

1 *The Lancet* 2020; 395:10221: 350-360 (DOI: [10.1016/S0140-6736\(19\)32998-8](https://doi.org/10.1016/S0140-6736(19)32998-8))

2

3 **Deep learning for prediction of colorectal cancer outcome: a discovery and validation**

4 **study**

5

6 Ole-Johan Skrede, M. Sc.^{1,2,*}, Sepp De Raedt, Ph. D.^{1,2,*}, Andreas Kleppe, Ph. D.^{1,2}, Tarjei S.

7 Hveem, Ph. D.¹, Prof. Knut Liestøl, Ph. D.^{1,2}, John Maddison, Ph. D.¹, Hanne A. Askautrud,

8 Ph. D.¹, Manohar Pradhan, Ph. D.¹, John Arne Nesheim, M. Sc.¹, Prof. Fritz Albrechtsen, M.

9 Sc.^{1,2}, Prof. Inger Nina Farstad, Ph. D.^{3,4}, Enric Domingo, Ph. D.⁵, David N. Church, D.

10 Phil.^{6,7}, Prof. Arild Nesbakken, Ph. D.^{4,8,9}, Prof. Neil A. Shepherd, D. M.¹⁰, Prof. Ian

11 Tomlinson, Ph. D.^{1,11}, Prof. Rachel Kerr, Ph. D.⁵, Prof. Marco Novelli, Ph. D.^{1,12}, Prof. David

12 J. Kerr, D. Sc.¹³, Prof. Håvard E. Danielsen, Ph. D.^{1,2,13**}

13

14 ¹Institute for Cancer Genetics and Informatics, Oslo University Hospital, Oslo, Norway

15 ²Department of Informatics, University of Oslo, Oslo, Norway

16 ³Department of Pathology, Division of Laboratory Medicine, Oslo University Hospital, Oslo,

17 Norway

18 ⁴Institute of Clinical Medicine, University of Oslo, Oslo, Norway

19 ⁵Department of Oncology, University of Oxford, Oxford, UK

20 ⁶NIHR Oxford Biomedical Research Centre, Oxford University Hospitals NHS Foundation

21 Trust, John Radcliffe Hospital, Oxford, UK

22 ⁷Wellcome Centre for Human Genetics, University of Oxford, Oxford, UK

23 ⁸Department of Gastrointestinal Surgery, Oslo University Hospital, Oslo, Norway

24 ⁹K.G. Jebsen colorectal cancer research centre, Oslo, Norway

25 ¹⁰Gloucestershire Cellular Pathology Laboratory, Cheltenham General Hospital, Cheltenham,
26 UK

27 ¹¹Edinburgh Cancer Research Centre, University of Edinburgh, Edinburgh, Scotland

28 ¹²Research Department of Pathology, University College London Medical School, London,
29 UK

30 ¹³Nuffield Division of Clinical Laboratory Sciences, University of Oxford, Oxford, UK

31

32 *Both authors contributed equally to this work.

33 **Corresponding author:

34 Prof Håvard E. Danielsen,

35 Institute for Cancer Genetics and Informatics,

36 Oslo University Hospital

37 Montebello, 0310, Oslo, Norway

38 Email: hdaniels@labmed.uio.no

39 Phone: +47 22782320

40

41 Words in abstract (not exceed 300): 297

42 Words in main text (up to 3500): 3889

43 Number of references (up to 30): 30

44 Number of figures: 2

45 Number of tables: 3

46 **Background:** Improved markers of prognosis are needed to stratify patients with early-stage
47 colorectal cancer to refine selection of adjuvant therapy. The aim of the present study was to
48 develop a biomarker of patient outcome after primary colorectal cancer resection by directly
49 analysing scanned conventional haematoxylin and eosin stained sections using deep learning.

50 **Methods:** More than 12,000,000 image tiles from 828 patients with distinctly good or poor
51 disease outcome were used to train a total of 10 convolutional neural networks, purpose-built
52 for classifying supersized heterogeneous images. A prognostic biomarker integrating the 10
53 networks were determined using 1645 patients with non-distinct outcome. The marker was
54 tested on 920 patients with slides prepared in UK, and finally independently validated
55 according to a pre-defined protocol in 1122 patients treated with single-agent capecitabine
56 using slides prepared in Norway. The primary outcome was cancer-specific survival.

57 **Findings:** The biomarker provided a hazard ratio for poor vs good prognosis of 3.84 (95%
58 confidence interval, 2.72-5.43; $p < 0.0001$) in the primary analysis of the validation cohort,
59 and 3.04 (95% confidence interval, 2.07-4.47; $p < 0.0001$) after adjusting for established
60 prognostic markers significant in univariable analyses of the same cohort; pN stage, pT stage,
61 lymphatic invasion, and venous vascular invasion.

62 **Interpretation:** It was possible to develop a clinically useful prognostic marker using deep
63 learning allied to digital scanning of conventional haematoxylin and eosin stained tumour
64 tissue sections. The assay has been extensively evaluated in large, independent patient
65 populations, correlates with and outperforms established molecular and morphological
66 prognostic markers, and gives consistent results across tumour and nodal stage. The
67 biomarker stratified stage II and III patients into sufficiently distinct prognostic groups that
68 these potentially could be used to guide selection of adjuvant treatment by avoiding therapy in
69 very low risk groups and identifying patients who would benefit from more intensive regimes.

70 **Funding:** The Research Council of Norway through its IKTPLUS Lighthouse program

71 (grant number 259204, project name DoMore!).

72

73 **Research in context**

74 **Evidence before this study**

75 Digital image analysis is one of the fields where the recent renaissance of deep learning has
76 achieved the most impressive results. We searched PubMed on June 12, 2019 without
77 language or time restrictions, using the terms “deep learning”, “prediction”, “survival”,
78 “cancer”, and “histology” (full specification of the search criteria is provided in the appendix
79 p 3). We systematically reviewed the 214 search results, and found 18 original research
80 studies which applied deep learning to predict patient outcome or related attributes using
81 histopathology images.

82

83 In 16 studies, the patient outcome was indirectly predicted by identifying attributes known to
84 correlate with patient outcome, e.g. stromal fraction, mitotic count, or Gleason pattern. Two
85 studies reported on direct prediction of survival, but neither presented a marker for automatic
86 prediction of patient outcome from scanned whole-slide sections; one required manual
87 annotation to locate interesting tissue regions, and the other classified tissue microarray spots.
88 Perhaps even more importantly, neither of these two studies evaluated their biomarker in
89 independent cohorts; the performance was instead estimated using cross-validation in the
90 same cohort as utilised for training, which can easily lead to overoptimistic estimates.

91

92 **Added value of the study**

93 We have applied deep learning to develop a biomarker for automatic prediction of cancer-
94 specific survival directly from scanned haematoxylin and eosin stained, formalin-fixed,
95 paraffin-embedded tumour tissue sections. Independent validation demonstrated that the

96 biomarker improved prediction of cancer-specific survival by stratifying stage II and III
97 colorectal cancer patients into distinct prognostic groups, supplementing established
98 prognostic markers, and outperforming most existing markers in terms of hazard ratios. The
99 marker could potentially be used to improve selection of adjuvant treatment after resection of
100 colorectal cancer by identifying patients at very low risk who may have been cured by surgery
101 alone, as well as patients at high risk who are much more likely to benefit from more
102 intensive regimes.

103

104 **Implications of all the available evidence**

105 It is possible to utilise deep learning to develop biomarkers for automatic prediction of patient
106 outcome directly from conventional histopathology images. In colorectal cancer, the marker
107 was found to be a clinically useful prognostic marker in analysis of a large series of patients
108 who received consistent, modern cancer treatment.

109

110 **Introduction**

111 Biomarkers are being used increasingly to match anticancer therapy to specific tumour
112 genotypes, protein, and RNA expression profiles, usually in patients with advanced disease.¹⁻³
113 One example of this is selection of *KRAS*-wild-type colorectal cancers (CRCs) for treatment
114 with epidermal growth factor receptor inhibitors.⁴ However, in the adjuvant setting for CRC,
115 the primary question is binary, whether to offer treatment at all, and subsequent selection of
116 drugs, dose, and schedule is predominantly driven by stage rather than by companion
117 diagnostics. If it were possible to further refine prognostic models, this could allow a more
118 targeted approach by defining subgroups in which the absolute benefits of adjuvant
119 chemotherapy are minimal, relative to surgery alone, and at the other end of the spectrum,
120 patients who might benefit from prolonged combination chemotherapy because of their poor
121 survival rate.⁵⁻⁸

122 More than two decades of adjuvant trials in patients with early-stage CRC using
123 fluoropyrimidines, in combination with cytotoxic agents like oxaliplatin, have yielded an
124 improved overall survival of around 3-5% for patients with stage II or IIIA CRC. Many
125 patients are cured by surgery alone, while around 25% will recur despite adjuvant
126 chemotherapy. There is likely to be a chemotherapy-associated death rate of 0.5-1%, and 20%
127 of patients will suffer significant side-effects. The risk-benefit ratio is therefore rather
128 marginal, but could potentially be much better if it were possible to define subgroups at
129 higher or lower risk of recurrence and cancer-specific death.⁹⁻¹²

130 Although clinically validated prognostic biomarkers would facilitate adjuvant therapeutic
131 decisions, very few have been sufficiently robustly validated for routine clinical application.
132 A case can be made for assessment of mismatch repair (MMR) status,^{13,14} as patients with
133 MMR-deficient tumours tend to have a good prognosis. We have recently reported that
134 measurement of tumour cellular DNA content (ploidy) in combination with stromal fraction

135 can stratify stage II patients into very good, intermediate, and poor prognostic groups.¹⁵
136 Interestingly, analysis of driver mutations and RNA signatures has shown them to be
137 individually weak prognostic markers and unable to guide clinical decision making.^{8,14}
138 Deep learning refers to the class of machine learning methods that make use of successively
139 more abstract representations of the input data to perform a specific task. These methods use a
140 training set to learn how these representations should be generated in a manner appropriate for
141 the given task. In contrast, traditional machine learning utilises handcrafted features to create
142 representations of the input data that are applied to perform the task. In many applications,
143 deep learning has been demonstrated to provide superior performance compared to other
144 machine learning techniques, and it is a growing expectation that deep learning will transform
145 current medical practice. Especially convolutional neural networks have excelled in many
146 image interpretation tasks, and could therefore be hypothesised to retrieve additional
147 information from histopathology images. The aim of the present study was to use deep
148 learning to analyse conventional whole-slide images (WSIs) in order to develop an automatic
149 prognostic biomarker for patients resected for primary CRC. The marker was trained using
150 828 patients with distinct prognosis from four cohorts, fine-tuned using 1645 other patients
151 from the same four cohorts, and tested on slides prepared at a different laboratory from 920
152 patients. Finally, the marker was independently validated according to the pre-defined
153 protocol (appendix pp 52-80) on 1122 patients analysed retrospectively from a trial
154 (QUASAR 2) of adjuvant therapy.¹⁶

155

156 **Methods**

157 **Training and Tuning Cohorts**

158 Four different cohorts were utilised for training and tuning to achieve a broad patient
159 representation and thereby improve the ability to generalise to new cohorts. Three cohorts

160 were consecutive series of stage I, II or III tumours from CRC patients treated at hospitals
161 with both rural and urban catchment areas: (i) 160 patients treated 1988-2000 at Akershus
162 University Hospital, Norway;¹⁷ (ii) 576 patients treated 1993-2003 at Aker University
163 Hospital, Norway;¹⁵ and (iii) 970 patients treated in Gloucester 1988-1996 and included in the
164 Gloucester Colorectal Cancer Study, UK.^{18,19} The fourth cohort were 767 stage II or III CRC
165 patients treated at 151 UK hospitals in 2002-2004 and included in the VICTOR trial (ISRCTN
166 registry number ISRCTN98278138).²⁰ Our cohorts included only patients with resectable
167 tumour, and a formalin-fixed, paraffin-embedded (FFPE) tumour tissue block available for
168 analysis.

169 To obtain clear ground-truth, we used as training cohort the 828 patients with so-called
170 distinct outcome, either good or poor. A patient was assigned to the good outcome group if
171 aged less than 85 years at surgery, had more than six years follow-up after surgery, and had
172 no record of recurrence or cancer-specific death. The poor outcome group consisted of those
173 aged less than 85 years at surgery and suffered cancer-specific death between 100 days
174 (inclusive) and 2.5 years (exclusive) after surgery. Patients not satisfying either of these group
175 criteria were defined as having non-distinct outcome, and these 1645 patients were used for
176 tuning. The protocol specifies additional cohort details, and demographics are summarised in
177 table 1.

178 **Test Cohort**

179 The test cohort consisted of 920 patients from the Gloucester Colorectal Cancer Study,
180 UK.^{18,19} WSIs were obtained from different FFPE tumour tissue blocks than those used in the
181 training and tuning cohorts.

182 **Validation Cohort**

183 The validation cohort consisted of 1122 patients from 170 hospitals in seven countries
184 recruited to the QUASAR 2 trial (ISRCTN registry number ISRCTN45133151).¹⁶ Inclusion

185 criteria were age 18 years or older, CRC adenocarcinoma histologically proven to be R0 M0
186 stage III or high-risk stage II, primary resection 4-10 weeks before randomisation, WHO
187 performance status score 0 or 1, and life expectancy (with comorbidities, but excluding cancer
188 risk) of at least five years. See protocol pp 22-25 for exclusion criteria and other details. All
189 patients received adjuvant therapy, either capecitabine plus bevacizumab or capecitabine
190 alone, with equal disease-free and overall survival in both trial arms.¹⁶

191 **Sample Preparation**

192 Slides in VICTOR cohort were prepared in Oxford, UK, while the other slides in the training
193 and tuning cohorts were prepared at the Institute for Cancer Genetics and Informatics (ICGI),
194 Norway. Introducing this variation in the development phase was hypothesised to increase the
195 robustness and generalisability of the trained marker. Slides in the test cohort were prepared
196 as a part of the routine histopathological examination in Cheltenham, UK, and the
197 performance in this cohort should thus indicate the prognostic ability when the marker is
198 assayed at a different laboratory using original slides. Slides in the validation cohort were
199 prepared at ICGI. All slides were made by staining a three μm FFPE tissue block section with
200 haematoxylin and eosin (H&E), and a pathologist (MP) ascertained that it contained tumour.
201 WSIs were acquired at the highest resolution available (referred to as 40x magnification by
202 the manufacturers) on two scanners, an Aperio AT2 (Leica Biosystems, Germany) and a
203 NanoZoomer XR (Hamamatsu Photonics, Japan).
204 Areas with high tumour content were identified using a segmentation network that was trained
205 on a subset of the training and tuning cohorts (protocol pp 6-10). A WSI with the so-called
206 40x resolution typically contained an order of 100,000x100,000 pixels, multiple orders of
207 magnitude larger than images currently feasible for classification by deep learning methods.
208 To preserve prognostic information contained at high-resolution, WSIs were partitioned into
209 multiple non-overlapping image regions called *tiles* at 10x and 40x resolutions, where each

210 pixel at 40x represents a physical size of approximately $0.24 \times 0.24 \mu\text{m}^2$. Patients without tiles
211 were excluded.

212 **Classification**

213 Five networks were trained on the 634,564 10x tiles and five networks on the 11,591,555 40x
214 tiles from the 1652 Aperio AT2 and NanoZoomer XR WSIs in the training cohort with the
215 patients' distinct outcomes as ground-truth. All networks were DoMore v1 networks, which
216 we designed for classifying supersized heterogeneous images. The DoMore v1 network was
217 built around multiple instance learning and comprised of a MobileNetV2²¹ representation
218 network, a Noisy-AND pooling function,²² and a fully-connected classification network
219 similar to the one used by Kraus et al²² (figure 1). Because of spatial heterogeneity, labelling a
220 tile with the label of its WSI might be problematic. Instead, the networks were trained on
221 labelled collections of tiles. A collection contained tiles from a single WSI, which label it
222 inherits. Collections of tiles were processed by the representation network before the resulting
223 tile representations were pooled and classified. The entire network was trained end-to-end, i.e.
224 directly from image to patient outcome, and each training iteration used a batch size of 32
225 collections with 64 tiles each. This many tiles were possible because we utilised a novel
226 gradient approximation technique which substantially reduce memory usage during training
227 (appendix pp 4-6). The Noisy-AND pooling function applied a trained non-linear function on
228 tile representation averages. This enhances robustness against tiles not representing the
229 ground-truth, and together with the large number of tiles, alleviates the issues of spatial
230 heterogeneity. During inference, the network processed all tiles in the WSI.

231 The networks were trained beyond apparent convergence using TensorFlow 1.10, and a
232 model was selected from each network training using the performance in the tuning cohort
233 with the c-index as metric, resulting in five models for each resolution (protocol pp 11-20).
234 Each of the five models provides a score reflecting the probability of poor outcome, and the

235 average was defined as the ensemble score. For use in categorical markers, suitable thresholds
236 for the 10x and the 40x ensemble scores were determined by evaluations in the tuning cohort
237 to define the ensemble classifiers (protocol pp 20-22). Furthermore, evaluations in the test
238 cohort indicated that combining 10x and 40x markers might be desirable, and two such
239 markers were defined, one continuous and one categorical. The continuous DoMore-v1-CRC
240 score was defined as the average of the 10x and the 40x ensemble scores. The categorical
241 DoMore-v1-CRC classifier assigned to good prognosis if both ensemble classifiers predicted
242 good outcome, uncertain if the ensemble classifiers predicted differently, and poor prognosis
243 if both predicted poor outcome. In a post-hoc analysis, the continuous DoMore-v1-CRC score
244 was categorised into five risk groups (appendix p 6).

245 Inception v3, a state-of-the-art convolutional neural network, was trained, tuned, and
246 evaluated with the same study setup as the DoMore v1 network (protocol pp 11-22), and
247 tested as a secondary analysis (protocol p 27). While the DoMore-v1-CRC marker was trained
248 using multiple instance learning, each single tile was labelled with the label of its WSI in
249 training the Inception v3 marker. The image distortion algorithm and network
250 hyperparameters were determined independently of the DoMore v1 network in the discovery
251 phase, resulting in slightly different choices for the Inception v3 network (protocol pp 15-16).

252 **Statistical Analysis**

253 This study conformed to the REMARK guideline²³ and relevant aspects of the guideline
254 proposed by Luo et al²⁴ (appendix pp 7-8). Primary and secondary analyses were planned in
255 advance of evaluations in the validation cohort and described in the protocol.

256 The pre-defined primary analysis for each scanner was univariable cancer-specific survival
257 (CSS) analysis of the DoMore-v1-CRC classifier; for simplicity, we first present results for
258 the Aperio AT2 scanner and in a separate paragraph address scanner differences. The
259 classifier was included as the only variable in a Cox model to compute the hazard ratio (HR)

260 with 95% confidence interval (CI) of patients with uncertain and poor prognosis relative to
261 patients with good prognosis. The proportional hazards assumption was found satisfactory
262 fulfilled using log-log plots (appendix p 26). The Mantel-Cox log-rank test was used to assess
263 whether the classifier predicted CSS.

264 Both the classifier and the continuous score were evaluated in multivariable Cox models as
265 secondary and post-hoc analyses, including markers available at the time of analysis (patients
266 with at least one missing value were excluded). To calculate classification metrics for 3-year
267 CSS, patients without event and less than 3-year follow-up were excluded and events after 3
268 years were ignored. Category-free net reclassification improvement (NRI) was computed
269 using the Kaplan-Meier estimates of five-year CSS. Two-sided $p < 0.05$ was considered
270 statistically significant. The confidence level of CIs is 95%. The bias-corrected and
271 accelerated bootstrap CI were computed for NRIs, c-indices and areas under the curves
272 (AUCs) using 10,000 bootstrap replicates and an acceleration constant estimated using leave-
273 one-out cross-validation. Time to CSS in the validation cohort was calculated from date of
274 randomisation to date of cancer-specific death or loss to follow-up. Survival analyses were
275 carried out in Stata/SE 15.1 (StataCorp, TX).

276 **Role of the funding source**

277 The funders had no role in study design, data collection, data analysis, data interpretation,
278 writing the report, or the decision to submit the paper for publication. The corresponding
279 author had full access to all data and the final responsibility to submit for publication.

280

281 **Results**

282 The DoMore-v1-CRC classifier was a strong predictor of CSS in the primary analysis of the
283 validation cohort (HR for uncertain vs good prognosis, 1.89; CI, 1.14-3.15; HR for poor vs
284 good prognosis, 3.84; CI, 2.72-5.43; figure 2A). The classifier remained strong in

285 multivariable analysis (HR for uncertain *vs* good prognosis, 1·56; CI, 0·92-2·65; HR for poor
286 *vs* good prognosis, 3·04; CI, 2·07-4·47; table 2) adjusting for established prognostic markers
287 significant in univariable analyses; pN stage, pT stage, lymphatic invasion, and venous
288 vascular invasion (appendix p 9).

289 The sensitivity was 52% (CI, 41%-63%), specificity 78% (CI, 75%-81%), positive predictive
290 value 19% (CI, 14%-25%), negative predictive value 94% (CI, 92%-96%), and correct
291 classification rate 76% (CI, 73%-79%) when comparing 3-year CSS to good prognosis *vs*
292 uncertain and poor prognosis. Compared to good and uncertain prognosis *vs* poor prognosis,
293 the sensitivity was 69% (CI, 58%-78%), specificity 66% (CI, 63%-69%), positive predictive
294 value 17% (CI, 13%-21%), negative predictive value 96% (CI, 94%-97%), and correct
295 classification rate 67% (CI, 63%-69%).

296 The constituents of the DoMore-v1-CRC classifier, the 10x and the 40x ensemble classifiers,
297 were strong predictors in univariable (appendix p 27) and multivariable analyses (appendix pp
298 10-11). The ensemble classifiers performed similarly as the best classifiers based on one of
299 the ten individual models that constituted the ensemble models (appendix pp 12 and 28-29).

300 The continuous ensemble scores were also strong predictors in univariable (appendix p 9) and
301 multivariable analyses (appendix pp 13-15). The DoMore-v1-CRC score associated strongly
302 with the patient outcome (appendix p 30), and provided a c-index of 0·674 (CI, 0·624-0·719;
303 appendix p 16) in all validation patients and an AUC of 0·713 (CI, 0·624-0·789; appendix p
304 31) in patients with distinct outcome. The c-index and AUC of the 10x ensemble score were
305 similar to the ones obtained for the DoMore-v1-CRC score (appendix pp 16 and 31).

306 The DoMore-v1-CRC classifier was a significant predictor of CSS in stage II (HR for poor *vs*
307 good prognosis, 2·71; CI, 1·25-5·86; figure 2C) and stage III (HR for poor *vs* good prognosis,
308 4·09; CI, 2·77-6·03; figure 2D), and this was confirmed in multivariable analysis (table 2) and
309 for the continuous score (appendix pp 9 and 13). The categorical marker identified patient

310 groups with substantially different CSS in stage IIIB and IIIC (appendix p 32), and was also
311 significant in pN stages (figures 2C, E, and F) and pT stages (pT1-3 vs pT4; appendix p 33).
312 The category-free NRI of supplementing substage with the DoMore-v1-CRC class for
313 prediction of five-year CSS was 61·6% (CI, 43·5%-79·3%); the event-NRI was 3·2% (CI, -
314 13·2%-20·0%), and the non-event-NRI was 58·3% (CI, 52·7%-63·8%).

315 The DoMore-v1-CRC classifier correlated with a number of factors such as age, pN stage, pT
316 stage, histological grade, location, tumour sidedness, *BRAF* mutation, and microsatellite
317 instability (table 3). Of special interest is the relation to the histopathological grading into
318 well, moderately, and poorly differentiated tumours. This was further studied in the test
319 cohort where all gradings were centrally reviewed by one highly experienced pathologist
320 (NAS).^{18,19} Among 133 tumours characterised as well differentiated, the DoMore-v1-CRC
321 classifier assigned 101 as good prognosis, 18 as uncertain and 14 as poor prognosis (appendix
322 p 17). The moderately differentiated tumours were distributed fairly evenly over the DoMore-
323 v1-CRC classes, while among 292 poorly differentiated tumours, the marker assigned 223 as
324 poor prognosis, 36 as uncertain, and 33 as good prognosis. Thus, the DoMore-v1-CRC class
325 was clearly associated to tumour differentiation. The large proportion of tumours classified as
326 moderately differentiated (e.g. 53% [489 of 920] in the test cohort and 75% [846 of 1122] in
327 the validation cohort) restricts the usefulness of this grading system, but also these patients
328 could be risk stratified by the DoMore-v1-CRC marker (appendix p 34).

329 Median processing time per patient for the entire classification pipeline, i.e. from scan to
330 predicted patient outcome, was 2·8 minutes (interquartile range, 1·8-3·9) in the validation
331 cohort on a computer with an NVIDIA GeForce RTX 2080 Ti and an Intel Core i7-7700K.
332 Inception v3 provided a marker of CSS with only slightly worse performance than the
333 DoMore-v1-CRC classifier (appendix pp 16 and 35-36).

334 In the test cohort with slides prepared at a different hospital, the classifier provided similar
335 HRs (appendix p 37) as in the validation cohort (figure 2), supporting that it is robust against
336 inter-laboratory differences in tissue preparation and staining.
337 When evaluated using another scanner (NanoZoomer XR), the DoMore-v1-CRC score tended
338 towards slightly higher values compared to when evaluated using the Aperio AT2 scanner,
339 resulting in a higher DoMore-v1-CRC class for some patients near the classification
340 thresholds (appendix p 38). However, the scores correlated strongly (Pearson's $r=0.956$; CI,
341 $0.951-0.961$), and the classifier provided similar prognostic information with both scanners
342 (see appendix pp 9, 16, 18-25, and 39-51 for results with NanoZoomer XR). Thus, the
343 classifier was also a strong predictor of CSS in the primary analysis of the validation cohort
344 when evaluated on NanoZoomer XR slide images (HR for uncertain vs good prognosis, 2.42 ;
345 CI, $1.45-4.03$; HR for poor vs good prognosis, 3.39 ; CI, $2.36-4.87$; appendix p 39).

346

347 **Discussion**

348 Building on recent developments in machine learning, we have developed a biomarker for
349 automatic prediction of the outcome of a patient resected for early-stage CRC which directly
350 analyse standard H&E stained histological sections. To assay the biomarker, one
351 convolutional neural network first automatically outlines cancerous tissue, and then a second
352 convolutional neural network stratifies the patients into prognostic categories. In the
353 validation, the good and poor prognosis groups included nearly 90% of the patients and
354 differed about 4 times in HR for CSS in univariable analysis and about 3 times in
355 multivariable analysis. The multivariable result indicated that the new biomarker will be a
356 useful supplement to the established markers and improve risk stratification.
357 Deep learning has already been shown to be suitable for detection and delineation of some
358 tumour types,²⁵ and various cancer classifications have been reported.²⁶ Recent studies have

359 suggested that deep learning could be used to develop markers which potentially utilise basic
360 morphology to predict the outcome of cancer patients, but these findings have not been
361 validated in independent cohorts.^{27,28} We have not yet seen independently validated markers
362 for directly predicting the outcome of cancer patients based on histological images.

363 We derived two markers using the same study setup, but different deep learning techniques.
364 In training the Inception v3 marker, each tile was labelled with the label of its WSI, while the
365 DoMore-v1-CRC marker was developed using multiple instance learning to allow training on
366 tile collections labelled with the label of its WSI. Both markers were strong predictors of CSS,
367 but the DoMore-v1-CRC marker performed slightly better and was the marker pre-selected
368 for independent validation in the QUASAR 2 cohort.

369 Automatic prognostication procedures reduce human intervention, and has the potential to
370 increase reproducibility of biomarkers. New procedures like the DoMore-v1-CRC markers
371 may initially be performed as services carried out at specialised laboratories with a high
372 degree of standardisation of procedure to avoid disparities in sample handling, including the
373 staining and scanning. Such centralised processing will also facilitate the collection of
374 information on new procedures and enable improvements in the decision support to
375 pathologists and clinicians. As an increasing number of laboratories are becoming digitalised,
376 accompanying decision support systems may include standardisation modules and facilitate a
377 more rapid spread of the automatic procedures. Moreover, supplemented by increased
378 robotisation of wet-lab procedures, the higher analytic throughput will allow decisions based
379 on multiple samples from a tumour. This may reduce the challenge of tumour heterogeneity,
380 which may be a key to improved accuracy of prognosis.

381 The DoMore-v1-CRC biomarker correlated with several recognised prognostic factors,
382 including the histological grading carried out by a specialised pathologist. The classifier
383 performed better than most other markers in terms of HRs in stage-specific multivariable

384 analyses, on a par with pN staging. As opposed to the grading system, the classifier had few
385 patients in the intermediate “uncertain” group.

386 The DoMore-v1-CRC classifier is technically simple to apply and can be delivered at
387 pathology laboratories everywhere. Although training the networks was resource demanding,
388 new patients can be assayed in a few minutes using consumer hardware.

389 Clinically, the marker will inform discussion with patients with stage II and III CRC on the
390 pros and cons of different adjuvant treatment options. Although the number of drugs used in
391 the adjuvant setting is limited to fluoropyrimidines ± oxaliplatin, recent data demonstrate that
392 three months treatment achieves approximately the same survival outcomes as six months for
393 the majority of stage III patients, while high risk patients (pT4 and pN2) might benefit from
394 prolonged therapy.^{29,30} It would be reasonable to hypothesise that stage III patients identified
395 as poor prognosis by the DoMore-v1-CRC classifier could benefit from prolonged
396 combination chemotherapy with oxaliplatin, or even consider experimental therapy
397 combining fluoropyrimidine + oxaliplatin + irinotecan as their high risk of cancer-specific
398 death should positively skew the risk-benefit ratio of more aggressive treatments (figures 2D
399 and F). At the other end, stage III patients with DoMore-v1-CRC good prognosis, the great
400 majority of whom are pN1, have very good survival with single-agent capecitabine (figure
401 2E), and good prognosis stage II patients have a very high chance of surgical cure, potentially
402 eliminating the need for adjuvant treatment.

403 We plan to undertake prospective adjuvant trials stratifying patients into different prognostic
404 groups using the DoMore-v1-CRC biomarker and randomising patients into observation, low
405 intensity and high intensity regimes depending on relative risk score. However, the currently
406 available data may also be used by clinicians and patients to make joint and more informed
407 decisions on adjuvant chemotherapy choices, as the proportional reduction in the HRs for
408 recurrence and death from CRC following adjuvant treatment is remarkably consistent at 20%

409 across most well-designed clinical trials, thus translating into quite different absolute survival
410 improvements for low and high risk subgroups.

411 Limitation of this study include that the DoMore-v1-CRC marker has not yet been tested
412 prospectively in clinical settings, and although we are planning a clinical trial with
413 randomisation, we at present only know the outcome of thorough retrospective testing. The
414 test and validation indicate good transferability between populations, but there are still
415 challenges related to standardisation, as illustrated by the differences between the tested
416 scanners. Differences between laboratories may also be seen for sample handling procedures,
417 and this is why the introduction into the clinic is suggested to be through services performed
418 at specialised laboratories. A well-known disadvantage of deep learning is its black-box
419 nature. The DoMore-v1-CRC marker is related to histological grading, but the marker is still
420 using small-scale features of the histological images with unknown biological correlates.
421 In summary, it has been possible to develop a clinically useful prognostic marker using deep
422 learning allied to digital scanning of conventional H&E stained, FFPE tumour tissue sections.
423 The assay has been extensively evaluated in large, independent patient populations, correlates
424 with and outperforms established molecular and morphological prognostic markers, gives
425 consistent results across tumour and nodal stage, and can potentially be used by clinicians to
426 improve decision making over adjuvant treatment choices.

427

428 **Contributors**

429 OJS, SDR, AK, TSH, KL, FA, DJK, and HED designed the study. HAA, JAN, AN, NAS, IT,
430 RK, MN, and DJK collected the samples and acquired the image data. MP, INF, ED, DNC,
431 AN, NAS, IT, RK, MN, and DJK provided clinical/pathological data and interpretations. OJS,
432 SDR, and JM performed the machine learning. AK performed the statistical analyses. OJS,
433 SDR, AK, TSH, KL, DJK, and HED interpreted the data and analyses. All authors vouch for

434 the data, analyses, and interpretations. OJS, SDR, AK, TSH, KL, DJK, and HED wrote the
435 first draft of the manuscript, and all authors reviewed, contributed to, and approved the
436 manuscript.

437

438 **Declaration of interests**

439 OJS, TSH, KL, JM, and HED report filing of a patent application entitled “Histological image
440 analysis” with International Patent Application Number PCT/EP2018/080828. The University
441 of Oxford (to DJK) received educational grants from Roche to support the QUASAR 2 trial
442 and from Merck to support the VICTOR trial. All other authors declare no competing
443 interests.

444

445 **Acknowledgements**

446 We thank Akershus University Hospital for access to their patient material, National Institute
447 for Health Research for funding support to Marco Novelli through Biomedical Research
448 Centres, Paul Callaghan for animating the appendix video, Marian Seiergren for creating
449 figure 1 and assembling figure 2, the laboratory and technical personnel at the Institute for
450 Cancer Genetics and Informatics for assistance, and the reviewers for valuable suggestions.
451 We also would like to thank the participating centres in the VICTOR and QUASAR 2 trials as
452 well as the staff at Akershus University Hospital, Aker University Hospital and the
453 Gloucestershire hospitals contributing to the Gloucester Colorectal Cancer Study, and last, but
454 not least all participating patients for making this study possible.

455

456 **References**

- 457 1. La Thangue NB, Kerr DJ. Predictive biomarkers: a paradigm shift towards
458 personalized cancer medicine. *Nat Rev Clin Oncol* 2011; **8**: 587–96.
- 459 2. Van Allen EM, Wagle N, Stojanov P, et al. Whole-exome sequencing and clinical
460 interpretation of formalin-fixed, paraffin-embedded tumor samples to guide precision cancer
461 medicine. *Nat Med* 2014; **20**: 682–88.
- 462 3. Moscow JA, Fojo T, Schilsky RL. The evidence framework for precision cancer
463 medicine. *Nat Rev Clin Oncol* 2018; **15**: 183–92.
- 464 4. Karapetis CS, Khambata-Ford S, Jonker DJ, et al. K-ras mutations and benefit from
465 cetuximab in advanced colorectal cancer. *N Engl J Med* 2008; **359**: 1757–65.
- 466 5. Kerr DJ, Shi Y. Biological markers: Tailoring treatment and trials to prognosis. *Nat*
467 *Rev Clin Oncol* 2013; **10**: 429–30.
- 468 6. Hutchins G, Southward K, Handley K, et al. Value of mismatch repair, KRAS, and
469 BRAF mutations in predicting recurrence and benefits from chemotherapy in colorectal
470 cancer. *J Clin Oncol* 2011; **29**: 1261–70.
- 471 7. Salazar R, Roepman P, Capella G, et al. Gene expression signature to improve
472 prognosis prediction of stage II and III colorectal cancer. *J Clin Oncol* 2011; **29**: 17–24.
- 473 8. Gray RG, Quirke P, Handley K, et al. Validation study of a quantitative multigene
474 reverse transcriptase-polymerase chain reaction assay for assessment of recurrence risk in
475 patients with stage II colon cancer. *J Clin Oncol* 2011; **29**: 4611–19.
- 476 9. QUASAR Collaborative Group. Comparison of fluorouracil with additional
477 levamisole, higher-dose folinic acid, or both, as adjuvant chemotherapy for colorectal cancer:
478 a randomised trial. *Lancet* 2000; **355**: 1588–96.
- 479 10. QUASAR Collaborative Group. Adjuvant chemotherapy versus observation in
480 patients with colorectal cancer: a randomised study. *Lancet* 2007; **370**: 2020–29.

- 481 11. Andre T, Boni C, Navarro M, et al. Improved overall survival with oxaliplatin,
482 fluorouracil, and leucovorin as adjuvant treatment in stage II or III colon cancer in the
483 MOSAIC trial. *J Clin Oncol* 2009; **27**: 3109–16.
- 484 12. Andre T, de Gramont A, Vernerey D, et al. Adjuvant Fluorouracil, Leucovorin, and
485 Oxaliplatin in Stage II to III Colon Cancer: Updated 10-Year Survival and Outcomes
486 According to BRAF Mutation and Mismatch Repair Status of the MOSAIC Study. *J Clin*
487 *Oncol* 2015; **33**: 4176–87.
- 488 13. Sinicrope FA. DNA mismatch repair and adjuvant chemotherapy in sporadic colon
489 cancer. *Nat Rev Clin Oncol* 2010; **7**: 174–77.
- 490 14. Mouradov D, Domingo E, Gibbs P, et al. Survival in stage II/III colorectal cancer is
491 independently predicted by chromosomal and microsatellite instability, but not by specific
492 driver mutations. *Am J Gastroenterol* 2013; **108**: 1785–93.
- 493 15. Danielsen HE, Hveem TS, Domingo E, et al. Prognostic markers for colorectal cancer:
494 estimating ploidy and stroma. *Ann Oncol* 2018; **29**: 616–23.
- 495 16. Kerr RS, Love S, Segelov E, et al. Adjuvant capecitabine plus bevacizumab versus
496 capecitabine alone in patients with colorectal cancer (QUASAR 2): an open-label, randomised
497 phase 3 trial. *Lancet Oncol* 2016; **17**: 1543–57.
- 498 17. Bondi J, Husdal A, Bukholm G, Nesland JM, Bakka A, Bukholm IR. Expression and
499 gene amplification of primary (A, B1, D1, D3, and E) and secondary (C and H) cyclins in
500 colon adenocarcinomas and correlation with patient outcome. *J Clin Pathol* 2005; **58**: 509–14.
- 501 18. Petersen VC, Baxter KJ, Love SB, Shepherd NA. Identification of objective
502 pathological prognostic determinants and models of prognosis in Dukes' B colon cancer. *Gut*
503 2002; **51**: 65–69.

- 504 19. Mitchard JR, Love SB, Baxter KJ, Shepherd NA. How important is peritoneal
505 involvement in rectal cancer? A prospective study of 331 cases. *Histopathology* 2010; **57**:
506 671–79.
- 507 20. Midgley RS, McConkey CC, Johnstone EC, et al. Phase III randomized trial assessing
508 rofecoxib in the adjuvant setting of colorectal cancer: final results of the VICTOR trial. *J Clin*
509 *Oncol* 2010; **28**: 4575–80.
- 510 21. Sandler M, Howard A, Zhu M, Zhmoginov A, Chen L. MobileNetV2: Inverted
511 Residuals and Linear Bottlenecks. *2018 IEEE/CVF Conference on Computer Vision and*
512 *Pattern Recognition* 2018: 4510–20.
- 513 22. Kraus OZ, Ba JL, Frey BJ. Classifying and segmenting microscopy images with deep
514 multiple instance learning. *Bioinformatics* 2016; **32**: i52–i59.
- 515 23. Altman DG, McShane LM, Sauerbrei W, Taube SE. Reporting recommendations for
516 tumor marker prognostic studies (REMARK): explanation and elaboration. *BMC Med* 2012;
517 **10**: 51.
- 518 24. Luo W, Phung D, Tran T, et al. Guidelines for Developing and Reporting Machine
519 Learning Predictive Models in Biomedical Research: A Multidisciplinary View. *J Med*
520 *Internet Res* 2016; **18**: e323.
- 521 25. Ehteshami Bejnordi B, Veta M, Johannes van Diest P, et al. Diagnostic Assessment of
522 Deep Learning Algorithms for Detection of Lymph Node Metastases in Women With Breast
523 Cancer. *JAMA* 2017; **318**: 2199–210.
- 524 26. Coudray N, Ocampo PS, Sakellaropoulos T, et al. Classification and mutation
525 prediction from non-small cell lung cancer histopathology images using deep learning. *Nat*
526 *Med* 2018; **24**: 1559–67.
- 527 27. Bychkov D, Linder N, Turkki R, et al. Deep learning based tissue analysis predicts
528 outcome in colorectal cancer. *Sci Rep* 2018; **8**: 3395.

- 529 28. Mobadersany P, Yousefi S, Amgad M, et al. Predicting cancer outcomes from
530 histology and genomics using convolutional networks. *Proc Natl Acad Sci U S A* 2018; **115**:
531 E2970–E79.
- 532 29. Grothey A, Sobrero AF, Shields AF, et al. Duration of Adjuvant Chemotherapy for
533 Stage III Colon Cancer. *N Engl J Med* 2018; **378**: 1177–88.
- 534 30. Iveson TJ, Kerr RS, Saunders MP, et al. 3 versus 6 months of adjuvant oxaliplatin-
535 fluoropyrimidine combination therapy for colorectal cancer (SCOT): an international,
536 randomised, phase 3, non-inferiority trial. *Lancet Oncol* 2018; **19**: 562–78.
- 537

538 **Figure Legends**

539

540 ***Figure 1: Pipeline of DoMore-v1-CRC classification***

541 Top: A whole-slide image (WSI) is segmented, and the segmented regions tiled at 40x
542 resolution and 10x resolution. For each resolution, the five trained models each produce one
543 score reflecting the probability of poor outcome. The average of those scores is the ensemble
544 score, one for 10x and one for 40x. If the ensemble score is above a certain threshold, the WSI
545 is classified as poor prognosis. The DoMore-v1-CRC class is determined by the agreement
546 between the two ensemble classifications. Bottom: The DoMore v1 network is comprised of a
547 representation network (MobileNetV2²¹), a pooling function (Noisy-AND²²), and a simple
548 fully-connected classification network. All components of the DoMore v1 network involve
549 trainable parameters, and the entire network is trained end-to-end. All tiles from a WSI are
550 processed by the representation network one by one, resulting in a collection of tile
551 representations. The pooling function reduces the representations into two numbers, which are
552 then processed by the classification network to produce the score outputted by the model.

553

554 **Figure 2: Kaplan-Meier analysis of cancer-specific survival by DoMore-v1-CRC class**
555 **evaluated on Aperio AT2 slide images in the QUASAR 2 validation cohort**
556 (A) The primary analysis; all patients evaluated with the pre-defined DoMore-v1-CRC
557 classifier. (B) A post-hoc analysis; all patients evaluated with the DoMore-v1-CRC classifier
558 variant with five categories. (C) A secondary analysis; stage II (equivalent to pN0) patients
559 evaluated with the pre-defined DoMore-v1-CRC classifier. (D) A secondary analysis; stage
560 III patients evaluated with the pre-defined DoMore-v1-CRC classifier. (E) A post-hoc
561 analysis; pN1 patients evaluated with the pre-defined DoMore-v1-CRC classifier. (F) A post-
562 hoc analysis; pN2 patients evaluated with the pre-defined DoMore-v1-CRC classifier.
563

Table 1: Patient characteristics in the training, tuning, test and validation cohorts

	Group	Training cohort (N=828)	Tuning cohort (N=1645)	Test cohort (N=920)	Validation cohort (N=1122)
Age, years		69 (61-75)	70 (61-77)	71 (64-78)	65 (59-71)
Sex					
	Female	402 (51%)	689 (42%)	421 (46%)	477 (43%)
	Male	426 (49%)	956 (58%)	499 (54%)	645 (57%)
Stage					
	I	101 (12%)	102 (6%)	70 (8%)	
	II	317 (38%)	797 (48%)	354 (38%)	402 (36%)
	III	410 (50%)	746 (45%)	496 (54%)	720 (64%)
pN stage					
	pN0	415 (50%)	891 (54%)	425 (46%)	402 (36%)
	pN1	241 (29%)	492 (30%)	258 (28%)	508 (45%)
	pN2	167 (20%)	239 (15%)	237 (26%)	183 (16%)
	Missing	5 (1%)	23 (1%)	0 (0%)	29 (3%)
pT stage					
	pT1	26 (3%)	30 (2%)	6 (1%)	17 (2%)
	pT2	110 (13%)	137 (8%)	65 (7%)	71 (6%)
	pT3	464 (56%)	1034 (63%)	411 (45%)	582 (52%)
	pT4	223 (27%)	423 (26%)	437 (48%)	404 (36%)
	Missing	5 (1%)	21 (1%)	1 (0%)	48 (4%)
Histological grade					
	1	77 (9%)	196 (12%)	134 (15%)	45 (4%)
	2	568 (69%)	1151 (70%)	489 (53%)	846 (75%)
	3	178 (21%)	280 (17%)	297 (32%)	168 (15%)
	Missing	5 (1%)	18 (1%)	0 (0%)	63 (6%)
Location					
	Rectum	222 (27%)	457 (28%)	311 (34%)	165 (15%)
	Distal colon	262 (32%)	533 (32%)	280 (30%)	451 (40%)
	Proximal colon	307 (37%)	505 (31%)	329 (36%)	453 (40%)
	Missing	37 (4%)	150 (9%)	0 (0%)	53 (5%)
Adjuvant treatment					
	No	467 (56%)	826 (50%)	538 (58%)	0 (0%)
	Chemotherapy	173 (21%)	397 (24%)	51 (6%)	1122 (100%)
	Radiotherapy	11 (1%)	6 (0%)	14 (2%)	0 (0%)
	Chemo- and radiotherapy	3 (0%)	9 (1%)	3 (0%)	0 (0%)
	Missing	174 (21%)	407 (25%)	314 (34%)	0 (0%)
Follow-up time, years		6.4 (1.7-8.2)	4.0 (2.2-5.2)	2.4 (1.0-4.6)	4.6 (3.3-5.1)

Data are median (IQR) or number (%). IQR=interquartile range.

Table 2: Multivariable cancer-specific survival analyses in the validation cohort; the multivariable model included the DoMore-v1-CRC class evaluated on Aperio AT2 slide images, and established prognostic markers that were significant in the corresponding stage-specific univariable analyses in the validation cohort

	Group	Stage II and III		Stage II		Stage III	
		HR (95% CI)	p	HR (95% CI)	p	HR (95% CI)	p
DoMore-v1-CRC			<0.0001		0.028		0.0001
	Good prognosis	ref.		ref.		ref.	
	Uncertain	1.56 (0.92-2.65)		1.22 (0.35-4.24)		2.14 (1.15-3.99)	
	Poor prognosis	3.04 (2.07-4.47)		2.71 (1.25-5.86)		2.95 (1.81-4.82)	
pN stage			<0.0001				<0.0001
	pN0	ref.					
	pN1	1.84 (1.13-2.98)				ref.	
	pN2	5.94 (3.71-9.52)				3.31 (2.14-5.13)	
pT stage			0.0058				0.014
	pT1	NA				NA	
	pT2	1.86 (0.90-3.86)				1.68 (0.64-4.45)	
	pT3	ref.				ref.	
	pT4	1.75 (1.22-2.51)				2.07 (1.33-3.22)	
Lymphatic invasion	Yes	1.66 (1.07-2.56)	0.023			1.98 (1.20-3.28)	0.0079
Venous vascular invasion	Yes	1.07 (0.76-1.51)	0.71			0.98 (0.64-1.52)	0.94
Sidedness	Right					1.09 (0.70-1.70)	0.69
BRAF	Mutated					1.39 (0.81-2.40)	0.24

Ref.=reference; NA=not available

Table 3: Associations between the DoMore-v1-CRC class evaluated on Aperio AT2 slide images and different patient characteristics in the validation cohort

	Group	DoMore-v1-CRC good prognosis	DoMore-v1-CRC uncertain	DoMore-v1-CRC poor prognosis	Spearman's correlation	
		(N=704)	(N=136)	(N=270)	ρ (95% CI)	p
Age (continuous), years		64 (58-71)	65 (60-71)	66 (60-72)	0.07 (0.01 to 0.13)	0.024
Age (dichotomous), years					0.03 (-0.03 to 0.09)	0.38
	≤72	568 (81%)	112 (82%)	209 (77%)		
	>72	136 (19%)	24 (18%)	61 (23%)		
Sex					-0.02 (-0.08 to 0.04)	0.59
	Female	297 (42%)	53 (39%)	122 (45%)		
	Male	407 (58%)	83 (61%)	148 (55%)		
Stage					0.04 (-0.02 to 0.10)	0.20
	II	261 (37%)	48 (35%)	88 (33%)		
	III	443 (63%)	88 (65%)	182 (67%)		
Stage with substage					0.15 (0.09 to 0.21)	<0.0001
	IIA	143 (21%)	19 (14%)	28 (11%)		
	IIB	110 (16%)	27 (20%)	54 (21%)		
	IIIA	67 (10%)	2 (2%)	6 (2%)		
	IIIB	269 (40%)	51 (38%)	104 (41%)		
	IIIC	83 (12%)	34 (26%)	64 (25%)		
pN stage					0.10 (0.04 to 0.16)	0.0008
	pN0	261 (38%)	48 (36%)	88 (33%)		
	pN1	339 (50%)	53 (39%)	111 (42%)		
	pN2	83 (12%)	34 (25%)	64 (24%)		
pT stage					0.26 (0.21 to 0.32)	<0.0001
	pT1	15 (2%)	0 (0%)	2 (1%)		
	pT2	61 (9%)	3 (2%)	6 (2%)		
	pT3	402 (60%)	75 (56%)	100 (39%)		
	pT4	194 (29%)	56 (42%)	148 (58%)		
Lymphatic invasion					0.04 (-0.02 to 0.10)	0.20
	No	599 (91%)	122 (92%)	220 (87%)		
	Yes	62 (9%)	10 (8%)	33 (13%)		
Venous vascular invasion					0.05 (-0.01 to 0.11)	0.11
	No	409 (61%)	74 (56%)	145 (56%)		
	Yes	257 (39%)	58 (44%)	112 (44%)		
Histological grade					0.14 (0.08 to 0.20)	<0.0001
	1	27 (4%)	7 (6%)	8 (3%)		
	2	565 (85%)	88 (69%)	186 (74%)		
	3	76 (11%)	32 (25%)	59 (23%)		
Location					0.15 (0.09 to 0.21)	<0.0001
	Rectum	118 (18%)	21 (16%)	23 (9%)		
	Distal colon	301 (45%)	46 (35%)	100 (38%)		
	Proximal colon	246 (37%)	64 (49%)	138 (53%)		
Sidedness					0.14 (0.08 to 0.20)	<0.0001
	Left	419 (63%)	67 (51%)	123 (47%)		
	Right	246 (37%)	64 (49%)	138 (53%)		
KRAS					-0.06 (-0.12 to 0.00)	0.069
	Wild-type	410 (65%)	86 (73%)	169 (70%)		
	Mutated	224 (35%)	32 (27%)	73 (30%)		
BRAF					0.22 (0.16 to 0.28)	<0.0001
	Wild-type	588 (93%)	89 (75%)	190 (77%)		

	Mutated	47 (7%)	29 (25%)	56 (23%)		
Microsatellite instability					-0.10 (-0.16 to -0.04)	0.0018
	Yes	66 (10%)	26 (21%)	40 (16%)		
	No	595 (90%)	99 (79%)	213 (84%)		
Follow-up time, years		4.8 (3.7-5.1)	4.9 (3.1-5.1)	4.1 (2.8-5.1)	-0.10 (-0.16 to -0.04)	0.0006

Data are median (IQR) or number (%). IQR=interquartile range.



