

1 **Deep learning-based Helicobacter pylori detection for histopathology: A diagnostic study**

2

3 Sharon Zhou^{1†}, Henrik Marklund^{1†}, Ondrej Blaha^{2,3}, Manisha Desai², Brock Martin⁴, David

4 Bingham⁴, Gerald J. Berry⁴, Ellen Gomulia⁴, Andrew Y. Ng^{1,5}, Jeanne Shen^{4,5*}

5

6 **Affiliations**

7 ¹ Department of Computer Science, Stanford University, Stanford, CA, USA

8 ² Quantitative Sciences Unit, Department of Medicine, Stanford University, Stanford, CA, USA

9 ³ Yale School of Medicine, New Haven, CT, USA

10 ⁴ Department of Pathology, Stanford University, Stanford, CA, USA

11 ⁵ Center for Artificial Intelligence in Medicine & Imaging, Stanford University, Stanford, CA, USA

12

13 † Equal contribution

14 * **Corresponding Author:** Jeanne Shen, Department of Pathology and Center for Artificial

15 Intelligence in Medicine & Imaging, Stanford University, USA (jeannes@stanford.edu)

16

17 **Running title:** Deep learning-assisted H. pylori detection

18 **Word count:** 3,958

19 **Conflicts of interest:** The authors declare no competing interests.

20

21 **Abstract**

22 **Aims:** Deep learning (DL), a sub-area of artificial intelligence, has demonstrated great promise
23 at automating diagnostic tasks in pathology, yet its translation into clinical settings has been
24 slow. Few studies have examined its impact on pathologist performance, when embedded into
25 clinical workflows. The identification of *H. pylori* on H&E stain is a tedious, imprecise task which
26 might benefit from DL assistance. Here, we developed a DL assistant for diagnosing *H. pylori* in
27 gastric biopsies and tested its impact on pathologist diagnostic accuracy and turnaround time.

28 **Methods and results:** H&E-stained whole-slide images (WSI) of 303 gastric biopsies with
29 ground truth confirmation by immunohistochemistry formed the study dataset; 47 and 126 WSI
30 were respectively used to train and optimize our DL assistant to detect *H. pylori*, and 130 were
31 used in a clinical experiment in which 3 experienced GI pathologists reviewed the same test set
32 with and without assistance. On the test set, the assistant achieved high performance, with a
33 WSI-level area-under-the-receiver-operating-characteristic curve (AUROC) of 0.965 (95% CI
34 0.934-0.987). On *H. pylori*-positive cases, assisted diagnoses were faster ($\hat{\beta}$, the fixed effect
35 size for assistance= -0.557, p=0.003) and much more accurate (OR=13.37, p<0.001) than
36 unassisted diagnoses. However, assistance increased diagnostic uncertainty on *H. pylori*-
37 negative cases, resulting in an overall decrease in assisted accuracy (OR=0.435, p=0.016) and
38 negligible impact on overall turnaround time ($\hat{\beta}$ for assistance=0.010, p=0.860).

39 **Conclusions:** DL can assist pathologists with *H. pylori* diagnosis, but its integration into clinical
40 workflows requires optimization to mitigate diagnostic uncertainty as a potential consequence of
41 assistance.

42

43 **Keywords:** Helicobacter, pathology, artificial intelligence, deep learning, machine learning,
44 stomach, imaging

45

46 **Introduction**

47 *Helicobacter pylori* is the most prevalent chronic bacterial infection worldwide, affecting an
48 estimated 4.4 billion individuals.¹ A well-established association exists between chronic *H.*
49 *pylori* infection and peptic ulcer disease, gastric cancer, and other gastric pathologies, as well
50 as evidence linking infection to iron deficiency anemia, colorectal cancer, and other extra-gastric
51 pathologies.²⁻⁴ Identifying and eradicating *H. pylori* in infected individuals reduces the risk of
52 progression to long-term complications.^{5,6} Every gastric biopsy received in the pathology lab is
53 evaluated for *H. pylori*, with the diagnosis resting on identification of one or more organisms
54 upon high-magnification examination. Although the bacteria are readily identifiable on H&E stain
55 in cases where high numbers of organisms are present, in other cases, making a diagnosis can
56 be time consuming, tedious, and subject to interobserver variability, with missed diagnoses not
57 uncommon.^{7,8} Although ancillary stains are available, these add significant cost and turnaround
58 time to diagnosis,^{8,9} are not recommended for reflex application,^{10,11} and may be unavailable in
59 low-resource settings. *H. pylori* diagnosis on H&E stain provides an ideal opportunity for deep
60 learning (DL) assistance, which has already shown promise at automating other pathology
61 tasks.¹²⁻¹⁶

62
63 Although DL (and other AI) models have demonstrated good performance on several
64 histopathologic tasks, most studies have retrospectively compared model performance to that of
65 a mixed group of diagnosticians of varying expertise levels. This may exaggerate the relative
66 performance and clinical utility of the model, as the true end-users may be different (often with
67 higher baseline diagnostic performance) from those against whom the model's performance
68 was compared. Furthermore, few studies have taken the next step of incorporating a model into
69 a clinical workflow and evaluating its impact on users.^{17,18}

70

71 In this study, we developed a DL ensemble of convolutional neural networks (CNNs) to assist
72 pathologists with *H. pylori* diagnosis on H&E-stained whole-slide images (WSI), and sought to
73 address the preceding gaps by testing the ensemble's impact on the diagnostic performance of
74 subspecialty GI pathologists.

75

76 **Materials and Methods**

77 Institutional Review Board approval was obtained (IRB #48684), with waived informed consent
78 for use of all patient material and data. The Standards for Reporting of Diagnostic Accuracy
79 Studies (STARD)¹⁹ guidelines were used.

80

81 **Dataset and reference standard annotations**

82 The study dataset consisted of H&E WSI of 311 gastric biopsies from 245 patients (160 *H.*
83 *pylori*-positive and 151 *H. pylori*-negative biopsies), all with diagnostic confirmation by both H&E
84 and *H. pylori* immunohistochemical (IHC) evaluation, obtained through stratified random
85 sampling (maintaining an approximate 50:50 class balance of *H. pylori*-positive and -negative
86 biopsies, for adequate representation of the morphologic range of each class) of all gastric
87 biopsies submitted to our institution from January 1, 2015-December 31, 2018. All WSI were
88 scanned at 40x magnification (0.25 micrometers per pixel) on an Aperio AT2 scanner (Leica
89 Biosystems, Germany).

90

91 As WSI are too large to directly input into CNNs, they are subdivided into smaller image patches
92 for input. Because not every patch in an *H. pylori*-positive WSI necessarily contains *H. pylori*,
93 reference standard patch-level annotations were generated using *H. pylori* IHC. A sequential
94 H&E de-stain/immunostain procedure was performed on the same tissue section to obtain
95 annotations for the *H. pylori*-positive biopsies (see Figure 1a and Supplementary methods for
96 details). This was done instead of using the original IHC slides, which contained sections cut at

97 different depths within the block from the H&E slides, resulting in differences in the tissue and *H.*
98 *pylori* content of the H&E and IHC slides (precluding accurate tissue co-registration and
99 generation of reference standard annotations). For *H. pylori*-negative biopsies, an original
100 diagnostic H&E slide was scanned from each biopsy.

101
102 During the de-stain/re-stain process, eight slides experienced focal tissue detachment and/or
103 incomplete immunostaining; these were excluded, resulting in a final dataset of 303 H&E WSI
104 (152 *H. pylori*-positive and 151 *H. pylori*-negative), which was grouped by patient, shuffled, and
105 randomly split (maintaining an approximate 50:50 positive:negative balance within each set) into
106 a training set of 47 WSI (27 positive and 20 negative), validation set of 126 WSI (60 positive and
107 66 negative) for internal validation and optimization, and test set of 130 WSI (65 positive and 65
108 negative), which was completely held out from model training and internal
109 validation/optimization, and used for the pathologist experiment. All WSI from the same patient
110 were assigned to the same set. Each WSI was subdivided into non-overlapping image patches
111 of size 1024 x 1024 pixels (256 x 256 micrometers) for model input, yielding, on average, 486
112 patches per WSI, and a total sample size of 147,258 patches. As multiple serial sections cut
113 from the same block might be present on a WSI, only one section per WSI (the first, or left-most
114 one on the slide) was used for training, yielding a training set of 24,786 patches. The 1024 x
115 1024 pixel patch size was chosen to be large enough to provide background tissue context for
116 the models to learn from, but small enough for pathologists to easily view the entire patch at
117 maximum resolution (40x magnification).

118

119 **Model development**

120 Two CNN architectures, ResNet²⁰ and DenseNet²¹, both of which have excelled on image
121 classification benchmarks, were leveraged by averaging their predictions in a model ensemble,
122 or set of CNNs whose outputs are combined to form a single prediction. Ensembling enables

123 more robust prediction that is less sensitive to outlier predictions, as the final prediction comes,
124 not from a single network, but from a collection of networks. Our ensemble consisted of three
125 ResNet-18 and three DenseNet-121 architectures, each of which input an image patch and
126 output a probability of *H. pylori* being present (for ensemble details, see Supplementary
127 methods).

128
129 Patch-level probabilities were aggregated into WSI-level probabilities by averaging the top 10
130 patch-level probabilities from a single serial section (the first section) per WSI. As each WSI
131 contained between 1-5 sections, only one section per WSI was used to avoid potential
132 prediction bias related to the number of sections on a slide. The patch number (10) was
133 determined empirically, based on the area under the receiver-operating-characteristic curve
134 (AUROC) performance of the ensemble on the validation set. This 10-patch probability average
135 was binarized using an optimal slide-level probability threshold (0.72) empirically determined
136 from the ensemble's performance on the validation set.

137

138 **Pathologist experiment**

139 In a diagnostic study designed to simulate the workload in a high-volume pathology practice, we
140 evaluated the ensemble's impact on the accuracy and turnaround time of three subspecialty GI
141 pathologists with 8, 9, and 27 years of respective practice experience, who reviewed the same
142 test set (65 *H. pylori* positive and 65 *H. pylori* negative biopsies) during a single session, without
143 time constraint. Half of each subgroup (positive and negative) was randomly assigned to review
144 with ensemble assistance, while the remaining half was reviewed unassisted. The WSI
145 sequence and assistance status were randomized for each pathologist (Figure 2a). The
146 pathologists were blinded to all patient identifying and clinical information.

147

148 QuPath²², an open-source digital pathology package, was used for WSI review. Prior to the
149 experiment, the ensemble was run on the 130 test WSI to generate patch-level *H. pylori*
150 probabilities for every tissue-containing patch across all sections in each WSI, with the 3-5
151 highest-probability patches from each WSI and their corresponding probabilities selected for
152 display. To avoid biasing the pathologists, binarized model predictions and the probability
153 threshold for positivity were not shown. The ensemble made predictions on all sections present
154 on the slide, reflective of a real-world practice scenario in which pathologists review all sections
155 on a slide. (However, for the purposes of evaluating machine-learning performance metrics for
156 the ensemble, only a single section per WSI was used.) For WSI containing only one section,
157 bounding boxes for the 3 top-ranked patches were displayed. For WSI containing two or more
158 sections, bounding boxes for the 5 top-ranked patches were displayed. This allowed the
159 ensemble to make predictions across all tissue present in a WSI, while displaying a relatively
160 limited number of patches, to avoid substantially slowing down the pathologists' slide review.

161
162 The ensemble's outputs were incorporated into the QuPath interface as follows: a window
163 adjacent to the main WSI-viewing window displayed the list of all bounding boxes for a given
164 WSI, with corresponding patch-level *H. pylori* probabilities. By clicking on a particular bounding
165 box, pathologists could jump directly to the corresponding region in the main WSI-viewing
166 window (see Figure 2b for an example). A numbered list of all 130 WSI was visible to the left of
167 the main WSI-viewing window. WSI assigned for assisted review were marked with a text
168 indicator ("Hierarchy"). On these cases, the pathologists were given discretion to review as few,
169 or as many, of the bounding boxes and probabilities as they felt were necessary (including
170 none). For WSI assigned to unassisted review, no bounding boxes or probabilities were
171 displayed, and the pathologists simply navigated the WSI on their own.

172

173 For each WSI, the pathologist entered one of four diagnoses into a data entry software
174 application which recorded a timestamp for computation of the diagnostic turnaround time:
175 Positive (P) = *H. pylori* evident on H&E stain, Uncertain Positive (UP) = features suggestive of
176 *H. pylori* infection, but would order IHC for confirmation, Uncertain Negative (UN) = most likely
177 negative for *H. pylori*, but would order IHC for confirmation, and Negative (N) = would
178 confidently diagnose as *H. pylori* negative on H&E stain. Use of the uncertain diagnostic choices
179 (UP/UN) was intended to reflect clinical practice, where, rather than making an immediate
180 choice between P or N, pathologists can be uncertain (resulting in a request for ancillary stains).
181 Diagnoses were binarized for statistical analysis as follows: P and N diagnoses concordant with
182 the reference standard annotation were considered correct, whereas UP and UN diagnoses,
183 and P and N diagnoses discordant with the reference standard, were treated as incorrect. The
184 rationale for treating uncertain diagnoses as incorrect was that uncertainty results in the
185 performance of ancillary stains, which increases diagnostic turnaround time and cost.

186
187 The experiment was administered to each pathologist by the same administrator, who logged
188 the time and duration of any interruptions (used to adjust the timestamp data for accurate
189 determination of the per-slide turnaround time). All experiments were performed using the same
190 workstation setup. Participants were given time at the beginning of each experiment to review a
191 tutorial on use of QuPath and the data entry software application, and to practice the experiment
192 workflow with 6 practice WSI that were not part of the 130 WSI test set.

193
194 **Statistical Analyses**

195 We evaluated both the WSI-level and patch-level diagnostic performance of our ensemble using
196 the AUROC, precision (positive predictive value), recall (sensitivity), specificity, accuracy, and
197 F1-score, based on the respective probability binarization thresholds of 0.504 and 0.72 for patch

198 and WSI-level predictions. Corresponding 95% confidence intervals for these metrics were
199 calculated by bootstrapping, with a replicate size of 2,000.

200
201 Patch-level model performance was assessed on 1024 x 1024 pixel patches sampled from 10
202 *H. pylori*-positive and 10 *H. pylori*-negative WSI randomly selected from the 130 WSI test set.
203 From each WSI, 100 patches were randomly sampled from a single section. If a WSI contained
204 fewer than 100 patches in a section, all patches from that section were sampled. This resulted
205 in a total of 871 patches from the 10 positive and 963 patches from the 10 negative WSI. Patch-
206 level reference standard diagnoses for positive cases were obtained using the de-stain/re-stain
207 method detailed previously, with confirmation by the reference pathologist. After excluding 87
208 patches with equivocal *H. pylori* status upon reference pathologist review, a total of 1,747
209 patches were used for patch-level performance evaluation.

210
211 The ensemble's WSI-level performance was evaluated on all 130 WSI in the test set, where the
212 average probability of *H. pylori* positivity across the 10 highest-probability patches from one
213 section per WSI was binarized using the probability threshold of 0.72 WSI (Figure 1b). The
214 same metrics used to evaluate patch-level performance were also calculated for WSI-level
215 performance.

216
217 Pathologist performance was reported using the diagnostic accuracy (with 95% Wilson score
218 confidence intervals²³) and per-slide turnaround time (with 95% t-score confidence intervals).
219 Our first objective was to investigate whether assistance was effective at increasing accuracy,
220 based on the definitions of correct and incorrect diagnoses used to binarize the pathologist
221 diagnoses. A generalized linear mixed model (GLMM) which included the assistance status
222 (with or without assistance) and pathologist as fixed effects, and WSI as a random effect, was
223 applied. The main effect of assistance was evaluated using a Wald z-test.

224
225 Our second objective was to investigate whether assistance affected the amount of time spent
226 reaching the diagnosis. A log-normal mixed effect model, which included the assistance status
227 and pathologist as fixed effects, and WSI as a random effect, was applied to estimate the effect
228 of these variables on diagnostic turnaround time. The main effect of assistance was evaluated
229 using a Wald t-test.

230
231 A significance level of $\alpha=0.05$ (two-tailed) was used for all statistical tests. The mixed effect
232 models were developed using the lme4²⁴ and MASS²⁵ packages in R.

233

234 **Results**

235 *Model performance*

236 On the validation set (126 WSI), our ensemble achieved a WSI-level AUROC=0.952 (95% CI
237 0.913-0.983), with precision=0.873 (95% CI 0.788-0.953), recall=0.917 (95% CI 0.844-0.982),
238 specificity=0.880 (95% CI 0.794-0.955), accuracy=0.897 (95% CI 0.841-0.944), and F1-
239 score=0.894 (95% CI 0.833-0.947). On the test set (130 WSI), the WSI-level AUROC=0.965
240 (95% CI 0.934-0.987), with precision=0.919 (95% CI 0.846-0.983), recall=0.877 (95% CI 0.788-
241 0.948), specificity=0.924 (95% CI 0.857-0.984), accuracy=0.900 (95% CI 0.846-0.946), and F1-
242 score=0.898 (95% CI 0.837-0.947). On the validation set, patch-level performance metrics
243 were: AUROC=0.877 (95% CI 0.863-0.881), accuracy=0.872 (95% CI 0.862-0.881), F1=0.604
244 (95% CI 0.574-0.633), precision=0.616 (95% CI 0.580-0.651), recall=0.593 (95% CI 0.558-
245 0.629), and specificity=0.927 (95% CI 0.918-0.935), using a probability threshold of 0.504 for *H.*
246 *pylori* positivity. On the test set, the patch-level AUROC=0.911 (95% CI 0.888-0.933), with
247 accuracy=0.892 (95% CI 0.878-0.907), F1=0.573 (95% CI 0.516-0.627), precision=0.488 (95%
248 CI 0.428-0.549), recall=0.692 (95% CI 0.623-0.760), and specificity=0.916 (95% CI 0.902-
249 0.929).

250

251 *Model impact on pathologist performance*

252 During the pathologist study, four reads were excluded due to data entry errors made by the
253 pathologists while recording diagnoses into the software application, resulting in a final set of
254 386 pathologist reads for analysis. The overall accuracy of the pathologists was 0.680 (95% CI
255 0.611-0.742) on unassisted cases and 0.573 (95% CI 0.502-0.641) on assisted cases. After
256 controlling for pathologist and WSI effects, we found that assistance had a significant positive
257 impact on *H. pylori*-positive cases, with assisted diagnoses being more accurate (OR=13.37,
258 95% CI 3.622-49.320, $p < 0.001$) than unassisted diagnoses. However, on *H. pylori*-negative
259 cases, assistance had a negative impact, with unassisted diagnoses being 2.30 times more
260 likely to be correct than assisted diagnoses (OR for correct diagnosis=0.435, 95% CI, 0.215-
261 0.844 $p=0.016$), resulting in a decrease in overall accuracy (Figure 3b).

262

263 The average per-slide turnaround time was 67.15 seconds (95% CI 59.11-75.20 s, range 6-408
264 s) in the unassisted state, and 70.69 seconds (95% CI 63.00-78.38 s, range 11-356 s) in the
265 assisted state. The average per-slide turnaround time for each diagnostic category was:
266 P=43.16 seconds (95% CI 35.60-50.73 s, range 6-356 s), N=51.87 seconds (95% CI 46.53-
267 57.21 s, range 16-168 s), UP=129.68 seconds (95% CI 115.30-144.07 s, range 47-408 s), and
268 UN=84.13 seconds (95% CI 75.74-92.53 s, range 39-221 s) (Figure 4a). After controlling for
269 pathologist and WSI effects, we found a significant reduction in diagnostic turnaround time for
270 *H. pylori*-positive cases ($\hat{\beta}$, the fixed effect size for assistance= -0.557, 95% CI -0.338-
271 -0.119, $p=0.003$), which was offset by a slower average turnaround time on negative cases,
272 resulting in an overall negligible change in turnaround time with assistance ($\hat{\beta}=0.010$, 95% CI -
273 0.097-0.116, $p=0.860$).

274

275 **Discussion**

276 In this study, we evaluated the impact of DL assistance on the diagnostic accuracy and
277 turnaround time of subspecialty GI pathologists, for the routine task of diagnosing *H. pylori* on
278 H&E-stained gastric biopsies. We observed a significant improvement in both metrics with
279 assistance on *H. pylori*-positive cases, but a detrimental effect of assistance on both metrics for
280 *H. pylori*-negative cases.

281
282 In particular, the negative impact of assistance when evaluating *H. pylori*-negative cases
283 resulted in a decrease in overall pathologist accuracy, which might be explained by increased
284 diagnostic uncertainty with assistance, with 81 versus 60 uncertain diagnoses made in the
285 assisted and unassisted states, respectively. Because the pathologists were not informed of the
286 probability threshold for positivity (to avoid biasing their diagnoses), subjectivity in their
287 interpretation of the probabilities could have contributed to increased uncertainty. Based on the
288 method used to binarize diagnoses into correct and incorrect categories (uncertain diagnoses
289 counted as incorrect, Figure 4b), 20 more incorrect diagnoses were made with assistance.
290 However, when uncertain diagnoses were excluded from analysis, more correct diagnoses were
291 made with, versus without, assistance, suggesting that the detrimental effect of assistance was
292 primarily attributable to increased diagnostic uncertainty, rather than to increased diagnostic
293 error. Increased uncertainty may be an unintended consequence of AI assistance, which should
294 be considered when designing and incorporating AI models into clinical practice. For example, a
295 follow-up study of the current model might evaluate the impact of assistance when binarized
296 outputs are displayed, rather than probabilities.

297
298 During the experiment, the pathologists were provided with the 3-5 top-ranked patches on each
299 assisted case, regardless of what the corresponding patch-level probabilities were. Although

300 they were allowed to review these patches at their discretion, they might have felt obligated to
301 view at least some, or even all, patches, simply because these were being presented. This
302 might have introduced diagnostic uncertainty where there initially was not, contributing to
303 additional turnaround time. Given the significant benefit of assistance on *H. pylori*-positive
304 slides, future user interfaces might be designed so that assistance is provided only when
305 patches exceed a probability threshold for positivity, and model outputs are hidden when a slide
306 is predicted to be negative.

307

308 While many studies of medical AI models have emphasized diagnostic accuracy, sensitivity, or
309 specificity as primary metrics, less attention has been devoted to examining the impact on other
310 practical considerations, such as turnaround time. In our study, the average per-slide
311 turnaround time was approximately 67 seconds unassisted and 71 seconds assisted, with no
312 statistically significant difference in turnaround time with assistance, after controlling for
313 pathologist and WSI effects. Given these relatively quick turnaround times, the impact of
314 assistance on this metric might be considered inconsequential, from a clinical standpoint.
315 However, the potential for human-computer interaction to increase turnaround time in
316 prospective diagnostic settings should be highlighted, as it may cause unintended harm to
317 patients when occurring in situations where diagnostic turnaround time is of utmost importance.
318 For the current application, the longer turnaround times for uncertain diagnoses (approximately
319 130 and 84 seconds for UP and UN diagnoses, respectively) suggest that a decrease in overall
320 turnaround time could be reached by reducing uncertainty in human-AI interaction during clinical
321 workflow integration.

322

323 Potential solutions for mitigating uncertainty might be to present the binarization threshold used
324 for *H. pylori* positivity, to display explicitly binarized model outputs, or, as previously discussed,

325 to set a probability threshold for presenting pathologists with model outputs. Another solution
326 might be to deploy the model, not as a primary diagnostic assistant, but as an automated pre-
327 screening tool, given its fast processing time, ability to correctly diagnose *H. pylori* on slides
328 where the bacteria are present, and lack of human diagnostic uncertainty. The model's accuracy
329 on the test set was 90%. In contrast, the pathologists' unassisted accuracy was 89.2%, when
330 P/UP diagnoses were counted as positive and N/UN diagnoses counted as negative (the
331 closest post-hoc approximation of a "no-uncertainty" scenario). The reported pathologist
332 sensitivity and specificity of *H. pylori* diagnosis on H&E stain ranges from 69-93% and 87-90%,
333 respectively.⁸ Our ensemble achieved a comparatively good sensitivity of 87.7% and a
334 specificity of 92.0% (higher than the upper range of reported pathologist specificities). If a
335 different operating point on the ensemble's ROC curve were selected to further maximize the
336 sensitivity while retaining acceptable specificity, cases flagged as negative by the ensemble
337 might be shifted to the end of the pathologist queue, while those flagged as positive might be
338 prioritized for review. Negative cases would no longer need to be painstakingly reviewed for *H.*
339 *pylori* or submitted for ancillary testing, while positive cases could be reviewed with the top-
340 ranked patches shown first, potentially reducing diagnostic turnaround time and ancillary
341 staining costs. Yet another deployment option might be to incorporate the ensemble into an
342 automated "double-check" tool, which could run in the background and alert the pathologist of
343 any disagreement between their diagnosis and the ensemble's prediction.

344
345 In this study, we tested the impact of the ensemble on subspecialty GI pathologists, as this is
346 the group that most commonly reviews gastric biopsies (and the only one that does so at our
347 institution). Given their high baseline unassisted accuracy (when uncertain diagnoses were
348 excluded, a total of only 2 incorrect diagnoses were made), it is somewhat expected that
349 assistance did not result in a statistically significant improvement in accuracy. While we did not
350 test the impact on non-GI pathologists, trainees, or other subgroups with less experience at the

351 task, it is possible that our ensemble could significantly improve accuracy and turnaround time
352 for these other subgroups. In practices where there is a shortage of GI-trained pathologists,
353 models such as the one in this study might provide value as a pre-screening or primary
354 diagnostic tool. To our knowledge, this is the first study to develop and test the impact of an AI
355 model for the direct detection of *H. pylori* organisms on H&E-stained whole-slide images.

356
357 Our study was subject to limitations. While the ensemble's accuracy was dependent on the
358 selected binarization threshold for *H. pylori* positivity, the pathologists' accuracy was dependent
359 on binarization of four possible diagnoses into correct and incorrect categories. Use of the UP
360 and UN categories precluded direct comparison of ensemble performance with that of the
361 pathologists. In theory, a direct comparison might be made if the ensemble's outputs were
362 thresholded into the same diagnostic choices available to the pathologists. In reality, the degree
363 of interobserver variability in the uncertainty threshold among different pathologists would make
364 it nearly impossible to establish a reliable threshold for the uncertain categories (whereas *H.*
365 *pylori*-positive and -negative diagnoses have a ground truth, there is none for uncertain
366 diagnoses).

367
368 Our ensemble was developed to help pathologists identify the presence of *H. pylori*, an
369 essential task which is performed, without exception, on every gastric biopsy. We acknowledge
370 that other concurrent pathologies may be present in the same biopsy, which are not currently
371 addressed by the ensemble, and which could be incorporated into future diagnostic suites
372 meant to assist with general gastric biopsy review.

373
374 Finally, our study was limited to data from a single pathology department serving a regional
375 healthcare system, with an imposed 50:50 class balance of positive and negative cases (chosen

376 to obtain broad representation of the morphologic range of cases, but which also reflects the
377 worldwide prevalence of *H. pylori* infection¹). Future studies incorporating datasets from multiple
378 institutions and regions, as well as more pathologists, are recommended to validate our
379 findings.

380

381 **Conclusions**

382 DL can diagnose *H. pylori* in H&E-stained gastric biopsies with high performance, and holds
383 potential for automating this common diagnostic task. However, contrary to prevailing
384 expectations regarding AI assistance, even a DL model with good performance metrics may fail
385 to improve human diagnostic performance, if it is not integrated into clinical workflows in an
386 optimal way. Although AI holds promise for improving healthcare quality and efficiency, its
387 ultimate clinical impact may be determined, not by model performance metrics, but by the
388 manner in which clinicians interact with these models. Our results suggest that increased
389 diagnostic uncertainty is an important unintended consequence of human-AI interaction, which
390 may decrease diagnostic accuracy and lead to longer case turnaround times. We hope that our
391 findings present a realistic picture of AI's impact on pathologists, while encouraging greater
392 attention toward addressing the various aspects of human-computer interaction that will
393 determine the ultimate real-world impact of these models.

394

395 **Author Contributions**

396 *Concept and design:* Zhou, Marklund, Ng, Shen

397 *Acquisition, analysis, or interpretation of data:* Zhou, Marklund, Blaha, Desai, Bingham, Martin,
398 Berry, Gomulia, Shen

399 *Drafting of the manuscript:* Zhou, Marklund, Blaha, Shen

400 *Manuscript revision and approval:* Zhou, Marklund, Blaha, Desai, Bingham, Martin, Berry,

401 Gomulia, Ng, Shen

402 *Statistical analysis:* Zhou, Marklund, Blaha, Desai, Shen

403 *Administrative, technical, or material support:* Zhou, Marklund, Gomulia, Ng, Shen

404 *Supervision:* Ng, Shen

405

406 **Acknowledgments**

407 We thank Norman L. Cyr (Graphic Arts & Imaging Services, Department of Pathology, Stanford
408 University) for his assistance in creating the figures for this article. This study was supported by
409 the Department of Pathology, Stanford University, with additional infrastructure provided by the
410 Stanford Machine Learning Group and the Stanford Center for Artificial Intelligence in Medicine
411 & Imaging.

412

413 **References**

- 414 1. Hooi JKY, Lai WY, Ng WK, et al. Global Prevalence of Helicobacter pylori Infection:
415 Systematic Review and Meta-Analysis. *Gastroenterology*. 2017;153(2):420-429.
- 416 2. Marshall BJ, Warren JR. Unidentified curved bacilli in the stomach of patients with gastritis
417 and peptic ulceration. *Lancet*. 1984;1:1311–1315.
- 418 3. International Agency for Research on Cancer. Schistosomes, liver flukes and Helicobacter
419 pylori, IARC Working Group on the Evaluation of Carcinogenic Risks to Humans, vol. 61. Lyon,
420 France: IARC; 1994.
- 421 4. Crowe SE. Indications and diagnostic tests for Helicobacter pylori infection. In Grover, S.
422 (Ed.), *UpToDate*. [https://www.uptodate.com/contents/indications-and-diagnostic-tests-for-](https://www.uptodate.com/contents/indications-and-diagnostic-tests-for-helicobacter-pylori-infection)
423 [helicobacter-pylori-infection](https://www.uptodate.com/contents/indications-and-diagnostic-tests-for-helicobacter-pylori-infection). Retrieved October 2, 2019.
- 424 5. Malfertheiner P, Megraud F, O’Morain CA, et al. Management of Helicobacter pylori infection
425 - the Maastricht V/Florence Consensus Report. *Gut*. 2017;66(1):6-30. doi: 10.1136/gutjnl-2016-
426 312288.

- 427 **6.** Chey WD, Wong BC, Practice Parameters Committee of the American College of
428 Gastroenterology. American College of Gastroenterology guideline on the management of
429 *Helicobacter pylori* infection. *Am J Gastroenterol*. 2007;102(8):1808.
- 430 **7.** Faigel DO, Childs M, Furth EE, Alavi A, Metz DC. New noninvasive tests for *Helicobacter*
431 *pylori* gastritis. Comparison with tissue-based gold standard. *Digest Dis Sci*. 1996;41(4):740-
432 748.
- 433 **8.** Lee JY, Kim N. Diagnosis of *Helicobacter pylori* by invasive test: Histology. *Ann Transl Med*.
434 2015;3(1):10.
- 435 **9.** 2019 Medicare Physician Fee Schedule. [https://documents.cap.org/documents/2019-final-](https://documents.cap.org/documents/2019-final-medicare-impact-table_181106_10474.pdf)
436 [medicare-impact-table_181106_10474.pdf](https://documents.cap.org/documents/2019-final-medicare-impact-table_181106_10474.pdf). Retrieved Jan 1, 2019.
- 437 **10.** Batts KP, Ketover S, Kakar S, et al. Appropriate use of special stains for identifying
438 *Helicobacter pylori*: Recommendations from the Rodger C. Haggitt Gastrointestinal Pathology
439 Society. *Am J Surg Pathol*. 2013;37(11):e12-e22. doi:10.1097/PAS.000000000000097
- 440 **11.** Pittman ME, Khararjian A, Wood LD, Montgomery EA, Voltaggio L. Prospective identification
441 of *Helicobacter pylori* in routine gastric biopsies without reflex ancillary stains is cost-efficient for
442 our health care system. *Hum Pathol*. 2016;58:90-96. doi:10.1016/j.humpath.2016.07.031
- 443 **12.** Salto-Tellez M, Maxwell P and Hamilton P. Artificial intelligence—the third revolution in
444 pathology. *Histopathology*. 2019;74: 372-376.
- 445 **13.** Bejnordi BE, Veta M, van Diest J, et al. Diagnostic assessment of deep learning algorithms
446 for detection of lymph node metastases in women with breast cancer. *JAMA*. 2017;318(22):
447 2199–2210.
- 448 **14.** Strom P, Kartasalo K, Olsson H, et al. Artificial intelligence for diagnosis and grading of
449 prostate cancer in biopsies: a population-based, diagnostic study. *Lancet Oncol*. 2020;21(2):
450 222-232.
- 451 **15.** Veta M, Heng YJ, Stathonikos N, et al. Predicting breast tumor proliferation from whole-slide
452 images: The TUPAC16 challenge. *Med Image Anal*. 2019;54:111-121.

- 453 **16.** Stalhammar G, Robertson S, Wedlund L, et al. Digital image analysis of Ki67 in hot spots is
454 superior to both manual Ki67 and mitotic counts in breast cancer. *Histopathology*. 2018;72(6):
455 974–989.
- 456 **17.** Steiner DF, MacDonald R, Liu Y, et al. Impact of deep learning assistance on the
457 histopathologic review of lymph nodes for metastatic breast cancer. *Am J Surg Pathol*. 2018;
458 42(12):1636–1646.
- 459 **18.** Kiani A, Uyumazturk B, Rajpurkar P, et al. Impact of a deep learning assistant on the
460 histopathologic classification of liver cancer. *NPJ Digital Med*. 2020;3:23; doi:10.1038/s41746-
461 020-0232-8
- 462 **19.** Bossuyt PM, Reitsma JB, Bruns DE, et al; STARD 2015: an updated list of essential items
463 for reporting diagnostic accuracy studies. *BMJ*. 2015;351:h5527. doi:10.1136/bmj.h5527.
- 464 **20.** He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. 2016 IEEE
465 Conference on Computer Vision and Pattern Recognition (CVPR). 2016:770-778.
- 466 **21.** Huang G, Liu Z, van der Maaten L, Weinberger KQ. Densely connected convolutional
467 networks. 2016: Preprint at <https://arxiv.org/abs/1608.06993>
- 468 **22.** Bankhead P, Loughrey MB, Fernandez JA, et al. QuPath: Open source software for digital
469 pathology image analysis. *Sci Rep*. 2017;7(1):16878.
- 470 **23.** Wilson EB. Probable inference, the law of succession, and statistical inference. *J AM Stat*
471 *Assoc*. 1927;22:209-212.
- 472 **24.** Bates D, Mächler M, Bolker B, Walker S. Fitting linear mixed-effects models using lme4. *J*
473 *Stat Softw*. 2015;67:1–48.
- 474 **25.** Venables WN, Ripley BD. *Modern Applied Statistics with S, Fourth edition*. New York, NY:
475 Springer; 2002.
- 476 **26.** Martin DR, Hanson JA, Gullapalli RR, Schultz FA, Sethi A, Clark DP. A deep learning
477 convolutional neural network can recognize common patterns of injury in gastric pathology. *Arch*
478 *Pathol Lab Med*. 2020;144(3):370-378.

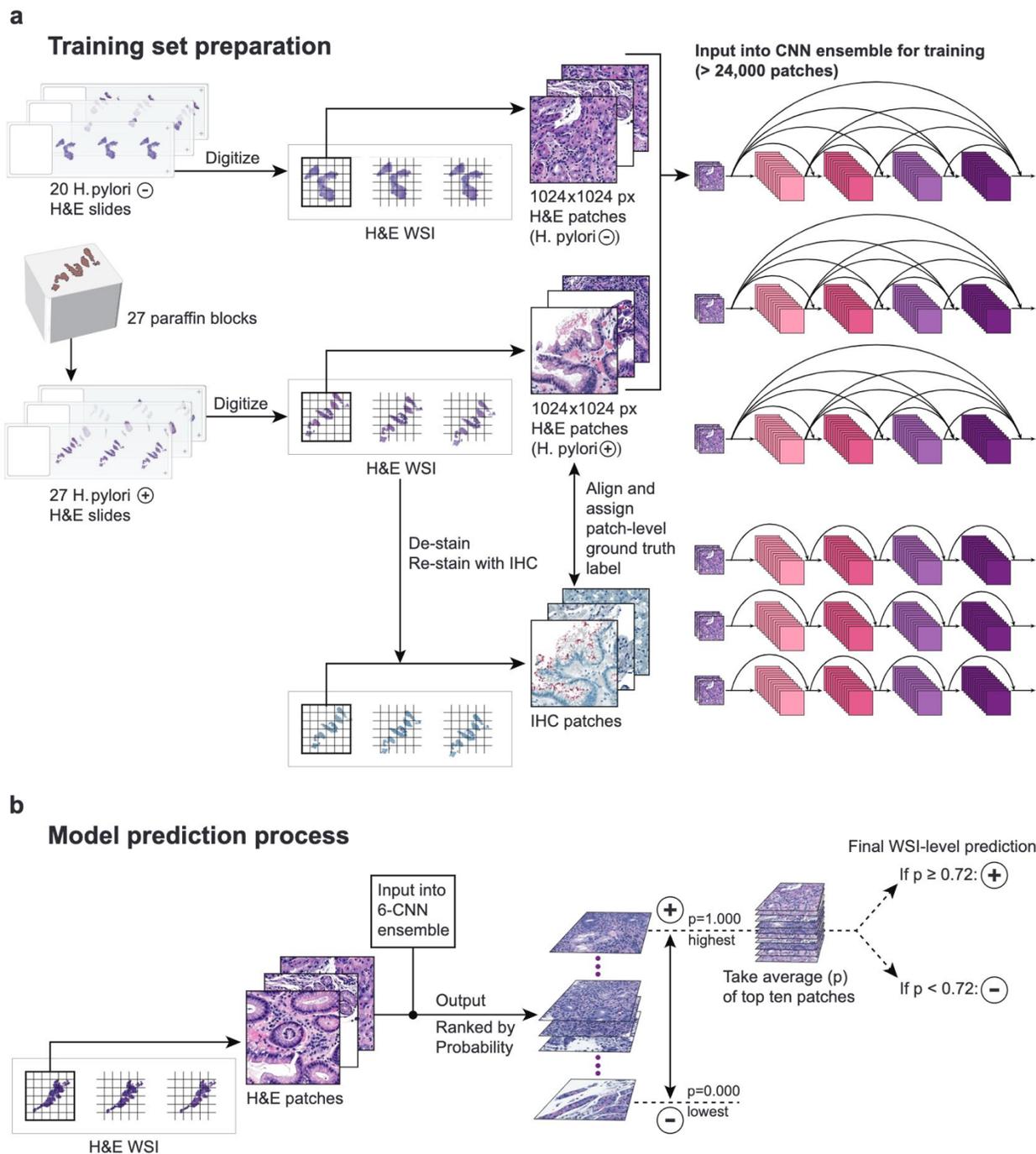
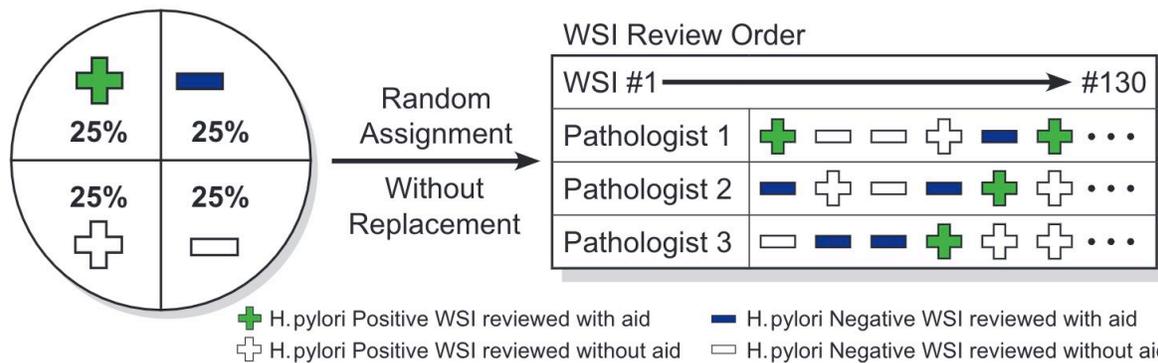


Figure 1. Training dataset preparation and model prediction process. **a**, The training dataset consisted of 24,786 non-overlapping image patches of size 1024 x 1024 pixels, extracted from one serial tissue section per H&E-stained slide (20 *H. pylori* negative and 27 *H. pylori* positive slides), input into a convolutional neural network (CNN) ensemble of 3 DenseNet-121 (top) and 3 ResNet-18 (bottom) architectures. To generate patch-level ground truth labels for the 27 *H. pylori*-positive slides, an H&E-stained slide was prepared from the paraffin block and digitized, then de-stained, re-stained with an *H. pylori* immunohistochemical (IHC) stain,

and digitized to generate a paired IHC WSI for tissue co-registration with the H&E WSI. **b**, The ensemble's WSI-level prediction of *H. pylori* status involved extracting all tissue-containing patches from a single level for input, computing the average of the patch-level probabilities for the top 10 highest probability patches, and binarizing this with a threshold of 0.72.

a
WSI assignment procedure



b
Whole slide image viewer interface

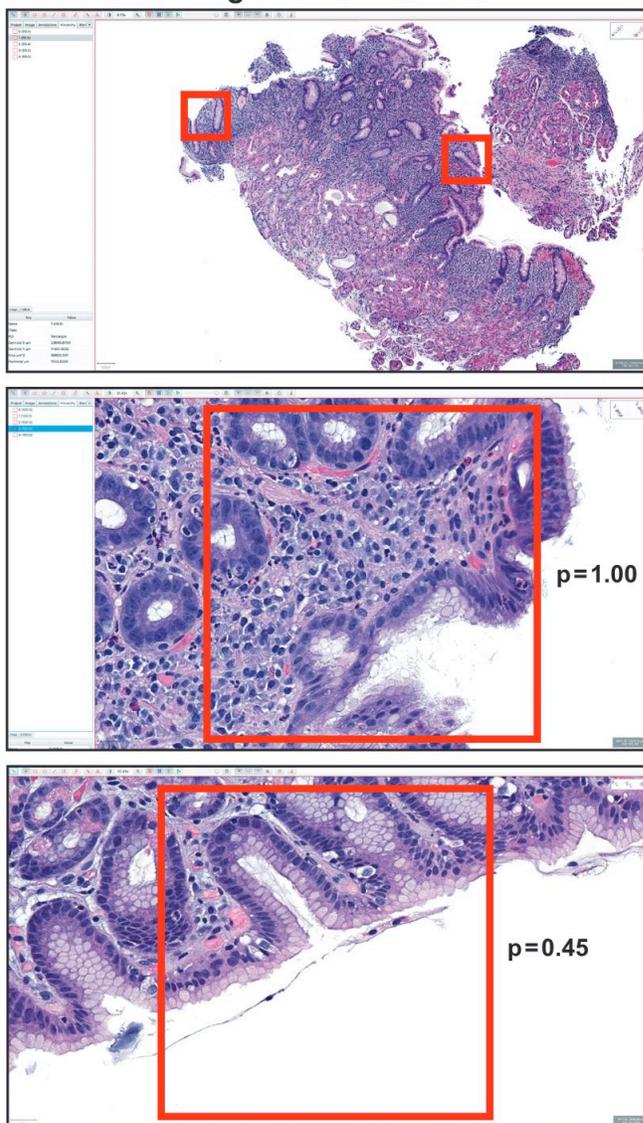


Figure 2. Design of the pathologist experiment and user interface. **a**, The 3 pathologists reviewed the same test set of 130 WSI (containing 65 positive and 65 negative WSI) during a single session, without time constraint. Half of the positive and negative WSI were randomly assigned to be reviewed with model assistance, with the remaining half assigned to be reviewed without assistance. Each pathologist received a unique randomized WSI review sequence and assignment of WSI to be reviewed with or without assistance. **b**, The WSI viewer interface consisted of a main WSI-viewing window (showing red bounding boxes around patches most likely to contain *H. pylori*) and a smaller window to the left displaying the list of all bounding boxes for that WSI, with their corresponding probabilities. By clicking on a bounding box in the smaller window, the user could automatically navigate to the corresponding region of the WSI in the main WSI-viewing window. The top panel shows a low magnification view, while the middle and bottom panels show higher magnification views with corresponding probabilities for two different bounding boxes from the same WSI. On the assisted cases, only the 3 to 5 bounding boxes that were most likely to contain *H. pylori* were displayed for each WSI. On unassisted cases, no bounding boxes or probabilities were displayed.

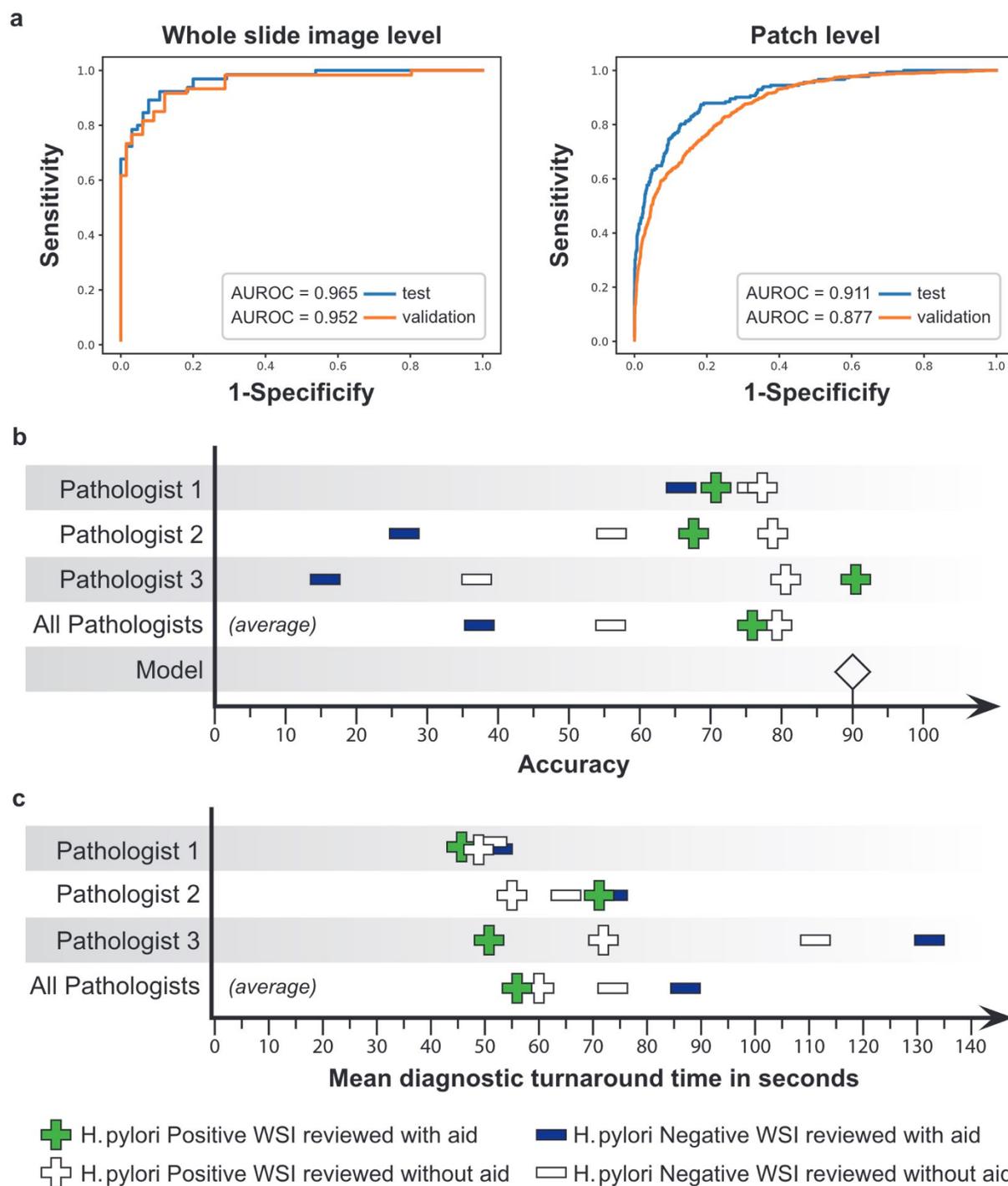
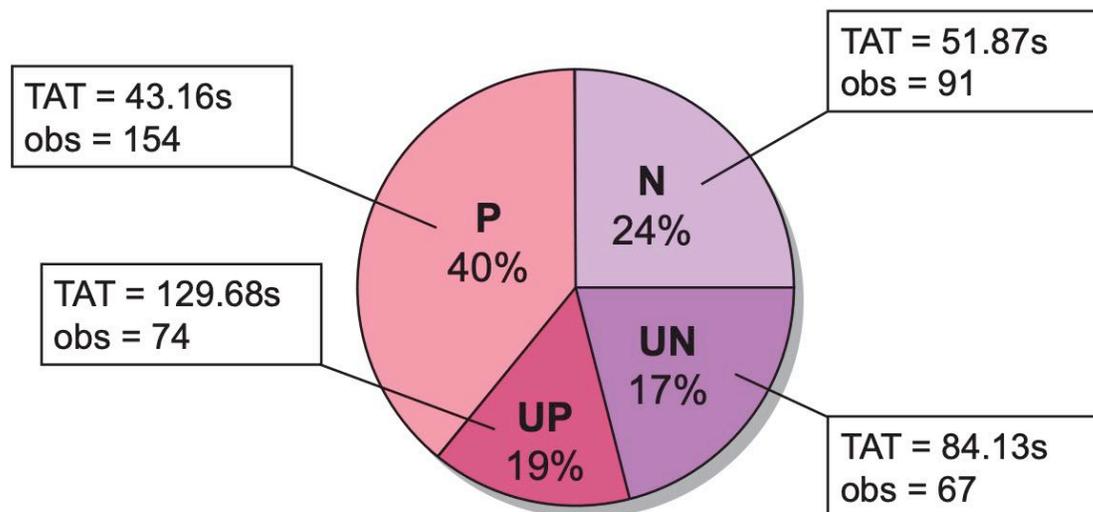


Figure 3. Performance of the model ensemble and impact of assistance on pathologist accuracy and diagnostic turnaround time. **a**, The WSI (left) and patch (right) level receiver-operating-characteristic (ROC) curves for the best-performing model ensemble on the validation (orange, n=126 WSI) and test (blue, n=130 WSI) datasets are shown, along with respective area under the curve (AUROC). **b**, The assisted and unassisted pathologist accuracies on *H. pylori* positive and negative WSI in the 130 WSI test set are shown. On *H. pylori*-positive WSI,

accuracies for individual pathologists ranged from 0.774-0.818 without assistance, and 0.688-0.906 with assistance. On *H. pylori* negative WSI, individual pathologist accuracies ranged from 0.375-0.758 unassisted and 0.161-0.677 assisted. The accuracy of the DL ensemble alone on the test set was 0.900. **c**, The average per-WSI diagnostic turnaround times (TAT) with and without assistance are shown. On *H. pylori*-positive WSI, individual pathologist TAT ranged from 49.7-73.7 seconds unassisted and 46.1-71.6 seconds assisted. On *H. pylori* negative WSI, individual TAT ranged from 47.8-110.8 seconds unassisted and 50.9-133.0 seconds assisted.

a



b

Definition of Correctness	With Aid		Without Aid		Total	
	Correct	Incorrect	Correct	Incorrect	Correct	Incorrect
1. Uncertain excluded	110 (99.1%)	1 (0.9%)	132 (98.5%)	2 (1.5%)	242 (98.8%)	3 (1.5%)
2. Uncertain = Correct	191 (99.5%)	1 (0.5%)	192 (99.0%)	2 (1.0%)	383 (99.2%)	3 (0.8%)
3. Positive = P, UP Negative = N, UN	162 (84.4%)	30 (15.6%)	173 (89.2%)	21 (10.8%)	335 (86.8%)	51 (13.2%)
4. Uncertain = Incorrect	110 (57.3%)	82 (42.7%)	132 (68.0%)	62 (32.0%)	242 (62.7%)	144 (37.3%)

Figure 4. Results of the pathologist experiment. **a**, The number and percentage of total observations (diagnoses made by the pathologists) and turnaround time (TAT) in seconds are shown for each diagnostic category, where P=Positive, N=Negative, UP=Uncertain Positive, UN=Uncertain Negative. **b**, The numbers and percentages of correct and incorrect diagnoses

made by the pathologists, given different possible definitions of correctness, are shown. For the analyses performed in the study, the last definition of correctness, where all Uncertain diagnoses were considered incorrect, was used.