

1
2 **Artificial Intelligence is Changing the Ethics of Medicine:**
3 **Reflections from the Australian Epilepsy Project**
4

5 Mangor Pedersen^{1,2*}, Heath R. Pardoe², Anton de Weger², Donna Hutchison²,
6 David F. Abbott^{2,4}, Karin Verspoor³ & Graeme D. Jackson^{2,4,5}
7

- 8 1. Department of Psychology and Neuroscience, Auckland University of Technology
9 (AUT), Auckland, New Zealand
10 2. The Florey Institute of Neuroscience and Mental Health, The University of
11 Melbourne, Australia.
12 3. School of Computing Technologies, RMIT University, Melbourne, Australia
13 4. Department of Medicine, Austin Health, The University of Melbourne, Australia
14 5. Department of Neurology, Austin Health, Melbourne, Australia.
15

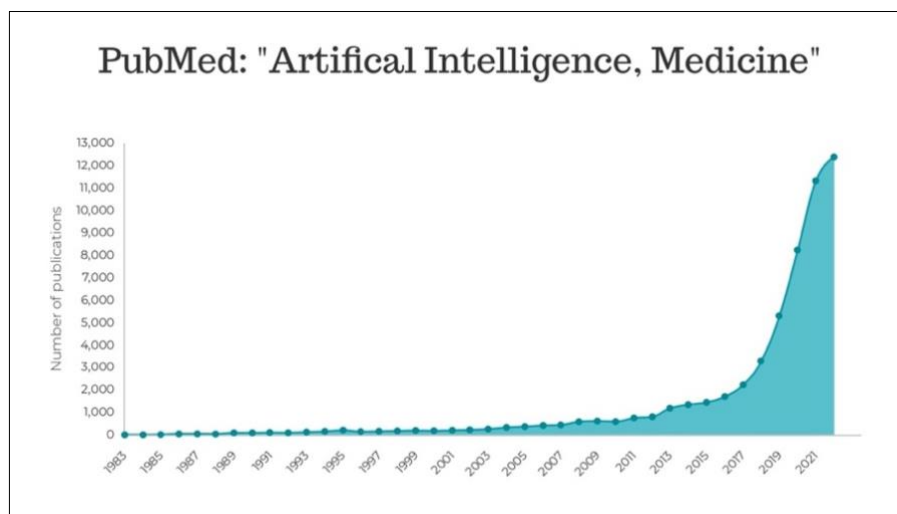
16 *Corresponding author: mangor.pedersen@aut.ac.nz
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33

1 **Abstract**

2 Artificial intelligence (AI) is a multidisciplinary scientific field that uses machines to solve
3 real-world problems and predict outcomes. Despite the current enthusiasm about AI's
4 potential as a clinical support tool, there is also a growing awareness and concern about the
5 potentially harmful effects of AI. Because AI will likely impact expert-based decision-
6 making in medicine, it is critical to consider the ethical issues that AI raises in medical
7 research. This paper outlines the AI ethics guidelines of the Australian Epilepsy Project. This
8 large-scale platform aims to democratise specialist care in epilepsy and use AI for clinical
9 decision support based on prospective multimodal datasets (MRI, genetic, clinical, and
10 cognitive data) from thousands of people with epilepsy. As AI develops rapidly, we focus on
11 key areas of medical AI ethics previously identified in the literature, including *Transparency*,
12 *Justice and Fairness*, *Non-maleficence*, *Responsibility*, and *Sustainability*. We believe AI is
13 changing the ethics of medicine, and it is imperative to advance and update ethical guidelines
14 adaptably while preparing for an era of augmented-intelligence-based medicine.

1 Background

2 Artificial intelligence (AI) uses computing storage, processing, and knowledge in
3 combination with interactions with environments to understand human intelligence ¹. AI is an
4 exciting, rapidly evolving, and multidisciplinary field of science focused on problem-solving
5 by machines. AI will most likely substantially impact expert-based decision-making in
6 medicine and increasingly realise the challenge of deep phenotyping, which means
7 understanding each individual's uniqueness through detailed characterization based on
8 multimodal data ². To process complex data meaningfully, clinicians and AI scientists will
9 benefit from working together to build and validate new models and demonstrate their utility
10 to improve health outcomes ³.



11
12 *Figure 1: In this figure, we display the number of publications in PubMed*
13 *(<https://pubmed.ncbi.nlm.nih.gov/>) for the keywords: Artificial Intelligence, Medicine. It*
14 *shows a rapid increase in medical publications referring to AI from 2017, at a time of*
15 *significant scaling of AI models and the computational feasibility of deep learning models.*
16 *Note that this is a rough approach to estimating the number of studies, not including papers*
17 *relating to AI (e.g., using terminologies such as machine learning and natural language*
18 *processing).*

19

20 Amongst the contemporary excitement around AI's potential, as well as its commercial
21 success (the expected worth of the medical AI market is projected to be US\$41.7 billion
22 before 2027*) and research success (see the recent rise in medical AI publications in Figure

* https://www.reportlinker.com/p06325450/Global-AI-in-Healthcare-Market-Analysis-By-Component-By-Algorithm-By-Application-By-End-User-By-Region-Size-Forecast-with-Impact-Analysis-of-COVID-19-and-Forecast-up-to.html?utm_source=GNW

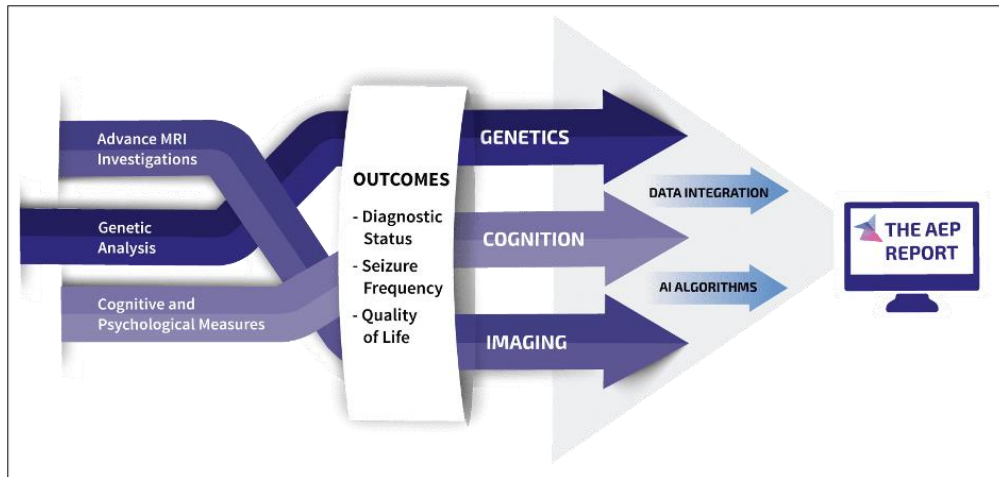
1) 1), there is an increasing awareness and apprehension about the potentially harmful effects of
2 AI. As new developments in AI are continuously emerging, there are concerns about the
3 potential detrimental impacts of AI: mainly how AI solutions are developed, shared,
4 monitored, and regulated in medicine ⁴⁻¹⁶.

5
6 Given that AI will lead to a change in medical practice, ethical guidelines need to advance to
7 consider ethical risks when developing, implementing, and governing AI algorithms in
8 clinical settings. Next, we outline our large-scale Australian Epilepsy Project that intends to
9 implement novel AI solutions for thousands of people living with epilepsy in Australia.

10 11 **The Australian Epilepsy Project and our approach to AI Ethics**

12 Epilepsy is more than just seizures. Reduced educational and occupational attainment,
13 decreased independence, lower quality of life, stigma, a higher risk of injury, cognitive
14 deficiencies, mental illness, and suicide are all connected with it ¹⁷⁻²². *The Australian*
15 *Epilepsy Project* ²³⁻²⁶ aims to give epilepsy specialists the information they need to make a
16 difference in people's lives, facilitate access to currently limited care, and answer questions
17 that are presently unanswerable. Will this person continue to suffer seizures? Which drug
18 works the best? Can brain surgery be a successful solution? These are fundamental concerns
19 that current medical practice frequently cannot answer.

20
21 The Australian Epilepsy Project has several aims that AI can address by integrating
22 multimodal data, including advanced brain imaging, cognitive, and genetics data (see Figure
23 2). The clinical purposes of the Australian Epilepsy Project include predicting whether a
24 person who has recently suffered their first seizure is likely to experience a second seizure
25 (epilepsy is diagnosed after more than one seizure). After people have been diagnosed with
26 epilepsy, more critical clinical decisions exist. Typically, the first decision for the treating
27 physician following diagnosis is the most appropriate antiseizure medication. In this context,
28 AI could help decide which antiseizure medication is most useful for each patient by
29 providing an individual-specific treatment plan. There is also a subset of patients where
30 seizures arise from a circumscribed part of the brain, and in one third of people with epilepsy
31 medication will not effectively control their seizures ²⁷. In these people, AI and advanced
32 MRI images can be used to identify subtle epileptogenic lesions not otherwise seen in
33 hospital-based brain imaging.



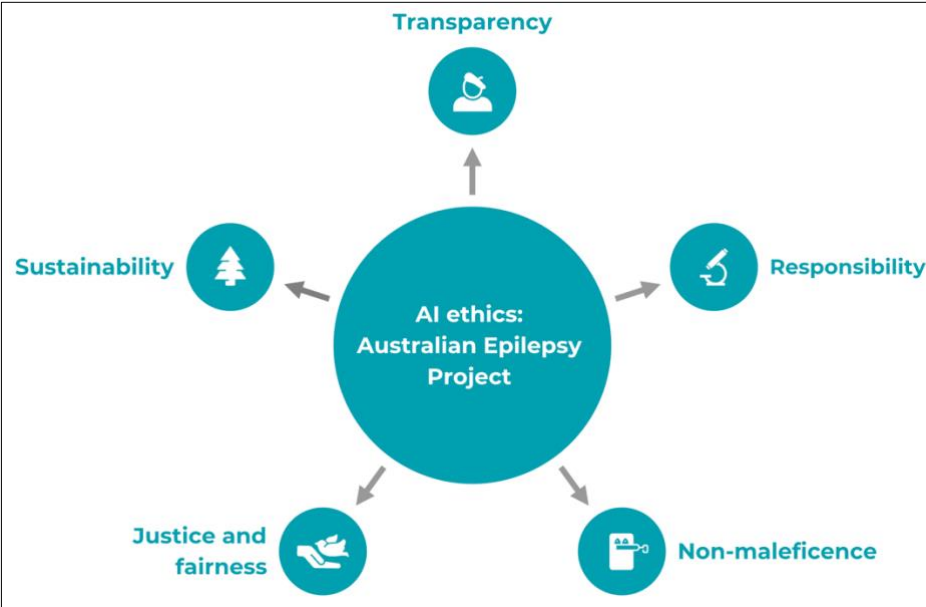
1
 2 *Figure 2: In this figure, we outline the structure of the Australian Epilepsy project, including*
 3 *the data collection and long-term outcomes integral to using AI to improve the lives of people*
 4 *with epilepsy.*

5
 6 A concerning trend within medical AI, identified in a study by Wu et al. ²⁸, highlights the
 7 need for large-scale clinical projects like the Australian Epilepsy Project. The authors
 8 reviewed the number of FDA-approved AI medical devices between 2017 and 2020. They
 9 found that the number of FDA-approved AI devices increased sharply in this period.
 10 However, the median sample size of the studies decreased over time, and single-site data was
 11 more common than multi-site studies. Most worrying was that most studies were also almost
 12 always retrospective. Obtaining quality annotated clinical information and outcome data is
 13 challenging. Nagendran et al. ²⁹ also acknowledges that few prospective AI studies exist in
 14 medical imaging, reinforcing the idea that the field of AI should move towards larger and
 15 prospective studies to enable reliable and generalizable science.

16
 17 The Australian Epilepsy Project is a *prospective* and *clinically-led* study with a target sample
 18 size of *thousands of people* living with epilepsy from *all over Australia*. Participants in the
 19 Australian Epilepsy Project will also be tracked two years after study inclusion to provide
 20 reliable and clinically valid outcome data for AI classification, such as medical and surgical
 21 outcomes. In addition to using AI to create predictions directly related to participant
 22 outcomes, the Australian Epilepsy Project intends to use AI in many aspects of the team’s
 23 workflow and planning. This could include using natural language processing to extract key
 24 terms from medical history and radiologist reports, estimating the “Failure to Attend”
 25 likelihood to help optimize bookings and minimize dropouts, and rating case similarity
 26 between participants to detect trends. Even if AI models are not implemented in production

1 workflows, data analysis using AI techniques may inform decision-making processes. This
2 use of pattern matching, analytics, and machine learning techniques will be undertaken with
3 due consideration of the ethical AI principles we discuss below. All data used to support
4 project workflow or additional research will be treated and secured like the core study data.

5
6 To develop an ethical AI framework for the Australian Epilepsy Project, we have used a data-
7 driven approach, using findings from a systematic review conducted by Jobin et al. ³⁰. Jobin
8 and others summarised information from 84 publications that covered AI ethics and showed
9 convergence around the following topics in AI ethics, *transparency, justice and fairness,*
10 *non-maleficence, and responsibility,* which we will review next.



12
13 *Figure 3: In this figure, we outline the five ethical AI domains that are important to the*
14 *Australian Epilepsy Project, all of which we take conscious steps to improve iteratively as the*
15 *field of AI advances.*

1 **Transparency:**

2 Transparency was deemed a critical AI ethical topic in 73/84 (87%) previous publications
3 with keywords such as explainability, interpretability, communication, and disclosure ³¹.

4
5 AI-based clinical decisions should be validated, explainable, and reproducible ³²⁻³⁴. In other
6 words, how does an algorithm reach its conclusion and prediction? Understanding how AI
7 systems arrive at their decisions or recommendations likely requires a shift into a
8 multidisciplinary area of inquiry. Many AI models, including deep learning, are often
9 conceived as a ‘black box’, and their outputs or results can be challenging to interpret ³⁵.

10

11 In epilepsy, identification of epileptogenic lesions increases the likelihood of obtaining a cure
12 for many people with focal epilepsy via brain surgery ³⁶. Contemporary AI solutions such as
13 deep learning have been explicitly developed for image-based applications and are well
14 suited for detecting neuroanatomical abnormalities ^{37,38}. An example of a transparent and
15 explainable lesion detection framework is proposed by Spitzer et al. ³⁹, namely the Multi-
16 centre Epilepsy Lesion Detection (MELD) framework. The MELD framework relies on
17 detailed structural information about the brain to obtain a probability of a seizure focus based
18 on structural brain imaging modalities commonly acquired in hospitals and research
19 institutions. In addition to lesion probability, MELD provides interpretable outputs regarding
20 the brain's cortical thickness, curvature, and white and grey matter imaging contrast
21 explaining the model’s output. Spitzer et al. ³⁹ validated this approach with 1015 participants
22 (epilepsy and controls) using a split-half training and testing paradigm. Patients with a visible
23 brain lesion –focal cortical dysplasia– were detected with 85% accuracy. The sensitivity and
24 specificity of the MELD approach were lower when the structural features of the seizure
25 focus was unclear. However, we anticipate rapid developments in lesion detection in epilepsy
26 over the next few years, given that AI and deep learning models are ideally suited to detect
27 abnormalities in images. Large-scale neuroimaging data collection efforts such as the
28 Australian Epilepsy Project will be critical to achieving this goal.

29

30 Another area in need of improvement in AI is research reporting. Several reporting checklists
31 encourage detailed reporting of AI methodologies, facilitating transparency and replicability.
32 One of the most common reporting frameworks in the field is the Transparent Reporting of
33 studies on prediction models for Individual Prognosis Or Diagnosis (TRIPOD ⁴⁰). A
34 systematic review by Nagendram et al. ²⁹ summarised adherence to the TRIPOD framework

1 in AI studies. The authors found that in most studies reviewed, less than half of the data items
2 in the TRIPOD framework were reported in peer-reviewed publications. These results show
3 that reporting AI standards are inadequate, but a potential issue with the TRIPOD reporting is
4 that it was made in an ‘older era’ of prediction modelling. The team is working on a
5 TRIPOD-AI checklist for newer AI models ⁴¹. In the meantime, other AI-specific reporting
6 checklists have also been created ⁴²⁻⁴⁴. One of the newer checklists is the Data, Optimization,
7 Model, and Evaluation (DOME) framework ⁴⁵. The DOME framework is developed
8 explicitly with contemporary AI and prediction models in mind. It includes a comprehensive
9 list of items in AI research design, train/test/validation dataset split, and model parameter
10 selection. Datasheets documenting key characteristics of datasets can improve transparency
11 and accountability of models built on top ⁴⁶.

12

13 In addition to reporting checklists, registering publications is a way to increase transparency
14 in research. In a registered study, a manuscript submitted in two stages (before and after
15 results are obtained) is pre-accepted for journal publication ⁴⁷. Emerging evidence suggests
16 that study bias is reduced in pre-registered reports. Evidence indicates that five times more
17 studies report no statistically significant findings if manuscripts are registered compared to
18 manuscripts that are not registered. This practice also allows for the publication of negative
19 findings ⁴⁸, providing a more reliable and comprehensive picture of relevant evidence. This is
20 in line with the benefits of clinical trial registration ⁴⁹ and peer-reviewed publication of
21 systematic review protocols via e.g., PROSPERO ⁵⁰ or Cochrane
22 (<https://www.cochrane.org/>), including to avoid selective outcome reporting ^{51,52}.

23

24 Non-transparent models are a concerning trend amongst large companies. For example, at the
25 time of writing, Open-AI’s large language model, GPT, has not been released for open
26 scientific exploration, nor have adequate details been provided about the resources and
27 models on which it is based to enable meaningful evaluation. This contradicts the idea of safe
28 and transparent AI. There is already substantial fear of AI in the community ⁵³. ⁵⁴ has
29 emphasized the critical need for developing adaptive guidelines that govern how AI
30 technologies are used and developed. Legal frameworks will likely struggle to keep up with
31 the pace of AI development. The European Union (EU) has proposed the 'AI-act' as a legal
32 AI framework, as a unified regulatory and legal framework for AI. It will likely change the
33 landscape of AI, ensuring that AI serves people and is a positive contribution to society
34 (<https://digital-strategy.ec.europa.eu/en/policies/european-approach-artificial-intelligence>).

1 **Justice and fairness:**

2 Justice and fairness were identified as essential topics in 68/84 (81%) publications in Jobin et
3 al.³⁰ review on ethical AI, with keywords including bias, discrimination, diversity, and
4 accessibility .

5

6 Study bias is a common and intrinsic issue in research studies. For instance, clinical trial
7 participation is generally lower among women, people of colour, and older people⁵⁶.

8 Genome-wide association studies have a strong Eurocentric bias, with 79% of GWAS

9 participants of European descent. Polygenic risk scores derived from this data have far lower
10 predictive value than non-European populations⁵⁷. Even in mega-sized community studies

11 such as the UK biobank project, there is evidence of study bias. The UK biobank study aims

12 to provide multimodal data from 500,000 people representing a ‘snapshot’ of the aging

13 population in the UK between 49-70 years. Nevertheless, in the UK biobank study, the

14 subjects are generally older, more often female, and have a higher socioeconomic status than

15 the general population in the UK⁵⁸. Subjects in the UK biobank study were also less likely to

16 be obese, to smoke, to drink alcohol daily, and to have fewer self-reported health conditions

17 when compared to the general population. The UK Biobank is an excellent initiative and an

18 invaluable resource for extensive population research. Still, given the ‘volunteer bias’

19 evidence, it may not represent the sampling population.

20

21 Biases in data acquisition also mean that AI algorithms will inevitably be biased. Medical AI

22 exhibits bias, too, often with insufficient generalisation and skewed samples regarding

23 gender, age, racial and ethnic backgrounds, hospitals, and methods related to data collection

24⁵⁹. Wang and colleagues demonstrated reduced bias if MRI-based AI models include

25 multisource data, including demographic, clinical, genetic, and cognitive scores. The U.S.

26 Food and Drug Administration (FDA) recently presented a draft document requiring overt

27 plans for having diversity in clinical trials (<https://www.fda.gov/media/106965/download>)

28 which signifies a step toward less biased and community-representative studies. This FDA

29 document outlines how a researcher should describe and explain the rationale behind the

30 anticipated enrolment of individuals from underrepresented groups.

31

32 The Australian population is ethnically diverse, which will likely benefit the generalisability

33 of AI models derived from the Australian Epilepsy Project. The project will also collect

34 demographic, gender, race, and other data often associated with biases in large datasets. It

1 will use this information to reduce potential acquisition bias when training and using AI
2 models. A priority for the Australian Epilepsy Project is to democratise clinical access for
3 people with epilepsy. This means bringing people to research centers with novel technologies
4 and capabilities, such as the ability to acquire advanced imaging including high-resolution
5 neuroanatomy, functional MRI and diffusion imaging scans to image the brain's white matter
6 pathways. Other data collection can now be done in people's homes, including saliva samples
7 for genetic tests which can be sent by mail, and tele-neuropsychology, which is delivered
8 online ²⁶.

9

10 **Non-maleficence:**

11 Non-maleficence was considered a topic of interest in 60/84 (71%) of previous studies ³⁰.
12 Keywords related to non-maleficence were security, safety, harm, and protection ⁶⁰.

13

14 A recent example of why governance and ethical regulations are needed in AI comes from a
15 recent report from Urbina et al. ⁶¹. The authors, part of a company developing drug
16 compounds to combat human disease, were invited to a conference hosted by the Swiss
17 government. This conference scrutinizes potential harm and adverse outcomes from AI
18 technologies. Based on this conference, Urbina et al. inverted their original 'good' AI
19 algorithm using chemical compounds to a new model designed to do 'harm'. By transforming
20 the positive algorithm into a harmful one, the authors created 40,000 compounds that could
21 be used in chemical warfare in less than six hours of data training. Many compounds were
22 deemed potentially more dangerous than current warfare agents. The authors acknowledge
23 that they had never previously considered that their AI models with chemical agents could be
24 used to cause harm. This example is an extreme scenario, as working with chemical
25 compounds is potentially risky. Nevertheless, this study highlights the need to consider the
26 known and potentially unknown ethical implications of AI and modelling, especially in
27 medicine, as these AI models influence people's lives and well-being.

28

29 A recent open letter signed by some key opinion leaders called for a six-month pause on giant
30 AI experiments that are yielding "unpredictable black-box models with emergent
31 capabilities", urging AI developers to "work with policymakers and regulators to
32 dramatically accelerate the development of robust AI governance systems" (see
33 <https://futureoflife.org/open-letter/pause-giant-ai-experiments/>). While the AI developments
34 planned by the Australian Epilepsy Project are more limited in scope, the need for oversight

1 of the use of resultant AI models is well recognized, including the need for appropriate
2 education of clinicians wishing to take advantage of AI-powered decision support tools.
3 Sharing of raw data and AI models are a double edge sword when it comes to AI ethics. It
4 enables transparency, which is good. But it also increases data safety and protection risks,
5 especially disclosing people's identities. Although most people are happy for their data to be
6 publicly shared ⁶², strategies are needed to ensure safe data sharing, as the re-identification of
7 individuals is a risk when openly releasing data. For example, if a person's whole genome is
8 freely available, re-identifying their identity is possible ⁶³. It is also possible to re-identify
9 individuals based on MRI of people's faces. Without de-facing MRI images, automated face
10 recognition was 97% accurate in matching participants' real photographs to the correct MRI,
11 and software packages exist to remove MRI features of people's faces before publicly
12 releasing data ⁶⁴.

13

14 Another way to safely release data to the public is to generate synthetic data. This is because
15 synthetic data has no identity and can cause similar features to the original, and real, data.
16 Latent Diffusion Models, such as stable diffusion ⁶⁵, have overtaken Generative Adversarial
17 Networks ⁶⁶ in the last year as the primary method for synthetic data generation. Pinaya et al.
18 ⁶⁷ recently showed that latent diffusion models could generate realistic MRI data using stable
19 diffusion algorithms and released 100,000 synthetic MRI datasets with varying ages and
20 ventricular volumes that can boost existing datasets, enabling large-scale AI approaches ⁶⁸.

21

22 The Australian Epilepsy Project is committed to publicly releasing data to the best standards
23 in the field. When sharing data, we will consider the intended use of the data and the scope of
24 the disclosure to determine how the data will be shared.

25

26 **Responsibility:**

27 Responsibility in AI was featured in 60/84 (71%) publications reviewed by Jobin et al ³⁰.
28 Terms around the topic of responsibility included accountability, liability, and acting with
29 integrity ⁶⁹⁻⁷¹

30

31 In medicine, we work with people's most private information and potentially influence life-
32 changing decisions. The argument of "who is responsible" for AI models and their decision-
33 making is not trivial. Consider a scenario where a novel (and nominally validated) deep
34 learning algorithm has detected a potential brain 'lesion' in a patient with refractory focal

1 epilepsy. The decision in this scenario was to use the AI-based lesion detection result as a
2 target for brain surgery which is often curative in epilepsy⁷². But, in this case, the surgery on
3 the brain area targeted by the AI algorithm was unsuccessful, and the patient still experienced
4 recurrent seizures. Who is responsible for this outcome? The medical professional or the AI
5 algorithm who ‘blindly’ identified a brain region that, in hindsight, was the wrong target for
6 surgery? Ultimately, treating physicians are responsible for patients, but it highlights an
7 important consideration as AI-based solutions become more integrated into clinical
8 workflows. Therefore, a pertinent current issue in AI is to determine how much we can trust
9 AI to correctly advise us of the ground truth.

10
11 Sand et al.¹⁵ highlights physicians' responsibility to understand AI models used for clinical
12 decision support. They list a set of criteria physicians should consider when employing AI
13 tools in clinical practice (building on previous work by Deitte et al.⁷³)

14 In the case of radiologists, required human knowledge includes*:

- 15 1. Reporting and informing about sensitivity rates and experimental performance
- 16 2. Understanding reasonable output
- 17 3. Understanding input data (e.g., relationship between image quality and accuracy rate)
- 18 4. Awareness of impact of utilizing medical AIs on one’s own skills and capacities
- 19 5. Awareness of task specificity of the medical AI
- 20 6. Assessing, monitoring and reporting of outputs over time

21
22 These guidelines empower human decision-making when developing, using, and interpreting
23 AI models and ensure that domain knowledge is imperative to successfully implementing AI
24 in medicine. This argument leads us onto the topic of augmented intelligence^{74,75}, an
25 emerging field in AI where people and computers symbiotically work together for the best
26 possible and explainable AI outcomes. Augmented intelligence leverages the best aspects of
27 human knowledge and computer algorithms. This includes the responsibility of those who
28 create and implement AI algorithms, as well as policymakers and regulators⁷⁶. It also implies
29 design of human-in-the-loop uses of AI in medicine, where AI is deployed to *support*
30 clinician decision making, rather than *replace* or *emulate* it⁷⁷.

31 The Australian Epilepsy Project is a patient-focused project led by clinicians, and clinical
32 expertise plays an integral role in AI model development and use. We hope our approach

* This list is verbatim from Sand et al.¹⁵.

1 toward augmented intelligence in the Australian Epilepsy Project can serve as a blueprint for
2 other projects that leverage human expertise to generate augmented intelligence-based
3 medicine. We envisage an environment in the Australian Epilepsy Project where AI provides
4 predictions and suggestions, and humans make the decisions.

6 **Other ethical considerations in AI**

7 Less discussed ethical issues in AI, according to Jobin et al.³⁰ review –highlighted in <50%
8 of AI ethics publications– included *privacy, beneficence, freedom and autonomy, trust,*
9 *sustainability, dignity, and solidarity.* We found it surprising that sustainability was only
10 considered an important AI ethics topic in 14/84 publications (17%). Although the hope is
11 that AI will ultimately play a significant role in reducing the effects of climate change and
12 other environmental impacts⁷⁸, currently, AI leaves a large carbon footprint on the planet^{79–}
13⁸¹. For example, ~15% of Google’s energy usage was dedicated to AI processing over the last
14 3 years – and the large language model used to train the GPT-3 model (the model behind
15 ChatGPT), with approximately 175 billion parameters, used 552.1 tons of CO₂-equivalent
16 emissions to train⁸². It is not only training AI algorithms that use large amounts of energy.
17 The majority of Facebook’s carbon emission is used for AI model inference –i.e., model
18 adaptation to new data– of their AI algorithms⁸³ and lifecycle of the Amazon Echo incurs a
19 massive toll on non-renewable materials, labor, and data⁸⁴. From an ethics perspective, is
20 also important to factor in that environmental impacts have disproportionately large effects
21 on marginalised communities⁸⁵.

22
23 There are several ways to reduce carbon emissions in AI. One of these includes accessing
24 computing resources in regions with cleaner energy. Accessing computing resources in
25 regions powered by more renewable energy sources can reduce emissions by up to 30 times
26⁸⁶. Another way to reduce carbon emission footprints is to explore whether a research
27 question can be answered with smaller AI models (as large AI models require the most
28 energy⁸²). Reporting the sensitivity of models to hyperparameters also allows for
29 understanding training resource requirements⁷⁹. It is furthermore good practice to use
30 frequent checkpoints when training AI models. Frequent checkpoints pause the training and
31 allow the researcher to estimate the model's performance. With checkpointing, errors or sub-
32 optimal performance can be detected early in the training of models, potentially saving time
33 and energy. It is worth noting that checkpointing can require a large amount of storage, and
34 consequent energy consumption, depending on the amount of information to evaluate. A way

1 to track the carbon footprint of AI models is the CodeCarbon Python package
 2 (<https://github.com/mlco2/codecarbon>)^{87,88}. CodeCarbon calculates the electricity
 3 consumption from GPU, CPU, and RAM, into a summary statistics of carbon footprint for
 4 the AI model.

5
 6 *Table 1. AI ethics goals in the Australian Epilepsy Project*

Category	Goals for AI ethics in the Australian Epilepsy Project
Transparency	<ul style="list-style-type: none"> • Encourage using explainable AI algorithms (e.g., MELD lesion detection approach³⁹). • Encourage using AI reporting checklists and registered reports to enhance transparency and reproducibility.
Justice and fairness	<ul style="list-style-type: none"> • Democratise clinical access for people with epilepsy in Australia. • Aim to monitor for, and remove, unacceptable demographic-related biases in AI models, including multimodal data.
Non-maleficence	<ul style="list-style-type: none"> • Sharing data to the best standards, with minimal probability of re-identification (e.g., operate within The Five Safes Framework⁸⁹)
Sustainability	<ul style="list-style-type: none"> • Encourage AI checkpoints to evaluate model performance. • Encourage tracking the carbon footprint of AI models (e.g., CodeCarbon Python package^{87,88})

7
 8 **Flexibility and adaptability in the fast-changing field of medical AI**
 9 The metrics that we use to judge and evaluate AI systems need careful consideration. While
 10 model accuracy, calibration, and robustness are common intrinsic evaluation metrics, other
 11 dimensions of AI ethics goals, including transparency, efficiency, fairness, and bias, should
 12 be quantified and reported. Various frameworks for AI model reporting have been proposed,
 13 such as model cards⁹⁰ to document the benchmark performance of models under various
 14 conditions. Proposed conditions include many directly relevant to medicine (e.g. cultural,
 15 demographic, or phenotypic groups). Medical applications of AI may require new metrics,
 16 and possibly additional dimensions – e.g. for generative language models, researchers have
 17 proposed that toxicity should be added to the list⁹¹. An additional concern for metric-based
 18 evaluation of AI models in a biomedical context is to balance the need for short-term

1 decisions about optimal model performance with the need for improved long-term health
2 outcomes in patients. Here the adage known as Goodhart’s law may be relevant, which is
3 often phrased as “when a measure becomes a target, it ceases to be a good measure”. Care
4 must be taken to avoid optimizing AI models for short-term metrics, such as performance on
5 an imaging-based lesion detection task, which may not reflect the ultimate patient-centred
6 goal of achieving improved quality of life. We must continuously mature our approach to
7 evaluation.

8

9 In 2023, AI is changing fast, and it sometimes seems futile to predict the next 5 and 10 years
10 of AI development and its impact on medicine. Therefore, we need to have a flexible and
11 adaptable approach to AI ethics that can help address the concerns of various stakeholders.
12 Viewpoints on AI ethics from patients, healthcare providers, regulators, and technology
13 developers are needed to generate adaptive frameworks that support safe AI to benefit health
14 and society. This paper is a starting point toward interpretable, secure, and sustainable AI in
15 the Australian Epilepsy Project (see Table 1). As AI is changing the ethics of medicine, we
16 aim to update ethics guidelines adaptively as we collect data from thousands of people with
17 epilepsy in the following years.

1 **Acknowledgements**

2 This work was supported by an Australian Government Medical Research Future Fund
3 Frontier Health Grant (RFRHPSI000008). MP acknowledges funding from Health Research
4 Council (HRC), New Zealand, Emerging Researcher Grant. D.F.A. acknowledges fellowship
5 funding from the Australian National Imaging Facility.

6

7 The Florey Institute of Neuroscience and Mental Health acknowledges the strong support
8 from the Victorian Government and, in particular, the funding from the Operational
9 Infrastructure Support Grant. The authors acknowledge the facilities and scientific and
10 technical assistance of the National Imaging Facility, a National Collaborative Research
11 Infrastructure Strategy (NCRIS) capability, at the Florey Institute of Neuroscience and
12 Mental Health.

13

1 **References**

- 2 1. McCarthy, J. What is Artificial Intelligence? (2004).
- 3 2. Weng, C., Shah, N. H. & Hripesak, G. Deep phenotyping: Embracing complexity and
4 temporality—Towards scalability, portability, and interoperability. *Journal of Biomedical*
5 *Informatics* **105**, 103433 (2020).
- 6 3. Acosta, J. N., Falcone, G. J., Rajpurkar, P. & Topol, E. J. Multimodal biomedical AI. *Nat*
7 *Med* **28**, 1773–1784 (2022).
- 8 4. Braun, M., Hummel, P., Beck, S. & Dabrock, P. Primer on an ethics of AI-based decision
9 support systems in the clinic. *Journal of Medical Ethics* **47**, e3–e3 (2021).
- 10 5. Geis, J. R. *et al.* Ethics of Artificial Intelligence in Radiology: Summary of the Joint
11 European and North American Multisociety Statement. *Radiology* **293**, 436–440 (2019).
- 12 6. Goldsmith, J. & Burton, E. Why Teaching Ethics to AI Practitioners Is Important.
13 *Proceedings of the AAAI Conference on Artificial Intelligence* **31**, (2017).
- 14 7. Gundersen, T. & Bærøe, K. The Future Ethics of Artificial Intelligence in Medicine:
15 Making Sense of Collaborative Models. *Sci Eng Ethics* **28**, 17 (2022).
- 16 8. Hamet, P. & Tremblay, J. Artificial intelligence in medicine. *Metabolism* **69**, S36–S40
17 (2017).
- 18 9. Martinho, A., Kroesen, M. & Chorus, C. A healthy debate: Exploring the views of
19 medical doctors on the ethics of artificial intelligence. *Artificial Intelligence in Medicine*
20 **121**, 102190 (2021).
- 21 10. McLennan, S. *et al.* Embedded ethics: a proposal for integrating ethics into the
22 development of medical AI. *BMC Med Ethics* **23**, 6 (2022).
- 23 11. Morley, J. *et al.* The ethics of AI in health care: A mapping review. *Social Science &*
24 *Medicine* **260**, 113172 (2020).
- 25 12. Muller, H., Mayrhofer, M. T., Van Veen, E.-B. & Holzinger, A. The Ten
26 Commandments of Ethical Medical AI. *Computer* **54**, 119–123 (2021).
- 27 13. Naik, N. *et al.* Legal and Ethical Consideration in Artificial Intelligence in Healthcare:
28 Who Takes Responsibility? *Front Surg* **9**, 862322 (2022).
- 29 14. Saheb, T., Saheb, T. & Carpenter, D. O. Mapping research strands of ethics of artificial
30 intelligence in healthcare: A bibliometric and content analysis. *Computers in Biology and*
31 *Medicine* **135**, 104660 (2021).
- 32 15. Sand, M., Durán, J. M. & Jongsma, K. R. Responsibility beyond design: Physicians’
33 requirements for ethical medical AI. *Bioethics* **36**, 162–169 (2022).

- 1 16. Vayena, E., Blasimme, A. & Cohen, I. G. Machine learning in medicine: Addressing
2 ethical challenges. *PLOS Medicine* **15**, e1002689 (2018).
- 3 17. Baker, G. A. Depression and suicide in adolescents with epilepsy. *Neurology* **66**, S5–S12
4 (2006).
- 5 18. Berg, A. T. Epilepsy, Cognition, and Behavior: The clinical picture. *Epilepsia* **52**, 7–12
6 (2011).
- 7 19. Elger, C. E., Helmstaedter, C. & Kurthen, M. Chronic epilepsy and cognition. *The Lancet*
8 *Neurology* **3**, 663–672 (2004).
- 9 20. Kerr, M. P. The impact of epilepsy on patients' lives. *Acta Neurologica Scandinavica*
10 **126**, 1–9 (2012).
- 11 21. Motamedi, G. & Meador, K. Epilepsy and cognition. *Epilepsy & Behavior* **4**, 25–38
12 (2003).
- 13 22. Smeets, V. M. J., van Lierop, B. A. G., Vanhoutvin, J. P. G., Aldenkamp, A. P. &
14 Nijhuis, F. J. N. Epilepsy and employment: Literature review. *Epilepsy & Behavior* **10**,
15 354–362 (2007).
- 16 23. Foster, E. *et al.* The costs of epilepsy in Australia: A productivity-based analysis.
17 *Neurology* **95**, e3221–e3231 (2020).
- 18 24. Pedersen, M. *et al.* Deep learning for reliable detection of epileptogenic lesions. in
19 *Augmenting Neurological Disorder Prediction and Rehabilitation Using Artificial*
20 *Intelligence* 163–175 (Elsevier, 2022).
- 21 25. Pedersen, M. *et al.* Artificial intelligence for clinical decision support in neurology. *Brain*
22 *Commun* **2**, (2020).
- 23 26. Tailby, C. *et al.* Teleneuropsychology in the time of COVID-19: The experience of The
24 Australian Epilepsy Project. *Seizure* **83**, 89–97 (2020).
- 25 27. Kwan, P. & Brodie, M. J. Early identification of refractory epilepsy. *N. Engl. J. Med.*
26 **342**, 314–319 (2000).
- 27 28. Wu, E. *et al.* How medical AI devices are evaluated: limitations and recommendations
28 from an analysis of FDA approvals. *Nat Med* **27**, 582–584 (2021).
- 29 29. Nagendran, M. *et al.* Artificial intelligence versus clinicians: systematic review of design,
30 reporting standards, and claims of deep learning studies. *BMJ* **368**, m689 (2020).
- 31 30. Jobin, A., Ienca, M. & Vayena, E. The global landscape of AI ethics guidelines. *Nature*
32 *Machine Intelligence* **1**, 389–399 (2019).
- 33 31. Afnan, M. A. M. *et al.* Interpretable, not black-box, artificial intelligence should be used
34 for embryo selection. *Human Reproduction Open* **2021**, hoab040 (2021).

- 1 32. Ehsan, U., Liao, Q. V., Muller, M., Riedl, M. O. & Weisz, J. D. Expanding
2 Explainability: Towards Social Transparency in AI systems. in *Proceedings of the 2021*
3 *CHI Conference on Human Factors in Computing Systems* 1–19 (Association for
4 Computing Machinery, 2021). doi:10.1145/3411764.3445188.
- 5 33. Felzmann, H., Fosch-Villaronga, E., Lutz, C. & Tamò-Larrieux, A. Towards
6 Transparency by Design for Artificial Intelligence. *Sci Eng Ethics* **26**, 3333–3361 (2020).
- 7 34. Larsson, S. & Heintz, F. Transparency in artificial intelligence. *Internet Policy Review* **9**,
8 1–16 (2020).
- 9 35. Castelvechi, D. Can we open the black box of AI? *Nature News* **538**, 20 (2016).
- 10 36. Baud, M. O. *et al.* European trends in epilepsy surgery. *Neurology* **91**, e96–e106 (2018).
- 11 37. Sinclair, B. *et al.* Machine learning approaches for imaging-based prognostication of the
12 outcome of surgery for mesial temporal lobe epilepsy. *Epilepsia* **63**, 1081–1092 (2022).
- 13 38. Suzuki, K. Overview of deep learning in medical imaging. *Radiol Phys Technol* **10**, 257–
14 273 (2017).
- 15 39. Spitzer, H. *et al.* Interpretable surface-based detection of focal cortical dysplasias: a
16 Multi-centre Epilepsy Lesion Detection study. *Brain* **145**, 3859–3871 (2022).
- 17 40. Collins, G. S., Reitsma, J. B., Altman, D. G. & Moons, K. G. Transparent reporting of a
18 multivariable prediction model for individual prognosis or diagnosis (TRIPOD): the
19 TRIPOD Statement. *BMC Medicine* **13**, 1 (2015).
- 20 41. Collins, G. S. *et al.* Protocol for development of a reporting guideline (TRIPOD-AI) and
21 risk of bias tool (PROBAST-AI) for diagnostic and prognostic prediction model studies
22 based on artificial intelligence. *BMJ Open* **11**, e048008 (2021).
- 23 42. Mongan, J., Moy, L. & Kahn, C. E. Checklist for Artificial Intelligence in Medical
24 Imaging (CLAIM): A Guide for Authors and Reviewers. *Radiology:*
25 *Artificial Intelligence* **2**, e200029 (2020).
- 26 43. Norgeot, B. *et al.* Minimum information about clinical artificial intelligence modeling:
27 the MI-CLAIM checklist. *Nat Med* **26**, 1320–1324 (2020).
- 28 44. Vasey, B. *et al.* Reporting guideline for the early stage clinical evaluation of decision
29 support systems driven by artificial intelligence: DECIDE-AI. *BMJ* **377**, e070904 (2022).
- 30 45. Walsh, I. *et al.* DOME: recommendations for supervised machine learning validation in
31 biology. *Nat Methods* **18**, 1122–1127 (2021).
- 32 46. Gebru, T. *et al.* Datasheets for datasets. *Commun. ACM* **64**, 86–92 (2021).
- 33 47. Chambers, C. D. & Tzavella, L. The past, present and future of Registered Reports. *Nat*
34 *Hum Behav* **6**, 29–42 (2022).

- 1 48. Allen, C. & Mehler, D. M. A. Open science challenges, benefits and tips in early career
2 and beyond. *PLOS Biology* **17**, e3000246 (2019).
- 3 49. Levin, L. A. & Palmer, J. G. Institutional Review Boards Should Require Clinical Trial
4 Registration. *Archives of Internal Medicine* **167**, 1576–1580 (2007).
- 5 50. Booth, A. *et al.* An international registry of systematic-review protocols. *The Lancet* **377**,
6 108–109 (2011).
- 7 51. Stewart, L., Moher, D. & Shekelle, P. Why prospective registration of systematic reviews
8 makes sense. *Systematic Reviews* **1**, 7 (2012).
- 9 52. Ge, L. *et al.* Association between prospective registration and overall reporting and
10 methodological quality of systematic reviews: a meta-epidemiological study. *Journal of*
11 *Clinical Epidemiology* **93**, 45–55 (2018).
- 12 53. Cugurullo, F. & Acheampong, R. A. Fear of AI: an inquiry into the adoption of
13 autonomous cars in spite of fear, and a theoretical framework for the study of artificial
14 intelligence technology acceptance. *AI & Soc* (2023) doi:10.1007/s00146-022-01598-6.
- 15 54. Sanderson, K. GPT-4 is here: what scientists think. *Nature* (2023) doi:10.1038/d41586-
16 023-00816-5.
- 17 55. Panch, T., Mattie, H. & Atun, R. Artificial intelligence and algorithmic bias: implications
18 for health systems. *J Glob Health* **9**, 020318.
- 19 56. Chastain, D. B. *et al.* Racial Disproportionality in Covid Clinical Trials. *N Engl J Med*
20 **383**, e59 (2020).
- 21 57. Martin, A. R. *et al.* Clinical use of current polygenic risk scores may exacerbate health
22 disparities. *Nat Genet* **51**, 584–591 (2019).
- 23 58. Alten, S. van, Domingue, B. W., Galama, T. & Marees, A. T. Reweighting the UK
24 Biobank to reflect its underlying sampling population substantially reduces pervasive
25 selection bias due to volunteering. 2022.05.16.22275048 Preprint at
26 <https://doi.org/10.1101/2022.05.16.22275048> (2022).
- 27 59. Wang, R., Chaudhari, P. & Davatzikos, C. Bias in machine learning models can be
28 significantly mitigated by careful training: Evidence from neuroimaging studies.
29 *Proceedings of the National Academy of Sciences* **120**, e2211613120 (2023).
- 30 60. Ellahham, S., Ellahham, N. & Simsekler, M. C. E. Application of Artificial Intelligence
31 in the Health Care Safety Context: Opportunities and Challenges. *Am J Med Qual* **35**,
32 341–348 (2020).
- 33 61. Urbina, F., Lentzos, F., Invernizzi, C. & Ekins, S. Dual use of artificial-intelligence-
34 powered drug discovery. *Nat Mach Intell* **4**, 189–191 (2022).

- 1 62. Mello, M. M., Lieou, V. & Goodman, S. N. Clinical Trial Participants' Views of the
2 Risks and Benefits of Data Sharing. *New England Journal of Medicine* **378**, 2202–2211
3 (2018).
- 4 63. Erlich, Y., Shor, T., Pe'er, I. & Carmi, S. Identity inference of genomic data using long-
5 range familial searches. *Science* **362**, 690–694 (2018).
- 6 64. Schwarz, C. G. *et al.* Changing the face of neuroimaging research: Comparing a new
7 MRI de-facing technique with popular alternatives. *NeuroImage* **231**, 117845 (2021).
- 8 65. Rombach, R., Blattmann, A., Lorenz, D., Esser, P. & Ommer, B. High-Resolution Image
9 Synthesis with Latent Diffusion Models. Preprint at
10 <https://doi.org/10.48550/arXiv.2112.10752> (2022).
- 11 66. Goodfellow, I. J. *et al.* Generative Adversarial Networks. Preprint at
12 <https://doi.org/10.48550/arXiv.1406.2661> (2014).
- 13 67. Pinaya, W. H. L. *et al.* Brain Imaging Generation with Latent Diffusion Models.
14 *arXiv.org* <https://arxiv.org/abs/2209.07162v1> (2022) doi:10.48550/arXiv.2209.07162.
- 15 68. Marquand, A. F., Rezek, I., Buitelaar, J. & Beckmann, C. F. Understanding
16 Heterogeneity in Clinical Cohorts Using Normative Models: Beyond Case-Control
17 Studies. *Biological Psychiatry* **80**, 552–561 (2016).
- 18 69. Coeckelbergh, M. Artificial Intelligence, Responsibility Attribution, and a Relational
19 Justification of Explainability. *Sci Eng Ethics* **26**, 2051–2068 (2020).
- 20 70. Constantinescu, M., Voinea, C., Uszkai, R. & Vică, C. Understanding responsibility in
21 Responsible AI. Dianoetic virtues and the hard problem of context. *Ethics Inf Technol* **23**,
22 803–814 (2021).
- 23 71. Tigard, D. W. Responsible AI and moral responsibility: a common appreciation. *AI*
24 *Ethics* **1**, 113–117 (2021).
- 25 72. Kral, T. *et al.* Outcome of epilepsy surgery in focal cortical dysplasia. *J Neurol*
26 *Neurosurg Psychiatry* **74**, 183–188 (2003).
- 27 73. Deitte, L. A. *et al.* Entrustable Professional Activities: Ten Things Radiologists Do. *Acad*
28 *Radiol* **23**, 374–381 (2016).
- 29 74. Gennatas, E. D. *et al.* Expert-augmented machine learning. *PNAS* **117**, 4571–4577
30 (2020).
- 31 75. Zheng, N. *et al.* Hybrid-augmented intelligence: collaboration and cognition. *Frontiers*
32 *Inf Technol Electronic Eng* **18**, 153–179 (2017).
- 33 76. Bazoukis, G. *et al.* The inclusion of augmented intelligence in medicine: A framework
34 for successful implementation. *Cell Reports Medicine* **3**, 100485 (2022).

- 1 77. Lederman, A., Lederman, R. & Verspoor, K. Tasks as needs: reframing the paradigm of
2 clinical natural language processing research for real-world decision support. *Journal of*
3 *the American Medical Informatics Association* **29**, 1810–1817 (2022).
- 4 78. Cowls, J., Tsamados, A., Taddeo, M. & Floridi, L. The AI gambit: leveraging artificial
5 intelligence to combat climate change—opportunities, challenges, and recommendations.
6 *AI & Soc* **38**, 283–307 (2023).
- 7 79. Strubell, E., Ganesh, A. & McCallum, A. Energy and Policy Considerations for Deep
8 Learning in NLP. in *Proceedings of the 57th Annual Meeting of the Association for*
9 *Computational Linguistics* 3645–3650 (Association for Computational Linguistics,
10 2019). doi:10.18653/v1/P19-1355.
- 11 80. Dhar, P. The carbon impact of artificial intelligence. *Nature Machine Intelligence* **2**, 423–
12 425 (2020).
- 13 81. Schwartz, R., Dodge, J., Smith, N. A. & Etzioni, O. Green AI. *Commun. ACM* **63**, 54–63
14 (2020).
- 15 82. Patterson, D. *et al.* The Carbon Footprint of Machine Learning Training Will Plateau,
16 Then Shrink. *Computer* **55**, 18–28 (2022).
- 17 83. Wu, C.-J. *et al.* Sustainable AI: Environmental Implications, Challenges and
18 Opportunities. (2021) doi:10.48550/arXiv.2111.00364.
- 19 84. Anatomy of an AI System. *Anatomy of an AI System* <http://www.anatomyof.ai>.
- 20 85. Bender, E. M., Gebru, T., McMillan-Major, A. & Shmitchell, S. On the Dangers of
21 Stochastic Parrots: Can Language Models Be Too Big? 🦜. in *Proceedings of the 2021*
22 *ACM Conference on Fairness, Accountability, and Transparency* 610–623 (Association
23 for Computing Machinery, 2021). doi:10.1145/3442188.3445922.
- 24 86. Henderson, P. *et al.* Towards the systematic reporting of the energy and carbon footprints
25 of machine learning. *J. Mach. Learn. Res.* **21**, 248:10039-248:10081 (2020).
- 26 87. Lottick, K., Susai, S., Friedler, S. A. & Wilson, J. P. Energy Usage Reports:
27 Environmental awareness as part of algorithmic accountability. Preprint at
28 <https://doi.org/10.48550/arXiv.1911.08354> (2019).
- 29 88. Lacoste, A., Luccioni, A., Schmidt, V. & Dandres, T. Quantifying the Carbon Emissions
30 of Machine Learning. Preprint at <https://doi.org/10.48550/arXiv.1910.09700> (2019).
- 31 89. Ritchie, F. *Five Safes: designing data access for research*. (2016).
32 doi:10.13140/RG.2.1.3661.1604.

- 1 90. Mitchell, M. *et al.* Model Cards for Model Reporting. in *Proceedings of the Conference*
2 *on Fairness, Accountability, and Transparency* 220–229 (Association for Computing
3 Machinery, 2019). doi:10.1145/3287560.3287596.
- 4 91. Liang, P. *et al.* Holistic Evaluation of Language Models. Preprint at
5 <http://arxiv.org/abs/2211.09110> (2022).
6