# DREAMS: Deep Read-level Error Model for Sequencing data applied to low-frequency variant calling and circulating tumor DNA detection

Mikkel H. Christensen [1,4] *, Simon Drue [1] *, Mads H. Rasmussen [1,4] *, Amanda Frydendahl [1,4] *, Iben Lyskjær [1,4], Christina Demuth [1], Jesper Nors [1,4], Kåre A. Gotschalck [2,4], Lene H. Iversen [3,4], Claus L. Andersen [1,4]# & Jakob Skou Pedersen [1,4]#

*Shared first author

#Shared senior authors / corresponding authors

[1]Department of Molecular Medicine, Aarhus University Hospital, Aarhus, Denmark.

[2]Department of Surgery, Horsens Regional Hospital, Horsens, Denmark.

[3]Department of Surgery, Aarhus University Hospital, Aarhus, Denmark.

[4]Department of Clinical Medicine, Faculty of Health, Aarhus University, Aarhus, Denmark

## Abstract

Circulating tumor DNA detection using Next-Generation Sequencing (NGS) data of plasma DNA is promising for cancer identification and characterization. However, the tumor signal in the blood is often low and difficult to distinguish from errors. We present DREAMS (**D**eep **Rea**d-level **M**odelling of **S**equencing-errors) for estimating error rates of individual read positions. Using DREAMS, we developed statistical methods for variant calling (DREAMS-*vc*) and cancer detection (DREAMS-*cc*).

For evaluation, we generated deep targeted NGS data of matching tumor and plasma DNA from 85 colorectal cancer patients. The DREAMS approach performed better than state-of-the-art methods for variant calling and cancer detection.

## Background

Degraded DNA fragments are released into the blood through apoptosis, necrosis and active secretion from a range of cell types and can be detected as circulating free DNA (cfDNA)[1]. Solid tumors also shed DNA into the bloodstream and cfDNA of cancer origin is called circulating tumor DNA (ctDNA)[2]. The ctDNA level in blood is reported to be positively associated with tumor burden[3, 4]. As the half-life of cfDNA is less than an hour, ctDNA measurements can be considered real-time assessments of tumor burden and studies have shown that ctDNA can be more sensitive than radiological imaging[5-7]. This makes ctDNA measurements a promising approach for detecting relapse in patients who have undergone curative surgery[6-10]. Other proposed applications include diagnosis and intervention planning, tracking therapeutic response, monitoring the development of treatment resistance, and ultimately early detection of cancer in screening programs[8, 11]. Since obtaining liquid biopsies, such as plasma from blood samples, is both cost-effective and minimally invasive, techniques for efficient ctDNA detection holds great promise for targeted treatment in precision medicine.

41    In clinical contexts with low tumor burden, e.g. detection of minimal residual disease after curative-

42    intended surgery and early detection of recurrence, the ctDNA constitute only a minor fraction of

43    the cfDNA, often less than 0.1%. Hence, the error rate of current sequencing methods is in the same

44    order of magnitude as the tumor signal[12], making it challenging to accurately distinguish errors

45    from true mutations in ctDNA applications. Errors can arise in several steps between the initial

46    shedding of cfDNA and the final generation of next-generation sequencing (NGS) reads (**Figure 1**).

47    DNA fragments may be damaged e.g. by deamination or oxidation[13, 14], during PCR amplification

48    of the sequencing library[13], and during sequencing from PCR amplification and/or sequencing

49    artefacts.{{Ma, 2019 #25}} For deep sequencing, some of the PCR and sequencing errors can be

50    rectified using unique molecular identifiers (UMIs). With the use of UMIs, each DNA fragment is

51    labeled with a unique "barcode" prior to PCR amplification, such that replicates of the same

52    fragment can be grouped together. Errors can then be eliminated by comparing the replicates within

53    a group, as errors from PCR amplification and sequencing are likely to be present in only a minority

54    of reads. However, some errors, such as DNA damage introduced prior to UMI labeling remains and

55    continue to challenge the discrimination of true low frequency mutational signal from these errors.

56    Several methods for detecting low frequency variants using NGS data have been developed. Most of

57    these establish a model for the expected frequency of errors and then assess the mutational signal

58    with a statistical test. They differ greatly in the required data prerequisites, how the errors are

59    modelled and handled, and the final assessment of the mutational signal.

60    Mutect2[15] and Shearwater[16] are examples of general somatic variant callers applicable for most

61    NGS data. Mutect2 realigns reads in regions with mutational signal and then calculates a log-odds for

62    the existence of the alternative allele using a statistical model in which the error rates are derived

63    from the PHRED scores. Shearwater is developed specifically for low-frequency somatic variant

64    detection for sub-clonal tumor mutations. It builds a position-specific error model based on the

65    observed rate of read alignment mismatches across a set of training samples. A mutation is called if

66    the observed signal exceeds what is expected from the error model. Additionally, this method can

67    incorporate prior knowledge about the probability of the mutations of interest.

68    Other methods, including MRDetect[17], INVAR [18]and iDES[12], have been specifically tailored to

69    detect ctDNA in NGS data. These methods build on the idea of aggregating the signal across multiple

70    mutations to classify a sample as ctDNA positive or negative, as opposed to calling each individual

71    mutation. For this purpose, a patient specific catalogue of mutations is generated from a matched

72    tumor sample. However, the enhanced performance of these methods come at the expense of

73    general applicability as they assume the presence of curated data from known ctDNA fragments or

74    specialized lab protocols.

75    Here we develop a generally applicable ctDNA detection method based on a detailed background

76    error model of individual read positions. This approach aims to capture general read-level error

77    behavior and thus be applicable even for genomic regions where training data is not available. Data

78    from reads known to come from ctDNA is not needed, and all data outside known mutated

79    positions, or from independent normal samples can be used as training data. However, training data

80    that was obtained similarly to the test data will provide the most precise model. Thus, severe

81    changes in laboratory protocols should optimally be accompanied by re-training of the model. Some

82    features such as the read position[19], proximity to fragment ends[14], UMI group size[12], GC-

83    content[20] and trinucleotide context[21] have been shown to affect the probability of errors at

84    individual read positions. By modelling their effect, the error rate of individual read positions may be

85    predicted. Thereby, a read alignment mismatch, i.e. a non-reference base, with a low predicted error

86    rate can provide more mutational evidence than a mismatch with a high error rate.  This allows for

87    improved cfDNA error modelling, which is key to develop accurate ctDNA applications.

88    In the following, we demonstrate how cfDNA errors can be modelled accurately using a neural

89    network, by combining read level features with information about the sequencing context. For this

90    we developed DREAMS (**Deep Rea**d-level **M**odelling of **S**equencing-errors) that incorporates both

91    read-level and local sequence-context features for positional error rate estimation. Based on

92    DREAMS, we developed a method for variant calling (DREAMS-*vc*) to accurately call individual cancer

93    mutations in cfDNA data. The method was generalized for cancer calling in DREAMS-*cc* that

94    aggregates the signal across a catalogue of mutations for accurate estimation of the tumor fraction

95    and sensitive determination of the overall cancer status. To evaluate the performance of DREAMS,

96    we performed deep-targeted sequencing of pre- and post-operative cfDNA samples from 85 stage I-

97    II colorectal cancer (CRC) patients and compared to state-of-the art methods Mutect2[15] and

98    Shearwater[16].

## Results

100    Plasma cfDNA was extracted from pre-operative (Pre-OP) and post-operative (Post-OP) blood draws

101    of 85 stage I-II CRC patients (**Table 1**) undergoing curative surgery. In addition, two stage III CRC

102    patients were used in the model training. A biopsy from the resected tumor and paired peripheral

103    blood cells was sequenced to generate a patient-specific mutational catalogue. Post-OP samples

104    were collected 2-4 weeks after surgical removal of the primary tumor (**Figure 2**). Each cfDNA sample

105    was sequenced using a custom hybrid-capture panel, designed to capture 41 exonic regions,

106    spanning 15.413 bp, frequently mutated in CRC (**Supplementary section 1** and **Supplementary table**

107    **1**). After UMI collapse the median of the average depths with corresponding interquartile range

108    (IQR) of samples were for Pre-OP; 3307 (IQR: 3560), Post-OP; 7143 (IQR: 8844), buffycoat; 1850 (IQR:

109    1468), and tumor samples; 2132 (IQR: 2145), no samples had an average read depth below 100. All

110    samples have been mapped and processed through the same pipeline (**Supplementary section 1).**

111

112

113

114

115

116

**Table 1:** Clinical characteristics

| Characteristic | Count or Median (percent or range) |
|---|---|
| Patients | 85 (100%) |
| Gender | |
| *Male* | 53 (62%) |
| *Female* | 32 (38%) |
| Age [years] | 71 (49-87) |
| Tumor location | |
| *Right colon* | 23 (27%) |
| *Left colon* | 26 (31%) |
| *Rectum* | 36 (42%) |
| Pathological T-stage | |
| *pT1* | 15 (18%) |
| *pT2* | 25 (29%) |
| *pT3* | 41 (48%) |
| *pT4* | 4 (4.7%) |
| UICC stage | |
| *I* | 40 (47%) |
| *II* | 45 (53%) |

117

118  We first identified features that are known or expected to affect the error rate (**Figure 3a**). In

119  general, they can be split into two types: local sequence-context features and read-level features.

120  The local sequence-context features capture the genomic sequence context, including the

121  trinucleotide context, information about the sequence complexity (Shannon entropy of nucleotide

122  frequency), and GC contents in an 11 bp window around the position of interest (**Methods**).

123    The read-level features capture the structural composition of the read, UMI characteristics and

124    sequencing information. The structural composition includes the strand a read aligns to (forward or

125    reverse), the number of insertions and deletions in the read, and the total size of the underlying

126    fragment. In the read pre-processing, UMIs were used to generate consensus reads with lowered

127    error rates (**Supplementary section 2**). For each consensus read, we extracted the UMI-group size,

128    the number of reads disagreeing with the consensus at the position, and the overall number of

129    mismatches outside the position of interest. As sequencing related features, we included the base

130    position in the read (read position) and whether the read is the first to be sequenced from the read-

131    pair. The read quality (PHRED score) was not included, as it had the same high value for all positions

132    in the UMI-collapsed consensus reads.

133    We evaluated the individual features association with the error rate by analyzing the total set of

134    read alignment mismatches (n=707,562) across all Post-OP samples (**Figure 3b-d**), after excluding

135    mutations and variants found in matching tumor and germline samples. The mismatches were

136    compared to an equal number of randomly sampled matches, to estimate the error rate for each

137    feature across its values (**Supplementary section 3**).

138    Since fragment lengths of cfDNA are influenced by nucleosome binding patterns, the fragment

139    length distribution have peaks at around 162 bp (mono-nucleosomal) and 340 bp (di-

140    nucleosomal)[22]. The error rate tended to be minimized in fragments of these lengths (**Figure 3b**).

141    As expected, we observed a lower error rate in consensus reads formed by larger UMI groups[12]

142    (**Figure 3c**).

143    The error distribution for the read position showed an increased error rate in the beginning of the

144    reads (**Figure 3d**).  We also observed a clear difference in error distribution along the read between

145    the first and second read of the pair. The 12 different nucleotide alterations showed widely different

146    error rates (**Figure 3e**), which is expected as error-induced mismatches are not equally likely, and the

147 rate further differed between the two strands. However, strand symmetric alterations were

148 generally similar, apart from the mismatches C→T/G→A and C→A/G→T.

149 Overall, we saw variation in the error rate for all the presented features (the remaining are shown in

150 **Supplementary section 3**). Thus, for a given genomic position, different reads may have different

151 error rates due to differences in read-level features. In the following, we present how this variation

152 can be captured and used to potentially improve detection of ctDNA.

153 **Neural network model and feature selection**

154 To predict the error rate at a given read position, we used a neural network model with the input

155 features described above (**Methods**). The predictive ability of individual features was evaluated

156 using a "leave-one-covariate-out" (LOCO) scheme[23] (**Supplementary section 4**). In short, we

157 evaluated the performance of a full model containing all features (baseline) and then the relative

158 performances of restricted models where each feature had been left out one by one. We used the

159 latter to measure and rank the importance of each feature (**Figure 4a**). When leaving out the

160 trinucleotide context, the reference base was provided instead to assess only the importance of the

161 two neighboring nucleotides.

162 We found the most informative feature for modelling the error rate to be the strand (**Figure 4a)**. The

163 second and thirds most informative features were whether the read is the first in a pair and the read

164 position. The trinucleotide context was fourth, indicating that there is a difference in error rate for

165 different contexts, as found by others[18]. The fragment length and the UMI group size also

166 contribute significantly to the model. The remaining features showed little to no effect on the model

167 performance.

168 An optimal subset of informative features was chosen using a stepwise procedure where features

169 were excluded in order of importance (**Methods**). The set of features chosen was the smallest model

170 that did not perform significantly worse than the full model (**Supplementary section 4)**. The four

171    least important features could be removed without any significant negative effect on the

172    performance (**Figure 4b**). Of the remaining ten features, eight were read-level features, namely the

173    features describing the UMI group, the number of errors in the UMI group, the number of deletions

174    in the read, the number of other errors in the read, the fragment length, read position, strand, and if

175    the read was first in pair. This showed that read-level features do contribute to accurate modelling

176    of the error rate.

177    The numerical and categorical variables are processed differently in the neural network prior to the

178    hidden layers (**Figure 4c**). The numerical features are batch normalized, the categorical features are

179    one-hot encoded, and the tri-nucleotide context is embedded in three dimensions to handle the

180    large number of possible contexts (**Methods**).

181     **Predictive performance in clinical data**

182    To validate the utilization of the DREAMS error model, we applied it in calling tumor variants

183    (DREAMS-*vc*) and cancer (DREAMS-*cc*) (**Methods)**. We assessed the performance using five repeats

184    of 2-fold cross-validation (5x2 CV) (**Figure 5a**). The model was trained on the Post-OP samples, and

185    Pre-OP samples were used for method validation. The split was done on patient level to ensure that

186    a model is not trained and tested on data from the same patient. This analysis was repeated with

187    five different randomized splits to control for split induced variation.

188    The performance of calling tumor mutations in the plasma samples was assessed by looking at the

189    area under Receiver Operating Characteristic curves (AUC). The performance of DREAMS-*vc* was

190    compared to state-of-the-art algorithms Mutect2 and Shearwater. Only positions with at least one

191    observed mismatch were included in the performance calculations (**Figure 5b**). Positions without

192    signal was called negative by any method, making them redundant for performance comparisons.

193    Using DREAMS-*vc*, we aimed to call the tumor mutations of each patient from their respective

194    mutation catalogue. As negative controls, we attempted to call cross-patient mutations, by

195    searching for the mutations found in other patients. Additionally, a validation set of 500 randomly

196    generated alterations within the covered sequencing panel was used as negative controls. Evaluating

197    across the combined negative set of both cross-patient mutations and validation alterations and

198    cancer stages, DREAMS-*vc* performs significantly better than both Shearwater and Mutect2 (**Figure**

199    **5b**). Additionally, the performance was assessed separately for stage I and stage II CRC patients. This

200    showed that superior performance of DREAMS-*vc* is predominantly due to the stage II CRC patients

201    (**Figure 5b**).  As expected, all models perform better on later stage patient samples as these are

202    expected to have a higher mutational signal in the cfDNA due to a higher tumor burden.

203    All methods perform similarly on stage I patients, however DREAMS-*vc* has marginally better

204    performance. Performance evaluations for each of the separate negative sets showed that DREAMS

205    performs better than Mutect2 with the cross-patient negative set and better than Shearwater with

206    the validation set as the negative set. The variation in performance of DREAMS-*vc* across splits and

207    folds is lower than for Mutect2 and Shearwater, which indicates that its variant calling is more stable

208    across patients and mutation types.

209    By maintaining the false positive rate at 5% for the alterations with signal in the validation set for

210    each model, we get comparable thresholds for the three confidence measures: p-values, Bayes

211    factor and TLOD for DREAMS-*vc*, Shearwater, and Mutect2, respectively. This allows for a

212    comparison of the sensitivity of the models at a pre-determined specificity of 95%. The model could

213    then be assessed across an alteration catalogue of 191 true positive mutations from the mutation

214    catalogue and 1290 cross-patient negative calls based on the mutation catalogue of the other

215    patients. Out of the alteration catalogue, 88 true mutations and 1100 cross-patient negative calls

216    had a signal for the alteration.

217    Using this threshold DREAMS-*vc* called 83% of the tumor mutations with signal, while Shearwater

218    and Mutect2 called 75% and 72.7%, respectively (Table 2). F1 and G-mean scores were calculated to

219    assess the performance of the models by using the cross-patient mutations as negative controls. G-

220    mean is the geometric mean of sensitivity and specificity, and F1 is the harmonic mean of precision

221    and sensitivity. For G-mean, DREAMS-*vc* performed better than Shearwater and Mutect2, however

222    the F1 score of Shearwater was very similar to DREAMS-*vc*, due to lower false-positive rate of

223    shearwater (Table 2). Considering all mutations observed in the tumors, including those without

224    signal in plasma, we found that about 38.2% could be recalled in Pre-OP liquid biopsy samples.

225

226

| Table 2 | Full alteration catalog[a] | | Catalogue alterations with signal[b] | | | |
|---|---|---|---|---|---|---|
| | Sensitivity | Specificity | Sensitivity | Specificity | F1 | G-mean |
| **DREAMS-*vc*** | 0.382 | 0.998 | 0.830 | 0.957 | 0.702 | **0.891** |
| **Shearwater** | 0.346 | 0.998 | 0.750 | 0.971 | **0.710** | 0.853 |
| **Mutect2** | 0.336 | 0.997 | 0.727 | 0.933 | 0.566 | 0.831 |

227

228    [a] Full alteration catalogue consisting of n=191 true positive mutations, and n=1290 potential cross-patient

229    negative calls.

230    [b] Catalogue of alterations with signal consisting of n=88 true positive mutations, and n=1100 potential cross-

231    patient negative calls.

232    By setting the threshold based on a 5% false positive rate in the cross-patient mutation set, the

233    validation mutation set can be used as negative controls. The true positives are still the same 191

234    mutations of which 88 has a signal for the alteration. The negatives are the 500 validation positions

235    multiplied with the 87 tested samples, giving a total of 43,500 possible alterations of which 1,350

236    had a signal. With this set we obtained an 83% true positive rate, compared to 77.3% for Shearwater

237    and 68.2% for Mutect2 (Table 3). DREAMS-*vc* scored highest in both F1 and G-mean scores. Here,

238    DREAMS-*vc* performed distinctly better than Shearwater, while Mutect2 had a more comparable F1

239    score.

| Table 3 | Full alteration catalog[a] | Catalogue alterations with signal[b] | |
|---|---|---|---|

| | Sensitivity | Specificity | Sensitivity | Specificity | F1 | G-mean |
|---|---|---|---|---|---|---|
| **DREAMS-*vc*** | 0.382 | 0.998 | 0.830 | 0.944 | **0.616** | **0.885** |
| **Shearwater** | 0.356 | 0.997 | 0.773 | 0.911 | 0.493 | 0.839 |
| **Mutect2** | 0.314 | 0.999 | 0.682 | 0.962 | 0.603 | 0.810 |

240

241  [a] Whole catalogue consisting of n=191 true positive mutations, and n=43500 potential validation set calls.

242  [b] Catalogue of positions with signal consisting of n=88 true positive mutations, and n=1350 potential validation

243  calls.

244  A common measure used to predict the presence of ctDNA is the estimated tumor fraction in

245  plasma. DREAMS-*cc* combines the mutational evidence across the mutation catalogue, to estimate

246  the tumor fraction with an accompanying p-value for the presence of cancer (**Methods**). We aimed

247  to detect cancer in the Pre-OP samples, since cancer is present and should, in theory, be detectable

248  given enough ctDNA is present in the blood. As a negative control, we attempted to detect cancer in

249  each Pre-OP sample (Tested Sample) with the mutation catalogue from all other patients (Candidate

250  patient) (**Figure 6a**). In case of shared mutations between the mutation catalogues, these were

251  eliminated to prevent false positives. As a benchmark, we constructed a cancer call score using the

252  product of the individual Bayes factors across the mutation catalogue from Shearwater, resulting in

253  a similar tendency (**Figure 6b**). The performance of calling cancer can be assessed by treating the

254  cross-patient mutation catalogues as expected negatives and calculate an AUC score. Performance

255  was compared using the 5x2 cross validation setup as above (**Figure 5a**). The AUC was very similar

256  between DREAMS-*cc* and Shearwater with respect to calling cancer, however DREAMS-*cc* showed a

257  slightly increased performance (p = 0.0343, one tailed t-test). As for variant calling, we only included

258  the samples with mutational signal to showcase and compare the performance of the different

259  methods in discriminating tumor from error signal.

260  For the patients with stage I and II CRC, we found tumor supporting reads in 47.5% (19/40) and 73%

261  (33/45) of the Pre-OP samples, respectively. We called cancer in 34% of the stage I CRC patients,

262  corresponding to 74% (14/19) of the patients with a mutational signal. We called cancer in 73% of

12

263   the stage II CRC patients, corresponding to 94% (31/33) of the patients with signal. These results

264   were obtained whilst still limiting the false positive rate to 5 % in cross-patient cancer calls with a

265   non-zero mutational signal.

266   Detailed analysis of the false positive cancer calls reveals that most are due to a specific KRAS G12V

267   variant: chr12:25245350 C>A. This variant is common in colon cancer, and it is therefore not

268   surprising to find in the patients [24]. However, the mutation was not found in the patient's

269   corresponding tumor or buffycoat samples. A possible explanation for this is that the mutation is not

270   detected in the tumor biopsy due to sub-clonality [25] or that there is an underlying germline signal

271   that was not caught in the buffycoat.

## 272   Discussion

273   We have developed DREAMS, as a new approach for modelling the error rates in sequencing data

274   that incorporates information from both the local sequence context and read-level information.

275   DREAMS is intended for settings that rely on accurate error identification and quantification. We

276   applied the error model for low-frequency ctDNA variant calling (DREAMS-*vc*) and cancer detection

277   (DREAMS-*cc*).

278   The error rate was found to vary depending on several of the proposed read-level features.

279   Surprisingly, fragment size was found to be correlated with the error rate, with the smallest error-

280   rates being observed for fragment sizes corresponding to the mono-nucleosomal and di-

281   nucleosomal lengths (**Figure 3b**). Fragments that deviate from these in length may have been

282   degraded in the blood for a longer time and thereby accumulated more errors. Fragments of ctDNA

283   are generally shorter and error rates are generally highest in short fragments, which shows the

284   importance of accurate error modelling[26, 27]. The error rate was also found to vary with the

285   strand, and symmetric mismatches occurred at different rates (**Figure 3e**). The G>T/C>A asymmetry

286   can be explained by the hybridization capture protocol only targeting one strand and thus only

287   capturing oxidative damage of that strand[14]. A similar mechanism might explain the C>T/G>A

288    asymmetry in the case of cytosine deamination. The error rate varied with the position in the read

289    and was especially increased in the beginning of reads (**Figure 3d**). This may be because ends of

290    fragments are prone to damage[14] and in thermodynamic equilibrium with being single stranded.

291    The error rate also varied depending on whether the read was the first or second in the pair (**Figure**

292    **3d**). Besides being intermitted by a PCR amplification step, the reads differ in composition and

293    length of adapters sequenced prior to the insert, which might cause this difference.

294    Training a background error model using DREAMS does not require known mutation sites in reads,

295    as it only models the errors found in aligned reads (BAM-files). These can originate from normal

296    samples or mutation filtered cancer samples, as in this study. Since error patterns are highly

297    dependent on laboratory procedures, the same protocol should be used for training samples and

298    subsequent testing samples. Training across multiple samples gathered over time, is expected to

299    learn the error patterns that are general across samples and batches. Conversely, if the amount of

300    data in a single sample is large, the error model can be trained on the sample itself, which

301    potentially yields a highly specific model that accounts for sample specific error patterns. The model

302    is built to be position agnostic and can therefore be used to predict error rates for positions for

303    which no training data is available. Furthermore, it is fit for both deep sequencing of panels and

304    shallow sequencing of whole genomes.

305    The error model has been implemented using a neural network, allowing the feature set to be

306    tailored to capture the relevant information of a specific setting. Analysis of the feature importance

307    revealed that several of the proposed read-level features are useful in predicting the error rate in

308    sequencing data (**Figure 4a**). Most features presented in this paper are general to NGS data,

309    however not all sequencing protocols use UMI based error correction, rendering UMI related

310    features redundant. In particular, UMI cannot be exploited for shallow whole-genome sequencing as

311    read groups cannot be formed. In such cases error rates would be increased, making accurate error

312    modelling as performed by DREAMS even more important.

313   Compared to simpler methods, the presented approach is more computationally demanding, due to

314   training of the neural network model and the use of complex data extracted from BAM-files. A

315   neural network is a simple and flexible approach for bridging the gap between a complex set of

316   contexts and read level features and the error rate of a given read position but might not be the

317   most efficient solution. The model can be trained on a regular laptop within a few hours, which

318   should only be done once, when the training dataset is defined. Using the trained model and the

319   statistical modules adds no significant computation time for calling mutations and cancer in the

320   current setting. However, very large mutation catalogues are expected to increase the computation

321   time for DREAMS-*cc*.

322   DREAMS was built to exploit read-level features under the assumption that these affect the error

323   rate in sequencing data. Thus, the power of this approach increases with the variability in the error

324   rate explained by read level features. Thereby, less emphasize is put on mismatches that are likely

325   errors, and more confidence in the potential tumor signal from other mismatches. Conversely, if

326   read level features are not improving error prediction, the performance is expected to be similar to

327   methods working with simpler summary data. Although DREAMS use information about the local

328   sequence-context, strong regional effects on the error rate are not expected to be captured by the

329   model.

330   In all performance comparisons DREAMS-*vc* performed better or equal to the other methods in

331   calling tumor mutations. This indicates that read-position level features can improve performance in

332   separating error from mutational signal. Similarly for cancer detection, DREAMS-*cc* performed equal

333   to calls based on Shearwater. Cancer was detected in most (73%) of stage II CRC cancer patients and

334   a third (34%) of stage I patients.

335   There are false positive cancer and mutation calls, some of which could potentially be explained by

336   clonal hematopoiesis of indeterminate potential (CHIP) or an unexpected error signal. To reduce the

337   signal from CHIP, we have excluded positions with significant presence of non-reference nucleotides,

338    found in the germline samples, however, a low signal might still be present. Remaining false positive

339    calls might be due to regional effects or sample specific artifacts. Many of the false positive mutation

340    calls in the Pre-OP samples were found to be a mutation leading to the KRAS G12V variant, and it

341    could therefore potentially be explained by a sub-clonal variant that was not identified in the tumor

342    sample or a germline signal of clonal hematopoiesis of indeterminate potential (CHIP) that was not

343    identified in the buffycoat samples.

344    Sensitive variant calling in liquid biopsies can provide non-invasive insight into tumor genetics, which

345    can potentially enable personalized treatment of patients and be a cost-effective approach for

346    cancer screening. DREAMS-*cc* integrates evidence across a mutation catalogue to increase sensitivity

347    in cancer detection. Cancer detection is expected to get more sensitive as the number of mutations

348    in the catalogue rises. A potential application of DREAMS-*cc* could be tumor agnostic cancer

349    detection based on a catalogue of commonly known tumor variants.

350    The approach presented in this paper does not utilize tumor specific signals such as the fragment

351    size distribution, fragmentation patterns, mutational signatures, expression information, etc.

352    However, together with the error characterizing properties of DREAMS-*cc*, this could potentially

353    refine the cancer calls. Addition of regional properties and positional information could potentially

354    further increase sensitivity. In this paper, we focus on the single nucleotide variants in the tumor,

355    but the model could be extended to be able to look for indels. The underlying ideas in DREAMS are

356    not restricted to variant calling and could be used in other tasks of sequencing data analysis such as

357    advanced error filtering.

## Conclusion

358

359    We have presented the DREAMS error rate model and demonstrated the importance of using read-

360    level features for modelling the errors in NGS data. The model was validated in a tumor informed

361    setting, using DREAMS-*vc* for variant calling and DREAMS-*cc* for cancer detection in patients with

362    CRC. DREAMS-*vc* allowed accurate detection of mutation signal in plasma samples extracted prior to

363     curative intended surgery with an improved performance compared to state-of-the-art methods.

364     This highlights the importance of including read-level information in modelling the background error

365     rate. Furthermore, DREAMS-*cc* demonstrated the ability to combine signal from multiple mutations

366     known from the tumor biopsy for improved cancer detection. DREAMS-*cc* was able to call cancer in

367     73 % of Pre-OP samples from CRC stage II patients, and 34 % of CRC stage I patients. Potential future

368     applications of DREAMS include analysis of WGS data and tumor agnostic cancer detection. The

369     approach presented with DREAMS is generally applicable across NGS applications that need accurate

370     handling and quantifications of errors, and the presented algorithms (DREAMS-*vc* and DREAMS-*cc*)

371     are only examples of how to exploit this. The specific application presented in this paper is

372     implemented as a user-friendly R package [https://github.com/JakobSkouPedersenLab/dreams].

## 373 Declarations

### 374 Ethics approval and consent to participate

375     The Committees on Biomedical Research Ethics in the Central Region of Denmark have approved the

376     study (J. No. 1-10-72-3-18). The study was performed in accordance with the Declaration of Helsinki

377     and all participants provided written informed consent.

### 378 Consent for publication

379     Not applicable.

### 380 Availability of data and materials

381     Sharing of sensitive patient specific clinical information and raw sequencing data is currently not

382     possible due to ethical and GDPR regulations.

### 383 Competing interests

384     The authors declare that they have no competing interests

### 385 Funding:

393  **Author contributions:**

394  MHC, SD, CLA, and JSP conceived and designed the study. MHC and SD developed the statistical

395  methods and the software under supervision by JSP with input from MHR and CLA. MHR, AF, IL, CD,

396  JN, KAG, and LHI acquired patient samples and generated patient data, including NGS data. MHC, SD,

397  MHR, AF, CLA, and JSP analyzed and interpreted the patient data. SD and MHC wrote the article

398  under supervision of CLA and JSP with revisions and suggestions from the other authors. All authors

399  read and approved the final manuscript.

# Methods

## Error rate prediction using read level information

In this study we present a method called DREAMS (**D**eep **Rea**d-level **M**odelling of **S**equencing-errors) for estimating the error rate at each read position using features of the individual read and the genomic context of the position. In practice, this is achieved by predicting the probability of observing each allele given the describing features of a position in a read and considering the probabilities of observing the alternative alleles as the error rates. The read specific features can include information such as the read position, the strand of the mapped read, the length of the fragment, and UMI-group size. The read position refers to the cycle number at which the position was sequenced starting with the first nucleotide of the fragment, thus disregarding cycles used for reading primers, adapters, unaligned ends etc. Context specific features contain information about the genomic sequence surrounding the position, including the neighboring bases (tri-nucleotide context), the complexity, and GC-content. The local complexity is calculated as the Shannon entropy for both single nucleotides and pairs. Similarly, the local GC content is calculated as the fraction of C and G nucleotides. In principle, any feature that can be thought to affect the error rate of a read position can be added to improve the error rate prediction. Another possible feature would be the positional read quality score given by the sequencing machine. However, the estimated quality for the collapsed consensus reads were all capped at the same high value and thus excluded as they do not include any information for further modelling.

## Data

Data for a read position can be extracted from a read mapping (BAM-file) with sequencing data from a next generation sequencing experiment. The training data for the model consists of a set of read positions from multiple samples, for which the observed allele is denoted together with the relevant features. This means that the training data includes both matches, where read positions where the observed allele is equal to the reference allele and mismatches where the observed and reference

19

433    allele differ. Mismatches that correspond to known single nucleotide polymorphisms found in the

434    germline samples are excluded from the training. Assuming that the training samples are non-

435    cancerous means that all remaining mismatches in the dataset can be assumed to be errors that

436    have occurred on a molecular level in the body or lab, or during sequencing of the sample.

437    The mismatches are extracted from the BAM-file using the mismatched positions annotated in the

438    MD-tag. The equivalent genomic position is found, and the 11- and 3-mer context is extracted from

439    the reference genome and used for calculation of local sequence-context features. The UMI errors

440    and UMI count are extracted from the cE and cD tags generated by the

441    CallMolecularConsensusReads from fgbio used for calling UMI consensus reads. Information about

442    the insertions and deletions is extracted from the cigar tag. The fragment size is the insert-size

443    (isize), and the read position is the position in the read sequence from the 5'-end of the read. Strand

444    and first in pair are extracted from BAM flag where this information is encoded in a bitwise fashion.

445    The model assumes that the input data for both training and testing is based on readings of unique

446    fragments, so each position in a fragment is only represented in one read. This can be assured using

447    unique molecular identifiers (UMIs) and by trimming overlapping read positions in the read pairs.

448    As training on every single read position in every single read is very demanding and inefficient, we

449    employ a methodology akin to importance sampling where we extract all the mismatches from the

450    data and randomly sample a subset of the non-mismatches. To account for this skew induced by

451    down-sampling one category of the training data a rescaling scheme inspired by[28] is used on the

452    predicted error rates. The method outlined in **Supplementary section 5**.

453    ## Neural network model

454    *Structure of the neural network*

455    To predict the error rate at a given read position we use a multilayer perceptron (MLP) which is a

456    simple neural network setup with multiple fully connected layers. The neural network allows us to

457  use the features without prior knowledge of how they interact amongst each other or how they

458  affect the error rate. The neural network is trained using a set of read positions where the features

459  describing the read positions are used as inputs and the observed allele as output.

460  For a given read position the possible observed outcomes are the alleles A, T, C or G. Interpreting

461  this as a random event, the observed allele can be seen as an outcome from a four-dimensional

462  multinomial distribution with one trial. Let $X_{ij}$ represents the observed allele in read $j$ at position $i$

463  and $D_{ij}$ be the set of observed features for that read position. For a non-mutated, homozygote

464  position the observed allele should predominantly be the reference allele, and any observations of

465  non-reference alleles, would be considered errors. In this situation $P(X_{ij} = A|D_{ij})$ would be close to

466  1 if $A$ was the reference allele for read position $(i, j)$, and $P(X_{ij} = x|D_{ij})$, $x \in \{T, C, G\}$ would be

467  the error rates for the remaining three alleles. Given a set of observations $\{(x_{ij}, D_{ij})\}_{i=1}^{N}$ it is then

468  possible to write the log-likelihood function for the observed data:

$$l\left(\left\{(x_{ij}, D_{ij})\right\}_{i,j}\right)$$

$$= \sum_{i,j} \log\left(P(X_{ij} = x_{ij}|D_{ij})\right)$$

$$= \sum_{i,j:x_{ij}=A} \log\left(P(X_{ij} = A|D_{ij})\right) + \sum_{i,j:x_{ij}=T} \log\left(P(X_{ij} = T|D_{ij})\right) +$$

$$\sum_{i,j:x_{ij}=C} \log\left(P(X_{ij} = C|D_{ij})\right) + \sum_{i,j:x_{ij}=G} \log\left(P(X_{ij} = G|D_{ij})\right)$$

469  The problem now becomes how to estimate the distribution $P(X_{ij}|D_{ij})$ above. To do this, start by

470  defining the probability functions via the SoftMax function:

$$P(X_{ij} = a|D_{ij}) = \frac{e^{f_a(D_{ij})}}{\sum_{a' \in \{A,T,C,G\}} e^{f_{a'}(D_{ij})}}$$

471  , where $f_a(D_{ij})$ is a predictor function for the allele $a$ using the observed information $D_{ij}$. As an

472  example, for classic multinomial logistic regression a linear predictor function is chosen such that

473 $f_a(X_i) = \beta_a \cdot X_i$, where $\beta_a$ is a vector of feature specific weights that can be found by maximizing

474 the log-likelihood function. To get a more flexible model, a neural network is chosen, since this can

475 approximate any arbitrary predictor function well including arbitrary interactions between input

476 features. To do this $P(X_{ij} = a | D_{ij})$ can be interpreted as the output from a neural network model

477 where SoftMax is used as the last activation function and $f_a(D_{ij})$ is the output from the last hidden

478 layer. To train such a model inspiration is drawn from likelihood theory and the negative log-

479 likelihood function is chosen as the loss function to minimize.

## Architecture

481 The neural network model allows for high flexibility in the choice of features and requires very

482 limited prior knowledge about the effect of the features on the error rate. The neural network was

483 selected to be a MLP with an input layer, three hidden layers and an output layer. The dimension of

484 the input layer depends on the selected input features, the hidden layers have a configuration of

485 128, 64, and 32 nodes with a ReLu activation function, and the output layer contains 4 nodes with

486 SoftMax activation, as explained above, corresponding to probability of observing each of the 4

487 alleles. The configuration of hidden layers can be varied, depending on the input data and the

488 available computational resources. The models were training using the Keras library (2.3.0) in R,

489 which is an interface that builds in Tensorflow (2.6.0) [29].

## Feature handling / embedding

491 The features are split into numeric, categorical, and embedded variables and handled accordingly.

492 Categorical features are one-hot encoded, and the numeric features are batch normalized. The

493 trinucleotide context can be seen as the three distinct features: reference allele and the two

494 neighboring bases. These can be handled as categorical features with individual one-hot encoded 4-

495 dimensional inputs using 12 (3x4) input nodes in total. Alternatively, a 64-dimensional (4x4x4) one-

496 hot encoded input of the entire trinucleotide context (TNC) can be used.  We will employ another

497 alternative that takes the 64-dimensional feature in the input layer and embeds it into a continuous

498   3-dimensional vector before including it in the model alongside the remaining input features.

499   Thereby, the model can learn the relationship between the contexts, and cluster contexts that have

500   a similar effect on the error rate close together and vice versa.

## Assessing cancer status across a catalogue of multiple mutation candidates

502   Based on the neural network error model developed above, it can now be assumed that the

503   individual error rates for a given position in each read is known. In this section the error rates will be

504   exploited to develop a statistical framework for estimating the tumor fraction in a sample based on a

505   catalogue of candidate mutations. This framework can ignore some mutation candidates if these are

506   not found in the sample, for example due to relatively low allelic frequency due to sub-clonality in

507   the tumor or due to little tumor in the circulation. Reduction in the candidate mutations allows for a

508   comprehensive mutation catalogue to be used, where mutation candidates with limited evidence

509   may be excluded. The subset of candidate mutations is selected statistically by finding mutations

510   with a consistently high mutational signal, and the tumor fraction is estimated based on these

511   candidates. This subset of mutations is then used in a statistical procedure for testing if the observed

512   mutational signal exceeds what we would expect if no mutated DNA were present, making it

513   possible to determine the cancer status of a patient based on the sample.

## The statistical model

515   Start by introducing $Z_i$ as a variable that controls the presence of a given mutation on the site $i$, such

516   that $Z_i = 1$ represent the case where the site is mutated, and $Z_i = 0$ when it is not. Furthermore

517   let:

$$Z_i \sim Bernulli(r)$$

518   Thus, given a catalogue of possible mutations, $r$ is the probability that each of them is present in the

519   sample. For site $i$ let $R$ be the germline reference allele and $M$ the alternative allele of interest.

520   Furthermore, it is assumed that the germline site is homozygote, such that any signal from non-

521   reference alleles must be due to errors or mutational signal from a tumor. To model the molecular

23

522     composition of the fragments covering site $i$ let $Y_{ij} \in \{R, M\}$ be the true error-free nucleotide of the

523     $j$'th fragment. If the $i$'th mutation is not present in the sample $(Z_i = 0)$, we are sure that the true

524     nucleotide of the fragment is the reference and thus the following distribution holds:

$$P\big(Y_{ij} = R | Z_i = 0\big) = 1, \qquad P\big(Y_{ij} = M | Z_i = 0\big) = 0$$

525     To model the mutational DNA present in the sample let $f > 0$ denote the tumor fraction. This is the

526     fraction of the DNA in the blood that originates from tumor cells. Assuming that the mutation of

527     interest is (sufficiently) clonal in the tumor, i.e. half of the DNA in the tumor has this mutation, the

528     probability of a given fragment having the mutation is $f/2$. Using this the following distribution for

529     $Y_{ij}$ can be assumed when the mutation is present in the sample $(Z_i = 1)$:

$$P\big(Y_{ij} = R | Z_i = 1\big) = 1 - \frac{f}{2}, \qquad P\big(Y_{ij} = M | Z_i = 1\big) = \frac{f}{2}$$

530     To model the errors that occur in NGS data let $X_{ij}$ be the observed nucleotide in fragment $j$ at

531     position $i$. Assume that the distribution of $X_{ij}$ depends only on the corresponding true nucleotide

532     $Y_{ij}$, in the sense that the event $X_{ij} \neq Y_{ij}$ corresponds to the observation being an error. This

533     distribution is exactly what the neural network model described above aims to approximate using

534     the observed features $D_{ij}$. To simplify notation the dependence of $X_{ij}$ on $D_{ij}$ will be omitted from

535     notation in the following. Note that observations $X_{ij}$ outside $\{R, M\}$ will have little information

536     about the true nucleotide $Y_{ij}$. Furthermore, since the error rates generally are low, the difference

537     between including interactions between all four possible alleles and only the two allele of interest is

538     negligible. Thus, to simplify the following calculations, we assume that $X_{ij} \in \{R, M\}$. In practice this

539     means that all fragments, $j'$, for which $x_{ij'} \notin \{R, M\}$ are eliminated from the analysis. Using this

540     assumption, we define the probability of observing the alternative allele in a reference allele

541     position as the following error rate:

$$e_{ij}^{R \to M} = P\big(X_{ij} = M | Y_{ij} = R, X_{ij} \in \{R, M\}\big) = \frac{P\big(X_{ij} = M | Y_{ij} = R\big)}{P\big(X_{ij} = R | Y_{ij} = R\big) + P\big(X_{ij} = M | Y_{ij} = R\big)}$$

542    Conversely, for a fragment that stems from a tumor cell and contains the mutated allele we define:

$$e_{ij}^{M \to R} = \frac{P(X_{ij} = R | Y_{ij} = M)}{P(X_{ij} = R | Y_{ij} = M) + P(X_{ij} = M | Y_{ij} = M)}$$

543    Estimating the tumor fraction and mutation presence

544    In this section we will develop a procedure for estimating the tumor fraction ($f$) and mutation

545    presence probability ($r$). For this, let $i \in \{1, \dots, K\}$ be the index of a catalogue of $K$ candidate

546    mutations, $N_i$ the corresponding number of covering reads and $\left\{ (x_{ij})_{j \in \{1,\dots,N\}} \right\}_{i \in \{1,\dots,K\}}$ all the

547    observed alleles. First, we write the likelihood function for $f$ and $r$:

$$L\left(f, r \middle| \{(x_{ij})\}_{i \in \{1,\dots,K\}, j \in \{1,\dots,N\}}\right)$$

$$= \prod_{i=1}^{K} P(Z_i = 0) \cdot$$

$$\prod_{j:x_{ij}=R} \left[ P(X_{ij} = R | Y_{ij} = R) P(Y_{ij} = R | Z_{ij} = 0) + P(X_{ij} = R | Y_{ij} = M) P(Y_{ij} = M | Z_{ij} = 0) \right] \cdot$$

$$\prod_{j:x_{ij}=M} \left[ P(X_{ij} = M | Y_{ij} = R) P(Y_{ij} = R | Z_{ij} = 0) + P(X_{ij} = M | Y_{ij} = M) P(Y_{ij} = M | Z_{ij} = 0) \right] +$$

$$P(Z_i = 1) \cdot$$

$$\prod_{j:x_{ij}=R} \left[ P(X_{ij} = R | Y_{ij} = R) P(Y_{ij} = R | Z_{ij} = 1) + P(X_{ij} = R | Y_{ij} = M) P(Y_{ij} = M | Z_{ij} = 1) \right] \cdot$$

$$\prod_{j:x_{ij}=M} \left[ P(X_{ij} = M | Y_{ij} = R) P(Y_{ij} = R | Z_{ij} = 1) + P(X_{ij} = M | Y_{ij} = M) P(Y_{ij} = M | Z_{ij} = 1) \right]$$

$$= \prod_{i=1}^{K} (1-r) \cdot \prod_{j:x_{ij}=R} \left(1 - e_{ij}^{R \to M}\right) \cdot \prod_{j:x_{ij}=M} e_{ij}^{R \to M} +$$

$$r \cdot \prod_{j:x_{ij}=R} \left[ \left(1 - e_{ij}^{R \to M}\right) \cdot \left(1 - \frac{f}{2}\right) + e_{ij}^{M \to R} \cdot \frac{f}{2} \right] \cdot \prod_{j:x_{ij}=M} \left[ e_{ij}^{R \to M} \cdot \left(1 - \frac{f}{2}\right) + \left(1 - e_{ij}^{M \to R}\right) \cdot \frac{f}{2} \right]$$

548    Getting a maximum likelihood estimate (MLE) of $f$ and $r$ by optimizing this expression analytically is

549    not tractable. However, by seeing $Y_{ij}$ and $Z_i$ as latent variables, estimates can be found by

550    employing an EM-algorithm, which will be developed in a **Supplementary section 6**. For now,

551    assume that $\hat{f}$ and $\hat{r}$ are a MLEs of $f$ and $r$ respectively.

552    To test if a sample has a significant content of mutational DNA, we focus on the parameter in the

553    model. By representing the hypothesis of a negative sample as a tumor fraction of 0 and no

554    mutations present $(H_0: f, r = 0$ ) and a positive sample as a positive tumor fraction and some

555    mutations present $\left(H_A: f > 0, r \geq \frac{1}{K}\right)$, a likelihood ratio test can be used to test for significance.

556    Note that $r \geq \frac{1}{K}$ in $H_A$ corresponds to at least one mutation being present in the sample. The LR-test

557    statistic for this test is:

$$Q = -2 \log \frac{L\left(0,0 \Big| \{(x_{ij})\}_{j=1}^{N}\right)}{L\left(\hat{f}, \hat{r} \Big| \{(x_{ij})\}_{j=1}^{N}\right)}$$

558    Since there are 2 free parameters in the model, it can be assumed that $Q$ is approximately $\chi^2(2)$-

559    distributed, and a p-value can be obtained as follows:

$$p_{val} = 1 - F_{\chi^2(2)}(Q)$$

560    Using this statistical model for cancer calling on top of the error rate predictions from DREAMS we

561    refer to as the DREAMS-cc.

## Calling individual mutations

563    In the special case where the number of mutations in the catalogue is $K = 1$, the algorithm outlined

564    above can be thought of as a regular variant caller. In this case the concept of some mutations not

565    being present in the sample is unnecessary, as the presence of the single mutations of interest can

566    be governed solely by the tumor fraction $f$. The algorithm above is easily modified to handle this by

567    assuming that $r = 1$, and using one degree of freedom for the $\chi^2$-distribution in the significance

568    test. The equations in the EM-algorithm can also be simplified by making this assumption. We refer

569    to the variant caller will be referred to as DREAMS-*vc*.

26

570 # Figure legends

571 **Figure 1:**

572 Error generation in Next Generation Sequencing data. Normal cells (grey) and tumor cells (blue) shed

573 DNA into the bloodstream. The tumor DNA (blue) contains a tumor mutation (yellow). The

574 circulating free DNA in the blood becomes damaged both *in vivo* and *in vitro* (green triangle). Errors

575 can be introduced at each PCR duplication during amplification (red circle). Further errors are

576 accumulated during sequencing and mapping (purple square). The final data contains mapped reads,

577 where some mismatches are errors, and others are mutation from tumor cells.

578 **Figure 2:**

579 The data collection setup for tumor-informed relapse detection in colon cancer patients. After the

580 patient is diagnosed with colorectal cancer a liquid biopsy is extracted prior to curative surgery (Pre-

581 OP). A biopsy is taken from the tumor. Following surgery liquid biopsies (Post-OP) can be collected to

582 monitor relapse. All collected samples are sequenced using Next-Generation Sequencing.

583 **Figure 3:**

584 a) Examples of local sequence-context features and read-level features extracted from a read for a

585 single position of interest in a read mapping. Centered at the position of interest, the trinucleotide

586 context is extracted, and the surrounding 11 bp region is used for calculating the regional features,

587 including GC content and K-mer complexity. The read pairs contain a forward and reverse read that

588 are enumerated as either the first or second of the pair according to the order of sequencing. Two

589 read pairs are used for illustration of the read-centric features in the panels on the right. The UMI

590 groups are shown to indicate the variation in the number of reads used for the consensus reads. The

591 read position and fragment size are shown for the consensus reads. b-e) Variation in observed error

592 rate for selected features based on their observed distribution: b) Fragment size, c) UMI group size,

593 d) Read position and the variation between the first and second read in a pair. e) Error type for each

594    strand (forward and reverse). For each feature the 95% confidence interval is indicated by the

595    shaded areas or error bars. See **Supplementary section 3** for how the error rates and confidence

596    intervals are calculated and similar plots of the remaining features.

597    **Figure 4:**

598    a) Features are individually removed one-by-one from the full model containing all features to

599    measure the decrease in validation error. The most important feature is then defined as the one that

600    decreases the validation error the most, and vice versa. The grey points show the mean decrease in

601    validation error for each fold of a 5-fold cross validation. The average of these is used to rank the

602    features by importance, indicated by the black points. b) Based on the importance ranking, the

603    features are cumulatively removed one-by-one to from a full model. If the decrease in validation

604    error compared to the full model is significant, the feature should not be removed from the model. A

605    feature is only kept if removing it worsen the performance in all folds of the 5-fold cross validation.

606    c) Structure of the neural network. The neural network uses three different types of input features:

607    numeric, categorical, and embedded. The input features are processed differently in each group. The

608    input features are then parsed through three hidden layers of decreasing width. The output contains

609    4 nodes representing the probability of observing each of the four based (A, T, C, G) at the given

610    read position.

611    **Figure 5:**

612    a) Illustration of 5x2-cross-validation procedure for the estimation of performance. The patients are

613    first split into two approximately equally sized folds. The neural network model is trained on the

614    Post-OP data of fold 1 and validated by testing the models on the Pre-OP samples of the other fold

615    (Test B). This is then repeated by swapping the data in fold 1 and 2. The whole process is repeated 5

616    times. b) Performance of variant calling using DREAMS-*vc* compared to state-of-the-art tools

617    Shearwater and Mutect2. The AUC is estimated based on the different negative sets: The cross-

618    patient calls, 500 random validation alterations and these sets combined (All). The AUC is also

619  estimated for the full group of patients (All), and the patients with stage I and stage II CRC,

620  individually (ns: p≥0.05, *: p<0.05, **: p<0.01, ***: p<0.001, ****: p<0.0001).

621  **Figure 6:**

622  Prediction of cancer using DREAMS-cc (a) and Shearwater (b). For each patient's LB-sample (y-axis)

623  the mutation catalogue (x-axis) for every candidate patient is used for calling cancer. The patients

624  are stratified into patients with stage I and stage II CRC, respectively. The diagonal is showing the

625  result of using a patient's own mutation catalogue for cancer calling and constitutes the expected

626  positives. The off diagonal is the cross-patient results, for which the mutation catalogue is filtered

627  with the patient's tumour and germline variants prior to cancer calling, and thus these are expected

628  to be negative. The colour scheme is chosen based on the matched quantiles from the p-value and

629  combined Bayes factors from DREAMS-cc (a) and Shearwater, respectively. The cancer predictions

630  show the results from one split in the 5x2 CV. c) AUC performance of DREAMS-cc and shearwater

631  with respect to calling cancer.

# References

632

633  1.      Hu Z, Chen H, Long Y, Li P, Gu Y: **The main sources of circulating cell-free DNA: Apoptosis,**

634          **necrosis and active secretion.** *Critical Reviews in Oncology/Hematology* 2021, **157**:103166.

635  2.      Bettegowda C, Sausen M, Leary RJ, Kinde I, Wang Y, Agrawal N, Bartlett BR, Wang H, Luber B,

636          Alani RM, et al: **Detection of circulating tumor DNA in early- and late-stage human**

637          **malignancies.** *Science translational medicine* 2014, **6**:224ra224-224ra224.

638  3.      Phallen J, Sausen M, Adleff V, Leal A, Hruban C, White J, Anagnostou V, Fiksel J, Cristiano S,

639          Papp E, et al: **Direct detection of early-stage cancers using circulating tumor DNA.** *Science*

640          *Translational Medicine* 2017, **9**:eaan2415.

641  4.      Abbosh C, Birkbak NJ, Wilson GA, Jamal-Hanjani M, Constantin T, Salari R, Le Quesne J,

642          Moore DA, Veeriah S, Rosenthal R, et al: **Phylogenetic ctDNA analysis depicts early-stage**

643          **lung cancer evolution.** *Nature* 2017, **545**:446-451.

644    5.    Coakley M, Garcia-Murillas I, Turner NC: **Molecular Residual Disease and Adjuvant Trial**

645          **Design in Solid Tumors.** *Clinical Cancer Research* 2019, **25**:6026-6034.

646    6.    Henriksen TV, Tarazona N, Frydendahl A, Reinert T, Gimeno-Valiente F, Carbonell-Asins JA,

647          Sharma S, Renner D, Hafez D, Roda D, et al: **Circulating Tumor DNA in Stage III Colorectal**

648          **Cancer, beyond Minimal Residual Disease Detection, toward Assessment of Adjuvant**

649          **Therapy Efficacy and Clinical Behavior of Recurrences.** *Clinical Cancer Research* 2022,

650          **28**:507-517.

651    7.    Øgaard N, Reinert T, Henriksen TV, Frydendahl A, Aagaard E, Ørntoft M-BW, Larsen MØ,

652          Knudsen AR, Mortensen FV, Andersen CL: **Tumour-agnostic circulating tumour DNA analysis**

653          **for improved recurrence surveillance after resection of colorectal liver metastases: A**

654          **prospective cohort study.** *European Journal of Cancer* 2022, **163**:163-176.

655    8.    Cescon DW, Bratman SV, Chan SM, Siu LL: **Circulating tumor DNA and liquid biopsy in**

656          **oncology.** *Nature Cancer* 2020, **1**:276-290.

657    9.    Diehl F, Schmidt K, Choti MA, Romans K, Goodman S, Li M, Thornton K, Agrawal N, Sokoll L,

658          Szabo SA, et al: **Circulating mutant DNA to assess tumor dynamics.** *Nature Medicine* 2008,

659          **14**:985-990.

660    10.   Garcia-Murillas I, Schiavon G, Weigelt B, Ng C, Hrebien S, Cutts RJ, Cheang M, Osin P,

661          Nerurkar A, Kozarewa I, et al: **Mutation tracking in circulating tumor DNA predicts relapse**

662          **in early breast cancer.** *Science Translational Medicine* 2015, **7**:302ra133-302ra133.

663    11.   Corcoran RB, Chabner BA: **Application of Cell-free DNA Analysis to Cancer Treatment.** *New*

664          *England Journal of Medicine* 2018, **379**:1754-1765.

665    12.   Newman AM, Lovejoy AF, Klass DM, Kurtz DM, Chabon JJ, Scherer F, Stehr H, Liu CL, Bratman

666          SV, Say C, et al: **Integrated digital error suppression for improved detection of circulating**

667          **tumor DNA.** *Nature Biotechnology* 2016, **34**:547-555.

668    13.    Ma X, Shao Y, Tian L, Flasch DA, Mulder HL, Edmonson MN, Liu Y, Chen X, Newman S,

669           Nakitandwe J, et al: **Analysis of error profiles in deep next-generation sequencing data.**

670           *Genome Biology* 2019, **20**.

671    14.    Chen L, Liu P, Evans Thomas C, Ettwiller Laurence M: **DNA damage is a pervasive cause of**

672           **sequencing errors, directly confounding variant identification.** *Science* 2017, **355**:752-756.

673    15.    Benjamin D, Sato T, Cibulskis K, Getz G, Stewart C, Lichtenstein L: **Calling Somatic SNVs and**

674           **Indels with Mutect2.** Cold Spring Harbor Laboratory; 2019.

675    16.    Gerstung M, Papaemmanuil E, Campbell PJ: **Subclonal variant calling with multiple samples**

676           **and prior knowledge.** *Bioinformatics* 2014, **30**:1198-1204.

677    17.    Zviran A, Schulman RC, Shah M, Hill STK, Deochand S, Khamnei CC, Maloney D, Patel K, Liao

678           W, Widman AJ, et al: **Genome-wide cell-free DNA mutational integration enables ultra-**

679           **sensitive cancer monitoring.** *Nature Medicine* 2020, **26**:1114-1124.

680    18.    Wan JCM, Heider K, Gale D, Murphy S, Fisher E, Mouliere F, Ruiz-Valdepenas A, Santonja A,

681           Morris J, Chandrananda D, et al: **ctDNA monitoring using patient-specific sequencing and**

682           **integration of variant reads.** *Science Translational Medicine* 2020, **12**:eaaz8084.

683    19.    Pfeiffer F, Gröber C, Blank M, Händler K, Beyer M, Schultze JL, Mayer G: **Systematic**

684           **evaluation of error rates and causes in short samples in next-generation sequencing.**

685           *Scientific Reports* 2018, **8**:10950.

686    20.    Huptas C, Scherer S, Wenning M: **Optimized Illumina PCR-free library preparation for**

687           **bacterial whole genome sequencing and analysis of factors influencing de novo assembly.**

688           *BMC Research Notes* 2016, **9**:269.

689    21.    Meacham F, Boffelli D, Dhahbi J, Martin DIK, Singer M, Pachter L: **Identification and**

690           **correction of systematic error in high-throughput sequence data.** *BMC Bioinformatics* 2011,

691           **12**:451.

692  22.  Fan HC, Blumenfeld YJ, Chitkara U, Hudgins L, Quake SR: **Analysis of the Size Distributions of**

693  **Fetal and Maternal Cell-Free DNA by Paired-End Sequencing.** *Clinical Chemistry* 2010,

694  **56**:1279-1286.

695  23.  Lei J, G'Sell M, Rinaldo A, Tibshirani RJ, Wasserman L: **Distribution-Free Predictive Inference**

696  **for Regression.** *Journal of the American Statistical Association* 2018, **113**:1094-1111.

697  24.  Hayama T, Hashiguchi Y, Okamoto K, Okada Y, Ono K, Shimada R, Ozawa T, Toyoda T,

698  Tsuchiya T, Iinuma H, et al: **G12V and G12C mutations in the gene KRAS are associated with**

699  **a poorer prognosis in primary colorectal cancer.** *International Journal of Colorectal Disease*

700  2019, **34**:1491-1496.

701  25.  Parikh A, Goyal L, Hazar-Rethinam M, Siravegna G, Blaszkowsky L, Russo M, Van Seventer E,

702  Nadres B, Shahzade H, Clark J, et al: **Systematic liquid biopsy identifies novel and**

703  **heterogeneous mechanisms of acquired resistance in gastrointestinal (GI) cancer patients.**

704  *Annals of Oncology* 2017, **28**:iii137.

705  26.  Cristiano S, Leal A, Phallen J, Fiksel J, Adleff V, Bruhm DC, Jensen S, Medina JE, Hruban C,

706  White JR, et al: **Genome-wide cell-free DNA fragmentation in patients with cancer.** *Nature*

707  2019, **570**:385-389.

708  27.  Mouliere F, Rosenfeld N: **Circulating tumor-derived DNA is shorter than somatic DNA in**

709  **plasma.** *Proceedings of the National Academy of Sciences* 2015, **112**:3178-3179.

710  28.  Pozzolo AD, Caelen O, Johnson RA, Bontempi G: **Calibrating Probability with Undersampling**

711  **for Unbalanced Classification.** In; *2015*. IEEE;

712  29.  Abadi M, Agarwal A, Barham P, Brevdo E, Chen Z, Citro C, Corrado GS, Davis A, Dean J, Devin

713  M: **Tensorflow: Large-scale machine learning on heterogeneous distributed systems.** *arXiv*

714  *preprint arXiv:160304467* 2016.

715

Shedding

Normal cell

Cancer cell

★ Mutation

Data

Reads from tumor tissue

Reads from normal tissue

Mutated position

▲ : DNA damage (*in vivo* + *in vitro*)

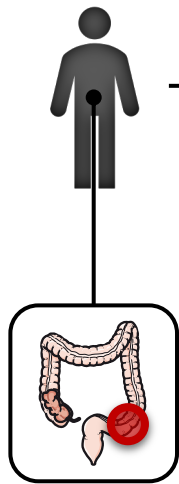● : PCR errors in amplification

■ : Read errors in sequencing

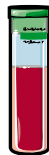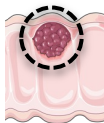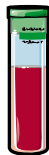Colorectal cancer diagnosis → Blood sample (Pre-OP) → Curative surgery → Blood sample (Post-OP)
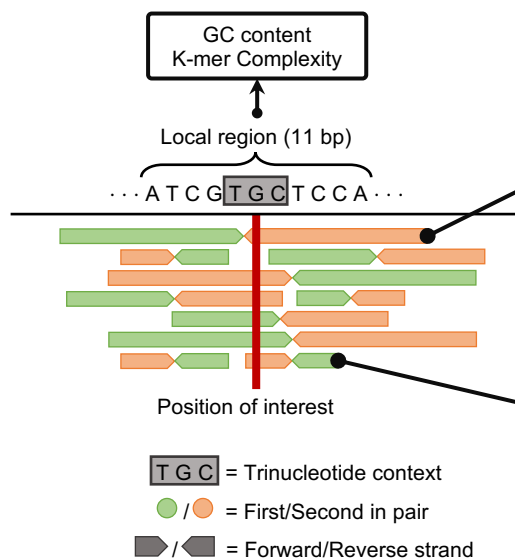
Patient sample set

Blood sample    Tumor biopsy    Blood sample
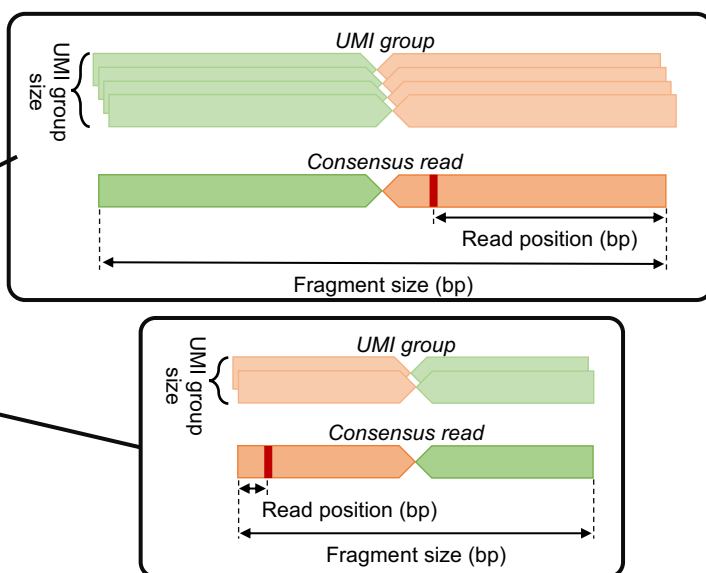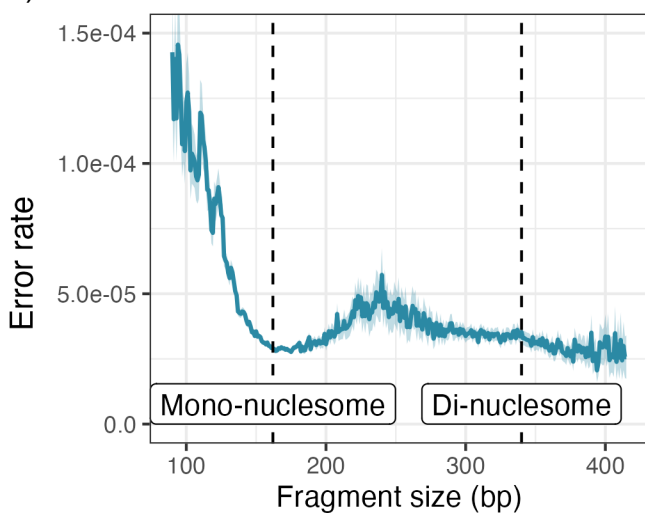
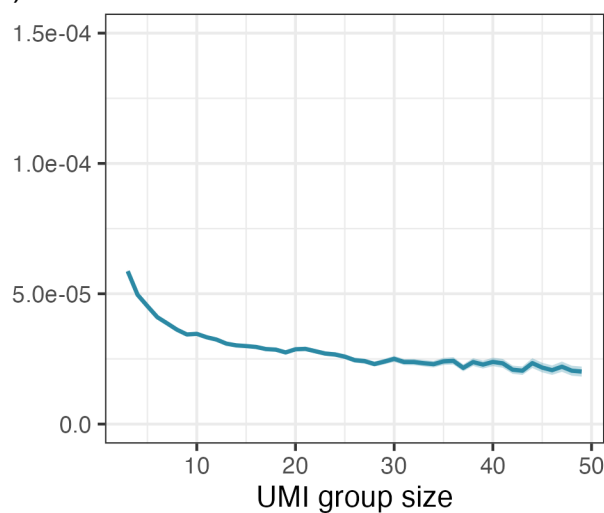Next-Generation Sequencing
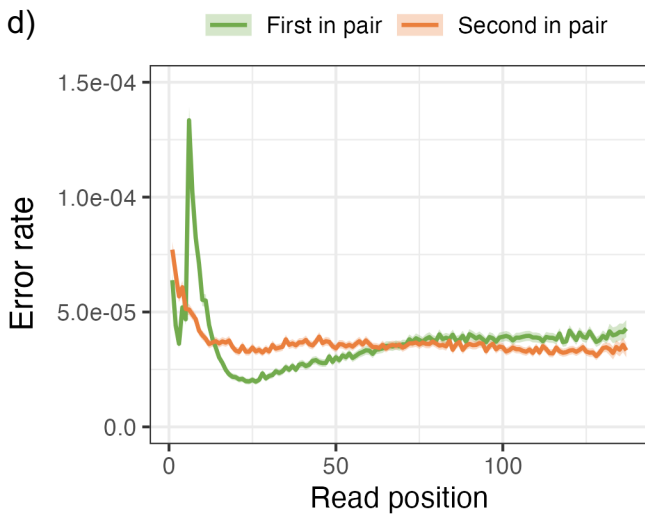
a) **Local sequence-context features** / **Read-level features**

GC content
K-mer Complexity

Local region (11 bp)

··· A T C G [T G C] T C C A ···

Position of interest

[T G C] = Trinucleotide context

● / ● = First/Second in pair

▶ / ◀ = Forward/Reverse strand

*UMI group*

UMI group size

*Consensus read*

Read position (bp)

Fragment size (bp)
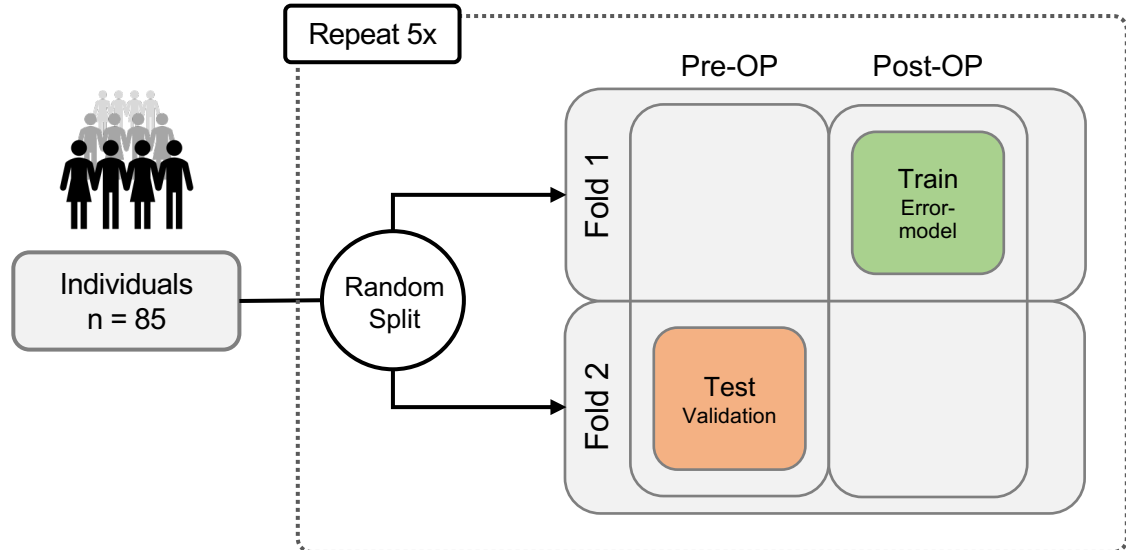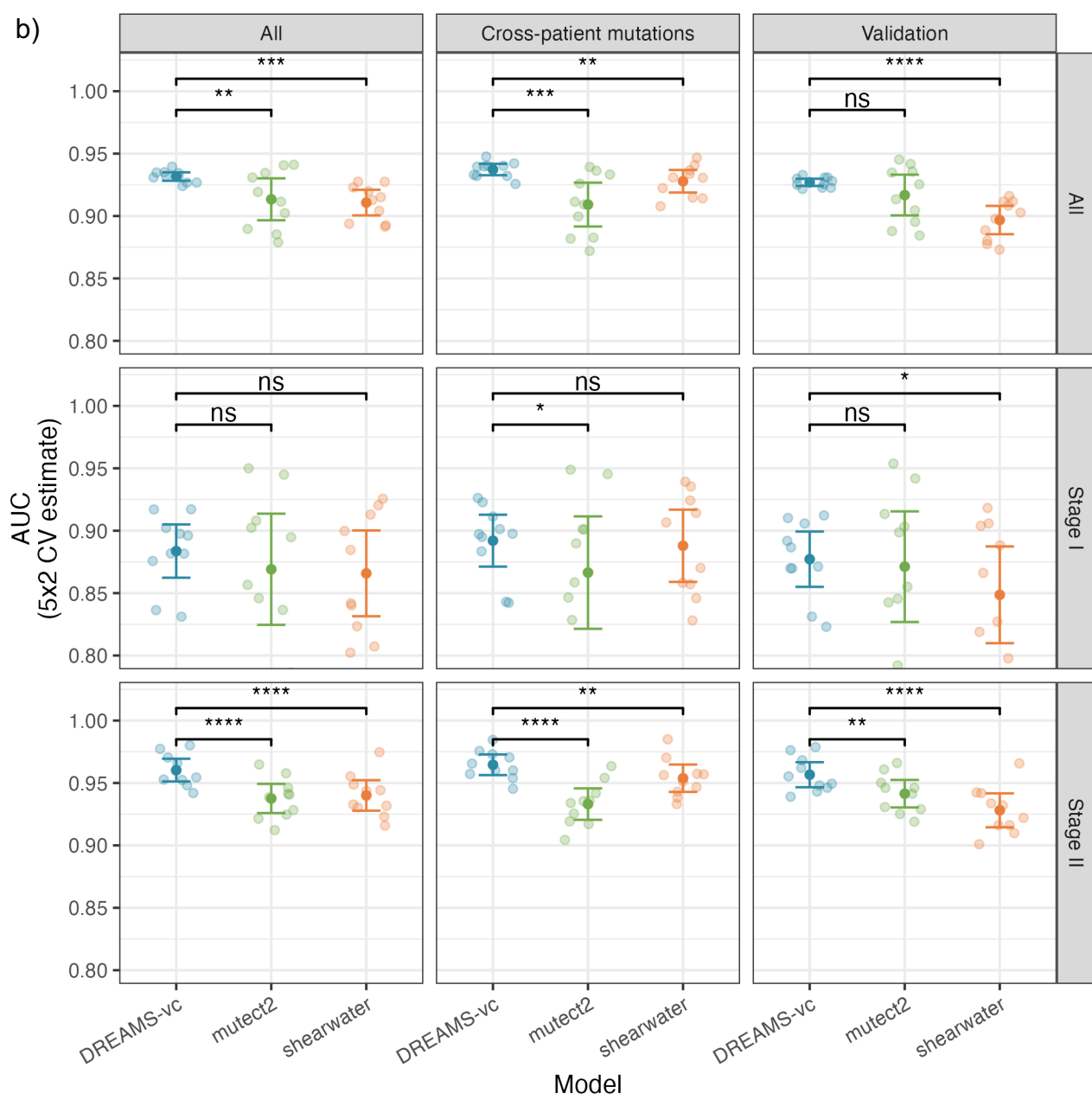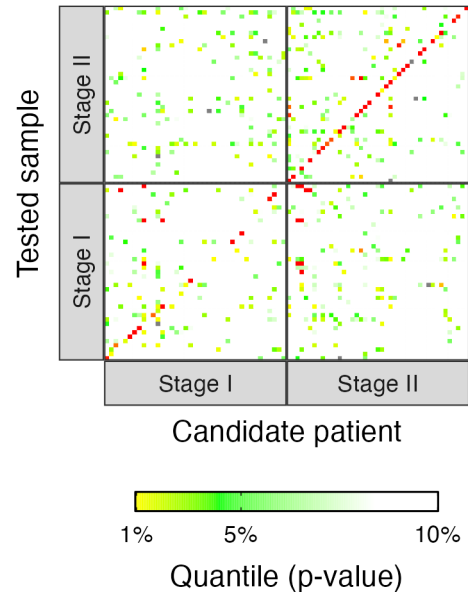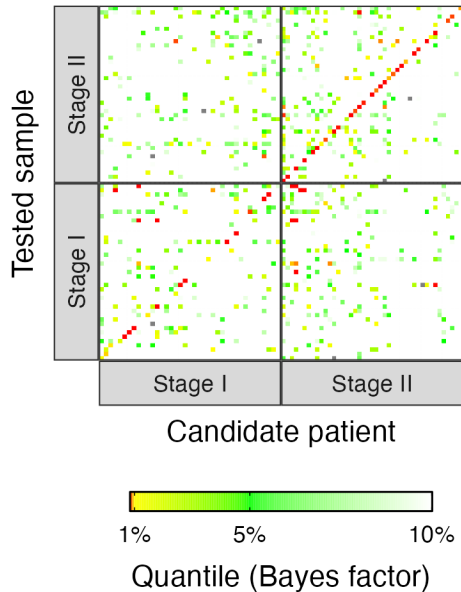
b) Mono-nucleosome / Di-nucleosome

c)

d) First in pair / Second in pair

e) Forward / Reverse

a) DREAMS-*cc*  b) Shearwater  c)

**a)** Tested sample (Stage II, Stage I) × Candidate patient (Stage I, Stage II)
Quantile (p-value): 1% — 5% — 10%

**b)** Tested sample (Stage II, Stage I) × Candidate patient (Stage I, Stage II)
Quantile (Bayes factor): 1% — 5% — 10%

**c)** AUC (5x2 CV estimate) vs Model (DREAMS-cc, shearwater)
p = 0.0343