

# Designing deep learning studies in cancer diagnostics

Andreas Kleppe<sup>1</sup>, Ole-Johan Skrede<sup>2</sup>, Sepp De Raedt<sup>3</sup>, Knut Liestøl<sup>4</sup>, David J. Kerr<sup>5</sup> and Håvard E. Danielsen<sup>6</sup>

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

Deep learning facilitates utilization of large data sets through direct learning of correlations between raw input data and target output, providing systems that may use intricate structures in high-dimensional input data to accurately model the association with the target output<sup>1,2</sup>. Numerous studies have reported on the applicability of deep learning in cancer diagnostics, including prediction of diagnosis, prognosis and treatment response<sup>3–5</sup>. Although many of these tools are claimed to perform comparably with or better than clinicians, few have yet demonstrated real-world medical utility<sup>6</sup>. This is partly a natural consequence of the time needed for evaluating and adapting systems affecting patient treatment. However, many studies evaluating apparently well-functioning systems are at high risk of bias<sup>6</sup>. Of particular concern is the frequent lack of stringent evaluation of external data<sup>7,8</sup> and that some systems are developed or evaluated on data that are too narrow or inappropriate for the intended medical setting<sup>9–12</sup>. Thus, the lack of a well-established sequence of evaluation steps for converting promising prototypes

into properly evaluated medical systems clearly limits the medical utilization of deep learning systems.

Whereas supervised machine learning techniques traditionally utilized carefully selected representations of the input data to predict the target output, modern deep learning techniques use highly flexible artificial neural networks to correlate input data directly to the target outputs<sup>1,2,13</sup>. The relations learnt by such direct correlation will often be true but may sometimes be spurious phenomena exclusive to the data utilized for learning. In fact, the millions of adjustable parameters make deep neural networks capable of performing perfectly in training sets even when the target outputs are randomly generated and, therefore, utterly meaningless<sup>14</sup>. Thus, the high capacity of neural networks induces serious challenges on how to design and develop deep learning systems, and on how to validate that such a system performs adequately in the intended medical setting<sup>15</sup>. Adequate clinical performance will only be possible if the system has good generalizability to subjects not included in the training data<sup>16,17</sup>.

The design challenge involves issues related to selection of appropriate training data, such as representativeness of the target population (BOX 1), as well as modelling questions such as how the variation of training data may be artificially increased without jeopardising the relationship between input data and target outputs in the training data<sup>18,19</sup>. The validation challenge includes verifying that the system generalizes well, for example performs satisfactorily when evaluated on relevant patient populations at new locations and when input data are obtained using differing laboratory procedures or alternative equipment<sup>15,16</sup>. Moreover, deep learning systems are typically developed iteratively, with repeated testing and often including various selection processes that may bias results<sup>20</sup>. Similar selection issues have been recognized as a general concern for the medical literature for many years<sup>21,22</sup>. Thus, when selecting design and validation processes for diagnostic deep learning systems, one will have to focus both on the generalization challenges and on preventing ‘classical’ pitfalls in data analysis. We will, however, argue that both sets of challenges may be diminished by adopting certain fairly simple principles partly borrowed from the drug clinical trial field.

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

## External cohort evaluation

Rigorous performance evaluation is particularly important due to the inherent high complexity of deep neural networks, as seemingly well-performing deep learning

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

As deep learning systems are developed by learning correlations between input data and target outcome directly from the training data, it is essential that the training data adequately represent the target population<sup>31,190</sup>. Otherwise, the system might learn correlations exclusive to the subpopulation represented in the training data and, consequently, perform worse on those not represented in the training data to a sufficient extent. Despite this, systems are often trained on data sets with prominent biases in demographic characteristics such as sex, race or ethnicity, with the consequence that many systems exhibit substantial discriminatory biases<sup>37,191,192</sup>. Restricting the target population to a particular sex, race or ethnicity would be appalling, and the medical application of any such deep learning system would systematically increase health-care disparities. It is therefore pivotal to utilize truly representative and unbiased data for training deep learning systems in cancer diagnostics. This extends beyond ensuring representative distributions of relevant demographics in the training data set. Racial bias may also be encoded into systems if the target outcome used in the training is affected by histories of unequal treatment of patients based on race or ethnicity<sup>193</sup>, or is a proxy such as health-care cost instead of health needs, which has been shown to be the reason why a widely used health-care prediction algorithm exhibited significant racial bias<sup>194</sup>. Researchers should strive to identify and compensate for any such biases in their data sets, as failure to do so might reinforce health inequities if the deep learning systems are applied in clinical practice<sup>195,196</sup>. Deficient deep learning systems might be identified through rigorous evaluations in external data sets truly representative of the target population, or representative of minority populations, as well as through comprehensive analyses of system explainability across different demographic characteristics.

systems might utilize unintentional and possibly false features<sup>10–12</sup> and respond unexpectedly to apparently irrelevant changes of the input data<sup>23</sup>. Failure to properly evaluate systems might have far-reaching consequences, including misdirection of further research, diminished credibility of research findings and, most importantly, being worthless or even harmful to patients if used to influence treatment<sup>24,25</sup>.

**The importance of an external cohort evaluation.** As an initial evaluation step, the cohort used for development of a deep learning system is often partitioned randomly into three distinct subsets, hereunder referred to as ‘training’, ‘tuning’ and ‘test’, where the training subset is applied to learn candidate deep learning models, the tuning subset to select the deep learning system that appears to perform best and the test subset applied to evaluate the performance of the selected system<sup>8</sup>. The evaluation of the test subset may provide unbiased estimation of the performance in the development cohort. It may also provide some information on the system’s ability to perform well in other populations by considering the extent to which the system performs better on the training subset than on the test subset, as this indicates the level of overfitting to the training data. Systems that are highly overfitted to the training data are likely not to perform well on other populations as the noise utilized to improve the performance on the training subset may negatively influence the performance on other populations. However, even a system

that performs similarly in training and test subsets might perform far from acceptably on cohorts distinct from the development cohort<sup>26,27</sup>. As discussed below and in BOX 1, this may be caused by the system utilizing data features that correlate with the target outcome only in the development cohort, which could be viewed as overfitting to the entire development cohort, or might also be caused by important predictive features not being adequately represented in the development cohort. Thus, using a random subset of the development cohort for testing does not imply that the results have external validity, that is, the performance of the system observed in the test subset may not generalize to patients external to the development cohort.

For example, Zech et al.<sup>11</sup> investigated a deep learning system for detection of pneumonia on chest X-ray images and found that it was not able to uphold the high discrimination performance achieved in the development cohort when applied to cohorts from different institutions. In this case, there was a substantially higher disease prevalence in one of the training cohorts, and it appears that the poor generalization was in part caused by utilization of cohort-specific characteristics. In particular, the system utilized metallic tokens that radiology technicians placed on patients to indicate laterality, as these often appeared differently in different cohorts. The authors further point out that the system might not even generalize well to other patients from the same institution as the development cohort, because some correlations between input data and target outcome in the development

cohort may not be present in new cohorts from the same institution. Winkler et al.<sup>12</sup> found that, for their system, visible surgical skin markings present in the image were associated with a higher prediction score for melanoma. Similarly, Narla et al.<sup>10</sup> reported that the presence of a ruler beside a lesion in an image was associated with a higher malignancy score for skin cancer. Of course, neither skin markings nor rulers are causing the skin cancer, but the apparent correlation present in the development cohort is sufficient for the deep learning system to make use of these associations. It could be argued that more thorough quality control on the training data could mitigate this, but it is highly unlikely that one is able to detect and control for all potential confounding factors present in the training set.

Thus, unbiased performance estimation in a real-world application of a deep learning system requires external cohorts representative of a target population<sup>22,28–30</sup>. In an external validation, no information from the external cohort should have influenced the design of the system or the estimation of any model parameter. Additionally, the external cohorts will implicitly define the patient population for which we have estimated the performance of the system. Thus, to know whether or not the results may be generalized to the entire target population, we need a broad validation where the cohorts may be regarded as representative of this desired target population, for example with respect to age, sex, ethnicity, geographical differences and disease prevalence<sup>31,32</sup>. Other types of evaluations may also be warranted prior to introducing the system in medical practice, including so-called domain validation to evaluate whether the system performs consistently across a range of laboratories and technical equipment (BOX 2).

Objective, non-random separation of patients from the same hospital or subjects from the same country — for example, distinguishing between patients treated before and after a certain date — allows using one cohort for training and tuning, and another for what has been denoted ‘narrow validation’<sup>22</sup> (BOX 2). Such evaluation might provide unbiased performance estimation for a particular hospital. However, the two cohorts should not simply be a non-random separation of an originally larger cohort but, instead, be processed separately when acquiring data and ascertaining target output<sup>33</sup>. Narrow validation is sometimes considered a limited type of external validation<sup>22</sup>.

**Prevalence in recent studies.** In order to investigate the prevalence of external cohort evaluation and other characteristics of recent studies on deep learning and cancer diagnostics, we searched PubMed on 21 April 2020 for original research articles published in 2015 or later (Supplementary Methods). The search provided 3,578 results, and the number of publications roughly doubled each year since 2016. To explore the use of external cohort evaluation and other characteristics in some of the most prominent and perhaps best studies, we restricted our evaluation to those with at least 20 citations per year or published in a journal with an impact factor of 10 or larger. Although studies satisfying either of these criteria are presumably quite influential, we acknowledge that some of the other studies might be equally good. In particular, recent studies may not have had time to accrue 20 citations even if they are currently of great interest, and such studies would only be included if published in a journal with an impact factor of 10 or larger. This will exclude most studies published in new journals that are expected to receive impact factors of 10 or larger when these become available. However, we consider the selected papers to be sufficient for the purposes of this discussion, as they show that some aspects of study design could be better even in some of the presumably best studies. Only 257 (7%) of the 3,578 search results satisfied at least 1 of these selection criteria, and another 43 search results were excluded because the document type in Web of Science indicated that these were not original research articles. The remaining 214 studies were manually evaluated (Supplementary Table 1). We further excluded 6 studies that were not original research articles and 102 studies where deep learning was not used to predict or classify features relevant for cancer diagnosis, prognosis or treatment response, or such potential utility of the deep learning system was not evaluated. After also excluding 14 studies without human subjects or only pertaining to cell biology, we ended up with 92 eligible studies<sup>34–125</sup>, of which 85 (92%) used images as input to the deep learning system<sup>34–57,59–64,66,67,69–93,95–99,101–121,123,125</sup>.

Among 516 original research articles on artificial intelligence for diagnostic analysis of medical images published in 2018, Kim et al.<sup>7</sup> found only 31 studies (6%) that evaluated an external cohort. By contrast, 50 (54%) of our 92 eligible studies evaluated the performance of the deep learning system on an external

cohort<sup>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</sup>.

This discrepancy is most likely mainly attributed to our selection of presumably influential studies and partly attributed to the increasing usage of external cohorts (FIG. 1a); 34 (72%) of the 47 eligible studies published in 2019 and 2020 evaluated an external cohort compared with 9 (39%) of the 23 eligible studies published in 2018 and 7 (32%) of the 22 eligible studies published before 2018.

Among studies satisfying both of our selection criteria, 79% (11 of 14) evaluated an external cohort, compared with 68% (25 of 37) for studies that satisfied only the impact factor criterion and 34% (14 of 41) for studies that satisfied only the citation frequency criterion. It thus appears that journals with a high impact factor have a preference for studies evaluating external cohorts. This is consistent with the call by editors of leading scientific journals for rigorous evaluation of artificial intelligence tools<sup>126,127</sup> and explicit prioritization of biomarker studies that evaluate external cohorts by some journals, for example *Journal of Clinical Oncology*.

## Generalizability

Although increased use of external cohorts is an important step towards proper validation of deep learning systems, one is still left with the challenge of ensuring that the results obtained for such a population provide a satisfactory measure of the performance within the entire intended target population. This target population may typically be patients who have a specific cancer type, and although often restricted, for example, to certain stages of the disease, the target population is normally broad. Although some studies may use more than one external cohort and some use trials with many centres distributed over several countries, it is difficult to obtain external cohorts that entirely cover the target population. Thus, successful application of a deep learning system will depend on good generalization properties, so that good performance on one population also indicates satisfactory performance on populations differing with respect to some properties. Fortunately, exploring generalization in deep learning is an active research area<sup>128</sup>, and by utilizing certain design principles, deep learning systems have shown remarkably good generalization performance on numerous tasks<sup>3–5</sup>.

### Box 2 | Approaches for evaluating a deep learning system

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

#### Internal validation

Internal validation is the evaluation of a deep learning system's performance in the development cohort. This can be done by evaluating the performance in a randomly sampled subset of the development cohort disjoint from the training and tuning subsets, or by using resampling techniques such as cross-validation or bootstrapping<sup>22</sup>.

#### Narrow validation

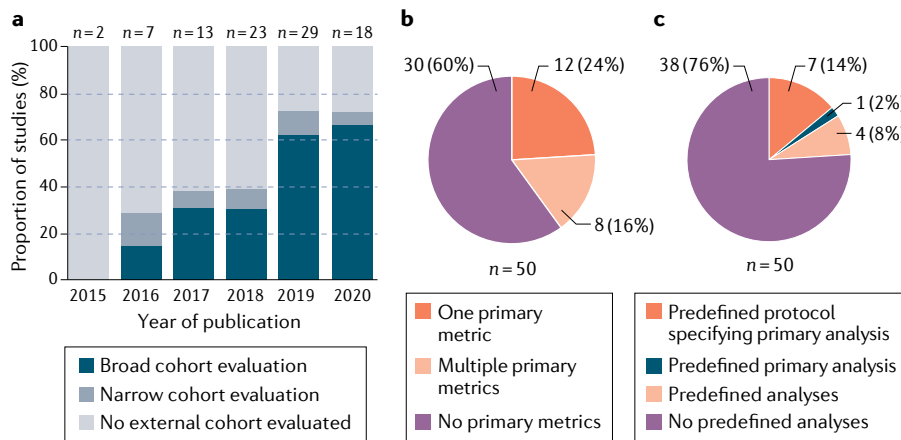
Narrow validation is the evaluation of a deep learning system's performance based on a cohort that is similar to but differs non-randomly from the development cohort, for example on a cohort from the same hospital as the development cohort but sampled in a time interval disjoint from the time interval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

#### Broad validation

Broad validation is the evaluation of a deep learning system's performance based on a cohort geographically separate from the development cohort, for example from a different hospital or country<sup>22</sup>. No information from the broad cohort should have influenced the development of the system.

#### Domain validation

Domain validation is the evaluation of a deep learning system's performance in a setting that is very different from the one in which the system was developed<sup>197</sup>. This includes validation in a cohort with characteristics not represented by the development cohort, for example developing a method on one type and stage of cancer and validating it on another type or stage of cancer. Other examples are when the validation data are obtained by equipment not used in the development, such as imaging systems from different vendors, or by sample preparation procedures intentionally different from those used for the development cohort. Domain validations should also be narrow or broad validations, and are typically performed after successful narrow or broad validations.



**Fig. 1 | Characteristics of recent, presumably influential, deep learning studies in cancer diagnostics.** **a** | Percentage of studies reporting on the evaluation of a broad or narrow cohort (BOX 2) by year of publication, for all 92 eligible studies. **b** | Percentage of studies specifying 1, multiple or no primary performance metrics in the analysis of the external cohort, for the 50 eligible studies that reported on the evaluation of an external cohort. **c** | Percentage of studies specifying a predefined analysis of the external cohort, for the 50 eligible studies that reported on the evaluation of an external cohort. Studies that specified predefined analyses of external cohorts without defining which was the primary, if any, were categorized as ‘predefined analyses’. Studies with a predefined primary analysis were categorized according to whether the primary analysis was pre-specified in a protocol or not.

One way of increasing generalization is to control the neural network’s capacity to express complex mappings, for example by limiting the number of adjustable parameters in the network, imposing various constraints on the network or regularizing the optimization<sup>129,130</sup>. Transfer learning could also increase generalization, particularly when training data for the task at hand are scarce<sup>131,132</sup>. In transfer learning, the network is initialized with parameters optimized using data for a different task, typically using large data sets such as ImageNet<sup>133,134</sup>, which may mitigate overfitting at the possible cost of introducing biases<sup>135–137</sup>. Making the training data set more diverse and more representative of the target population is another way of increasing generalization<sup>138</sup>. Of particular importance is to ensure adequate and unbiased representation across demographic characteristics such as sex, race and ethnicity (BOX 1). In addition to expanding the natural training data set, that is, the set of training data acquired from a range of patient samples with associated target outcome, one may artificially augment the training data set by applying smaller transformations on the inputs while maintaining their relationship to the target output<sup>18,139</sup>. This can reduce the network’s ability to memorize details of the training data and thereby increase generalization, especially in situations where the availability of training data is limited. The transforms can randomly change, often called ‘distort’, the input data by, for example, adding noise, erasing parts,

shifting and scaling colours or altering the image geometry<sup>19</sup>. Artificially diversifying the training data may increase generalization by enabling the resulting system to ignore vagaries of the measurement process and even become applicable to multiple data acquisition procedures, for example different acquisition equipment<sup>140,141</sup>. Other augmentation techniques include those that generate artificial input data, for example by mixing multiple data inputs<sup>19</sup>. The value

of augmentation techniques has been observed in various application domains<sup>19</sup>, including use on images obtained from radiology<sup>38,142–144</sup> and histopathology<sup>141,145</sup>.

To illustrate the importance of the amount of and variation in training data, and more specifically show how data distortion may work to improve deep learning systems in cancer diagnostics, we show this type of analysis here using data from a previously published study<sup>113</sup>. This previous study applied deep learning to predict colorectal cancer-specific survival directly from conventional haematoxylin and eosin-stained sections, with training and tuning data derived from 2,473 patients from four cohorts. The performance was evaluated on an external cohort consisting of 1,122 patients from a randomized controlled trial of a drug that was observed to not affect survival<sup>146</sup>. We applied the convolutional neural network called Inception-v3 (REF.<sup>147</sup>), which is a network commonly used in medical image diagnostics<sup>8</sup>, in both the previously published analyses and the new analyses presented here.

Initially, we applied the same distortion process as in our published analyses<sup>113</sup>. This process artificially increased the variation of the training images by randomly distorting their colours, which is an augmentation technique that appears crucial when training deep learning systems in histopathology<sup>145</sup>. Initially, the maximum amount of distortion we allowed was quite modest (FIG. 2a). To illustrate the effect of reducing the

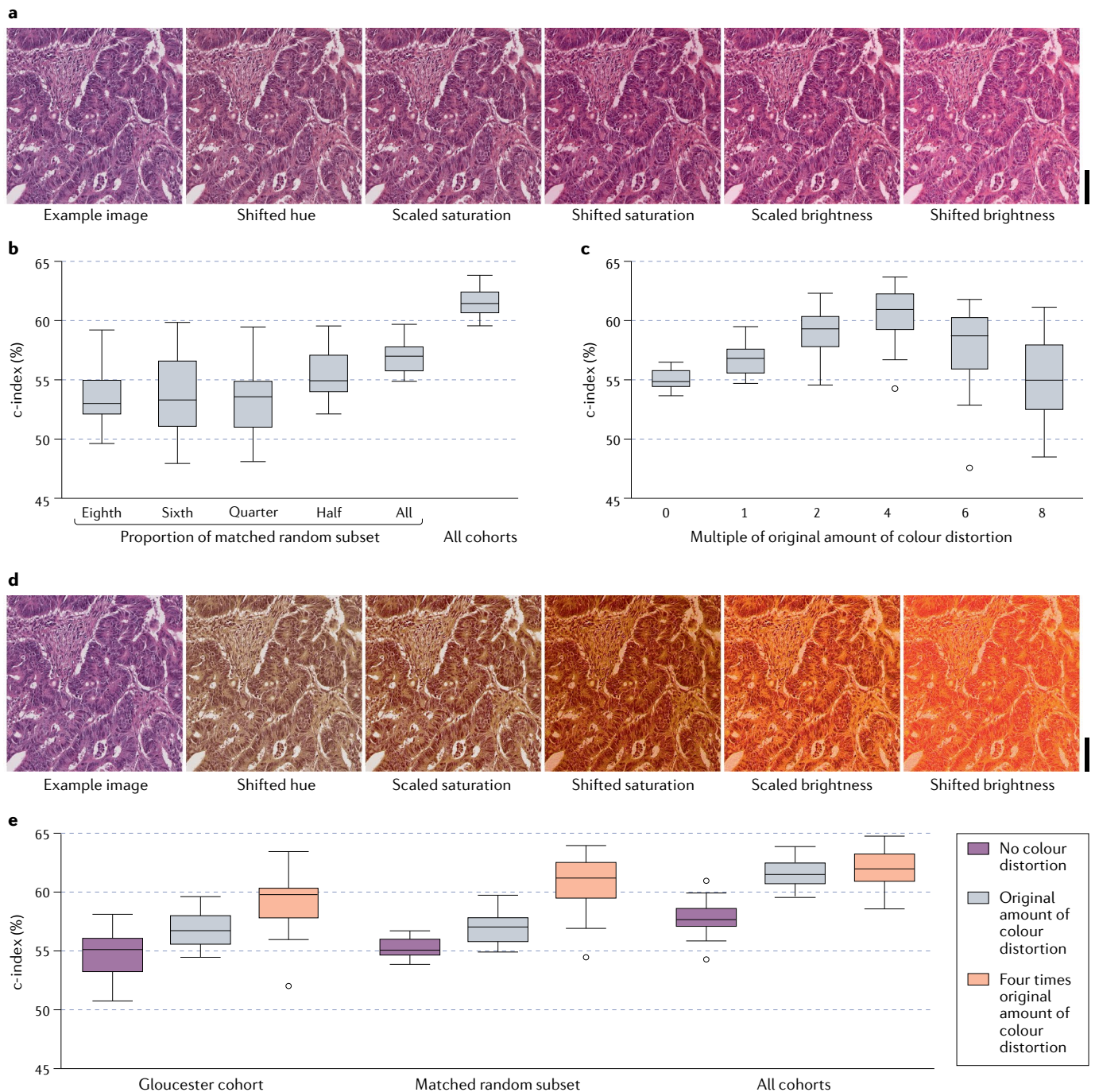
**Fig. 2 | Effect of data variation when training deep learning systems.** For each analysis set-up, 20 individual deep learning systems were trained and tuned for prediction of colorectal cancer-specific survival using images of haematoxylin and eosin-stained sections acquired by both Aperio AT2 (Leica Biosystems, Germany) and NanoZoomer XR (Hamamatsu Photonics, Japan), as in the previously published analyses<sup>113</sup>. The individual systems were applied to evaluate the external cohort using NanoZoomer XR slide images, and the concordance index (c-index) of the system’s binary output was computed. Standard box plots were made using Stata/SE 16.1 (StataCorp, USA). The matched random subset contained the same number of training and tuning patients with and without cancer-specific death as in the Gloucester cohort, in total 979 patients. **a** | An example image from the training data set and the results of applying the maximum possible amount of colour distortion at each step in the random distortion process used in the published Inception-v3 analyses<sup>113</sup>. Generally, the distortion process applies random colour distortions to an image by converting the image to HSV colour space, adding a random value between –0.05 and 0.05 to the hue, scaling the saturation by a random value between 1/1.1 and 1.1, adding a random value between –0.1 and 0.1 to the saturation, scaling the brightness (or, technically, the value channel in the HSV colour space) by a random value between 1/1.1 and 1.1, adding a random value between –0.1 and 0.1 to the brightness and converting back to RGB colour space. Intuitively, the leftmost and rightmost images represent the range of the random colour distortion, that is, the minimum and maximum possible amount of colour distortion for the applied distortion process, where the minimum is no colour distortion. Scale bar, 100 µm. **b** | Effect of changing the number of patients in the training and tuning subsets when using the original amount of colour distortion, as depicted in part **a**. **c** | Effect of changing the amount of colour distortion when training and tuning using the matched random subset. Label ‘0’ on the horizontal axis identifies deep learning systems trained without any colour distortion, label ‘1’ identifies systems trained with the colour distortion process depicted in part **a** and label ‘4’ identifies systems trained with the colour distortion process depicted in part **d**. **d** | Similar to part **a**, but four times the amount of colour distortion was used at each step in the distortion process. Scale bar, 100 µm. **e** | Effect of changing the amount of colour distortion and the number of patients and cohorts in the training and tuning subsets.

number of patients while keeping the patient heterogeneity implied by having data from four cohorts, we randomly sampled 979 patients in such a manner that the data had the same number of training and tuning patients with and without cancer-specific death as in the cohort from the Gloucester Colorectal Cancer Study, UK (the largest of the four training and tuning cohorts). The decreased performance of the resulting deep learning system when evaluated on the external cohort (FIG. 2b) exemplifies the importance of a large natural training data set and its intrinsic variation<sup>138</sup>.

Further reduction of the number of patients decreased the performance further; training and tuning on a quarter of the 979 patients or fewer (that is, fewer than 250 patients) provided systems that did not perform substantially better than random guessing (FIG. 2b).

We then showed that modifying the distortion process may mitigate for the performance loss observed when reducing the number of patients in training and tuning. Compared with using all 2,473 patients for training and tuning, using 979 randomly selected patients and

four times the original amount of colour distortion provided similar performance on the external cohort (FIG. 2c). For this modified distortion process we allowed quite substantial colour distortions (FIG. 2d), and the results showed that artificial augmentation may, in some cases, compensate for limited natural training and tuning data. However, increasing the amount of colour distortion further provided worse performance (FIG. 2c), illustrating the trade-off between preventing overfitting through random distortions and occluding relevant information for the prediction task.



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

In general, the most suitable distortion process will depend on the particular medical prediction task because the involved data will tolerate different amounts of the various types of distortions before true correlations between input and target output are occluded. For instance, deep learning systems that classify based on images of skin lesions or tumour sections are likely to benefit from being invariant to rotations, whereas systems aimed at supporting radiology might rely on the orientation in images of larger organ structures and, thereby, perform worse if forced to be rotation invariant. Thus, the distortion process needs to be fine-tuned to the particular application, as findings about which distortion process appears most beneficial in one scenario — for example, findings from the example presented in FIG. 2 — are not necessarily directly applicable to other scenarios. However, the general principle is that including much and varied training data is important. As the importance of artificial augmentation decreases with the amount and diversity in the natural training data, prediction tasks where the true correlations between input data and target output are easily obscured by distortion warrants a more comprehensive natural training data set.

### Predefined primary analysis

In the development of a deep learning system, researchers will often evaluate different systems sequentially, each time having the possibility to learn from interpreting the previous evaluations and adapt the system to the specific data used for evaluation. Such repeated evaluations will bias the estimates, and their dependence on previous evaluations makes established statistical approaches for adjusting for multiple comparisons not applicable<sup>148,149</sup>. Similar reanalysis issues may arise if the

initial analysis of a specific deep learning system reveals issues that are then corrected and the performance is re-evaluated. Such problems of repeated or multiple evaluations are well-known from examinations of the data analysis in various types of published medical studies, and have been identified as important contributors to biased inference and irreproducible results<sup>20,150</sup>.

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

### Prevalence of predefined primary analysis.

In our evaluation of recent, presumably influential, deep learning studies in cancer diagnostics, all studies performed multiple analyses of the external cohort typically in the form of evaluating multiple systems, analysing multiple subpopulations or using various analysis methods. Only 3 (6%) of the 50 eligible studies that evaluated an external cohort used one of the well-established methods for adjustment for multiple comparisons<sup>51,62,114</sup>, for example Bonferroni correction. This implies that most studies should have specified which analysis was considered the primary analysis prior to evaluation of the external cohort, if such a decision was made, in order to inform the reader which analysis was not affected by selection bias and to help distinguish studies with a predefined primary analysis from those that repeatedly evaluated the external cohort and might have ended up reporting severely biased performance estimates. Although the principle of using an external data set only once to evaluate the final hypothesis should be well-known in the machine learning community<sup>151,152</sup>, it seems, currently, that there is no tradition for specifying the predefined primary analysis in deep learning publications other than those reporting on clinical trials. In our evaluation, 20 (40%) of the 50 studies evaluating an external cohort specified one or more primary performance metrics<sup>55,60,73,82,83,85,86,93,98,102,105,108–110,113,115,116,120,121,125</sup> (FIG. 1b), but only 8 (16%) of the 50 studies specified a predefined primary analysis<sup>73,83,102,105,109,113,120,121</sup> (FIG. 1c).

Pre-specification of the primary analysis has previously been advocated in diagnostic and prognostic research<sup>153,154</sup>, but this is

unfortunately still not common practise despite being the only direct protection against selection bias<sup>20</sup>. To ensure unbiased estimation, the primary analysis should be unequivocally specified prior to all investigations that could reveal correlations between input data and target output in the external cohort. This would require the researchers to define all relevant aspects of the validation prior to analysing the cohort, including the deep learning system, target output, and patient and input data in the external cohort. Predefining the primary analysis will entail a commitment to the main analysis, which implies that the analysis should be carefully planned in advance and that researchers will be discouraged from performing creative data dredging<sup>155</sup>.

**Choosing the primary metric.** Many medical questions are categorical in nature, for example whether tumour or not, whether mutated or not and whether to offer treatment or not. However, deep learning models often output continuous values reflecting the predicted probability of each possible outcome. In such cases, the predefined primary analysis should preferably evaluate a categorization of the model output aimed at answering the medical question. The primary analysis will then be comparing predicted and target outcome in the external cohort, for example by measuring the so-called balanced accuracy<sup>156</sup>. Measuring the performance using categorical outputs often provides more conservative estimates<sup>157</sup> and avoids issues with metrics frequently applied to measure the performance using continuous outputs. For instance, the area under the receiver operating characteristic curve (AUC)<sup>158</sup> and the concordance index (c-index)<sup>159</sup> are only affected by the ranking of the continuous outputs, not the prediction scores themselves<sup>160</sup>. Thus, such metrics may indicate that a deep learning system performs well even if it predicts markedly too high probabilities for all patients in a specific cohort, provided that the continuous outputs of the system rank the patients in a fairly correct order. In another cohort, the same system may similarly appear to perform well even if it predicts markedly too low probabilities for all of those patients. The generalizability of such a system is poor, yet this would not be evident from the AUC and c-index of the continuous outputs but would be evident from the AUC and c-index of a categorization defined irrespective of the external cohorts. The categorization may be defined by, for example, determining suitable

thresholds during tuning or selecting the outcome with the highest prediction score as the predicted outcome. Defining the categorization using the external cohort, even at predefined levels of, for example, sensitivity, adapts the categorical marker to the specific external cohort and may occlude shifts in the prediction scores as with the AUC and c-index of the continuous outputs.

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

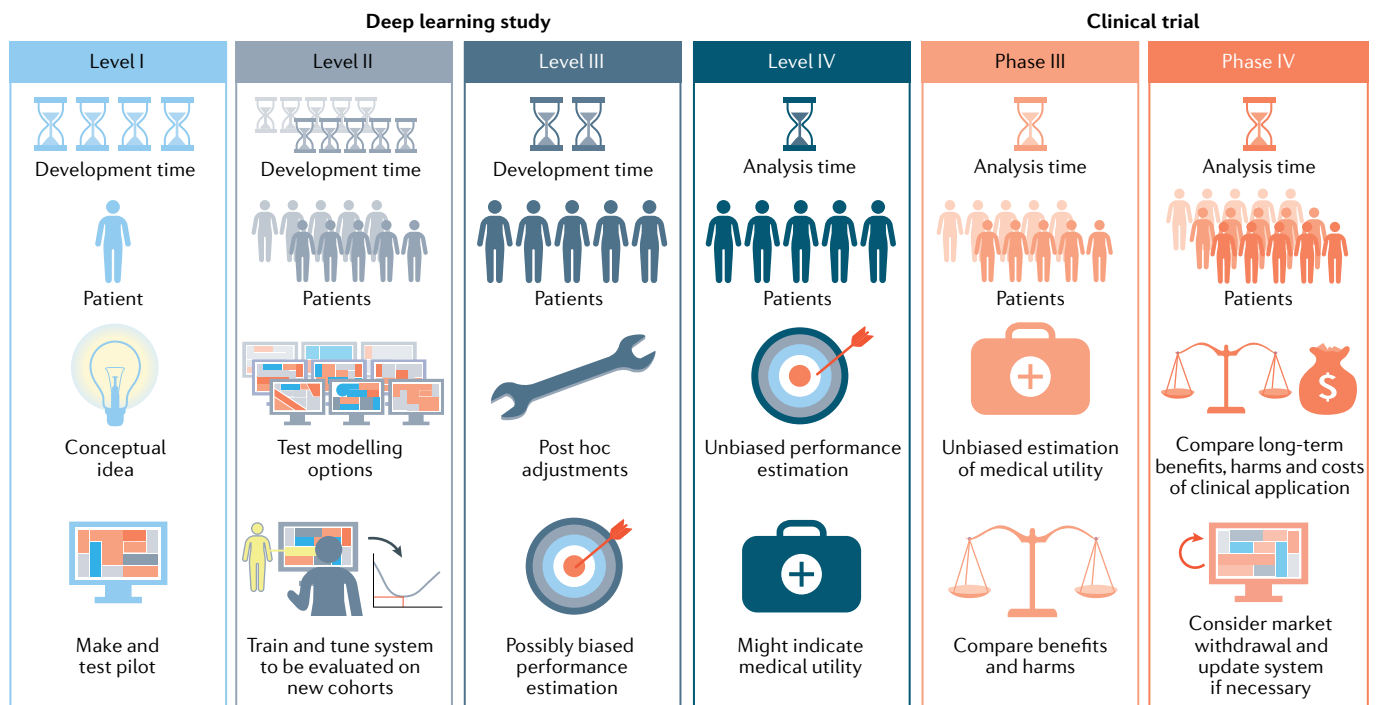
For certain deep learning systems, the intended medical application directly utilizes the system's continuous output, for example to triage patients for further examinations, and in such cases the continuous output should be evaluated in the primary analysis. This may warrant additional analyses to reveal generalization issues that might be occluded by the selected performance metric, for example to consider a calibration plot in addition to the c-index when evaluating a clinical decision support system for predicting patient outcome<sup>22,26</sup>. In general, the metric chosen for the primary analysis should be one that measures how well the deep learning system performs in the intended medical application. For instance, the overall performance in a classification task could be measured using the balanced accuracy.

### From conception to application

All research with the potential to influence patient treatment should undergo careful evaluation sequences and be driven by protocols with a predefined statistical

analysis plan<sup>153</sup>. FIGURE 3 illustrates what we consider natural and important steps in the development and evaluation of deep learning systems for medical applications.

The initial exploratory studies aim to answer whether deep learning appears suitable for the task at hand or whether further investigations based on deep learning are not warranted at this time, usually because the hypothesis seems ill-founded or the available data are not expected to provide a system with adequate performance. The performance estimates obtained in such pilot studies are frequently inflated by the use of a limited development cohort, but promising findings may motivate further investigations. After a series of explorations, and possibly expansions, of the development cohort, the development should conclude by deciding which system appears to perform best on the intended medical task, considering also the sensitivity to vagaries of the measurement process. Of particular importance to prevent selection of a system that performs much worse on patients outside the development



**Fig. 3 | Development and evaluation of deep learning systems.** A deep learning project often begins with testing a conceptual idea using pilot software based on a related open source implementation and data easily available to the researchers. Successful level I studies will typically evolve into explorative testing of different modelling options that might be more suitable for the particular task. The system that appears to perform best should be determined in a level II study that includes sufficient amount and variation in the natural training data set. Although performance estimates obtained in such studies are often inflated by the use of a subset that closely resembles the training subset, level II is an important step in the evaluation sequence that could motivate investigators to pursue evaluation of external cohorts and attract collaborators.

As the lack of predefined primary analysis often entails post hoc adjustments influenced by the performance in the external cohort, we distinguish between studies without (level III) and with (level IV) a primary analysis unequivocally specified prior to all investigations that could reveal correlations between input data and target output in the external cohort. If the medical validity of a deep learning system is established in level IV studies, the indicated medical utility should be prospectively evaluated in randomized phase III clinical trials where the system directly intervenes with the current standard of care. If medical utility is demonstrated and necessary governmental agencies approve routine medical application, the system can be applied in medical practice while monitoring the long-term benefits, harms and costs of its application.

cohort, the study could include a sufficient amount and variation in the natural training data set and use techniques such as data distortion to increase the variation artificially.

There is growing interest in explainable deep learning systems<sup>161–163</sup>, including the creation of inherently more explainable systems and post hoc explanations of existing systems<sup>164</sup>. For image classification tasks in particular, so-called saliency maps visualize the contribution of each pixel to the final prediction score and can be created using numerous different techniques<sup>165–167</sup>. By increasing the transparency, the more explained systems might have more predictable generalizing abilities. This may be used to identify target populations within which the system is expected to generalize well or settings where the system is prone to fail. For example, Winkler et al.<sup>12</sup> used such a technique to support their finding that surgical skin markings unduly increased the system's prediction score for melanoma. Although current explainability techniques might suggest generalizability, and thereby suggest suitable target populations or influence the selection of which system to evaluate further, they will only provide indications and, thus, not reduce the need for proper validation.

Whereas efficacy studies of pharmaceutical products are usually preceded by prospective trials to estimate basic features such as safety and dosing<sup>168</sup>, deep learning systems for diagnostic purposes can to a larger extent utilize retrospective cohorts, for example from earlier clinical trials or medical practice. Given the risks, time frame and costs of interventional research<sup>168–170</sup>, we recommend rigorous, retrospective analyses to evaluate the medical validity of a deep learning system by conducting an external validation

according to a predefined primary analysis. The results of such studies provide valuable information to direct further research, thus warranting publication regardless of the significance of the findings, which would also mitigate publication bias.

Rigorous, retrospective analyses of a deep learning system might warrant conducting a prospective, randomized phase III clinical trial where the system directly intervenes with the current standard of care in order to evaluate the system's medical utility in a specific real-world application, considering both benefits and harms for patients in the target population<sup>30,171</sup>. Systems demonstrated to have medical utility and approved by necessary governmental agencies can be applied in medical practice while monitoring the long-term benefits, harms and costs for each specific real-world medical application in phase IV clinical trials. Such surveillance might eventually indicate that the system needs to be updated because of changes in medical practice or data acquisition<sup>172</sup>.

The levels of deep learning studies depicted in FIG. 3 and the phases of clinical trials were used to categorize recent, presumably influential, deep learning studies in cancer diagnostics in relation to the reliability of the performance estimation approach and the demonstrated applicability of the system in medical practice. Although some group sizes are very small, there appear to be notable differences between research fields defined by the input to the deep learning system (FIG. 4). The proportion of studies evaluating an external cohort was lowest for the 7 studies with only non-image inputs, such as omics data (29%; 2 of 7 studies), and highest for the 22 studies with images other than histopathology section and radiology images as input, for example from gastrointestinal endoscopic

examinations or dermoscopic images (64%; 14 of 22 studies). Five (23%) of the 22 studies with other images as input even had a predefined primary analysis of the external cohort<sup>73,102,105,109,121</sup>, which included the 3 studies reporting on a randomized clinical trial, all of which evaluated a deep learning system to aid gastrointestinal examinations<sup>102,105,121</sup>.

**Recommended protocol items.** When planning to evaluate the medical validity of a deep learning system through rigorous, retrospective analyses, we recommend the unequivocal specification of the predefined primary analysis to be documented in a study protocol. Relevant items in such protocols would differ from clinical trial protocols, which are the target of guidelines such as SPIRIT (Standard Protocol Items: Recommendations for Interventional Trials)<sup>173</sup> and its extension to artificial intelligence<sup>174</sup>. Protocols should be developed before conducting the validation, and relevant items would therefore also differ from those in original research articles, which are the target of many reporting guidelines such as CONSORT (Consolidated Standards of Reporting Trials)<sup>175</sup> and TRIPOD (Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis)<sup>22</sup> as well as their extension or anticipated adaption to machine learning<sup>176,177</sup>. There is therefore a need to establish guidelines dedicated to study protocols describing validations of deep learning systems. We propose a non-exhaustive list of items that we consider essential in such protocols, termed Protocol Items for External Cohort Evaluation of a deep learning System (PIECES) in cancer diagnostics.

In order to be sufficiently concrete about the predefined primary analysis,

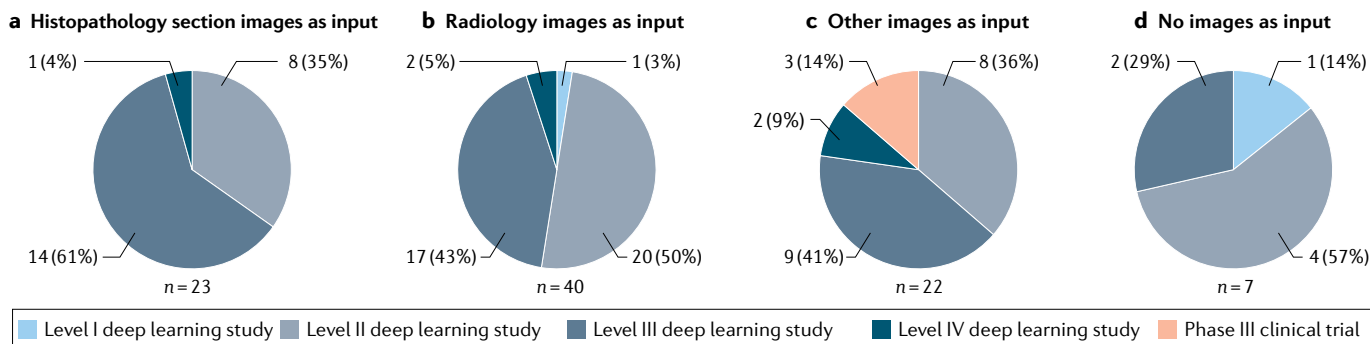


Fig. 4 | Reliability of performance estimations in recent, presumably influential, deep learning studies in cancer diagnostics. Percentage of studies categorized in the different levels of deep learning studies or phases of clinical trials depicted in FIG. 3 for all 92 eligible studies separated by type of input to the neural network. The input was histopathology section images in 23 (25%) of the studies (part a), radiology images in 40 (43%) of the studies (part b), other images in 22 (24%) of the studies (part c) and other types of input in 7 (8%) of the studies (part d).



the protocol needs to describe the deep learning system and how it will be assayed; define the external cohort, including its origin, what it represents in terms of medical setting and target population, and input data and target output; and clearly specify the performance evaluation. These three parts of the protocol form the basis of our PIECES recommendations together with a declaration of status (BOX 3). The status declaration should scrupulously elucidate any investigations performed before finalising the protocol that could

reveal correlations between input data and target output in the external cohort, or state that no such investigations were performed.

The PIECES recommendations are designed to facilitate identification of ambiguities and disagreements between the researchers planning to conduct an external validation as well as to provide a clear description of the predefined primary analysis as a reference for all readers, which may aid medical professionals in identifying well-designed studies and their

applicability to their own clinical practice. The thought and work that should go into making such a protocol could also allow the researchers to make appropriate changes prior to performing the external validation. For instance, considering what the external cohort is intended to represent and how the deep learning system is envisioned to be applied in practice could affect the inclusion and exclusion criteria for patients and samples as well as the metric or statistical test applied in the primary analysis.

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

**Study registration.** We recommend registration of the study protocol in an online repository before analysing the external cohort. Most major trial registries, for example [ClinicalTrials.gov](https://www.clinicaltrials.gov) and the [International Standard Randomized Controlled Trial Number \(ISRCTN\) registry](https://www.isrctn.com/), accept registration of diagnostic accuracy studies<sup>154</sup>. These registries can be used to record external validation studies in deep learning, but some items will not be relevant and some important items, such as defining the deep learning system, will not be encouraged. A dedicated repository to register the study protocol describing the external validation of a deep learning system is therefore warranted. We recognize that it may be undesirable to publish a detailed study protocol in an online repository prior to conclusion of the study as this would reveal novel work prior to publication of the results and perhaps in some rare cases

## Glossary

### Area under the receiver operating characteristic curve

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

### Artificial neural networks

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

### Balanced accuracy

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

### Capacity

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

### Concordance index

(c-index). A performance metric measuring the concordance between a target outcome, usually defined by time to event data, and the ranking of subjects provided by a continuous or categorical marker. A c-index of 50% indicates random guessing and 100% indicates perfect prediction. For dichotomous outcomes, the c-index and the area under the receiver operating characteristic curve are equivalent.

### Deep learning

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

### Deep learning models

Computational models obtained by training deep neural networks. Note that a single training of a neural network produces a sequence of models as each new optimization iteration produces a model slightly different from the previous one. A tuning data set may be used to select among these models.

### Deep learning systems

Systems utilizing one or more deep learning models to make predictions. A system's output may be a function of the outputs of the models, for example by averaging and thresholding the model outputs.

### Development cohort

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

### External cohorts

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

### External validation

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

### Generalizability

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

### Overfitting

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

### Supervised machine learning

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

### Test

Although frequently used by the machine learning community to refer to an evaluation of a system's performance, we use 'test' to refer to evaluations other than external validations, for example internal validations.

### Training

Optimization of model parameters based on data.

### Tuning

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

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

**Status**

- Specify the date the protocol was last modified.
- Scrupulously elucidate any investigations performed before finalising the protocol that could reveal correlations between input data and target output in the external cohort, or state that no such investigations were performed.

**System**

- Describe the development of the deep learning system, including utilized cohorts, network architecture, hyperparameters and any categorization of the neural network model's output.
- Unequivocally specify how to assay the deep learning system in a blinded fashion for a single, new subject, including what the system receives as input and what it directly outputs.

**External cohort**

- Describe the origin of the cohort, and explain why it should be regarded as external to the development cohort.
- Precisely define criteria for inclusion and exclusion of subjects and samples, preferably starting from a consecutive series of subjects.
- Clearly state the medical setting and target population that the cohort represents.
- Specify the acquisition of input data, including whether it was acquired blinded to the deep learning system and target output. Note the expertise of any humans involved in the process, for example that a pathologist annotated the regions of interest in slide images.

- Specify the ascertainment of target output, including whether it was ascertained blinded to the deep learning system.
- If multiple external cohorts are planned to be analysed as a pooled cohort, then the preceding five protocol items should be completed for the pooled cohort and differences between the individual cohorts should be stated. If multiple external cohorts are to be analysed independently, the five preceding protocol items should be completed for each cohort, as well as subsequent protocol items if the predefined analyses differ between cohorts.

**Analyses**

- Unequivocally specify the primary analysis, including the target output and the performance metric and/or statistical test with interpretation.
- If the chosen metric or statistical test depends on other markers, describe how these markers were assayed and whether done blinded to the deep learning system and target output, and specify how missing values will be handled.
- If the deep learning system was designed to evolve upon usage, for example by learning from unlabelled data or adapting to a cohort, specify that this will not be done when evaluating the external cohort. The system's prediction should thus not depend on the order in which a set of patients is evaluated and should also be identical if the same patient is evaluated multiple times.
- If additional analyses will be performed and reported in disseminations, for example of other deep learning systems, target outputs, metrics or statistical tests, or in specific patient subgroups, specify these analyses in the same manner as the primary analysis and identify them as secondary analyses.

jeopardise publication. In a dedicated repository, a submission could be partially or completely invisible to the public and the protocol encrypted until the authors choose to reveal the submission and provide the required decryption key, thus facilitating preregistration of study protocols without requiring authors to reveal novel ideas prematurely.

Registration of observational studies has been advocated by editors of major clinical journals<sup>178,179</sup>, many editorial board members<sup>180</sup> and researchers<sup>181,182</sup>, and the criticism this has received from epidemiologists in relation to the exploratory nature of epidemiology<sup>183–185</sup> does not apply to external validation studies. For diagnostic and prognostic biomarker studies in particular, the registration of a study protocol with a predefined analysis plan has been recommended by several researchers<sup>153,154,186–188</sup>, provided that it precedes the onset of the study<sup>189</sup>. This would facilitate a more balanced evaluation of the proposed marker, identification and prevention of selective reporting, increased transparency, reduced proportion of false positive findings, mitigation of publication bias through identification of unpublished studies, and prevention of unnecessary duplication of research while facilitating collaboration between researchers and identification of research gaps. Consequently, widespread preregistration

of detailed study protocols for deep learning systems might translate into more rapid identification of promising systems and thereby expedite progression of the research field. It would also communicate a study to peers without disclosing the findings and interpretations prior to editorial and peer review, thus providing some of the benefits of preprint archiving while allowing critical appraisal of the findings and interpretations before publication.

Amendments of clinical trial protocols are common but should be tracked and dated<sup>173</sup>. Whereas clinical trials often take years to conduct due to patient recruitment and follow-up, most external validations of deep learning systems use retrospective data, and the analysis part of the validation may be performed in a matter of days. Consequently, it should rarely be necessary to modify the study protocol describing the external validation of a deep learning system after initiating the validation. We therefore generally discourage protocol amendments, but if found necessary for a particular study, we recommend amendments to be included as postscripts to the study protocol, leaving the original protocol unaltered. Both the postscript and disseminations of the validation results should concretely specify what was changed as well as describe the motivation and rationale for the change.

**Conclusions**

Including much natural and artificial data variation when training rigorous deep learning systems appears pivotal, as analyses indicate its instrumental role in increasing the performance and generalizability of systems. Utilizing multiple sets of patients, samples and data acquisition procedures will diversify the training data, whereas augmentation techniques artificially enhance the variation further. The resulting systems may be capable of handling the diversity in routine medical practice and, in some cases, even generalize to completely new settings.

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

Andreas Kleppe<sup>1,2</sup>, Ole-Johan Skrede<sup>1,2</sup>,  
Sepp De Raedt<sup>1,2</sup>, Knut Liestøl<sup>1,2</sup>, David J. Kerr<sup>1,2,3</sup>  
and Håvard E. Danielsen<sup>1,2,3</sup>✉

<sup>1</sup>Institute for Cancer Genetics and Informatics,  
Oslo University Hospital, Oslo, Norway.

<sup>2</sup>Department of Informatics, University of Oslo,  
Oslo, Norway.

<sup>3</sup>Nuffield Division of Clinical Laboratory Sciences,  
University of Oxford, Oxford, UK.

✉e-mail: hdaniels@ifi.uio.no

<https://doi.org/10.1038/s41568-020-00327-9>

Published online 29 January 2021

- Schmidhuber, J. Deep learning in neural networks: an overview. *Neural Netw.* **61**, 85–117 (2015).
- LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436–444 (2015).
- Hosny, A., Parmar, C., Quackenbush, J., Schwartz, L. H. & Aerts, H. J. W. L. Artificial intelligence in radiology. *Nat. Rev. Cancer* **18**, 500–510 (2018).
- Vamathevan, J. et al. Applications of machine learning in drug discovery and development. *Nat. Rev. Drug Discov.* **18**, 463–477 (2019).
- Bera, K., Schalper, K. A., Rimm, D. L., Velcheti, V. & Madabhushi, A. Artificial intelligence in digital pathology — new tools for diagnosis and precision oncology. *Nat. Rev. Clin. Oncol.* **16**, 703–715 (2019).
- Nagendran, M. et al. Artificial intelligence versus clinicians: systematic review of design, reporting standards, and claims of deep learning studies. *BMJ* **368**, m689 (2020).
- Kim, D. W., Jang, H. Y., Kim, K. W., Shin, Y. & Park, S. H. Design characteristics of studies reporting the performance of artificial intelligence algorithms for diagnostic analysis of medical images: results from recently published papers. *Korean J. Radiol.* **20**, 405–410 (2019).
- Liu, X. et al. A comparison of deep learning performance against health-care professionals in detecting diseases from medical imaging: a systematic review and meta-analysis. *Lancet Digit. Health* **1**, e271–e297 (2019).
- Ross, C. & Swettlitz, I. IBM's Watson supercomputer recommended 'unsafe and incorrect' cancer treatments, internal documents show. *STAT* <https://www.statnews.com/2018/07/25/ibm-watson-recommended-unsafe-incorrect-treatments/> (2018).
- Narla, A., Kuprel, B., Sarin, K., Novoa, R. & Ko, J. Automated classification of skin lesions: from pixels to practice. *J. Invest. Dermatol.* **138**, 2108–2110 (2018).
- Zech, J. R. et al. Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: a cross-sectional study. *PLoS Med.* **15**, e1002683 (2018).
- Winkler, J. K. et al. Association between surgical skin markings in dermoscopic images and diagnostic performance of a deep learning convolutional neural network for melanoma recognition. *JAMA Dermatol.* **155**, 1135–1141 (2019).
- Rueckert, D. & Schnabel, J. A. Model-based and data-driven strategies in medical image computing. *Proc. IEEE* **108**, 110–124 (2020).
- Zhang, C., Bengio, S., Hardt, M., Recht, B. & Vinyals, O. Understanding deep learning requires rethinking generalization. *Proc. Int. Conf. Learn. Represent.* <https://arxiv.org/abs/1611.03530> (2017).
- Liu, Y., Chen, P.-H. C., Krause, J. & Peng, L. How to read articles that use machine learning: users' guides to the medical literature. *JAMA* **322**, 1806–1816 (2019).
- Ransohoff, D. F. Bias as a threat to the validity of cancer molecular-marker research. *Nat. Rev. Cancer* **5**, 142–149 (2005).
- Moons, K. G. M. et al. PROBAST: a tool to assess risk of bias and applicability of prediction model studies: explanation and elaboration. *Ann. Intern. Med.* **170**, W1–W33 (2019).
- Simard, P., Victorri, B., LeCun, Y. & Denker, J. Tangent Prop — a formalism for specifying selected invariances in an adaptive network. *Adv. Neural Inf. Process. Syst.* **4**, 895–903 (1992).
- Shorten, C. & Khoshgoftaar, T. M. A survey on image data augmentation for deep learning. *J. Big Data* **6**, 60 (2019).
- Ioannidis, J. P. A. What have we (not) learnt from millions of scientific papers with *P* values? *Am. Stat.* **73**, 20–25 (2019).
- Ioannidis, J. P. A. Why most published research findings are false. *PLoS Med.* **2**, e124 (2005).
- Moons, K. G. M. et al. Transparent Reporting of a Multivariable Prediction Model for Individual Prognosis Or Diagnosis (TRIPOD): explanation and elaboration. *Ann. Intern. Med.* **162**, W1–W73 (2015).
- Heaven, D. Why deep-learning AIs are so easy to fool. *Nature* **574**, 163–166 (2019).
- Ioannidis, J. P. A. Evolution and translation of research findings: from bench to where? *PLoS Clin. Trials* **1**, e36 (2006).
- Topol, E. J. High-performance medicine: the convergence of human and artificial intelligence. *Nat. Med.* **25**, 44–56 (2019).
- Justice, A. C., Covinsky, K. E. & Berlin, J. A. Assessing the generalizability of prognostic information. *Ann. Intern. Med.* **130**, 515–524 (1999).
- Subbaswamy, A. & Saria, S. From development to deployment: dataset shift, causality, and shift-stable models in health AI. *Biostatistics* **21**, 345–352 (2020).
- Ioannidis, J. P. A. & Khoury, M. J. Improving validation practices in "omics" research. *Science* **334**, 1230–1232 (2011).
- Obermeyer, Z. & Emanuel, E. J. Predicting the future — big data, machine learning, and clinical medicine. *N. Engl. J. Med.* **375**, 1216–1219 (2016).
- Keane, P. A. & Topol, E. J. With an eye to AI and autonomous diagnosis. *NPJ Digit. Med.* **1**, 40 (2018).
- Gianfrancesco, M. A., Tamang, S., Yazdany, J. & Schmajuk, G. Potential biases in machine learning algorithms using electronic health record data. *JAMA Intern. Med.* **178**, 1544–1547 (2018).
- Noor, P. Can we trust AI not to further embed racial bias and prejudice? *BMJ* **368**, m363 (2020).
- Luo, W. et al. Guidelines for developing and reporting machine learning predictive models in biomedical research: a multidisciplinary view. *J. Med. Internet Res.* **18**, e323 (2016).
- Hua, K. L., Hsu, C. H., Hidayati, S. C., Cheng, W. H. & Chen, Y. J. Computer-aided classification of lung nodules on computed tomography images via deep learning technique. *Onco Targets Ther.* **8**, 2015–2022 (2015).
- Ciampi, F. et al. Automatic classification of pulmonary peri-fissural nodules in computed tomography using an ensemble of 2D views and a convolutional neural network out-of-the-box. *Med. Image Anal.* **26**, 195–202 (2015).
- Arevalo, J., González, F. A., Ramos-Pollán, R., Oliveira, J. L. & Guevara-Lopez, M. A. Representation learning for mammography mass lesion classification with convolutional neural networks. *Comput. Methods Prog. Biomed.* **127**, 248–257 (2016).
- Setio, A. A. A. et al. Pulmonary nodule detection in CT images: false positive reduction using multi-view convolutional networks. *IEEE Trans. Med. Imaging* **35**, 1160–1169 (2016).
- Roth, H. R. et al. Improving computer-aided detection using convolutional neural networks and random view aggregation. *IEEE Trans. Med. Imaging* **35**, 1170–1181 (2016).
- Kallenberg, M. et al. Unsupervised deep learning applied to breast density segmentation and mammographic risk scoring. *IEEE Trans. Med. Imaging* **35**, 1322–1331 (2016).
- Litjens, G. et al. Deep learning as a tool for increased accuracy and efficiency of histopathological diagnosis. *Sci. Rep.* **6**, 26286 (2016).
- Huynh, B. O., Li, H. & Giger, M. L. Digital mammographic tumor classification using transfer learning from deep convolutional neural networks. *J. Med. Imaging* **3**, 034501 (2016).
- Nie, K. et al. Rectal cancer: assessment of neoadjuvant chemoradiation outcome based on radiomics of multiparametric MRI. *Clin. Cancer Res.* **22**, 5256–5264 (2016).
- Kooi, T. et al. Large scale deep learning for computer aided detection of mammographic lesions. *Med. Image Anal.* **35**, 303–312 (2017).
- Esteva, A. et al. Dermatologist-level classification of skin cancer with deep neural networks. *Nature* **542**, 115–118 (2017).
- Dhungel, N., Carneiro, G. & Bradley, A. P. A deep learning approach for the analysis of masses in mammograms with minimal user intervention. *Med. Image Anal.* **37**, 114–128 (2017).
- Yu, L., Chen, H., Dou, Q., Qin, J. & Heng, P. Automated melanoma recognition in dermoscopy images via very deep residual networks. *IEEE Trans. Med. Imaging* **36**, 994–1004 (2017).
- Sun, W., Tseng, T. B., Zhang, J. & Qian, W. Enhancing deep convolutional neural network scheme for breast cancer diagnosis with unlabeled data. *Comput. Med. Imaging Graph.* **57**, 4–9 (2017).
- Cruz-Roa, A. et al. Accurate and reproducible invasive breast cancer detection in whole-slide images: a deep learning approach for quantifying tumor extent. *Sci. Rep.* **7**, 46450 (2017).
- Ciampi, F. et al. Towards automatic pulmonary nodule management in lung cancer screening with deep learning. *Sci. Rep.* **7**, 46479 (2017).
- Araújo, T. et al. Classification of breast cancer histology images using convolutional neural networks. *PLoS ONE* **12**, e0177544 (2017).
- Becker, A. S. et al. Deep learning in mammography: diagnostic accuracy of a multipurpose image analysis software in the detection of breast cancer. *Invest. Radiol.* **52**, 434–440 (2017).
- Dou, Q., Chen, H., Yu, L., Qin, J. & Heng, P. Multilevel contextual 3-D CNNs for false positive reduction in pulmonary nodule detection. *IEEE Trans. Biomed. Eng.* **64**, 1558–1567 (2017).
- Lao, J. et al. A deep learning-based radiomics model for prediction of survival in glioblastoma multiforme. *Sci. Rep.* **7**, 10353 (2017).
- Setio, A. A. A. et al. Validation, comparison, and combination of algorithms for automatic detection of pulmonary nodules in computed tomography images: the LUNA16 challenge. *Med. Image Anal.* **42**, 1–13 (2017).
- Ehteshami Bejnordi, B. et al. Diagnostic assessment of deep learning algorithms for detection of lymph node metastases in women with breast cancer. *JAMA* **318**, 2199–2210 (2017).
- Mohamed, A. A. et al. A deep learning method for classifying mammographic breast density categories. *Med. Phys.* **45**, 314–321 (2018).
- Khosravi, P., Kazemi, E., Imielinski, M., Elemento, O. & Hajirasouliha, I. Deep convolutional neural networks enable discrimination of heterogeneous digital pathology images. *EBioMedicine* **27**, 317–328 (2018).
- Xiao, Y., Wu, J., Lin, Z. & Zhao, X. A deep learning-based multi-model ensemble method for cancer prediction. *Comput. Methods Prog. Biomed.* **153**, 1–9 (2018).
- Marchetti, M. A. et al. Results of the 2016 International Skin Imaging Collaboration International Symposium on Biomedical Imaging challenge: comparison of the accuracy of computer algorithms to dermatologists for the diagnosis of melanoma from dermoscopic images. *J. Am. Acad. Dermatol.* **78**, 270–277.e1 (2018).
- Chen, P.-J. et al. Accurate classification of diminutive colorectal polyps using computer-aided analysis. *Gastroenterology* **154**, 568–575 (2018).
- Bychkov, D. et al. Deep learning based tissue analysis predicts outcome in colorectal cancer. *Sci. Rep.* **8**, 3395 (2018).
- Yasaka, K., Akai, H., Abe, O. & Kiryu, S. Deep learning with convolutional neural network for differentiation of liver masses at dynamic contrast-enhanced CT: a preliminary study. *Radiology* **286**, 887–896 (2018).
- Chang, K. et al. Residual convolutional neural network for the determination of IDH status in low- and high-grade gliomas from MR imaging. *Clin. Cancer Res.* **24**, 1073–1081 (2018).
- Ribli, D., Horváth, A., Unger, Z., Pollner, P. & Csabai, I. Detecting and classifying lesions in mammograms with deep learning. *Sci. Rep.* **8**, 4165 (2018).
- Chaudhary, K., Poirion, O. B., Lu, L. & Garmire, L. X. Deep learning-based multi-omics integration robustly predicts survival in liver cancer. *Clin. Cancer Res.* **24**, 1248–1259 (2018).
- Mobadersany, P. et al. Predicting cancer outcomes from histology and genomics using convolutional networks. *Proc. Natl. Acad. Sci. USA* **115**, E2970–E2979 (2018).
- Saltz, J. et al. Spatial organization and molecular correlation of tumor-infiltrating lymphocytes using deep learning on pathology images. *Cell Rep.* **23**, 181–193.e7 (2018).
- van de Goor, R., van Hooren, M., Dingemans, A.-M., Kremer, B. & Kross, K. Training and validating a portable electronic nose for lung cancer screening. *J. Thorac. Oncol.* **13**, 676–681 (2018).
- Chang, H., Han, J., Zhong, C., Snijders, A. M. & Mao, J. Unsupervised transfer learning via multi-scale convolutional sparse coding for biomedical applications. *IEEE Trans. Pattern Anal. Mach. Intell.* **40**, 1182–1194 (2018).

70. Han, S. S. et al. Classification of the clinical images for benign and malignant cutaneous tumors using a deep learning algorithm. *J. Invest. Dermatol.* **138**, 1529–1538 (2018).
71. Hirasawa, T. et al. Application of artificial intelligence using a convolutional neural network for detecting gastric cancer in endoscopic images. *Gastric Cancer* **21**, 653–660 (2018).
72. Chang, P. et al. Deep-learning convolutional neural networks accurately classify genetic mutations in gliomas. *Am. J. Neuroradiol.* **39**, 1201–1207 (2018).
73. Haenssle, H. A. et al. Man against machine: diagnostic performance of a deep learning convolutional neural network for dermoscopic melanoma recognition in comparison to 58 dermatologists. *Ann. Oncol.* **29**, 1836–1842 (2018).
74. Coudray, N. et al. Classification and mutation prediction from non-small cell lung cancer histopathology images using deep learning. *Nat. Med.* **24**, 1559–1567 (2018).
75. Wang, P. et al. Development and validation of a deep-learning algorithm for the detection of polyps during colonoscopy. *Nat. Biomed. Eng.* **2**, 741–748 (2018).
76. Urban, G. et al. Deep learning localizes and identifies polyps in real time with 96% accuracy in screening colonoscopy. *Gastroenterology* **155**, 1069–1078.e8 (2018).
77. Rajpurkar, P. et al. Deep learning for chest radiograph diagnosis: a retrospective comparison of the CheXNeXt algorithm to practicing radiologists. *PLoS Med.* **15**, e1002686 (2018).
78. Hosny, A. et al. Deep learning for lung cancer prognostication: a retrospective multi-cohort radiomics study. *PLoS Med.* **15**, e1002711 (2018).
79. Nam, J. G. et al. Development and validation of deep learning-based automatic detection algorithm for malignant pulmonary nodules on chest radiographs. *Radiology* **290**, 218–228 (2019).
80. Byrne, M. F. et al. Real-time differentiation of adenomatous and hyperplastic diminutive colorectal polyps during analysis of unaltered videos of standard colonoscopy using a deep learning model. *Gut* **68**, 94–100 (2019).
81. Horie, Y. et al. Diagnostic outcomes of esophageal cancer by artificial intelligence using convolutional neural networks. *Gastrointest. Endosc.* **89**, 25–32 (2019).
82. Kather, J. N. et al. Predicting survival from colorectal cancer histology slides using deep learning: a retrospective multicenter study. *PLoS Med.* **16**, e1002730 (2019).
83. Rodríguez-Ruiz, A. et al. Detection of breast cancer with mammography: effect of an artificial intelligence support system. *Radiology* **290**, 305–314 (2019).
84. Li, X. et al. Diagnosis of thyroid cancer using deep convolutional neural network models applied to sonographic images: a retrospective, multicohort, diagnostic study. *Lancet Oncol.* **20**, 193–201 (2019).
85. Wang, S. et al. Predicting EGFR mutation status in lung adenocarcinoma on CT image using deep learning. *Eur. Respir. J.* **53**, 1800986 (2019).
86. Brinker, T. J. et al. A convolutional neural network trained with dermoscopic images performed on par with 145 dermatologists in a clinical melanoma image classification task. *Eur. J. Cancer* **111**, 148–154 (2019).
87. Kickingereder, P. et al. Automated quantitative tumour response assessment of MRI in neuro-oncology with artificial neural networks: a multicentre, retrospective study. *Lancet Oncol.* **20**, 728–740 (2019).
88. Brinker, T. J. et al. Deep learning outperformed 136 of 157 dermatologists in a head-to-head dermoscopic melanoma image classification task. *Eur. J. Cancer* **113**, 47–54 (2019).
89. Choi, K. S., Choi, S. H. & Jeong, B. Prediction of IDH genotype in gliomas with dynamic susceptibility contrast perfusion MR imaging using an explainable recurrent neural network. *Neuro Oncol.* **21**, 1197–1209 (2019).
90. Ardila, D. et al. End-to-end lung cancer screening with three-dimensional deep learning on low-dose chest computed tomography. *Nat. Med.* **25**, 954–961 (2019).
91. Yala, A., Lehman, C., Schuster, T., Portnoi, T. & Barzilay, R. A deep learning mammography-based model for improved breast cancer risk prediction. *Radiology* **292**, 60–66 (2019).
92. Kather, J. N. et al. Deep learning can predict microsatellite instability directly from histology in gastrointestinal cancer. *Nat. Med.* **25**, 1054–1056 (2019).
93. Liu, Y. et al. Artificial intelligence-based breast cancer nodal metastasis detection: insights into the black box for pathologists. *Arch. Pathol. Lab. Med.* **143**, 859–868 (2019).
94. Kehl, K. L. et al. Assessment of deep natural language processing in ascertaining oncologic outcomes from radiology reports. *JAMA Oncol.* **5**, 1421–1429 (2019).
95. Campanella, G. et al. Clinical-grade computational pathology using weakly supervised deep learning on whole slide images. *Nat. Med.* **25**, 1301–1309 (2019).
96. Chen, P.-H. C. et al. An augmented reality microscope with real-time artificial intelligence integration for cancer diagnosis. *Nat. Med.* **25**, 1453–1457 (2019).
97. Hu, L. et al. An observational study of deep learning and automated evaluation of cervical images for cancer screening. *J. Natl Cancer Inst.* **111**, 923–932 (2019).
98. Rodríguez-Ruiz, A. et al. Stand-alone artificial intelligence for breast cancer detection in mammography: comparison with 101 radiologists. *J. Natl Cancer Inst.* **111**, 916–922 (2019).
99. Wang, X. et al. Weakly supervised deep learning for whole slide lung cancer image analysis. *IEEE Trans. Cybern.* **50**, 3950–3962 (2019).
100. Jurmeister, P. et al. Machine learning analysis of DNA methylation profiles distinguishes primary lung squamous cell carcinomas from head and neck metastases. *Sci. Transl. Med.* **11**, eaaw8513 (2019).
101. Courtiol, P. et al. Deep learning-based classification of mesothelioma improves prediction of patient outcome. *Nat. Med.* **25**, 1519–1525 (2019).
102. Wang, P. et al. Real-time automatic detection system increases colonoscopic polyp and adenoma detection rates: a prospective randomised controlled study. *Gut* **68**, 1813–1819 (2019).
103. Liao, F., Liang, M., Li, Z., Hu, X. & Song, S. Evaluate the malignancy of pulmonary nodules using the 3-D deep leaky noisy-OR network. *IEEE Trans. Neural Netw. Learn. Syst.* **30**, 3484–3495 (2019).
104. Luo, H. et al. Real-time artificial intelligence for detection of upper gastrointestinal cancer by endoscopy: a multicentre, case-control, diagnostic study. *Lancet Oncol.* **20**, 1645–1654 (2019).
105. Wu, L. et al. Randomised controlled trial of WISENSE, a real-time quality improving system for monitoring blind spots during esophagogastroduodenoscopy. *Gut* **68**, 2161–2169 (2019).
106. Shkolyar, E. et al. Augmented bladder tumor detection using deep learning. *Eur. Urol.* **76**, 714–718 (2019).
107. Yamamoto, Y. et al. Automated acquisition of explainable knowledge from unannotated histopathology images. *Nat. Commun.* **10**, 5642 (2019).
108. McKinney, S. M. et al. International evaluation of an AI system for breast cancer screening. *Nature* **577**, 89–94 (2020).
109. Hollon, T. C. et al. Near real-time intraoperative brain tumor diagnosis using stimulated Raman histology and deep neural networks. *Nat. Med.* **26**, 52–58 (2020).
110. Haenssle, H. A. et al. Man against machine reloaded: performance of a market-approved convolutional neural network in classifying a broad spectrum of skin lesions in comparison with 96 dermatologists working under less artificial conditions. *Ann. Oncol.* **31**, 137–143 (2020).
111. Ström, P. et al. Artificial intelligence for diagnosis and grading of prostate cancer in biopsies: a population-based, diagnostic study. *Lancet Oncol.* **21**, 222–232 (2020).
112. Bulten, W. et al. Automated deep-learning system for Gleason grading of prostate cancer using biopsies: a diagnostic study. *Lancet Oncol.* **21**, 233–241 (2020).
113. Skrede, O.-J. et al. Deep learning for prediction of colorectal cancer outcome: a discovery and validation study. *Lancet* **395**, 350–360 (2020).
114. Saillard, C. et al. Predicting survival after hepatocellular carcinoma resection using deep-learning on histological slides. *Hepatology* **72**, 2000–2013 (2020).
115. Jin, E. H. et al. Improved accuracy in optical diagnosis of colorectal polyps using convolutional neural networks with visual explanations. *Gastroenterology* **158**, 2169–2179.e8 (2020).
116. de Groof, A. J. et al. Deep-learning system detects neoplasia in patients with Barrett's esophagus with higher accuracy than endoscopists in a multistep training and validation study with benchmarking. *Gastroenterology* **158**, 915–929.e4 (2020).
117. Bangalore Yogananda, C. G. et al. A novel fully automated MRI-based deep-learning method for classification of IDH mutation status in brain gliomas. *Neuro Oncol.* **22**, 402–411 (2020).
118. Zheng, X. et al. Deep learning radiomics can predict axillary lymph node status in early-stage breast cancer. *Nat. Commun.* **11**, 1236 (2020).
119. Galateau Salle, F. et al. Comprehensive molecular and pathologic evaluation of transitional mesothelioma assisted by deep learning approach: a multi-institutional study of the International Mesothelioma Panel from the MESOPATH Reference Center. *J. Thorac. Oncol.* **15**, 1037–1053 (2020).
120. Baldwin, D. R. et al. External validation of a convolutional neural network artificial intelligence tool to predict malignancy in pulmonary nodules. *Thorax* **75**, 306–312 (2020).
121. Wang, P. et al. Effect of a deep-learning computer-aided detection system on adenoma detection during colonoscopy (CADE-DB trial): a double-blind randomised study. *Lancet Gastroenterol. Hepatol.* **5**, 343–351 (2020).
122. Song, Q., Zheng, Y., Sheng, W. & Yang, J. Tridirectional transfer learning for predicting gastric cancer morbidity. *IEEE Trans. Neural Netw. Learn. Syst.* <https://doi.org/10.1109/TNNLS.2020.2979486> (2020).
123. Dong, D. et al. Deep learning radiomic nomogram can predict the number of lymph node metastasis in locally advanced gastric cancer: an international multicenter study. *Ann. Oncol.* **31**, 912–920 (2020).
124. Shin, H. et al. Early-stage lung cancer diagnosis by deep learning-based spectroscopic analysis of circulating exosomes. *ACS Nano* **14**, 5435–5444 (2020).
125. Kann, B. H. et al. Multi-institutional validation of deep learning for pretreatment identification of extranodal extension in head and neck squamous cell carcinoma. *J. Clin. Oncol.* **38**, 1304–1311 (2020).
126. [No authors listed] AI diagnostics need attention. *Nature* **555**, 285 (2018).
127. [No authors listed] Is digital medicine different? *Lancet* **392**, 95 (2018).
128. Kawaguchi, K., Kaelbling, L. P. & Bengio, Y. Generalization in deep learning. *arXiv* <https://arxiv.org/abs/1710.05468> (2017).
129. LeCun, Y. in *Connectionism in Perspective* (eds Pfeifer, R., Schreier, Z., Fogelman, F., & Steels, L.) 143–156 (Elsevier, 1989).
130. Neyshabur, B., Bhojanapalli, S., McAllester, D. & Srebro, N. Exploring generalization in deep learning. *Adv. Neural Inf. Process. Syst.* **30**, 5947–5956 (2017).
131. Pan, S. J. & Yang, Q. A survey on transfer learning. *IEEE Trans. Knowl. Data Eng.* **22**, 1345–1359 (2010).
132. Weiss, K., Khoshgoftar, T. M. & Wang, D. A survey of transfer learning. *J. Big Data* **3**, 9 (2016).
133. Deng, J. et al. ImageNet: a large-scale hierarchical image database. *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* <https://doi.org/10.1109/CVPR.2009.5206848> (2009).
134. Russakovsky, O. et al. ImageNet large scale visual recognition challenge. *Int. J. Comput. Vis.* **115**, 211–252 (2015).
135. Shankar, S. et al. No classification without representation: assessing geodiversity issues in open data sets for the developing world. *NIPS Workshop Mach. Learn. Dev. World* <https://arxiv.org/abs/1711.08536> (2017).
136. Geirhos, R. et al. ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. *Proc. Int. Conf. Learn. Represent.* <https://arxiv.org/abs/1811.12231> (2019).
137. Beyer, L., Hénaff, O. J., Kolesnikov, A., Zhai, X. & van den Oord, A. Are we done with ImageNet? *arXiv* <https://arxiv.org/abs/2006.07159> (2020).
138. Sun, C., Shrivastava, A., Singh, S. & Gupta, A. Revisiting unreasonable effectiveness of data in deep learning era. *Proc. IEEE Int. Conf. Comput. Vis.* <https://doi.org/10.1109/ICCV.2017.97> (2017).
139. Simard, P. Y., Steinkraus, D. & Platt, J. C. Best practices for convolutional neural networks applied to visual document analysis. *Proc. 7th Int. Conf. Anal. Recognit.* <https://doi.org/10.1109/ICDAR.2003.1227801> (2003).
140. Baird, H. S. Document image defect models and their uses. *Proc. 2nd Int. Conf. Doc. Anal. Recognit.* <https://doi.org/10.1109/ICDAR.1993.395781> (1993).
141. Stacke, K., Eilertsen, G., Unger, J. & Lundström, C. Measuring domain shift for deep learning in histopathology. *IEEE J. Biomed. Health Inform.* <https://doi.org/10.1109/JBHI.2020.3032060> (2020).

142. Lakhani, P. & Sundaram, B. Deep learning at chest radiography: automated classification of pulmonary tuberculosis by using convolutional neural networks. *Radiology* **284**, 574–582 (2017).
143. Hussain, Z., Gimenez, F., Yi, D. & Rubin, D. Differential data augmentation techniques for medical imaging classification tasks. *AMIA Annu. Symp. Proc.* **2017**, 979–984 (2018).
144. Sajjad, M. et al. Multi-grade brain tumor classification using deep CNN with extensive data augmentation. *J. Comput. Sci.* **30**, 174–182 (2019).
145. Tellez, D. et al. Quantifying the effects of data augmentation and stain color normalization in convolutional neural networks for computational pathology. *Med. Image Anal.* **58**, 101544 (2019).
146. Kerr, R. S. et al. Adjuvant capecitabine plus bevacizumab versus capecitabine alone in patients with colorectal cancer (QUASAR 2): an open-label, randomised phase 3 trial. *Lancet Oncol.* **17**, 1543–1557 (2016).
147. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J. & Wojna, Z. Rethinking the inception architecture for computer vision. *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.* <https://doi.org/10.1109/CVPR.2016.308> (2016).
148. Miller, R. G. J. *Simultaneous Statistical Inference* 2nd edn (Springer, 1981).
149. Hochberg, Y. & Tamhane, A. C. *Multiple Comparison Procedures* (Wiley, 2009).
150. Michiels, S., Koscielny, S. & Hill, C. Prediction of cancer outcome with microarrays: a multiple random validation strategy. *Lancet* **365**, 488–492 (2005).
151. Hastie, T., Tibshirani, R. & Friedman, J. *The Elements of Statistical Learning* 2nd edn (Springer-Verlag, 2009).
152. Russell, S. & Norvig, P. *Artificial Intelligence: A Modern Approach* 3rd edn (Prentice Hall, 2010).
153. Hemingway, H., Riley, R. D. & Altman, D. G. Ten steps towards improving prognosis research. *BMJ* **339**, b4184 (2009).
154. Korevaar, D. A. et al. Facilitating prospective registration of diagnostic accuracy studies: a STARD initiative. *Clin. Chem.* **63**, 1331–1341 (2017).
155. Ioannidis, J. P. A. The importance of predefined rules and prespecified statistical analyses: do not abandon significance. *JAMA* **321**, 2067–2068 (2019).
156. Brodersen, K. H., Ong, C. S., Stephan, K. E. & Buhmann, J. M. The balanced accuracy and its posterior distribution. *Proc. 20th Int. Conf. Pattern Recognit.* <https://doi.org/10.1109/ICPR.2010.764> (2010).
157. van den Hout, W. B. The area under an ROC curve with limited information. *Med. Decis. Mak.* **23**, 160–166 (2003).
158. Fawcett, T. An introduction to ROC analysis. *Pattern Recognit. Lett.* **27**, 861–874 (2006).
159. Harrell, F. E. Jr, Califf, R. M., Pryor, D. B., Lee, K. L. & Rosati, R. A. Evaluating the yield of medical tests. *J. Am. Med. Assoc.* **247**, 2543–2546 (1982).
160. Lobo, J. M., Jiménez-Valverde, A. & Real, R. AUC: a misleading measure of the performance of predictive distribution models. *Glob. Ecol. Biogeogr.* **17**, 145–151 (2008).
161. Voosen, P. How AI detectives are cracking open the black box of deep learning. *Science* <https://www.sciencemag.org/news/2017/07/how-ai-detectives-are-cracking-open-black-box-deep-learning> (2017).
162. Adadi, A. & Berrada, M. Peeking inside the black-box: a survey on explainable artificial intelligence (XAI). *IEEE Access.* **6**, 52138–52160 (2018).
163. Barredo Arrieta, A. et al. Explainable artificial intelligence (XAI): concepts, taxonomies, opportunities and challenges toward responsible AI. *Inf. Fusion.* **58**, 82–115 (2020).
164. Montavon, G., Samek, W. & Müller, K.-R. Methods for interpreting and understanding deep neural networks. *Digit. Signal. Process.* **73**, 1–15 (2018).
165. Simonyan, K., Vedaldi, A. & Zisserman, A. Deep inside convolutional networks: visualising image classification models and saliency maps. *Proc. Int. Conf. Learn. Represent.* <https://arxiv.org/abs/1312.6034> (2014).
166. Bach, S. et al. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS ONE* **10**, e0130140 (2015).
167. Sundararajan, M., Taly, A. & Yan, Q. Axioomatic attribution for deep networks. *Proc. 34th Int. Conf. Mach. Learn.* **70**, 3319–3328 (2017).
168. Friedman, L. M., Furberg, C. D., DeMets, D. L., Reboussin, D. M. & Granger, C. B. *Fundamentals of Clinical Trials* 5th edn (Springer, 2015).
169. van Luijn, H. E. M., Musschenga, A. W., Keus, R. B., Robinson, W. M. & Aaronson, N. K. Assessment of the risk/benefit ratio of phase II cancer clinical trials by Institutional Review Board (IRB) members. *Ann. Oncol.* **13**, 1307–1313 (2002).
170. Martin, L., Hutchens, M., Hawkins, C. & Radnov, A. How much do clinical trials cost? *Nat. Rev. Drug Discov.* **16**, 381–382 (2017).
171. Teutsch, S. M. et al. The evaluation of genomic applications in practice and prevention (EGAPP) initiative: methods of the EGAPP Working Group. *Genet. Med.* **11**, 3–14 (2009).
172. Vollmer, S. et al. Machine learning and artificial intelligence research for patient benefit: 20 critical questions on transparency, replicability, ethics, and effectiveness. *BMJ* **368**, i6927 (2020).
173. Chan, A.-W. et al. SPIRIT 2013 explanation and elaboration: guidance for protocols of clinical trials. *BMJ* **346**, e7586 (2013).
174. Cruz Rivera, S. et al. Guidelines for clinical trial protocols for interventions involving artificial intelligence: the SPIRIT-AI extension. *Nat. Med.* **26**, 1351–1363 (2020).
175. Moher, D. et al. CONSORT 2010 explanation and elaboration: updated guidelines for reporting parallel group randomised trials. *BMJ* **340**, c869 (2010).
176. Collins, G. S. & Moons, K. G. M. Reporting of artificial intelligence prediction models. *Lancet* **393**, 1577–1579 (2019).
177. Liu, X. et al. Reporting guidelines for clinical trial reports for interventions involving artificial intelligence: the CONSORT-AI extension. *Nat. Med.* **26**, 1364–1374 (2020).
178. [No authors listed] Should protocols for observational research be registered? *Lancet* **375**, 348 (2010).
179. Loder, E., Groves, T. & MacAuley, D. Registration of observational studies. *BMJ* **340**, c950 (2010).
180. Chambers, C. & Munafo, M. Trust in science would be improved by study pre-registration. *The Guardian* <https://www.theguardian.com/science/blog/2013/jun/05/trust-in-science-study-pre-registration> (2013).
181. Williams, R. J., Tse, T., Harlan, W. R. & Zarin, D. A. Registration of observational studies: is it time? *Can. Med. Assoc. J.* **182**, 1638–1642 (2010).
182. Gill, J. & Prasad, V. Improving observational studies in the era of big data. *Lancet* **392**, 716–717 (2018).
183. Sørensen, H. T. & Rothman, K. J. The prognosis for research. *BMJ* **340**, c703 (2010).
184. Vandenbroucke, J. P. Registering observational research: second thoughts. *Lancet* **375**, 982–983 (2010).
185. [No authors listed] The registration of observational studies — when metaphors go bad. *Epidemiology* **21**, 607–609 (2010).
186. Andre, F. et al. Biomarker studies: a call for a comprehensive biomarker study registry. *Nat. Rev. Clin. Oncol.* **8**, 171–176 (2011).
187. Hooft, L. & Bossuyt, P. M. Prospective registration of marker evaluation studies: time to act. *Clin. Chem.* **57**, 1684–1686 (2011).
188. Altman, D. G. The time has come to register diagnostic and prognostic research. *Clin. Chem.* **60**, 580–582 (2014).
189. Rifai, N. et al. Registering diagnostic and prognostic trials of tests: is it the right thing to do? *Clin. Chem.* **60**, 1146–1152 (2014).
190. Rajkumar, A., Dean, J. & Kohane, I. Machine learning in medicine. *N. Engl. J. Med.* **380**, 1347–1358 (2019).
191. Zou, J. & Schiebinger, L. AI can be sexist and racist — it's time to make it fair. *Nature* **559**, 324–326 (2018).
192. Adamson, A. S. & Smith, A. Machine learning and health care disparities in dermatology. *JAMA Dermatol.* **154**, 1247–1248 (2018).
193. Vyas, D. A., Eisenstein, L. G. & Jones, D. S. Hidden in plain sight — reconsidering the use of race correction in clinical algorithms. *N. Engl. J. Med.* **383**, 874–882 (2020).
194. Obermeyer, Z., Powers, B., Vogeli, C. & Mullainathan, S. Dissecting racial bias in an algorithm used to manage the health of populations. *Science* **366**, 447–453 (2019).
195. Rajkumar, A., Hardt, M., Howell, M. D., Corrado, G. & Chin, M. H. Ensuring fairness in machine learning to advance health equity. *Ann. Intern. Med.* **169**, 866–872 (2018).
196. Owens, K. & Walker, A. Those designing healthcare algorithms must become actively anti-racist. *Nat. Med.* **26**, 1327–1328 (2020).
197. Moons, K. G. M. et al. Risk prediction models: II. External validation, model updating, and impact assessment. *Heart* **98**, 691–698 (2012).

#### Acknowledgements

The authors thank M. Seiergren for assembling all figures, T. S. Hveem for discussions, T. Ystanes, H. A. Inderhaug and B. M. Sannes for setting up and maintaining our computer network and computational infrastructure, and the authors of Inception-v3 for making their code freely available under an open source licence (Apache License, version 2.0). The authors of this Perspective acknowledge funding from the Research Council of Norway through its IKTPLUSS Lighthouse programme (project number 259204).

#### Author contributions

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. evaluated the recent, presumably influential, deep learning studies in cancer diagnostics. S.D.R. executed the training, tuning and evaluation of Inception-v3 systems. A.K. drafted the manuscript, and all authors contributed to reviewing and editing the manuscript.

#### Competing interests

The authors declare no competing interests.

#### Peer review information

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

#### Publisher's note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

#### Supplementary information

Supplementary information is available for this paper at <https://doi.org/10.1038/s41568-020-00327-9>.

#### RELATED LINKS

ClinicalTrials.gov registry: <https://www.clinicaltrials.gov>  
International Standard Randomized Controlled Trial Number (ISRCTN) registry: <https://www.isrctn.com>  
Journal Policies | *Journal of Clinical Oncology*: <https://ascopubs.org/jco/authors/journal-policies>

© Springer Nature Limited 2021