

1 **Designing deep learning studies in cancer diagnostics**

2

3 Andreas Kleppe^{1,2}, Ole-Johan Skrede^{1,2}, Sepp De Raedt^{1,2}, Knut Liestøl^{1,2}, David J. Kerr³, and Håvard E.

4 Danielsen^{1,2,3†}

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

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

7 ³Nuffield Division of Clinical Laboratory Sciences, University of Oxford, Oxford, United Kingdom

8

9 [†]Corresponding author: E-mail: hdaniels@ifi.uio.no

10

11 **Abstract**

12 The number of publications on deep learning for cancer diagnostics is rapidly increasing, and systems are
13 frequently claimed to perform comparable to or better than clinicians. However, few systems have yet
14 demonstrated real-world medical utility. In this Perspective, we discuss reasons for the moderate progress, and
15 describe remedies designed to facilitate transition to the clinic. Recent, presumably influential deep learning
16 studies in cancer diagnostics, of which the vast majority used images as input to the system, are reviewed to
17 reveal the status of the field. By manipulating real data, we then exemplify that much and varied training data
18 facilitates the generalisability of neural networks, and thus the ability to use them clinically. To reduce the risk of
19 biased performance estimation of deep learning systems, we advocate evaluation in external cohorts, and
20 strongly advise that the planned analyses, including a predefined primary analysis, are described in a protocol
21 preferentially stored in an online repository. Recommended protocol items should be established for the field,
22 and we present our suggestions.

23

24 [H1] Introduction

25 Deep learning [G] facilitates utilisation of large datasets through direct learning of correlations between raw
26 input data and target output, providing systems that may use intricate structures in high-dimensional input data to
27 accurately model the association with the target output^{1,2}. A number of studies have reported on the applicability
28 of deep learning in cancer diagnostics, including prediction of diagnosis, prognosis and treatment response³⁻⁵.
29 While a large number of these tools are claimed to perform comparably or better than clinicians, few have yet
30 demonstrated real-world medical utility⁶. This is partly a natural consequence of the time needed for evaluating
31 and adapting systems affecting patient treatment. However, many studies evaluating apparently well-functioning
32 systems are at high risk of bias⁶. Of particular concern is the frequent lack of stringent evaluation on external
33 data^{7,8} and that some systems are developed or evaluated on data that are too narrow or inappropriate for the
34 intended medical setting⁹⁻¹². Thus, the lack of a well-established sequence of evaluation steps for converting
35 promising prototypes into properly evaluated medical systems clearly limits the medical utilisation of deep
36 learning systems [G].

37

38 While supervised machine learning [G] techniques traditionally utilised carefully selected representations of the
39 input data to predict the target output, modern deep learning techniques use highly flexible artificial neural
40 networks [G] to correlate input data directly to the target outputs^{1,2,13}. The relations learned by such direct
41 correlation will often be true but may sometimes be spurious phenomena exclusive to the data utilised for
42 learning. In fact, the millions of adjustable parameters make deep neural networks capable of performing
43 perfectly in training [G] sets even when the target outputs are randomly generated and therefore utterly
44 meaningless¹⁴. Thus, the high capacity [G] of neural networks induces serious challenges on how to design and
45 develop deep learning systems, and on how to validate that such a system performs adequately in the intended
46 medical setting¹⁵. Adequate clinical performance will only be possible if the system has good generalisability
47 [G] to subjects not included in the training data^{16,17}.

48

49 The design challenge involves issues related to selection of appropriate training data, such as representativeness
50 of the target population (BOX 1), as well as modelling questions such as how the variation of training data may
51 be artificially increased without jeopardising the relationship between input data and target outputs in the

52 training data^{18,19}. The validation challenge includes verifying that the system generalises well, e.g. performs
53 satisfactorily when evaluated on relevant patient populations at new locations and when input data are obtained
54 using differing laboratory procedures or alternative equipment^{15,16}. Moreover, deep learning systems are
55 typically developed iteratively, with repeated testing and often including various selection processes that may
56 bias results²⁰. Similar selection issues have been recognised as a general concern for the medical literature for
57 many years^{21,22}. Thus, when selecting design and validation processes for diagnostic deep learning systems, one
58 will have to focus both on the generalisation challenges and on preventing ‘classical’ pitfalls in data analysis. We
59 will, however, argue that both sets of challenges may be diminished by adopting certain fairly simple principles
60 partly borrowed from the drug clinical trial field.

61

62 In this Perspective, we first describe the validation challenges with focus on the use of **external cohorts [G]**. An
63 evaluation of presumably influential deep learning studies is used to reveal the status of the field particularly
64 with respect to validation procedures. We then consider generalisation issues, especially looking at the
65 importance of both natural and artificially induced variations in training datasets. In the last part, we highlight
66 the importance of evaluating an external cohort according to a predefined primary analysis to reduce selection
67 bias, and we outline a suggested sequence of evaluation steps for deep learning studies in cancer diagnostics,
68 including the use of protocols with predefined analysis plans.

69

70 **[H1] External cohort evaluation**

71 Rigorous performance evaluation is particularly important due to the inherent high complexity of deep neural
72 networks, as seemingly well-performing deep learning systems might utilise unintentional and possibly false
73 features¹⁰⁻¹² and respond unexpectedly to apparently irrelevant changes of the input data²³. Failure to properly
74 evaluate systems might have far-reaching consequences, including misdirection of further research, diminished
75 credibility of research findings and, most importantly, being worthless or even harmful to patients if used to
76 influence treatment^{24,25}.

77

78 ***[H2] The importance of an external cohort evaluation***

79 As an initial evaluation step, the cohort used for development of a deep learning system is often partitioned
80 randomly into three distinct subsets hereunder referred to as ‘training’, ‘tuning’ [G] and ‘test’ [G], where the
81 training subset is applied to learn candidate deep learning models [G], the tuning subset to select the deep
82 learning system that appears to perform best, and the test subset to evaluate the performance of the selected
83 system⁸. The evaluation on the test subset may provide unbiased estimation of the performance in the
84 development cohort [G]. It may also provide some information on the system’s ability to perform well in other
85 populations by considering the extent to which the system performs better on the training subset than on the test
86 subset, as this indicates the level of overfitting [G] to the training data. Systems that are highly overfitted to the
87 training data are likely not to perform well on other populations as the noise utilised to improve the performance
88 on the training subset may negatively influence the performance on other populations. However, even a system
89 that performs similarly in training and test subsets might perform far from acceptably on cohorts distinct from
90 the development cohort^{26,27}. As discussed below and in BOX 1, this may be caused by the system utilising data
91 features that correlate with the target outcome only in the development cohort, which could be viewed as
92 overfitting to the entire development cohort, or it might also be caused by important predictive features not being
93 adequately represented in the development cohort. Thus, using a random subset of the development cohort for
94 testing does not imply that the results have external validity, i.e. the performance of the system observed in the
95 test subset may not generalise to patients external to the development cohort.

96

97 For example, Zech, Badgeley and colleagues¹¹ investigated a deep learning system for detection of pneumonia in
98 chest X-rays, and found that it was not able to uphold the high discrimination performance achieved in the
99 development cohort when applied to cohorts from different institutions. In this case there was a substantially
100 higher disease prevalence in one of the training cohorts, and it appears that the poor generalisation was in part
101 caused by utilisation of cohort-specific characteristics. In particular, the system utilised metallic tokens that
102 radiology technicians placed on patients to indicate laterality, as these often appeared differently in different
103 cohorts. The authors further point out that the system might not even generalise well to other patients from the
104 same institution as the development cohort, because some correlations between input data and target outcome in
105 the development cohort may not be present in new cohorts from the same institution. Winkler and colleagues¹²
106 found that for their system, visible surgical skin markings present in the image were associated with higher
107 prediction score for melanoma. Similarly, Narla and colleagues¹⁰ reported that the presence of a ruler beside a
108 lesion in an image was associated with a higher malignancy score. Of course, neither skin markings nor rulers

109 are causing the skin disease, but the apparent correlation present in the development cohort is sufficient for the
110 deep learning system to make use of these associations. It could be argued that a more thorough quality control
111 on the training data could mitigate this, but it is highly unlikely that one is able to detect and control for all
112 potential confounding factors present in the training set.

113

114 Thus, unbiased performance estimation in a real-world application of a deep learning system requires external
115 cohorts representative for a target population^{22,28-30}. In an **external validation [G]**, no information from the
116 external cohort should have influenced the design of the system or the estimation of any model parameter.
117 Additionally, the external cohorts will implicitly define the patient population for which we have estimated the
118 performance of the system. Thus, to know whether or not the results may be generalised to the entire target
119 population, we need a broad validation where the cohorts may be regarded as representative of this desired target
120 population, e.g. with respect to age, sex, ethnicity, geographical differences and disease prevalence^{31,32}. Other
121 types of evaluations may also be warranted prior to introducing the system in medical practice, including so-
122 called domain validation to evaluate whether the system performs consistently across a range of laboratories and
123 technical equipment (BOX 2).

124

125 Objective, non-random separation of patients from the same hospital or subjects from the same country, e.g.
126 distinguishing between patients treated before and after a certain date, allows using one cohort for training and
127 tuning and the other for what has been denoted ‘narrow validation’ (BOX 2)²². Such evaluation might provide
128 unbiased performance estimation for a particular hospital. However, the two cohorts should not simply be a non-
129 random separation of an originally larger cohort but instead be processed separately when acquiring data and
130 ascertaining target output³³. Narrow validation is sometimes considered a limited type of external validation²².

131

132 ***[H2] Prevalence in recent studies***

133 In order to investigate the prevalence of external cohort evaluation and other characteristics of recent studies on
134 deep learning and cancer diagnostics, we searched PubMed on 21st of April 2020 for original research articles
135 published in 2015 or later (Supplementary Methods). The search provided 3,578 results, and the number of

136 publications roughly doubled each year since 2016. To explore the use of external cohort evaluation and other
137 characteristics in some of the most prominent and perhaps best studies, we restricted our evaluation to those with
138 at least 20 citations per year or published in a journal with impact factor 10 or larger. Although studies satisfying
139 either of these criteria are presumably quite influential, we acknowledge that some of the other studies might be
140 equally good. In particular, recent studies may not have had time to accrue 20 citations even if they are currently
141 of great interest, and such studies would only be included if published in a journal with impact factor 10 or
142 larger. This will exclude most studies published in new journals that are expected to receive impact factors 10 or
143 larger when this becomes available. However, we consider the selected papers to be sufficient for the purposes of
144 this discussion, as they show that some aspects of study design could be better even in some of the presumably
145 best studies. Only 257 (7%) of the 3,578 search results satisfied at least one of these selection criteria, and
146 another 43 search results were excluded because the document type in Web of Science indicated that these were
147 not original research articles. The remaining 214 studies were manually evaluated (Supplementary Table 1). We
148 further excluded 6 studies that were not original research articles and 102 studies where deep learning was not
149 used to predict or classify features relevant for cancer diagnosis, prognosis or treatment response, or such
150 potential utility of the deep learning system was not evaluated. After also excluding 14 studies without human
151 subjects or only pertaining cell biology, we ended up with 92 eligible studies³⁴⁻¹²⁵, of which 85 (92%) used
152 images as input to the deep learning system^{34-57,59-64,66,67,69-93,95-99,101-121,123,125}.

153

154 Among 516 original research articles on artificial intelligence for diagnostic analysis of medical images
155 published in 2018, Kim, Jang and colleagues⁷ found only 31 (6%) studies that evaluated an external cohort. In
156 contrast, 50 (54%) of our 92 eligible studies evaluated the performance of the deep learning system on an
157 external cohort^{37,40,48,49,51,53,55,60,62,63,65,70,73-75,78-80,82-87,90,92,93,95,96,98,100-102,104-116,120,121,123,125}. This discrepancy is most
158 likely mainly attributed to our selection of presumably influential studies, and partly attributed to the increasing
159 usage of external cohorts (FIG. 1a); 34 (72%) of the 47 eligible studies published in 2019 and 2020 evaluated an
160 external cohort compared to 9 (39%) of the 23 eligible studies published in 2018 and 7 (32%) of the 22 eligible
161 studies published before 2018.

162

163 Among studies satisfying both our selection criteria, 79% (11 of 14) evaluated an external cohort, compared to
164 68% (25 of 37) for studies that satisfied only the impact factor criterion and 34% (14 of 41) for studies that

165 satisfied only the citation frequency criterion. It thus appears that journals with high impact factor have a
166 preference for studies evaluating external cohorts. This is consistent with the call by editors of leading scientific
167 journals for rigorous evaluation of artificial intelligence tools^{126,127} and explicit prioritisation of biomarker
168 studies that evaluate external cohorts by some journals, e.g. the Journal of Clinical Oncology
169 (<https://ascopubs.org/jco/authors/journal-policies>).

170

171 **[H1] Generalisability**

172 While increased use of external cohorts is an important step towards proper validation of deep learning systems,
173 one is still left with the challenge of ensuring that the results obtained on such a population provides a
174 satisfactory measure of the performance within the entire intended target population. This target population may
175 typically be patients who have a specific cancer type, and although often restricted e.g. to certain stages of the
176 disease, the target population is normally broad. Although some studies may use more than one external cohort
177 and some use trials with many centres distributed over several countries, it is difficult to obtain external cohorts
178 that entirely cover the target population. Thus, successful application of a deep learning system will depend on
179 good generalisation properties, so that good performance on one population also indicate satisfactory
180 performance on populations differing with respect to some properties. Fortunately, exploring generalisation in
181 deep learning is an active research area¹²⁸, and by utilising certain design principles, deep learning systems have
182 shown remarkably good generalisation performance on a number of tasks²⁻⁵.

183

184 One way of increasing generalisation is to control the neural network's capacity to express complex mappings,
185 e.g. by limiting the number of adjustable parameters in the network, imposing various constraints on the network
186 or regularising the optimisation^{129,130}. Transfer learning could also increase generalisation, particularly when
187 training data for the task at hand is scarce^{131,132}. In transfer learning, the network is initialised with parameters
188 optimised using data for a different task, typically using large datasets such as ImageNet^{133,134}, which may
189 mitigate overfitting at the possible cost of introducing biases¹³⁵⁻¹³⁷. Making the training dataset more diverse and
190 more representative of the target population is another way of increasing generalisation¹³⁸. Of particular
191 importance is to ensure adequate and unbiased representation across demographic characteristics such as sex,
192 race and ethnicity (BOX 1). In addition to expanding the natural training dataset, i.e. the set of training data

193 acquired from a range of patient samples with associated target outcome, one may artificially augment the
194 training dataset by applying smaller transformations on the inputs while maintaining their relationship to the
195 target output^{18,139}. This can reduce the network's ability to memorise details of the training data and thereby
196 increase generalisation, especially in situations where the availability of training data is limited. The transforms
197 can randomly change, often called 'distort', the input data by e.g. adding noise, erasing parts, shifting and scaling
198 colours or altering the image geometry¹⁹. Artificially diversifying the training data may increase generalisation
199 by enabling the resulting system to ignore vagaries of the measurement process and even become applicable to
200 multiple data acquisition procedures, e.g. different acquisition equipment^{140,141}. Other augmentation techniques
201 include those that generate artificial input data, e.g. by mixing multiple data inputs¹⁹. The value of augmentation
202 techniques has been observed in various application domains¹⁹, including the use on images obtained in
203 radiology^{38,142-144} and histopathology^{141,145}.

204

205 To illustrate the importance of the amount and variation in training data, and more specifically show how data
206 distortion may work to improve deep learning systems in cancer diagnostics, we show this type of analysis here
207 using data from a previously published study¹¹³. This previous study applied deep learning to predict colorectal
208 cancer-specific survival directly from conventional haematoxylin and eosin stained sections, with training and
209 tuning data derived from 2,473 patients from four cohorts. The performance was evaluated on an external cohort
210 consisting of 1,122 patients from a randomised controlled trial on a drug that was observed to not affect
211 survival¹⁴⁶. We applied the convolutional neural network called Inception-v3¹⁴⁷, which is a commonly used
212 network in medical image diagnostics⁸, in both the previously published analyses and the new analyses presented
213 here.

214

215 Initially, we applied the same distortion process as in our published analyses¹¹³. This process artificially
216 increased the variation of the training images by randomly distorting their colours, which is an augmentation
217 technique that appears crucial when training deep learning systems in histopathology¹⁴⁵. Initially, the maximum
218 amount of distortion we allowed was quite modest (FIG. 2a). To illustrate the effect of reducing the number of
219 patients while keeping the patient heterogeneity implied by having data from four cohorts, we randomly sampled
220 979 patients in such a manner that the data had the same number of training and tuning patients with and without
221 cancer-specific death as in the cohort from the Gloucester Colorectal Cancer Study, UK (the largest of the four

222 training and tuning cohorts). The decreased performance of the resulting deep learning system when evaluated
223 on the external cohort (FIG. 2b) exemplifies the importance of a large natural training dataset and its intrinsic
224 variation¹³⁸. Further reduction of the number of patients decreased the performance further; training and tuning
225 on a quarter of the 979 patients or less (that is, less than 250 patients) provided systems that did not perform
226 substantially better than random guessing (FIG. 2b).

227

228 We then showed that modifying the distortion process may mitigate for the performance loss observed when
229 reducing the number of patients in training and tuning. Compared to using all 2,473 patients for training and
230 tuning, using 979 randomly selected patients and four times the original amount of colour distortion provided
231 similar performance on the external cohort (FIG. 2c). For this modified distortion process we allowed quite
232 substantial colour distortions (FIG. 2d), and the results showed that artificial augmentation may in some cases
233 compensate for limited natural training and tuning data. However, increasing the amount of colour distortion
234 further provided worse performance (FIG. 2c), illustrating the trade-off between preventing overfitting through
235 random distortions and occluding relevant information for the prediction task.

236

237 Randomly sampling 979 patients from all four cohorts maintained much of the variation in the natural training
238 and tuning data. If we instead used only the Gloucester cohort, which contained the same number of training and
239 tuning patients with and without cancer-specific death as in the random sample, we obtained worse performance
240 on the external cohort, most clearly when including more colour distortion in training (FIG. 2e). This underlines
241 the importance of designing studies such that the natural training data is diverse, and FIG. 2e additionally
242 illustrates that natural and artificial variation works well together to increase generalisability.

243

244 In general, the most suitable distortion process will depend on the particular medical prediction task because the
245 involved data will tolerate different amounts of the various types of distortions before true correlations between
246 input and target output are occluded. For instance, deep learning systems that classify based on images of skin
247 lesions or tumour sections are likely to benefit from being invariant to rotations, while systems aimed at
248 supporting radiology might rely on the orientation in images of larger organ structures and thereby perform
249 worse if forced to be rotation invariant. Thus, the distortion process needs to be fine-tuned to the particular

250 application, as findings about which distortion process appears most beneficial in one scenario, e.g. findings
251 from the example presented in FIG. 2, are not necessarily directly applicable to other scenarios. However, the
252 general principle is that including much and varied training data is important. As the importance of artificial
253 augmentation decreases with the amount and diversity in the natural training data, prediction tasks where the true
254 correlations between input data and target output are easily obscured by distortion warrants a more
255 comprehensive natural training dataset.

256

257 **[H1] Predefined primary analysis**

258 In the development of a deep learning system, researchers will often evaluate different systems sequentially,
259 each time having the possibility to learn from interpreting the previous evaluations and adapt the system to the
260 specific data used for evaluation. Such repeated evaluations will bias the estimates, and their dependence on
261 previous evaluations makes established statistical approaches for adjusting for multiple comparisons not
262 applicable^{148,149}. Similar re-analysis issues may arise if the initial analysis of a specific deep learning system
263 reveals issues that are then corrected and the performance is re-evaluated. Such problems of repeated or multiple
264 evaluations are well-known from examinations of the data analysis in various types of published medical studies,
265 and have been identified as important contributors to biased inference and irreproducible results^{20,150}.

266

267 As discussed above, evaluation on an external cohort is required for unbiased performance estimation in a real-
268 world application of the deep learning system, but it is only a prerequisite as multiple or repeated evaluations
269 may cause bias even if evaluating an external cohort. Great caution would therefore be needed when interpreting
270 studies that report multiple analyses without specifying which was initially planned to be the primary analysis, if
271 any.

272

273 ***[H2] Prevalence of predefined primary analysis***

274 In our evaluation of recent, presumably influential deep learning studies in cancer diagnostics, all studies
275 performed multiple analyses of the external cohort in the form of either evaluating multiple systems, analysing
276 multiple subpopulations or using various analysis methods. Only 3 (6%) of the 50 eligible studies that evaluated

277 an external cohort used one of the well-established methods for adjustment for multiple comparisons^{51,62,114}, e.g.
278 Bonferroni correction. This implies that most studies should have specified which analysis was considered the
279 primary analysis prior to evaluation of the external cohort, if such a decision was made, in order to inform the
280 reader which analysis was not affected by selection bias and to help distinguish studies with a predefined
281 primary analysis from those that repeatedly evaluated the external cohort and might have ended up reporting
282 severely biased performance estimates. Although the principle of using an external dataset only once to evaluate
283 the final hypothesis should be well-known in the machine learning community^{151,152}, it seems currently that there
284 is no tradition for specifying the predefined primary analysis in deep learning publications other than those
285 reporting on clinical trials. In our evaluation, 20 (40%) of the 50 studies evaluating an external cohort specified
286 one or more primary performance metrics (FIG. 1b)^{55,60,73,82,83,85,86,93,98,102,105,108-110,113,115,116,120,121,125}, but only 8
287 (16%) of the 50 studies specified a predefined primary analysis (FIG. 1c)^{73,83,102,105,109,113,120,121}.

288

289 Prespecification of the primary analysis has previously been advocated in diagnostic and prognostic
290 research^{153,154}, but this is unfortunately still not common practise despite being the only direct protection against
291 selection bias²⁰. To ensure unbiased estimation, the primary analysis should be unequivocally specified prior to
292 all investigations that could reveal correlations between input data and target output in the external cohort. This
293 would require the researchers to define all relevant aspects of the validation prior to analysing the cohort,
294 including the deep learning system, target output, and patient and input data in the external cohort. Predefining
295 the primary analysis will entail a commitment to the main analysis, which implies that the analysis should be
296 carefully planned in advance and that researchers will be discouraged from performing creative data dredging¹⁵⁵.

297

298 *[H2] Choosing the primary metric*

299 Many medical questions are categorical in nature, e.g. whether tumour or not, whether mutated or not, and
300 whether to offer treatment or not. However, deep learning models often output continuous values reflecting the
301 predicted probability of each possible outcome. In such cases, the predefined primary analysis should preferably
302 evaluate a categorisation of the model output aimed at answering the medical question. The primary analysis will
303 then be comparing predicted and target outcome in the external cohort, e.g. by measuring the so-called **balanced**
304 **accuracy** [G]¹⁵⁶. Measuring the performance using categorical outputs often provides more conservative

305 estimates¹⁵⁷ and avoids issues with metrics frequently applied to measure the performance using continuous
306 outputs. For instance, the **area under the receiver operating characteristic curve [G]** (AUC)¹⁵⁸ and **concordance**
307 **index [G]** (c-index)¹⁵⁹ are only affected by the ranking of the continuous outputs, not the prediction scores
308 themselves¹⁶⁰. Thus, such metrics may indicate that a deep learning system performs well even if it predicts
309 markedly too high probabilities for all patients in a specific cohort, provided that the continuous outputs of the
310 system rank the patients in a fairly correct order. In another cohort, the same system may similarly appear to
311 perform well even if it predicts markedly too low probabilities for all those patients. The generalisability of such
312 a system is poor, yet this would not be evident from the AUC and c-index of the continuous outputs, but it would
313 be evident from the AUC and c-index of a categorisation defined irrespective of the external cohorts. The
314 categorisation may be defined by e.g. determining suitable thresholds during tuning or selecting the outcome
315 with highest prediction score as the predicted outcome. Defining the categorisation using the external cohort,
316 even at predefined levels of e.g. sensitivity, adapts the categorical marker to the specific external cohort and may
317 occlude shifts in the prediction scores as with the AUC and c-index of the continuous outputs.

318

319 In our evaluation of recent, presumably influential deep learning studies in cancer diagnostics, we found that 34
320 (68%) of the 50 studies evaluating an external cohort reported the estimated performance of a categorical marker
321 on the external cohort, with a categorisation defined irrespective of the external cohort<sup>48,49,53,55,60,62,63,65,73,75,78-
322 80,82,85,87,90,98,100,102,104-106,108-111,113-116,120,121,125</sup>. The proportion was lower for studies reporting on deep learning
323 systems that used histopathology section images as input, with only 6 (40%) of 15 studies evaluating a fixed
324 categorical marker on the external cohort^{48,55,82,111,113,114}, which is surprising since most histopathological
325 evaluations provide categorical values.

326

327 For certain deep learning systems, the intended medical application directly utilises the system's continuous
328 output, e.g. to triage patients for further examinations, and in such cases the continuous output should be
329 evaluated in the primary analysis. This may warrant additional analyses to reveal generalisation issues that might
330 be occluded by the selected performance metric, e.g. to consider a calibration plot in addition to the c-index
331 when evaluating a clinical decision support system for predicting patient outcome^{22,26}.

332

333 **[H1] From conception to application**

334 All research with the potential to influence patient treatment should undergo careful evaluation sequences and be
335 driven by protocols with a predefined statistical analysis plan¹⁵³. FIG. 3 illustrates what we consider as natural
336 and important steps in the development and evaluation of deep learning systems for medical applications.

337

338 The initial exploratory studies aim to answer whether deep learning appears suitable for the task at hand or
339 whether further investigations based on deep learning are not warranted at this time, usually because the
340 hypothesis seems ill-founded or the available data is not expected to provide a system with adequate
341 performance. The performance estimates obtained in such pilot studies are frequently inflated by the use of a
342 limited development cohort, but promising findings may motivate further investigations. After a series of
343 explorations and possibly expansions of the development cohort, the development should conclude by deciding
344 which system appears to perform best on the intended medical task, considering also the sensitivity to vagaries
345 of the measurement process. Of particular importance to prevent selection of a system that performs much worse
346 on patients outside the development cohort, the study could include sufficient amount and variation in the natural
347 training dataset and use techniques like data distortion to increase the variation artificially.

348

349 There is a growing interest in explainable deep learning systems¹⁶¹⁻¹⁶³, including the creation of inherently more
350 explainable systems and post-hoc explanations of existing systems¹⁶⁴. For image classification tasks in particular,
351 so-called saliency maps visualise the contribution of each pixel to the final prediction score and can be created
352 using a number of different techniques¹⁶⁵⁻¹⁶⁷. By increasing the transparency, the more explained systems might
353 have more predictable generalising abilities. This may be used to identify target populations within which the
354 system is expected to generalise well or settings where the system is prone to fail. For example, Winkler and
355 colleagues¹² used such a technique to support their finding that surgical skin markings unduly increased the
356 system's prediction score for melanoma. While current explainability techniques might suggest generalisability
357 and thereby suggest suitable target populations or influence the selection of which system to evaluate further,
358 they will only provide indications and thus not reduce the need for proper validation.

359

360 While efficacy studies of pharmaceutical products are usually preceded by prospective trials to estimate basic
361 features such as safety and dosing¹⁶⁸, deep learning systems for diagnostic purposes can to a larger extent utilise
362 retrospective cohorts, e.g. from earlier clinical trials or medical practice. Given the risks, timeframe and costs of
363 interventional research¹⁶⁸⁻¹⁷⁰, we recommend rigorous, retrospective analyses to evaluate the medical validity of
364 a deep learning system by conducting an external validation according to a predefined primary analysis. The
365 results of such studies provide valuable information to direct further research, thus warranting publication
366 regardless of the significance of the findings, which would also mitigate publication bias.

367

368 Rigorous, retrospective analyses of a deep learning system might warrant conducting a prospective, randomised
369 phase III clinical trial where the system directly intervenes with the current standard of care in order to evaluate
370 the system's medical utility in a specific real-world application, considering both benefits and harms for patients
371 in the target population^{30,171}. Systems demonstrated to have medical utility and approved by necessary
372 governmental agencies can be applied in medical practice while monitoring the long-term benefits, harms and
373 costs for each specific real-world medical application in phase IV clinical trials. Such surveillances might
374 eventually indicate that the system needs to be updated because of changes in medical practice or data
375 acquisition¹⁷².

376

377 The levels of deep learning studies depicted in FIG. 3 and the phases of clinical trials were used to categorise
378 recent, presumably influential deep learning studies in cancer diagnostics in relation to the reliability of the
379 performance estimation approach and the demonstrated applicability of the system in medical practice. Although
380 some group sizes are very small, there appears to be notable differences between research fields defined by the
381 input to the deep learning system (FIG. 4). The proportion of studies evaluating an external cohort was lowest
382 for the 7 studies with only non-image inputs such as omics data (29%; 2 of 7 studies), while highest for 22
383 studies with images other than histopathology section and radiology images as input, e.g. from gastrointestinal
384 endoscopic examinations or dermoscopic images (64%; 14 of 22 studies). Five (23%) of the 22 studies with
385 other images as input even had a predefined primary analysis of the external cohort^{73,102,105,109,121}, which included
386 the 3 studies reporting on a randomised clinical trial, all of which evaluated a deep learning system to aid
387 gastrointestinal examinations^{102,105,121}.

388

389 ***[H2] Recommended protocol items***

390 When planning to evaluate the medical validity of a deep learning system through rigorous, retrospective
391 analyses, we recommend the unequivocal specification of the predefined primary analysis to be documented in a
392 study protocol. Relevant items in such protocols would differ from clinical trial protocols, which are the target of
393 guidelines such as SPIRIT (Standard Protocol Items: Recommendations for Interventional Trials)¹⁷³ and its
394 extension to artificial intelligence¹⁷⁴. Protocols should be developed before conducting the validation, and
395 relevant items would therefore also differ from those in original research articles, which are the target of many
396 reporting guidelines such as CONSORT (Consolidated Standards of Reporting Trials)¹⁷⁵ and TRIPOD
397 (Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis)²² as well as
398 their extension or anticipated adaption to machine learning^{176,177}. It is therefore a need to establish guidelines
399 dedicated to study protocols describing validations of deep learning systems. We propose a non-exhaustive list
400 of items that we consider essential in such protocols, termed Protocol Items for External Cohort Evaluation of a
401 deep learning System (PIECES) in cancer diagnostics.

402

403 In order to be sufficiently concrete about the predefined primary analysis, the protocol needs to describe the deep
404 learning system and how it will be assayed, define the external cohort, including its origin, what it represents in
405 terms of medical setting and target population, input data and target output, and clearly specify the performance
406 evaluation. These three parts of the protocol form the basis of our PIECES recommendations together with a
407 declaration of status (BOX 3). The status declaration should scrupulously elucidate any investigations performed
408 before finalising the protocol that could reveal correlations between input data and target output in the external
409 cohort, or state that no such investigations were performed.

410

411 The PIECES recommendations are designed to facilitate identification of ambiguities and disagreements
412 between the researchers planning to conduct an external validation as well as to provide a clear description of the
413 predefined primary analysis as reference for all readers, which may aid medical professionals in identifying well-
414 designed studies and their applicability to their own clinical practice. The thought and work that should go into
415 making such a protocol could also allow the researchers to make appropriate changes prior to performing the

416 external validation. For instance, considering what the external cohort is intended to represent and how the deep
417 learning system is envisioned to be applied in practice, could affect the inclusion and exclusion criteria for
418 patients and samples as well as the metric or statistical test applied in the primary analysis.

419

420 Researchers conducting an external validation would often like to perform multiple, related analyses to elucidate
421 the performance of the deep learning system. To separate preplanned analyses from exploratory, post hoc
422 analyses, the PIECES recommendation encourages specification of predefined secondary analyses that the
423 researchers would like to commit themselves to report on publication of their findings. Such secondary analyses
424 would be affected by the multiple comparisons problem but predefining and reporting all secondary analyses
425 would provide a transparency that would substantially increase the credibility of the results. Importantly, the
426 specification of predefined secondary analyses does not diminish the validity of the predefined primary analysis.
427 Any analyses the researchers consider reporting, but do not wish to commit themselves to report, should not be
428 specified as secondary analyses in the protocol and therefore should be reported as exploratory analyses, even
429 though they might be thought of prior to analysing the external cohort.

430

431 *[H2] Study registration*

432 We recommend registration of the study protocol in an online repository before analysing the external cohort.
433 Most major trial registries, e.g. ClinicalTrials.gov (<https://www.clinicaltrials.gov>) and the International Standard
434 Randomised Controlled Trial Number (ISRCTN) registry (<https://www.isrctn.com>), accept registration of
435 diagnostic accuracy studies¹⁵⁴. These registries can be used to record external validation studies in deep learning,
436 but some items will not be relevant, while some important items such as defining the deep learning system will
437 not be encouraged. A dedicated repository to register the study protocol describing the external validation of a
438 deep learning system is therefore warranted. We recognise that it may be undesirable to publish a detailed study
439 protocol in an online repository prior to conclusion of the study as it would reveal novel work prior to
440 publication of the results and perhaps in some rare cases jeopardise publication. In a dedicated repository, a
441 submission could be partially or completely invisible to the public and the protocol encrypted until the authors
442 choose to reveal the submission and provide the required decryption key, thus facilitating preregistration of study
443 protocols without requiring authors to reveal novel ideas prematurely.

444

445 Registration of observational studies has been advocated by editors of major clinical journals^{178,179}, many
446 editorial board members¹⁸⁰ and researchers^{181,182}, and the criticism it has received from epidemiologists in
447 relation to the exploratory nature of epidemiology¹⁸³⁻¹⁸⁵ does not apply to external validation studies. For
448 diagnostic and prognostic biomarker studies in particular, the registration of a study protocol with a predefined
449 analysis plan has been recommended by several researchers^{153,154,186-188}, provided that it precedes the onset of the
450 study¹⁸⁹. This would facilitate a more balanced evaluation of the proposed marker, identification and prevention
451 of selective reporting, increased transparency, reduced proportion of false positive findings, mitigation of
452 publication bias through identification of unpublished studies, and prevention of unnecessary duplication of
453 research while facilitating collaboration between researchers and identification of research gaps. Consequently,
454 widespread preregistration of detailed study protocols for deep learning systems might translate into more rapid
455 identification of promising systems and thereby expedite progression of the research field. It would also
456 communicate a study to peers without disclosing the findings and interpretations prior to editorial and peer
457 review, thus providing some of the benefits of preprint archiving while allowing critical appraisal of the findings
458 and interpretations before publication.

459

460 Amendments of clinical trial protocols are common but should be tracked and dated¹⁷³. While clinical trials often
461 take years to conduct due to patient recruitment and follow-up, most external validations of deep learning
462 systems use retrospective data and the analysis part of the validation may be performed in a matter of days.
463 Consequently, it should rarely be necessary to modify the study protocol describing the external validation of a
464 deep learning system after initiating the validation. We therefore generally discourage protocol amendments, but
465 if found necessary for a particular study, we recommend amendments to be included as postscripts to the study
466 protocol, leaving the original protocol unaltered. Both the postscript and disseminations of the validation results
467 should concretely specify what was changed as well as describe the motivation and rationale for the change.

468

469 **[H1] Conclusions**

470 Including much natural and artificial data variation when training rigorous deep learning systems appears
471 pivotal, as analyses indicate its instrumental role in increasing the performance and generalisability of systems.

472 Utilising multiple sets of patients, samples and data acquisition procedures will diversify the training data, while
473 augmentation techniques artificially enhance the variation further. The resulting systems may be capable of
474 handling the diversity in routine medical practice and in some cases even generalise to completely new settings.

475

476 Going forward, the medical validity of a deep learning system should be evaluated according to a preregistered
477 study protocol specifying the primary analysis and using an external cohort representative of the intended
478 medical setting and target population. This facilitates balanced performance evaluations by reducing selection
479 bias and increasing transparency, and helps medical professionals distinguish rigorous, retrospective validation
480 studies from studies that repeatedly evaluated the external cohort and might end up reporting severely biased
481 performance estimates. It would therefore assist in identifying deep learning systems that warrant prospective
482 evaluations in randomised clinical trials and ultimately drive the development of systems that could transform
483 current medical practice.

484

485 **References**

- 486 1 Schmidhuber, J. Deep learning in neural networks: An overview. *Neural Netw.* **61**, 85-117 (2015).
- 487 2 LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436-444 (2015).
- 488 3 Hosny, A., Parmar, C., Quackenbush, J., Schwartz, L. H. & Aerts, H. J. W. L. Artificial intelligence in
489 radiology. *Nat. Rev. Cancer* **18**, 500-510 (2018).
- 490 4 Vamathevan, J. *et al.* Applications of machine learning in drug discovery and development. *Nat. Rev.*
491 *Drug Discov.* **18**, 463-477 (2019).
- 492 5 Bera, K., Schalper, K. A., Rimm, D. L., Velcheti, V. & Madabhushi, A. Artificial intelligence in digital
493 pathology — new tools for diagnosis and precision oncology. *Nat. Rev. Clin. Oncol.* **16**, 703-715
494 (2019).
- 495 6 Nagendran, M. *et al.* Artificial intelligence versus clinicians: systematic review of design, reporting
496 standards, and claims of deep learning studies. *BMJ* **368**, m689 (2020).
- 497 7 Kim, D. W., Jang, H. Y., Kim, K. W., Shin, Y. & Park, S. H. Design Characteristics of Studies
498 Reporting the Performance of Artificial Intelligence Algorithms for Diagnostic Analysis of Medical
499 Images: Results from Recently Published Papers. *Korean J. Radiol.* **20**, 405-410 (2019).

500 8 Liu, X. *et al.* A comparison of deep learning performance against health-care professionals in detecting
501 diseases from medical imaging: a systematic review and meta-analysis. *Lancet Digit. Health* **1**, e271-
502 e297 (2019).

503 9 Ross, C. & Sweltitz, I. *IBM's Watson supercomputer recommended 'unsafe and incorrect' cancer*
504 *treatments, internal documents show.* STAT. [https://www.statnews.com/2018/07/25/ibm-watson-](https://www.statnews.com/2018/07/25/ibm-watson-recommended-unsafe-incorrect-treatments/)
505 [recommended-unsafe-incorrect-treatments/](https://www.statnews.com/2018/07/25/ibm-watson-recommended-unsafe-incorrect-treatments/) (2018).

506 10 Narla, A., Kuprel, B., Sarin, K., Novoa, R. & Ko, J. Automated Classification of Skin Lesions: From
507 Pixels to Practice. *J. Invest. Dermatol.* **138**, 2108-2110 (2018).

508 11 Zech, J. R. *et al.* Variable generalization performance of a deep learning model to detect pneumonia in
509 chest radiographs: A cross-sectional study. *PLoS Med.* **15**, e1002683 (2018).

510 12 Winkler, J. K. *et al.* Association Between Surgical Skin Markings in Dermoscopic Images and
511 Diagnostic Performance of a Deep Learning Convolutional Neural Network for Melanoma Recognition.
512 *JAMA Dermatol.* **155**, 1135-1141 (2019).

513 13 Rueckert, D. & Schnabel, J. A. Model-Based and Data-Driven Strategies in Medical Image Computing.
514 *Proc. IEEE* **108**, 110-124 (2020).

515 14 Zhang, C., Bengio, S., Hardt, M., Recht, B. & Vinyals, O. Understanding deep learning requires
516 rethinking generalization. *Proc. Int. Conf. Learn. Represent.* (2017).

517 15 Liu, Y., Chen, P.-H. C., Krause, J. & Peng, L. How to Read Articles That Use Machine Learning:
518 Users' Guides to the Medical Literature. *JAMA* **322**, 1806-1816 (2019).

519 16 Ransohoff, D. F. Bias as a threat to the validity of cancer molecular-marker research. *Nat. Rev. Cancer*
520 **5**, 142-149 (2005).

521 17 Moons, K. G. M. *et al.* PROBAST: A Tool to Assess Risk of Bias and Applicability of Prediction
522 Model Studies: Explanation and Elaboration. *Ann. Intern. Med.* **170**, W1-W33 (2019).

523 18 Simard, P., Victorri, B., LeCun, Y. & Denker, J. Tangent Prop - A formalism for specifying selected
524 invariances in an adaptive network. *Adv. Neural Inf. Process. Syst.* **4**, 895-903 (1992).

525 19 Shorten, C. & Khoshgoftaar, T. M. A survey on Image Data Augmentation for Deep Learning. *J. Big*
526 *Data* **6**, 60 (2019).

527 20 Ioannidis, J. P. A. What Have We (Not) Learnt from Millions of Scientific Papers with P Values? *Am.*
528 *Stat.* **73**, 20-25 (2019).

529 21 Ioannidis, J. P. A. Why Most Published Research Findings Are False. *PLoS Med.* **2**, e124 (2005).

530 22 Moons, K. G. M. *et al.* Transparent Reporting of a multivariable prediction model for Individual
531 Prognosis Or Diagnosis (TRIPOD): Explanation and Elaboration. *Ann. Intern. Med.* **162**, W1-W73
532 (2015).

533 23 Heaven, D. Why deep-learning AIs are so easy to fool. *Nature* **574**, 163-166 (2019).

534 24 Ioannidis, J. P. A. Evolution and translation of research findings: from bench to where? *PLoS Clin.*
535 *Trials* **1**, e36-e36 (2006).

536 25 Topol, E. J. High-performance medicine: the convergence of human and artificial intelligence. *Nat.*
537 *Med.* **25**, 44-56 (2019).

538 26 Justice, A. C., Covinsky, K. E. & Berlin, J. A. Assessing the Generalizability of Prognostic Information.
539 *Ann. Intern. Med.* **130**, 515-524 (1999).

540 27 Subbaswamy, A. & Saria, S. From development to deployment: dataset shift, causality, and shift-stable
541 models in health AI. *Biostatistics* **21**, 345-352 (2020).

542 28 Ioannidis, J. P. A. & Khoury, M. J. Improving Validation Practices in “Omics” Research. *Science* **334**,
543 1230-1232 (2011).

544 29 Obermeyer, Z. & Emanuel, E. J. Predicting the Future — Big Data, Machine Learning, and Clinical
545 Medicine. *N. Engl. J. Med.* **375**, 1216-1219 (2016).

546 30 Keane, P. A. & Topol, E. J. With an eye to AI and autonomous diagnosis. *npj Digit. Med.* **1**, 40 (2018).

547 31 Gianfrancesco, M. A., Tamang, S., Yazdany, J. & Schmajuk, G. Potential Biases in Machine Learning
548 Algorithms Using Electronic Health Record Data. *JAMA Intern. Med.* **178**, 1544-1547 (2018).

549 32 Noor, P. Can we trust AI not to further embed racial bias and prejudice? *BMJ* **368**, m363 (2020).

550 33 Luo, W. *et al.* Guidelines for Developing and Reporting Machine Learning Predictive Models in
551 Biomedical Research: A Multidisciplinary View. *J. Med. Internet Res.* **18**, e323 (2016).

552 34 Hua, K. L., Hsu, C. H., Hidayati, S. C., Cheng, W. H. & Chen, Y. J. Computer-aided classification of
553 lung nodules on computed tomography images via deep learning technique. *Onco Targets Ther.* **8**,
554 2015-2022 (2015).

555 35 Ciompi, F. *et al.* Automatic classification of pulmonary peri-fissural nodules in computed tomography
556 using an ensemble of 2D views and a convolutional neural network out-of-the-box. *Med. Image Anal.*
557 **26**, 195-202 (2015).

558 36 Arevalo, J., González, F. A., Ramos-Pollán, R., Oliveira, J. L. & Guevara Lopez, M. A. Representation
559 learning for mammography mass lesion classification with convolutional neural networks. *Comput.*
560 *Methods Programs Biomed.* **127**, 248-257 (2016).

561 37 Setio, A. A. A. *et al.* Pulmonary Nodule Detection in CT Images: False Positive Reduction Using
562 Multi-View Convolutional Networks. *IEEE Trans. Med. Imaging* **35**, 1160-1169 (2016).

563 38 Roth, H. R. *et al.* Improving Computer-Aided Detection Using Convolutional Neural Networks and
564 Random View Aggregation. *IEEE Trans. Med. Imaging* **35**, 1170-1181 (2016).

565 39 Kallenberg, M. *et al.* Unsupervised Deep Learning Applied to Breast Density Segmentation and
566 Mammographic Risk Scoring. *IEEE Trans. Med. Imaging* **35**, 1322-1331 (2016).

567 40 Litjens, G. *et al.* Deep learning as a tool for increased accuracy and efficiency of histopathological
568 diagnosis. *Sci. Rep.* **6**, 26286 (2016).

569 41 Huynh, B. Q., Li, H. & Giger, M. L. Digital mammographic tumor classification using transfer learning
570 from deep convolutional neural networks. *J. Med. Imaging* **3**, 034501 (2016).

571 42 Nie, K. *et al.* Rectal Cancer: Assessment of Neoadjuvant Chemoradiation Outcome based on Radiomics
572 of Multiparametric MRI. *Clin. Cancer Res.* **22**, 5256-5264 (2016).

573 43 Kooi, T. *et al.* Large scale deep learning for computer aided detection of mammographic lesions. *Med.*
574 *Image Anal.* **35**, 303-312 (2017).

575 44 Esteva, A. *et al.* Dermatologist-level classification of skin cancer with deep neural networks. *Nature*
576 **542**, 115-118 (2017).

577 45 Dhungel, N., Carneiro, G. & Bradley, A. P. A deep learning approach for the analysis of masses in
578 mammograms with minimal user intervention. *Med. Image Anal.* **37**, 114-128 (2017).

579 46 Yu, L., Chen, H., Dou, Q., Qin, J. & Heng, P. Automated Melanoma Recognition in Dermoscopy
580 Images via Very Deep Residual Networks. *IEEE Trans. Med. Imaging* **36**, 994-1004 (2017).

581 47 Sun, W., Tseng, T. B., Zhang, J. & Qian, W. Enhancing deep convolutional neural network scheme for
582 breast cancer diagnosis with unlabeled data. *Comput. Med. Imaging Graph.* **57**, 4-9 (2017).

583 48 Cruz-Roa, A. *et al.* Accurate and reproducible invasive breast cancer detection in whole-slide images: A
584 Deep Learning approach for quantifying tumor extent. *Sci. Rep.* **7**, 46450 (2017).

585 49 Ciompi, F. *et al.* Towards automatic pulmonary nodule management in lung cancer screening with deep
586 learning. *Sci. Rep.* **7**, 46479 (2017).

587 50 Araújo, T. *et al.* Classification of breast cancer histology images using Convolutional Neural Networks.
588 *PLoS One* **12**, e0177544 (2017).

589 51 Becker, A. S. *et al.* Deep Learning in Mammography: Diagnostic Accuracy of a Multipurpose Image
590 Analysis Software in the Detection of Breast Cancer. *Invest. Radiol.* **52**, 434-440 (2017).

591 52 Dou, Q., Chen, H., Yu, L., Qin, J. & Heng, P. Multilevel Contextual 3-D CNNs for False Positive
592 Reduction in Pulmonary Nodule Detection. *IEEE Trans. Biomed. Eng.* **64**, 1558-1567 (2017).

593 53 Lao, J. *et al.* A Deep Learning-Based Radiomics Model for Prediction of Survival in Glioblastoma
594 Multiforme. *Sci. Rep.* **7**, 10353 (2017).

595 54 Setio, A. A. A. *et al.* Validation, comparison, and combination of algorithms for automatic detection of
596 pulmonary nodules in computed tomography images: The LUNA16 challenge. *Med. Image Anal.* **42**, 1-
597 13 (2017).

598 55 Ehteshami Bejnordi, B. *et al.* Diagnostic Assessment of Deep Learning Algorithms for Detection of
599 Lymph Node Metastases in Women With Breast Cancer. *JAMA* **318**, 2199-2210 (2017).

600 56 Mohamed, A. A. *et al.* A deep learning method for classifying mammographic breast density categories.
601 *Med. Phys.* **45**, 314-321 (2018).

602 57 Khosravi, P., Kazemi, E., Imielinski, M., Elemento, O. & Hajirasouliha, I. Deep Convolutional Neural
603 Networks Enable Discrimination of Heterogeneous Digital Pathology Images. *EBioMedicine* **27**, 317-
604 328 (2018).

605 58 Xiao, Y., Wu, J., Lin, Z. & Zhao, X. A deep learning-based multi-model ensemble method for cancer
606 prediction. *Comput. Methods Programs Biomed.* **153**, 1-9 (2018).

607 59 Marchetti, M. A. *et al.* Results of the 2016 International Skin Imaging Collaboration International
608 Symposium on Biomedical Imaging challenge: Comparison of the accuracy of computer algorithms to
609 dermatologists for the diagnosis of melanoma from dermoscopic images. *J. Am. Acad. Dermatol.* **78**,
610 270-277.e271 (2018).

611 60 Chen, P.-J. *et al.* Accurate Classification of Diminutive Colorectal Polyps Using Computer-Aided
612 Analysis. *Gastroenterology* **154**, 568-575 (2018).

613 61 Bychkov, D. *et al.* Deep learning based tissue analysis predicts outcome in colorectal cancer. *Sci. Rep.*
614 **8**, 3395 (2018).

615 62 Yasaka, K., Akai, H., Abe, O. & Kiryu, S. Deep Learning with Convolutional Neural Network for
616 Differentiation of Liver Masses at Dynamic Contrast-enhanced CT: A Preliminary Study. *Radiology*
617 **286**, 887-896 (2018).

618 63 Chang, K. *et al.* Residual Convolutional Neural Network for the Determination of IDH Status in Low-
619 and High-Grade Gliomas from MR Imaging. *Clin. Cancer Res.* **24**, 1073-1081 (2018).

620 64 Ribli, D., Horváth, A., Unger, Z., Pollner, P. & Csabai, I. Detecting and classifying lesions in
621 mammograms with Deep Learning. *Sci. Rep.* **8**, 4165 (2018).

622 65 Chaudhary, K., Poirion, O. B., Lu, L. & Garmire, L. X. Deep Learning–Based Multi-Omics Integration
623 Robustly Predicts Survival in Liver Cancer. *Clin. Cancer Res.* **24**, 1248-1259 (2018).

624 66 Mobadersany, P. *et al.* Predicting cancer outcomes from histology and genomics using convolutional
625 networks. *Proc. Natl. Acad. Sci. U. S. A.* **115**, E2970-E2979 (2018).

626 67 Saltz, J. *et al.* Spatial Organization and Molecular Correlation of Tumor-Infiltrating Lymphocytes
627 Using Deep Learning on Pathology Images. *Cell Rep.* **23**, 181-193.e187 (2018).

628 68 van de Goor, R., van Hooren, M., Dingemans, A.-M., Kremer, B. & Kross, K. Training and Validating
629 a Portable Electronic Nose for Lung Cancer Screening. *J. Thorac. Oncol.* **13**, 676-681 (2018).

630 69 Chang, H., Han, J., Zhong, C., Snijders, A. M. & Mao, J. Unsupervised Transfer Learning via Multi-
631 Scale Convolutional Sparse Coding for Biomedical Applications. *IEEE Trans. Pattern Anal. Mach.*
632 *Intell.* **40**, 1182-1194 (2018).

633 70 Han, S. S. *et al.* Classification of the Clinical Images for Benign and Malignant Cutaneous Tumors
634 Using a Deep Learning Algorithm. *J. Invest. Dermatol.* **138**, 1529-1538 (2018).

635 71 Hirasawa, T. *et al.* Application of artificial intelligence using a convolutional neural network for
636 detecting gastric cancer in endoscopic images. *Gastric Cancer* **21**, 653-660 (2018).

637 72 Chang, P. *et al.* Deep-Learning Convolutional Neural Networks Accurately Classify Genetic Mutations
638 in Gliomas. *Am. J. Neuroradiol.* **39**, 1201-1207 (2018).

639 73 Haenssle, H. A. *et al.* Man against machine: diagnostic performance of a deep learning convolutional
640 neural network for dermoscopic melanoma recognition in comparison to 58 dermatologists. *Ann. Oncol.*
641 **29**, 1836-1842 (2018).

642 74 Coudray, N. *et al.* Classification and mutation prediction from non–small cell lung cancer
643 histopathology images using deep learning. *Nat. Med.* **24**, 1559-1567 (2018).

644 75 Wang, P. *et al.* Development and validation of a deep-learning algorithm for the detection of polyps
645 during colonoscopy. *Nat. Biomed. Eng.* **2**, 741-748 (2018).

646 76 Urban, G. *et al.* Deep Learning Localizes and Identifies Polyps in Real Time With 96% Accuracy in
647 Screening Colonoscopy. *Gastroenterology* **155**, 1069-1078.e1068 (2018).

648 77 Rajpurkar, P. *et al.* Deep learning for chest radiograph diagnosis: A retrospective comparison of the
649 CheXNeXt algorithm to practicing radiologists. *PLoS Med.* **15**, e1002686 (2018).

650 78 Hosny, A. *et al.* Deep learning for lung cancer prognostication: A retrospective multi-cohort radiomics
651 study. *PLoS Med.* **15**, e1002711 (2018).

652 79 Nam, J. G. *et al.* Development and Validation of Deep Learning–based Automatic Detection Algorithm
653 for Malignant Pulmonary Nodules on Chest Radiographs. *Radiology* **290**, 218-228 (2019).

654 80 Byrne, M. F. *et al.* Real-time differentiation of adenomatous and hyperplastic diminutive colorectal
655 polyps during analysis of unaltered videos of standard colonoscopy using a deep learning model. *Gut*
656 **68**, 94-100 (2019).

657 81 Horie, Y. *et al.* Diagnostic outcomes of esophageal cancer by artificial intelligence using convolutional
658 neural networks. *Gastrointest. Endosc.* **89**, 25-32 (2019).

659 82 Kather, J. N. *et al.* Predicting survival from colorectal cancer histology slides using deep learning: A
660 retrospective multicenter study. *PLoS Med.* **16**, e1002730 (2019).

661 83 Rodríguez-Ruiz, A. *et al.* Detection of Breast Cancer with Mammography: Effect of an Artificial
662 Intelligence Support System. *Radiology* **290**, 305-314 (2019).

663 84 Li, X. *et al.* Diagnosis of thyroid cancer using deep convolutional neural network models applied to
664 sonographic images: a retrospective, multicohort, diagnostic study. *Lancet Oncol.* **20**, 193-201 (2019).

665 85 Wang, S. *et al.* Predicting EGFR Mutation Status in Lung Adenocarcinoma on CT Image Using Deep
666 Learning. *Eur. Respir. J.*, 1800986 (2019).

667 86 Brinker, T. J. *et al.* A convolutional neural network trained with dermoscopic images performed on par
668 with 145 dermatologists in a clinical melanoma image classification task. *Eur. J. Cancer* **111**, 148-154
669 (2019).

670 87 Kickingereder, P. *et al.* Automated quantitative tumour response assessment of MRI in neuro-oncology
671 with artificial neural networks: a multicentre, retrospective study. *Lancet Oncol.* **20**, 728-740 (2019).

672 88 Brinker, T. J. *et al.* Deep learning outperformed 136 of 157 dermatologists in a head-to-head
673 dermoscopic melanoma image classification task. *Eur. J. Cancer* **113**, 47-54 (2019).

674 89 Choi, K. S., Choi, S. H. & Jeong, B. Prediction of IDH genotype in gliomas with dynamic susceptibility
675 contrast perfusion MR imaging using an explainable recurrent neural network. *Neuro Oncol.* **21**, 1197-
676 1209 (2019).

677 90 Ardila, D. *et al.* End-to-end lung cancer screening with three-dimensional deep learning on low-dose
678 chest computed tomography. *Nat. Med.* **25**, 954-961 (2019).

679 91 Yala, A., Lehman, C., Schuster, T., Portnoi, T. & Barzilay, R. A Deep Learning Mammography-based
680 Model for Improved Breast Cancer Risk Prediction. *Radiology* **292**, 60-66 (2019).

681 92 Kather, J. N. *et al.* Deep learning can predict microsatellite instability directly from histology in
682 gastrointestinal cancer. *Nat. Med.* **25**, 1054-1056 (2019).

683 93 Liu, Y. *et al.* Artificial Intelligence–Based Breast Cancer Nodal Metastasis Detection: Insights Into the
684 Black Box for Pathologists. *Arch. Pathol. Lab. Med.* **143**, 859-868 (2019).

685 94 Kehl, K. L. *et al.* Assessment of Deep Natural Language Processing in Ascertaining Oncologic
686 Outcomes From Radiology Reports. *JAMA Oncol.* **5**, 1421-1429 (2019).

687 95 Campanella, G. *et al.* Clinical-grade computational pathology using weakly supervised deep learning on
688 whole slide images. *Nat. Med.* **25**, 1301-1309 (2019).

689 96 Chen, P.-H. C. *et al.* An augmented reality microscope with real-time artificial intelligence integration
690 for cancer diagnosis. *Nat. Med.* **25**, 1453-1457 (2019).

691 97 Hu, L. *et al.* An Observational Study of Deep Learning and Automated Evaluation of Cervical Images
692 for Cancer Screening. *J. Natl. Cancer Inst.* **111**, 923-932 (2019).

693 98 Rodriguez-Ruiz, A. *et al.* Stand-Alone Artificial Intelligence for Breast Cancer Detection in
694 Mammography: Comparison With 101 Radiologists. *J. Natl. Cancer Inst.* **111**, 916-922 (2019).

695 99 Wang, X. *et al.* Weakly Supervised Deep Learning for Whole Slide Lung Cancer Image Analysis. *IEEE*
696 *Trans. Cybern.*, 1-13 (2019).

697 100 Jurmeister, P. *et al.* Machine learning analysis of DNA methylation profiles distinguishes primary lung
698 squamous cell carcinomas from head and neck metastases. *Sci. Transl. Med.* **11**, eaaw8513 (2019).

699 101 Courtiol, P. *et al.* Deep learning-based classification of mesothelioma improves prediction of patient
700 outcome. *Nat. Med.* **25**, 1519-1525 (2019).

701 102 Wang, P. *et al.* Real-time automatic detection system increases colonoscopic polyp and adenoma
702 detection rates: a prospective randomised controlled study. *Gut* **68**, 1813-1819 (2019).

703 103 Liao, F., Liang, M., Li, Z., Hu, X. & Song, S. Evaluate the Malignancy of Pulmonary Nodules Using
704 the 3-D Deep Leaky Noisy-OR Network. *IEEE Trans. Neural Netw. Learn. Syst.* **30**, 3484-3495 (2019).

705 104 Luo, H. *et al.* Real-time artificial intelligence for detection of upper gastrointestinal cancer by
706 endoscopy: a multicentre, case-control, diagnostic study. *Lancet Oncol.* **20**, 1645-1654 (2019).

707 105 Wu, L. *et al.* Randomised controlled trial of WISENSE, a real-time quality improving system for
708 monitoring blind spots during esophagogastroduodenoscopy. *Gut* **68**, 2161-2169 (2019).

709 106 Shkolyar, E. *et al.* Augmented Bladder Tumor Detection Using Deep Learning. *Eur. Urol.* **76**, 714-718
710 (2019).

711 107 Yamamoto, Y. *et al.* Automated acquisition of explainable knowledge from unannotated histopathology
712 images. *Nat. Commun.* **10**, 5642 (2019).

713 108 McKinney, S. M. *et al.* International evaluation of an AI system for breast cancer screening. *Nature*
714 **577**, 89-94 (2020).

715 109 Hollon, T. C. *et al.* Near real-time intraoperative brain tumor diagnosis using stimulated Raman
716 histology and deep neural networks. *Nat. Med.* **26**, 52-58 (2020).

717 110 Haenssle, H. A. *et al.* Man against machine reloaded: performance of a market-approved convolutional
718 neural network in classifying a broad spectrum of skin lesions in comparison with 96 dermatologists
719 working under less artificial conditions. *Ann. Oncol.* **31**, 137-143 (2020).

720 111 Ström, P. *et al.* Artificial intelligence for diagnosis and grading of prostate cancer in biopsies: a
721 population-based, diagnostic study. *Lancet Oncol.* **21**, 222-232 (2020).

722 112 Bulten, W. *et al.* Automated deep-learning system for Gleason grading of prostate cancer using
723 biopsies: a diagnostic study. *Lancet Oncol.* **21**, 233-241 (2020).

724 113 Skrede, O.-J. *et al.* Deep learning for prediction of colorectal cancer outcome: a discovery and
725 validation study. *Lancet* **395**, 350-360 (2020).

726 114 Saillard, C. *et al.* Predicting survival after hepatocellular carcinoma resection using deep-learning on
727 histological slides. *Hepatology* **Advance online publication**, <https://doi.org/10.1002/hep.31207> (2020).

728 115 Jin, E. H. *et al.* Improved Accuracy in Optical Diagnosis of Colorectal Polyps Using Convolutional
729 Neural Networks with Visual Explanations. *Gastroenterology* **158**, 2169-2179.e2168 (2020).

730 116 de Groof, A. J. *et al.* Deep-Learning System Detects Neoplasia in Patients With Barrett's Esophagus
731 With Higher Accuracy Than Endoscopists in a Multistep Training and Validation Study With
732 Benchmarking. *Gastroenterology* **158**, 915-929.e914 (2020).

733 117 Bangalore Yogananda, C. G. *et al.* A novel fully automated MRI-based deep-learning method for
734 classification of IDH mutation status in brain gliomas. *Neuro Oncol.* **22**, 402-411 (2020).

735 118 Zheng, X. *et al.* Deep learning radiomics can predict axillary lymph node status in early-stage breast
736 cancer. *Nat. Commun.* **11**, 1236 (2020).

737 119 Galateau Salle, F. *et al.* Comprehensive Molecular and Pathologic Evaluation of Transitional
738 Mesothelioma Assisted by Deep Learning Approach: A Multi-Institutional Study of the International
739 Mesothelioma Panel from the MESOPATH Reference Center. *J. Thorac. Oncol.* **15**, 1037-1053 (2020).

740 120 Baldwin, D. R. *et al.* External validation of a convolutional neural network artificial intelligence tool to
741 predict malignancy in pulmonary nodules. *Thorax* **75**, 306-312 (2020).

742 121 Wang, P. *et al.* Effect of a deep-learning computer-aided detection system on adenoma detection during
743 colonoscopy (CADE-DB trial): a double-blind randomised study. *Lancet Gastroenterol. Hepatol.* **5**,
744 343-351 (2020).

745 122 Song, Q., Zheng, Y., Sheng, W. & Yang, J. Tridirectional Transfer Learning for Predicting Gastric
746 Cancer Morbidity. *IEEE Trans. Neural Netw. Learn. Syst.*, 1-14 (2020).

747 123 Dong, D. *et al.* Deep learning radiomic nomogram can predict the number of lymph node metastasis in
748 locally advanced gastric cancer: an international multicenter study. *Ann. Oncol.* **31**, 912-920 (2020).

749 124 Shin, H. *et al.* Early-Stage Lung Cancer Diagnosis by Deep Learning-Based Spectroscopic Analysis of
750 Circulating Exosomes. *ACS Nano* **14**, 5435-5444 (2020).

751 125 Kann, B. H. *et al.* Multi-Institutional Validation of Deep Learning for Pretreatment Identification of
752 Extranodal Extension in Head and Neck Squamous Cell Carcinoma. *J. Clin. Oncol.* **38**, 1304-1311
753 (2020).

754 126 Nature. AI diagnostics need attention. *Nature* **555**, 285 (2018).

755 127 The Lancet. Is digital medicine different? *Lancet* **392**, 95 (2018).

756 128 Kawaguchi, K., Kaelbling, L. P. & Bengio, Y. Generalization in Deep Learning. Preprint at
757 <https://arxiv.org/abs/1710.05468> (2017).

758 129 LeCun, Y. in *Connectionism in perspective* (ed. Pfeifer, R., Schreter, Z., Fogelman, F., & Steels, L.)
759 143-156 (Elsevier, Zürich, Switzerland, 1989).

760 130 Neyshabur, B., Bhojanapalli, S., Mcallester, D. & Srebro, N. Exploring Generalization in Deep
761 Learning. *Adv. Neural Inf. Process. Syst.* **30**, 5947-5956 (2017).

762 131 Pan, S. J. & Yang, Q. A Survey on Transfer Learning. *IEEE Trans. Knowl. Data Eng.* **22**, 1345-1359
763 (2010).

764 132 Weiss, K., Khoshgoftaar, T. M. & Wang, D. A survey of transfer learning. *J. Big Data* **3**, 9 (2016).

765 133 Deng, J. *et al.* ImageNet: A large-scale hierarchical image database. *Proc. IEEE Conf. Comput. Vis.*
766 *Pattern Recognit.*, 248-255 (2009).

767 134 Russakovsky, O. *et al.* ImageNet Large Scale Visual Recognition Challenge. *Int. J. Comput. Vis.* **115**,
768 211-252 (2015).

769 135 Shankar, S. *et al.* No Classification without Representation: Assessing Geodiversity Issues in Open
770 Data Sets for the Developing World. *NIPS Workshop Mach. Learn. Dev. World* (2017).

771 136 Geirhos, R. *et al.* ImageNet-trained CNNs are biased towards texture; increasing shape bias improves
772 accuracy and robustness. *Proc. Int. Conf. Learn. Represent.* (2019).

773 137 Beyer, L., Hénaff, O. J., Kolesnikov, A., Zhai, X. & van den Oord, A. Are we done with ImageNet?
774 Preprint at <https://arxiv.org/abs/2006.07159> (2020).

775 138 Sun, C., Shrivastava, A., Singh, S. & Gupta, A. Revisiting Unreasonable Effectiveness of Data in Deep
776 Learning Era. *Proc. IEEE Int. Conf. Comput. Vis.*, 843-852 (2017).

777 139 Simard, P. Y., Steinkraus, D. & Platt, J. C. Best practices for convolutional neural networks applied to
778 visual document analysis. *Proc. 7th Int. Conf. Doc. Anal. Recognit.*, 958-963 (2003).

779 140 Baird, H. S. Document image defect models and their uses. *Proc. 2nd Int. Conf. Doc. Anal. Recognit.*,
780 62-67 (1993).

781 141 Stacke, K., Eilertsen, G., Unger, J. & Lundstrom, C. Measuring Domain Shift for Deep Learning in
782 Histopathology. *IEEE J. Biomed. Health Inform.* **Advance online publication**,
783 <https://doi.org/10.1109/JBHI.2020.3032060> (2020).

784 142 Lakhani, P. & Sundaram, B. Deep Learning at Chest Radiography: Automated Classification of
785 Pulmonary Tuberculosis by Using Convolutional Neural Networks. *Radiology* **284**, 574-582 (2017).

786 143 Hussain, Z., Gimenez, F., Yi, D. & Rubin, D. Differential Data Augmentation Techniques for Medical
787 Imaging Classification Tasks. *AMIA Annu. Symp. Proc.* **2017**, 979-984 (2018).

788 144 Sajjad, M. *et al.* Multi-grade brain tumor classification using deep CNN with extensive data
789 augmentation. *J. Comput. Sci.* **30**, 174-182 (2019).

790 145 Tellez, D. *et al.* Quantifying the effects of data augmentation and stain color normalization in
791 convolutional neural networks for computational pathology. *Med. Image Anal.* **58**, 101544 (2019).

792 146 Kerr, R. S. *et al.* Adjuvant capecitabine plus bevacizumab versus capecitabine alone in patients with
793 colorectal cancer (QUASAR 2): an open-label, randomised phase 3 trial. *Lancet Oncol.* **17**, 1543-1557
794 (2016).

795 147 Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J. & Wojna, Z. Rethinking the Inception Architecture for
796 Computer Vision. *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2818-2826 (2016).

797 148 Miller, R. G. J. *Simultaneous Statistical Inference*, 2nd edn. (Springer New York, 1981).

798 149 Hochberg, Y. & Tamhane, A. C. *Multiple Comparison Procedures*. (Wiley, 2009).

799 150 Michiels, S., Koscielny, S. & Hill, C. Prediction of cancer outcome with microarrays: a multiple
800 random validation strategy. *Lancet* **365**, 488-492 (2005).

801 151 Hastie, T., Tibshirani, R. & Friedman, J. *The Elements of Statistical Learning*, 2nd edn. (Springer-
802 Verlag New York, 2009).

803 152 Russell, S. & Norvig, P. *Artificial Intelligence: A Modern Approach*, 3rd edn. (Prentice Hall, 2010).

804 153 Hemingway, H., Riley, R. D. & Altman, D. G. Ten steps towards improving prognosis research. *BMJ*
805 **339**, b4184 (2009).

806 154 Korevaar, D. A. *et al.* Facilitating Prospective Registration of Diagnostic Accuracy Studies: A STARD
807 Initiative. *Clin. Chem.* **63**, 1331-1341 (2017).

808 155 Ioannidis, J. P. A. The Importance of Predefined Rules and Prespecified Statistical Analyses: Do Not
809 Abandon Significance. *JAMA* **321**, 2067-2068 (2019).

810 156 Brodersen, K. H., Ong, C. S., Stephan, K. E. & Buhmann, J. M. The Balanced Accuracy and Its
811 Posterior Distribution. *Proc. 20th Int. Conf. Pattern Recognit.*, 3121-3124 (2010).

812 157 van den Hout, W. B. The Area under an ROC Curve with Limited Information. *Med. Decis. Making* **23**,
813 160-166 (2003).

814 158 Fawcett, T. An introduction to ROC analysis. *Pattern Recognit. Lett.* **27**, 861-874 (2006).

815 159 Harrell, F. E., Jr, Califf, R. M., Pryor, D. B., Lee, K. L. & Rosati, R. A. Evaluating the yield of medical
816 tests. *J. Am. Med. Assoc.* **247**, 2543-2546 (1982).

817 160 Lobo, J. M., Jiménez-Valverde, A. & Real, R. AUC: a misleading measure of the performance of
818 predictive distribution models. *Glob. Ecol. Biogeogr.* **17**, 145-151 (2008).

819 161 Voosen, P. *How AI detectives are cracking open the black box of deep learning*. Science.
820 [https://www.sciencemag.org/news/2017/07/how-ai-detectives-are-cracking-open-black-box-deep-](https://www.sciencemag.org/news/2017/07/how-ai-detectives-are-cracking-open-black-box-deep-learning)
821 [learning](https://www.sciencemag.org/news/2017/07/how-ai-detectives-are-cracking-open-black-box-deep-learning) (2017).

822 162 Adadi, A. & Berrada, M. Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence
823 (XAI). *IEEE Access* **6**, 52138-52160 (2018).

824 163 Barredo Arrieta, A. *et al.* Explainable Artificial Intelligence (XAI): Concepts, taxonomies,
825 opportunities and challenges toward responsible AI. *Inf. Fusion* **58**, 82-115 (2020).

826 164 Montavon, G., Samek, W. & Müller, K.-R. Methods for interpreting and understanding deep neural
827 networks. *Digit. Signal Process.* **73**, 1-15 (2018).

828 165 Simonyan, K., Vedaldi, A. & Zisserman, A. Deep Inside Convolutional Networks: Visualising Image
829 Classification Models and Saliency Maps. *Proc. Int. Conf. Learn. Represent.* (2014).

830 166 Bach, S. *et al.* On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise
831 Relevance Propagation. *PLoS One* **10**, e0130140 (2015).

832 167 Sundararajan, M., Taly, A. & Yan, Q. Axiomatic attribution for deep networks. *Proc. 34th Int. Conf.*
833 *Mach. Learn.* **70**, 3319-3328 (2017).

834 168 Friedman, L. M., Furberg, C. D., DeMets, D. L., Reboussin, D. M. & Granger, C. B. *Fundamentals of*
835 *Clinical Trials*, 5th edn. (Springer, 2015).

836 169 van Luijn, H. E. M., Musschenga, A. W., Keus, R. B., Robinson, W. M. & Aaronson, N. K. Assessment
837 of the risk/benefit ratio of phase II cancer clinical trials by Institutional Review Board (IRB) members.
838 *Ann. Oncol.* **13**, 1307-1313 (2002).

839 170 Martin, L., Hutchens, M., Hawkins, C. & Radnov, A. How much do clinical trials cost? *Nat. Rev. Drug*
840 *Discov.* **16**, 381-382 (2017).

841 171 Teutsch, S. M. *et al.* The Evaluation of Genomic Applications in Practice and Prevention (EGAPP)
842 initiative: methods of the EGAPP Working Group. *Genet. Med.* **11**, 3-14 (2009).

843 172 Vollmer, S. *et al.* Machine learning and artificial intelligence research for patient benefit: 20 critical
844 questions on transparency, replicability, ethics, and effectiveness. *BMJ* **368**, l6927 (2020).

845 173 Chan, A.-W. *et al.* SPIRIT 2013 explanation and elaboration: guidance for protocols of clinical trials.
846 *BMJ* **346**, e7586 (2013).

847 174 Cruz Rivera, S. *et al.* Guidelines for clinical trial protocols for interventions involving artificial
848 intelligence: the SPIRIT-AI extension. *Nat. Med.* **26**, 1351-1363 (2020).

849 175 Moher, D. *et al.* CONSORT 2010 Explanation and Elaboration: updated guidelines for reporting
850 parallel group randomised trials. *BMJ* **340**, c869 (2010).

851 176 Collins, G. S. & Moons, K. G. M. Reporting of artificial intelligence prediction models. *Lancet* **393**,
852 1577-1579 (2019).

853 177 Liu, X. *et al.* Reporting guidelines for clinical trial reports for interventions involving artificial
854 intelligence: the CONSORT-AI extension. *Nat. Med.* **26**, 1364-1374 (2020).

855 178 The Lancet. Should protocols for observational research be registered? *Lancet* **375**, 348 (2010).

856 179 Loder, E., Groves, T. & MacAuley, D. Registration of observational studies. *BMJ* **340**, c950 (2010).

857 180 Chambers, C. & Munafo, M. *Trust in science would be improved by study pre-registration*. The
858 Guardian. [https://www.theguardian.com/science/blog/2013/jun/05/trust-in-science-study-pre-](https://www.theguardian.com/science/blog/2013/jun/05/trust-in-science-study-pre-registration)
859 registration (2013).

860 181 Williams, R. J., Tse, T., Harlan, W. R. & Zarin, D. A. Registration of observational studies: Is it time?
861 *Can. Med. Assoc. J.* **182**, 1638-1642 (2010).

862 182 Gill, J. & Prasad, V. Improving observational studies in the era of big data. *Lancet* **392**, 716-717 (2018).

863 183 Sørensen, H. T. & Rothman, K. J. The prognosis for research. *BMJ* **340**, c703 (2010).

864 184 Vandenbroucke, J. P. Registering observational research: second thoughts. *Lancet* **375**, 982-983 (2010).

865 185 Epidemiology. The Registration of Observational Studies—When Metaphors Go Bad. *Epidemiology*
866 **21**, 607-609 (2010).

867 186 Andre, F. *et al.* Biomarker studies: a call for a comprehensive biomarker study registry. *Nat. Rev. Clin.*
868 *Oncol.* **8**, 171-176 (2011).

869 187 Hooft, L. & Bossuyt, P. M. Prospective Registration of Marker Evaluation Studies: Time to Act. *Clin.*
870 *Chem.* **57**, 1684-1686 (2011).

871 188 Altman, D. G. The Time Has Come to Register Diagnostic and Prognostic Research. *Clin. Chem.* **60**,
872 580-582 (2014).

873 189 Rifai, N. *et al.* Registering Diagnostic and Prognostic Trials of Tests: Is It the Right Thing to Do? *Clin.*
874 *Chem.* **60**, 1146-1152 (2014).

875 190 Rajkomar, A., Dean, J. & Kohane, I. Machine Learning in Medicine. *N. Engl. J. Med.* **380**, 1347-1358
876 (2019).

877 191 Zou, J. & Schiebinger, L. AI can be sexist and racist - it's time to make it fair. *Nature* **559**, 324-326
878 (2018).

879 192 Adamson, A. S. & Smith, A. Machine Learning and Health Care Disparities in Dermatology. *JAMA*
880 *Dermatol.* **154**, 1247-1248 (2018).

- 881 193 Vyas, D. A., Eisenstein, L. G. & Jones, D. S. Hidden in Plain Sight — Reconsidering the Use of Race
882 Correction in Clinical Algorithms. *N. Engl. J. Med.* **383**, 874-882 (2020).
- 883 194 Obermeyer, Z., Powers, B., Vogeli, C. & Mullainathan, S. Dissecting racial bias in an algorithm used to
884 manage the health of populations. *Science* **366**, 447-453 (2019).
- 885 195 Rajkomar, A., Hardt, M., Howell, M. D., Corrado, G. & Chin, M. H. Ensuring Fairness in Machine
886 Learning to Advance Health Equity. *Ann. Intern. Med.* **169**, 866-872 (2018).
- 887 196 Owens, K. & Walker, A. Those designing healthcare algorithms must become actively anti-racist. *Nat.*
888 *Med.* **26**, 1327-1328 (2020).
- 889 197 Moons, K. G. M. *et al.* Risk prediction models: II. External validation, model updating, and impact
890 assessment. *Heart* **98**, 691-698 (2012).

891

892 **Acknowledgements**

893 We thank Marian Seiergren for assembling all figures, Tarjei Sveinngjerd Hveem for discussions, Trond Ystanes,
894 Haakon André Inderhaug and Bjørn Morten Sannes for setting up and maintaining our computer network and
895 computational infrastructure, and the authors of Inception-v3 for making their code freely available under an
896 open source licence (Apache License, Version 2.0). We acknowledge funding from the Research Council of
897 Norway through its IKTPLUSS Lighthouse program (project number 259204).

898

899 **Author contributions**

900 H.E.D and D.J.K initiated the project. All authors researched data for the article. A.K., O.J.S. and K.L. assessed
901 papers in the review of recent, presumably influential deep learning studies in cancer diagnostics. S.D.R.
902 executed the training, tuning and evaluation of Inception-v3 systems. A.K. drafted the manuscript, and all
903 authors contributed to reviewing and editing the manuscript.

904

905 **Competing interests**

906 The authors declare no competing interests.

907

908 **Peer review information**

909 *Nature Reviews Cancer* thanks J. Kather and the other, anonymous, reviewer(s) for their contribution to the peer
910 review of this work.

911

912 **Related links**

913 Journal Policies in the Journal of Clinical Oncology: <https://ascopubs.org/jco/authors/journal-policies>

914 ClinicalTrials.gov registry: <https://www.clinicaltrials.gov>

915 International Standard Randomised Controlled Trial Number (ISRCTN) registry: <https://www.isrctn.com>

916

917 **Supplementary information**

918 Supplementary information is available for this paper at <https://doi.org/10.1038/s415XX-XXX-XXXX-X>

919

920 **Glossary**

921 **Artificial neural networks**

922 Mathematical functions mapping input data to output representations, structured as a directed graph of nodes and
923 edges.

924 **Deep learning**

925 A class of machine learning methods that make use of successively more abstract representations of the input
926 data to perform a specific task, typically implemented using artificial neural networks. They also consist of an
927 objective function that compares the final output with a target output as well as an optimisation method that is
928 used to optimise the objective function.

929 **Deep learning models**

930 Computational models obtained by training deep neural networks. Note that a single training of a neural network
931 produces a sequence of models since each new optimisation iteration produces a model slightly different from
932 the previous one. A tuning dataset may be used to select among these models.

933 **Deep learning systems**

934 Systems utilising one or more deep learning models to make predictions. A system's output may be a function of
935 the outputs of the models, e.g. by averaging and thresholding the model outputs.

936 **Supervised machine learning**

937 A methodology in which learning occurs by mimicking the mapping of input data to target output labels. In
938 contrast, the input data are not associated with any output labels in unsupervised learning.

939 **Capacity**

940 The ability of a model class, e.g. a particular network architecture, to express complicated correlations between
941 input data and target output. Model classes with high capacity have the potential to produce models that are able
942 to map training data to target outputs with a high degree of accuracy, but are also more prone to overfitting.

943 **Development cohort**

944 A cohort used for training and sometimes tuning and internal validation of a system.

945 **External cohorts**

946 Also known as independent cohorts, these differ non-randomly from the development cohort. In cancer
947 diagnostics, the external cohorts will often contain patients suspected of having the same disease or disease
948 attribute, at risk of developing the same event or suspected to respond to the same treatment as patients in the
949 development cohort. However, external cohorts may be intentionally more different from the development
950 cohort.

951 **Training**

952 Optimisation of model parameters based on data.

953 **Tuning**

954 Informed selection of hyperparameter values (parameters not optimised during training) based on data. Examples
955 include network architecture, optimisation method and threshold for a model's continuous output. The
956 nomenclature in machine learning is to use 'validation' instead of 'tuning'.

957 **Test**

958 While frequently used by the machine learning community to refer to an evaluation of a system's performance,
959 we use 'test' to refer to evaluations other than external validations, e.g. internal validations.

960 **External validation**

961 An evaluation of a system's performance on an external cohort that did not influence the development of the
962 system.

963 **Overfitting**

964 Utilising noise or features in the training data that are not generally relevant for the prediction task but cause the
965 system to perform better on the training sample.

966 **Generalisability**

967 The ability of a system to perform similarly on subjects not included in training as on those included in the
968 training. Poor generalisability can be caused by overfitting to the training data or by the lack of generally
969 relevant features in the training data.

970 **Balanced accuracy**

971 A classification performance metric calculated by averaging the proportion of true predicted outcomes across all
972 possible outcomes. For dichotomous outcomes, this reduces to the average between the sensitivity and
973 specificity.

974 **Area under the receiver operating characteristic curve (AUC)**

975 A performance metric measuring the concordance between a dichotomous outcome and the ranking of subjects
976 provided by a continuous or categorical marker. An AUC of 50% indicates random guessing and 100% indicates
977 perfect prediction. For dichotomous markers, AUC and balanced accuracy are equivalent.

978 **Concordance index (c-index)**

979 A performance metric measuring the concordance between a target outcome, usually defined by time-to-event
980 data, and the ranking of subjects provided by a continuous or categorical marker. A c-index of 50% indicates
981 random guessing and 100% indicates perfect prediction. For dichotomous outcomes, c-index and AUC are
982 equivalent.

983 **Figure legends**

984

985 **Fig. 1 | Characteristics of recent, presumably influential deep learning studies in cancer diagnostics. a |**
986 Percentage of studies reporting on the evaluation of a broad or narrow cohort (BOX 2) by year of publication, for
987 all 92 eligible studies. **b |** Percentage of studies specifying one, multiple or no primary performance metrics in
988 the analysis of the external cohort, for the 50 eligible studies that reported on the evaluation of an external
989 cohort. **c |** Percentage of studies specifying a predefined analysis of the external cohort, for the 50 eligible studies
990 that reported on the evaluation of an external cohort. Studies that specified predefined analyses of external
991 cohorts without defining which one was the primary, if any, were categorised as 'Predefined analyses'. Studies
992 with a predefined primary analysis were categorised according to whether the primary analysis was prespecified
993 in a protocol or not.

994

995 **Fig. 2 | Effect of data variation when training deep learning systems.** For each analysis setup, 20 individual
996 deep learning systems were trained and tuned for prediction of colorectal cancer-specific survival using images
997 of haematoxylin and eosin stained sections acquired by both Aperio AT2 (Leica Biosystems, Germany) and
998 NanoZoomer XR (Hamamatsu Photonics, Japan), as in the previously published analyses¹¹³. The individual
999 systems were applied to evaluate the external cohort using NanoZoomer XR slide images, and the c-index of the
1000 system's binary output was computed. Standard box plots were made using Stata/SE 16.1 (StataCorp, TX). The
1001 matched random subset contained the same number of training and tuning patients with and without cancer-
1002 specific death as in the Gloucester cohort, in total 979 patients. **a |** An example image from the training dataset
1003 and the results of applying the maximum possible amount of colour distortion at each step in the random
1004 distortion process used in the published Inception-v3 analyses¹¹³. Generally, the distortion process applies
1005 random colour distortions to an image by: 1) converting the image to HSV colour space, 2) adding a random
1006 value between -0.05 and 0.05 to the hue, 3) scaling the saturation by a random value between 1/1.1 and 1.1, 4)
1007 adding a random value between -0.1 and 0.1 to the saturation, 5) scaling the brightness (or technically the value
1008 channel in the HSV colour space) by a random value between 1/1.1 and 1.1, 6) adding a random value between -
1009 0.1 and 0.1 to the brightness, and 7) converting back to RGB colour space. Intuitively, the leftmost and rightmost
1010 images represent the range of the random colour distortion, i.e. the minimum and maximum possible amount of
1011 colour distortion for the applied distortion process, where the minimum is no colour distortion. Scale bar, 100

1012 μm . **b** | Effect of changing the number of patients in training and tuning when using the original amount of
1013 colour distortion, as depicted in figure part **a**. **c** | Effect of changing the amount of colour distortion when
1014 training and tuning using the matched random subset. Label '0' on the horizontal axis identifies deep learning
1015 systems trained without any colour distortion, label '1' identifies systems trained with the colour distortion
1016 process depicted in figure part **a**, and label '4' identifies systems trained with the colour distortion process
1017 depicted in figure part **d**. **d** | Similar to figure part **a**, but four times the amount of colour distortion was used at
1018 each step in the distortion process. Scale bar, 100 μm . **e** | Effect of changing the amount of colour distortion and
1019 the number of patients and cohorts in training and tuning.

1020

1021 **Fig. 3 | Development and evaluation of deep learning systems.** A deep learning project often begins with
1022 testing a conceptual idea using a pilot software based on a related open source implementation and data easily
1023 available to the researchers. Successful level I studies will typically evolve into explorative testing of different
1024 modelling options that might be more suitable for the particular task. The system that appears to perform best
1025 should be determined in a level II study that includes sufficient amount and variation in the natural training
1026 dataset. Although performance estimates obtained in such studies are often inflated by the use of a subset that
1027 closely resembles the training subset, level II is an important step in the evaluation sequence that could motivate
1028 investigators to pursue evaluation on external cohorts and attract collaborators. As the lack of predefined primary
1029 analysis often entails post hoc adjustments influenced by the performance in the external cohort, we distinguish
1030 between studies without (level III) and with (level IV) a primary analysis unequivocally specified prior to all
1031 investigations that could reveal correlations between input data and target output in the external cohort. If the
1032 medical validity of a deep learning system is established in level IV studies, the indicated medical utility should
1033 be prospectively evaluated in randomised phase III clinical trials where the system directly intervenes with the
1034 current standard of care. If medical utility is demonstrated and necessary governmental agencies approve routine
1035 medical application, the system can be applied in medical practice while monitoring the long-term benefits,
1036 harms and costs of its application.

1037

1038 **Fig. 4 | Reliability of performance estimations in recent, presumably influential deep learning studies in**
1039 **cancer diagnostics.** Percentage of studies categorised in the different levels of deep learning studies or phases of
1040 clinical trials depicted in FIG. 3, for all 92 eligible studies separated by type of input to the neural network. The

1041 input was histopathology section images in 23 (25%) of the studies **(a)**, radiology images in 40 (43%) of the
1042 studies **(b)**, other images in 22 (24%) of the studies **(c)** and other types of input in 7 (8%) of the studies **(d)**.

1043

1044 **Boxes**

1045

1046 **Box 1 | Representation and biases in training data**

1047 As deep learning systems are developed by learning correlations between input data and target outcome directly
1048 from the training data, it is essential that the training data adequately represents the target population^{31,190}.

1049 Otherwise, the system might learn correlations exclusive to the subpopulation represented in the training data
1050 and consequently perform worse on those not represented in the training data to a sufficient extent. Despite this,
1051 systems are often trained on datasets with prominent biases in demographic characteristics such as sex, race or
1052 ethnicity, with the consequence that many systems exhibit substantial discriminatory biases^{32,191,192}. Restricting
1053 the target population to a particular sex, race or ethnicity would be appalling, and the medical application of any
1054 such deep learning system would systematically increase health care disparities. It is therefore pivotal to utilise
1055 truly representative and unbiased data for training deep learning systems in cancer diagnostics. This extends
1056 beyond ensuring representative distributions of relevant demographics in the training dataset. Racial bias may
1057 also be encoded into systems if the target outcome used in the training is affected by histories of unequal
1058 treatment of patients based on race or ethnicity¹⁹³ or is a proxy such as health care cost instead of health needs,
1059 which has been shown to be the reason why a widely used health care prediction algorithm exhibited significant
1060 racial bias¹⁹⁴. Researchers should strive to identify and compensate for any such biases in their datasets, as
1061 failure to do so might reinforce health inequities if the deep learning systems are applied in clinical
1062 practice^{195,196}. Deficient deep learning systems might be identified through rigorous evaluations in external
1063 datasets truly representative of the target population, or representative of minority populations, as well as
1064 through comprehensive analyses of system explainability across different demographic characteristics.

1065

1066 **Box 2 | Approaches for evaluating a deep learning system**

1067 Different approaches for estimating the performance of a deep learning system provide indications of the
1068 system's ability to make accurate predictions in different scenarios. Even if successful, internal and narrow
1069 validations do not indicate a general medical validity in themselves. Successful broad or domain validations
1070 might warrant assessment of the system's medical utility in prospective, randomised phase III clinical trials.

1071 **[bH1] Internal validation**

1072 Internal validation is evaluation of a deep learning system's performance in the development cohort. This can be
1073 done by evaluating the performance in a randomly sampled subset of the development cohort disjoint from the
1074 training and tuning subsets, or by using resampling techniques such as cross-validation or bootstrapping²².

1075 **[bH1] Narrow validation**

1076 Narrow validation is evaluation of a deep learning system's performance based on a cohort that is similar but
1077 differs non-randomly from the development cohort, e.g. on a cohort from the same hospital as the development
1078 cohort, but sampled in a time interval disjoint from the time interval where the development cohort was sampled.
1079 No information from the narrow cohort should have influenced the development of the system, including that it
1080 should be collected and handled separately from the development cohort.

1081 **[bH1] Broad validation**

1082 Broad validation is evaluation of a deep learning system's performance based on a cohort geographically
1083 separate from the development cohort, e.g. from a different hospital or country²². No information from the broad
1084 cohort should have influenced the development of the system.

1085 **[bH1] Domain validation**

1086 Domain validation is evaluation of a deep learning system's performance in a setting that is very different from
1087 the one where the system was developed¹⁹⁷. This includes validation in a cohort with characteristics not
1088 represented by the development cohort, e.g. developing a method on one type and stage of cancer and validating
1089 it on another type or stage of cancer. Other examples are when the validation data are obtained by equipment not
1090 used in the development such as imaging systems from different vendors, or by sample preparation procedures
1091 intentionally different from the ones used for the development cohort. Domain validations should also be narrow
1092 or broad validations, and are typically performed after successful narrow or broad validations.

1094 **Box 3 | Recommended Protocol Items for External Cohort Evaluation of a deep learning System**
1095 **(PIECES) in cancer diagnostics**

1096 **[bH1] Status**

1097 [b1] Specify the date the protocol was last modified.

1098 [b1] Scrupulously elucidate any investigations performed before finalising the protocol that could reveal
1099 correlations between input data and target output in the external cohort, or state that no such investigations were
1100 performed.

1101 **[bH1] System**

1102 [b1] Describe the development of the deep learning system, including utilised cohorts, network architecture,
1103 hyperparameters and any categorisation of the neural network model's output.

1104 [b1] Unequivocally specify how to assay the deep learning system in a blinded fashion for a single, new subject,
1105 including what the system receives as input and what it directly outputs.

1106 **[bH1] External cohort**

1107 [b1] Describe the origin of the cohort, and explain why it should be regarded as external to the development
1108 cohort.

1109 [b1] Precisely define criteria for inclusion and exclusion of subjects and samples, preferably starting from a
1110 consecutive series of subjects.

1111 [b1] Clearly state the medical setting and target population that the cohort represents.

1112 [b1] Specify the acquisition of input data, including whether it was acquired blinded to the deep learning system
1113 and target output. Note the expertise of any humans involved in the process, e.g. that a pathologist annotated the
1114 regions of interest in slide images.

1115 [b1] Specify the ascertainment of target output, including whether it was ascertained blinded to the deep learning
1116 system.

1117 [b1] If multiple external cohorts are planned to be analysed as a pooled cohort, then the preceding five protocol
1118 items should be completed for the pooled cohort and differences between the individual cohorts should be stated.

1119 If multiple external cohorts are to be analysed independently, the five preceding protocol items should be
1120 completed for each cohort, as well as subsequent protocol items if the predefined analyses differ between
1121 cohorts.

1122 **[bH1] Analyses**

1123 [b1] Unequivocally specify the primary analysis, including the target output and the performance metric and/or
1124 statistical test with interpretation.

1125 [b1] If the chosen metric or statistical test depends on other markers, describe how these markers were assayed
1126 and whether done blinded to the deep learning system and target output, and specify how missing values will be
1127 handled.

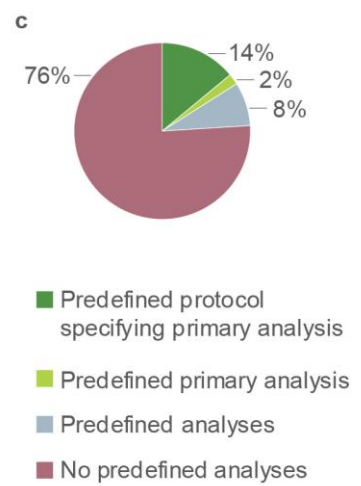
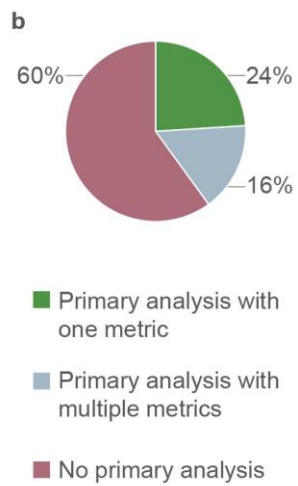
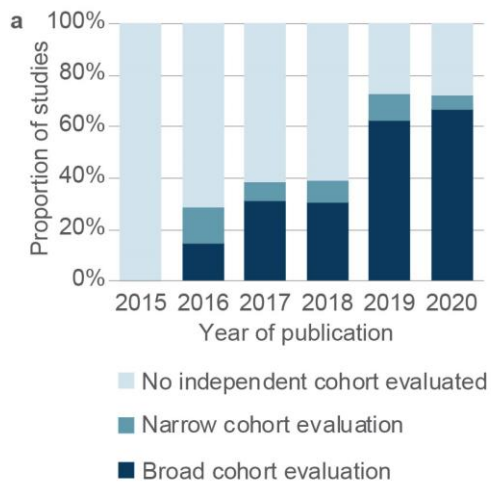
1128 [b1] If the deep learning system was designed to evolve upon usage, e.g. by learning from unlabelled data or
1129 adapting to a cohort, specify that this will not be done when evaluating the external cohort. The system's
1130 prediction should thus not depend on the order in which a set of patients is evaluated and also be identical if the
1131 same patient is evaluated multiple times.

1132 [b1] If additional analyses will be performed and reported in disseminations, e.g. of other deep learning systems,
1133 target outputs, metrics or statistical tests or in specific patient subgroups, specify these analyses in the same
1134 manner as the primary analysis and identify them as secondary analyses.

1135

1136 **Table of Contents Summary**

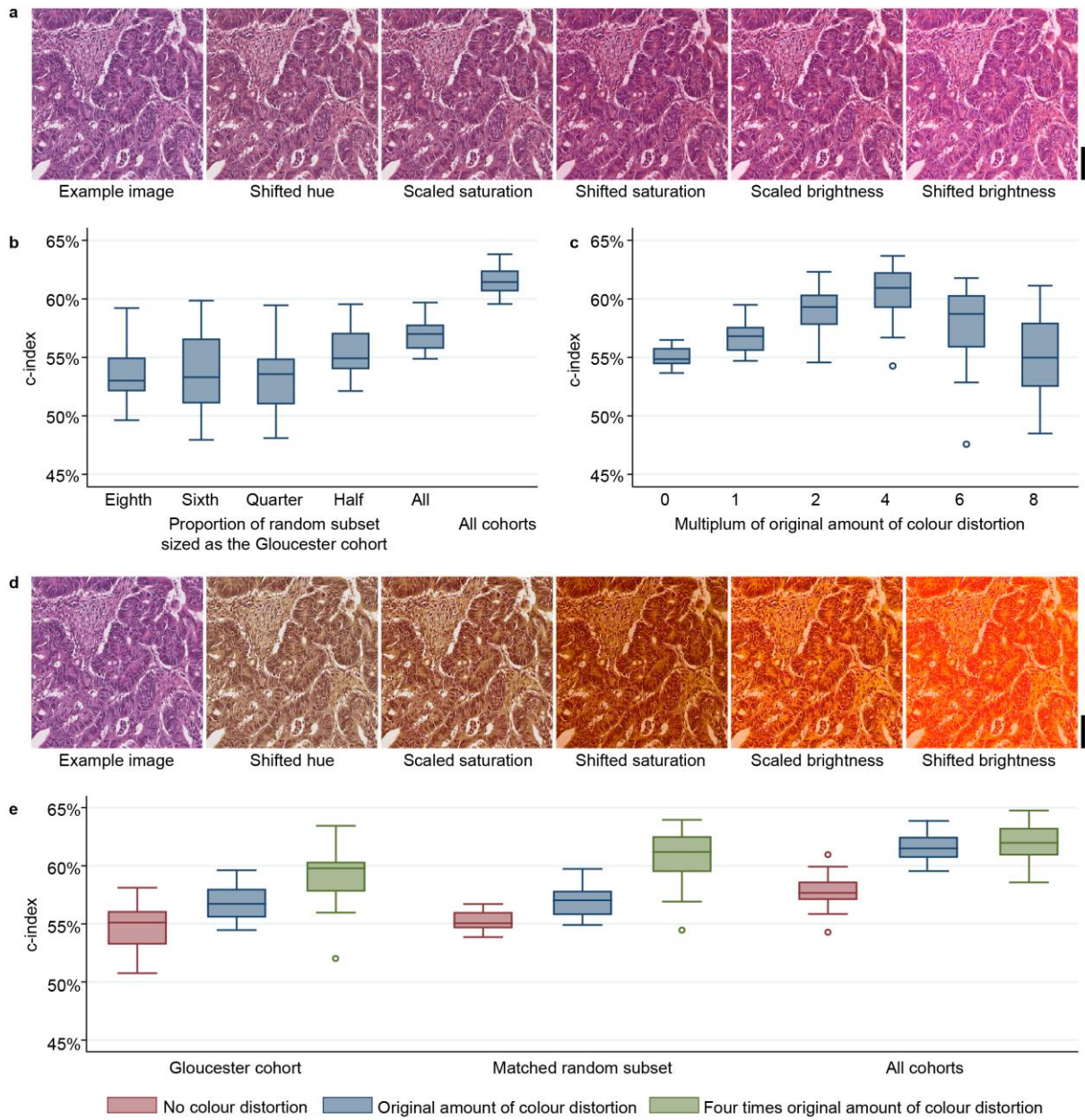
1137 The number of publications on deep learning for cancer diagnostics is rapidly increasing, but clinical translation
1138 is slow. This Perspective advocates performance estimation in external cohorts, and strongly advises that a
1139 primary analysis is predefined in a standardised protocol preferentially stored in an online repository.



1140

1141 Figure 1.

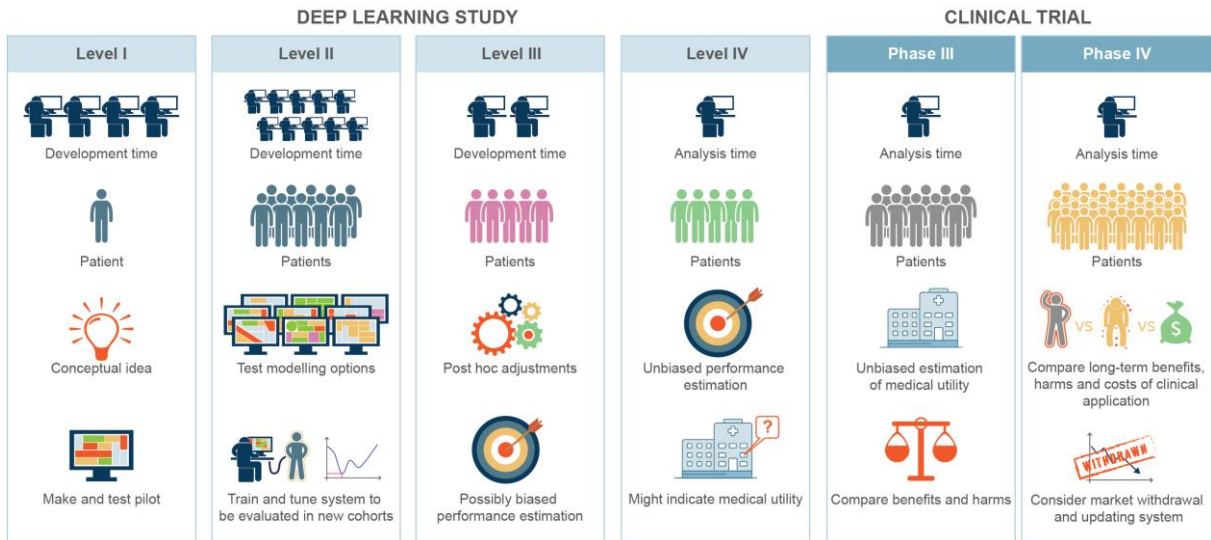
1142



1143

1144 Figure 2.

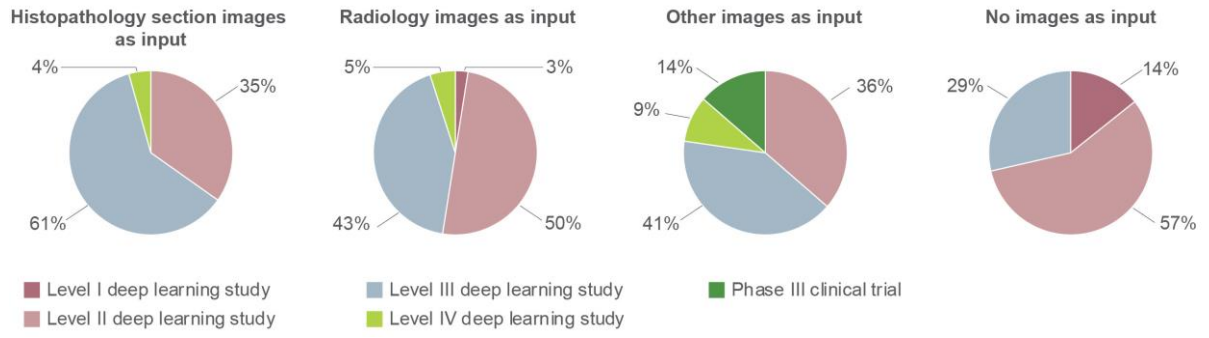
1145



1146

1147 Figure 3.

1148



1149

1150 Figure 4.