

1 **TITLE**

2
3 **'Deep learning of HIV field-based rapid tests'**

4
5 **AUTHOR LIST**

6
7 Valérian Turbé¹, Carina Herbst², Thobeka Mngomezulu², Sepehr Meshkinfamfard¹,
8 Nondumiso Dlamini², Themrani Mhlongo², Theresa Smit², Valeriia Cherepanova³, Koki
9 Shimada³, Jobie Budd^{1,4}, Nestor Arsenov¹, Steven Gray⁵, Deenan Pillay^{2,6}, Kobus Herbst^{2,7},
10 Maryam Shahmanesh^{2,8}, Rachel A. McKendry^{1,4}

11
12 Corresponding authors:

13 R. A. McKendry (r.a.mckendry@ucl.ac.uk)

14 Valérian Turbé (v.turbe@ucl.ac.uk)

15 Maryam Shahmanesh (m.shahmanesh@ucl.ac.uk)

16 Kobus Herbst (Kobus.Herbst@ahri.org)

17
18 **AFFILIATIONS**

19
20 ¹London Centre for Nanotechnology, University College London, 17-19 Gordon Street,
21 London WC1H 0AH, UK

22 ²Africa Health Research Institute, K-RITH Tower Building, Nelson R. Mandela Medical
23 School, 719 Umbilo Rd, Umbilo, Durban, 4001, South Africa

24 ³Department of Computer Science, University College London, Gower St, Bloomsbury,
25 London WC1E 6EA, UK

26 ⁴Division of Medicine, Rayne Building, University College London, 5 University Street,
27 London, WC1E 6JF, UK

28 ⁵UCL Centre for Advanced Spatial Analysis, Gower Street, London, WC1E 6BT, UK

29 ⁶Division of Infection and Immunity, UCL Cruciform Building, University College London,
30 Gower Street, London, WC1E 6BT, UK

31 ⁷DSI-MRC South African Population Research Infrastructure Network, 491 Peter Mokaba
32 Ridge Road, Durban, South Africa

33 ⁸Institute for Global Health, University College London, Mortimer Market Centre, off Capper
34 Street, London WC1E 6JB, UK

35
36 **Key words:** Rapid diagnostic test, deep learning, HIV

37
38 **INTRODUCTORY PARAGRAPH**

39
40 Although deep learning algorithms show increasing promise for disease diagnosis, their use
41 with rapid diagnostic tests performed in the field has not been extensively tested. Here, we
42 use deep learning to classify images of rapid HIV tests acquired in rural South Africa. Using
43 newly developed image capture protocols with the Samsung SM-P585 tablet, 60
44 fieldworkers routinely collected images of HIV lateral-flow tests. From a library of 11,374
45 images, deep learning algorithms were trained to classify tests as positive or negative. A
46 pilot field study of the algorithms deployed as a mobile application demonstrated high levels
47 of sensitivity (97.8%) and specificity (100%), compared to traditional visual interpretation by
48 humans - experienced nurses and newly trained community health worker staff - and
49 reduced the number of false positives and false negatives. Our findings lay the foundations
50 for a new paradigm of deep learning-enabled diagnostics in low- and middle-income
51 countries, termed REASSURED diagnostics¹, for Real-time connectivity, Ease of specimen
52 collection, Affordable, Sensitive, Specific, User-friendly, Rapid, Equipment-free, and

53 Deliverable. Such diagnostics have the potential to provide a platform for workforce training,
54 quality assurance, decision support, and mobile connectivity to inform disease control
55 strategies, strengthen healthcare system efficiency, and improve patient outcomes and
56 outbreak management of emerging infections.

57
58

MAIN TEXT

59 Rapid diagnostic tests (RDTs) save lives by informing case management, treatment,
60 screening, disease control and elimination programmes¹. Lateral flow tests are among the
61 most common RDTs and hundreds of millions of these tests are performed worldwide each
62 year. They have the potential to support near person testing and decentralised management
63 of a range of clinically important diseases (including malaria, HIV, syphilis, tuberculosis,
64 influenza and non-communicable diseases²), making it convenient for the end-user and
65 more affordable for health systems³. RDT also present some issues, namely: errors in
66 performing the test and interpreting the result^{4,5}, quality control, and lack of electronic data
67 capture records of the test and results within health systems and surveillance. Many of these
68 would be overcome with the 'R' in REASSURED - the new criteria for an ideal test to reflect
69 the importance of digital connectivity, coined by Peeling and coworkers¹. The 'R' stands for
70 'real-time connectivity' using mobile phone connected RDTs. To date there have been few
71 peer reviewed studies or evaluations of the effectiveness of connected lateral flow tests at
72 scale in populations in need in low- and middle-income countries.

73 Recent studies that compare the human interpretation of a HIV RDT to various gold
74 standards, such as Western Blot⁶⁻⁹, Enzyme Immunoassay^{7,9-11}, standardised test panels¹²
75 or different HIV RDTs¹³⁻¹⁵, have highlighted the common issue of subjective interpretation of
76 the test result, which can lead to incorrect diagnosis. User error (especially in the case of
77 weak reactive lines) and inadequate supervision of testers were identified as prime factors
78 for misinterpretation¹⁶. In a study of differently experienced users interpreting results of HIV
79 RDTs by looking at pictures of tests¹⁷, the accuracy of interpretation varied between 80%
80 and 97%. This highlights the importance of experience in reading the test, as well as the
81 subjectivity involved in reading a weak test line. Evidence also suggests that some
82 fieldworkers struggle to interpret RDTs because of colour blindness or short-sightedness.¹⁸
83 Another study used photographs of HIV RDTs to quantify the subtle difference in tests with
84 faint lines declared as True- or False-positive by a panel of human users¹⁹. While these were
85 small-scale studies (N = 148 and 8, respectively), both highlighted the potential for
86 photographs to improve quality control and decision-making.

87

88 Deep learning algorithms, harnessing advances in large data sets and processing power,
89 have recently shown the ability to exceed human performance in a plethora of visual
90 tasks, including cell-based diagnostics²⁰, interpreting dermatology²¹, ophthalmology²² and
91 radiography images²³, playing strategic games²⁴, and in clinical medicine when used
92 alongside appropriate guidelines^{25,26}. While some studies are emerging looking at
93 applying deep learning to the interpretation of RDT^{27,28}, little is known about the ability of
94 machine learning models to analyse field-acquired diagnostic test data, with concerns about
95 the potential uniformity of images (e.g. focus, tilt), harsh environmental factors such as
96 lighting, and the variety of test types. In addition, there is a general lack of large real-world
97 datasets available to successfully train deep learning classifiers, particularly from low- and
98 middle-income countries. Recent advances in consumer electronic devices and deep
99 learning, have the potential to improve RDT quality assurance, staff training and
100 connectivity, eventually supporting self-testing, such as HIV-self-testing, which has been
101 shown to be cost-effective²⁹, to appeal to young people³⁰ and help reduce anxiety³¹.

102

103 Mobile health (mHealth) approaches, which marry RDTs with widely available mobile
104 phones, take advantage of inbuilt sensors (e.g. cameras) found in the phones, battery life,

105 processing power, screens to display results, and connectivity to send results to health
106 databases. A recent field study has shown high levels of acceptability for a device
107 sending HIV RDT results to online data bases in real-time³². An array of approaches have
108 been piloted at small scales ($N \leq 283$) and have shown good performance. However,
109 most require a physical attachment, such as a dongle (92-100% sensitivity, 97-100%
110 specificity)³³, a cradle³⁴, or a portable reader (97-98% sensitivity)³⁵, which increases cost
111 and complexity, and typically rely on simple image analysis software.

112
113 We explore the potential of deep learning algorithms to classify field-based RDT images as
114 either positive or negative, focusing on HIV as an exemplar, and piloting at scale in
115 population 'test beds' in KwaZulu-Natal, typical of semi-rural settings in Sub-Saharan Africa.
116 Figure 1 shows the concept of our deep learning-enabled REASSURED diagnostic system
117 to capture and interpret RDT results. Our approach first involved building a large image
118 library of field-acquired test images as training data set, optimising algorithms for high
119 sensitivity and specificity, and then to deploy our classifier in a pilot study to assess its
120 performance compared to traditional visual interpretation with a range of end users with
121 varying levels of training.

122
123 Our standard image collection protocol (Figure 2a) and library are described in the Methods
124 section. In brief, 11,374 photographs of HIV RDT were captured by over 60 fieldworkers
125 using Samsung tablets (SM-P585, 8Megapixel camera, f1/9, with autofocus capability).
126 Embedding routine image collection into staff workflows was acceptable and feasible, and
127 participant consent rate was 96%. We optimised our mHealth system for the two different
128 HIV RDTs used in the study as part of routine household population surveillance. At first
129 glance these RDTs appear similar but have different features and number of test lines. To
130 reduce the number of variables, we cropped the images around the region of interest (ROI)
131 (Figure 2b). Figure 2c shows a snapshot of the very diverse real-world field conditions where
132 the images were captured (indoors, outdoors, in the shade and in direct sunlight).

133
134 Each image was labelled (see Online Methods) according to the test result. Figure 3a details
135 the number of images used to train classifiers to automatically read the result of HIV RDT
136 images. The training process is described in the Online Methods section. In order to test the
137 reproducibility of the process, we performed a 10-fold cross validation. As can be seen in
138 Figure 3b, the average sensitivity (95.9% \pm 5.1 for type A, 98.7% \pm 1.7 for type B) and
139 specificity (99.0% \pm 0.6 for type A, 99.8% \pm 0.2 for type B) achieved across the 10 folds was
140 high and consistent for both types of HIV RDT. We therefore used all the available data to
141 train a final classifier for each type of test, which were used in our field study. We
142 investigated different common classification methods being used for clinical diagnostic
143 (Support Vector Machine³⁶ (SVM) and Convolutional Neural Networks (CNN)) including 3
144 different CNN architectures (ResNet50³⁷, MobileNetV2^{38,39} and MobileNetV3⁴⁰), and found
145 MobileNetV2 was the most appropriate for our task, as can be seen in Figure 3c.

146
147 We then conducted a field pilot study in rural South Africa to assess the performance of our
148 mHealth system compared to visual interpretation with a range of end-users with varying
149 levels of training (see Online Methods). Five participants (2 nurses, 3 newly trained
150 community healthworkers) were each asked to give their interpretation of 40 HIV RDTs and
151 to acquire a photograph of the RDT via the app. All five participants (100%) were able to use
152 our mHealth system without training, demonstrating its feasibility and acceptability. The
153 photographs were then evaluated by an expert RDT interpreter, followed by our deep
154 learning algorithms on a secure server. The results were not fed back to the study
155 participants to avoid confirmation bias. The performance results can be seen in Figure 4.

156
157 When comparing the traditional visual interpretation of the RDTs, we observed varied levels
158 of agreement between participants, (61-100%) as can be seen in Figure 4a. As expected,
159 agreement between nurses (N1 & N2: 100% and 94.4% agreement for test types A and B

160 respectively) was greater than between newly trained community health workers (C1, C2 &
161 C3: 80-90% and 61.1-94.4% for test types A and B, respectively). Test type B showed the
162 lower level of agreement. The low level of agreement between participants, and variability
163 due to the type of HIV RDT, were of concern and highlighted the need for a more objective
164 and consistent method to interpret HIV RDTs in the field. The confusion matrices in Figure
165 4b, demonstrate our mHealth system reduced the number of errors in reading RDTs. The
166 number of False Positive results from our mHealth system was found to be significantly
167 lower than for the traditional visual interpretation (0 compared to 11 – the largest variation
168 being observed for community health workers, 10), which translates as an improvement in
169 specificity from 89% to 100%, and an improvement in Positive Predictive Value from 88.7%
170 to 100%. Similarly, the number of False Negative results was just two in our mHealth
171 system, compared to four in traditional visual interpretation, which translates as an
172 improvement in sensitivity from 95.6% to 97.8%, and an improvement in Negative Predictive
173 Value from 95.7% to 98%. We plotted the ratio of our mHealth system performance to the
174 participant performance, both for sensitivity and specificity (Figure 4c). All participants had a
175 sensitivity index equal or greater than one for test type A; four out of five participants (N1,
176 N2, C1, C2) also did for test type B, demonstrating our mHealth system was better than
177 those participants at reading positive test results. Our system was also more reliable at
178 reading negative tests, as all participants had a specificity index equal or greater than one
179 for both types of HIV RDTs.

180

181 We acknowledge the following limitations of our study. Firstly, our pilot study involved a
182 relatively small number of participants (five), although we note this is comparable to other
183 similar pilot studies reported in the field. In future, larger evaluation studies and clinical trials
184 are needed to assess the performance of the system, involving participants with a broader
185 range of demographics including age, gender and different levels of digital literacy, as well
186 as more expert readers. In addition, future studies would benefit from including an invalid
187 test classifier and different mobile phone types with varying camera specifications. The
188 images were analysed on a secure server, however, future analysis could be on-device
189 overcoming the need to upload images. We are also currently investigating a picture
190 segmentation approach using deep learning for the next iteration of the smartphone
191 application.

192

193 To conclude, we demonstrated the potential of deep learning to accurately classify RDT
194 images, with an overall performance of 98.9% accuracy, significantly higher than traditional
195 visual interpretation of study participants (92.1%), which are comparable with reports of 80-
196 97% accuracy¹⁷. Given that over 100 million HIV tests are performed annually, even a small
197 improvement in quality assurance could impact the lives of millions of people by reducing the
198 risk of false positives and negatives. To the best of our knowledge our real-world image
199 library is the first of its kind at this scale and we demonstrate that deep learning models can
200 be deployed in mobile devices in the field, without the need for cradles, dongles or other
201 attachments. It lays the foundation for deep learning enabled REASSURED diagnostics,
202 demonstrating that RDTs linked to a mobile device could standardise capture and
203 interpretation of test results for decision-makers, reducing interpretation and transcription
204 errors and workforce training. Our findings are based on HIV testing decision support for
205 fieldworkers, nurses and community health workers, but in future could be applicable to
206 decision support for self-testing. We focused on HIV as an exemplar, but the capacity of the
207 classifier to adapt to two different test types suggests that it is amenable to a large range of
208 RDTs spanning communicable and non-communicable diseases. This platform could be
209 utilised for workforce training, quality assurance, decision support, and mobile connectivity to
210 inform disease control strategies, strengthen healthcare systems efficiency, and improve
211 patient outcomes, and outbreak management. The ideal connected system would link to
212 connected RDTs to laboratory systems, whereby remote monitoring of RDT functionality and
213 utilisation could also allow health programmes to optimise testing deployment and supply
214 management to deliver the Sustainable Development Goals and ensure no one is left

215 behind. The real-time alerting capability of connected RDTs could also support public health
216 outbreak management, by mapping ‘hotspots’ for epidemics including COVID-19 to protect
217 populations.

218

219 REFERENCES

220

- 221 1. Land, K. J., Boeras, D. I., Chen, X.-S., Ramsay, A. R. & Peeling, R. W. REASSURED
222 diagnostics to inform disease control strategies, strengthen health systems and
223 improve patient outcomes. *Nat. Microbiol.* **4**, 46–54 (2019).
- 224 2. World Health Organization. *Second WHO Model List of Essential In Vitro Diagnostics.*
225 (2019).
- 226 3. Peeling, R. W. Diagnostics in a digital age: an opportunity to strengthen health
227 systems and improve health outcomes. (2015). doi:10.1093/inthealth/ihv062
- 228 4. Ghani, A. C., Burgess, D. H., Reynolds, A. & Rousseau, C. Expanding the role of
229 diagnostic and prognostic tools for infectious diseases in resource-poor settings.
230 *Nature* **528**, S50–S52 (2015).
- 231 5. Figueroa, C. *et al.* Reliability of HIV rapid diagnostic tests for self-testing compared
232 with testing by health-care workers: a systematic review and meta-analysis. *Lancet*
233 *HIV* **5**, e277–e290 (2018).
- 234 6. Klarkowski, D. B. *et al.* The evaluation of a rapid in situ HIV confirmation test in a
235 programme with a high failure rate of the WHO HIV two-test diagnostic algorithm.
236 *PLoS One* **4**, e4351 (2009).
- 237 7. Gray, R. H. *et al.* Limitations of rapid HIV-1 tests during screening for trials in Uganda:
238 diagnostic test accuracy study. *BMJ* **335**, 188 (2007).
- 239 8. Martin, E. G., Salaru, G., Paul, S. M. & Cadoff, E. M. Use of a rapid HIV testing
240 algorithm to improve linkage to care. *J. Clin. Virol.* **52**, S11–S15 (2011).
- 241 9. Cham, F. *et al.* The World Health Organization African region external quality
242 assessment scheme for anti-HIV serology. *Afr. J. Lab. Med.* **1**, 39 (2012).
- 243 10. Galiwango, R. M. *et al.* Evaluation of current rapid HIV test algorithms in Rakai,
244 Uganda. *J. Virol. Methods* **192**, 25–7 (2013).
- 245 11. Louis, F. J. *et al.* Evaluation of an external quality assessment program for HIV testing
246 in Haiti, 2006-2011. *Am. J. Clin. Pathol.* **140**, 867–71 (2013).
- 247 12. Peck, R. B. *et al.* What Should the Ideal HIV Self-Test Look Like? A Usability Study of
248 Test Prototypes in Unsupervised HIV Self-Testing in Kenya, Malawi, and South Africa.
249 *AIDS Behav.* **18**, 422–432 (2014).
- 250 13. Baveewo, S. *et al.* Potential for false positive HIV test results with the serial rapid HIV
251 testing algorithm. *BMC Res. Notes* **5**, 154 (2012).
- 252 14. Crucitti, T., Taylor, D., Beelaert, G., Fransen, K. & Damme, L. Van. Performance of a
253 Rapid and Simple HIV Testing Algorithm in a Multicenter Phase III Microbicide Clinical
254 Trial. *Clin. Vaccine Immunol.* **18**, 1480 (2011).
- 255 15. Tegbaru, B. *et al.* Assessment of the implementation of HIV-rapid test kits at different
256 levels of health institutions in Ethiopia. *Ethiop. Med. J.* **45**, 293–9 (2007).
- 257 16. Johnson, C. C. *et al.* To err is human, to correct is public health: a systematic review
258 examining poor quality testing and misdiagnosis of HIV status. *J. Int. AIDS Soc.* **20**,
259 21755 (2017).
- 260 17. Learmonth, K. M. *et al.* Assessing proficiency of interpretation of rapid human
261 immunodeficiency virus assays in nonlaboratory settings: ensuring quality of testing.
262 *J. Clin. Microbiol.* **46**, 1692–7 (2008).
- 263 18. Garcia, P. J. *et al.* Rapid Syphilis Tests as Catalysts for Health Systems
264 Strengthening: A Case Study from Peru. *PLoS One* **8**, e66905 (2013).
- 265 19. Sacks, R., Omodele-Lucien, A., Whitbread, N., Muir, D. & Smith, A. Rapid HIV testing
266 using Determine™ HIV 1/2 antibody tests: is there a difference between the visual
267 appearance of true- and false-positive tests? *Int. J. STD AIDS* **23**, 644–646 (2012).
- 268 20. Doan, M. & Carpenter, A. E. Leveraging machine vision in cell-based diagnostics to
269 do more with less. *Nat. Mater.* **18**, 414–418 (2019).

- 270 21. Esteva, A. *et al.* Dermatologist-level classification of skin cancer with deep neural
 271 networks. *Nature* **542**, 115–118 (2017).
- 272 22. De Fauw, J. *et al.* Clinically applicable deep learning for diagnosis and referral in
 273 retinal disease. *Nat. Med.* **24**, 1342–1350 (2018).
- 274 23. Xu, Y. *et al.* Deep Learning Predicts Lung Cancer Treatment Response from Serial
 275 Medical Imaging. *Clin. Cancer Res.* **25**, 3266–3275 (2019).
- 276 24. Silver, D. *et al.* A general reinforcement learning algorithm that masters chess, shogi,
 277 and Go through self-play. *Science (80-.)*. **362**, 1140–1144 (2018).
- 278 25. Ascent of machine learning in medicine. *Nat. Mater.* **18**, 407–407 (2019).
- 279 26. Ching, T. *et al.* Opportunities and obstacles for deep learning in biology and medicine.
 280 *bioRxiv* 142760 (2018). doi:10.1101/142760
- 281 27. Zeng, N., Wang, Z., Zhang, H., Liu, W. & Alsaadi, F. E. Deep Belief Networks for
 282 Quantitative Analysis of a Gold Immunochromatographic Strip. *Cognit. Comput.* **8**,
 283 684–692 (2016).
- 284 28. Carrio, A., Sampedro, C., Sanchez-Lopez, J. L., Pimienta, M. & Campoy, P.
 285 Automated low-cost smartphone-based lateral flow saliva test reader for drugs-of-
 286 abuse detection. *Sensors (Switzerland)* **15**, 29569–29593 (2015).
- 287 29. Neuman, M. *et al.* The effectiveness and cost-effectiveness of community-based lay
 288 distribution of HIV self-tests in increasing uptake of HIV testing among adults in rural
 289 Malawi and rural and peri-urban Zambia: protocol for STAR (self-testing for Africa)
 290 cluster randomized evaluations. *BMC Public Health* **18**, 1234 (2018).
- 291 30. Aicken, C. R. H. *et al.* Young people’s perceptions of smartphone-enabled self-testing
 292 and online care for sexually transmitted infections: qualitative interview study. *BMC*
 293 *Public Health* **16**, 974 (2016).
- 294 31. Witzel, T. C., Weatherburn, P., Rodger, A. J., Bourne, A. H. & Burns, F. M. Risk,
 295 reassurance and routine: a qualitative study of narrative understandings of the
 296 potential for HIV self-testing among men who have sex with men in England. *BMC*
 297 *Public Health* **17**, 491 (2017).
- 298 32. Nsabimana, A. P. *et al.* Bringing Real-Time Geospatial Precision to HIV Surveillance
 299 Through Smartphones: Feasibility Study. *JMIR public Heal. Surveill.* **4**, e11203
 300 (2018).
- 301 33. Laksanasopin, T. *et al.* A smartphone dongle for diagnosis of infectious diseases at
 302 the point of care. *Sci. Transl. Med.* **7**, 273re1 (2015).
- 303 34. Mudanyali, O. *et al.* Integrated rapid-diagnostic-test reader platform on a cellphone.
 304 *Lab Chip* **12**, 2678 (2012).
- 305 35. Allan-Blitz, L.-T. *et al.* Field evaluation of a smartphone-based electronic reader of
 306 rapid dual HIV and syphilis point-of-care immunoassays. *Sex. Transm. Infect.* **94**,
 307 589–593 (2018).
- 308 36. S, F. *et al.* Immunochromatographic Diagnostic Test Analysis Using Google Glass.
 309 *ACS Nano* **8**, (2014).
- 310 37. Guan, Q. *et al.* Diagnose like a Radiologist: Attention Guided Convolutional Neural
 311 Network for Thorax Disease Classification. (2018).
- 312 38. Sandler, M., Howard, A., Zhu, M., Zhmoginov, A. & Chen, L.-C. MobileNetV2: Inverted
 313 Residuals and Linear Bottlenecks. (2018).
- 314 39. Chaturvedi, S. S., Gupta, K. & Prasad, P. S. Skin Lesion Analyser: An Efficient
 315 Seven-Way Multi-class Skin Cancer Classification Using MobileNet. in 165–176
 316 (Springer, Singapore, 2021). doi:10.1007/978-981-15-3383-9_15
- 317 40. Howard, A. *et al.* Searching for MobileNetV3. (2019).

318 319 320 **ACKNOWLEDGEMENTS**

321
322 We thank the community of the uMkhanyakude district and the study participants, as well as
 323 the AHRI team of fieldworkers and their supervisors. We thank A. Koza, Z. Thabethe, T.
 324 Madini, N. Okesola and S. Msane for their help with the pilot study; D. Gareta and J. Dreyer

325 for IT support; V. Lampos and I. J. Cox for useful discussions; and E. Manning and J.
326 McHugh for their help with editing and project management. This research was funded by
327 the m-Africa Medical Research Council GCRF Global Infections Foundation Award (no.
328 MR/P024378/1, to C.H., D.P., K.H., M.S., R.A.M. and V.T.) and is part of the EDCTP2
329 program supported by the European Union, i-sense Engineering and Physical Sciences
330 Research Council Interdisciplinary Research Collaboration (EPSRC IRC) in Early Warning
331 Sensing Systems for Infectious Disease (no. EP/K031953/1, to R.A.M., V.T., D.P., S.M.,
332 S.G., N.A. and M.S.), the i-sense: EPSRC IRC in Agile Early Warning Sensing Systems for
333 Infectious Diseases and Antimicrobial Resistance (no. EP/R00529X/1, to R.A.M., V.T., D.P.,
334 S.G., N.A. and S.M.) and supported by the National Institute for Health Research University
335 College London Hospitals Biomedical Research Centre (R.A.M. and S.M.). We thank the m-
336 Africa and i-sense investigators and advisory boards. The AHRI is supported by core funding
337 from the Wellcome Trust (core grant no. 082384/Z/07/Z, to T.S., D.P. and K.H.).
338

339 AUTHOR CONTRIBUTION STATEMENT

340
341 VT and RAM wrote the manuscript with input from co-authors; VT, CH, TMn, ND and TMh
342 collected the field data; VT and SM developed the machine learning models
343 with contributions from VC, KS SG and RAM; VT, NA and JB were involved in manual data
344 pre-processing; KH oversaw the data collection and data management; TS and MS provided
345 access to anonymised blood samples used in the pilot study. RAM, VT, MS, KH and DP
346 conceived the overall project, designed the study and secured the funding. RAM was the PI
347 with overall responsibility for the i-sense EPSRC IRC and the m-Africa programmes. She
348 was the supervisor of the Research Associates (VT, SM and NA) and students (VC, KS and
349 JB) involved in this study.
350

351 COMPETING INTERESTS STATEMENT

352
353 The authors declare no competing interests.
354

355 FIGURE LEGENDS

356
357 **Figure 1: Infographic to illustrate the benefits of data capture to support field**
358 **decisions.** *In blue, the current workflow used by fieldworkers. In orange, our proposed*
359 *mHealth system of automated RDT classifier plus data capture and transmission to a secure*
360 *mHealth database. In green, the benefits arising from deploying the proposed system. The*
361 *black rectangle represents a tablet or smartphone.*
362

363 **Figure 2: Standardisation of image capture, image pre-processing and training library.**
364 **a)** *Fieldworker capturing a photograph of two HIV RDTs at the time of interpretation, in the*
365 *field in rural South Africa (photo credit: Africa Health Research Institute). The two HIV RDTs*
366 *are fitted in a plastic tray designed to standardise image capture and facilitate image pre-*
367 *processing. b)* *Interpretation process, starting from the original picture of HIV RDTs used*
368 *during the study, pre-processing to select the region of interest (ROI), then interpretation of*
369 *the test result. If two lines (control + test) are present on the paper strip at the time of*
370 *interpretation, the test result is positive. Note: for the ABON HIV RDT, one or two different*
371 *test lines can appear (T1 and T2) depending on the type of HIV infection (HIV-1 and HIV-2,*
372 *respectively). The test result is positive regardless of which test line is present, or if both test*
373 *lines are present on the paper strip at the time of interpretation. If only the top line (control) is*
374 *present, the test is negative. If no control line can be seen, the test is deemed invalid. c)*
375 *Snapshot of the image library of HIV RDTs collected in the field in rural South Africa (162*
376 *randomly selected images out of 11374), illustrating the diversity of the colour, background*
377 *and brightness.*
378

379 **Figure 3: Algorithm training and performance. a)** Table showing the number of images in
380 the training library, divided in two labels categories ('positive' and 'negative') as well as two
381 sub-categories corresponding to the test type. **b)** Table to summarise the training process
382 using cross-validation, with a training set of $N = 3998$ (test type A) and $N = 6221$ (test type
383 B). The sensitivity and specificity were obtained using a hold-out testing dataset of $N = 445$
384 (test type A) and $N = 693$ (test type B). **c)** Barplots showing the average performance
385 (sensitivity and specificity) of 4 classification methods trained on our dataset, using cross
386 validation (the error bars represent the standard deviation from the mean). The three CNN
387 pretrained on the ImageNet dataset (ResNet50, MobileNetV2 and MobileNetV3) were
388 retrained and tested using our dataset. The SVM was trained using features extracted by
389 Histogram of Oriented Gradients. All four classifiers were trained using the same training set
390 described in panel b). The sensitivity and specificity were obtained using the hold-out testing
391 dataset described in panel b).
392

393 **Figure 4. Performance evaluation of our mHealth system compared to traditional**
394 **visual interpretation, field pilot study. a)** Graphics showing the agreement (%) between
395 pairs of study participants, when asked to interpret HIV RDTs results using traditional visual
396 interpretation. Participants are divided between experienced nurses (N1, N2) and community
397 health workers (C1, C2, C3). For each pair of participants, the number of HIV RDTs was $N =$
398 38 . The observations are separated according to the two types of HIV RDTs used in the
399 study. The purple square on both graphics highlights the agreement between the two
400 experienced nurses, while the orange polygon highlights the agreement between the three
401 pairs of community health workers. **b)** Confusion matrices showing the number of True
402 Negative, False Positive, False Negative and True Positive results, when comparing the
403 interpretation of our mHealth system (top row) and traditional visual interpretation (bottom
404 row) to the groundtruth. Red matrices on the left include the results for all study participants,
405 which are broken down into experienced nurses (orange matrices) and community health
406 workers (purple matrices). **c)** Barplots showing the performance index for individual
407 participants. Participants are divided between experienced nurses (N1, N2) and Community
408 health workers (C1, C2, C3). The performance index is the ratio of the performance of our
409 mHealth system over that of traditional visual interpretation. A performance index greater (or
410 equal) to one indicates our mHealth system performed better than (or as well as) traditional
411 visual interpretation. The observations are separated according to the two types of HIV
412 RDTs used in the study.
413

414 METHODS

415 Ethics

416 Ethical approval for the demographic surveillance study was granted by the Biomedical
417 Research Ethics Committee of the University of KwaZulu-Natal, South Africa, Reference
418 Number BE435/17. Separate informed consent is required for the main household survey,
419 for the HIV sero-survey, the HIV point of care test and the photographs of the HIV test.
420
421
422

423 Ethical approval for the collection of human blood samples used in the pilot study was
424 granted by the Biomedical Research Ethics Committee of the University of KwaZulu-Natal,
425 South Africa, Reference Number BFCJ 11/18.
426

427 Recruitment of participants to AHRI Population Implementation Platform for the image library

428 Eligible participants are all individuals age 15 years and older resident within the geographic
429 boundaries of the AHRI population intervention programme surveillance area (Cohort profile:
430 Africa Centre demographic information system (ACDIS) and population-based HIV survey.
431 International journal of epidemiology. 2007 Nov 12;37(5):956-62.). Individuals who have died
432 or outmigrated prior to the surveillance visit are no longer eligible. There are three contact
433

434 attempts by the fieldworker team and a further three contact attempts by a tracking team
435 before the individual is considered to be uncontactable. All individuals in the study gave
436 informed consent. Specifically, all contacted eligible individuals who gave informed consent
437 for this study were offered a rapid HIV test if they were not currently on anti-retroviral
438 therapy. For children under the age of 18, written consent for Rapid HIV testing was
439 obtained for the parent or guardian and assent from the participant.

441 *HIV RDT Image library collection*

442
443 The original RDT images library was collected in rural South Africa by a team of 60
444 fieldworkers (between 2017 and 2019). AHRI fieldworkers survey a population of 170,000
445 people in rural KwaZulu-Natal. Participants were visited at their home, those giving informed
446 consent were tested for HIV using a combination of two HIV RDTs, and upon further
447 consent, a picture of their two HIV RDTs was captured by the fieldworker on a tablet at the
448 time of interpretation. Both HIV RDTs were used as part of routine demographic surveillance
449 in Africa Health Research Institute. The test type continued to change during this study
450 following recommendations by the South African government, exemplifying the need for
451 robust systems to read multiple test formats.

452
453 While the two HIV RDTs used in this study have their own instructions for use (see
454 manufacturer's instructions), they all generally follow the same principle of collecting a drop
455 of blood from the participant's fingertip, delivering that drop of blood to the sample pad and
456 using a drop of chase buffer to help the blood sample flow through the length of the paper
457 strip. The result (a combination of one or two lines appearing on the paper strip) is then read
458 out after a period of 10 to 40 min, depending on the type of HIV RDT used.

459
460 In order to least disturb the fieldworker's workflow, a plastic tray designed to hold both HIV
461 RDT was given to each fieldworker. A picture of the tray can be seen in Figure 2a. This
462 ensured the fieldworkers only had to capture one picture per participant. The tasks of
463 separating the two HIV RDT and isolating the ROI used to train the classifier were
464 conducted down the line as part of data pre-processing.

465
466 A standard operating procedure (SOP) on how to capture the image was co-created and
467 optimised with the team of fieldworkers. A copy of the SOP can be found in the Extended
468 Data section (Extended Data Figure 1). The SOP was designed to minimise the impact of
469 environmental factors, as well as to ensure a standard way of capturing the pictures. All
470 fieldworkers attended a two-day initial training programme during which the objectives of the
471 data collection and design of the plastic tray were clearly explained, and each fieldworker
472 was personally trained and given feedback on how to capture valid photographs. A training
473 protocol was also established, in order to ensure newly enrolled fieldworkers who did not
474 attend the initial training session could also be trained to capture pictures for the project.
475 Finally, picture quality assessment sessions were conducted in order to give the fieldworkers
476 team feedback, and to ensure most pictures were of high enough quality to be used for
477 training the classifier.

478
479 All pictures were captured using Samsung tablets (SM-P585, 8MPixels camera, f1/9, with
480 autofocus capability) using the native Android camera application, stored on the device until
481 the end of the day when they were transferred to a secure database at AHRI. Our mHealth
482 system only allows one picture per test and per participant to be saved to the tablet and
483 uploaded to the AHRI database. After anonymisation (including stripping geo-coordinates
484 from the picture EXIF data), batches of 2000-3000 pictures were securely transferred to UCL
485 team members on a quarterly basis, and stored securely in a 'Data Safe Haven' managed by
486 the university.

487

488 Both the feasibility (93%) and acceptability (98%) of the system used to capture the HIV
489 RDTs pictures were high, according to a survey taken by the fieldworkers involved in the
490 study.

491
492 For the purpose of this study, an initial batch of 11, 374 images were used. As only very few
493 invalid results were obtained from the field, it was decided, for the purpose of this proof of
494 concept study, to focus on training the classifier to distinguish between positive and negative
495 results. In order to optimise this task, the ROI around each HIV RDT was isolated and used
496 to train the classifier.

497 498 *Image labelling*

499
500 All pre-processed images were labelled by a group of three RDT experts (99.2% agreement
501 with fieldworkers labelling). Labelling is the process of sorting the images into categories,
502 which are then used to train the classifier. The categories chosen here correspond to the
503 possibilities for the HIV RDT result, i.e. 'positive' and 'negative'. We recognise that a third
504 outcome, 'invalid', is also possible and needs to be considered when using the system to
505 provide a confident diagnostic. However, the absence of invalid test results in our library of
506 images collected by fieldworkers did not allow us to train the classifier on this third category
507 in this study. We therefore focused the training on the two main categories ('positive' and
508 'negative'), and are exploring other ways to incorporate the 'invalid' outcome in our mHealth
509 system. This could mean either using data augmentation techniques on the low numbers of
510 invalid test results images, or adding a pre-processing step to detect the presence of a
511 control line on the image before deciding to feed it (or not, in case the control line is absent)
512 to the classifier.

513 514 *Training library*

515
516 The labelled images were divided into two sub-categories corresponding to the HIV RDT
517 type. The two types of tests in our library are:

- 518 • **Type A:** ABON™ HIV 1/2/O Tri-Line Human Immunodeficiency Virus Rapid Test Device
519 (Whole Blood/Serum/Plasma) (ABON Biopharm (Hangzhou) Co.,Ltd)
- 520 • **Type B:** ADVANCED QUALITY™ ONE STEP Anti-HIV (1&2) Test (InTec PRODUCTS,
521 INC)

522
523 While there are two tests per patient, herein in this study we treat each test individually since
524 the tests are from different manufacturers and therefore could respond differently to the
525 same blood sample. The collection system design also guaranteed that there was never
526 more than one image of a given test per participant.

527 528 *Image normalisation*

529
530 Before being used for training, each image was resized to the dimensions of the input layer
531 then standardised. Standardisation of the data was performed using equation (1) below,
532 where x_s is the standardised pixel value, x_o the original pixel value, μ and σ are the mean
533 and standard deviation of all pixels in the image, respectively.

$$534$$
$$535$$
$$536 \quad x_s = \frac{x_o - \mu}{\sigma} \quad \text{Equation (1)}$$

537

538 *Cross-validation*

539
540 Each dataset (one for each type of HIV RDT) was randomly divided into 10 equal folds.
541 Using the leave-one-out method, 10 classifiers were trained using nine folds as the training

542 set (further randomly divided into 80% training and 20% validation). To account for
543 imbalanced datasets (roughly 13:1 negative:positive ratio), we forced every batch during
544 training to contain 50% positive images and 50% negative images using random sampling.
545 Each model was then optimised by creating a ROC curve using the validation set. This
546 yielded an optimal threshold which was used to evaluate the model performance model on
547 the testing set (remaining 10th fold). The deployment models were obtained by retraining
548 using all the available data, for each type of HIV RDT. All training and evaluation were
549 conducted using the scikit-learn and Tensorflow libraries in Python.

550

551 *Comparison with established classification methods*

552

553 The SVM was trained using pre-processed features extracted using Histogram of Oriented
554 Gradients (HOG), with Principal Component Analysis used to filter out less significant
555 features. The three CNN (ResNet50, MobileNetV2 and MobileNetV3) were pre-trained using
556 the ImageNet dataset, and re-trained using our dataset. For all four methods, training and
557 evaluation was conducted using the scikit-learn and Tensorflow libraries in Python.

558

559

560 *Android application*

561

562 We developed a smartphone/tablet Android application designed for end-users to capture a
563 picture of their HIV RDT, at the time of reading the test result. Together with end users, we
564 optimised the design so as to maximise the simplicity of the process, in order to make our
565 mHealth system accessible to end users with a broad range of digital literacy. All that is
566 required from the end user is to roughly align a semi-transparent template of the HIV RDT
567 with their HIV RDT and press a button to capture a picture. Cropping around the ROI was
568 then performed automatically in the background (using the pixel coordinates of the template
569 overlay), as was the process of sending the ROI to our classifier and receiving our mHealth
570 system's result. For the purpose of this pilot study, participants were not made aware of our
571 mHealth system's interpretation of the test results, so as to avoid bias for their own
572 interpretation. Screenshots of the application can be found in the Extended Data section
573 (Extended Data Figure 2).

574

575 *Field pilot study protocol*

576

577 The Android application was deployed in a field pilot study in KwaZulu Natal, South Africa.
578 Five participants were randomly selected from the staff at AHRI – two experienced nurses
579 and three community healthworkers. 40 HIV RDT (20 of type A, 20 of type B) were
580 performed following manufacturer's guidelines using discarded anonymised human blood
581 samples (10 positive, 10 negative according to ELISA). For each of the 40 HIV RDTs, each
582 participant was asked to record their visual interpretation of the test result, then use our
583 mHealth system on a tablet to capture a photograph of the HIV RDT. The system consisted
584 of our Android app (described above), installed on a single Samsung SM-P585 tablet,
585 identical to the ones used by fieldworkers for data collection. Participants were not shown
586 the automated interpretation of the test result provided by our mHealth system in order to
587 avoid confirmation bias. The field pilot study took place at the AHRI rural site at the heart of
588 the community (Mtubatuba, KwaZulu-Natal), under lighting conditions identical to the ones
589 the mHealth system is intended to be used. A short (10 minutes) demonstration on how to
590 use the smartphone application was given to all participants, who were then left on their own
591 to proceed with the task of reading the HIV RDTs and capturing pictures.

592

593 *Field pilot study data analysis*

594

595 The data analysis consisted of the comparison of three datasets:

596

- 597 i) Traditional visual interpretation by study participants
- 598 ii) Independent expert interpretation of the images captured by study participants
- 599 iii) Automated machine learning interpretation by our classifier

600

601 Traditional visual interpretation was recorded on the tablet by each study participant
602 immediately after being shown the HIV RDTs. Only two of the 40 HIV RDTs (corresponding
603 to 10 images out of 200) had to be discarded from the analysis, as one participant took a
604 photograph of the wrong HIV RDTs and it was therefore not possible to compare
605 interpretation results across all five participants.

606

607 An independent RDT expert subsequently visually interpreted all 190 HIV RDTs images. The
608 independent RDT expert had significant experience conducting performance evaluations of
609 lateral flow rapid tests for ocular and genital *Chlamydia trachomatis* in The Philippines, The
610 Gambia and Senegal. The visual interpretation occurred 1-5 hours after sample addition.
611 The independent expert certified that none of the HIV RDT results had changed during this
612 time frame.

613

614 The automated machine learning interpretation by our classifiers occurred on our secured
615 server. The results were compared to traditional visual interpretation and the independent
616 RDT expert, shown in the confusion matrices in Figure 4, then analysed using the
617 performance indicators described below.

618

619 *Performance indicators*

620

621 The four indicators of performance investigated were sensitivity, specificity, positive
622 predictive value (PPV) and negative predictive value (NPV). For each image, the classifier
623 produces an outcome that belongs to either of the four categories: True Positive (TP), True
624 Negative (TN), False Positive (FP) and False Negative (FN). Whether the outcome is True
625 or False depends on the comparison with the gold standard chosen.

626

627 The sensitivity is the ability of the classifier to correctly detect a positive result, by measuring
628 the ratio $\frac{TP}{TP+FN}$, while the specificity is the ratio $\frac{TN}{TN+FP}$ and translates the ability of the
629 classifier to correctly detect a negative result. The PPV is the ratio $\frac{TP}{TP+FP}$, the NPV is the
630 ratio $\frac{TN}{TN+FN}$. They indicate the proportion of positive and negative results (respectively) by a
631 diagnostic test that are true positives and true negatives (respectively).

632

633 *Data availability*

634

635 The datasets generated during and/or analysed during the current study are available from
636 the AHRI data repository:
637 Herbst, K., & McKendry, R. (2019). *m-Africa: Building mobile phone-connected diagnostics
638 and online care pathways for optimal delivery of population HIV testing, prevention and care
639 in decentralised settings* (Version 1) [Data set]. Africa Health Research Institute (AHRI).
640 <https://doi.org/10.23664/AHRI.M-AFRICA.2019.V1>

641

642 *Code availability*

643

644 Custom code used in this study is available on the public repository:
645 https://xip.uclb.com/product/classify_ai

646