

1 **Title Page**

2 Title: Enhanced physician performance when using an artificial intelligence model to detect
3 ischemic stroke on computed tomography

4

5 Authors:

6 James M Hillis, MBBS DPhil*^{1,2,3}

7 Bernardo C Bizzo, MD PhD*^{1,3,4}

8 Romane Gauriau, PhD¹

9 Christopher P Bridge, DPhil^{1,3,5}

10 John K Chin, MD¹

11 Buthaina Hakamy, MBBS¹

12 Sarah Mercaldo, PhD^{1,3,4}

13 John Conklin, MD^{3,4}

14 Sayon Dutta, MD MPH^{3,6}

15 William A Mehan, MD MBA^{1,3,4}

16 Robert W Regenhardt, MD PhD^{2,3}

17 Ajay Singh, MD^{3,4}

18 Aneesh B Singhal, MBBS MD^{2,3}

19 Jonathan D Sonis, MD MHCM^{3,6}

20 Marc D Succj, MD^{3,4}

21 Tianhao Zhang, PhD⁷

22 Bin Xing, PhD⁷

23 John F Kalafut, PhD⁷

24 Keith J Dreyer, DO PhD^{1,3,4}

25 Michael H Lev, MD^{3,4}

26 R Gilberto González, MD PhD^{3,4,5}

27

28 *These authors contributed equally to the work.

29

30 Author Affiliations:

31 ¹Data Science Office, Mass General Brigham, Boston, MA, USA

32 ²Department of Neurology, Massachusetts General Hospital, Boston, MA, USA

33 ³Harvard Medical School, Boston, MA, USA

34 ⁴Department of Radiology, Massachusetts General Hospital, Boston, MA, USA

35 ⁵Athinoula A. Martinos Center for Biomedical Imaging, Massachusetts General Hospital

36 ⁶Department of Emergency Medicine, Massachusetts General Hospital, Boston, MA, USA

37 ⁷GE Healthcare

38

39 Corresponding author:

40 Dr. James Hillis

41 Digital CRO, Data Science Office, Mass General Brigham

42 Suite 1303, Floor 13, 100 Cambridge St, Boston, MA 02114

43 james.hillis@mgh.harvard.edu

44 617-726-2000

45

46 Date of revision: January 16 2023

47 **Abstract (word count: 55)**

48 Acute ischemic stroke can be subtle to detect on non-contrast computed tomography imaging.

49 We show that a novel artificial intelligence model significantly improves the performance of

50 physicians, including ED physicians, neurologists and radiologists, in identifying and quantifying

51 the volume of acute ischemic stroke lesions. This model may lead to improved clinical decision-

52 making for stroke patients.

53 **Article (word count: 1,118)**

54 The early imaging features of acute ischemic stroke on non-contrast computed tomography
55 (CT) can be subtle. We previously reported the development of a deep learning model that
56 detects ischemic core on CT studies and was superior to three experienced neuroradiologists.¹
57 As part of its training, the model utilized segmentations obtained from the region of acute
58 infarct on paired magnetic resonance imaging (MRI) studies that were then registered onto the
59 CT studies. This design had recognized that MRI better detects early acute ischemic stroke but
60 CT is a cheaper, quicker and more widely available imaging modality.

61
62 In further evaluating this model, we wished to see how its use impacted physicians who were
63 interpreting non-contrast CT studies. We therefore designed a multi-reader multi-case study
64 whereby physicians interpreted 180 CT cases both with and without the use of the model
65 (Supplementary Figure 1). There were 8 physicians (2 emergency physicians, 2 emergency
66 radiologists, 2 neurologists, 2 neuroradiologists) who each interpreted 90 cases with the model
67 output and 90 cases without the model output (Figure 1A). After a four-week washout period,
68 they interpreted the cases in the opposite manner (i.e., the cases that they had previously
69 interpreted without the model output were now interpreted with the model output).

70
71 The model output included the binary classification of whether ischemic core $\geq 5\text{mL}$ was
72 present. If there was ischemic core $\geq 5\text{mL}$, the output also included the volume of the ischemic
73 core and the segmented region of the ischemic core. The physicians were asked up to three
74 multiple choice questions: the presence of any ischemic core (positive, negative), their level of

75 certainty for the presence or absence (scale of 1-7; see Supplementary Table 1 for details), and,
76 if they stated the presence of ischemic core, their volume estimate (0-20mL, 20-50mL, >50mL;
77 upper bound of range considered inclusive). Their responses were compared to the ground
78 truth interpretations demonstrated by diffusion imaging on a paired MRI case.

79
80 The physicians performed significantly better at binary detection when using the model
81 compared with not using the model. Their area under the receiver operating characteristic
82 curve (AUC) improved from 0.696 to 0.836 (difference 0.141; 95% CI: 0.081-0.200; $p < 0.001$;
83 Figure 1B and Table 1). Their sensitivity improved from 53.9% to 66.9% (difference 13.1%; 95%
84 CI: 6.5-19.6%; $p < 0.001$). Their specificity improved from 78.9% to 89.2% (difference 10.3%; 95%
85 CI: 6.6-14.0%; $p < 0.001$). Their interpretation time also improved from 62.30 seconds to 43.92
86 seconds (difference 18.38 seconds; 95% CI: 15.59-21.18 seconds; $p < 0.001$).

87
88 These improvements were maintained across most specialties as part of a subgroup analysis
89 (Table 1). The physicians with the greatest improvement were emergency physicians and
90 emergency radiologists, which is consistent with them having the least experience in
91 interpreting brain imaging and therefore having the greatest opportunity to benefit. The
92 physicians with the least improvement were neuroradiologists, who have the most experience.

93
94 The improvement in sensitivity was most pronounced for larger infarct volumes (Table 1). The
95 physicians' sensitivity improved from 54.2% without the model output to 84.6% with the model
96 output for 20-50mL infarcts, and 61.7% to 84.2% for >50mL infarcts. These large infarcts are

97 more likely to involve a large vessel occlusion, which means the patients could benefit from
98 treatment with endovascular thrombectomy.²⁻⁸ The increased ability to detect these infarcts
99 could prompt physicians to obtain CT angiography to detect a large vessel occlusion and cue the
100 physicians interpreting the CT angiography to the likely vessel involved. This increased ability to
101 detect may also lead to improved triage and sooner evaluation for endovascular
102 thrombectomy. A decreased time to thrombectomy has previously been shown to improve
103 outcomes.⁹

104
105 There was, however, a decline in sensitivity for infarct volumes 0-5mL (Table 1). This decreased
106 performance was expected given that the model only outputs whether it has detected ischemic
107 core $\geq 5\text{mL}$ (i.e., so the model should classify these cases as negative). The reason for the model
108 using this volume threshold is to avoid false positive interpretations from small regions of noise
109 on a CT study; the clinical concern is for such interpretations to increase MRI utilization to
110 evaluate for infarct more definitively. While we acknowledge the subsequent decreased
111 sensitivity amongst physicians for these infarct volumes, it is important to recognize that this
112 study occurred outside of the clinical environment where other factors, especially the acute
113 onset of neurologic symptoms, could alert physicians to the occurrence of ischemic stroke.

114
115 The physicians also performed better at volume quantification when using the model (Figure
116 1C). They correctly identified the volume range 73.8% of the time for $>50\text{mL}$ infarcts with the
117 model output compared to 27.1% without the model output. They correctly identified the

118 volume range 58.8% of the time for 20-50mL infarcts with the model output compared to
119 30.4% without the model output.

120
121 A key consideration moving forward is how the model might impact clinical workflow. Currently
122 the clinical paradigm is that an ischemic stroke should be assumed when a patient presents
123 with stroke-like symptoms and a non-contrast CT does not reveal an abnormality; the key
124 reason for obtaining the CT is to exclude intracranial hemorrhage. This study demonstrates the
125 model improves physicians' detection of ischemic core. Further research should be conducted
126 to determine whether the CT could be used more than it currently is for confirmation of
127 ischemic core. For instance, confirmation of ischemic core on non-contrast CT could be
128 particularly helpful when a patient's symptoms are ambiguous and ischemic stroke has not
129 been considered as a likely differential diagnosis. The presence of ischemic core could help
130 trigger and triage the next management steps such as CT angiography, administration of
131 thrombolytic medication and consideration of endovascular thrombectomy.¹⁰

132
133 The model most helped physicians with less experience in interpreting brain imaging and may
134 be most beneficial in rural areas with fewer subspecialty physicians. The model may separately
135 be most helpful at hospitals that can perform non-contrast CT but do not have emergent access
136 to advanced imaging modalities like CT perfusion or MRI. It may therefore assist in reducing the
137 urban-rural inequities in acute stroke care.¹¹

138

139 A limitation of this study was that the cases were taken from the test set from model
140 development. While these cases were sequestered and not exposed to the model during
141 development, they are likely to be the most similar to the training cases and therefore provide
142 the best model performance. The assessment of generalizability of the model will benefit from
143 evaluation on a more diverse dataset prior to clinical use. We note that the physicians were not
144 aware of the standalone model performance on these cases during this study.

145

146 Overall, this study demonstrates how the use of an artificial intelligence model enhances
147 physicians' identification and volume quantification of ischemic core on non-contrast CT. It
148 suggests that the model could provide benefit in the acute stroke clinical environment.

149 **Methods**

150 **Study design**

151 This retrospective multi-reader multi-case study was conducted using radiology cases from
152 hospitals within the Mass General Brigham network. It was approved by the Mass General
153 Brigham Institutional Review Board with waiver of informed consent per the Common Rule. It
154 was conducted in accordance with relevant guidelines and regulations including the Health
155 Insurance Portability and Accountability Act (HIPAA). This report followed the Standards for
156 Reporting Diagnostic Accuracy (STARD 2015) reporting guideline.

157

158 **Case selection and model inference**

159 The 180 cases were selected from the test set that had been used at the time of model
160 development.¹ As described previously, they had accompanying MRI studies that were used to
161 establish the ground truth interpretations (within 3 hours of the CT for positive cases and 5
162 days for negative cases). They were selected using stratified randomization from the entire
163 primary test set such that there were 30 cases with ischemic core 0-20mL, 30 cases with
164 ischemic core 20-50mL, 30 cases with ischemic core >50mL and 90 cases without any ischemic
165 core. This randomization ensured there were cases on which the model performed accurately
166 and inaccurately (see Supplementary Table 2 for standalone model performance on these
167 cases).

168

169 For the studies without the model output, the physicians were provided with only the axial
170 5mm series. For the studies with the model output, the physicians were provided with the axial

171 5mm series and an identical series with the model output incorporated into the imaging pixel
172 data (Figure 1A). The physicians were able to visualize these two series simultaneously in
173 adjacent window panes. The outputs included the binary classification of whether ischemic core
174 $\geq 5\text{mL}$ was present, and, if it was present, the volume of the ischemic core and the segmented
175 region of the ischemic core. The volume threshold of 5mL is proposed for future clinical use
176 given it optimizes specificity by avoiding false positive interpretations through incorrect
177 interpretation of small regions of noise on CT.

178

179 **Physicians**

180 The physicians were chosen to ensure representation of likely future clinical users including
181 emergency physicians, emergency radiologists, neurologists and neuroradiologists. They were
182 all board-certified for their relevant specialty. They were trained on the annotation tasks and
183 completed twelve training cases. They had not been involved with model development and
184 were not aware of prior model performance results. They were informed that specificity was
185 prioritized over sensitivity as part of the model design given that future clinical users should
186 similarly be aware of this fact.

187

188 **Reader study process**

189 The multi-reader multi-case study design involved the physicians interpreting all radiology
190 studies twice: both with and without the model output (Supplementary Figure 1). The
191 interpretations were performed as part of two sessions that were separated by at least four
192 weeks as recommended by the US Food and Drug Administration¹²; each study was interpreted

193 once in each session. Within each session, half of the studies were interpreted with the model
194 output and half of the studies without the model output; the studies were interpreted in the
195 opposite manner for the other session. The studies were split into two batches to facilitate
196 these interpretations. The order of studies within each session was randomized and differed for
197 each physician. The interpretation of the first batch of studies with or without model output for
198 each physician was also randomized.

199
200 The physicians were assisted in working through the cases by an internal web-based annotation
201 system that required them to interpret the cases in the defined order. This annotation system
202 incorporated the FDA-cleared eUnity image visualization software (Version 6 or higher). The
203 physicians were firstly asked about the presence of any ischemic core (positive, negative). They
204 were then asked their level of certainty for this ischemic core on a scale of 1-7 (see
205 Supplementary Table 1 for options). If they stated that ischemic core was present, they were
206 also asked to estimate the volume (0-20mL, 20-50mL, >50mL; the upper bound of these ranges
207 was considered inclusive). The annotation system also recorded the start and stop times for the
208 interpretation of each case. When submitting the annotations for each case, the physicians
209 could opt to move to the next case or exit; they did not need to interpret all cases in a single
210 sitting.

211 212 **Statistical techniques**

213 This study was a pilot multi-reader multi-case study for this model. The predefined primary
214 assessment was comparison of the AUC with and without the model outputs. The predefined

215 secondary assessments were comparison of sensitivity, specificity and interpretation time with
216 and without the model outputs. The specialty subgroup analysis, infarct volume subgroup
217 analysis and volume quantification analysis were calculated as exploratory assessments. Given
218 the pilot nature of this study, there was not an estimated effect size and powering was not
219 performed. There were no missing data with the exception of the excluded interpretation times
220 as described below.

221
222 The selection of cases involved multiple randomization procedures including for the cases
223 selected in each batch, the interpretation order for each physician, and whether the first batch
224 for a physician was with or without model output. This randomization used the Latin squares
225 methodology.

226
227 The analysis of the AUC was based on the 7-point scale assessments from each physician and
228 the ground truth interpretations based on the MRI. The overall comparison between the two
229 methods (with and without model outputs) was derived from the difference between the mean
230 AUCs¹³ via analysis of variance (ANOVA), taking into account the variability components:
231 methods, physicians, cases and their interactions. The analysis was implemented by using the
232 OR-DBM MRMC software package version 2.5 or higher.^{14,15} The ANOVA model treated both
233 physicians and cases as random samples, to be able to draw inferences for the whole
234 population of physicians and cases. Two-sided p-values were based on 95% confidence intervals
235 of the difference in AUCs.

236

237 The analysis of sensitivity and specificity used the binary interpretations from each physician
238 and the ground truth interpretations based on the MRI. The comparison between the two
239 methods (with and without model outputs) was based on the generalized estimating equations
240 (GEE) model¹⁶ by using SAS procedure PROC GENMOD, taking into account repeated
241 observations from the MRMC design. Based on SAS GENMOD, case was treated as the repeated
242 subject, accounting for correlations of physician and method within case, by using an
243 exchangeable covariance structure. Two-sided p-values were based on 95% confidence
244 intervals of the difference in sensitivity and specificity.

245
246 The analysis of interpretation time was performed in SAS using a repeated measures ANOVA.
247 Two-sided p-values were based on 95% confidence intervals of the difference in interpretation
248 time. Thirteen interpretation times (from 2880 interpretations) were excluded; the paired study
249 for the same physician (i.e., the equivalent case with or without model output) was also
250 excluded to ensure balance of such exclusions. These exclusions occurred for two reasons.
251 Firstly, the physicians could notify study management if an interpretation time should be
252 excluded (e.g., they were interrupted). Secondly, the annotation system initially recorded the
253 start time for each case incorrectly; to ensure consistency across the entire cohort, the
254 interpretation time was rederived by calculating the difference in start time with the
255 subsequent case or the difference in stop time with the previous case (the minimum of these
256 times was taken); the interpretation times were excluded if this derived number did not appear
257 consistent with an expected duration (e.g., it appeared a physician only interpreted one case in
258 a sitting so the difference in times with the previous and subsequent cases would be incorrect).

259

260 The analysis of volume outputs used the volume estimates from each physician and the ground
261 truth interpretations based on the MRI. The frequencies of how the volume estimates matched
262 with ranges of ground truth volumes of the two methods (with and without model outputs)
263 were calculated using Microsoft Excel.

264 **Acknowledgements**

265 The authors thank the broader Mass General Brigham Data Science Office and GE Healthcare
266 teams for their assistance with this project.

267

268 **Funding**

269 This study was funded by GE Healthcare. JMH, BCB, RG, CPB, JKC, BH, SM, JC, SD, WAH, RWR,
270 AS, ABS, JDS, MDS, KJD, MHL, RGG were employees of Mass General Brigham and/or
271 Massachusetts General Hospital at the time of this study, which had received institutional
272 funding from GE Healthcare for the study. TZ, BX, JFK were employees of GE Healthcare at the
273 time of this study.

274

275 **Competing interests**

276 RWR reports the following competing interests: Rapid Medical – Clinical Trial DSMB;
277 Microvention – Site PI; Penumbra – Site PI; National Institute of Neurological Disorders and
278 Stroke – Research Grant; Society of Vascular and Interventional Neurology – Research Grant;
279 Heitman Foundation – Research Grant.

280

281 **Data availability**

282 The data used were obtained from hospitals within the Mass General Brigham network. Data
283 use was approved by relevant Institutional Review Board. The data are not publicly available
284 and restrictions apply to their use.

285

286 **Author contributions**

287 JMH, BCB, JKC, TZ, JFK, KJD, MHL and RGG conceptualized the study. JMH, BCB, RG, CPB, JKC

288 and BH conducted the study including producing the model output. JC, SD, WAH, RWR, AS, ABS,

289 JDS and MDS were physician readers. SM, TZ and BX performed the randomization and

290 statistical analysis. JMH drafted the manuscript text and figures. All authors reviewed the

291 manuscript.

292 References

- 293 1 Bizzo, B. C. *et al.* *Head CT Deep Learning Model for Early Stroke Identification*
294 *Outperforms Human Experts* (Research Square, 2021).
- 295 2 Berkhemer, O. A. *et al.* A randomized trial of intraarterial treatment for acute ischemic
296 stroke. *N Engl J Med* **372**, 11-20, doi:10.1056/NEJMoa1411587 (2015).
- 297 3 Goyal, M. *et al.* Randomized assessment of rapid endovascular treatment of ischemic
298 stroke. *N Engl J Med* **372**, 1019-1030, doi:10.1056/NEJMoa1414905 (2015).
- 299 4 Saver, J. L. *et al.* Stent-retriever thrombectomy after intravenous t-PA vs. t-PA alone in
300 stroke. *N Engl J Med* **372**, 2285-2295, doi:10.1056/NEJMoa1415061 (2015).
- 301 5 Campbell, B. C. *et al.* Endovascular therapy for ischemic stroke with perfusion-imaging
302 selection. *N Engl J Med* **372**, 1009-1018, doi:10.1056/NEJMoa1414792 (2015).
- 303 6 Jovin, T. G. *et al.* Thrombectomy within 8 hours after symptom onset in ischemic stroke.
304 *N Engl J Med* **372**, 2296-2306, doi:10.1056/NEJMoa1503780 (2015).
- 305 7 Nogueira, R. G. *et al.* Thrombectomy 6 to 24 Hours after Stroke with a Mismatch
306 between Deficit and Infarct. *N Engl J Med* **378**, 11-21, doi:10.1056/NEJMoa1706442
307 (2018).
- 308 8 Albers, G. W. *et al.* Thrombectomy for Stroke at 6 to 16 Hours with Selection by
309 Perfusion Imaging. *N Engl J Med* **378**, 708-718, doi:10.1056/NEJMoa1713973 (2018).
- 310 9 Jahan, R. *et al.* Association Between Time to Treatment With Endovascular Reperfusion
311 Therapy and Outcomes in Patients With Acute Ischemic Stroke Treated in Clinical
312 Practice. *JAMA* **322**, 252-263, doi:10.1001/jama.2019.8286 (2019).
- 313 10 Powers, W. J. *et al.* Guidelines for the Early Management of Patients With Acute
314 Ischemic Stroke: 2019 Update to the 2018 Guidelines for the Early Management of
315 Acute Ischemic Stroke: A Guideline for Healthcare Professionals From the American
316 Heart Association/American Stroke Association. *Stroke* **50**, e344-e418,
317 doi:10.1161/STR.0000000000000211 (2019).
- 318 11 Hammond, G., Luke, A. A., Elson, L., Towfighi, A. & Joynt Maddox, K. E. Urban-Rural
319 Inequities in Acute Stroke Care and In-Hospital Mortality. *Stroke* **51**, 2131-2138,
320 doi:10.1161/STROKEAHA.120.029318 (2020).
- 321 12 US Food and Drug Administration. *Clinical Performance Assessment: Considerations for*
322 *Computer-Assisted Detection Devices Applied to Radiology Images and Radiology Device*
323 *Data in Premarket Notification (510(k)) Submissions*,
324 <https://www.fda.gov/media/77642/download> (2022).
- 325 13 Hillis, S. L., Obuchowski, N. A., Scharzt, K. M. & Berbaum, K. S. A comparison of the
326 Dorfman-Berbaum-Metz and Obuchowski-Rockette methods for receiver operating
327 characteristic (ROC) data. *Stat Med* **24**, 1579-1607, doi:10.1002/sim.2024 (2005).
- 328 14 Hillis, S. L., Berbaum, K. S. & Metz, C. E. Recent developments in the Dorfman-Berbaum-
329 Metz procedure for multireader ROC study analysis. *Acad Radiol* **15**, 647-661,
330 doi:10.1016/j.acra.2007.12.015 (2008).
- 331 15 Hillis, S. L. & Scharzt, K. M. Multireader sample size program for diagnostic studies:
332 demonstration and methodology. *J Med Imaging (Bellingham)* **5**, 045503,
333 doi:10.1117/1.JMI.5.4.045503 (2018).

334 16 *Generalized Estimating Equations,*
335 <https://support.sas.com/rnd/app/stat/topics/gee/gee.pdf> (<
336

337 **Figure Legends**

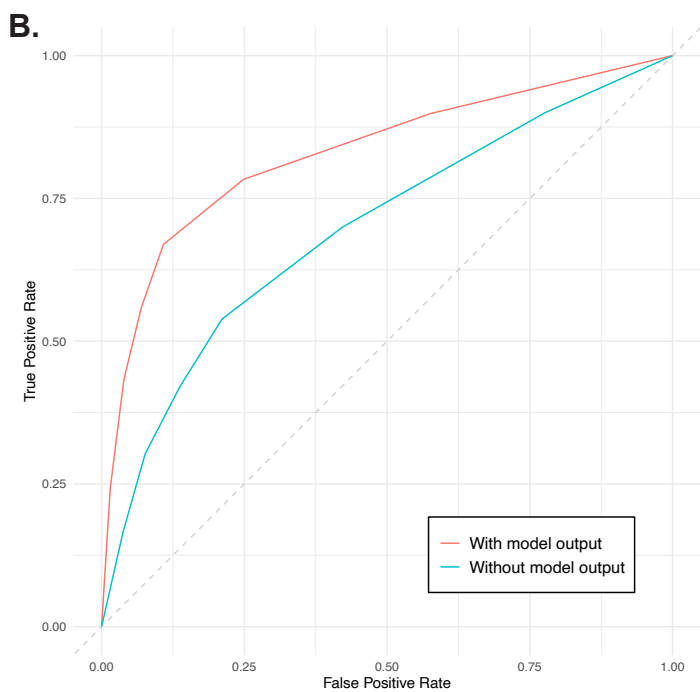
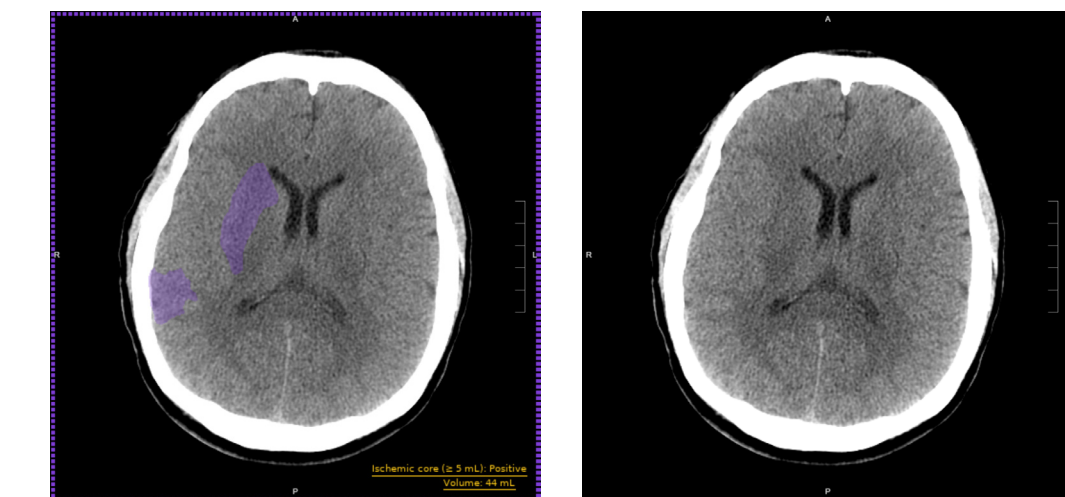
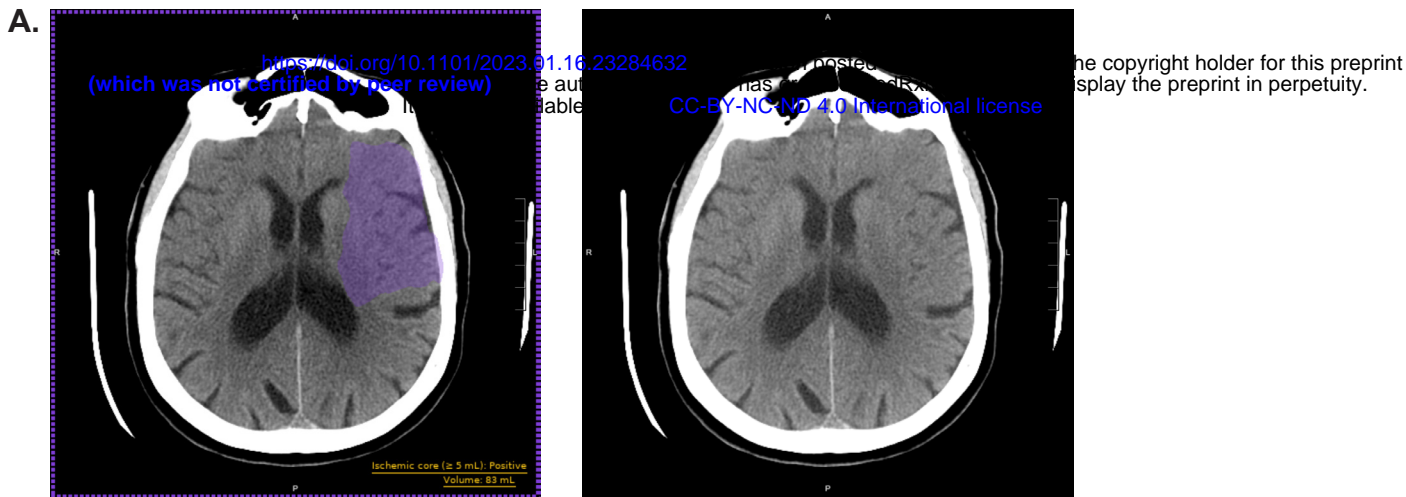
338 **Figure 1: A.** Example of two cases with and without model output showing text incorporating
339 binary classification for ischemic core $\geq 5\text{mL}$ and volume, and purple segmented region. The first
340 case, with infarct volume 77mL on MRI, was detected by 1 out of the 8 physicians without
341 model output and 6 with model output. The second case, with infarct volume 45mL on MRI,
342 was detected by 2 without model output and 8 with model output. **B.** Receiver operating
343 characteristic curves for both with model output and without model output. **C.** Confusion
344 matrices comparing physician volume with ground truth volume; a correct volume estimate
345 occurs on the diagonal; a volume of 0mL reflects a negative case.

346 **Tables**

347 **Table 1:** Results summary for physician performance comparing interpretations performed with
 348 or without the model outputs.
 349

	Metric _{with}	Metric _{without}	Difference (95% CI)	p value
AUC				
Overall	0.836	0.696	0.141 (0.081 to 0.200)	<0.001
ED Physician	0.820	0.632	0.188 (0.119 to 0.256)	<0.001
ED Radiologist	0.805	0.617	0.188 (0.098 to 0.279)	0.003
Neurologist	0.854	0.738	0.116 (0.081 to 0.151)	0.001
Neuroradiologist	0.866	0.796	0.070 (-0.337 to 0.477)	0.328
Sensitivity				
Overall	66.9%	53.9%	13.1% (6.5 to 19.6%)	<0.001
ED Physician	64.4%	43.9%	20.6% (10.4% to 30.8%)	<0.001
ED Radiologist	74.4%	48.9%	25.6% (16.0% to 35.1%)	<0.001
Neurologist	70.6%	69.4%	1.1% (-8.0% to 10.3%)	0.814
Neuroradiologist	58.3%	53.3%	5.0% (-1.0% to 11.4%)	0.123
Specificity				
Overall	89.2%	78.9%	10.3% (6.6 to 14.0%)	<0.001
ED Physician	91.1%	78.3%	12.8% (6.1% to 19.4%)	<0.001
ED Radiologist	76.7%	72.8%	3.9% (-4.7% to 12.5%)	0.376
Neurologist	91.1%	69.4%	21.7% (15.5% to 27.8%)	<0.001
Neuroradiologist	97.8%	95.0%	2.8% (-0.4% to 6.0%)	0.090
Interpretation time				
Overall	43.92s	62.30s	-18.38s (-21.18s to -15.59s)	<0.001
ED Physician	40.74s	61.91s	-21.17s (-25.95s to -16.38s)	<0.001
ED Radiologist	42.01s	70.58s	-28.57s (-37.03s to -20.10s)	<0.001
Neurologist	40.59s	60.79s	-20.20s (-24.32s to -16.07s)	<0.001
Neuroradiologist	52.33s	55.94s	-3.60s (-7.12s to -0.09s)	0.044
Sensitivity for different infarct volumes (based on MRI ground truth)				
0-5mL	16.0%	34.7%	-18.8% (-28.2% to -9.3%)	<0.001
5-20mL	56.3%	62.5%	-6.3% (-21.1% to 8.6%)	0.410
20-50mL	84.6%	54.2%	30.4% (23.6% to 37.3%)	<0.001
>50mL	84.2%	61.7%	22.5% (14.1% to 30.9%)	<0.001

350



C.

With model output

		N	Physician volume			
			0mL	0-20mL	20-50mL	>50mL
Ground truth (MRI) volume	0mL	720	89.2%	6.4%	3.8%	0.7%
	0-20mL	240	67.9%	20.4%	9.6%	2.1%
	20-50mL	240	15.4%	16.7%	58.8%	9.2%
	>50mL	240	15.8%	2.5%	7.9%	73.8%

Without model output

		N	Physician volume			
			0mL	0-20mL	20-50mL	>50mL
Ground truth (MRI) volume	0mL	720	78.9%	10.7%	6.9%	3.5%
	0-20mL	240	54.2%	26.7%	16.3%	2.9%
	20-50mL	240	45.8%	17.1%	30.4%	6.7%
	>50mL	240	38.3%	16.7%	17.9%	27.1%