

# Transparency and reproducibility in artificial intelligence

Citation for published version (APA):

Haibe-Kains, B., Adam, G. A., Hosny, A., Khodakarami, F., Shraddha, T., Shraddha, T., Kusko, R., Sansone, S.-A., Tong, W., Wolfinger, R. D., Mason, C. E., Jones, W., Dopazo, J., Furlanello, C., Waldron, L., Wang, B., McIntosh, C., Goldenberg, A., Kundaje, A., ... Aerts, H. J. W. L. (2020). Transparency and reproducibility in artificial intelligence. *Nature*, 586(7829), E14-E16. <https://doi.org/10.1038/s41586-020-2766-y>

## Document status and date:

Published: 15/10/2020

## DOI:

[10.1038/s41586-020-2766-y](https://doi.org/10.1038/s41586-020-2766-y)

## Document Version:

Publisher's PDF, also known as Version of record

## Document license:

Taverne

## Please check the document version of this publication:

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

## General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

[www.umlib.nl/taverne-license](http://www.umlib.nl/taverne-license)

## Take down policy

If you believe that this document breaches copyright please contact us at:

[repository@maastrichtuniversity.nl](mailto:repository@maastrichtuniversity.nl)

providing details and we will investigate your claim.

# Transparency and reproducibility in artificial intelligence

<https://doi.org/10.1038/s41586-020-2766-y>

Received: 1 February 2020

Accepted: 10 August 2020

Published online: 14 October 2020

 Check for updates

Benjamin Haibe-Kains<sup>1,2,3,4,5</sup>✉, George Alexandru Adam<sup>3,5</sup>, Ahmed Hosny<sup>6,7</sup>, Farnoosh Khodakarami<sup>1,2</sup>, Massive Analysis Quality Control (MAQC) Society Board of Directors\*, Levi Waldron<sup>8</sup>, Bo Wang<sup>2,3,5,9,10</sup>, Chris McIntosh<sup>2,5,9</sup>, Anna Goldenberg<sup>3,5,11,12</sup>, Anshul Kundaje<sup>13,14</sup>, Casey S. Greene<sup>15,16</sup>, Tamara Broderick<sup>17</sup>, Michael M. Hoffman<sup>1,2,3,5</sup>, Jeffrey T. Leek<sup>18</sup>, Keegan Korthauer<sup>19,20</sup>, Wolfgang Huber<sup>21</sup>, Alvis Brazma<sup>22</sup>, Joelle Pineau<sup>23,24</sup>, Robert Tibshirani<sup>25,26</sup>, Trevor Hastie<sup>25,26</sup>, John P. A. Ioannidis<sup>25,26,27,28,29</sup>, John Quackenbush<sup>30,31,32</sup> & Hugo J. W. L. Aerts<sup>6,7,33,34</sup>

ARISING FROM S. M. McKinney et al. *Nature* <https://doi.org/10.1038/s41586-019-1799-6> (2020)

Breakthroughs in artificial intelligence (AI) hold enormous potential as it can automate complex tasks and go even beyond human performance. In their study, McKinney et al.<sup>1</sup> showed the high potential of AI for breast cancer screening. However, the lack of details of the methods and algorithm code undermines its scientific value. Here, we identify obstacles that hinder transparent and reproducible AI research as faced by McKinney et al.<sup>1</sup>, and provide solutions to these obstacles with implications for the broader field.

The work by McKinney et al.<sup>1</sup> demonstrates the potential of AI in medical imaging, while highlighting the challenges of making such work reproducible. The authors assert that their system improves the speed and robustness of breast cancer screening, generalizes to populations beyond those used for training, and outperforms radiologists in specific settings. Upon successful prospective clinical validation and approval by regulatory bodies, this new system holds great potential for streamlining clinical workflows, reducing false positives, and improving patient outcomes. However, the absence of sufficiently documented methods and computer code underlying the study effectively undermines its scientific value. This shortcoming limits the evidence required for others to prospectively validate and clinically implement such technologies. By identifying obstacles hindering transparent and reproducible AI research as faced by McKinney et al.<sup>1</sup>, we provide potential solutions with implications for the broader field.

Scientific progress depends on the ability of independent researchers to scrutinize the results of a research study, to reproduce the study's main results using its materials, and to build on them in future studies (<https://www.nature.com/nature-research/editorial-policies/>

reporting-standards). Publication of insufficiently documented research does not meet the core requirements underlying scientific discovery<sup>2,3</sup>. Merely textual descriptions of deep-learning models can hide their high level of complexity. Nuances in the computer code may have marked effects on the training and evaluation of results<sup>4</sup>, potentially leading to unintended consequences<sup>5</sup>. Therefore, transparency in the form of the actual computer code used to train a model and arrive at its final set of parameters is essential for research reproducibility. McKinney et al.<sup>1</sup> stated that the code used for training the models has “a large number of dependencies on internal tooling, infrastructure and hardware”, and claimed that the release of the code was therefore not possible. Computational reproducibility is indispensable for high-quality AI applications<sup>6,7</sup>; more complex methods demand greater transparency<sup>8</sup>. In the absence of code, reproducibility falls back on replicating methods from textual description. Although, McKinney and colleagues<sup>1</sup> claim that all experiments and implementation details were described in sufficient detail in the supplementary methods section of their Article<sup>1</sup> to “support replication with non-proprietary libraries”, key details about their analysis are lacking. Even with extensive description, reproducing complex computational pipelines based purely on text is a subjective and challenging task<sup>9</sup>.

In addition to the reproducibility challenges inherent to purely textual descriptions of methods, the description by McKinney et al.<sup>1</sup> of the model development as well as data processing and training pipelines lacks crucial details. The definitions of several hyperparameters for the model's architecture (composed of three networks referred to as the breast, lesion and case models) are missing (Table 1). In their

<sup>1</sup>Princess Margaret Cancer Centre, University Health Network, Toronto, Ontario, Canada. <sup>2</sup>Department of Medical Biophysics, University of Toronto, Toronto, Ontario, Canada. <sup>3</sup>Department of Computer Science, University of Toronto, Toronto, Ontario, Canada. <sup>4</sup>Ontario Institute for Cancer Research, Toronto, Ontario, Canada. <sup>5</sup>Vector Institute for Artificial Intelligence, Toronto, Ontario, Canada. <sup>6</sup>Artificial Intelligence in Medicine (AIM) Program, Brigham and Women's Hospital, Harvard Medical School, Boston, MA, USA. <sup>7</sup>Radiation Oncology and Radiology, Dana-Farber Cancer Institute, Brigham and Women's Hospital, Harvard Medical School, Boston, MA, USA. <sup>8</sup>Department of Epidemiology and Biostatistics and Institute for Implementation Science in Population Health, CUNY Graduate School of Public Health and Health Policy, New York, NY, USA. <sup>9</sup>Peter Munk Cardiac Centre, University Health Network, Toronto, Ontario, Canada. <sup>10</sup>Department of Laboratory Medicine and Pathobiology, University of Toronto, Ontario, Canada. <sup>11</sup>SickKids Research Institute, Toronto, Ontario, Canada. <sup>12</sup>Child and Brain Development Program, CIFAR, Toronto, Ontario, Canada. <sup>13</sup>Department of Genetics, Stanford University School of Medicine, Stanford, CA, USA. <sup>14</sup>Department of Computer Science, Stanford University, Stanford, CA, USA. <sup>15</sup>Dept. of Systems Pharmacology and Translational Therapeutics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, USA. <sup>16</sup>Childhood Cancer Data Lab, Alex's Lemonade Stand Foundation, Philadelphia, PA, USA. <sup>17</sup>Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA, USA. <sup>18</sup>Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA. <sup>19</sup>Department of Statistics, University of British Columbia, Vancouver, British Columbia, Canada. <sup>20</sup>BC Children's Hospital Research Institute, Vancouver, British Columbia, Canada. <sup>21</sup>European Molecular Biology Laboratory, Genome Biology Unit, Heidelberg, Germany. <sup>22</sup>European Molecular Biology Laboratory, European Bioinformatics Institute, EMBL-EBI, Hinxton, UK. <sup>23</sup>McGill University, Montreal, Quebec, Canada. <sup>24</sup>Montreal Institute for Learning Algorithms, Quebec, Canada. <sup>25</sup>Department of Statistics, Stanford University School of Humanities and Sciences, Stanford, CA, USA. <sup>26</sup>Department of Biomedical Data Science, Stanford University School of Medicine, Stanford, CA, USA. <sup>27</sup>Department of Medicine, Stanford University School of Medicine, Stanford, CA, USA. <sup>28</sup>Meta-Research Innovation Center at Stanford (METRICS), Stanford, CA, USA. <sup>29</sup>Department of Epidemiology and Population Health, Stanford University School of Medicine, Stanford, CA, USA. <sup>30</sup>Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA, USA. <sup>31</sup>Channing Division of Network Medicine, Brigham and Women's Hospital, Boston, MA, USA. <sup>32</sup>Department of Data Science, Dana-Farber Cancer Institute, Boston, MA, USA. <sup>33</sup>Radiology and Nuclear Medicine, Maastricht University, Maastricht, The Netherlands. <sup>34</sup>Cardiovascular Imaging Research Center, Massachusetts General Hospital, Harvard Medical School, Boston, MA, USA. \*A list of authors and their affiliations appears at the end of the paper. ✉e-mail: bhaibeka@uhnresearch.ca

**Table 1 | Essential hyperparameters for reproducing the study for each of the three models**

	Lesion	Breast	Case
Learning rate	Missing	0.0001	Missing
Learning rate schedule	Missing	Stated	Missing
Optimizer	Stochastic gradient descent with momentum	Adam	Missing
Momentum	Missing	Not applicable	Not applicable
Batch size	4	Unclear	2
Epochs	Missing	120,000	Missing

publication, McKinney et al.<sup>1</sup> did not disclose the settings for the augmentation pipeline; the transformations used are stochastic and can considerably affect model performance<sup>10</sup>. Details of the training pipeline were also missing. Without this key information, independent reproduction of the training pipeline is not possible.

Numerous frameworks and platforms exist to make artificial intelligence research more transparent and reproducible (Table 2). For the sharing of code, these include Bitbucket, GitHub and GitLab, among others. The many software dependencies of large-scale machine learning applications require appropriate control of the software environment, which can be achieved through package managers including Conda, as well as container and virtualization systems, including Code Ocean, Gigantum, Colaboratory and Docker. If virtualization of the McKinney et al.<sup>1</sup> internal tooling proved to be difficult, they could have released the computer code and documentation. The authors could have also created small artificial examples or used small public datasets<sup>11</sup> to show how new data must be processed to train the model and generate predictions. Sharing the fitted model (architecture along with learned parameters) should be simple aside from privacy concerns that the model may reveal sensitive information about the set of patients used to train it. Nevertheless, techniques for achieving differential privacy exist to alleviate such concerns. Many platforms allow sharing of deep learning models, including TensorFlow Hub, ModelHub.ai, ModelDepot and Model Zoo with support for several frameworks such as PyTorch and Caffe, as well as the TensorFlow library used by the authors. In addition to improving accessibility and transparency, such resources can considerably accelerate model development, validation and transition into production and clinical implementation.

Another crucial aspect of ensuring reproducibility lies in access to the data the models were derived from. In their study, McKinney et al.<sup>1</sup> used two large datasets under license, properly disclosing this limitation in their publication. The sharing of patient health information is highly regulated owing to privacy concerns. Despite these challenges, the sharing of raw data has become more common in biomedical literature, increasing from under 1% in the early 2000s to 20% today<sup>12</sup>. However, if the data cannot be shared, the model predictions and data labels themselves should be released, allowing further statistical analyses. Above all, concerns about data privacy should not be used as a way to distract from the requirement to release code.

Although sharing of code and data are widely seen as a crucial part of scientific research, the adoption varies across fields. In fields such as genomics, complex computational pipelines and sensitive datasets have been shared for decades<sup>13</sup>. Guidelines related to genomic data are clear, detailed and, most importantly, enforced. It is generally accepted that all code and data are released alongside a publication. In other fields of medicine and science as a whole, this is much less common, and data and code are rarely made available. For scientific efforts in which a clinical application is envisioned and human

**Table 2 | Frameworks to share code, software dependencies and deep-learning models**

Resource	URL
<b>Code</b>	
BitBucket	<a href="https://bitbucket.org">https://bitbucket.org</a>
GitHub	<a href="https://github.com">https://github.com</a>
GitLab	<a href="https://about.gitlab.com">https://about.gitlab.com</a>
<b>Software dependencies</b>	
Conda	<a href="https://conda.io">https://conda.io</a>
Code Ocean	<a href="https://codeocean.com">https://codeocean.com</a>
Gigantum	<a href="https://gigantum.com">https://gigantum.com</a>
Colaboratory	<a href="https://colab.research.google.com">https://colab.research.google.com</a>
<b>Deep-learning models</b>	
TensorFlow Hub	<a href="https://www.tensorflow.org/hub">https://www.tensorflow.org/hub</a>
ModelHub	<a href="http://modelhub.ai">http://modelhub.ai</a>
ModelDepot	<a href="https://modeldepot.io">https://modeldepot.io</a>
Model Zoo	<a href="https://modelzoo.co">https://modelzoo.co</a>
<b>Deep-learning frameworks</b>	
TensorFlow	<a href="https://www.tensorflow.org/">https://www.tensorflow.org/</a>
Caffe	<a href="https://caffe.berkeleyvision.org/">https://caffe.berkeleyvision.org/</a>
PyTorch	<a href="https://pytorch.org/">https://pytorch.org/</a>

lives would be at stake, we argue that the bar of transparency should be set even higher. If a dataset cannot be shared with the entire scientific community, because of licensing or other insurmountable issues, at a minimum a mechanism should be set so that some highly-trained, independent investigators can access the data and verify the analyses.

The lack of access to code and data in prominent scientific publications may lead to unwarranted and even potentially harmful clinical trials<sup>14</sup>. These unfortunate lessons have not been lost on journal editors and their readers. Journals have an obligation to hold authors to the standards of reproducibility that benefit not only other researchers, but also the authors themselves. Making one's methods reproducible may surface biases or shortcomings to authors before publication<sup>5</sup>. Preventing external validation of a model will likely reduce its impact, as it also prevents other researchers from using and building upon it in future studies. The failure of McKinney et al. to share key materials and information transforms their work from a scientific publication open to verification and adoption by the scientific community into a promotion of a closed technology.

We have high hopes for the utility of AI methods in medicine. Ensuring that these methods meet their potential, however, requires that these studies be scientifically reproducible. The recent advances in computational virtualization and AI frameworks are greatly facilitating the implementations of complex deep neural networks in a more structured, transparent, and reproducible way. Adoption of these technologies will increase the impact of published deep-learning algorithms and accelerate the translation of these methods into clinical settings.

## Reporting summary

Further information on research design is available in the Nature Research Reporting Summary linked to this paper.

## Data availability

No data have been generated as part of this manuscript.

1. McKinney, S. M. et al. International evaluation of an AI system for breast cancer screening. *Nature* **577**, 89–94 (2020).
2. Bluemke, D. A. et al. Assessing radiology research on artificial intelligence: a brief guide for authors, reviewers, and readers—from the *Radiology* editorial board. *Radiology* **293**, 315–316 (2019).
3. Gundersen, O. E., Gil, Y. & Aha, D. W. On reproducible AI: towards reproducible research, open science, and digital scholarship in AI publications. *AI Mag.* **39**, 56–68 (2018).
4. Crane, M. Questionable answers in question answering research: reproducibility and variability of published results. *Trans. Assoc. Comput. Linguist.* **6**, 241–252 (2018).
5. Sculley, D. et al. in *Advances in Neural Information Processing Systems 28* (eds Cortes, C. et al.) 2503–2511 (Curran Associates, Inc., 2015).
6. Stodden, V. et al. Enhancing reproducibility for computational methods. *Science* **354**, 1240–1241 (2016).
7. Hutson, M. Artificial intelligence faces reproducibility crisis. *Science* **359**, 725–726 (2018).
8. Bzdok, D. & Ioannidis, J. P. A. Exploration, inference, and prediction in neuroscience and biomedicine. *Trends Neurosci.* **42**, 251–262 (2019).
9. Gundersen, O. E. & Kjensmo, S. State of the art: Reproducibility in artificial intelligence. In *Thirty-second AAAI Conference on Artificial Intelligence (AAAI-18)* 1644–1651 (2018).
10. Shorten, C. & Khoshgoftaar, T. M. A survey on Image Data Augmentation for Deep Learning. *J. Big Data* **6**, 60 (2019).
11. Lee, R. S. et al. A curated mammography data set for use in computer-aided detection and diagnosis research. *Sci. Data* **4**, 170177 (2017).
12. Wallach, J. D., Boyack, K. W. & Ioannidis, J. P. A. Reproducible research practices, transparency, and open access data in the biomedical literature, 2015–2017. *PLoS Biol.* **16**, e2006930 (2018).
13. Amann, R. I. et al. Toward unrestricted use of public genomic data. *Science* **363**, 350–352 (2019).
14. Carlson, B. Putting oncology patients at risk. *Biotechnol. Healthc.* **9**, 17–21 (2012).

**Acknowledgements** We thank S. McKinney and colleagues for their prompt and open communication regarding the materials and methods of their study. This work was supported in part by the National Cancer Institute (R01 CA237170).

**Author contributions** B.H.-K. and G.A.A. wrote the first draft of the manuscript. B.H.-K. and H.J.W.L.A. designed and supervised the study. A.H., F.K., T.S., R.K., S.-A.S., W.T., R.D.W.,

C.E.M., W.J., J.D., C.F., L.W., B.W., C. McIntosh, A.G., A.K., C.S.G., T.B., M.M.H., J.T.L., K.K., W.H., A.B., J.P., R.T., T.H., J.P.A.I. and J.Q. contributed to the writing of the manuscript.

**Competing interests** A.H. is a shareholder of and receives consulting fees from Altis Labs. M.M.H. received a GPU Grant from Nvidia. H.J.W.L.A. is a shareholder of and receives consulting fees from Onc.AI. B.H.K. is a scientific advisor for Altis Labs. C.M. holds an equity position in Bridge7Oncology and receives royalties from RaySearch Laboratories. A.K. is on the SAB of ImmuneAI Inc, a consultant for Biogen Inc., a scientific co-founder of RavelBio Inc. and a shareholder of Freenome Inc. G.A.A., F.K., L.W., B.W., C.S.G., J.T.L., W.H., A.B., J.P., R.T., T.H., J.P.A.I. and J.Q. declare no other competing interests related to the manuscript.

#### Additional information

**Supplementary information** is available for this paper at <https://doi.org/10.1038/s41586-020-2766-y>.

**Correspondence and requests for materials** should be addressed to B.H.-K.

**Reprints and permissions information** is available at <http://www.nature.com/reprints>.

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

© The Author(s), under exclusive licence to Springer Nature Limited 2020

#### Massive Analysis Quality Control (MAQC) Society Board of Directors

**Thakkar Shraddha<sup>35</sup>, Rebecca Kusko<sup>36</sup>, Susanna-Assunta Sansone<sup>37</sup>, Weida Tong<sup>35</sup>, Russ D. Wolfinger<sup>38</sup>, Christopher E. Mason<sup>39</sup>, Wendell Jones<sup>40</sup>, Joaquin Dopazo<sup>41</sup> & Cesare Furlanello<sup>42</sup>**

<sup>35</sup>National Center for Toxicological Research, US Food and Drug Administration, Jefferson, AR, USA. <sup>36</sup>Immuneering Corporation, Cambridge, MA, USA. <sup>37</sup>Engineering Science Department, Oxford e-Research Centre, University of Oxford, Oxford, UK. <sup>38</sup>SAS Institute Inc, Cary, NC, USA. <sup>39</sup>Weill Cornell Medicine, New York, NY, USA. <sup>40</sup>Q2 Solutions, Morrisville, NC, USA. <sup>41</sup>Hospital Virgen del Rocío, Sevilla, Spain. <sup>42</sup>Fondazione Bruno Kessler, Trento, Italy.

## Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

### Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

- |                                     |   |
|-------------------------------------|---|
| n/a                                 | Confirmed   |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> The exact sample size ( <i>n</i> ) for each experimental group/condition, given as a discrete number and unit of measurement   |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly  |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> The statistical test(s) used AND whether they are one- or two-sided<br><i>Only common tests should be described solely by name; describe more complex techniques in the Methods section.</i>   |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> A description of all covariates tested   |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons  |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals) |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> For null hypothesis testing, the test statistic (e.g. <i>F</i> , <i>t</i> , <i>r</i> ) with confidence intervals, effect sizes, degrees of freedom and <i>P</i> value noted<br><i>Give P values as exact values whenever suitable.</i>                     |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings   |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes   |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Estimates of effect sizes (e.g. Cohen's <i>d</i> , Pearson's <i>r</i> ), indicating how they were calculated   |

*Our web collection on [statistics for biologists](#) contains articles on many of the points above.*

### Software and code

Policy information about [availability of computer code](#)

Data collection

Data analysis

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research [guidelines for submitting code & software](#) for further information.

### Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

### Field-specific reporting

## Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	<input type="text" value="n/a"/>
Data exclusions	<input type="text" value="n/a"/>
Replication	<input type="text" value="n/a"/>
Randomization	<input type="text" value="n/a"/>
Blinding	<input type="text" value="n/a"/>

## Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

### Materials & experimental systems

n/a	Included in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> Antibodies
<input checked="" type="checkbox"/>	<input type="checkbox"/> Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/> Palaeontology and archaeology
<input checked="" type="checkbox"/>	<input type="checkbox"/> Animals and other organisms
<input checked="" type="checkbox"/>	<input type="checkbox"/> Human research participants
<input checked="" type="checkbox"/>	<input type="checkbox"/> Clinical data
<input checked="" type="checkbox"/>	<input type="checkbox"/> Dual use research of concern

### Methods

n/a	Included in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> ChIP-seq
<input checked="" type="checkbox"/>	<input type="checkbox"/> Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/> MRI-based neuroimaging

# Reply to: Transparency and reproducibility in artificial intelligence

<https://doi.org/10.1038/s41586-020-2767-x>

Published online: 14 October 2020

 Check for updates

Scott Mayer McKinney<sup>1</sup>✉, Alan Karthikesalingam<sup>2</sup>, Daniel Tse<sup>1</sup>, Christopher J. Kelly<sup>2</sup>, Yun Liu<sup>1</sup>, Greg S. Corrado<sup>1</sup> & Shrivya Shetty<sup>1</sup>✉

REPLYING TO B. Haibe-Kains et al. *Nature* <https://doi.org/10.1038/s41586-020-2766-y> (2020)

We thank the authors of the accompanying Comment<sup>1</sup> for their interest in our work<sup>2</sup> and their thoughtful contribution. We agree that transparency and reproducibility are paramount for scientific progress. In keeping with this principle, the largest data source used in our publication is available to the academic community. Any researcher can apply for access to the OPTIMAM database (<https://medphys.royalsurrey.nhs.uk/omidb/getting-access>), which our institution helped fund. The broad accessibility of the database was part of the reason we pursued this collaboration. In fact, since our article came out, another group has already published results on this very dataset<sup>3</sup>.

The other dataset, from the United States, was shared with our research team after approval from the hospital system's Institutional Review Board (IRB). The IRB judged that the potential benefits of the research outweighed the minimal privacy risks associated with sharing de-identified data with a trusted party capable of and committed to safeguarding these data. As the authors understand, we are not at liberty to share data that we do not own. More generally, widely releasing data considerably alters the risk–benefit calculus for patients, so institutions must be thoughtful about how and when they do this. Because of these considerations, large medical image datasets with associated breast cancer outcomes are rarely made openly available<sup>3–5</sup>. However, as our support for the OPTIMAM database demonstrates, we endorse such efforts where practical. Although there are some small, publicly available mammography datasets<sup>6</sup>, restricting published research to such datasets would provide an extremely limited picture of an algorithm's clinical applicability.

The commenters<sup>1</sup> asked for more information concerning the training of our deep learning models. We strove to document all relevant machine learning methods while keeping the paper accessible to a clinical and general scientific audience. We thank the authors for highlighting the omission of some hyperparameters. We have supplied the requested methodological details and further elaborated on our data augmentation strategies in an Addendum<sup>7</sup> to our original Article<sup>2</sup>.

The authors of the Comment<sup>1</sup> suggest open-sourcing all the code associated with this project. Most of our work builds on open-source implementations, such as ResNet ([https://github.com/tensorflow/models/blob/master/research/slim/nets/resnet\\_v1.py](https://github.com/tensorflow/models/blob/master/research/slim/nets/resnet_v1.py)), MobileNet ([https://github.com/tensorflow/models/blob/master/research/slim/nets/mobilenet/mobilenet\\_v2.py](https://github.com/tensorflow/models/blob/master/research/slim/nets/mobilenet/mobilenet_v2.py)), multidimensional image augmentation (<https://github.com/deepmind/multidim-image-augmentation>), and the Tensorflow Object Detection API ([https://github.com/tensorflow/models/tree/master/research/object\\_detection](https://github.com/tensorflow/models/tree/master/research/object_detection)), all of which were released by our institution. Much of the remaining code concerns data input–output and the orchestration of the training process across internal compute clusters, both of which are of scant scientific value and limited utility to researchers outside our organization. Given the

extensive textual description in the supplementary information of our Article<sup>2</sup>, we believe that investigators proficient in deep learning should be able to learn from and expand upon our approach.

The authors<sup>1</sup> further suggest releasing a containerized version of our model for others to apply to new images. It is important to note that regulators commonly classify technologies such as the one proposed here as ‘medical device software’ or ‘software as a medical device’. Unfortunately, the release of any medical device without appropriate regulatory oversight could lead to its misuse. As such, doing so would overlook material ethical concerns. Because liability issues surrounding artificial intelligence in healthcare remain unresolved<sup>8</sup>, providing unrestricted access to such technologies may place patients, providers, and developers at risk. In addition, the development of impactful medical technologies must remain a sustainable venture to promote a vibrant ecosystem that supports future innovation. Parallels to hardware medical devices and pharmaceuticals may be useful to consider in this regard. Finally, increasing evidence suggests that a model's learned parameters may inadvertently expose properties of its training set to attack; how to safeguard potentially susceptible models is the subject of active research<sup>9</sup>. As our training data are private or under restricted access, sharing the model openly seems premature and may introduce risks that are not well characterized. On the basis of these concerns, we deliberately approach sharing artefacts derived from patient data (even if de-identified) with an abundance of caution.

No doubt the commenters<sup>1</sup> are motivated by protecting future patients as much as scientific principle. We share that sentiment. This work serves as an initial proof of concept, and is by no means the end of the story. We intend to subject our software to extensive testing before its use in a clinical environment, working alongside patients, providers and regulators to ensure efficacy and safety.

1. Haibe-Kains, B. et al. Transparency and reproducibility in artificial intelligence. *Nature* <https://doi.org/10.1038/s41586-020-2766-y> (2020).
2. McKinney, S. M. et al. International evaluation of an AI system for breast cancer screening. *Nature* **577**, 89–94 (2020).
3. Kim, H.-E. et al. Changes in cancer detection and false-positive recall in mammography using artificial intelligence: a retrospective, multireader study. *Lancet Digital Health* **2**, e138–e148 (2020).
4. Wu, N. et al. Deep neural networks improve radiologists' performance in breast cancer screening. *IEEE Trans. Med. Imaging* **39**, 1184–1194 (2019).
5. Rodriguez-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).
6. Lee, R. S. et al. A curated mammography data set for use in computer-aided detection and diagnosis research. *Sci. Data* **4**, 170177 (2017).
7. McKinney, S. M. et al. Addendum: International evaluation of an AI system for breast cancer screening. *Nature* <https://doi.org/10.1038/s41586-020-2679-9> (2020).
8. Price, W. N., II, Gerke, S. & Cohen, I. G. Potential liability for physicians using artificial intelligence. *J. Am. Med. Assoc.* **322**, 1765–1766 (2019).
9. Abadi, M. et al. Deep learning with differential privacy. In *Proc. 2016 ACM SIGSAC Conference Computer Communications Security CCS'16* 308–318 (2016).

<sup>1</sup>Google Health, Palo Alto, CA, USA. <sup>2</sup>Google Health, London, UK. ✉e-mail: scottmayer@google.com; sshetty@google.com

# Matters arising

---

**Acknowledgements** We thank A. Dai and E. Gabrilovich for comments.

**Author contributions** This Reply was prepared by a subset of the authors of the original Article in addition to Y.L., all of whom have expertise related to this exchange. S.M.M., A.K., D.T., C.J.K., Y.L., G.S.C. and S.S. wrote and revised this Reply.

**Competing interests** This study was funded by Google LLC. S.M.M., A.K., D.T., C.J.K., Y.L., G.S.C. and S.S. are employees of Google and own stock as part of the standard compensation package. The authors have no other competing interests to disclose.

## Additional information

**Correspondence and requests for materials** should be addressed to S.M.M. or S.S.

**Reprints and permissions information** is available at <http://www.nature.com/reprints>.

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

© The Author(s), under exclusive licence to Springer Nature Limited 2020