

# The importance of transparency and reproducibility in artificial intelligence research

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>, MAQC Society Board<sup>8,\*</sup>, Levi Waldron<sup>9</sup>, Bo Wang<sup>2,5,10</sup>, Chris McIntosh<sup>2,5,10</sup>, Anshul Kundaje<sup>11</sup>, Casey S. Greene<sup>12,13</sup>, Michael M. Hoffman<sup>1,2,3</sup>, Jeffrey T. Leek<sup>14</sup>, Wolfgang Huber<sup>15</sup>, Alvis Brazma<sup>16</sup>, Joelle Pineau<sup>17,18</sup>, Robert Tibshirani<sup>19,20</sup>, Trevor Hastie<sup>19,20</sup>, John P.A. Ioannidis<sup>19,20,21,22</sup>, John Quackenbush<sup>24,25,26</sup>, Hugo J.W.L. Aerts<sup>6,7,27,20</sup>

## Affiliations

<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

<sup>8</sup> Massive Analysis Quality Control (MAQC) Society Board of Directors

<sup>9</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>10</sup> Peter Munk Cardiac Centre, University Health Network, Toronto, Ontario, Canada  
Hospital, Harvard Medical School, Boston, MA, USA

<sup>11</sup> Department of Genetics, Stanford University School of Medicine, Stanford, CA, USA

<sup>12</sup> Dept. of Systems Pharmacology and Translational Therapeutics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, USA

<sup>13</sup> Childhood Cancer Data Lab, Alex's Lemonade Stand Foundation, Philadelphia, PA, USA

<sup>14</sup> Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health, Baltimore MD, USA

<sup>15</sup> European Molecular Biology Laboratory, Genome Biology Unit, Heidelberg, Germany

<sup>16</sup> European Molecular Biology Laboratory, European Bioinformatics Institute, EMBL-EBI, Hinxton, UK

<sup>17</sup> McGill University, Montreal, QC, Canada

<sup>18</sup> Montreal Institute for Learning Algorithms, QC, Canada

<sup>19</sup> Meta-Research Innovation Center at Stanford (METRICS), Stanford, CA, USA

<sup>20</sup> Department of Biomedical Data Science, Stanford University School of Medicine, Stanford, CA, USA

<sup>21</sup> Departments of Medicine, of Health Research and Policy, and of Biomedical Data Science, Stanford University School of Medicine, Stanford, CA, USA

<sup>22</sup> Department of Statistics, Stanford University School of Humanities and Sciences, Stanford, CA, USA

<sup>23</sup> Department of Epidemiology and Population Health, Stanford University School of Medicine, Stanford, CA, USA

<sup>24</sup> Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA, USA

<sup>25</sup> Channing Division of Network Medicine, Brigham and Women's Hospital, Boston, MA, USA

<sup>26</sup> Department of Data Science, Dana-Farber Cancer Institute, Boston, MA, USA

<sup>27</sup> Radiology and Nuclear Medicine, Maastricht University, Maastricht, Netherlands

<sup>28</sup> Cardiovascular Imaging Research Center, Massachusetts General Hospital, Harvard Medical School, Boston, MA, USA

**§ Corresponding Author**

Benjamin Haibe-Kains: [bhaibeka@uhnresearch.ca](mailto:bhaibeka@uhnresearch.ca)

**\* Massive Analysis Quality Control (MAQC) Society Board of Directors**

Thakkar Shraddha, Rebecca Kusko, Susanna-Assunta Sansone, Weida Tong, Russ D. Wolfinger, Christopher Mason, Wendell Jones, Joaquin Dopazo, Cesare Furlanello

**Abstract:**

In their study, McKinney et al. showed the high potential of artificial intelligence for breast cancer screening. However, the lack of detailed methods and computer code undermines its scientific value. We identify obstacles hindering transparent and reproducible AI research as faced by McKinney et al and provide solutions with implications for the broader field.

**Main Text:**

The evaluation of deep learning for the detection of breast cancer from mammograms by McKinney and colleagues<sup>1</sup> showed promising improvements in screening performance, while highlighting challenges around the reproducibility and transparency of artificial intelligence (AI) research. They 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 validation, 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. Here, we identify obstacles hindering transparent and reproducible AI research as faced by McKinney et al. and provide potential solutions with implications for the broader field.

Scientific progress depends upon the ability of independent researchers to (1) scrutinize the results of a research study, (2) reproduce the study's main results using its materials, and (3) build upon them in future studies<sup>2</sup>. Publication of insufficiently documented research violates the core principles underlying scientific discovery<sup>3,4</sup>. The authors state *"The code used for training the models has a large number of dependencies on internal tooling, infrastructure and hardware, and its release is therefore not feasible"*. Computational reproducibility is indispensable for robust AI applications<sup>5,6</sup> more complex methods demand greater transparency<sup>7</sup>. In the absence of code, reproducibility falls back on replicating methods from textual description. Although, the authors claim that *"all experiments and implementation details are described in sufficient detail in the Supplementary Methods section to support replication with non-proprietary libraries"*, key details about their analysis are lacking. Even with sufficient description, reproducing complex computational pipelines based purely on text is a subjective and challenging task<sup>8,9</sup>.

More specifically, the authors' description of the model development as well as data processing and training pipelines lacks critical details. The definition of multiple hyperparameters for the model's architecture (composed of three networks referred to as the Breast, Lesion, and Case models) is missing (Table 1). The authors did not disclose the parameters used for data augmentation; the transformations used are stochastic and can significantly affect model performance<sup>10</sup>. Details of the training pipeline were also missing. For instance, they state that the mini-batches were sampled to contain an equal proportion of negative and positive examples, potentially leading to multiple instances of the same patients in a given epoch. Deep learning optimization algorithms such as stochastic gradient descent typically operate under the

assumption that a given sample is provided to the model exactly once per epoch. The lack of detail regarding the per-batch balancing of classes prevents replicating the training pipeline.

There exist numerous frameworks and platforms 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 multiple software dependencies of large-scale machine learning applications require appropriate control of software environment, which can be achieved through package managers including Conda, as well as container and virtualization systems, including Code Ocean, Gigantum and Colaboratory. If virtualization of the McKinney et al. internal tooling proved to be difficult, they could have released the computer code and documentation. The authors could also have created toy examples to show how new data must be processed to generate predictions. As for the trained model, many platforms allow sharing of deep learning models, including TensorFlow Hub, ModelHub.ai, ModelDepot, and Model Zoo with support for multiple frameworks such as PyTorch, Caffe, and TensorFlow. Since the authors created their model with the publicly available TensorFlow library, sharing of the model should be trivial. In addition to improving accessibility and transparency, such tools can significantly 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. used two large datasets under license, properly disclosing this limitation in their publication. Sharing of patient health information is highly regulated due to privacy concerns. Despite these challenges, sharing of raw data has become more common in biomedical literature, increasing from under 1% in the early 2000s to 20% today<sup>11</sup>. However, if the data cannot be shared, the model predictions and data labels themselves should be released, allowing further statistical analyses.

Although sharing of code and data is 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>12</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 where a clinical application is envisioned and human lives would be at stake, we argue that the bar of transparency should be set even higher. If data 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. This would allow a truly adequate peer-review of the study and its evidence before moving into clinical implementation.

Many egregious failures of science were due to lack of public access to code and data used in the discovery process<sup>13,14</sup>. These unfortunate lessons should not be lost on either journal editors or its readers. Journals have an obligation to hold authors to the standards of reproducibility that

benefit not only other researchers, but also the creators of a given method. Making one's methods reproducible may surface biases or shortcomings to authors before publication<sup>15</sup>. Preventing external validation of a model will likely reduce its impact and could lead to unintended consequences<sup>15</sup>. The failure of McKinney et al. to share key materials and information transforms their work from a scientific publication open to verification 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 reproducible. Unfortunately, the biomedical literature is littered with studies that have failed the test of reproducibility, and many of these can be tied to methodologies and experimental practices that could not be investigated due to failure to fully disclose software and data. This is even more important for applications intended for use in the diagnosis or treatment of human disease.

### **Competing Interests**

AH is a shareholder of and receives consulting fees from Altis Labs. MMH received a GPU Grant from Nvidia. HJWLA is a shareholder of and receives consulting fees from Onc.AI. BHK is a scientific advisor for Altis Labs. GAA, FK, LW, BW, CM, AK, CSG, JTL, WH, AB, JP, RT, TH, JPAI and JQ declare no other competing interests related to the manuscript.

### **Author Contributions**

BHK and GAA wrote the first draft of the manuscript. BHK and HJWLA designed and supervised the study. AH, FK, LW, BW, CM, AK, CSG, MMH, JTL, WH, AB, JP, RT, TH, JPAI and JQ contributed to the writing of the manuscript.

### **References**

1. McKinney, S. M., Sieniek, M., Godbole, V. & Godwin, J. International evaluation of an AI system for breast cancer screening. *Nature* (2020).
2. Nature Research Editorial Policies. Reporting standards and availability of data, materials, code and protocols. *Springer Nature*  
<https://www.nature.com/nature-research/editorial-policies/reporting-standards>.
3. 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* 192515 (2019) doi:10.1148/radiol.2019192515.
4. Gundersen, O. E., Gil, Y. & Aha, D. W. On reproducible AI: Towards reproducible research, open science, and digital scholarship in AI publications. *AI Magazine* **39**, 56–68 (2018).
5. Stodden, V. *et al.* Enhancing reproducibility for computational methods. *Science* **354**, 1240–1241 (2016).
6. Hutson, M. Artificial intelligence faces reproducibility crisis. *Science* **359**, 725–726 (2018).
7. Bzdok, D. & Ioannidis, J. P. A. Exploration, Inference, and Prediction in Neuroscience and Biomedicine. *Trends Neurosci.* **42**, 251–262 (2019).
8. Gundersen, O. E. & Kjensmo, S. State of the art: Reproducibility in artificial intelligence. in

*Thirty-second AAAI conference on artificial intelligence* (2018).

9. Beaulieu-Jones, B. K. & Greene, C. S. Reproducibility of computational workflows is automated using continuous analysis. *Nat. Biotechnol.* **35**, 342–346 (2017).
10. Shorten, C. & Khoshgoftaar, T. M. A survey on Image Data Augmentation for Deep Learning. *Journal of Big Data* **6**, 60 (2019).
11. 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).
12. Amann, R. I. *et al.* Toward unrestricted use of public genomic data. *Science* **363**, 350–352 (2019).
13. Iqbal, S. A., Wallach, J. D., Khoury, M. J., Schully, S. D. & Ioannidis, J. P. A. Reproducible Research Practices and Transparency across the Biomedical Literature. *PLoS Biol.* **14**, e1002333 (2016).
14. Carlson, B. Putting oncology patients at risk. *Biotechnol. Healthc.* **9**, 17–21 (2012).
15. Sculley, D. *et al.* Hidden Technical Debt in Machine Learning Systems. in *Advances in Neural Information Processing Systems 28* (eds. Cortes, C., Lawrence, N. D., Lee, D. D., Sugiyama, M. & Garnett, R.) 2503–2511 (Curran Associates, Inc., 2015).

**Table 1: Essential hyperparameters for reproducing the study for each of the three models (Lesion, Breast, and Case), including those missing from the description in Mckinney et al.**

	<b>Lesion</b>	<b>Breast</b>	<b>Case</b>
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

**Table 2: Frameworks and platforms to share code, software dependencies and deep learning models to make artificial intelligence research more transparent and reproducible.**

Resource		URL
<i>Code</i>		
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>
<i>Software dependencies</i>		
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>
<i>Deep learning models</i>		
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>
<i>Deep learning frameworks</i>		
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>



<b>Author Name</b>	<b>ORCID</b>
Benjamin Haibe-Kains	0000-0002-7684-0079
George-Alexandru Adam	0000-0001-9084-3703
Ahmed Hosny	0000-0002-1844-481X
Levi Waldron	0000-0003-2725-0694
Bo Wang	/0000-0002-9620-3413
Chris McIntosh	0000-0003-1371-1250
Anshul Kundaje	0000-0003-3084-2287
Casey S. Greene	0000-0001-8713-9213
Michael M. Hoffman	0000-0002-4517-1562
Jeffrey T. Leek	0000-0002-2873-2671
Wolfgang Huber	0000-0002-0474-2218
Alvis Brazma	0000-0001-5988-7409
Joelle Pineau	0000-0003-0747-7250
Robert Tibshirani	0000-0003-0553-5090
Trevor Hastie	0000-0002-0164-3142
John P. A. Ioannidis	0000-0003-3118-6859
John Quackenbush	0000-0002-2702-5879
Hugo JWL Aerts	0000-0002-2122-2003

<b>MAQC Society Board Members</b>	<b>ORCID</b>
Thakkar Shraddha	0000-0002-2920-7713
Rebecca Kusko	0000-0001-5331-5119

Susanna-Assunta Sansone	0000-0001-5306-5690
Weida Tong	0000-0003-3488-6148
Russ D. Wolfinger	0000-0001-8575-0537
Christopher Mason	0000-0002-1850-1642
Wendell Jones	0000-0002-9676-5387
Joaquin Dopazo	0000-0003-3318-120X
Cesare Furlanello	0000-0002-5384-3605