

109 The Methods section includes a description of procedures to augment the training data or
110 improve Clair’s network architecture that we tested but that did not improve precision and
111 recall of variant calling. Developers working on further improving Clair’s performance can
112 save time by avoiding the same methods, or the same settings in a method.

113

114 Performance on ONT

115 ONT datasets are currently available for two GIAB samples, HG001 and HG002. The HG001
116 rel6 dataset generated by the Nanopore WGS Consortium¹⁴ contains approximately 44.3-
117 fold coverage of human genome (the dataset is also referred to as 1:44x, where '1' means
118 the sample suffix and '44x' means the coverage). The rel6 dataset was base-called with
119 Guppy 2.3.8, using the HAC (High-ACcuracy) model. In addition to the rel6 dataset, we
120 obtained a separate 124.1-fold coverage dataset for HG001 (1:124x) directly from Oxford
121 Nanopore (Philipp Rescheneder, personal communication). That dataset was base-called
122 with Guppy 2.2.3 using the Flip-Flop model. In some experiments, we combined 1:44x and
123 1:124x to form a new dataset 1:168x to maximize the coverage. For HG002, we used a
124 dataset with ~64-fold coverage (2:64x) from the GIAB consortium, which was base-called
125 with Guppy 2.3.5 using the Flip-Flop model. The links to the datasets are available in the
126 Supplementary Notes. The details about "the GIAB truth variant datasets", and "the
127 benchmarking methods and metrics" are available in "Methods – Benchmarking". Please
128 note that we have removed the low-complexity regions defined by GA4GH (The Global
129 Alliance for Genomics and Health)⁶ from benchmarking and we suggest removing these
130 regions as a common practice in calling small variants using the current ONT data. The
131 details of the regions removed, and the performance differences before and after removing
132 the low complexity regions using different sequencing technologies are available in
133 "Methods – Benchmarking".

134

135 **Figure 2** shows the precision and recall of Clair and other variant callers on SNPs and indels
136 in multiple experiments with ONT data. Supplementary Table 1 contains more details,
137 including precision, recall and F1-score in five categories, including overall, SNP, indel,

138 insertion, and deletion. Our results show that Clair not only outperformed other variant
139 callers, including Clairvoyante, Longshot, and Medaka, but also ran much faster. Using
140 1:168x|2:64x (i.e., test variant calling using HG002 with 64-fold coverage against a model
141 trained using HG001 with 168-fold coverage) as Clair’s primary result, Clair achieved 98.36%
142 precision, 96.46% recall, and 97.40% F1-score overall performance. In terms of SNPs, the
143 three metrics were 99.29%, 97.78% and 98.53%, respectively. For indels, they were
144 somewhat lower at 81.15%, 73.88%, and 77.34%. Clair significantly outperformed its
145 predecessor Clairvoyante on both SNP and indel calling (overall F1-score 97.40% versus
146 93.45%). Clair had a slightly higher F1-score on SNPs than Longshot (98.53% versus 98.41%),
147 but Longshot detects only SNPs, and Clair ran five times faster than Longshot (320 versus
148 1,797 minutes). Clair had a better performance than Medaka (overall F1-score 97.40%
149 versus 94.81%) and ran 30 times faster (320 versus 10,817 minutes). It is worth mentioning
150 that we didn’t benchmark Nanopolish¹⁹, which is also capable of variant calling on ONT data,
151 because it also requires raw signals as input, which are not publicly available for HG002.

152

153 We ran further experiments to answer five additional questions about Clair, as follows.

154

155 **Is the Clair model reference-genome specific?** In our experiments, performance did not
156 depend on whether we used GRCh37 or GRCh38. The performance of 1:168x|2:64x and
157 1:168x|2:64x(b37) was similar; the latter experiment tested HG002 GRCh37 read alignments
158 on a model trained using HG001 GRCh38 read alignments. Actually, 1:168x|2:64x(b37)
159 performed slightly better than 1:168x|2:64x, with a 0.18% better F1-score on SNPs, and
160 1.4% on indels.

161

162 **Does higher coverage in the test sample helps improve variant calling performance?** Yes,
163 but improvement seems to asymptote at ~60-fold coverage. In a comparison of
164 1:168x|2:64x to 1:168x|2:32x, the overall F1-score increased from 94.10% to 97.40%
165 (+3.30%), the SNP from 95.51% to 98.53% (+3.02%), and the indel from 68.87% to 77.34%
166 (+8.47%). Further increasing the coverage in the test sample will not significantly increase
167 the variant calling performance as we discuss below.

168

169 **Does higher coverage for model training help improve variant calling performance?** Yes,
170 but it depends on the coverage of the test sample. In a comparison of 1:124x|2:64x to
171 1:44x|2:64x, the overall F1-score increased from 96.84% to 97.51% (+0.67%), the SNP from
172 98.01% to 98.54% (+0.53%), and the indel from 75.78% to 78.44% (+2.66%). In a comparison
173 of 1:168x|2:64x to 1:124x|2:64x, the performance was similar, or even slightly dropped
174 from 97.51% to 97.40% overall. One possible reason is that the lower coverage test sample
175 cannot benefit from the much higher coverage used for model training. We propose how to
176 deal with excessively high coverage in test samples (i.e., coverage exceeding that used in
177 model training) in the Discussion section below.

178

179 **Does multiple subsampled coverage for model training improve variant calling**
180 **performance?** Yes. in a comparison of 1:44x|2:64x to '1:44x (single cov.)|2:64x', the latter
181 used only the full coverage 44-fold in model training; the overall F1-score increased from
182 95.47% to 96.84% (+1.37%), the SNP from 96.94% to 98.01% (+1.07%), and the indel from
183 75.78% to 78.44% (+2.86%). The results show that even without sufficient coverage for
184 model training, using multiple subsampled coverage still improved the variant calling
185 performance significantly.

186

187 **What is the upper bound on performance?**

188 To determine Clair's performance limit on the current ONT data, we intentionally overfitted
189 Clair by adding the samples we are going to test to the model training. That is, we
190 performed a biased test by exposing the test samples to model training, and if a true variant
191 is not called even after a biased training, it suggests the input signal is simply too weak
192 against the noise. Theoretically, a valid biased test requires the training data to be flawless
193 and the model design to be perfect, which are neither the case in our study, nor realistic in
194 real-world problems. However, in terms of training samples, GIAB has improved the quality
195 continuously. In the latest version v3.3.2, 99.5% of the variants have been correctly
196 phased¹³, suggesting their unprecedented quality. In terms of model design, we admit that
197 Clair will still have rooms to be improved, but the improvement is likely to be insignificant if
198 we do not deviate from using pileup data because, in this study, we have systematically
199 optimized the method using the techniques laid out in the Methods section. That is to say,
200 as we expect the real performance cap will be higher, the gap between it and using a biased
201 test in our study is small. Thus, it is appropriate to use a biased test to explore the limit of
202 Clair.

203

204 The two tests we did were 1:168x+2:64x|2:64x and 1:168x+2:64x|1:168x. Although the test
205 sample coverage in the first test was much lower than that in the second (64-fold against
206 168-fold), their performance was similar, with the overall F1-score at 97.77% and 97.82%,
207 SNP at 98.75% and 98.77%, and indel at 79.92% and 81.37%. The biased test
208 1:168x+2:64x|2:64x did not significantly outperform 1:168x|2:64x; the overall F1-score
209 increased from 97.40% to 97.77% (+0.33%), SNP from 98.53% to 98.75% (+0.22%), and indel

210 from 77.34% to 79.92% (+2.58%). Even with this biased experiment, we observed that the
211 performance of using Clair on the current ONT data was capped at about 97.8% F1-score
212 overall, 98.8% on SNPs, and 80% on indels. We consider how the new ONT chemistry that
213 provides a lower base error rate can raise the upper bound of Clair's variant calling
214 performance in the Discussion section below.

215

216 We analyzed and categorized the FP and FN results of Clair on ONT data. We randomly
217 extracted 100 FPs and 100 FNs from the 1:168x|2:64x experiment. **Figure 3** shows a
218 summary and examples of different categories, and Supplementary Table 2 shows a detailed
219 analysis of each FP and FN. Within the 100 FPs, the three largest categories are "Incorrect
220 allele with $AF \geq 0.2$ " (41/100), "Homopolymer" (25/100), and "Tandem repeat" (11/100).
221 "Incorrect allele with $AF \geq 0.2$ " means that at the FP variant, an incorrect allele dominates
222 other alleles in the read alignments (including the correct one), and the incorrect allele has a
223 frequency $\geq 20\%$. "Homopolymer", "Tandem repeat", and "Low complexity region" mean
224 that the FP variant is in a repetitive region, which remains difficult for ONT base-calling. It is
225 worth mentioning that these repetitive regions are ≤ 10 bp because we removed all GA4GH
226 low-complexity regions longer than 10bp from benchmarking. It may not be possible to
227 perfectly resolve these three categories for FP variants using pileup data for variant calling,
228 although complete read alignments might help to provide better precision. Three out of 100
229 FPs had "Incorrect insertion bases", while two out of 100 were categorized as "Overlapping
230 insertions", which means that the alleles of two consecutive insertions overlapped each
231 other in an input tensor; thus, the correct allele cannot be resolved for both insertions.
232 These two categories of errors can be resolved using the '--pysam_for_all_indel' option in
233 Clair, but this slows down Clair for ONT data by a factor of up to ten times. Other errors,

234 including "Incorrect indel length" and "Incorrect zygosity", are errors made by Clair's neural
235 network. In the 100 FNs, the three major categories are "Correct allele with AF<0.25"
236 (54/100), "Homopolymer" (18/100), and "Tandem repeat" (7/100). "Correct allele with
237 AF<0.25" means that at the location of the missed (FN) variant, the signal of the correct
238 allele is rather weak, with allele frequency lower than 25%. One FN categorized as "More
239 than two possible alternative alleles" is an error due to an alignment error in segmental
240 duplications, in which more than two alternative alleles seem correct.

241

242 [Performance on PacBio CCS](#)

243 In early 2019¹⁷, PacBio developed a protocol based on single-molecule, circular consensus
244 sequencing (CCS) to generate highly accurate (99.8%) long reads averaging as much as
245 13.5kb. PacBio published CCS datasets for HG001 (in this section also referred to as 1:30x; 1
246 as the sample suffix and 30x means 30-fold coverage), HG002 (2:33x) and HG005 (5:33x). All
247 three samples are involved in model training. To demonstrate a possible overfitting
248 phenomenon on deep learning based variant callers, both HG002 and HG005 are used in
249 benchmarking. To align with ONT's benchmarking results, the low-complexity regions
250 defined by GA4GH were removed from benchmarking. But noteworthy, PacBio CCS data is
251 less erroneous and the performance in the low-complexity regions are not significantly
252 degraded. The performance differences before and after removing the low complexity
253 regions using different sequencing technologies are available in "Methods –
254 Benchmarking".

255

256 Supplementary Table 3 shows the results of Clair and three other variant callers:
257 Clairvoyante, Longshot, and DeepVariant. Testing on HG002, DeepVariant performed the

258 best, with an overall F1-score of 99.96%, SNP of 99.97%, and indel of 99.92%. The primary
259 result of Clair 1:30x+5:33x|2:33x had an overall F1-score of 99.83%, which was 0.13% lower
260 than DeepVariant, but outperformed both Clairvoyante and Longshot. On SNP,
261 1:30x+5:33x|2:33x had an F1-score of 99.88%, which was 0.09% lower than DeepVariant,
262 0.43% higher than Longshot, and 0.17% higher than Clairvoyante. On indel,
263 1:30x+5:33x|2:33x had an F1-score at 99.07%, which was 0.85% lower than DeepVariant,
264 but 19.17% higher than Clairvoyante, showing that the new methods applied to Clair have
265 effective solved the indel-calling problem in Clairvoyante. In terms of speed, Clair (147
266 minutes) is slightly faster than Longshot (206 minutes), and about seven times faster than
267 DeepVariant (1,072 minutes). We also tested HG005. Interestingly, while Clair, Clairvoyante,
268 and Longshot all performed better on HG005 than HG002, DeepVariant performed worse.
269 Comparing 1:30x|2:33x to 1:30x|5:33x, Clair's overall F1-score increased from 99.77% to
270 99.80%. Clairvoyante's overall F1-score increased from 98.61% to 98.70%. Longshot's SNP
271 F1-score increased from 99.45% to 99.46%. The performance of the three callers verifies the
272 quality of the HG005 dataset. However, DeepVariant's F1-score dropped from 99.96% to
273 99.92%, the SNP F1-score decreased from 99.97% to 99.93%, and the indel F1-score
274 dropped most significantly, from 99.92% to 99.78%. The most probable reason is that,
275 DeepVariant's current PacBio CCS model was trained completely using HG002²⁰. We suggest
276 using DeepVariant's result on HG005 as its real performance on PacBio CCS data. The biased
277 test 1:30x+2:33x+5:33x|2:33x found the performance cap of Clair at 99.88% on SNP, which
278 was the same as 1:30x+5:33x|2:33x, and 99.28% on indel, which was 0.21% higher than
279 1:30x+5:33x|2:33x. While in 1:30x+5:33x|2:33x, the highest coverage used for model
280 training was only 33x, we expect to fill the performance gap on indel calling by using higher
281 coverage for model training. The performance gap between Clair and DeepVariant on

282 HG005 (99.28% against 99.78%, -0.5%) is the result of Clair using pileup data, while
283 DeepVariant uses complete read alignments that contain information at a per-read level.
284 This is also a reason DeepVariant runs slower than Clair. We discuss the possibility of
285 improving Clair to use complete read alignments without slowing down performance in the
286 Discussion section below.

287

288 [Performance on Illumina](#)

289 Approximately 300x coverage in 148-bp Illumina paired-end read data is available for five
290 GIAB samples, including HG001, HG002, HG003, HG004 and HG005¹¹. We used HG001,
291 HG003, HG004, HG005 for model training, and HG002 for benchmarking. To resemble the
292 typical coverage in whole genome sequencing, we used full coverage of HG001 (306-fold)
293 and HG005 (352-fold), but down-sampled HG002, HG003 and HG004 to 52-, 57-, and 66-
294 fold. To align with ONT's benchmarking results, the low-complexity regions defined by
295 GA4GH were removed from benchmarking. But we observed that Illumina's results were not
296 very much affected by removing the low-complexity regions.

297

298 Supplementary Table 4 shows the results of Clair and DeepVariant. DeepVariant performed
299 better, with an overall F1-score of 99.94%. The primary result of Clair
300 1:306x+3:57x+4:66x+5:352x|2:52x was an overall F1-score of 99.83%, which was 0.11%
301 lower than DeepVariant's. For SNPs, the F1-score of Clair was 0.09% lower than that of
302 DeepVariant (99.85% versus 99.94%). For Indel, the F1-score of Clair was 0.42% lower than
303 DeepVariant's (99.48% versus 99.90%). In terms of speed, Clair was about seven times faster
304 than DeepVariant (77 versus 537 minutes). The biased test
305 1:306x+2:52x+3:57x+4:66x+5:352x|2:52x found the performance cap of Clair to be 99.87%

306 for SNPs, which was 0.02% higher than the primary result, but 0.07% lower than that of
307 DeepVariant, and 99.57% for indels, which was 0.09% higher than the primary result, but
308 0.33% lower than that of DeepVariant. Similar to the ONT and PacBio CCS experiments, we
309 expect to fill in the performance gap through partially making use of complete read
310 alignments, as discussed in the Discussion section.

311 Discussion

312 In this paper we present Clair, a germline small variant caller for single molecule sequencing
313 data. The name Clair means 'clear' in French, echoing its predecessor, named Clairvoyante,
314 meaning 'clear seeing'. Clair adds new methods to solve problems that Clairvoyante had
315 trouble with, including multiallelic variant calling and long indel calling. In our experiments
316 on ONT data, Clair outperformed all existing tools in terms of precision, recall and speed. On
317 PacBio CCS and Illumina data, Clair performed slightly worse than DeepVariant, but ran
318 about an order of magnitude faster. Looking closer at the FP and FN variants shows that
319 Clair is approaching the limit on how accurately it can call variants using pileup data. Some
320 of the erroneous variant calls can be corrected using complete read alignments instead of
321 pileup data. However, dealing with complete read alignments requires a more powerful
322 neural network design with much greater computational demands. In the future, we will
323 explore using an ensemble method to handle the majority of the variants using Clair, while
324 for the extremely tricky ones we will use a new, more sophisticated method.

325

326 The quality and sufficiency of training data is key to the performance of Clair, as well as
327 other deep learning based variant callers, such as DeepVariant. To train a model for
328 production purposes, we used five samples (HG001 to 5) for Illumina data, but only two

329 samples (HG001 and HG002) for ONT, due to the limited availability of public high-coverage
330 whole genome sequencing datasets for the GIAB samples. ONT sequencing of the other
331 GIAB samples is ongoing, and more data will be available in the near future. With additional
332 datasets, we expect to see even higher performance in Clair on ONT data.

333

334 On ONT data, although Clair performed the best, its indel calling precision and recall were
335 only about 80%, even excluding GA4GH low-complexity regions, which leaves substantial
336 room for improvement. While the precision can be further improved by considering
337 complete read alignments, the recall is bounded by input and can be improved only with a
338 lower read-level base-calling error rate. Future improvements in ONT technology offer the
339 possibility of reducing the error rate to 2-3%, which in turn should improve Clair's ability to
340 detect indels in these data.

341

342 The GIAB datasets we used for model training have moderate whole-genome sequencing
343 coverage. Although we can use samples with very high coverage (over 300-fold, which is
344 sometimes seen in amplicon sequenced data) with Clair for variant calling, such samples
345 might show degraded performance because very high coverage variants were not
346 adequately observed in model training. To solve this problem, we propose two methods.

347 One method is to do transfer learning using a trained model on additional datasets with
348 very high coverage. Clair supports transfer learning and can be applied to additional
349 datasets instantly. Another method is an ensemble method, which generates multiple
350 copies of randomly subsampled read alignments at a candidate variant for Clair to call
351 variant. A majority vote or a decision tree can be used to make the final decision, using the
352 results of each copy.

353

354 A limitation of Clair is that it cannot be applied to polyploid species, which are inconsistent
355 with its neural network design. For the same reasons, Clair is not applicable to somatic
356 variant calling, where a single sample might hold multiple distinct populations of cells. Our
357 next steps include extending Clair to support polyploid species and somatic variant calling.

358 Method

359

360 Clair's input/output

361 Input

362 For a truth variant for training or a candidate variant for calling, the read alignments that
363 overlap or are adjacent to the variant are summarized (i.e. pile-up data) into a three-
364 dimensional tensor of shape 33 by 8 by 4, comprising 1056 integer numbers. The three
365 dimensions correspond to the position, the count of four possible bases from two different
366 strands, and four different ways of counting. In the first dimension, 33 positions include the
367 starting position of a variant at the center and 16 flanking bases on both sides. The second
368 dimension corresponds to the count of 'A+', 'A-', 'C+', 'C-', 'G+', 'G-', 'T+' or 'T-', with the
369 symbols +/- denoting the count from the forward/reverse strand. The third dimension
370 replicates the first two dimensions with four different ways of counting to highlight 1) the
371 allelic count of the reference allele, 2) insertions, 3) deletions and 4) single nucleotide
372 alternative alleles. "Supplementary Note – Pseudocode for generating the input tensor"
373 shows the pseudo code of the exact algorithm of how the input tensor is generated.
374 Supplementary Figure 1 demonstrates how the tensors look like for ONT data at a random
375 'non-variant', a 'SNP', an 'Insertion', and a 'Deletion'.

376

377 Output

378 The output of Clair has four tasks (a.k.a. four output components, in total 90 probabilities),
379 including 1) the 21-genotype probabilistic model (21 probabilities); 2) zygosity (3
380 probabilities); 3) the length of the first indel allele (33 probabilities); and 4) the length of the
381 second indel allele (33 probabilities). One of the breakthroughs in Clair is the invention of
382 the 21-genotype probabilistic model. It comprises all of the possible genotypes of a diploid
383 sample at a genome position, including 'AA', 'AC', 'AG', 'AT', 'CC', 'CG', 'CT', 'GG', 'GT', 'TT',
384 'AI', 'CI', 'GI', 'TI', 'AD', 'CD', 'GD', 'TD', 'II', 'DD', and 'ID', where 'A', 'C', 'G', 'T', 'I' (insertion)
385 and 'D' (deletion) denote the six possible alleles. The new model covers variants with two
386 alternative alleles, which could not be called in Clairvoyante. The zygosity task outputs the
387 probability of the input being 1) a homozygous reference (0/0); 2) heterozygous with 1 or 2
388 alternative alleles (0/1 or 1/2); or 3) a homozygous variant (1/1). The zygosity task is
389 partially redundant to the 21-genotype task, but it makes decisions independently, and it
390 crosschecks the decision made by the 21-genotype task. Tasks three and four have the same
391 design. They output the length of up to two indel alleles. Each task outputs 33 probabilities,
392 including the likelihood of 1) more than 15bp deleted (<-15bp); 2) any number between -
393 15bp and 15bp, including 0bp, and; 3) more than 15bp inserted (>15bp). In training, the
394 indel allele with a smaller number is set as the first indel allele. For example, for a
395 heterozygous 1bp deletion, the first indel allele is set as -1bp, the second as 0bp (-1bp/0bp).
396 For a heterozygous 1bp insertion, 0bp/1bp is set. This design makes the non-0bp training
397 variants for both tasks balanced. For a heterozygous indel with two alternative alleles, say,
398 one -2bp and one 5bp, -2bp/5bp are set. For a homozygous indel, two indel alleles are set to
399 the same value. For indels longer than 15bp, the exact length is determined using an

400 additional step (Method – New methods used in Clair – Dealing with indels longer than
401 15bp). The output of the two indel allele tasks are also used for crosschecking with the 21-
402 genotype task, with 0bp supporting an SNP allele, and non-0bp supporting an indel allele.
403 More details about how the four tasks crosscheck each other to come up with a result
404 coherently are in "Method – New methods used in Clair – Determining the most probable
405 variant type using the four tasks of Clair".

406

407 [New methods used in Clair](#)

408 Clair has been fully revamped while a few basic deep-learning techniques in Clairvoyante
409 have been retained, including 1) model initialization; 2) activation function; 3) optimizer; 4)
410 dropout; 5) L2 regularization; and 6) combining multiple samples for model training. The
411 parameters in 1, 2, and 3 remain default. The dropout rate of each layer is depicted in
412 Figure 1. For L2 regularization, we tested and found a constant 1e-3 worked best with
413 cyclical learning rate, which will be introduced in a following section. Below we discussed
414 the new methods we have applied in Clair.

415

416 [Moving from convolutional to recurrent neural network](#)

417 In Clairvoyante, we observed that the size of the convolutional kernel in the three
418 convolutional layers had limited the performance. Increasing the kernel size will increase
419 the performance, at a cost that the number of model parameters and running time will
420 increase exponentially. When designing Clair, we tried two strategies, including 1) multiple
421 dilated kernels, and 2) recurrent neural network (RNN) with bi-directional LSTM, for
422 substituting the three convolutional layers in Clairvoyante. The performance of using dilated
423 kernels was good. But to achieve the same performance as using RNN, six three times n

424 dilated kernels are needed for our input for each convolutional layer, which increased the
425 computation significantly. An RNN layer is usually slower than a convolutional layer with the
426 same number parameters, but it outperformed dilated kernel because to achieve the same
427 performance, RNN with bi-directional LSTM requires over 50% fewer parameters in our new
428 design. We used two RNN layers of in Clair in contrast to three convolutional layers in
429 Clairvoyante.

430

431 [Dealing with indels longer than 15bp](#)

432 For each candidate variant, Clair directly outputs the length of up to two alternative indel
433 alleles. However, if an insertion goes beyond 15bp, or a deletion goes below -15bp, Clair
434 runs an additional step to decide its exact length and allele. In the additional step, Clair
435 gathers all possible insertion/deletion alleles longer than 15bp at a genome position
436 through pysam (a wrapper around htslib and the samtools²¹ package). Depending on the
437 genotype concluded by Clair, we choose 1) the insertion/deletion with the highest allelic
438 count for 'AI', 'CI', 'GI', 'TI', 'AD', 'CD', 'GD' and 'TD'; 2) the insertions with the highest and/or
439 the second-highest allelic count for 'II'; 3) the deletions with highest and/or the second-
440 highest allelic count for 'DD', or; 4) both the insertion and deletion with the highest allelic
441 count for 'ID'. The additional step is slow, but it is required only for indels longer than 15bp.
442 We investigated HG001 and found 570,367 indels in its truth variant set; only 10,672
443 (1.87%) were >15bp. In our experiments, we found the slowdown was acceptable. Users can
444 set an option in Clair to enable this additional step for all indels, but our experiments found
445 that while the improvement in precision is small, it slows down Clair by about two times
446 with Illumina and PacBio CCS data, and by more than 10 times on ONT data.

447

448 Determining the most probable variant type using the four Clair tasks
449 Clair outputs data on four tasks. With an independent penultimate layer (Figure 1, FC5
450 layer) immediately before each task, the output of each task is considered independent. We
451 made two observations from our experiments: 1) for true positive variants, a random task
452 or two will make a mistake occasionally, but usually, the best and the second-best
453 probabilities are near and can be disambiguated if considered with other tasks; 2) for false
454 positive variants, the tasks do not usually agree well with each other, leading to two or
455 more possible decisions with similar probabilities. Thus, in Clair, we implemented a method
456 as a submodule for making a decision using the output of all four tasks. Variants are divided
457 into 10 categories: 1) a homozygous reference allele; 2) a homozygous 1 SNP allele; 3) a
458 heterozygous 1 SNP allele, or heterozygous 2 SNP alleles; 4) a homozygous 1 insertion allele;
459 5) a heterozygous 1 insertion allele, or heterozygous 1 SNP and 1 insertion alleles; 6)
460 heterozygous 2 insertion alleles; 7) a homozygous 1 deletion allele; 8) a heterozygous 1
461 deletion allele, or heterozygous 1 SNP and 1 deletion alleles; 9) heterozygous 2 deletion
462 alleles; and 10) a heterozygous 1 insertion and 1 deletion alleles. The likelihood value of the
463 10 categories is calculated for each candidate variant, and the category with the largest
464 likelihood value is chosen (Pseudocode in "Supplementary Note – Pseudo code for
465 determining the most probable variant type"). The variant quality is calculated as the square
466 of the Phred score of the distance between the largest and the second-largest likelihood
467 values.

468

469 Cyclical learning rate

470 The "initial learning rate" and "how the learning rate decays" are two critical
471 hyperparameters in training a deep neural network model. A model might be stuck at a local

472 optimum (i.e. unable to achieve the best precision and recall) if the initial learning rate is
473 too large, or the decay is too fast. But a large initial learning rate, and a slow decay rate
474 make the training process either unstable or take too long to finish. So in common practice,
475 a tediously long grid search that is very costly is needed to find the best hyperparameters.
476 Furthermore, through a grid search, we found that different sequencing technologies differ
477 in their best hyperparameters. This problem makes model training too complicated and
478 largely impedes Clair from being applied to new datasets and sequencing technologies. To
479 solve the problem, we implemented Cyclical Learning Rate (CLR)²² in Clair. CLR is a new deep
480 learning technique that eliminates the need to find the best values of the two
481 hyperparameters. CLR gives a way to schedule the learning rate in an efficient way during
482 training, by cyclically varying between a lower and higher threshold. Following the CLR
483 paper, we determined the higher threshold to be 0.03 and the lower threshold to be 0.0001.
484 The two thresholds worked well on the training variants of all three sequencing
485 technologies (Illumina, PacBio CCS and ONT). In terms of which CLR scheduler to use, we
486 chose the triangular schedule with exponential decay. In our experiments, on PacBio CCS
487 and Illumina datasets, CLR decreased model training time by about 1–3 times, while often
488 outperforming the three-step decay method introduced in Clairvoyante for both precision
489 and recall. However, on ONT datasets, CLR has a lower, but almost negligible, performance
490 than the three-step decay. We provide both CLR and three-step decay options in Clair. To
491 train a model for production, we suggest users try both options and choose the best
492 through benchmarking. In our results, we used CLR for PacBio CCS and Illumina datasets,
493 and the three-step decay method for ONT datasets.
494

495 Focal loss

496 Our training data uses the truth variants from the GIAB consortium and is unbalanced in
497 terms of variant type. For example, the number of heterozygous variants is nearly twice that
498 of the homozygous variants. SNPs are about five times more numerous than indels. Worst
499 of all, only ~1.1% (39,898 of 3,619,471 in HG001) of variants have two or more alternative
500 alleles. And among them, only 884 (~0.024%) are multiallelic SNPs. This problem leads to
501 degenerate models, as the numerous easy variants contribute no useful learning signals and
502 overwhelm training. In our practice, if we leave the problem unaddressed, we observe a
503 significant drop in recall for the underrepresented variant types. For multiallelic SNPs, the
504 recall dropped to zero. To solve this problem, we used the "Focal loss" technique²³, which
505 applies a modulating term to the cross-entropy loss in Clair's output to focus training on
506 underrepresented hard variants and down-weight the numerous easy variants. Focal loss
507 calculates the loss as $(1 - p_t)^\gamma \times \alpha_t \times -\log(p_t)$, where $p_t = p$, $\alpha_t = \alpha$, if the prediction
508 matches the truth, or $p_t = (1 - p)$, $\alpha_t = (1 - \alpha)$ otherwise. In addition to the traditional
509 cross entropy loss, focal loss uses two more parameters: γ (the focusing parameter) to
510 differentiate easy/hard training examples, and α (the balancing parameter) to balance the
511 importance of positive/negative training examples. We have tested the combinations of $\gamma =$
512 1, 2, 4, and $\alpha = 0.25, 0.5$. We determined $\gamma = 2$ and $\alpha = 0.25$ work best for the GIAB truth
513 variants with a 1:2 ratio of truth variant and non-variant. The use of focal loss significantly
514 increases the performance of underrepresented variant types. It also allows us to be more
515 lenient on variant type balance when augmenting the training data.

516

517 Training data augmentation using subsampled coverage

518 Lower coverage usually leads to lower precision and recall in variant calling. To train Clair to
519 achieve better performance on variants with lower coverages, we subsampled each dataset
520 into four or nine additional datasets with lower coverages. The subsampling factors f are
521 determined as $(\sqrt[h]{4 \div c})^n$, where c is full coverage of each sample, 4 is the minimal
522 coverage, h is either 4 or 9, and n is from 1 to h . Using HG002 as an example, its full
523 coverage is 63.68-fold, and the nine subsampled coverages are 46.82-, 34.43-, 25.31-,
524 18.61-, 13.69-, 10.06-, 7.40-, 5.44- and 4.00-fold. If variant samples were lower than 4x after
525 subsampling, we removed them from training. We used the command "samtools view -s f "
526 to generate a subsampled BAM. A different seed counting from zero for random number
527 generation was set for each coverage. The use of subsampled coverages improved the recall
528 on indel significantly.

529

530 Methods tested but showed no improvement to accuracy

531 In this section we discuss methods we tested that had no effect on Clair's performance. For
532 researchers working on further improving the performance of Clair, these methods could be
533 avoided or revised. This section is elaborated in detail in the Supplementary Notes.

534

535 Benchmarking

536 The GIAB truth variant datasets

537 We used the GIAB version 3.3.2 datasets as our truth variants. Depending on the availability
538 of deep sequencing data, our ONT experiments used samples HG001 or HG001+HG002 for
539 model training, our PacBio CCS experiments used HG001 or HG001+HG005, and our Illumina

540 experiments used HG001 or HG001+HG003+HG004+HG005. For benchmarking, ONT, PacBio
541 CCS and Illumina experiments have used HG002, HG005, and HG002, respectively. The links
542 to the truth variants and high-confidence regions are available in “Methods – Data sources –
543 Truth variants”. Depending on the reference genome used in the already available read
544 alignments, we used GRCh38 for our ONT and Illumina experiments, and GRCh37 for our
545 PacBio CCS experiments. The links to the reference genomes we used are available in
546 “Methods – Data sources – Reference genomes”

547

548 [Removing GA4GH low-complexity regions from benchmarking](#)

549 Krusche et al.⁶ from the GA4GH benchmarking team and the GIAB consortium published the
550 low-complexity regions, including homopolymers, STRs, VNTRs, and other repetitive
551 sequences for stratifying variants in their paper titled "Best practices for benchmarking
552 germline small-variant calls in human genomes". All low-complexity regions defined by
553 GA4GH are longer than 10bp. The performance difference between before and after
554 removing the low-complexity regions are in Supplementary Table 5. ONT's performance
555 degraded significantly (precision -11.41%, recall -55.33%), while that of PacBio CCS and
556 Illumina dropped only 0.99–1.67% in precision and recall. Thus, when computing variant
557 calling using ONT, we suggest removing the variants called in the low-complexity regions. In
558 our benchmarks for all datasets, in addition to using the high-confidence regions of each
559 sample provided by GIAB, we removed the low-complexity regions. The procedures are
560 available in "Supplementary Note – Commands – Remove GA4GH low complexity regions
561 from GIAB's high-confidence regions". There was retention of 92.61–93.47% high-
562 confidence regions in GRCh38, and 94.40–95.05% in GRCh37 of the five samples HG001 to 5
563 after removing the low-complexity regions (Supplementary Table 8).

564

565 [Benchmarking methods and metrics](#)

566 Clair trains a model either for 30 epochs, using the Cyclical Learning Rate (used for PacBio
567 CCS and Illumina datasets), or by decaying the learning rate three times (by one tenth each
568 time) until the validation losses converge (used for ONT datasets). While the performance of
569 last few epochs are generally similar, the best-performing one will be chosen for
570 benchmarking. We did not run replications of model training because choosing from the
571 best epoch actually resembles the process of having multiple replications. In ONT and
572 Illumina experiments, the GRCh38 reference genome was used, while in PacBio CCS
573 experiments, GRCh37 was used. For each variant calling experiment, we used the
574 submodule vcfEval in RTG Tools²⁴ version 3.9 to generate three metrics, 'Precision', 'Recall',
575 and 'F1-score', for five categories of variants: 'Overall', 'SNP', 'Indel', 'Insertion', and
576 'Deletion'. All time consumptions were gauged on two 12-core Intel Xeon Silver 4116 (in
577 total 24 cores), with 12 concurrent Clair processes, each with 4 Tensorflow threads. As Clair
578 has some serial steps that use only one thread, we observed our setting sufficient to
579 maximize the utilization of all 24 cores. For other variant callers, including DeepVariant,
580 Longshot and Medaka, options were to set to use all 24 cores for the best speed.

581

582 [Computational performance](#)

583 Clair requires Python3, Pypy3 and Tensorflow. Variant calling using Clair requires only a
584 CPU. For a typical 30-fold human WGS sample, Clair takes about an hour for Illumina data
585 and PacBio CCS data, and five hours on ONT data, using two 12-core Intel Xeon Silver 4116
586 processors. Memory consumption depends on both input data and concurrency. ONT data
587 has a higher memory footprint than Illumina and PacBio CSS, while Clair is capped at 7GB

588 per process (helper scripts at 4.5GB and Tensorflow at 2.5GB). Model training requires a
589 high-end GPU; we used the Nvidia Titan RTX 24GB in our experiment. Using Clair’s default
590 parameters, generating 1 million training samples takes about 38 seconds. For example, the
591 Illumina model with four samples (HG001, 3, 4, 5) and 30 coverages in total (10 for 1 and 5,
592 5 for 2 and 3) has 284,367,735 training samples and takes about 11,000 seconds per epoch.
593 Training the 1:168x|2:64x ONT model used 246,099s (2.85 days) on the Titan RTX. In
594 comparison, the Nvidia RTX 2080 Ti 11GB is about 15% slower, and the Nvidia GTX 1080 Ti
595 11GB is about 35% slower.

596

597 [Code availability](#)

598 Clair is open source, available at <https://github.com/HKU-BAL/Clair>. Clair is licensed under
599 the BSD 3-Clause license.

600

601 [Data availability](#)

602 The details of and links to the 1) reference genomes, 2) truth variants, 3) Oxford Nanopore
603 Technologies (ONT) data, 4) Pacific Bioscience (PacBio) CCS data, and 5) Illumina data that
604 are supporting the findings of this study are available in the “Data Sources” section in the
605 Supplementary Notes. The VCF files generated by Clair in this study are available at
606 http://www.bio8.cs.hku.hk/clair_models/VCFBenchmarked/.

607 [Acknowledgements](#)

608 We thank Steven Salzberg, Mike Schatz, and Fritz Sedlazeck for their valuable comments. R.
609 L. was supported by the ECS (Grant No. 27204518) of the HKSAR government, and the URC

610 fund at HKU. T. L., C. W., Y. W., C. T., C. Li. and C. Le. were supported by the ITF (Grant No.
611 ITF/331/17FP) from the Innovation and Technology Commission, HKSAR government.

612 [Author contributions](#)

613 R. L. and T. L. conceived the study. R. L, C. W., Y. W., C. T., C. Li. and C. Le. analyzed the data
614 and wrote the paper.

615 [Competing interests](#)

616 The authors declare no competing interests

617

618 References

- 619 1 Goodwin, S., McPherson, J. D. & McCombie, W. R. Coming of age: ten years of next-
620 generation sequencing technologies. *Nat Rev Genet* **17**, 333-351,
621 doi:10.1038/nrg.2016.49 (2016).
- 622 2 Ashley, E. A. Towards precision medicine. *Nat Rev Genet* **17**, 507-522,
623 doi:10.1038/nrg.2016.86 (2016).
- 624 3 Li, H. Toward better understanding of artifacts in variant calling from high-coverage
625 samples. *Bioinformatics* **30**, 2843-2851, doi:10.1093/bioinformatics/btu356 (2014).
- 626 4 Luo, R., Schatz, M. C. & Salzberg, S. L. 16GT: a fast and sensitive variant caller using a
627 16-genotype probabilistic model. *GigaScience* (2017).
- 628 5 Van der Auwera, G. A. *et al.* From FastQ data to high confidence variant calls: the
629 Genome Analysis Toolkit best practices pipeline. *Curr Protoc Bioinformatics* **43**, 11 10
630 11-33, doi:10.1002/0471250953.bi1110s43 (2013).
- 631 6 Krusche, P. *et al.* Best practices for benchmarking germline small-variant calls in
632 human genomes. *Nat Biotechnol* **37**, 555-560, doi:10.1038/s41587-019-0054-x
633 (2019).
- 634 7 The long view on sequencing. *Nat Biotechnol* **36**, 287, doi:10.1038/nbt.4125 (2018).
- 635 8 Sedlazeck, F. J., Lee, H., Darby, C. A. & Schatz, M. C. Piercing the dark matter:
636 bioinformatics of long-range sequencing and mapping. *Nat Rev Genet*,
637 doi:10.1038/s41576-018-0003-4 (2018).
- 638 9 Ameer, A., Kloosterman, W. P. & Hestand, M. S. Single-Molecule Sequencing:
639 Towards Clinical Applications. *Trends Biotechnol* **37**, 72-85,
640 doi:10.1016/j.tibtech.2018.07.013 (2019).
- 641 10 Luo, R., Sedlazeck, F. J., Lam, T. W. & Schatz, M. C. A multi-task convolutional deep
642 neural network for variant calling in single molecule sequencing. *Nat Commun* **10**,
643 998, doi:10.1038/s41467-019-09025-z (2019).
- 644 11 Zook, J. M. *et al.* Extensive sequencing of seven human genomes to characterize
645 benchmark reference materials. *Sci Data* **3**, 160025, doi:10.1038/sdata.2016.25
646 (2016).
- 647 12 Zook, J. M. *et al.* Integrating human sequence data sets provides a resource of
648 benchmark SNP and indel genotype calls. *Nat Biotechnol* **32**, 246-251,
649 doi:10.1038/nbt.2835 (2014).
- 650 13 Zook, J. M. *et al.* An open resource for accurately benchmarking small variant and
651 reference calls. *Nature Biotechnology* **37**, 561-566, doi:10.1038/s41587-019-0074-6
652 (2019).
- 653 14 Jain, M. *et al.* Nanopore sequencing and assembly of a human genome with ultra-
654 long reads. *Nat Biotechnol* **36**, 338-345, doi:10.1038/nbt.4060 (2018).
- 655 15 Edge, P. & Bansal, V. Longshot enables accurate variant calling in diploid genomes
656 from single-molecule long read sequencing. *Nat Commun* **10**, 4660,
657 doi:10.1038/s41467-019-12493-y (2019).
- 658 16 *medaka: Sequence correction provided by ONT Research.*
659 <https://github.com/nanoporetech/medaka>, accessed Nov 17 2019.
- 660 17 Wenger, A. M. *et al.* Accurate circular consensus long-read sequencing improves
661 variant detection and assembly of a human genome. *Nat Biotechnol* **37**, 1155-1162,
662 doi:10.1038/s41587-019-0217-9 (2019).

663 18 Poplin, R. *et al.* A universal SNP and small-indel variant caller using deep neural
664 networks. *Nature Biotechnology* **2018/09/24/online**, doi:10.1038/nbt.4235 (2018).

665 19 Simpson, J. T. *et al.* Detecting DNA cytosine methylation using nanopore sequencing.
666 *Nature methods* **14**, 407 (2017).

667 20 Poplin, R. *et al.* *DeepVariant training data*.
668 [https://github.com/google/deepvariant/blob/r0.9/docs/deepvariant-details-](https://github.com/google/deepvariant/blob/r0.9/docs/deepvariant-details-training-data.md)
669 [training-data.md](https://github.com/google/deepvariant/blob/r0.9/docs/deepvariant-details-training-data.md), accessed Nov 22 2019.

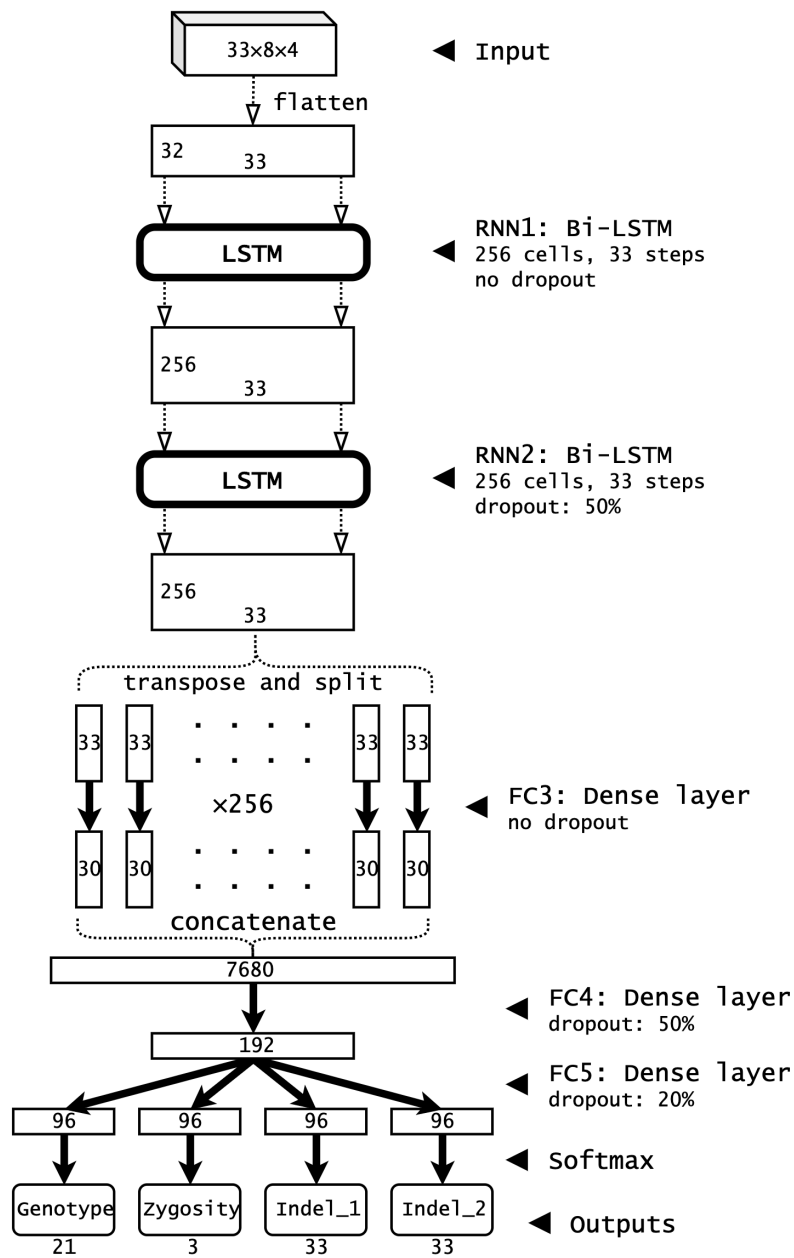
670 21 Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,
671 2078-2079, doi:10.1093/bioinformatics/btp352 (2009).

672 22 Smith, L. N. in *2017 IEEE Winter Conference on Applications of Computer Vision*
673 (WACV). 464-472 (IEEE).

674 23 Lin, T.-Y., Goyal, P., Girshick, R., He, K. & Dollár, P. in *Proceedings of the IEEE*
675 *international conference on computer vision*. 2980-2988.

676 24 Cleary, J. G. *et al.* Joint variant and de novo mutation identification on pedigrees
677 from high-throughput sequencing data. *J Comput Biol* **21**, 405-419,
678 doi:10.1089/cmb.2014.0029 (2014).

679

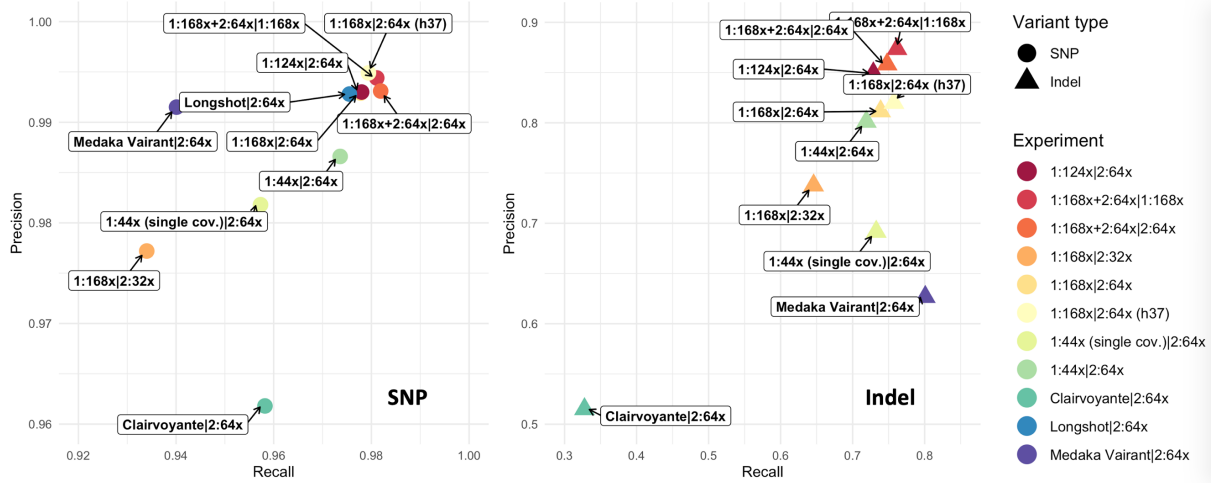


681

682 Figure 1. Clair network architecture and layer details. RNN: Recurrent Neural Network. FC:

683 Fully Connected layer. Bi-LSTM: Bi-directional Long Short-Term Memory layer.

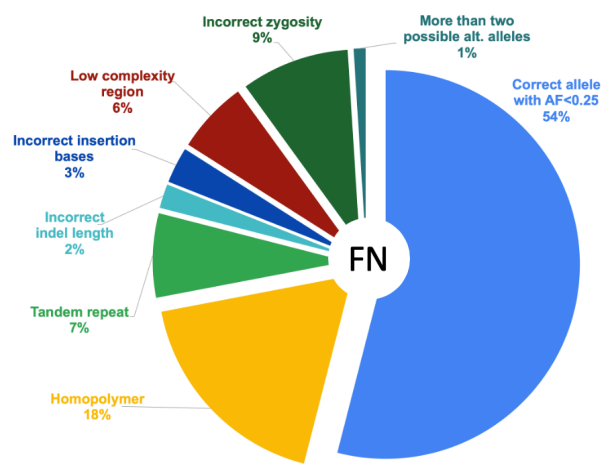
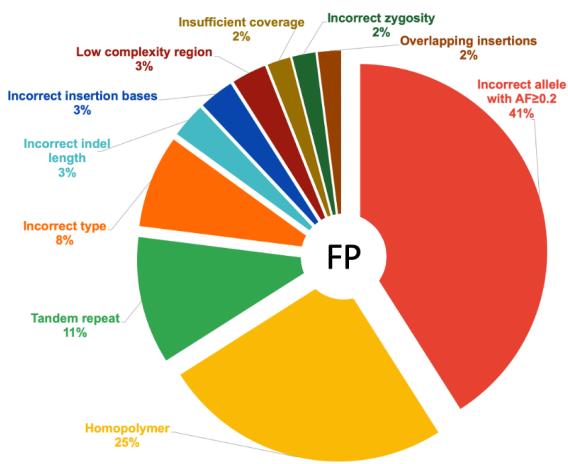
684



685

686 Figure 2. ONT benchmarking results. For Clair, the datasets used for model training and
 687 testing are separated with a vertical bar '|', and are written as 'a:bx', where a denotes the
 688 suffix of the GIAB sample ID (e.g., 1 means HG001), and b denotes the coverage of the
 689 dataset. Longshot calls only SNP variants, so it is not shown in the indel results.

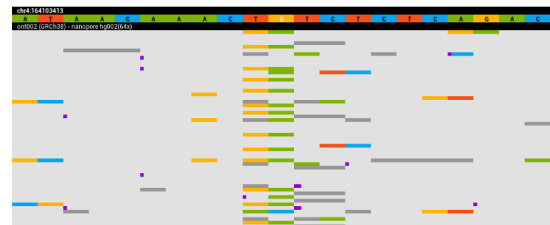
690



Homopolymer FP - chr2:207,157,429



Tandem repeat FN - chr4:164,103,413



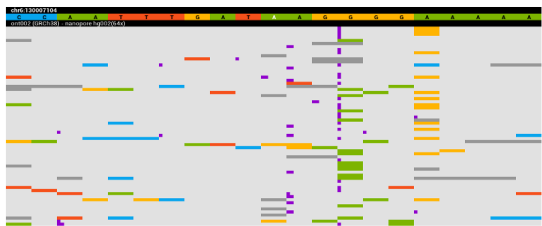
Insufficient coverage FP - chr2:93,219,242



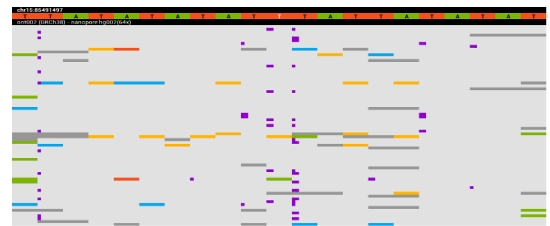
Low complexity region FN - chr1:119,620,624



More than two possible alt. alleles FN - chr6:130,007,104



Overlapping insertions FP - chr15:85,491,497



691

692 Figure 3. The category distribution of FPs and FNs made by Clair in the 1:168x|2:64x

693 experiment on ONT data, and six genome browser screen captures showing examples of

694 different categories. In the screen captures, bases A, C, G, and T are green, blue, yellow, and

695 red, respectively. Gaps (i.e., deletions) are dark gray. Insertions are purple dots between

696 two bases and are wider when the insertion is longer.