

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15

Triage-driven diagnosis for early detection of esophageal cancer using deep learning

Marcel Gehrung,^{1,2} Mireia Crispin-Ortuzar,¹ Adam G. Berman,¹
Maria O'Donovan,^{3,4} Rebecca C. Fitzgerald,^{3,a} Florian Markowetz^{1,a*}

¹Cancer Research UK Cambridge Institute, University of Cambridge, UK

²The Alan Turing Institute, London, UK

³MRC Cancer Unit, University of Cambridge, UK

⁴Department of Pathology, Cambridge University Hospitals NHS Trust, UK

^aThese authors share senior authorship.

*To whom correspondence should be addressed; E-mail: florian.markowetz@cruk.cam.ac.uk

Abstract

6
7
8
9
10
11
12
13
14
15

Deep learning methods have been shown to achieve excellent performance on diagnostic tasks, but it is still an open challenge how to optimally combine them with expert knowledge and existing clinical decision pathways. This question is particularly important for the early detection of cancer, where high volume workflows might potentially benefit substantially from automated analysis. Here, we present a deep learning framework to analyse samples of the Cytosponge®-TFF3 test, a minimally invasive alternative to endoscopy, for detecting Barrett's Esophagus, the main precursor of esophageal cancer. We trained and independently validated the framework on data from two clinical trials, analysing a combined total of 4,662 pathology slides from 2,331

16 **patients. Our approach exploits screening patterns of expert gastrointestinal**
17 **pathologists and established decision pathways to define eight triage classes of**
18 **varying priority for manual expert review. By substituting manual review with**
19 **automated review in low-priority classes, we can reduce pathologist workload**
20 **by up to 66% while matching the diagnostic performance of expert patholo-**
21 **gists. These results lay the foundation for tailored, semi-automated decision**
22 **support systems embedded in clinical workflows.**

23 **Introduction**

24 Early detection of cancer often leads to better survival (1), because pre-malignant lesions and
25 early stage tumors can be more effectively treated (2). Most pre-malignant lesions amenable to
26 early detection rely on targeted sampling and show only minor tissue changes on pathology as-
27 sessment (3–5). In addition, pathology procedures often involve laborious and time-consuming
28 steps which can lead to errors and adversely affect patient care (6). Recent developments in Ar-
29 tificial Intelligence (AI) have achieved excellent performance on diagnostic tasks (7–9). How-
30 ever, understanding how these techniques can be integrated into clinical workflows most effi-
31 ciently and to assess the actual benefits they bring remains a challenge. The design of a clinical
32 decision support system needs to balance its performance against workload reduction and po-
33 tential economic impact. Replacing pathologists entirely could lead to substantial workload
34 reduction, but such an approach would only be viable if performance remains comparable to
35 that of human experts. Between a fully automated approach and the *status quo* of fully manual
36 review lies a semi-automated approach, which uses computational methods to triage patients
37 and only presents pathologists with difficult cases. A semi-automated approach will not reduce
38 workload as much as a fully automated approach, but its performance benefits from existing
39 expert knowledge and heuristics. Here we present such a semi-automated triage system using

40 deep learning for the early detection of esophageal cancer.

41 Esophageal cancer is the sixth most common cause for cancer related deaths (10). Patients
42 usually present at an advanced stage with dysphagia and weight loss, and the 5-year overall
43 survival of esophageal adenocarcinoma (EAC), one of two pathological subtypes, is 13% (11).
44 EAC can arise from a precursor lesion called Barrett's Esophagus (BE) (12, 13), providing an
45 effective starting point for early detection. BE occurs in patients with Gastresophageal Reflux
46 Disease (GERD), a digestive disorder where acid and bile from the stomach return into the
47 esophagus leading to heartburn symptoms. In Western countries, 10 to 15% of the adult popu-
48 lation are affected by GERD (14) and, therefore, at an increased risk of having BE. The pathog-
49 nomonic feature of BE is intestinal metaplasia (IM), a process whereby the stratified squamous
50 epithelial lining localized in the lower esophagus is replaced with columnar epithelium con-
51 taining goblet cells (15, 16). The conventional diagnosis of BE requires an invasive endoscopic
52 procedure of the upper gastrointestinal tract. However, there is no routine endoscopic screening
53 of the GERD population and thus the vast majority of BE patients are undiagnosed (14).

54 Cytosponge-TFF3 is a non-endoscopic, minimally invasive diagnostic test for BE (17–19).
55 It is a cell collection device consisting of a compressed sponge on a string inside a gelatin
56 capsule. The capsule is swallowed by the patient and the gelatin dissolves in the stomach after
57 a few minutes, allowing the sponge to expand. The sponge is then withdrawn from the stomach
58 by the attached string, sampling superficial epithelial cells from the top of the stomach, the
59 esophagus, and the oropharynx (Figure 1a). Therefore, the cellular composition of the sample
60 is dominated by squamous cells, gastric columnar epithelium, and respiratory epithelium as well
61 as any Barrett's cells, if present. Following removal, the device is placed in a container with
62 preservative solution and the sampled cells are processed, embedded in paraffin and stained with
63 Hematoxylin & Eosin (H&E) as well as immunohistochemically stained with Trefoil Factor 3
64 (TFF3) (20). H&E stains allow the identification and quantification of cellular phenotypes,

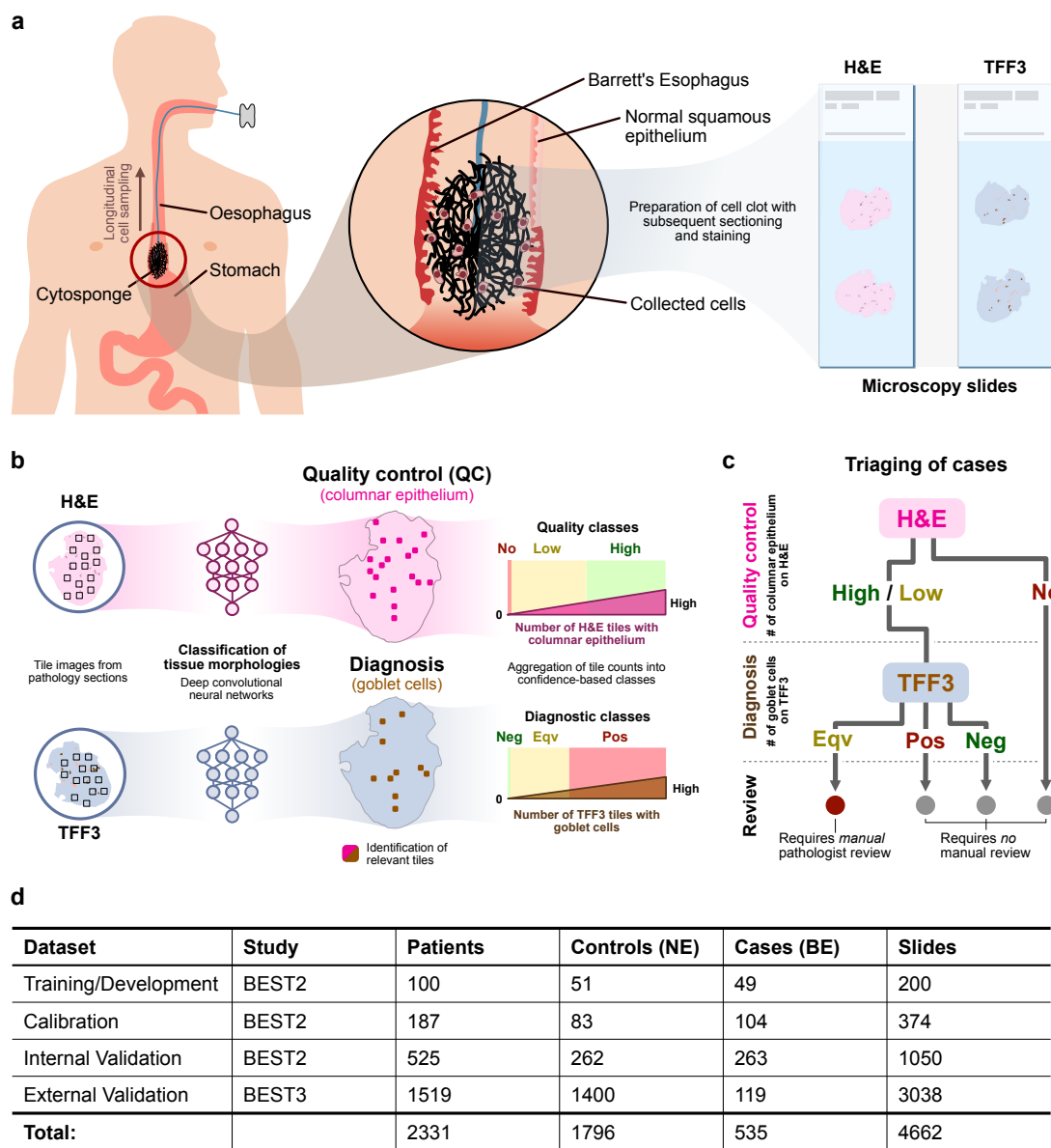


Figure 1: Cytosponge procedure with conceptual patient triage scheme and data summary. **a** During withdrawal the sponge samples superficial epithelial cells from the top of the stomach and the esophagus. These cells are processed into a cell block, then sectioned and stained with Hematoxylin & Eosin (H&E) and Trefoil Factor 3 (TFF3). **b** Convolutional neural networks, trained on an independent training dataset, are used for inference of H&E and TFF3 stains. The resulting tile maps are analysed for relevant regions (columnar epithelium on H&E and goblet cells on TFF3 stain) and aggregated into quality control and diagnostic classes based on tile detections. **c** Quality and diagnostic classes are mapped to a conceptualised pathway for sample stratification. The review layer (bottom) describes to what extent a human pathologist has to review the microscopy slides. (Pos = Positive, Neg = Negative) **d** Overview of data used in this study.

65 which is critical for quality control. TFF3 is over-expressed in mucin-producing goblet cells
66 which are a key feature of BE. TFF3 also functions as a protector of the mucosa from insults,
67 stabilizes the mucus layer, and promotes healing of the epithelium (21). TFF3 stains allow the
68 identification and quantification of goblet cells, which are indicative of IM. Therefore, TFF3 is
69 the key diagnostic biomarker for BE (20).

70 The Cytosponge-TFF3 approach has profound and well-tested clinical significance. It of-
71 fers, with substantial clinical trial data underpinning its efficacy, a long-awaited diagnostic al-
72 ternative to endoscopy (BEST1 (17), BEST2 (18), BEST3 (22)). The BEST3 study found that
73 the Cytosponge-TFF3 test had in excess of a 10-fold increase in detection of Barrett’s compared
74 to usual clinical care in which patients with heartburn receive medication and an endoscopy if
75 deemed necessary. This performance makes the Cytosponge a major advance in patient man-
76 agement. The BEST3 study also concluded that the pathology assessment is a major bottleneck
77 for scaling the test to large patient populations. Since the analysis of Cytosponge-TFF3 pathol-
78 ogy slides is a very laborious process due to the large amount of sampled cellular material. It
79 comprises several time-consuming tasks such as assessing the amount of sampled material and
80 checking the presence of gastric-type columnar epithelium to confirm that the capsule reached
81 the stomach, followed by assessment for the presence of goblet cells indicative of BE. Though
82 effective, the laboriousness of this process gives rise to a major opportunity for a clinical deci-
83 sion support system to improve analysis and scalability of the Cytosponge-TFF3 test.

84 Here, we use a deep learning approach for quality control and diagnosis of pathology slides
85 for the Cytosponge-TFF3 test (Figure 1b). We propose a triage-driven approach, which retains
86 diagnostic accuracy by leveraging the decision-making rules of expert gastrointestinal patholo-
87 gists (Figure 1c). We train, calibrate, and internally validate our approach on data of the BEST2
88 multi-centre clinical trial (18) and externally validate it in an independent cohort from the recent
89 BEST3 multi-centre trial (22) (Figure 1d). Additionally, we explore in a simulation study how

90 well our results generalise to more general populations.

91 **Results**

92 **Deep learning models achieve high performance for tile-level classifications**

93 The first step of our approach is based on the tile-level detection of different classes of cells
94 relevant for quality control and diagnosis of BE. For model development and internal valida-
95 tion, we used 812 Cytosponge-TFF3 patient samples with paired pathology and endoscopy data
96 from the BEST2 clinical case-control study (18). Samples were randomly divided into train-
97 ing/development (n=100), calibration (n=187) and internal validation (n=525) sets (Figure 1d).
98 An additional independent dataset (n=1,519) from the BEST3 study was used for external vali-
99 dation of the developed approach.

100 Training sets of larger size did not improve tile-level accuracy (Figure S1). Training, cali-
101 bration, and validation sets were kept separate. Endoscopic as well as Cytosponge pathology
102 diagnoses were only unblinded after tile-wise tissue classification models were calibrated and
103 validated, respectively. All training slides were tessellated prior to training: For H&E we de-
104 rived 193,734 tiles from 100 slides and for TFF3 we derived 235,932 tiles from 100 slides (based
105 on the size of annotated areas, see Methods). All tiles were 200-by-200 μm and all labels were
106 taken from expert slide annotations.

107 For both quality control (H&E) and diagnostic (TFF3) tasks, we trained several state-of-the-
108 art networks (AlexNet (23), DenseNet (24), Inception v3 (25), ResNet-18 (26), SqueezeNet (27),
109 and VGG-16 (28)) and evaluated their performance on the development dataset. Using indi-
110 vidual tiles, we compared tile-level precision and recall for classifying columnar epithelium
111 using the presence of gastric-type cells (on H&E) and positive goblet cells (on TFF3) (Ta-
112 ble S1, description in Methods): For gastric-type columnar epithelium, VGG-16, DenseNet and
113 Inception v3 achieved the highest recalls (0.950, 0.947, 0.940, respectively) with consistent pre-

114 cisions (0.843, 0.865, 0.857). For goblet cells, VGG-16, Inception v3, and ResNet-18 achieved
115 the highest recalls (0.919, 0.919, 0.912) with consistent precisions (0.856, 0.856, 0.827). Ex-
116 amples for whole slide images classified positive and negative for quality control and diagnosis
117 are shown in Figure 2a. We also observed a relationship between the tile-level results and the
118 complexity of the applied architectures (Table S1).

119 **Saliency maps agree with pathologist criteria for classification of tissue tiles**

120 To understand which characteristics of the tile images were relevant to our models' classifica-
121 tions, we generated saliency maps using Gradient-weighted Class Activation Mapping (Grad-
122 CAM) (29). These maps highlight the local regions of an image most relevant to a model's iden-
123 tification of a particular class. We generated saliency maps for classes in one H&E-based model
124 (VGG-16) and one TFF3-based model (VGG-16) (Figure 2b). For the gastric-type columnar ep-
125 ithelium class of the H&E-based model, the saliency maps highlight gastric cells by both the
126 linear organisation of their nuclei as well as the presence of a straight border between the cells
127 and the lumen. For the positive class of the TFF3-based model, we found that the saliency
128 maps highlighted the mucin-containing goblet cells that characterise IM with high precision. In
129 addition to the three representative examples in Figure 2b, we compared landmarks selected by
130 an expert pathologist with tile images and respective saliency maps (Figure S2). The saliency
131 maps confirm that the models learned features are similar to those used by pathologists to iden-
132 tify different tissue classes.

133 **Fully automated approach shows suboptimal performance**

134 Tile-level classifications were aggregated into patient-level classifications using tile counts above
135 thresholds determined by the specificity of expert pathologists on the calibration cohort (Meth-
136 ods, Table S2, Figure S4). We then performed Receiver Operating Characteristics (ROC) anal-
137 ysis with matched Cytosponge pathology and endoscopy ground truth on the internal validation

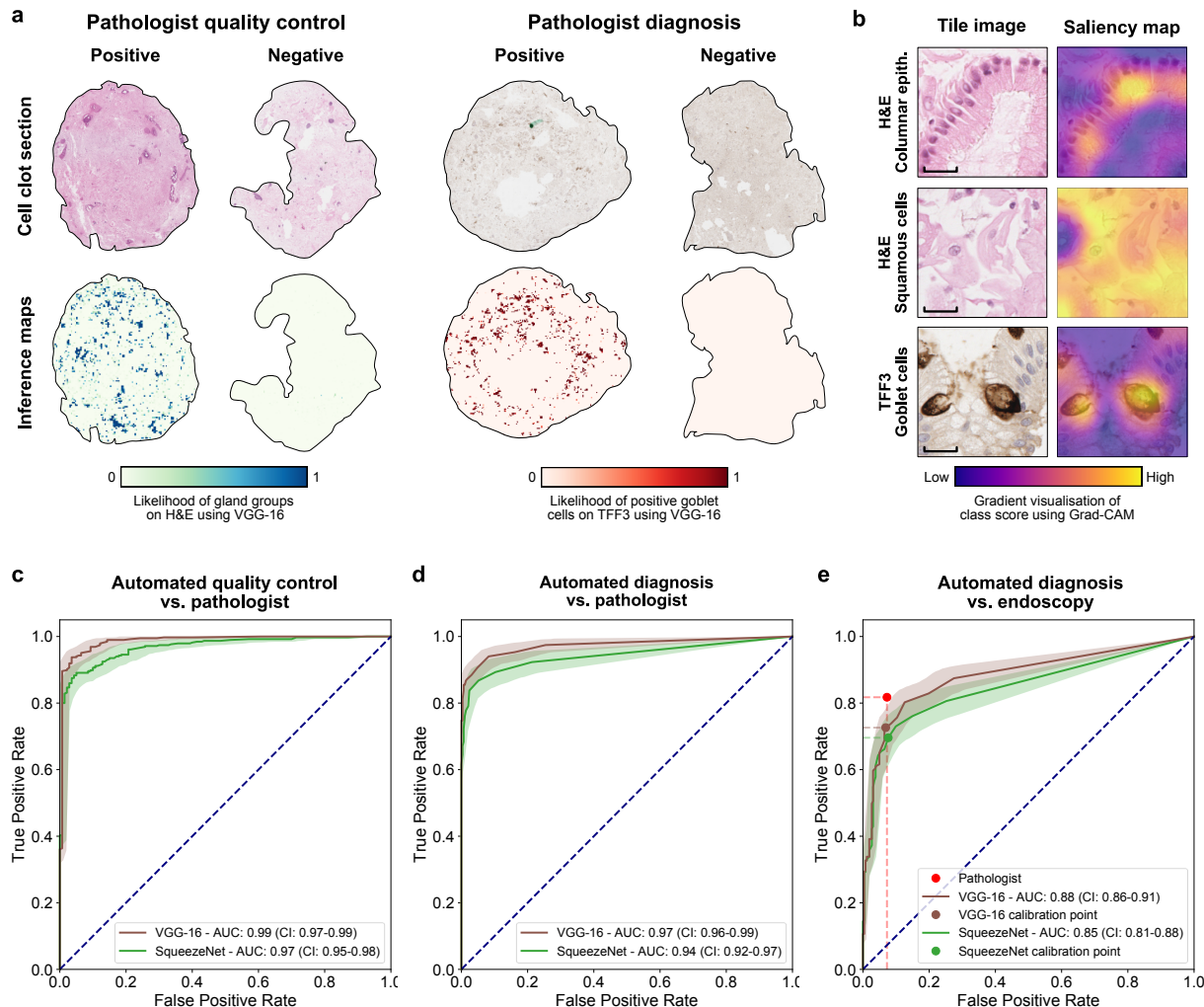


Figure 2: Tile and patient-level classification of Cytosponge-TFF3 samples. **a** Examples of tile-level inference maps for samples, which were classified by a pathologist as positive and negative for quality control (H&E) and diagnosis (TFF3), respectively. **b** Comparison of two tile images from H&E and one tile image from TFF3 with their respective Grad-CAM saliency maps. *Top*: Columnar epithelium (H&E) of gastric type with clear focus on columnar arrangement in saliency map. *Middle*: Squamous cells (H&E) with distributed focus in saliency map. *Bottom*: TFF3-positive goblet cells with localisation in saliency maps. Scale bar = 100 μ m. **c** ROC-AUC internal validation cohort analysis of automated tile counts of columnar epithelium on H&E with pathologist ground truth. **d** ROC-AUC internal validation cohort analysis of automated tile counts of positive goblet cells on TFF3 with pathologist ground truth. **e** ROC-AUC internal validation cohort analysis of pathologist and automated tile counts of positive goblet cells on TFF3 with endoscopy ground truth (BE patients defined according to the Prague criteria (Methods) with confirmed IM on biopsy). Pathologist performance: Sensitivity: 81.749% (CI: 76.597% - 85.951%), Specificity: 92.748% (CI: 89.444% - 95.769%). VGG-16 performance (calibration point determined as in Methods): Sensitivity: 72.624% (CI: 67.424% - 78.213%), Specificity: 93.130% (CI: 90.038% - 96.133%). SqueezeNet performance (calibration point, Methods): Sensitivity: 69.582% (CI: 63.915% - 75.522%), Specificity: 92.366% (88.468% - 95.518%). CI = 95% bootstrap confidence interval. Shaded areas show CIs.

138 cohort (Figure 2c-e).

139 First, the patient-level scores were compared against the binary Cytosponge-TFF3 ground
140 truth by the pathologist on the internal validation set. For quality control, VGG-16 ranked high-
141 est for detecting columnar epithelium in H&E stains (ROC-AUC: 0.99 (CI 95%: 0.98 - 0.99)).
142 SqueezeNet, the least complex architecture we trained, ranked lowest (ROC-AUC: 0.97 (CI
143 95%: 0.95 - 0.98), Figure 2c). For diagnosis, VGG-16 ranked highest for detecting goblet cells
144 in TFF3 stains (ROC-AUC: 0.97 (CI 95%: 0.96 - 0.99), Figure 2d). Again, SqueezeNet ranked
145 lowest (ROC-AUC: 0.94 (CI 95%: 0.92 - 0.96)). Confidence intervals were derived by boot-
146 strapping (Methods). Results for all architectures are presented in table S3, and fig. S5a/b. In
147 summary, for both quality control and diagnosis in comparison to Cytosponge-TFF3 pathology
148 ground truth, VGG-16 provided the highest performance, and SqueezeNet the lowest.

149 Next, patient-level counts were compared to endoscopy ground truth for detecting BE on
150 the internal validation set (Methods). This ground truth was defined according to the Prague
151 criteria (Methods) with confirmed IM on endoscopy biopsies (30). To calculate sensitivity and
152 specificity for the fully automated method on the internal validation cohort, we used operating
153 points determined on the calibration cohort (Table S2). VGG-16 ranked highest for detecting
154 patients with BE from TFF3 stains (ROC-AUC: 0.88 (CI 95%: 0.85 - 0.91), Sensitivity: 72.62%
155 (CI: 67.42% - 78.21%), Specificity: (93.13% (CI: 90.04% - 96.13%)), Figure 2e). SqueezeNet
156 ranked lowest for detecting patients with BE from TFF3 stains (ROC-AUC: 0.85 (CI 95%: 0.81
157 - 0.88), Sensitivity: 69.58% (CI: 63.92% - 75.52%), Specificity: 92.37% (88.47% - 95.52%),
158 Figure 2e). For comparison, the pathologists achieve a sensitivity of 81.7% (CI 95%: 77.4% -
159 86.5%) and a specificity of 92.7% (CI 95%: 89.6% - 95.6%). Performances of all architectures
160 are presented in table S3, and fig. S5c. In summary, results for the fully automated approach
161 on the internal validation cohort showed a loss of sensitivity of 9.1% for BE detection when
162 compared to an expert pathologist.

163 **Triage-driven approach selects patients for manual review**

164 We then explored whether a different modelling approach based on established decision path-
165 ways could boost performance. We developed a triage-driven, semi-automated approach as an
166 alternative to the fully automated approach described above. Both approaches use the same
167 patient-level aggregations as input, but their outputs are different: the fully automated approach
168 tries to directly mimic pathology assessment by classifying patients as positive or negative for
169 BE. In contrast, the triage approach defines different quality and diagnostic confidence classes
170 to select challenging patient samples for manual review. Although it cannot reduce workload as
171 much as a fully automated approach, a triage approach keeps sample stratification more inter-
172 pretable and transparent.

173 We first selected deep learning architectures and defined cut-offs for different quality and
174 diagnostic confidence classes based on thresholds determined by two expert observers on the
175 calibration cohort (Figure S6, Methods). For quality confidence classes, pathologists conclude
176 that the sponge reached the stomach if they observe columnar epithelial groups (18, 20). We en-
177 coded these subjective metrics in a quantitative scheme where the number of tiles detected with
178 gastric-type columnar epithelium on H&E were classified as no confidence, low confidence, or
179 high confidence (Figure S6a, Table S4). For diagnostic confidence classes, the number of tiles
180 detected with TFF3-positive goblet cells were classified as high confidence negative, low confi-
181 dence equivocal, or high confidence positive (Figure S6b, Table S4). On the internal validation
182 cohort, we observed a visual agreement between these confidence classes and pathology and
183 endoscopy ground truths (Figure 3, Table S5).

184 We then combined the quality and diagnostic classes into eight triage classes of varying
185 priority for manual review (Figure 4a). The relative priority of each class was determined by
186 expert pathologists: Cases with low confidence in sample quality (none or few columnar ep-
187 ithelium detected on H&E) or low confidence in diagnosis (few goblet cells detected on TFF3)

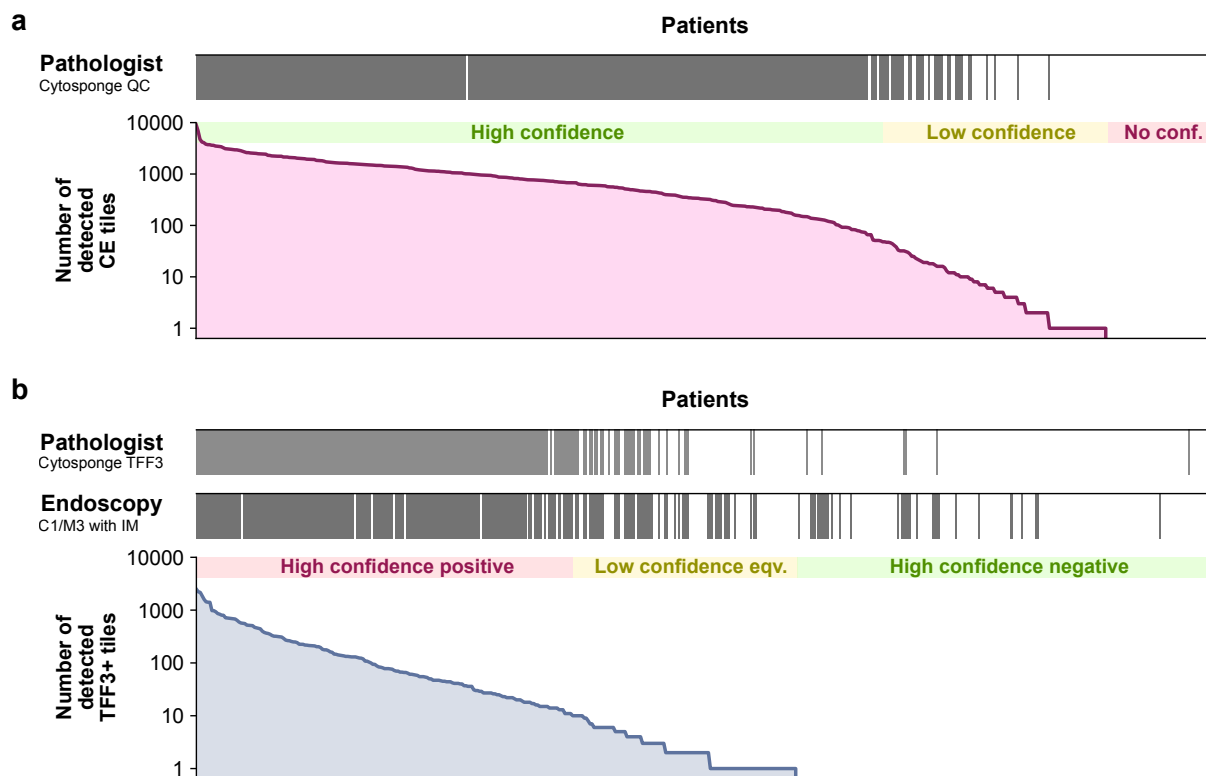


Figure 3: **Application of quality control and diagnostic confidence class scheme to internal validation cohort.** **a** Quality ground truth by pathologist from Cytosponge (top) compared with number of detected columnar epithelium (CE) tiles on H&E detected by VGG-16 (bottom). **b** Diagnosis ground truth by pathologist from Cytosponge (top), Endoscopy (with confirmed IM on biopsy) ground truth (middle) compared with number of detected TFF3-positive tiles on TFF3 detected by ResNet-18 (bottom) / eqv. = equivocal.

188 should be prioritised for human expert assessment over cases with high-confidence positive or
189 negative evidence. In our internal validation cohort, we find that only 13.0% of patients fall
190 into the triage classes with high priority (4 and 5), while 87.0% fall into the other six classes
191 (Figure 4a).

192 We next asked which classes can be substituted by automated review while retaining the
193 accuracy of full manual review by a human pathologist (sensitivity: 81.7%; specificity: 92.7%).
194 We applied a cumulative substitution scheme and started by substituting class 1 with automated
195 review, then classes 1 and 2, then classes 1, 2, and 3, and so on. In the validation cohort, we
196 found that sensitivity and specificity remain stable if classes 1, 2, and 3 are substituted, but
197 decrease with the substitution of class 4, 5, and 6 (Figure 4b). Repeating this procedure starting
198 with class 8 shows that sensitivity and specificity are stable if classes 8 or 7 are substituted, but
199 decrease with the substitution of classes 6, 5, and 4 (Figure 4c). These results show that five of
200 the eight classes (1, 2, 3, 7, 8) can be substituted by automated review while three classes (4,
201 5, 6) should be manually reviewed by a pathologist. This substitution scheme would result in
202 similar performance (sensitivity: 82.5% (CI 95%: 77.3% - 87.2%); specificity: 92.7% (CI 95%:
203 89.6% - 95.9%)) as fully manual review by a pathologist. These classes cover the majority
204 of patients (66.3% (CI 95%: 62.7% - 70.1%) in validation cohort) and triage-driven, semi
205 automated review would thus save 66% of the pathologists' workload (Methods) by enabling
206 them to focus on difficult cases while leaving easy cases for automated review.

207 **Simulation of varying cohort composition corroborates reduction in expected workload**

208 Our case-control cohort is not representative of a real-world population eligible for Cytosponge-
209 TFF3 testing. In our internal validation set we had a disease prevalence of 50.0%, while the
210 prevalence expected in a real-world population with GERD symptoms ranges from 3.0% to
211 7.5% (17, 31–33). Additionally, the allocation of samples to triage classes depends directly on

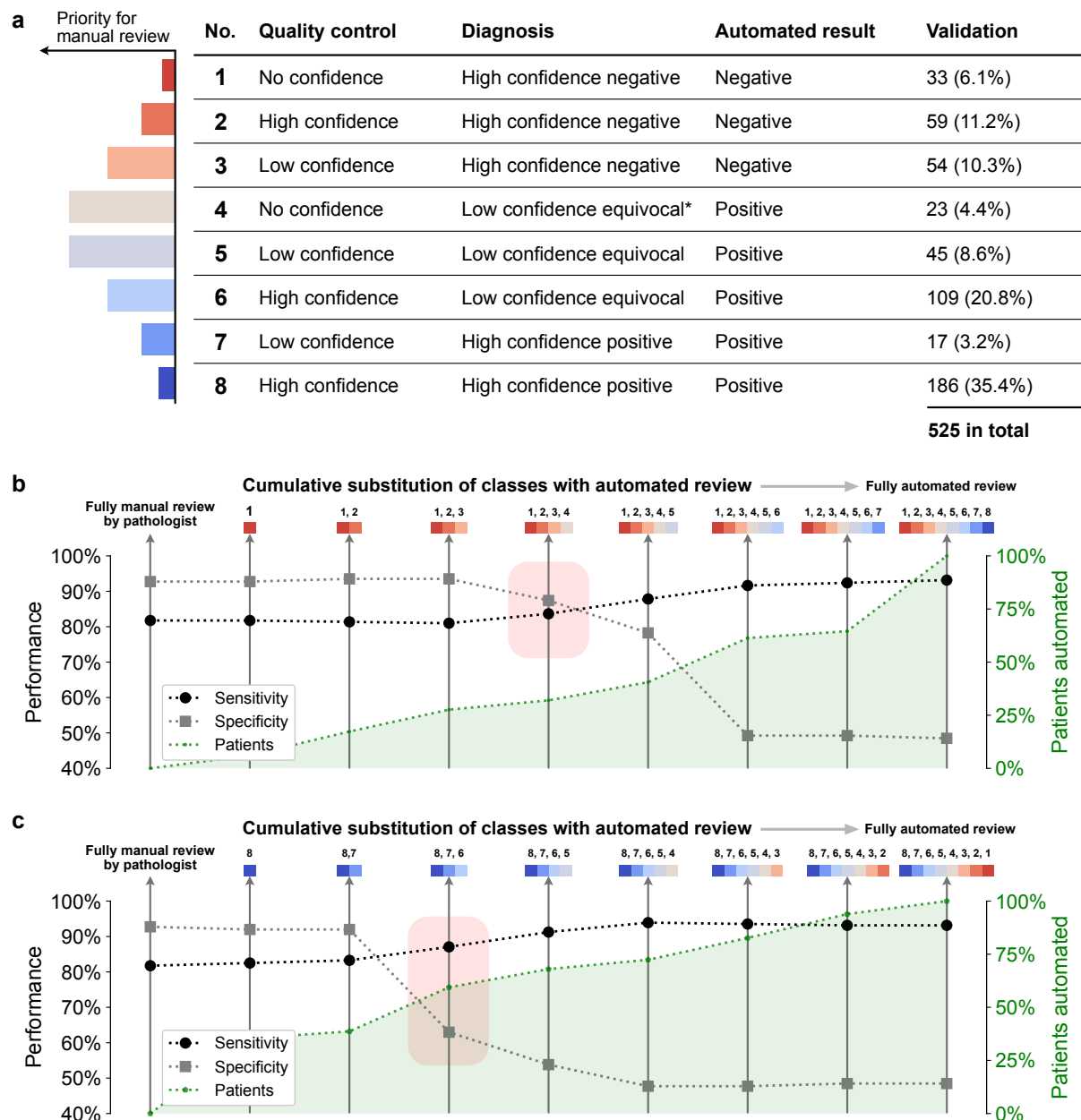


Figure 4: Triage-driven approach with incremental triage class substitution scheme on internal validation set

a Table of quality control and diagnosis classes. Each class has been assigned a qualitative priority for manual review. Column ‘Automated result’ refers to the label a sample would be assigned if all samples of this class were automatically reviewed. Asterisk (*): includes combination of no confidence (quality control) and high confidence positive (diagnosis) despite minimal likelihood of occurrence. **b** Cumulative substitution scheme starting with fully manual review, followed by substitution with automated review of class no. 1, then 1 and 2, etc. Red rectangle indicates a drop of performance at substitution stage. **c** Cumulative substitution scheme starting with fully manual review, followed by substitution with automated review of class no. 8, then 8 and 7, etc. Red rectangle indicates a drop of performance at substitution stage.

212 the amount of sampled cellular material and the resulting sample confidence, which can vary
213 widely and might improve with future refinements of the device administration procedure.

214 To understand how our results generalize, we devised a simulation approach to vary how
215 many samples have BE and how many samples are allocated to high/low confidence triage
216 classes (Methods). To simulate the change in workload over a range of possible prevalences
217 of BE, we first determined the proportion of patients with and without BE in each triage class
218 and then weighted each vector of proportions by a new prevalence ranging from 0 to 55%. To
219 simulate the effect that relative changes in overall sample confidence have on the workload,
220 we first determined the proportion of patients in triage classes with highest sample confidence
221 (determined by quality control and diagnostic class: 2 and 8) and lower sample confidence (1,
222 3, 4, 5, 6, and 7). We then modified the proportion of high confidence samples and inversely
223 adapted the proportion of lower confidence samples within a range from -25% to 25%.

224 Over a fine grid of varying disease prevalence and changes in sample confidence, we ob-
225 served a negative impact of decreasing cohort BE prevalence and a positive impact of sample
226 confidence on the potential workload reduction (Figure 5a). According to this simulation, in
227 a realistic cohort with a BE prevalence of 7%, we would still be able to reduce the pathology
228 workload by 57%. In order to retain the same workload reduction we observed in the validation
229 cohort, the proportion of samples with high confidence in a realistic cohort would need to be
230 increased by 15%.

231 **External validation of triage-driven approach**

232 Finally, we tested the validity of our results and the extrapolation in the simulation study in
233 an independent test set of 3038 slides from 1519 patients from from 109 primary care sites
234 in the UK (BEST3 trial) (22). All slides were processed in the same way and with the same
235 model parameters as the BEST2 validation cohort (fig. S7, table S6). Following the method

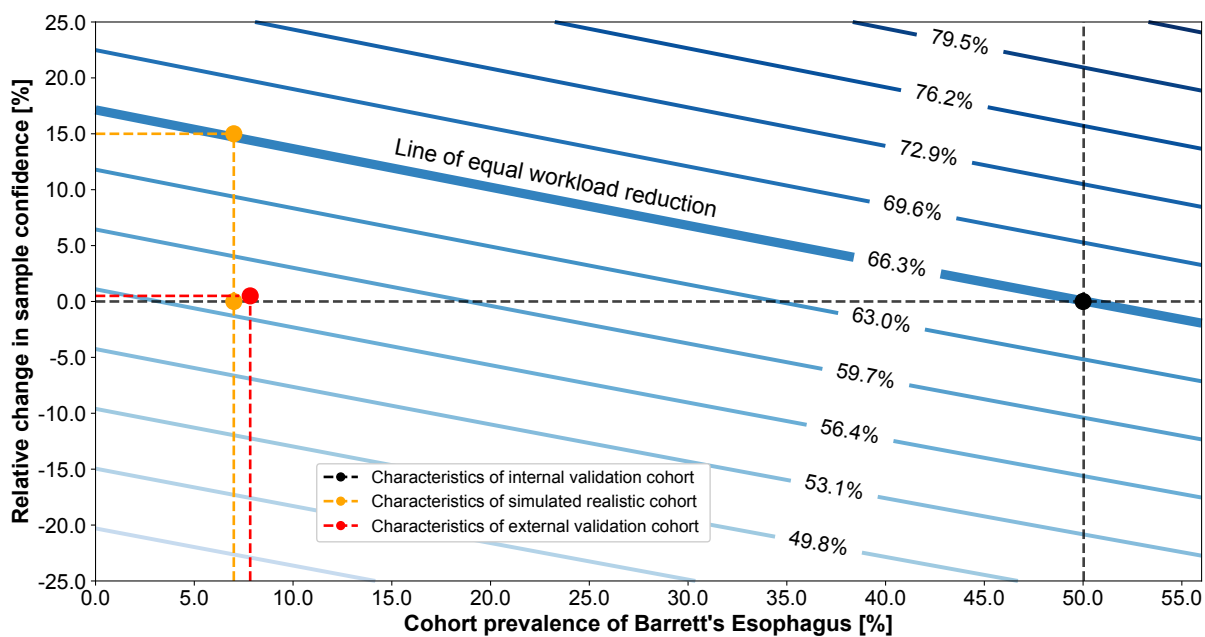


Figure 5: Triage model applied to external validation cohort and simulation of cohort variation
 Simulation of changes in cohort prevalence of BE and sample confidence with impact on workload reduction. Every contour line (blue) represents the same level of workload reduction as indicated by the percentages. Solid black lines indicate the workload reduction of the validation cohort. The dotted yellow line illustrates the workload reduction of a realistic primary care referral cohort (with 7% prevalence) with no change in sample confidence classes (lower yellow marker) and the confidence change required to match the same amount of workload reduction as in the validation cohort (upper yellow marker). The results from the external validation cohort are shown in red.

236 described in the previous section, we used manual pathologist reviews for samples that fell into
237 triage classes 4, 5 and 6. In the BEST3 trial, endoscopy data was only available for positive
238 Cytosponge patients and those who had Barrett's diagnosed at follow-up as a result of standard
239 of care. In addition, the trial was not designed to investigate sensitivity or specificity but positive
240 predictive value (PPV) instead. We also calculated the negative predictive value (NPV) based
241 on findings aggregated through the primary endpoint analysis (coded BE diagnosis in patient
242 records). For this external validation cohort, fully manual review by pathologists resulted in a
243 PPV of 56.08% and NPV of 99.02%. After application of the triage-driven, semi-automated
244 approach the PPV of the overall cohort was 53.37% and the NPV 99.39% (fig. S7). Using this
245 approach in a realistic primary care setting would have resulted in the following key results: In
246 total 872 patients out of 1519 patients (57.41%) would have been reviewed automatically while
247 42.59% would have had to be reviewed manually. This agrees with the simulated, expected
248 value (fig. 5) of workload reduction given the prevalence (7.8%) of BE in this external validation
249 cohort. Six additional patients would have been diagnosed with BE while being missed by the
250 pathologist at the cost of 19 additional endoscopies when compared to fully manual review.
251 One patient would have received an automated negative diagnosis even though the pathologist
252 scored it as positive with BE finding at endoscopy.

253 Discussion

254 We have presented a triage-driven approach that analyses samples of the Cytosponge-TFF3 test
255 using deep learning for the early detection of esophageal cancer. Our approach combines quality
256 control and diagnostic metrics of pathology slides to stratify patients into 8 triage classes which
257 determine whether a patient sample requires manual or if automated review would suffice.

258 For the analysis of Cytosponge-TFF3 samples, our triaging approach has several benefits:
259 We are able to substantially reduce workload and match the sensitivity and specificity of expert
260 pathologists. In our internal validation cohort, fully manual review by a pathologist achieves
261 81.7% sensitivity and 92.7% specificity. In a fully automated approach, we observed a sensitiv-
262 ity of 72.6% and a specificity of 93.1%. With our triage-driven approach, we demonstrate that
263 up to 66% of cases can be reviewed automatically while achieving a sensitivity of 82.5% and
264 specificity of 92.7%, a performance marginally superior to fully manual review by pathologists.
265 Further, in an external validation cohort from a large randomised controlled trial we observed
266 a PPV of 53.37% and NPV of 99.39%. For comparison, pathologist review resulted in very
267 similar values with a PPV of 56.08% and NPV of 99.02%. While a small number of additional
268 endoscopies would have been triggered, they would have also yielded more positive diagnoses.
269 In this more realistic cohort, 57.41% workload for the pathologists would have been reduced.
270 These results (Figure 5) have several implications: First, a fully automated review would re-
271 duce sensitivity (at fixed specificity) and therefore suffer from a loss of clinical utility. Second,
272 while a triage-driven approach is not able to reduce workload as much as a fully automated ap-
273 proach, the described triage classes provide a logical way for stage-wise clinical adoption and
274 performance testing in routine practice.

275 Another benefit of our approach is that we were able to directly adopt heuristics applied
276 by pathologists familiar with Cytosponge-TFF3 samples in our algorithmic design process. As

277 a result, our approach demonstrates traceability and interpretability (8): First, we mimicked
278 the screening process of samples observed by expert pathologists by replicating their decision-
279 making scheme (Figure 1c). Second, the saliency maps we generated from deep learning mod-
280 els to visualize learned features in the pathology images show strong agreement with manual
281 landmarks placed by pathologists (Figure S2).

282 As a further benefit, our triage approach achieves strong performance from only 287 sam-
283 ples for training and calibration by incorporating informative prior knowledge about biological
284 and clinical test characteristics, followed by rigorous testing in independent cohorts. This per-
285 formance compares favorably to previous fully automated approaches reporting expert-level
286 performance that relied on very large datasets with training set sizes ranging from 10,000 to
287 more than 100,000 examples (34, 35) - dataset sizes that cannot be expected for most applica-
288 tions.

289 Finally, a quantitative analysis of workload reduction across varying disease prevalences and
290 sample confidences shows that our approach is expected to generalize well to a real-world pop-
291 ulation. A more general population would have a lower disease prevalence than a case-control
292 study, which would cause a larger workload due to the distribution of BE/non-BE patients within
293 the individual triage classes. We were further able to confirm this simulation with an external
294 validation cohort. These findings provide realistic expectations of how clinical decision-making
295 systems are affected by bias in cohort composition.

296 Our approach has several limitations: First, while samples used in this work were gen-
297 erated at multiple centres they were processed at only a single site (Addenbrookes Hospital,
298 Cambridge, UK). Thus, our data might not fully reflect the variation in histology sectioning
299 and staining across different laboratories (36). We compensated for this limitation through
300 data augmentation by spatial and color profile distortion. Additionally, our data are not too
301 far from future real-world applications, because for large-scale rollout of the Cytosponge test a

302 centralised laboratory is envisaged to ensure processing as well as analysis with proper quality
303 assurance. In future work, we plan to test whether the superiority of the triage-driven approach
304 over fully manual pathologist review will generalize by incorporating multi-centre data from
305 ongoing and future Cytosponge-TFF3 studies to evaluate this effect more extensively.

306 Second, the underlying machine learning model could be further optimized. For example,
307 instead of using a transfer learning model based on pre-training with a primary dataset, we could
308 train a model from scratch, which has proven to improve results in some CNN applications (37).
309 In addition, the tile size needs further investigation because it determines the receptive field in
310 which the CNNs build feature representations of images. Our tile size was chosen by expert
311 pathologists to capture relevant structures like columnar epithelium and goblet cells. Although
312 good performance was observed, a refined multi-scale classification with several magnifications
313 might be necessary to achieve better classification of tissue types. Further improvements might
314 be realised from using attention-based models to reduce the laborious annotation steps required
315 for expanding the training data (38) or aggregating tiles to patient level with more sophisticated
316 approaches based on sequence models (34).

317 Third, a major determinant of workload reduction is the quality and therefore diagnostic
318 confidence attributed to a sample. However, what determines the amount of columnar mate-
319 rial sampled is unknown. One hypothesis is that the strength of esophageal peristalsis, which
320 can be influenced by variations in device ingestion, may be associated with the likelihood of
321 the Cytosponge reaching the stomach. We plan to investigate determinants of sample quality
322 by comparing the data generated by the trained deep learning models with patient and device
323 administrator profiles.

324 In summary, our triage approach differs from previous applications of deep learning to medi-
325 cal images (7,34) which used fully automated approaches on extremely large datasets. We show
326 that for a modest dataset size, leveraging existing heuristics of pathologist decision-making in a

327 triage-based approach is a powerful alternative to fully automated classification models, which
328 generalises well to an independent validation cohort. These results lay the foundation for tai-
329 lored, semi-automated decision support systems embedded in clinical workflows.

330 **Methods**

331 **Study design and dataset**

332 The multicentre Barrett's Esophagus Screening Trial 2 (BEST2) (18) case-control study (study
333 registration: ISRCTN12730505) investigates the automated analysis of Cytosponge-TFF3 sam-
334 ples as a secondary objective. Ethics approval was obtained from the East of England - Cam-
335 bridge Central Research Ethics Committee (number 10/H0308/71) and registered in the UK
336 Clinical Research Network Study Portfolio (9461). Patients enrolled underwent a Cytosponge
337 procedure followed by an endoscopy with biopsies where required. The objective of this work
338 was the comparison of: fully manual review of Cytosponge-TFF3 pathology slides by human
339 experts, fully automated review of Cytosponge-TFF3 pathology slides by a deep learning-based
340 method, and triage-driven, semi-automated review of Cytosponge-TFF3 pathology by a hybrid
341 method relying on deep learning methods and human experts.

342 812 patients were randomly selected from the entire BEST2 cohort (from 11 hospitals in the
343 UK) for digitisation of their respective H&E and TFF3 pathology slides (1624 in total) on an
344 Aperio AT2 digital whole-slide scanner (Leica Biosystems Nussloch GmbH, Germany) at 40x
345 magnification.

346 BEST2 patients were randomly partitioned into three distinct subsets: 100 patients for train-
347 ing/development (labels unblinded for training purposes), 187 patients for calibration (labels
348 unblinded for calibration), and 525 patients as an internal validation set (labels unblinded af-
349 ter validation). The distribution of patients with or without Barrett's Esophagus (BE) for each
350 partition is shown in Figure 1d.

351 For independent external validation we used data from the Barrett's Esophagus Screen-
352 ing Trial 3 (BEST3) REF randomised controlled trial (study registration: ISRCTN68382401).
353 Ethics approval was obtained from the East of England - Cambridge Central Research Ethics

354 Committee (number 16/EE/0546). Patients enrolled either were invited to a Cytosponge proce-
355 dure or received standard of care. Both arms were followed up after 8 to 18 months (weighted
356 overall average of approx. 12 months). Only patients who underwent a Cytosponge procedures
357 or were referred as part of usual care received an endoscopy. A patient was considered as pos-
358 itive for Barrett's Oesophagus if they either had a diagnosis at endoscopy or as a result of a
359 coded search in records from the primary care site.

360 1519 patients were randomly selected from the entire BEST3 cohort (from 109 primary care
361 sites in the UK) for digitisation of their respective H&E and TFF3 pathology slides (1638 in
362 total) on Hamamatsu S60 and S210 whole-slide scanners (Hamamatsu, Japan) at 40x magnifi-
363 cation. For each patient, the repeat test was used if one as performed due to inadquace of the
364 baseline test.

365 All BEST3 patients were processed using the fully automated and triage-driven, semi-
366 automated approach presented in this work. Labels were unblinded after validation.

367 Confidence intervals in this work were defined as the 2.5th and 97.5th percentiles on distri-
368 butions of 500 samples (with replacement) of the respective dataset size.

369 **Cytosponge-TFF3 procedure**

370 The Cytosponge-TFF3 is a non-endoscopic diagnostic modality for BE. It is a cell collection
371 device, consisting of a mesh sphere on a string inside a gelatine capsule, coupled with an im-
372 munohistochemical biomarker called Trefoil Factor 3 (TFF3).

373 The capsule is swallowed by the patient, and passes to the stomach, where the gelatine dis-
374 solves allowing the mesh sphere to expand to a diameter of 3 cm. After 5 to 7.5 minutes, the
375 sponge is withdrawn from the stomach by the attached string, sampling superficial epithelial
376 cells from the top of the stomach, the esophagus, and the oropharynx. The removed device
377 is placed in a container with preservative solution (SurePath Preservative Fluid, BD) and pro-

378 cessed in a laboratory for histochemical (Hematoxylin & Eosin) and immunohistochemical
379 (TFF3) staining. The stained pathology slides are then screened by a pathologist. The primary
380 objective of the Cytosponge-TFF3 test is the detection of columnar epithelium of intestinal type
381 (with TFF3-positive goblet cells) in the squamous oesophagus which is indicative of the patient
382 having Barrett's Esophagus (BE). These TFF3-positive patients can then be referred for an up-
383 per gastrointestinal endoscopy to confirm the diagnosis. Previous studies (17–19) have shown a
384 consistent sensitivity (73.3 % and 79.9 %) and specificity (93.8 % and 92.4 %) for the diagnosis
385 of BE using the Cytosponge coupled with TFF3.

386 **Endoscopy procedure**

387 Esophago-gastroduodenoscopies were carried out by an endoscopist after the Cytosponge test.
388 BE was defined as endoscopically visible columnar-lined esophagus that measured at least 1 cm
389 circumferentially or at least 3 cm in non-circumferential tongues according to the Prague criteria
390 ($\geq C1$ or $\geq M3$ (39)). An additional criterion for BE was histopathological evidence of intestinal
391 metaplasia (IM) on at least one endoscopy biopsy. For cases with suspected BE, diagnostic
392 biopsies were collected following the recommended Seattle surveillance protocol (40). When
393 reviewing the biopsy data, all of the pathologists were blinded to the result of the Cytosponge-
394 TFF3 test.

395 **Whole-slide image annotation for training**

396 One H&E- and one TFF3-stained slide for each of the 100 BEST2 patients from the training
397 set were manually annotated and reviewed by an expert pathologist (MO) using the ASAP soft-
398 ware (41). Regions of interest (ROIs) were selected in the digitised pathology slides at a mag-
399 nification of 40x. Each of these ROIs was labeled with a class for training. For the H&E-based
400 quality control model, four different classes were identified: gastric-type columnar epithelium,

401 respiratory-type columnar epithelium, intestinal metaplasia, and background (including other
402 cellular material such as squamous cells and slide artefacts). Gastric-type columnar epithelial
403 cells were considered as the marker for quality control, as their presence confirms that the Cy-
404 tospane has reached the stomach. For the TFF3-based diagnostic model, three classes were
405 identified: TFF3-positive regions (darkly stained goblet cells), TFF3-equivocal regions (regions
406 of ambiguous staining that may be goblet cells), and background. TFF3-positive cells were con-
407 sidered as the marker for the presence of IM, as they indicate that the patient might have BE.
408 All slides were annotated using the existing patient-level ground truth data for comparison. We
409 aimed for a representative fraction of available material on each slide to be labelled.

410 **Tesselation of whole-slide images for training**

411 Tesselation, or tiling, of whole-slide images was performed in order to prepare data prior to
412 model training. A custom tiling method was developed to optimise the yield and coverage of
413 annotated cellular material in the images. Whereas packing problems of squares in polygons
414 can be neglected for large annotations, optimal coverage for tiles in combination with small
415 annotation sizes is not straightforward and requires a tailored solution. Annotations with an
416 area of $1.5 * \text{tile area}$ or larger were cropped into tiles by taking the top-left coordinate of
417 the enveloping bounding box and iterating tiles along the x- and y-axis of the image. Tiles
418 with an intersection of less than 0.33 (for H&E) or 0.66 (for TFF3) with their corresponding
419 annotation were rejected. Annotations with an area smaller than $1.5 * \text{tile area}$ were treated as
420 single examples and a tile was placed in the center-of-mass of the respective annotation. Tiles
421 with sufficient annotation coverage (determined by intersection) were extracted and labelled
422 according to the class of their parent annotation. For this work, a tile size of 400-by-400 pixels
423 (corresponding to 200-by-200 μm at a magnification of 40x) was selected in accordance with
424 sizes of relevant tissue features. Tiles were extracted from whole-slide images as JPEG images

425 with minimal compression.

426 **Model training using deep learning**

427 We implemented two different deep learning frameworks: one for performing quality control on
428 H&E-stained slides, and a second one for performing automated BE diagnosis from the TFF3-
429 stained slide images. Both deep learning frameworks for quality control and diagnosis were
430 created by comparative transfer learning of multiple convolutional neural network architectures:
431 AlexNet (23), DenseNet (24), Inception v3 (25), ResNet-18 (26), SqueezeNet (27), and VGG-
432 16 (28). All architectures were initialised with the best parameter set that was achieved on
433 the ImageNet competition. Training tile images were resized as required for the individual
434 architectures, resulting in a change of effective magnification from 22x to 30x. We then unfroze
435 all layers to enable fine-tuning of the entire network. For all models, training continued on two
436 NVIDIA GTX 1080Ti graphics cards for 25 epochs with an architecture-specific batch size
437 (ResNet-18: 128, VGG-16: 48, Inception v3: 48, AlexNet: 64, SqueezeNet: 256, DenseNet:
438 84) and a learning rate that decayed by a factor of 0.1 every 7 epochs. All models used cross-
439 entropy loss. To account for slight variations in the training data, random vertical/horizontal flip,
440 random rotation, and random color jitter (variation in hue, contrast, brightness, and saturation)
441 were introduced for data augmentation. Differences in tile class sizes were accounted for by
442 using a modified imbalanced dataset sampler, a function which oversamples from minority
443 classes and undersamples from majority classes. The parameter set of epoch with the highest
444 accuracy on the development subset was selected for further use. All models were trained using
445 the PyTorch deep learning framework (42). Final model versions used a split of 85:15 patients
446 for training and development subset. We further investigated the effect of increased training set
447 sizes by incrementally increasing the training subset while fixing the development subset size
448 (Figure S1).

449 **Evaluation of tile-level performance**

450 In order to compare the performance of all six deep learning architectures, we calculated class-
451 specific performance in the quality control and diagnosis frameworks (Table S1). To obtain
452 these numbers, we selected the epochs with the best weighted accuracy score on the develop-
453 ment subset for each training run. We then calculated precision and recall of all four classes in
454 the H&E-based model and all three classes in the TFF3-based model in the selected epoch. For
455 visual comparison, we also created 2D inference maps of samples which were classified as
456 positive or negative by a pathologist for quality control and diagnosis, respectively. Tile-level
457 results were not used to select architectures for the fully automated or semi-automated, triage-
458 driven approach. The best performing architectures according to relevant class precision and
459 recall on tile level for quality control and diagnosis were selected for saliency map generation.

460 **Generation of saliency maps using Grad-CAM**

461 Gradient-weighted Class Activation Mapping (Grad-CAM) class localisation maps are created
462 by visualising the gradients flowing into the final convolutional layer of the network, just before
463 the fully-connected layers (29). Since convolutional layers contain class-specific spatial infor-
464 mation from the input image which is lost in the fully connected layers, this is the optimal point
465 for map generation. Unlike conventional class-activation maps (CAMs), Grad-CAM has the
466 benefit of not requiring any modifications to the existing model architecture, nor does it require
467 any retraining of the model (29). In order to create the class-specific Grad-CAM localisation
468 map for class c , $L_{\text{Grad-CAM}}^c$, it is first necessary to compute the gradient $\frac{\partial y^c}{\partial A^k}$ of the score y^c for
469 class c with respect to the feature map A^k of the final convolutional layer (29). Once $\frac{\partial y^c}{\partial A^k}$ has
470 been computed for each feature map k , these backward-flowing gradients are global-average-
471 pooled across the width and height of the network (indexed by i and j) to yield α_k^c , the weights
472 of neuron importance for each of the feature maps k (29):

$$\alpha_k^c = \frac{1}{Z} \sum_i \sum_j \frac{\partial y^c}{\partial A_{ij}^k}$$

473 α_k^c , the neuron-importance weights for each feature map k , therefore estimate the saliency
474 of each feature map to the prediction of class c (29). Finally, to get class c -specific Grad-CAM
475 localisation map $L_{\text{Grad-CAM}}^c$, we take the *ReLU* of the weighted sum of the feature maps A^k ,
476 where each feature map k 's weight is α_k^c (29):

$$L_{\text{Grad-CAM}}^c = \text{ReLU} \left(\sum_k \alpha_k^c A^k \right) \quad (1)$$

477 Note that the *ReLU* operation is used to retain only the features which have a positive
478 influence on the prediction of class c , and that the resulting localisation map will be the same
479 size as the feature maps of the last convolutional layer (29).

480 We generated saliency maps for both models trained on H&E and TFF3, respectively. The
481 target layer from the VGG-16 architecture was the last feature layer (no. 30) before several
482 stacked fully connected layers. Tiles were randomly selected from the development subset.
483 For qualitative comparison between saliency maps and manual landmarks, we asked one expert
484 pathologist (MO) to highlight important areas. Areas highlighted by the pathologist provide
485 a representation of features which a human observer uses for classification of tile images. To
486 investigate qualitative agreement of landmarks by the pathologist with generated saliency maps,
487 a side-by-side comparison of tile images and respective saliency maps was prepared (Figure S2).

488 **Model inference on calibration and validation cohort**

489 All six deep learning architectures trained separately for quality control and diagnosis tasks
490 were applied to pathology slides randomised to calibration and validation cohort. Whole-slide
491 images were tessellated on the fly as described above. Detection of tissue was achieved by
492 luminance thresholding of tile values in the LAB colour space. Tiles were forward-passed

493 through the trained deep learning architectures and softmax probabilities were aggregated for
494 each tile position.

495 **Aggregation of classifications on tile level to the patient level**

496 We explored two different aggregation approaches based on propagation of the individual tile-
497 level classifications to patient-level classifications for quality control and diagnosis: a fully au-
498 tomated approach which operates on the basis of a single operating point, and a semi-automated,
499 triage-driven approach which leverages two operating points. For the former approach, perfor-
500 mance was assessed using sensitivity and specificity; for the latter, performance was assessed
501 using an incremental substitution scheme with simultaneous analysis of sensitivity and speci-
502 ficity. For both approaches, tile-level probabilities had to be thresholded to obtain the number
503 of positive tiles per slide for quality control and diagnosis. In the following section, we describe
504 how tile-level probabilities were thresholded and how the operating points on the resulting num-
505 bers of positive tiles (quality control and diagnosis) were then calibrated and evaluated as part
506 of each approach.

507 **Determination of tile-level probability thresholds**

508 In order to generalise the tile-level probabilities to the number of positive tiles per patient,
509 we determined thresholds for each model and endpoint (quality control and diagnosis). The
510 probability threshold of individual tiles for quality control and diagnosis had to be determined,
511 then, the resulting number of positive tiles per threshold was assessed against the best ROC-
512 AUC on the calibration cohort (Figure S3, Table S2).

513 To achieve the best-performing threshold for individual tile probabilities and subsequent
514 aggregation, we iterated over a range of tile thresholds on a fine grid from 0 to 1 (in 0.005
515 steps and inclusive of 0.999, 0.9999, and 0.99999). For the quality control model on H&E,

516 the relevant class was gastric-type columnar epithelium. For the diagnosis model on TFF3, the
517 relevant class was TFF3-positive goblet cells.

518 In order to determine the resulting number of positive tiles per threshold, probability thresh-
519 olds for quality control were compared (ROC-AUC) to the pathologist ground truth of H&E
520 slide analysis. Probability thresholds for diagnosis were compared (ROC-AUC) to endoscopy
521 (confirmation of BE presence by endoscopist and IM on endoscopy biopsy by pathologist)
522 ground truth. This step was required to determine the optimal threshold for individual tile clas-
523 sification. This threshold was then used in the calibration and validation of the fully automated
524 and semi-automated, triage-driven model as described in the next section.

525 **Calibration of fully automated model**

526 All six deep learning architectures trained for quality control and diagnosis were applied to the
527 whole-slide images from the calibration cohort (see Model inference). The number of positive
528 tiles per sample for quality control and diagnosis was determined as described above. To de-
529 termine an adequate operating point for the fully automated patient-level model, ROC analysis
530 was performed on the number of detected tiles (quality control and diagnosis) per patient. On
531 the same set of patients, we calculated the performance by an expert pathologist. In order to
532 determine the ideal cut-off for number of detected tiles, we fixed the specificity of each model to
533 the performance of an expert pathologist on the calibration cohort. The resulting operating point
534 was then chosen for validation of the fully automated model in the validation cohort (Table S2).
535 Tile-level thresholds which yielded the best sensitivity on the calibration cohort were used for
536 evaluating all approaches on the validation cohort. The best-performing architecture (assessed
537 by sensitivity) on the calibration cohort was considered the representative model for application
538 on the validation cohort. However, due to the simplicity of operating point determination, the
539 performance of all other architectures on the validation cohort was also investigated.

540 **Evaluation of fully automated model using ROC analysis**

541 All six deep learning architectures trained for quality control and diagnosis were applied to the
542 whole-slide images from the validation cohort (see Model inference). The number of positive
543 tiles per sample for quality control and diagnosis was determined as described above. Subse-
544 quently, the previously determined operating point (calibration) for each of the deep learning
545 architectures was applied. The binary results were then compared against ground truth of the
546 quality control and diagnosis models. For quality control on H&E, the results were compared to
547 the ground truth of the pathologist who was reading the H&E slide of the Cytosponge test. For
548 diagnosis on TFF3, the results were compared with endoscopy ground truth (with confirmation
549 of BE presence by endoscopist and IM on endoscopy biopsy by pathologist). Sensitivities and
550 specificities on the validation cohort were calculated for all models with an additional presen-
551 tation of ROCs for visualisation (Table S3, Figure S5). For comparison with other approaches,
552 performance metrics of the architecture selected during calibration of the fully automated model
553 were used.

554 **Calibration of triage-driven, semi-automated model**

555 All six deep learning architectures trained for quality control and diagnosis were applied to the
556 whole-slide images from the calibration cohort (see Model inference). For calibration, only
557 the best model (according to ROC-AUC) was presented to two expert observers to determine
558 operating points. The number of positive tiles per sample for quality control and diagnosis was
559 determined as described above (Figure S6). The objective of this approach was a more granular
560 classification of patients into three classes for quality control and diagnosis and subsequent
561 stratification by different class combinations. Therefore, two operating points were determined
562 for each model, instead of one.

563 Both observers were presented with the number of detected tiles and relevant ground truth

564 (Cytosponge pathology and endoscopy) for quality control and diagnosis models. They were in-
565 structed to choose two operating points for each task: First, an operating point which optimises
566 sensitivity with a low number of false positives. Second, an operating point which separates the
567 intermediate region of the first and second operating point from samples with optimised speci-
568 ficity and a low number of false negatives. The resulting operating points were then chosen for
569 validation of the semi-automated, triage-driven model in the validation cohort (Table S5).

570 The two operating points for quality control and diagnosis resulted in three tiers per frame-
571 work and were labelled as follows: for quality control, samples above the first operating point
572 were to be considered as high confidence, samples between the first and second operating point
573 as low confidence, and samples below the second operating point as no confidence. For diagno-
574 sis, samples above the first operating point were to be considered as high confidence positive,
575 samples between the first and second operating points as low confidence equivocal, and samples
576 below the second operating point as high confidence negative. Eight triage classes (number 1
577 to 8) were composed by all possible combinations of quality control and diagnosis classes. The
578 combination (no confidence in quality and high confidence in diagnosis) is likely artifactual and
579 was therefore merged (with no confidence in quality and equivocal in diagnosis) to form triage
580 class 4. Two expert observers then ranked all eight classes from lowest to highest likelihood for
581 patients having BE. They further assigned a qualitative rank for priority of manual review based
582 on the subjective difficulty to review samples that are part of specific triage classes.

583 **Evaluation of triage-driven model on internal validation cohort**

584 The triage-driven, semi-automated model was evaluated by applying a cumulative substitution
585 scheme on the internal validation cohort. The base scenario for all cumulative substitutions
586 was the performance of the pathologists on the entire validation cohort. At every substitution,
587 the pathologists' Cytosponge-TFF3 results were substituted with automated review in the re-

588 spectively triage classes. Then, sensitivity, specificity, and proportion of patients substituted with
589 automated review were calculated and compared against the previous substitution steps. The
590 substitution scheme was applied starting from both ends of the triage class list. First, class 1
591 was substituted with automated review, then classes 1 and 2, then classes 1, 2, and 3, and so on.
592 Second, class 8 was substituted with automated review, then classes 8 and 7, then classes 8, 7,
593 and 6, and so on. We then analysed the sensitivity and specificity curves for deviations from
594 their previous values for each step in both applications of the scheme. Classes which caused a
595 drop in sensitivity or specificity on substitution were considered as ‘difficult’ and we retained
596 manual review by a pathologist for associated samples. For each of the difficult classes we then
597 summed up the number of patients that fell into these classes and divided by the total number
598 of patient in the validation cohort. This ratio was to be considered as the potential workload
599 reduction which this substitution scheme could achieve without notable loss in performance.

600 **Simulation of cohort variation and impact on workload reduction**

601 In order to assess workload reduction in cohorts with different compositions, we simulated the
602 distribution of patients within triage classes with varying BE prevalences and sample confi-
603 dences. Let P be a set of all patients with two subsets: $Q \subseteq P$ contains all patients with BE and
604 its complement $R = P \setminus Q$ contains all patients without BE. We count the proportions of pa-
605 tients in each triage class in each of the sets P , Q , R as vectors \mathbf{c}^P , \mathbf{c}^Q and \mathbf{c}^R , respectively. Our
606 simulation consists in re-weighting these vectors to reflect different BE prevalences and sample
607 confidences. For each element of a range of BE prevalences ($\mathbf{s}_{\text{prev}} = \{0.00, 0.01, \dots, 0.55\}$) we
608 multiply \mathbf{c}^Q by $s \in \mathbf{s}_{\text{prev}}$ and \mathbf{c}^R by $1 - s$. At the same time, for each element of a range of rel-
609 ative sample confidences ($\mathbf{t}_{\text{conf}} = \{-0.25, -0.24, \dots, 0.25\}$) we shift proportions of \mathbf{c}^P between
610 triage classes $\{1, 3, 4, 5, 6, 7\}$ and $\{2, 8\}$ by adding $t \in \mathbf{t}_{\text{conf}}$ to one set of classes and subtracting
611 it from the other. Reduction of workload (W) at every simulation step was defined as \mathbf{c}^P for

612 classes 4, 5, and 6 over classes 1, 2, 3, 7, and 8:

$$W = \frac{c_4^P + c_5^P + c_6^P}{c_1^P + c_2^P + c_3^P + c_7^P + c_8^P}$$

613 **Evaluation of triage-driven model on external validation cohort**

614 The triage-driven, semi-automated model was further evaluated applying it with frozen model
615 parameters on the external validation cohort. Processing of images was performed as described
616 on the internal validation cohort above. The trial from the data originates was investigating
617 real-world implementation of the Cytosponge device technology. Therefore, endoscopy data
618 endoscopy data was only available for positive Cytosponge patients and those who had Barrett's
619 diagnosed at follow-up as a result of standard of care. This resulted in a difference of available
620 data as the study was designed for PPV instead of sensitivity and specificity. The NPV was also
621 calculated by using aggregated findings from the primary trial endpoint. An analysis according
622 to the presented substitution scheme was additionally performed ()

623 **Code availability**

624 The source code of this work is freely available at a public repository:

625 <https://github.com/markowetzlab/cytosponge-triage>.

626 **Data availability**

627 The dataset is governed by data usage policies specified by the data controller (University of
628 Cambridge, Cancer Research UK). We are committed to complying with Cancer Research UK's
629 Data Sharing and Preservation Policy. Whole-slide images used in this study will be available
630 for non-commercial research purposes upon approval by a Data Access Committee due to in-
631 stitutional requirements. Applications for data access should be directed to rcf29@cam.ac.uk.
632 Data derived from the raw images are freely available at a public repository:

633 <https://github.com/markowetzlab/cytosponge-triage>. The code and included
634 data enable replication of the results and figures in this manuscript.

635 **References**

- 636 1. Hawkes, N. Cancer survival data emphasise importance of early diagnosis (2019).
- 637 2. Schiffman, J. D., Fisher, P. G. & Gibbs, P. Early detection of cancer: past, present, and
638 future. *American Society of Clinical Oncology Educational Book* **35**, 57–65 (2015).
- 639 3. Nanda, K. *et al.* Accuracy of the papanicolaou test in screening for and follow-up of
640 cervical cytologic abnormalities: a systematic review. *Annals of internal medicine* **132**,
641 810–819 (2000).
- 642 4. CyR, P. R. Atypical moles. *American family physician* **78** (2008).
- 643 5. Talbot, I., Price, A. & Salto-Tellez, M. *Biopsy pathology in colorectal disease* (CRC Press,
644 2006).
- 645 6. Maung, R. Pathologists' workload and patient safety. *Diagnostic Histopathology* **22**, 283–
646 287 (2016).
- 647 7. Esteva, A. *et al.* A guide to deep learning in healthcare. *Nature medicine* **25**, 24 (2019).
- 648 8. Miotto, R., Wang, F., Wang, S., Jiang, X. & Dudley, J. T. Deep learning for healthcare:
649 review, opportunities and challenges. *Briefings in bioinformatics* **19**, 1236–1246 (2018).
- 650 9. Yu, K.-H., Beam, A. L. & Kohane, I. S. Artificial intelligence in healthcare. *Nature biomedical
651 engineering* **2**, 719–731 (2018).

- 652 10. Bray, F. *et al.* Global cancer statistics 2018: Globocan estimates of incidence and mortality
653 worldwide for 36 cancers in 185 countries. *CA: a cancer journal for clinicians* **68**, 394–424
654 (2018).
- 655 11. Pohl, H., Sirovich, B. & Welch, H. G. Esophageal adenocarcinoma incidence: are we reach-
656 ing the peak? *Cancer Epidemiology and Prevention Biomarkers* **19**, 1468–1470 (2010).
- 657 12. Smyth, E. C. *et al.* Oesophageal cancer. *Nature reviews Disease primers* **3**, 17048 (2017).
- 658 13. Peters, Y. *et al.* Barrett oesophagus. *Nature Reviews Disease Primers* **5** (2019). URL
659 <https://doi.org/10.1038/s41572-019-0086-z>.
- 660 14. El-Serag, H. B., Sweet, S., Winchester, C. C. & Dent, J. Update on the epidemiology of
661 gastro-oesophageal reflux disease: a systematic review. *Gut* **63**, 871–880 (2014).
- 662 15. Spechler, S. J. & Souza, R. F. Barrett’s esophagus. *The New England journal of medicine*
663 **371**, 836–845 (2014).
- 664 16. Odze, R. Histology of barretts metaplasia: Do goblet cells matter? *Digestive diseases and*
665 *sciences* **63**, 2042–2051 (2018).
- 666 17. Kadri, S. R. *et al.* Acceptability and accuracy of a non-endoscopic screening test for barretts
667 oesophagus in primary care: cohort study. *Bmj* **341**, c4372 (2010).
- 668 18. Ross-Innes, C. S. *et al.* Evaluation of a minimally invasive cell sampling device coupled
669 with assessment of trefoil factor 3 expression for diagnosing barrett’s esophagus: a multi-
670 center case–control study. *PLoS medicine* **12**, e1001780 (2015).
- 671 19. Freeman, M., Offman, J., Walter, F. M., Sasieni, P. & Smith, S. G. Acceptability of the
672 cytosponge procedure for detecting barrett’s oesophagus: a qualitative study. *BMJ open* **7**,
673 e013901 (2017).

- 674 20. Paterson, A. L., Gehrung, M., Fitzgerald, R. C. & O'Donovan, M. Role of tff3 as an
675 adjunct in the diagnosis of barrett's esophagus using a minimally invasive esophageal
676 sampling devicethe cytospongetm. *Diagnostic Cytopathology* (2019). URL [https://](https://onlinelibrary.wiley.com/doi/abs/10.1002/dc.24354)
677 onlinelibrary.wiley.com/doi/abs/10.1002/dc.24354. [https://](https://onlinelibrary.wiley.com/doi/pdf/10.1002/dc.24354)
678 onlinelibrary.wiley.com/doi/pdf/10.1002/dc.24354.
- 679 21. Lao-Sirieix, P. *et al.* Non-endoscopic screening biomarkers for barretts oesophagus: from
680 microarray analysis to the clinic. *Gut* (2009).
- 681 22. Fitzgerald, R. *et al.* Cytosponge-trefoil factor 3 versus usual care to identify barretts oe-
682 sophagus in a primary care setting: a prospective, multicentre, pragmatic, randomised con-
683 trolled trial. *The Lancet* (2020).
- 684 23. Krizhevsky, A., Sutskever, I. & Hinton, G. E. Imagenet classification with deep convolu-
685 tional neural networks. In *Advances in neural information processing systems*, 1097–1105
686 (2012).
- 687 24. Huang, G., Liu, Z., Van Der Maaten, L. & Weinberger, K. Q. Densely connected convo-
688 lutional networks. In *Proceedings of the IEEE conference on computer vision and pattern*
689 *recognition*, 4700–4708 (2017).
- 690 25. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J. & Wojna, Z. Rethinking the inception
691 architecture for computer vision. In *Proceedings of the IEEE conference on computer*
692 *vision and pattern recognition*, 2818–2826 (2016).
- 693 26. He, K., Zhang, X., Ren, S. & Sun, J. Deep residual learning for image recognition. In
694 *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778
695 (2016).

- 696 27. Iandola, F. N. *et al.* Squeezenet: Alexnet-level accuracy with 50x fewer parameters and;
697 0.5 mb model size. *arXiv preprint arXiv:1602.07360* (2016).
- 698 28. Simonyan, K. & Zisserman, A. Very deep convolutional networks for large-scale image
699 recognition. *arXiv preprint arXiv:1409.1556* (2014).
- 700 29. Selvaraju, R. R. *et al.* Grad-cam: Visual explanations from deep networks via gradient-
701 based localization. In *Proceedings of the IEEE International Conference on Computer*
702 *Vision*, 618–626 (2017).
- 703 30. Fitzgerald, R. C. *et al.* British society of gastroenterology guidelines on the diagnosis and
704 management of barrett’s oesophagus. *Gut* **63**, 7–42 (2014).
- 705 31. Fan, X. & Snyder, N. Prevalence of barretts esophagus in patients with or without gerd
706 symptoms: role of race, age, and gender. *Digestive diseases and sciences* **54**, 572–577
707 (2009).
- 708 32. Rex, D. K. *et al.* Screening for barretts esophagus in colonoscopy patients with and without
709 heartburn. *Gastroenterology* **125**, 1670–1677 (2003).
- 710 33. Elizondo, J. H. *et al.* Prevalence of barrett’s esophagus: An observational study from
711 a gastroenterology clinic. *Revista de Gastroenterología de México (English Edition)* **82**,
712 296–300 (2017).
- 713 34. Campanella, G. *et al.* Clinical-grade computational pathology using weakly supervised
714 deep learning on whole slide images. *Nature medicine* **25**, 1301–1309 (2019).
- 715 35. Esteva, A. *et al.* Dermatologist-level classification of skin cancer with deep neural net-
716 works. *Nature* **542**, 115–118 (2017).

- 717 36. Bejnordi, B. E., Timofeeva, N., Otte-Höller, I., Karssemeijer, N. & van der Laak, J. A.
718 Quantitative analysis of stain variability in histology slides and an algorithm for standard-
719 ization. In *Medical Imaging 2014: Digital Pathology*, vol. 9041, 904108 (International
720 Society for Optics and Photonics, 2014).
- 721 37. Kieffer, B., Babaie, M., Kalra, S. & Tizhoosh, H. R. Convolutional neural networks for
722 histopathology image classification: Training vs. using pre-trained networks. In *2017 Sev-*
723 *enth International Conference on Image Processing Theory, Tools and Applications (IPTA)*,
724 1–6 (IEEE, 2017).
- 725 38. Courtiol, P. *et al.* Deep learning-based classification of mesothelioma improves prediction
726 of patient outcome. *Nature medicine* **25**, 1519–1525 (2019).
- 727 39. Sharma, P. *et al.* The development and validation of an endoscopic grading system for
728 barretts esophagus: the prague c & m criteria. *Gastroenterology* **131**, 1392–1399 (2006).
- 729 40. Levine, D. S. *et al.* An endoscopic biopsy protocol can differentiate high-grade dysplasia
730 from early adenocarcinoma in barrett’s esophagus. *Gastroenterology* **105**, 40–50 (1993).
- 731 41. Computation Pathology Group, part of the Diagnostic Image Analysis Group, at
732 the Radboud University Medical Center. *Asap*. URL [https://github.com/](https://github.com/computationalpathologygroup/ASAP)
733 [computationalpathologygroup/ASAP](https://github.com/computationalpathologygroup/ASAP).
- 734 42. Paszke, A. *et al.* Pytorch: An imperative style, high-performance deep learning library. In
735 Wallach, H. *et al.* (eds.) *Advances in Neural Information Processing Systems 32*, 8024–
736 8035 (Curran Associates, Inc., 2019). URL [http://papers.neurips.cc/paper/](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
737 [9015-pytorch-an-imperative-style-high-performance-deep-learning-library](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
738 [pdf](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf).

739 **Acknowledgments**

740 This research has been supported by Cancer Research UK (FM: C14303/A17197), the Medi-
741 cal Research Council (RCF: RG84369) and Cambridge University Hospitals NHS Foundation
742 Trust. BEST2 was funded by Cancer Research UK (12088 and 16893). MG acknowledges
743 support from an Enrichment Fellowship from the Alan Turing Institute. MCO acknowledges
744 support from a Borysiewicz Fellowship from the University of Cambridge and Junior Research
745 Fellowship from Trinity College, Cambridge. FM is a Royal Society Wolfson Research Merit
746 Award holder. We thank Michael Schneider, Ruben Drews, Paula Martinez-Gonzalez, and Tris-
747 tan Whitmarsh for valuable input on this work. The authors thank the Cambridge Biomedical
748 Research Centre and the Experimental Cancer Medicine Centre for their support and for provid-
749 ing the infrastructure for the research procedures in Cambridge. Further, we thank the Human
750 Research Tissue Bank, which is supported by the UK National Institute for Health Research
751 (NIHR) Cambridge Biomedical Research Centre, from Addenbrookes Hospital. Last, we thank
752 the BEST2 trial team, the Histopathology core facility at the Cancer Research UK Cambridge
753 Institute, and Pathognomics Ltd for their support.

754 **Author contributions**

755 MG conceived and led the analysis; MCO and AB contributed to the analysis; MG and AB
756 wrote the code for analysis; MO and RCF were involved in collection and labelling of the data;
757 RCF conceived the study; RCF and FM directed the project; MG and FM wrote the manuscript
758 with the assistance and feedback of all other co-authors.

759 **Declaration of interests**

760 The Cytosponge® device technology and the associated TFF3 biomarker have been licensed
761 to Covidien GI solutions (now owned by Medtronic) by the Medical Research Council. MG,
762 MCO, and FM are named inventors on a patent pertaining to technology applied in this work.
763 RCF and MO are named inventors on patents pertaining to the Cytosponge and associated tech-
764 nology. MG, MO, and RCF are shareholders of Cyted Ltd, a company working on early detec-
765 tion technology.

766 Supplementary materials

767 Figures

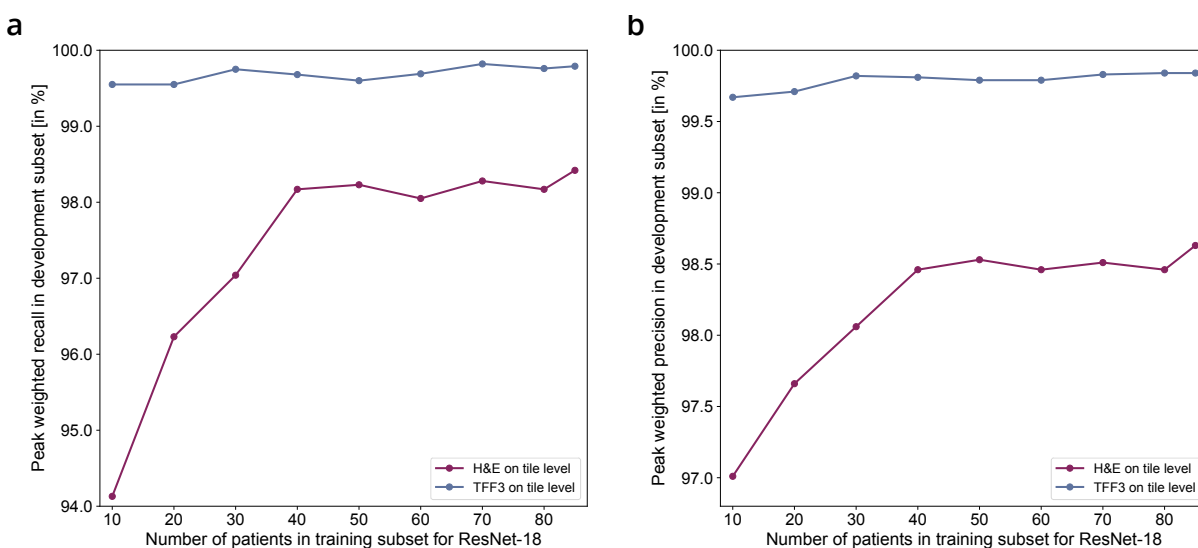


Figure S1: Differential increase of training partition size for ResNet-18. Training subset refers to the relative proportion of the training partition used in the model training phase. Development subset refers to the relative proportion of the training partition used in the model development phase. The peak development weighted recall (a) and precision (b) correspond to the best performing cohort for each training run. The size of the development set was fixed at 15 patients. For each patient, an average of 3,500 tiles was used. For both H&E and TFF3 no substantial increase in performance metrics could be observed after a training subset size of 50 patients. H&E benefited more from an increased number of patients than the TFF3 model. This difference is associated with the increased complexity of detecting different tissue morphologies on H&E vs. brown goblet cells on TFF3.

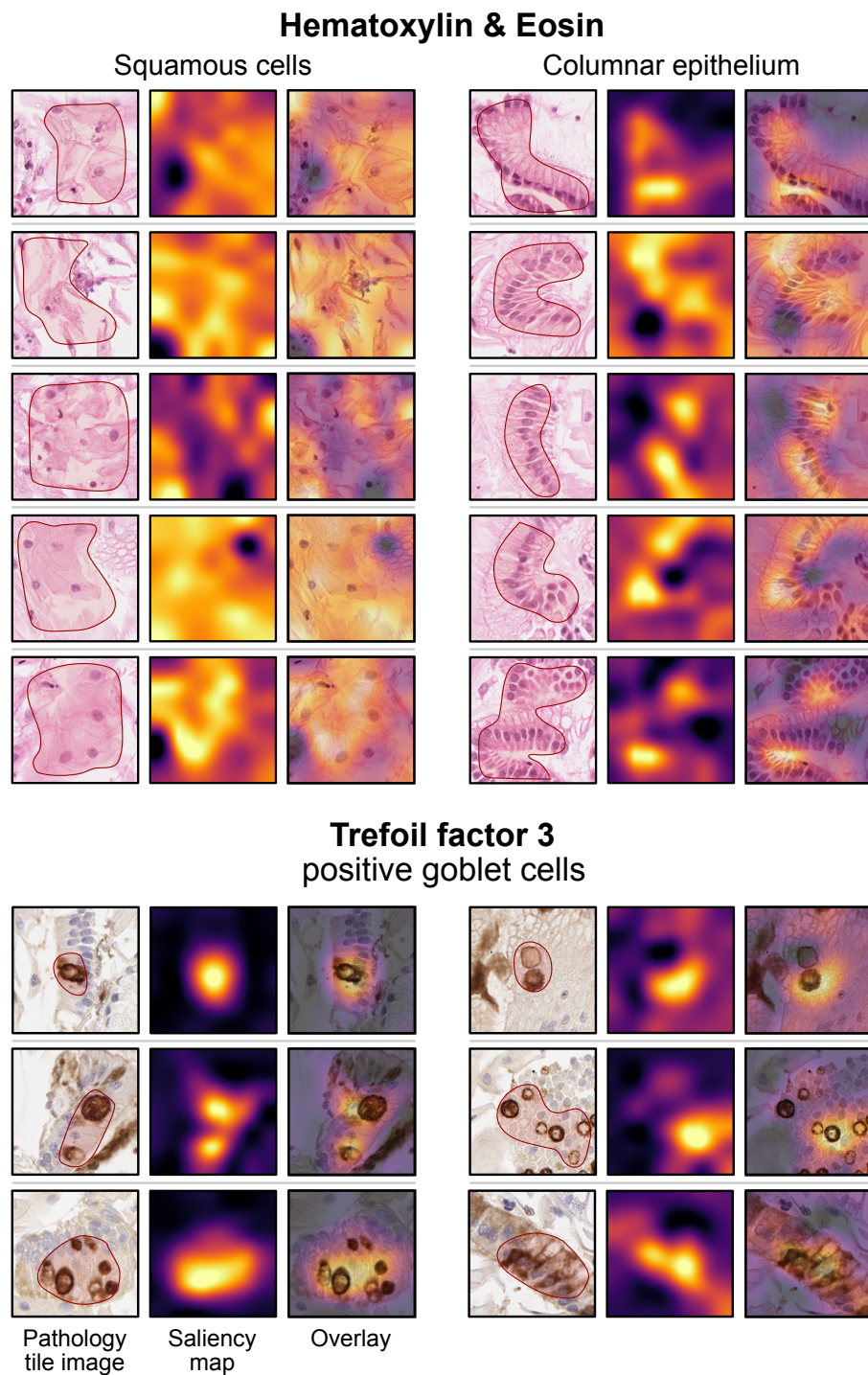


Figure S2: **Comparison of pathologist landmarks with saliency maps extracted from VGG-16 architectures.** Additional examples of saliency maps for Hematoxylin & Eosin stain (squamous cells and columnar epithelium) and Trefoil factor 3 (positive goblet cells). Landmarks selected by an expert pathologist are shown as overlays with red borders on pathology tile images. For all classes, there was visual agreement between highlighted areas by the pathologist and saliency map activations.

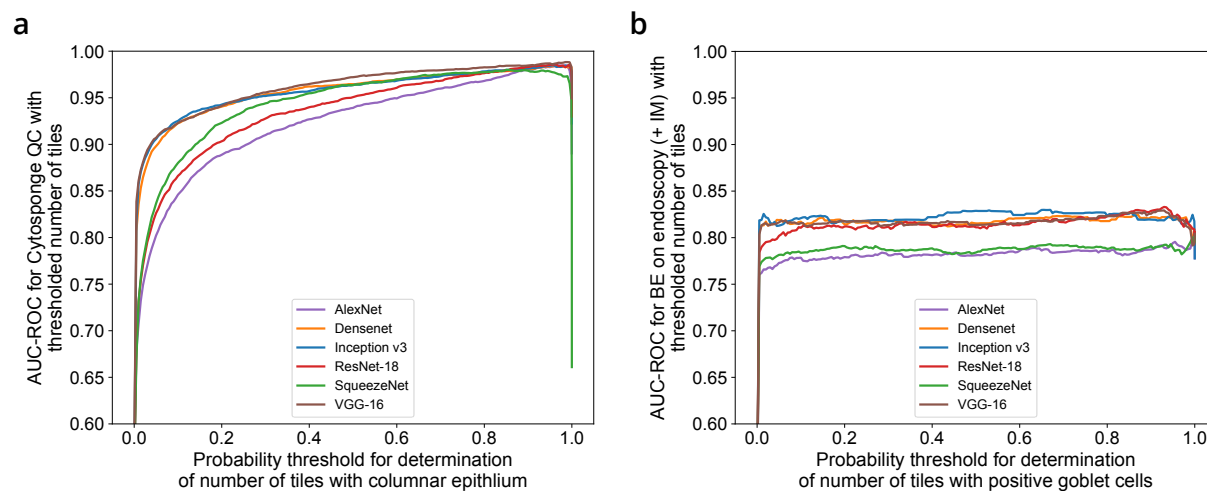


Figure S3: Determination of probability thresholds in order to obtain number of tiles. Both plots show the AUC-ROC for individual probability thresholds (after softmax) which are used to decide whether a tile falls into the relevant class. (a) AUC-ROC for quality control (QC) ground truth determined by the pathologist compared with number of tiles containing columnar epithelium at individual probability thresholds. (b) AUC-ROC for diagnosis ground truth determined by the endoscopy (with confirmed IM on pathology) compared with number of tiles containing positive goblet cells at individual probability thresholds.

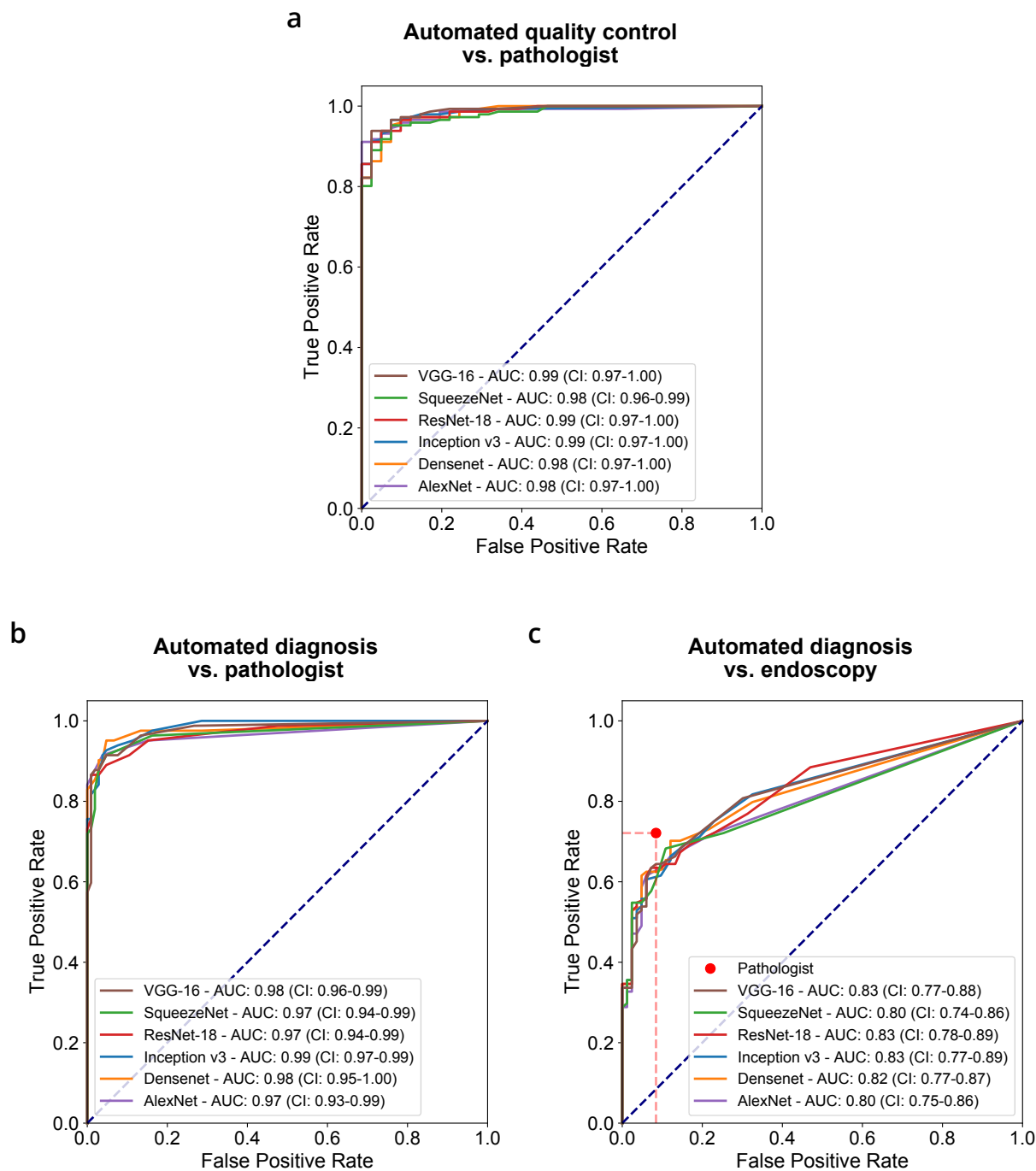


Figure S4: Performance of all deep learning architectures on the calibration cohort. (a) ROC analysis of number of tiles containing columnar epithelium on H&E compared with pathologist ground truth from Cytosponge (b) ROC analysis of number of tiles containing positive goblet cells on TFF3 compared with pathologist ground truth from Cytosponge (c) ROC analysis of number of tiles containing positive goblet cells on TFF3 compared with endoscopy (with confirmed IM) ground truth. A weak AUC dependency on architecture complexity can be observed.

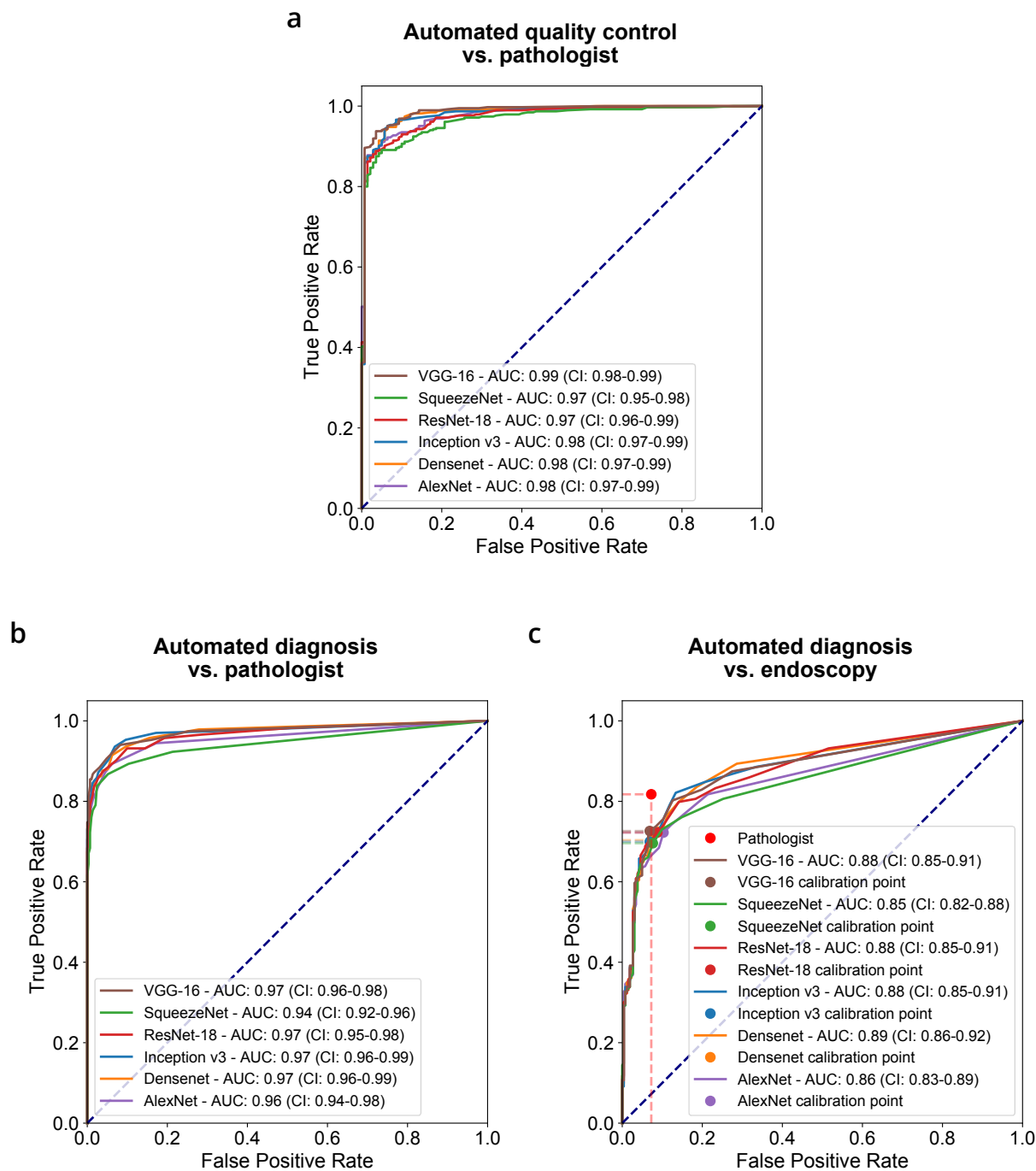


Figure S5: Performance of all deep learning architectures on the validation cohort. (a) ROC analysis of number of tiles containing columnar epithelium on H&E compared with pathologist ground truth from Cytosponge (b) ROC analysis of number of tiles containing positive goblet cells on TFF3 compared with pathologist ground truth from Cytosponge (c) ROC analysis of number of tiles containing positive goblet cells on TFF3 compared with endoscopy (with confirmed IM) ground truth. As in the calibration cohort, a weak AUC dependency on architecture complexity can be observed.

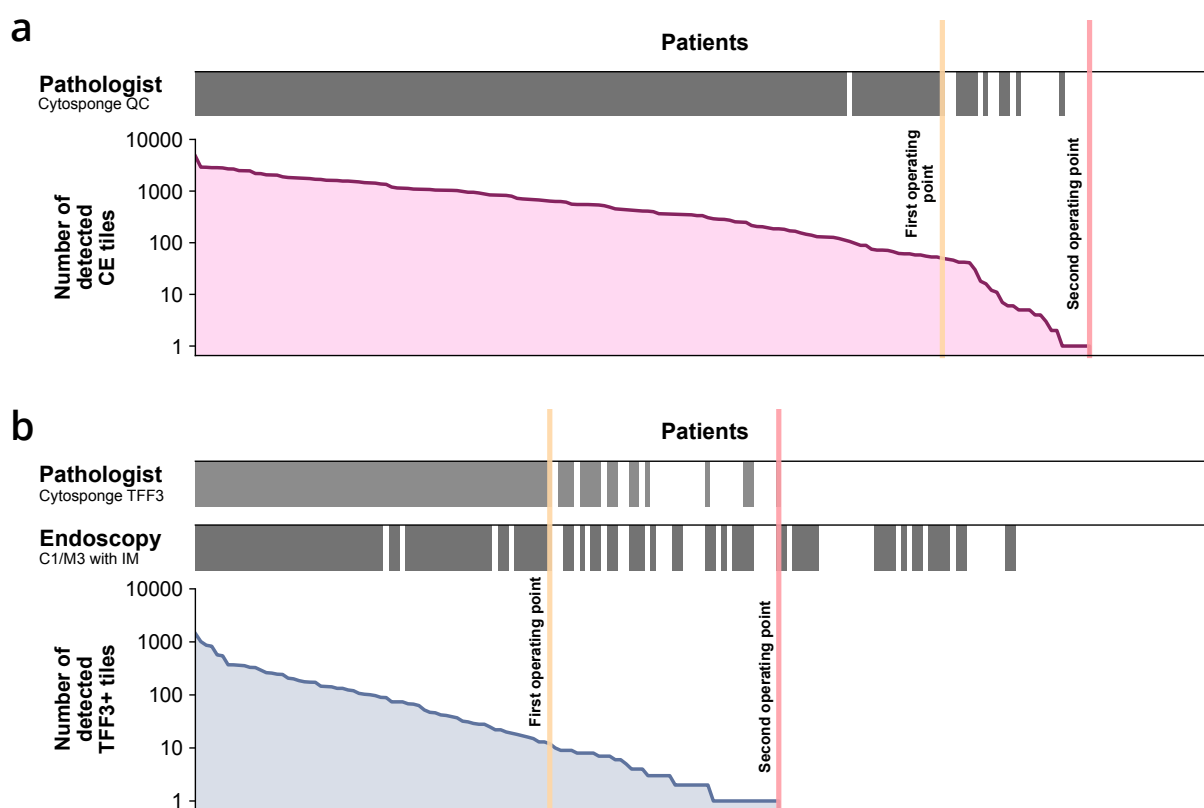


Figure S6: **Application of quality control and diagnostic confidence class scheme to calibration cohort.** **a** Quality ground truth by pathologist from Cytosponge (top) compared with number of detected columnar epithelium (CE) tiles on H&E detected by VGG-16 (bottom). **b** Diagnosis ground truth by pathologist from Cytosponge (top), Endoscopy (with confirmed IM on biopsy) ground truth (middle) compared with number of detected TFF3-positive tiles on TFF3 detected by ResNet-18 (bottom) / eqv. = equivocal.

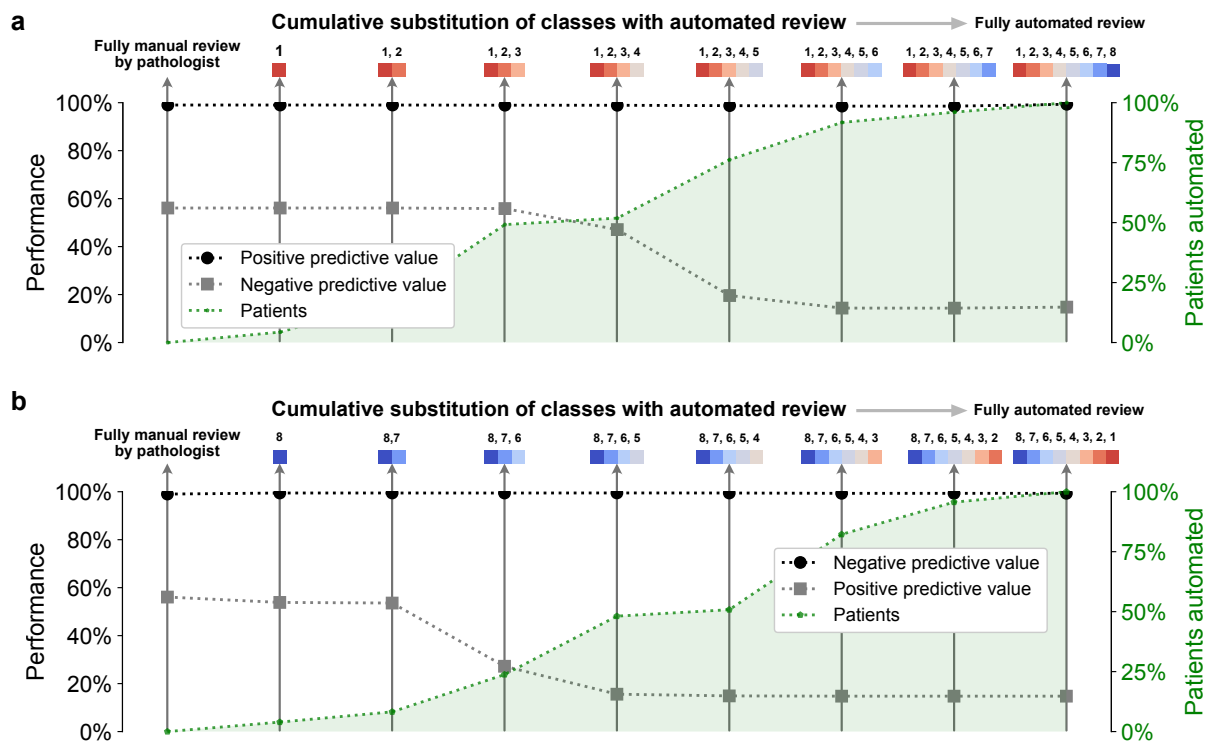


Figure S7: Performance of semi-automated, triage-driven model on external validation cohort **a** Cumulative substitution scheme starting with fully manual review, followed by substitution with automated review of class no. 1, then 1 and 2, etc. **b** Cumulative substitution scheme starting with fully manual review, followed by substitution with automated review of class no. 8, then 8 and 7, etc.

768 **Tables**

	AlexNet	DenseNet	Inception	ResNet	SqueezeNet	VGG
H&E						
Overall accuracy	0.977	0.990	0.989	0.984	0.959	0.988
Precision						
Background	0.999	0.999	0.999	0.999	0.999	0.999
CE (gastric type)	0.791	0.865	0.857	0.807	0.763	0.843
CE (respiratory type)	0.389	0.750	0.895	0.667	0.241	0.741
Intestinal Metaplasia	0.393	0.688	0.609	0.518	0.215	0.640
Recall						
Background	0.984	0.995	0.996	0.991	0.963	0.995
CE (gastric type)	0.893	0.947	0.940	0.921	0.935	0.950
CE (respiratory type)	0.802	0.779	0.588	0.794	0.832	0.634
Intestinal Metaplasia	0.606	0.610	0.629	0.606	0.643	0.568
TFF3						
Overall accuracy	0.996	0.999	0.998	0.998	0.999	0.998
Precision						
Positive	0.752	0.903	0.856	0.827	0.589	0.856
Equivocal	0.233	0.513	0.533	0.385	0.133	0.404
Negative	1.000	1.000	1.000	1.000	1.000	1.000
Recall						
Positive	0.912	0.890	0.919	0.912	0.897	0.919
Equivocal	0.465	0.465	0.372	0.465	0.767	0.442
Negative	0.997	1.000	1.000	0.999	0.991	0.999

Table S1: **Tile-level precision and recall for all classes from H&E and TFF3 models.** This data is derived from the tiles in the development set. (DenseNet = DenseNet-121, Inception = Inception v3, ResNet = ResNet-18, VGG = VGG-16). The highest value(s) per row is/are highlighted in bold.

	AlexNet	DenseNet	Inception	ResNet	SqueezeNet	VGG
Quality control						
Probability threshold	0.97	0.96	0.995	0.96	0.85	0.99
AUC	0.985	0.984	0.986	0.986	0.980	0.988
Diagnosis						
Probability threshold	0.9999	0.87	0.655	0.93	0.99999	0.93
AUC	0.80	0.82	0.83	0.83	0.80	0.83
Sensitivity at fixed specificity (91.57%)	63.4%	62.5%	61.5%	63.5%	60.6%	64.4%
Tile number threshold	3	8	10	9	4	6

Table S2: Individual probability threshold calibration with associated performance based on differential ROC analysis for quality control and diagnosis. The AUC for quality control relates to the performance on the calibration cohort at the given probability threshold for individual tiles containing columnar epithelium on H&E. The AUC for diagnosis relates to the performance on the calibration cohort at the given probability threshold for individual tiles containing positive goblet cells on TFF3. Sensitivity is based on a fixed value of specificity derived from the pathologist performance on the calibration cohort. The tile number threshold is the resulting cutoff from the fixed specificity.

	AUC (CI 95%) vs. pathologist	AUC (CI 95%) vs. endoscopy	Sensitivity (CI 95%)	Specificity (CI 95%)
Quality control				
AlexNet	0.98 (0.97-0.99)	n/a	n/a	n/a
DenseNet	0.98 (0.97-0.99)	n/a	n/a	n/a
Inception v3	0.98 (0.97-0.99)	n/a	n/a	n/a
ResNet-18	0.97 (0.96-0.99)	n/a	n/a	n/a
SqueezeNet	0.97 (0.95-0.98)	n/a	n/a	n/a
VGG-16	0.99 (0.98-0.99)	n/a	n/a	n/a
Diagnosis				
Pathologist	n/a	n/a	81.75% (76.67%-85.92%)	92.75% (89.37%-95.51%)
AlexNet	0.96 (0.94-0.98)	0.86 (0.83-0.89)	72.24% (66.98%-77.37%)	89.70% (85.80%-92.97%)
DenseNet	0.97 (0.96-0.99)	0.89 (0.86-0.91)	70.34% (64.84%-76.24%)	92.75% (89.84%-95.85%)
Inception v3	0.97 (0.96-0.99)	0.88 (0.85-0.91)	69.96% (64.71%-75.65%)	93.13% (89.74%-96.03%)
ResNet-18	0.97 (0.95-0.98)	0.88 (0.85-0.91)	72.24% (66.67%-77.18%)	91.22% (87.72%-94.64%)
SqueezeNet	0.94 (0.92-0.96)	0.85 (0.82-0.88)	69.58% (63.59%-74.54%)	92.37% (88.85%-95.42%)
VGG-16	0.97 (0.96-0.99)	0.88 (0.85-0.91)	72.62% (66.72%-77.64%)	93.13% (89.75%-96.05%)

Table S3: **Performance of all architectures after application on the validation cohort.** Quality control models relied on pathologist calls on sample quality. Sensitivities or specificities were not determined due to irrelevance in the fully automated model approach. Diagnosis models relied on thresholds determined on the calibration cohort.

Quality classes	No confidence	Low confidence	High confidence
No. of patients	22	27	138
Proportion	11.8%	14.4%	73.8%
QC positive (path)	0	9	137
QC negative (path)	22	18	1
Diagnostic classes			
	High conf. negative	Low conf. equivocal	High conf. positive
No. of patients	56	59	72
Proportion	30.0%	31.5%	38.5%
TFF3 positive (path)	1	10	71
TFF3 negative (path)	55	49	1
Barrett's esophagus	12	26	66
No Barrett's esophagus	44	33	6

Table S4: Characteristics of patients in quality control and diagnosis classes from calibration cohort. For each of the three quality control and diagnosis classes, the number of patients within the class and the paired ground truth is shown.

Quality classes	No confidence	Low confidence	High confidence
No. of patients	55	116	354
Proportion	10.5%	22.1%	67.4
QC positive (path)	0	35	350
QC negative (path)	55	81	4
Diagnostic classes			
	High conf. negative	Low conf. equivocal	High conf. positive
No. of patients	145	177	203
Proportion	27.6%	33.7%	38.7%
TFF3 positive (path)	4	33	197
TFF3 negative (path)	141	144	6
Barrett's esophagus	18	61	184
No Barrett's esophagus	127	116	19

Table S5: Characteristics of patients in quality control and diagnosis classes from validation cohort. For each of the three quality control and diagnosis classes, the number of patients within the class and the paired ground truth is shown.

Quality classes	No confidence	Low confidence	High confidence
No. of patients	107	912	500
Proportion	7.1%	60.0%	32.9
QC positive (path)	38	733	350
QC negative (path)	69	179	4
Diagnostic classes			
Diagnostic classes	High conf. negative	Low conf. equivocal	High conf. positive
No. of patients	747	646	126
Proportion	49.2%	42.5%	8.3%
TFF3 positive (path)	1	83	105
TFF3 negative (path)	746	563	21
Barrett's esophagus	5	38	76
No Barrett's esophagus	742	608	50

Table S6: Characteristics of patients in quality control and diagnosis classes from external validation cohort. For each of the three quality control and diagnosis classes, the number of patients within the class and the paired ground truth is shown.