# A reusable benchmark of brain-age prediction from M/EEG resting-state signals

Denis A. Engemann*[a,b,c], Apolline Mellot[b], Richard Höchenberger[b], Hubert Banville[b,h], David Sabbagh[b,d], Lukas Gemein[e,f], Tonio Ball[e,g], Alexandre Gramfort[b]

[a] Roche Pharma Research and Early Development, Neuroscience and Rare Diseases, Roche Innovation Center Basel, F. Hoffmann–La Roche Ltd., Basel, Switzerland
[b] Université Paris-Saclay, Inria, CEA, Palaiseau, France
[c] Max Planck Institute for Human Cognitive and Brain Sciences, Department of Neurology, D-04103, Leipzig, Germany
[d] Inserm, UMRS-942, Paris Diderot University, Paris, France
[e] Neuromedical AI Lab, Department of Neurosurgery, Medical Center – University of Freiburg, Faculty of Medicine, University of Freiburg, Engelbergerstr. 21, 79106, Freiburg, Germany
[f] Neurorobotics Lab, Computer Science Department – University of Freiburg, Faculty of Engineering, University of Freiburg, Georges-Köhler-Allee 80, 79110, Freiburg, Germany
[g] BrainLinks-BrainTools Cluster of Excellence, University of Freiburg, Freiburg, Germany
[h] InteraXon Inc., Toronto, Canada

Correspondence: *denis.engemann@roche.com

## Keywords

clinical neuroscience, brain age, electroencephalography, magnetoencephalography, machine learning, population modeling, Riemannian geometry, , random forests, deep learning

## Highlights

- We provide systematic reusable benchmarks for brain age from M/EEG signals

- The benchmarks were carried out on M/EEG from four countries > 2500 recordings

- We compared machine learning pipelines capable of handling the non-linear regression task of relating biomedical outcomes to M/EEG dynamics, based on classical machine learning and deep learning

- Next to data-driven methods we benchmarked template-based source localization as a practical tool for generating features less affected by electromagnetic field spread

- The benchmarks are built on top of the MNE ecosystem and the braindecode package and can be applied on any M/EEG dataset presented in the BIDS format

## Abstract

Population-level modeling can define quantitative measures of individual aging by applying machine learning to large volumes of brain images. These measures of brain age, obtained from the general population, helped characterize disease severity in neurological populations, improving estimates of diagnosis or prognosis. Magnetoencephalography (MEG) and Electroencephalography (EEG) have the potential to further generalize this approach towards prevention and public health by enabling assessments of brain health at large scales in socioeconomically diverse environments. However, more research is needed to define methods that can handle the complexity and diversity of M/EEG signals across diverse real-world contexts. To catalyse this effort, here we propose reusable benchmarks of competing machine learning approaches for brain age modeling. We benchmarked popular classical machine learning pipelines and deep learning architectures previously used for pathology decoding or brain age estimation in 4 international M/EEG cohorts from diverse countries and cultural contexts, including recordings from more than 2500 participants. Our benchmarks were built on top of the M/EEG adaptations of the BIDS standard, providing tools that can be applied with minimal modification on any M/EEG dataset provided in the BIDS format. Our results suggest that, regardless of whether classical machine learning or deep learning was used, the highest performance was reached by pipelines and architectures involving spatially aware representations of the M/EEG signals, leading to $R^2$ scores between 0.60-0.71. Hand-crafted features paired with random forest regression provided robust benchmarks even in situations in which other approaches failed. Taken together, this set of benchmarks, accompanied by open-source software and high-level Python scripts, can serve as a starting point and quantitative reference for future efforts at developing M/EEG-based measures of brain aging. The generality of the approach renders this benchmark reusable for other related objectives such as modeling specific cognitive variables or clinical endpoints.

## Introduction

1  Aging-related disorders of the central nervous system affect hundreds of millions of patients,
2  their caregivers and national health services. Over the past decades, important progress has
3  been made in clinical neuroscience, resulting in improvements to clinical diagnosis and
4  treatment (Walhovd et al. 2010; Ewers et al. 2011). Backed by increasingly advanced
5  analytical methods, this has enabled fine-grained characterization of neurodegenerative
6  conditions (Gaubert et al. 2019; Schumacher et al. 2021; Güntekin et al. 2021). Yet, from a
7  public-health perspective, rather than focusing on pathology, it is essential to detect risk

8  factors early within the general population in order to provide actionable feedback for
9  preventive medicine, e.g., by targeting life-style changes. Such predictions are still
10  challenging. Could it be helpful to look at biological rather than chronological age to better
11  estimate the risk of declining brain health?

12  Recently, brain age has emerged as a concept for estimating biological aging in the general
13  population (James H. Cole and Franke 2017; Liem et al. 2017; Dosenbach et al. 2010).
14  Biological aging can be inferred from the genome via telomere length, mitochondrial function,
15  epigenetics and other cellular features (Ferrucci et al. 2020; Mather et al. 2011). Yet, the age
16  of a person is only a noisy measure of these cellular processes (people of the same
17  chronological age can have different biological ages). At the same time, biological aging
18  affects brain structure and function (K. S. King et al. 2014), inducing loss of brain volume
19  (Driscoll et al. 2009; Scahill et al. 2003) and characteristic changes in neuronal activity
20  (Cabeza et al. 2002; Damoiseaux et al. 2008; Babiloni et al. 2006). A proxy of biological aging
21  can, thus, be obtained by mapping chronological age to brain data from large populations of
22  subjects using machine learning (Liem et al. 2017; Dadi et al. 2021). The resulting models
23  can be used to compute an expectation of a person's age given her brain data. This is
24  achieved by quantitatively comparing that person's brain data to the distribution of brain data
25  across different ages within the general population. This statistical expectation can tell how
26  old (or young) a brain "looks" (Spiegelhalter 2016), hence, predicting the risk of neurological
27  complications potentially more precisely than the chronological age.

28  This empirical measure of biological aging derived from the general population has proven a
29  useful marker of neurodegeneration and cognitive decline in clinical populations (Cole et al.
30  2018; Raffel et al. 2017; Denissen et al. 2021; Gonneaud et al. 2021). In these cohorts,
31  patients typically appear to have older brains than their chronological age would suggest.
32  Importantly, similar trends emerge when evaluating brain age in the general population where
33  elevated brain age, compared to chronological age, has been associated with lower cognitive
34  capacity, well-being, and general health (Dadi et al. 2021; Cole 2020; Wrigglesworth et al.
35  2021). Yet, so far, this approach has mainly been based on anatomical brain scans and
36  hemodynamic signals obtained from magnetic resonance imaging (MRI). This limits the broad
37  utility of brain age for public health, as cerebral MRI scans are usually collected when there is
38  an indication, which can be too late. Even when people from the general population are
39  motivated to participate in brain research, this only concerns a small fraction of society: MRI
40  devices and neuroscientific studies are not equally accessible in all regions of the world and
41  do not attract all people equally from within society, potentially leading to selection bias (Fry
42  et al. 2017).

43  New hope to generalize this approach has been sparked by advances in large-scale modeling
44  of biomedical outcomes from non-invasive electrophysiological data including

magnetoencephalography (MEG) and electroencephalography (EEG) (Gaubert et al. 2019; Engemann et al. 2018). This line of research in clinical neurology may help develop assessments of brain health in many additional contexts in which MRI cannot be applied. First MEG-based brain-age models have allowed to validate MEG-derived brain age against MRI-derived brain age. Results from several studies have shown that the MEG- and MRI-derived brain age are statistically related, leading to overlapping correlations between ensuing brain age estimates (Engemann et al. 2020; Sabbagh et al. 2020; Xifra-Porxas et al. 2021) and individual differences in cognition and health. This overlap can be explained by electromagnetic field spread, independently of neuronal activity: As brain structure changes due to aging, cortical activity, even if unchanged, will project differently onto the M/EEG sensor array, making age indirectly decodable (Sabbagh et al. 2020). Importantly, multiple articles have found that neuronal activity captured by MEG adds specific information not present in MRI-derived brain age (Engemann et al. 2020; Xifra-Porxas et al. 2021), leading to improved prediction performance and richer neurocognitive characterization (Engemann et al. 2020).

While MEG can provide an important discovery context, it is unlikely to be the right instrument for addressing the availability issues of MRI-based brain age as MEG scanners are even rarer than MRI scanners. In this context, EEG can make a true difference as EEG is economical and allows for flexible instrumentation for neural assessments in a wide range of clinical and real-world situations including at-home assessments. First evidence suggests that MEG-based strategies for brain-age modeling can be translated to EEG. In an earlier publication (Engemann et al. 2020) we found that among many alternative features of varying data-processing complexity, the spatial distribution of cortical power spectra in the beta (13-30Hz) and alpha (8-13Hz) frequency band explained most of the MEG's performance as brain-age regressor. This type of information can be well accessed without source localization from the sensor-space covariance using spatial filtering approaches or Riemannian geometry (Sabbagh et al. 2020; D. Sabbagh et al. 2019), which has led to successful translation of this MEG-derived strategy to clinical EEG with around 20 electrodes (David Sabbagh et al. 2020). In clinical and real-world contexts in which EEG is frequently collected, fine-grained spatial information may not be present as only a few electrodes are used. This has favored alternative EEG-derived brain-age models focusing on a wealth of spectral and temporal features (Al Zoubi et al. 2018) which may perform better on sparse EEG-montages and has enabled sleep-based brain age measures (Sun et al. 2019; Ye et al. 2020).

These results provide a sense of the flexibility and future potential of EEG-based brain age as a widely applicable real-world measure of brain health. Yet, to fully develop this research program, more and richer evidence is desirable. At this point, comparisons between different machine learning strategies are difficult. Most models were not only developed and validated in one specific context, but their implementations and data-processing routines are dataset-

82   specific. Moreover, general machine learning approaches successful at pathology decoding
83   should be well-suited for brain age modeling too, yet they have never been tested for that
84   purpose (Gemein et al. 2020; Banville et al. 2020; Engemann et al. 2018). This makes it hard
85   to know whether any strategy is globally optimal and where specific strategies have their
86   preferred niche. As a result, uncertainty is added to comparisons between MEG, EEG and
87   MRI, slowing down efforts of validating M/EEG-based brain age. Finally, to mitigate the impact
88   of selection bias concerning the subjects investigated, it will be crucial to analyze many,
89   socially and culturally diverse M/EEG datasets and find representations that are invariant to
90   confounding effects that can raise issues of fairness and racial bias if remaining unaddressed
91   (Choy, Baker, and Stavropoulos 2021). To develop the next generation of M/EEG-derived
92   brain age models, to facilitate processing of larger numbers of diverse M/EEG-data resources
93   and to avoid fragmentation of research efforts, standardized software and reusable
94   benchmarks are needed.

95   In this paper we wish to make a first step in that direction. We provide reusable brain-age-
96   prediction benchmarks for different machine learning strategies validated on multiple M/EEG
97   datasets from different countries. The benchmarks are built on top of highly standardized
98   dataset-agnostic code enabled by the BIDS standard (Gorgolewski et al. 2016; Niso et al.
99   2018; Appelhoff et al. 2019). This makes the benchmarks easy to extend in the future for
100  additional datasets. The paper is organized as follows. The method section motivates the
101  choice of the different machine learning benchmarks. The general data processing approach
102  and software developed for this contribution are presented in the context of the benchmark.
103  The selection of datasets is motivated, and datasets are then described in detail and
104  compared regarding key figures that could provoke differences between benchmarks.
105  Dataset-specific processing steps and peculiarities are highlighted. Then a model validation
106  strategy is developed. The results section presents benchmarks on prediction performance
107  across machine learning models and datasets and different performance metrics. The
108  discussion inspects differences between models, modalities, and datasets, identifying unique
109  niches, safe bets as well as unresolved challenges. The work concludes with practical
110  suggestions on additional benchmarks that can be readily explored using the proposed tools
111  and resources. The scripts and library code for this benchmark are publicly available on
112  GithHub[1].

---

[1] https://github.com/meeg-ml-benchmarks/meeg-brain-age-benchmark-paper

## Methods

### Brain age benchmarks

113  Many different approaches exist for ML in neuroscience, and it can be hard to select among

114  them. The following categorization may help orient practical reasoning and study design. What

115  varies in the taxonomy of methods discussed below is how much M/EEG data are statistically

116  summarized before being presented to the learning algorithm. In other words, ML methods

117  vary with respect to the extent to which compression and summary of the M/EEG signals is

118  performed by the learning algorithm vs. feature-defining procedures performed before and

119  independently of the machine learning algorithm.

#### *A-priori defined, a.k.a. handcrafted, features*

120  The first category represents approaches in which features are inspired by theoretical and

121  empirical results in neuroscience or neural engineering. Here, M/EEG is summarized in a rigid

122  fashion by global aggregation across sensors, time, and frequencies or by visiting specific

123  regions of interest (Gemein et al. 2020; Sitt et al. 2014; Engemann et al. 2018). A meaningful

124  composition of features requires prior knowledge of the (clinical) neuroscience literature,

125  especially when interpretation of the model is a priority. In practice, it is convenient to extract

126  all or the most relevant features discussed in a given field, apply multiple spatial and temporal

127  aggregation strategies, and then bet on the capacity of the learning algorithm to ignore

128  irrelevant features (Sitt et al. 2014). This motivates the use of tree-based algorithms like

129  random forests (Breiman 2001) that are easy to tune, can fit nonlinear functions (higher-order

130  interaction effects), and are relatively robust to the presence of uninformative features. As

131  local methods that can be seen as adaptive nearest neighbors (Hastie et al. 2005), the

132  predictions of random forests and related methods are bounded by the minimum and

133  maximum of the outcome in the training distribution. For clinical neuroscience applications,

134  this has proven to yield robust off-the-shelf prediction models that are relatively unaffected by

135  noise in the data and in the outcome (Engemann et al. 2018). This approach is also a natural

136  choice when using sparse EEG-montages with few electrodes.

137  Here we implemented a strategy pursued in (Gemein et al. 2020) and (Banville et al.

138  2020), aiming at a broad set of different summary statistics of the time-series or the power

139  spectrum. This approach has turned out useful for a pathology detection task in which the

140  labeling of EEG as pathological can be due to different clinical reasons, hence, affecting many

141  different EEG signatures in potentially diffuse ways. Features were computed using the MNE-

142  features package (Schiratti, Le Douget, Van Quyen, et al. 2018). More specifically we used

143  as features (each computed for individual channels and concatenated across channels, and

144  then averaged across epochs): the standard-deviation, the kurtosis, the skewness, the

145     different quantiles (10%, 25%, 75%, 90%), the peak-to-peak amplitude, the mean, the power

146     ratios in dB among all frequency bands (0 to 2Hz, 2 to 4Hz, 4 to 8Hz, 8 to 13Hz, 13 to 18Hz,

147     18 to 24Hz, 24 to 30Hz and 30Hz to 49Hz), the spectral entropy (Inouye et al. 1991), the

148     approximate and sample entropy (Richman and Moorman 2000), the temporal complexity

149     (Roberts, Penny, and Rezek 1999), the Hurst exponent as used in (Devarajan et al. 2014),

150     the Hjorth complexity and mobility as used in (Päivinen et al. 2005), the line length (Esteller,

151     Echauz, et al. 2001), the energy of wavelet decomposition coefficients as proposed in

152     (Teixeira et al. 2011), the Higuchi fractal dimension as used in (Esteller, Vachtsevanos, et al.

153     2001), the number of zero crossings and the SVD Fisher Information (per channel) (Roberts,

154     Penny, and Rezek 1999).

### *Covariance-based filterbank approaches*

155     This category represents approaches in which the spatial dimension of M/EEG is fully exposed

156     to the model, whereas temporal or spectral aspects of the signal are to some extent

157     summarized before modeling. As M/EEG signals reflect linear superposition of neuronal

158     activity projected to the sensors through linear field/potential spread, it is natural to use linear

159     (additive) models for adaptively summarizing the spatial dimension of M/EEG signals (King et

160     al. 2018; Stokes, Wolff, and Spaak 2015; King and Dehaene 2014). This intuition is driving

161     the success of linear decoders for evoked response analysis but faces additional challenges

162     when applied to power spectra (Sabbagh et al. 2020). Computing power features on M/EEG

163     sensor-space signals renders the regression task a non-linear problem for which linear models

164     will provide sub-optimal results (Sabbagh et al. 2019). In practice, this can be overcome by

165     extracting nonlinear features like spectral power after anatomy-based source localization, or

166     in a data-driven fashion that does not require availability of individual MRI scans. Spatial

167     filtering techniques provide unmixing of brain sources based on statistical criteria without

168     using explicit anatomical information, which has led to supervised spatial filtering pipelines

169     (de Cheveigné and Parra 2014; Dähne et al. 2014). Another related strategy consists in

170     computing features that are invariant to field spread. This can be achieved by Riemannian

171     geometry, an approach first applied to M/EEG in the context of brain computer interfaces but

172     that has also proven effective for biomarker learning (Barachant et al. 2012; Yger, Berar, and

173     Lotte 2017; Rodrigues, Jutten, and Congedo 2019). These approaches have in common to

174     favor the covariance of M/EEG sensors as a practical representation of the signals.

175     Manipulating the covariance allows one to suppress the effects of linear mixing while, at the

176     same time, exposing the power spectrum and the spatial structure of neuronal activity in each

177     frequency band (Sabbagh et al. 2020). To scan along the entire power spectrum, one

178     computes covariances from several narrow-band signals covering low to high frequencies

179    (Sabbagh et al. 2020). This provides spatially fine-grained information of frequency-specific

180    neuronal activity, hence the term *filterbank*.

181        Here we implemented the filterbank models from (Sabbagh et al. 2020; Sabbagh et al.

182    2019) based on Riemannian geometry that were found to provide a practical alternative to

183    MRI-based source localization, although falling slightly behind in terms of performance. This

184    may be explained by the model violations arising from computing the Riemannian embedding

185    across multiple participants. The Riemannian embedding assumes linear field spread but

186    each recording comes from a different head and different sensor locations, which is explicitly

187    modeled when computing individual-specific source estimates. It is an open question whether

188    template-based source localization can improve upon the Riemannian pipeline, observing that

189    in the case of MEG such a procedure would be informed by the head position in the MEG

190    dewar. Both average brain templates and Riemannian embeddings mitigate field spread in a

191    global way with the difference that the average template uses some anatomical information

192    and approximate sensor locations in the context of MEG, whereas Riemannian embeddings

193    are purely a data-driven procedure with some whitening based on the average covariance

194    (across subjects).

195        To evaluate the benefit of a template-based anatomy, we included a filterbank model

196    using source localization based on the *fsaverage* subject from FreeSurfer (Fischl 2012). The

197    forward model was computed with a 3-layer Boundary Element Method (BEM) model. Source

198    spaces were equipped with a set of 4098 candidate dipole locations per hemisphere. Source

199    points closer than 5mm from the inner skull surface were excluded. The noise covariance

200    matrices used along with forward solutions to compute minimum-norm estimates inverse

201    operators were taken as data-independent diagonal matrices. Diagonal values defaulted to

202    the M/EEG-specific expected scale of noise (obtained via the "make_ad_hoc_cov" function

203    from MNE-Python). All computations were done with MNE (Gramfort et al. 2014, 2013). For

204    computational efficiency, source power estimates were obtained by applying the inverse

205    operators to the subjects' covariance data (MNE-Python function "apply_inverse_cov").

206    Dimensionality reduction was carried out with a parcellation containing 448 ROIs (Khan et al.

207    2018). This procedure closely followed the one from (Engemann et al. 2020), with the

208    difference that here an MRI template was used instead of subject-specific MRIs. Finally, the

209    448 ROI-wise source power estimates represented as diagonal matrices were the inputs of

210    the log-diag pipeline from (Sabbagh et al. 2020; D. Sabbagh et al. 2019). Features were

211    computed using the coffeine package[2].

---

[2] https://github.com/coffeine-labs/coffeine

### *Deep learning approaches*

212  This category concerns modeling strategies in which the outcome is mapped directly from the
213  raw signals without employing separate a priori feature-defining procedures. Instead, multiple
214  layers of nonlinear but parametric transformations are estimated end-to-end to successively
215  summarize and compress the input data. This process is controlled by supervision and
216  enabled by a coherent single optimization objective. In many fields, emerging deep learning
217  methods keep defining the state of the art in generalization performance, often outperforming
218  humans. Deep learning models are however greedy for data, and it may take hundreds of
219  thousands if not millions of training examples until these models show a decisive advantage
220  over classical machine-learning pipelines. Applied to neuroscience, where the bulk of datasets
221  is small to medium-sized, deep learning models may or may not outperform classical
222  approaches (Poldrack, Huckins, and Varoquaux 2020; Schulz et al. 2020; Roy et al. 2019; He
223  et al. 2020). The success of using a deep-learning model may, eventually, depend on the
224  amount of energy and resources invested in its development (Gemein et al. 2020).

225  Apart from high performance on standard laboratory M/EEG datasets and decoding
226  tasks, deep learning models are attractive for other reasons. First, when very specific
227  hypotheses about data generators or noise generators are available (Kietzmann, McClure,
228  and Kriegeskorte 2019). In this setting, the model architecture can be designed to implement
229  this knowledge, e.g. to explicitly extract band power features in a motor decoding task.
230  Second, these models have a strategic advantage when the data generating mechanism is
231  not known at all, hence, few hypotheses about classes of features are available (Schirrmeister
232  et al. 2017). In this setting, models with a generic architecture can learn and identify relevant
233  features themselves without requiring expert knowledge of the researcher. With neural
234  architecture search and automated hyperparameter optimization, there is also intense
235  research to even reduce the amount of expert knowledge needed to create the network
236  architecture itself. This flexibility has led neuroscientists to discover the framework as a vector
237  for hypothesis-driven research probing brain functions and neural computation (Yamins and
238  DiCarlo 2016; Bao et al. 2020). At the same time, this flexibility is equally beneficial under
239  complex environmental conditions that degrade the quality of M/EEG recordings (e.g. real-
240  world recordings outside of controlled laboratory conditions), in which the classes of relevant
241  features are not a priori known and deep learning models can exploit the structure of the data
242  and noise sources to provide robust predictions. (Banville et al. 2021).

243  Based on prior work, here we benchmarked two battle-tested general architectures
244  (Gemein et al. 2020) implemented using the Braindecode package[3] (Schirrmeister et al. 2017;
245  Gramfort et al. 2013). Braindecode is an open-source library for end-to-end learning on EEG

---

[3] https://braindecode.org

246  signals. It is closely intertwined with other libraries. One of them is Mother of all BCI

247  Benchmarks (MOABB) (Jayaram and Barachant 2018), which allows for convenient EEG-

248  data fetching, MNE (Gramfort et al. 2013, 2014), implements well established data structures,

249  preprocessing functionality, and more. A second key dependency is Skorch (Tietz et al. 2017),

250  which implements the commonly known scikit-learn (Pedregosa et al. 2011) API for neural

251  network training (Buitinck et al. 2013). For these reasons, Braindecode is equally useful for

252  EEG researchers who desire to apply deep learning as well as for deep learning researchers

253  who desire to work with EEG data. Braindecode builds on PyTorch (Paszke et al. 2019) and

254  comprises a zoo of decoding models that were already successfully applied to a wide variety

255  of EEG decoding classification and regression tasks, such as motor (imagery) decoding

256  (Schirrmeister et al. 2017; Kostas and Rudzicz 2020), pathology decoding (Gemein et al.

257  2020; van Leeuwen et al. 2019; Tibor Schirrmeister et al. 2017), error decoding (Völker et al.

258  2018), sleep staging (Chambon et al. 2018; Perslev et al. 2021), and relative positioning

259  (Banville et al. 2020).

260  For this benchmark and the task of age regression we used two Convolutional Neural

261  Networks (ConvNets, sometimes abbreviated CNNs) (LeCun et al. 1999) namely

262  ShallowFBCSPNet (BD-Shallow) and Deep4Net (BD-Deep) (Schirrmeister et al. 2017). BD-

263  Shallow was inspired by the famous filter bank common spatial pattern (FBCSP) (Ang et al.

264  2008) algorithm. Initially, it has two layers that represent a temporal convolution as well as a

265  spatial filter. Together with a squaring and logarithmic non-linearity it was designed to

266  specifically extract bandpower features. Of note, in the present context this architecture is

267  closely related to SPoC (Dähne et al 2014) and, in therefore, in principle, has the capacity to

268  deliver consistent regression models as was formally proven in previous work (Sabbagh et al

269  2020).

270  In contrast, BD-Deep is a much more generic architecture. In total, it has four blocks of

271  convolution-max-pooling and is therefore not restricted to any specific features. While BD-

272  Deep has around 276k trainable parameters and has therefore more learning capacity, BD-

273  Shallow has only about 36k parameters.

274  It is important to note, that we did neither adjust the model architectures (apart from those

275  changes required by the regression task) nor run task-specific hyperparameter optimization.

276  Both ConvNets were used as implemented in Braindecode with hyperparameters that were

277  already successfully applied to pathology decoding from the TUH Abnormal EEG Corpus

278  (Gemein et al. 2020; van Leeuwen et al. 2019; Tibor Schirrmeister et al. 2017). For more

279  information on Braindecode or the ConvNets, please refer to the original publication

280  (Schirrmeister et al. 2017). For decoding, we converted the MEG input data from Tesla to

281  Femtotesla, the EEG input data from Volts to Microvolts, and additionally rescaled the data,

282    such that it has roughly zero mean and unit variance by dividing by the standard deviation of

283    each dataset (see Section Datasets).

**General data processing strategy using BIDS and the MNE-BIDS pipeline**

284    Neuroimaging and behavioral data are stored in many different complex formats, potentially

285    hampering efforts of building widely usable methods, hence, impeding reproducible research.

286    Our goal was to provide brain-age prediction models that can be directly applicable to any

287    new electrophysiological dataset. For this purpose, we used the Brain Imaging Data Structure

288    (BIDS) (Gorgolewski et al. 2016) which allows us to organize neuroimaging data in a

289    standardized way supporting interoperability between programming languages and software

290    tools. We used the MNE-BIDS software (Appelhoff et al. 2019) for programmatically

291    converting M/EEG datasets into the BIDS format (Pernet et al. 2019; Niso et al. 2018). This

292    has allowed us to access all datasets included in this work in the same way, enabling data

293    analysis for all these datasets with the same code. We will now summarize the general

294    workflow (cf. *Fig. 1*).

295         For this study, we used the MNE-BIDS-Pipeline for automatic preprocessing of MEG

296    and EEG data stored in BIDS format[4] (Jas et al. 2018). Its main advantage is that we can

297    implement various custom analyses for different datasets without having to write any

298    elaborate code. Modifying the overall processing pipeline or adapting a given pipeline to a

299    new dataset only requires few edits. Controlling the pipeline is achieved through dataset-

300    specific configuration files that specify the desired processing steps and options of the MNE-

301    BIDS-Pipeline while dealing with the peculiarities of the data. The MNE-BIDS-Pipeline scripts

302    themselves do not need to be modified and are readily applicable on diverse datasets.
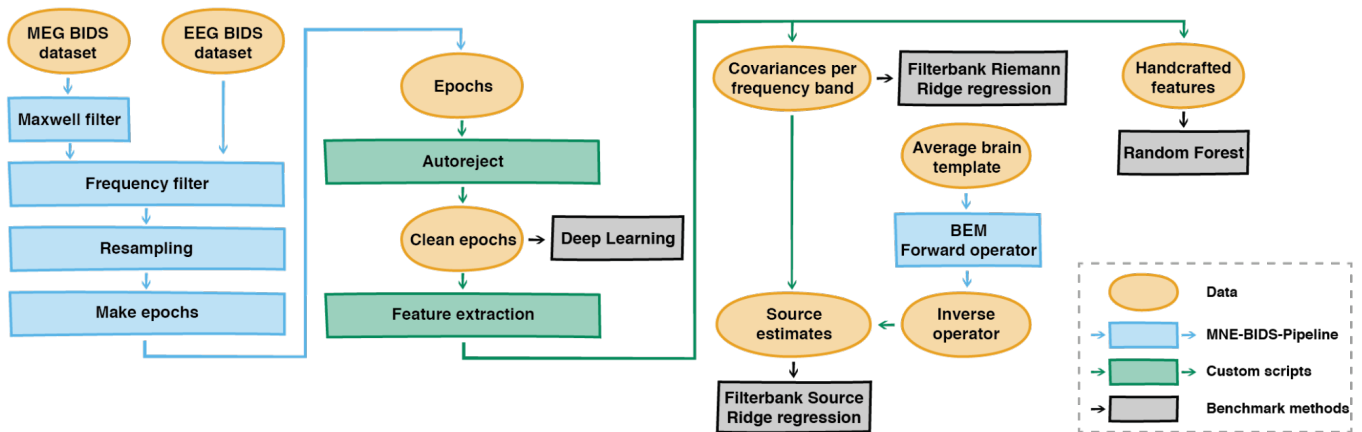
303         We designed configuration files to implement data processing steps common to all

304    datasets analyzed in this benchmark while handling dataset-specific details. Raw signals

305    bandpass-filtered between 0.1 and 49Hz using a zero-phase finite impulse response (FIR)

306    filter with Hamming window. Window length and transition bandwidth were automatically

307    controlled by default settings of MNE-Python (v0.24). We considered epochs of 10-second

308    length without overlap. These epochs coincided with eyes-closed or eyes-open resting-state

309    conditions in some of the datasets. As additional channels measuring ocular and cardiac

310    activity were not consistently available across datasets, we only implemented amplitude-

311    based artifact rejection using the local autoreject method (Jas et al. 2017). Through 5-fold

312    cross-validation, autoreject chose channel-specific rejection peak-to-peak-amplitude

313    thresholds and then decided if a given epoch could be repaired using interpolation, or if it

314    should be rejected to obtain clean data. We kept the default grid of candidate values for the

---

[4] https://github.com/mne-tools/mne-bids-pipeline

315    hyperparameters 'rho' (the consensus proportion of bad channels leading to rejection of an

316    epoch) and 'kappa' (maximum number of channels allowed to be interpolated). For 'rho' we

317    considered a linearly spaced grid of 11 points between 0 and 1. For 'kappa' we considered 1,

318    4, or 32 channels. As the local autoreject is not yet supported in the MNE-BIDS pipeline, this

319    step was implemented in a custom script (see the "compute_autoreject.py" in the code

320    repository). Apart from preprocessing, we also made use of the MNE-BIDS-Pipeline to

321    generate forward solutions and inverse operators for the source localization approach based

322    on template MRI (see section *Covariance-based filterbank approaches* for detailed

323    explanations).

324         Each model of the benchmark is based on features extracted from clean epochs.

325    Again, the conversion of datasets to BIDS has enabled feature extraction using one general

326    script for all datasets ("compute_features.py" in the code repository).



**Figure 1: Data processing, feature extraction and model construction based on the BIDS standard.** This benchmark project provides a common data processing and feature extraction code allowing comparisons of different classical and deep learning-based machine learning models across different M/EEG datasets. Support for new datasets can be added with minimal modifications. For a detailed description consider the main text and the open-source code repository supporting this article[5].

## Datasets

327    Large datasets and biobanks are the backbone of population modeling. In the past 10 years,

328    this has led to a wealth of publications in cognitive neuroscience on modeling biomedical

329    outcomes and individual differences in cognition from MRI data (Kernbach et al. 2018; James

330    H. Cole 2020; Smith et al. 2015). This has been enabled by consortia and large-scale

331    institutional collaborations (Bycroft et al. 2018; Van Essen et al. 2013) that aim at

332    recontextualizing existing data for open-ended future usage (Leonelli 2016). More recently,

333    the first M/EEG datasets have emerged with a focus on characterizing populations (Taylor et

334    al. 2017; Larson-Prior et al. 2013; Babayan et al. 2019; Obeid and Picone 2016; Niso et al.

---

[5] https://github.com/meeg-ml-benchmarks/meeg-brain-age-benchmark-paper

335   2016; Valdes-Sosa et al. 2021; Bosch-Bayard et al. 2020). The selection of datasets for the

336   present study did not aim at comprehensiveness but represents an attempt to secure a

337   minimum degree of diversity. Social bias and fairness are important challenges, not only in

338   the field of machine learning but also in biomedical research. It has been shown for modern

339   biobanks that the sample deviates from the general population in important ways,

340   oversampling Caucasian people with higher education degrees (Fry et al. 2017; Henrich and

341   Heine 2010). For deployment of predictive biomarkers, this can have tragic consequences as

342   clinical utility may depend on sex and ethnicity (Duncan et al. 2019). As a result, in EEG

343   research, specific risks of racial bias have been recognized lately, highlighting the risk of

344   selection bias and confounding, e.g., due to culture-specific hair style (Choy, Baker, and

345   Stavropoulos 2021). Taken together, this emphasizes the importance of benchmarking on

346   socially and culturally different datasets. Our selection includes M/EEG datasets from four

347   different countries representing culturally and socioeconomically diverse contexts. In the

348   following we will provide a high-level introduction to the datasets, highlighting characteristic

349   differences, challenges and opportunities for unique benchmarks.

### *Cam-CAN MEG data.*

350   The Cambridge Centre of Ageing and Neuroscience (Cam-CAN) dataset (Taylor et al. 2017;

351   Shafto et al. 2014) has been the starting point of our efforts in building brain age models

352   (Engemann et al. 2020; David Sabbagh et al. 2020) and we like to see it as a discovery

353   context. The combination of a wide, almost uniformly distributed age range and MEG data

354   alongside MRI and fine-grained neurobehavioral results make it a rich resource for exploring

355   aging-related cortical dynamics. On the other hand, models developed on this dataset may

356   not be generalizable to real-world contexts in which EEG is operated. The following two

357   sections are based on the methods description from our previous publications (Engemann et

358   al. 2020; Sabbagh et al. 2020).

359   *Sample description.* The present work was based on the latest BIDS release of the

360   Cam-CAN dataset (downloaded February 2021). We included resting-state MEG recordings

361   from 646 participants (female = 319, male = 327). The age of the participants ranged from

362   18.5 to 88.9 years with a mean age of 54.9 (female = 54.5, male = 55.4) and a standard

363   deviation of 18.4 years. Data is provided in Tesla and has a standard deviation of 369.3

364   Femtotesla. We did not apply any data exclusion. Final numbers of samples reflect successful

365   preprocessing and feature extraction. For technical details regarding the MEG instrumentation

366   and data acquisition, please consider the reference publications by the Cam-CAN (Taylor et

367   al. 2017; Shafto et al. 2014). In the following we highlight a few points essential for

368   understanding our benchmarks on the Cam-CAN MEG data.

369     *Data acquisition and processing.* MEG was recorded with a 306 VectorView system
370     (Elekta Neuromag, Helsinki). This system allowed measuring magnetic fields with 102
371     magnetometers and 204 orthogonal planar gradiometers inside a light magnetically shielded
372     room. During acquisition, an online filter was applied between around 0.03Hz and 1000Hz.
373     After bandpass filtering (0.1 - 49Hz), we applied decimation by a factor of 5, leading to a
374     sample frequency of 200Hz (at the epoching stage). To mitigate the contamination of the MEG
375     signal by environmental magnetic interference, we applied the temporal signal-space-
376     separation (tSSS) method (Taulu, Simola, and Kajola 2005). Default settings were applied for
377     the harmonic decomposition (8 components of the internal sources, 3 for the external sources)
378     on a 10-s sliding window. To discard segments for which inner and outer signal components
379     were poorly distinguishable, we applied a correlation threshold of 98%. As a result of this
380     procedure, the signal was high pass filtered at 0.1Hz and the dimensionality of the data was
381     reduced to 65, approximately. It is worthwhile to note that Maxwell filtering methods like tSSS
382     merge the signal from magnetometers and gradiometers into one common low-rank
383     representation. As a result, after tSSS, the signal displayed on magnetometers becomes a
384     linear transformation of the signals displayed on the gradiometers. This leads to virtually
385     identical results when conducting analyses exclusively on magnetometers versus
386     gradiometers (Garcés et al. 2017). To reduce computation time, we analyzed the
387     magnetometers for our benchmark. To deal with the reduced data rank, a PCA projection to
388     the common rank of 65 was applied whenever the machine learning pipeline was sensitive to
389     the rank (e.g., Riemannian filterbank models). For the full specification of the preprocessing,
390     please refer to the "config_camcan_meg.py" file in the code repository.

### LEMON EEG data.

391     The Leipzig Mind-Brain-Body (LEMON) dataset offers rich multimodal EEG, MRI and fMRI
392     data for a well characterized group of young and elderly adults sampled from the general
393     population (Babayan et al. 2019). As it was the case for the Cam-CAN data, here the research
394     was conducted in a research context using high-end equipment accompanied by rich and fine-
395     grained neurocognitive and behavioral assessments.

396     *Sample description.* EEG resting-state data from 227 healthy individuals from the
397     LEMON dataset were included in this study. This sample contains 82 females (mean age =
398     44.2) and 145 males (mean age = 36), representing a clearly visible difference in the
399     composition of the sample (*Fig.* 2). Their age distribution went from 20 to 77 years old with an
400     average of 38.9 +- 20.3 years. Our sample covers the whole available dataset (downloaded
401     September 2021) as we did not apply any exclusion criteria. It is a peculiarity of this dataset
402     is that it is divided into 2 distinct age subpopulations, one between 20-35, the second between
403     55-77 (*Fig.* 2), rendering the mean a bad representation of the age distribution. Moreover, the

404 public version of the datasets only provides ages in a granularity of 5 years to mitigate the risk
405 of identifying participants. For the purpose of this study, we included the precise ages obtained
406 through institutional collaboration. The impact on average modeling results turned out
407 negligible, however. Data is provided in Volts and has a standard deviation of 9.1 Microvolts.

408 *Data acquisition and processing.* EEG was recorded with 62-channel active ActiCAP
409 electrodes and a bandpass filter between 0.015Hz and 1kHz. We applied additional bandpass
410 filtering between 0.1Hz and 49Hz. The channel placement implemented the 10-5 system
411 (Oostenveld and Praamstra 2001). EEG data were sampled at 2500Hz. After bandpass
412 filtering (0.1 - 49Hz), data were decimated by a factor of 5, yielding a final sampling frequency
413 of 500Hz. As a peculiarity of the dataset, resting-state recordings encompass samples from
414 two conditions: eyes-closed and eyes-open. Our pipeline explicitly respected these different
415 conditions. To include a maximum of data and, potentially, a larger set of distinguishable EEG
416 sources, we pooled the data prior to feature extraction. For the full specification of the
417 preprocessing, please refer to the "config_lemon_eeg.py" file in the code repository.

### CHBP EEG data.

418 The Cuban Human Brain Mapping Project (CHBP) provides rich multimodal EEG and MRI
419 data sampled from young to middle-aged adults from the general population (Valdes-Sosa et
420 al. 2021; Hernandez-Gonzalez et al. 2011; Bosch-Bayard et al. 2020). As for the Cam-CAN
421 and LEMON data, research was carried out using high-end electrophysiological equipment in
422 a biomedical research context. However, the data was collected in a Latin American mid-
423 income country, (Valdes-Sosa et al. 2021), adding a much-needed opportunity for increasing
424 the diversity in population-level neuroscience datasets. This diversity expresses itself in the
425 composition of EEG protocols which contain elements of real-world neurology exams, e.g., a
426 hyperventilation task.

427 *Sample description.* EEG resting-state data from 282 healthy individuals from the
428 CHBP dataset were included in this study. The sample contained 87 females (mean age =
429 36.7) and 195 males (mean age = 29.9), representing a clearly visible difference in the
430 composition of the sample (*Fig*. 2). The overall age distribution went from 18 to 68 years with
431 an average of 32 +/- 9.3 years. Data is provided in Volts and has a standard deviation of 6.6
432 Microvolts. Our sample covers the whole available dataset (download June 2021) as we did
433 not apply any exclusion criteria. Final numbers reflect successful processing of the data.

434 *Data acquisition and processing.* EEG data were recorded using a MEDICID 5 system
435 and two different electrode caps of either 64 or 128 channels. The channel placement
436 implemented the 10-5 system (Oostenveld and Praamstra 2001). Here we focused the
437 analysis on the subset of common channels present in all recordings, leading to 53 channels.
438 We applied additional bandpass filtering between 0.1Hz and 49Hz. As in the LEMON dataset,

439 resting-state recordings encompassed samples from eyes-closed and eyes-open conditions.

440 Again, we pooled both conditions prior to feature extraction. Note that for the data release

441 (downloaded July 2021) used in this work, we could not benefit from the expert-based

442 annotations of clean data. The results obtained on this dataset may therefore be impacted by

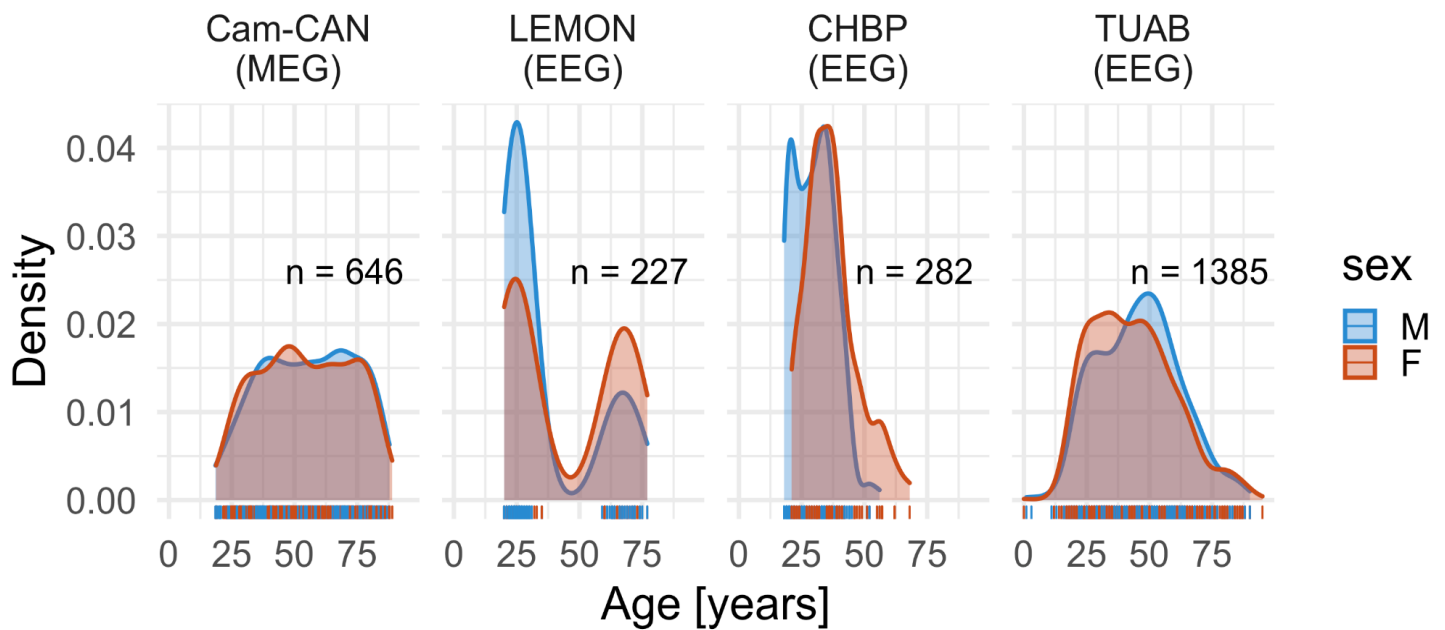443 quality issues to unknown extents.

444 For the full specification of the preprocessing, please refer to the "config_chbp_eeg.py"

445 file in the code repository.

### *TUAB EEG data.*

446 The Temple University Hospital Abnormal EEG Corpus (TUAB) provides socially and

447 ethnically heterogeneous clinical EEG data (Obeid and Picone 2016) mostly from Latin-

448 American and African American participants (personal communication, Joseph Picone). As a

449 peculiarity, the EEG data is obtained from an archival effort of recovering different EEG exams

450 from the Temple University Hospital in Philadelphia. The clinical and social diversity render

451 the TUAB dataset an important resource for electrophysiological population modeling

452 (Gemein et al. 2020; David Sabbagh et al. 2020).

453 *Sample description.* Here, we focused exclusively on the EEG recordings labeled as

454 not pathological by medical experts comprising a subsample of 1385 subjects (female = 775

455 and males = 610). This sample contained individuals ranging from newborn children (min age

456 = 0 for female and min age = 1 for male) to elderly (max age = 95 for female and 90 for male)

457 people (*Fig*. 2). The average age is 44.4 +/- 16.5 years. Data is provided in Volts and has a

458 standard deviation of 9.7 Microvolts. The data processing closely followed our previous work

459 on the TUAB data (Sabbagh et al. 2020). For further details about the dataset, please refer to

460 the reference publications (Harati et al. 2014; Obeid and Picone 2016).

461 *Data acquisition and processing.* EEG data were recorded using different Nicolet EEG

462 devices (Natus Medical Inc.), equipped with between 24 and 36 channels. For channel

463 placement, the 10-5 system was applied (Oostenveld and Praamstra 2001). All sessions have

464 been recorded with an average reference. Here we considered a subset of 21 common

465 channels. As channel numbers differed across recordings, re-referencing was necessary. For

466 consistency, we also applied re-referencing with an average reference on all other EEG

467 datasets. As sampling frequencies were inconsistent across recordings, we resampled the

468 data to 200Hz. For many patients, multiple recordings were available. For simplicity we only

469 considered the first recording. For the full specification of the preprocessing, please refer to

470 the "config_tuab_eeg.py" file in the code repository.

**Figure 2: Age distributions by gender by dataset.** The kernel density (y axis) is plotted across the age range (x axis) for all four M/EEG datasets included in the study, separately for male (blue) and female (red) participants. Individual observations are displayed by rug plots at the bottom of each panel. The Cam-CAN data (MEG) show a wide age range with a quasi-uniform distribution and no obvious sex imbalance. This situation poses no a priori challenges for age prediction while, at the same time, analysis of MEG data may be more complex. The LEMON dataset included a group of young participants and a group of old participants, leading to a characteristic bi-modal distribution. Sex imbalance is clearly visible with more male participants in the group of young participants and fewer male participants in the group of older participants. This may lead to potential sex differences in prediction success and renders the average age a bad summary of the age distribution. The CHBP data shows a rather reduced age range with a right-skewed age distribution and some sex imbalance (again more young male participants). Predicting the age can be expected to turn out more difficult on this dataset for the implied lack of density along the age range. Finally, the TUAB data present a symmetric age distribution with minor sex differences, however, a less uniform age distribution. This may lead to more pronounced errors in young and elderly participants. This may, however, be compensated for by the more generous sample size. To summarize, the four datasets investigated here pose unique challenges for M/EEG brain age modeling.

## Model evaluation and comparison

471    To gauge model performance, we first defined a baseline model that should not provide any

472    intelligent prediction. As in previous work (Sabbagh et al. 2020; Sabbagh et al. 2019;

473    Engemann et al. 2020), we employed a dummy regressor model as a low-level baseline in

474    which the outcome is guessed from the average of the outcome on the training data. This

475     approach is fast and typically converges with more computationally demanding procedures

476     based on permutation testing that we shall briefly outline.

477     This is particularly relevant for the present benchmark where the combinatorial matrix

478     of machine learning models (including deep learning) versus datasets would lead to

479     unpleasant computation times when applying tens of thousands of permutations. The same

480     can be said for other approximations focusing on ranking statistics across hundreds of Monte

481     Carlo cross-validation iterations (Sabbagh et al. 2019). Finally, another approach relies on

482     large left-out datasets, entirely independent from model construction, in which predictions can

483     be treated like random variables, hence, classical inferential statistics are valid. In previous

484     work (Dadi et al. 2021), permutation tests and the non-parametric bootstrap were employed

485     on more than 4000 left-out data points to assess performance above chance and pairwise

486     differences between models. Such generous held-out datasets are not available in the present

487     setting, nor can we readily compute statistics across folds, as cross-validation iterations are

488     not statistically independent. We therefore implemented a less formal approach comparing

489     competing models against dummy regressors and against each other based on standard 10-

490     fold cross-validation based on fixed random seeds. This ensured that for any model under

491     consideration, identical data splits were used. Of note, our reusable benchmark code allows

492     interested readers to implement more exhaustive model comparison strategies.

493     For scoring prediction performance, we focused on two complementary metrics. The

494     coefficient of determination ($R^2$) score and the mean absolute error (MAE). Considering the

495     dummy regressor, the $R^2$ score is a natural choice as it quantifies the incremental success of

496     a model over a regressor returning the average of the training-data as a guess for the

497     outcome. Compared to Pearson correlations that are sometimes used in applied neuroscience

498     studies, the $R^2$ metric is more rigorous as it is sensitive to the scale of the error and the

499     location: Predictions that are entirely biased, e.g, shifted by a large offset, could still be

500     correlated with the outcome. In contrast, the $R^2$ metric clearly penalizes systematically wrong

501     predictions by assigning scores smaller than 0. Positive predictive success thus falls into a

502     range of $R^2$ between 0 and 1. This facilitates comparisons across models within the same

503     dataset while posing challenges when comparing models across datasets.

504     We therefore considered the MAE which has the benefit of expressing prediction errors

505     at the scale of the outcome. This is particularly convenient for scientific interpretation when

506     the outcome has some practical meaning as is the case in the present benchmarks on age

507     prediction. Importantly, the MAE does not per se resolve the problem of comparisons across

508     datasets as the meaning of errors entirely depend on the distribution of the outcome: Small

509     errors in years are good for datasets with wide age distributions but bad in datasets with

510     narrow age distributions. This obviously calls for contextualizing the MAE against a dummy

511     baseline regression model. While this does not necessarily facilitate comparisons across

512  datasets, it helps make visible situations in which one cannot rely solely on the $R^2$ for model

513  comparisons.

**Computational considerations and software**

514  *M/EEG data processing.* BIDS conversion and subsequent data analysis steps were carried

515  out in Python 3.7.1, the MNE-Python software (v0.24, Gramfort et al. 2014, 2013), the MNE-

516  BIDS package (v0.9, Appelhoff et al. 2019) and the MNE-BIDS-pipeline on a 48-core Linux

517  high-performance server with 504 GB RAM. The joblib library (v1.0.1) was used for parallel

518  processing. For artifact removal, the latest development version (v0.3dev) of the autoreject

519  package (Jas et al. 2017) was used.

520        *Classical machine learning benchmarks.* For future computation, the mne-features

521  (0.2, Schiratti, Le Douget, Le Van Quyen, et al. 2018), PyRiemann (v0.2.6) and the coffeine

522  (0.1, Sabbagh et al. 2020) libraries were used. Analyses were composed in custom scripts

523  and library functions based on the Scientific Python Stack with NumPy (v1.19.5, Harris et al.

524  2020), SciPy (v1.6.3, Virtanen et al. 2020) and pandas (v.1.2.4, McKinney and Others 2011).

525  Machine-learning specific computation was composed using the scikit-learn package

526  (Pedregosa et al. 2011). Analysis was carried out on a 48-core Linux high-performance server

527  with 504 GB RAM. Feature extraction, depending on the dataset, completed within several

528  hours to days. Model training and evaluation completed within a few minutes to hours.

529  However, feature computation could last several days, depending on the dataset and the

530  types of features.

531        *Deep learning benchmarks.* A high-performance Linux server with 72 cores, 376 GB

532  RAM and 1 or 2 Nvidia Tesla V100 or P4 GPUs was used. Code was implemented using the

533  PyTorch (Paszke et al. 2019) and braindecode (Schirrmeister et al. 2017) packages. Model

534  training and evaluation completed within 2-3 days.

535        *Data visualization.* Graphical displays and tables were composed on an Apple Silicon

536  M1 Macbook Pro (space gray) in R (v4.0.3 "Bunny-Wunnies Freak Out") using the ggplot2

537  (v3.3.5, Wickham 2011), patchwork (v1.1.1, Pedersen 2019), ggthemes (v4.2.4) and scales

538  (v1.1.1, Arnold 2017) packages with their respective dependencies.
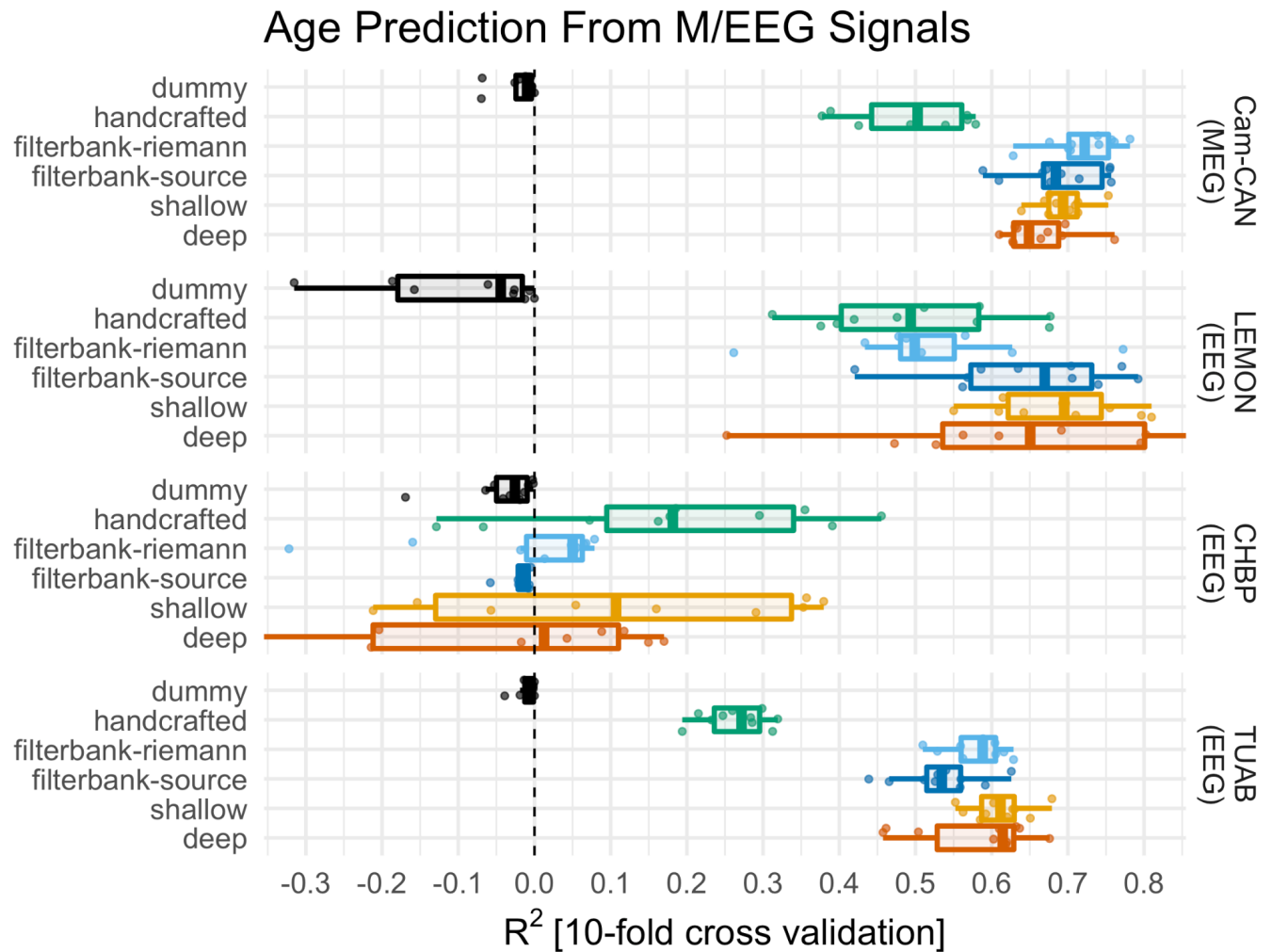
## Results

539  For the age prediction benchmark, we considered five alternative approaches: heterogeneous

540  handcrafted features & random forest ('handcrafted'), filterbank features based on

541  Riemannian embeddings & ridge regression ('filterbank-riemann'), filterbank features based

542     on source localization with MRI-average template & ridge regression ('filterbank-source'), a

543     shallow deep learning architecture ("shallow") and a 4-layer deep-learning architecture

544     ('deep'). These approaches were benchmarked across four M/EEG datasets: The Cambridge

545     Centre of Ageing and Neuroscience (Cam-CAN) dataset (Taylor et al. 2017), the Cuban

546     Human Brain Mapping Project (CHBP) dataset (Valdes-Sosa et al. 2021), the Leipzig Mind-

547     Brain-Body (LEMON) dataset (Babayan et al. 2019) and the Temple University Hospital

548     Abnormal EEG Corpus (TUAB) dataset (Obeid and Picone 2016). Generalization

549     performance was estimated using 10-fold cross validation after shuffling the samples (fixed

550     random seed). The coefficient of determination ($R^2$) was used as a metric enabling

551     comparisons between datasets independently of the age distribution, mathematically

552     quantifying the additional variance explained by predicting better than the average age. A

553     dummy model empirically quantifies chance-level prediction by returning the average age of

554     the training data as prediction. The results are displayed in Fig. 3. One can see that on most

555     of the datasets all machine learning models achieved $R^2$ scores well beyond the dummy

556     baseline. The highest scores were observed on the Cam-CAN MEG dataset, followed by the

557     LEMON EEG dataset. Caution is warranted though to avoid premature conclusions: The $R^2$

558     offers a common scale that explicitly compares the incremental model performance over the

559     average predictor. This is achieved by dividing the sum of squares of the model's prediction

560     by the sum of squares of the average predictor but, in turn, depends on the distribution of age.

561     As a result, this can be misleading in cross-dataset comparisons when the variance of the

562     outcome is not the same, which is the case here (cf. Fig. 2). We therefore also computed

563     results using the mean absolute error as a performance metric (Fig 4). One can now see that

564     the overall distribution of scores, including the scores of the dummy model, depend not only

565     on the dataset but also on its age range. Where the range is small, improvements over the

566     baseline models are harder to observe. Moreover, comparing MAE scores across datasets

567     without taking into account the baseline can yield misleading conclusions. For example, the

568     same score of *e.g.* an MAE = 10 can be way above chance in one dataset (Cam-CAN) but

569     below chance in another dataset (CHBP). To alleviate this problem, normalized MAE scores

570     have been suggested in which the MAE scores are related to the range of the age distribution

571     (Cole, Franke, and Cherbuin 2019). This does not come without its own problems, as then

572     outliers in non-uniform distributions could drive the scores. As research keeps evolving on this

573     topic and the community has not yet agreed on the best metric, we recommend considering

574     multiple classical machine learning metrics when comparing model performance – in critical

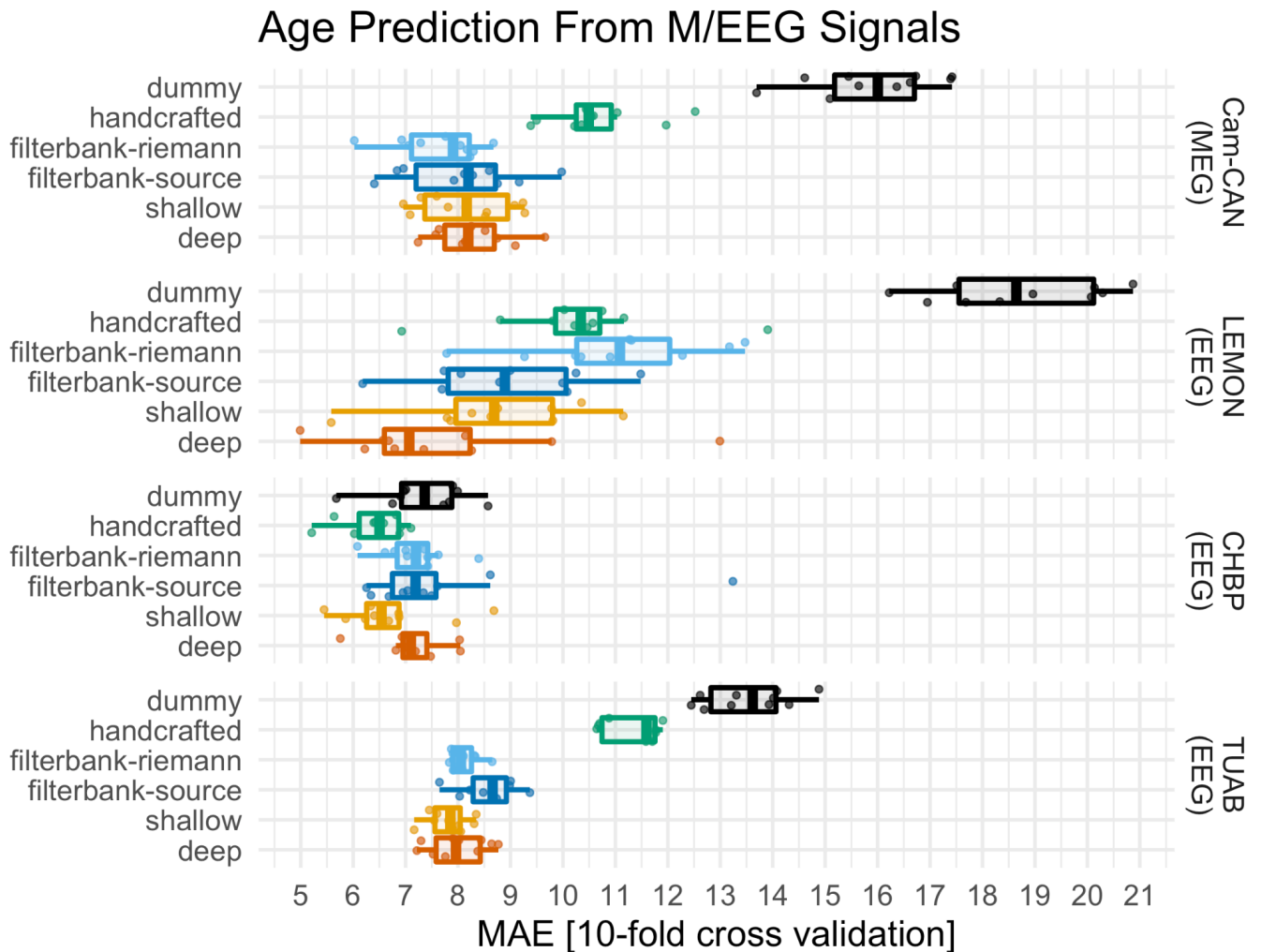575     awareness of their respective limitations.

576         Confronting the relative performances of models to the dummy baseline in Fig. 3 and

577     Fig. 4, one can see overall similar performance rankings between the models, regardless of

578     the metric. See Table 1 for side-by-side comparisons of the aggregated cross-validation

579    distributions. The big-picture results argue in favor of the importance of fine-grained spatial

580    features for M/EEG prediction while considering important between-dataset heterogeneity.

581    Both filterbank pipelines provide features based on spatially aware representations of the

582    M/EEG signals, which either explicitly or implicitly deal with the spatial spread of electrical

583    potentials and fields characteristic for M/EEG signals. The source-level filterbank

584    approximates source localization using the average MRI template, whereas Riemannian

585    embeddings provide non-linear spectral features that are affine invariant, hence, independent

586    of linear mixing. The deep benchmarks, on the other hand, implied spatial-filtering layers

587    capable of mimicking source localization by learning an unmixing function. Surprisingly, using

588    the average MRI template instead of the Riemannian embedding to construct a filterbank

589    model did not lead to consistent improvements across datasets, suggesting that both

590    approaches may be equally effective in practice. We would have conjectured that even an

591    imprecise biophysical head model would provide inverse solutions leading to more accurate

592    unmixing of M/EEG sources. Compared to our previous benchmarks (Engemann et al. 2020;

593    Sabbagh et al. 2020) favoring filterbank models based on source-localization, one has to point

594    out that this finding may reflect at least two differences: The use of an MRI template instead

595    of individual co-registration and the use of empty-room-based suppression of environmental

596    noise. The second factor may be less relevant for EEG though where empty room recordings

597    are not available and data-based covariances are more common in event-related studies

598    where brain activity induced by stimuli is compared against the background resting-state

599    activity. As a practical implication, and if inspection of the brain sources is not a priority, the

600    purely data-driven pipelines may be more practical as no additional MRI-based data

601    processing is needed (cf. Fig. 1).

**Fig. 3. Age prediction benchmarks across M/EEG datasets ($R^2$ score).** Generalization performance was assessed by 10-fold cross-validation and the $R^2$ score for five machine learning strategies compared against a dummy model (rows) and four datasets (panels). Across datasets, dummy models were mostly well-calibrated with $R^2$ scores close to zero. The LEMON dataset was one exception as dummy scores were systematically worse than chance, which can be explained by the bimodal age distribution (cf. Fig. 2), rendering the average age a bad guess for the age. The 'handcrafted' benchmark delivered moderate but systematic prediction success across all datasets. The two filterbank models performed well across datasets with similar performances, markedly higher than for the 'handcrafted' approach. The only exception was the CHBP benchmark for which neither the filterbank nor the deep models delivered useful predictions. Note that here, for the 'filterbank-source', a single fold with an abysmal $R^2$ score of -15 was obtained (x limits constrained to a range between -.3 and 1.0). Overall, the deep learning benchmarks performed similarly to the filterbank models.

## Age Prediction From M/EEG Signals



**Fig. 4. Age prediction benchmarks across M/EEG datasets (mean absolute error).** Same visual conventions as in Fig. 3. As the mean absolute error (MAE) is sensitive to the scale and distribution of the outcome, one can see characteristic differences across datasets. The distribution of the dummy scores provides an estimate of the random guessing. As before, in all but the Cuban datasets all benchmarks achieved MAE scores markedly better than the dummy with no overlap between model and dummy distributions. Model rankings resemble the ones obtained using the $R^2$. On the LEMON data, the deep benchmark now presented a slight advantage over all other benchmarks.

602    Interestingly, none of these approaches involving spatially fine-grained
603    representations of the M/EEG signal worked well on the CHBP data, whereas the random
604    forest on top of hand-crafted features scored systematically better than the dummy baseline.
605    This may be related to three factors that come together in the CHBP benchmark dataset: Like
606    the LEMON dataset, the sample size is relatively small. Second, the age distribution is far less
607    uniform, leading to underrepresentation of elderly participants. This makes the learning task

608    at hand harder as models have fewer training examples from elderly populations. These

609    challenges apply equally to all machine learning benchmarks, hence, do not explain why the

610    random forest model on hand-crafted features is working to some extent. In this context, it

611    may be worthwhile to consider that the CHBP uses two different EEG montages, one with 60,

612    one with 120 electrodes, which may induce strong difference in the covariance structure of

613    the signals due to montage-specific noise structure related to the number of electrodes. This

614    may have affected the random-forest pipeline less strongly as the hand-crafted features

615    extracted marginal channel-wise summary statistics of the time-series or the power spectrum

616    rather than pairwise interactions. Progress on this specific benchmark may therefore involve

617    explicit consideration of the montage when selecting samples for cross-validation or even at

618    the level of the machine learning model (e.g., by including the number of electrodes or

619    montage type as covariate). Moreover, future availability of samples from older populations in

620    the CHBP dataset will help disambiguate this point. Finally, once the expert-based quality-

621    control annotations are considered for epochs-selection, the results obtained in this

622    benchmark may change (see section Datasets/CHBP EEG data for details).

623        A different type of challenge is illustrated by the benchmarks on the LEMON dataset.

624    As the age distribution is bimodal here (Fig. 2), the $R^2$ score is not well calibrated as the

625    average predictor will not provide a reasonable summary of the distribution. This is not

626    automatically mitigated by considering the MAE as a metric. On the other hand, it will not

627    affect the ranking of the machine learning models, which compare overall well to results

628    obtained on the Cam-CAN and the TUAB datasets. To obtain a more rigorous baseline, one

629    could envision a group-wise average predictor that, depending on the age group, would return

630    the groups' respective average age from the training data. We did not implement such a

631    custom baseline here as it was our goal to stick to standard routines provided by the software

632    libraries our benchmarks were based on. Second, it was our intention to expose such issues

633    as this may stimulate future research and development.

**Table 1.** *Aggregate cross-validation results across benchmarks and datasets*

| dataset | benchmark | $R^2_{(M)}$ | $R^2_{(SD)}$ | $MAE_{(M)}$ | $MAE_{(SD)}$ |
|---|---|---|---|---|---|
| Cam-CAN (MEG) | deep | 0.66 | 0.05 | 8.29 | 0.74 |
| Cam-CAN (MEG) | shallow | 0.69 | 0.03 | 8.14 | 0.90 |
| Cam-CAN (MEG) | filterbank-source | 0.69 | 0.06 | 8.10 | 1.11 |
| Cam-CAN (MEG) | filterbank-riemann | 0.72 | 0.05 | 7.65 | 0.81 |
| Cam-CAN (MEG) | handcrafted | 0.49 | 0.07 | 10.65 | 0.98 |
| Cam-CAN (MEG) | dummy | -0.02 | 0.03 | 15.90 | 1.22 |
| LEMON (EEG) | deep | 0.65 | 0.20 | 7.78 | 2.25 |
| LEMON (EEG) | shallow | 0.69 | 0.08 | 8.80 | 1.58 |
| LEMON (EEG) | filterbank-source | 0.65 | 0.12 | 8.93 | 1.56 |

| | | | | | |
|---|---|---|---|---|---|
| LEMON (EEG) | filterbank-riemann | 0.51 | 0.13 | 11.00 | 1.73 |
| LEMON (EEG) | handcrafted | 0.50 | 0.13 | 10.26 | 1.76 |
| LEMON (EEG) | dummy | -0.13 | 0.17 | 18.70 | 1.60 |
| CHBP (EEG) | deep | -0.10 | 0.28 | 7.14 | 0.65 |
| CHBP (EEG) | shallow | 0.03 | 0.38 | 6.74 | 0.96 |
| CHBP (EEG) | filterbank-source | -1.49 | 4.67 | 7.76 | 2.05 |
| CHBP (EEG) | filterbank-riemann | -0.01 | 0.13 | 7.17 | 0.63 |
| CHBP (EEG) | handcrafted | 0.19 | 0.19 | 6.40 | 0.61 |
| CHBP (EEG) | dummy | -0.04 | 0.05 | 7.33 | 0.83 |
| TUAB (EEG) | deep | 0.58 | 0.08 | 7.99 | 0.55 |
| TUAB (EEG) | shallow | 0.61 | 0.04 | 7.82 | 0.38 |
| TUAB (EEG) | filterbank-source | 0.53 | 0.06 | 8.58 | 0.51 |
| TUAB (EEG) | filterbank-riemann | 0.58 | 0.04 | 8.10 | 0.26 |
| TUAB (EEG) | handcrafted | 0.26 | 0.04 | 11.32 | 0.53 |
| TUAB (EEG) | dummy | -0.01 | 0.01 | 13.55 | 0.82 |

## Discussion

634  In this study, we proposed empirical benchmarks for age prediction comparing distinct
635  machine learning approaches across diverse M/EEG datasets comprising, in total, more than
636  2500 recordings. The benchmarks were implemented in Python based on the MNE-software
637  ecosystem, the Braindecode package and the BIDS data standard. The explicit reliance on
638  the BIDS standard renders these pipelines applicable to any M/EEG data presented in the
639  BIDS format. This enabled coherent side-by-side comparisons of classical machine learning
640  models and deep learning methods across M/EEG datasets recorded in different research or
641  medical contexts.

642      Our cross-dataset and cross-model benchmarks pointed out stable ranking of model
643  performance across two metrics, the $R^2$ score and mean absolute error (MAE). $R^2$ scores
644  have been less consistently reported in the literature, however, the top MAE scores observed
645  across datasets in this benchmark of 7 to 8 years are well in line with reports from previous
646  publications (Sun et al. 2019; Sabbagh et al. 2020; Xifra-Porxas et al. 2021). While direct
647  comparisons against MRI were not performed in this study, the present benchmarks would be
648  compatible with the impression that for what concerns the overall performance of age
649  prediction, M/EEG features are slightly weaker than MRI features (Engemann et al. 2020;
650  Xifra-Porxas et al. 2021). We found that, overall, Riemannian filterbank models and deep
651  learning models achieved the highest scores, whereas random forests based on hand-crafted
652  features delivered robust performance. In line with previous work (Gemein et al. 2020), these
653  results suggest that deep learning methods do not necessarily show a consistent advantage
654  over classical pipeline models: Similar performance may be explained by the fact that our

655    filterbank models and the deep models imply similar spatially aware representations of the

656    M/EEG data (see results section for detailed discussion in context). Moreover, given the

657    relatively small training datasets, it can be considered good news that these parameter-rich

658    models did not seem to overfit as was evidenced by comparisons against simpler classical

659    models. Yet, it may be simply a matter of collecting larger samples until deep learning

660    approaches may reveal their advantage at extracting more elaborate representations of

661    M/EEG signals. This may lead to positioning M/EEG-based brain age prediction on par with

662    MRI-based brain age prediction just as MRI-based deep learning models of brain age have

663    defined state-of-the-art performance on large datasets (Cole et al. 2017; Bashyam et al. 2020;

664    Jonsson et al. 2019). However, more importantly, the value of M/EEG-derived brain age

665    models should not be defined in terms of incremental improvement over MRI-based models

666    as M/EEG-based models may enhance MRI-derived information (Engemann et al. 2020) or

667    may be the only option available (Sun et al. 2019).

668    Our results nicely demonstrate a second critical merit of cross-model and cross-

669    dataset benchmarking. It was sufficient to analyze four different sources of data until we found

670    a perfectly legitimate EEG dataset from an academic research context (CHBP) in which our

671    previously favored modeling techniques developed on the Cam-CAN and the TUAB data did

672    not perform well by default. There may be good reasons for these discrepancies related to the

673    age distribution found in the CHBP data and the fact that multiple different EEG montages

674    were used in that dataset (see results section for detailed discussion in context). But more

675    importantly, we did not anticipate this to happen and would have never learned about it had

676    we confined the scope to previously analyzed datasets. Such discoveries are favored by

677    systematic benchmarks with dataset-independent code implementation, which has the

678    potential to lower the burden threshold for including always more datasets into model

679    development. In the long run, we hope that this effort will stimulate new research leading to

680    more generalizable models.

681    This brings us to some limitations of this work. Our work has been motivated by the

682    absolute necessity to diversify datasets for development of M/EEG-based measures of brain

683    health. This has led us to analyzing more than 2500 M/EEG recordings and, yet we only

684    included four datasets. Other M/EEG datasets come to mind that would have been potentially

685    relevant. The Human Connectome Project MEG data (Larson-Prior et al. 2013) includes MEG

686    recordings from less than 100 participants, which we deemed insufficient for predictive

687    regression modeling. The OMEGA data resource (Niso et al. 2016) was not accessible at the

688    time of this investigation but would have been a good match for this study. Finally, the LIFE

689    cohort (Loeffler et al. 2015) includes a large number of EEG recordings of participants

690    sampled from the general population yet follows a closed / controlled access scheme. The

691    Healthy-Brain-Network EEG data (Alexander et al. 2017) concerns a developmental cohort.

692    Despite potentially relevant similarities between brain development and aging, age prediction
693    in developmental cohorts would have exceeded the scope of the present study. Even if we
694    had integrated these resources in the present benchmark, this may have only marginally
695    enhanced the diversity covered by the current selection datasets as most public neuroscience
696    datasets come from the wealthiest nations. We hope that this situation will improve as new
697    promising international consortia and efforts emerge that focus on curating large EEG
698    datasets from diverse national and cultural contexts (Ibanez et al. 2021; Shekh Ibrahim et al.
699    2020; "Global Brain Consortium Homepage"). A second limitation of the present study
700    concerns the depth of validation. To advance our understanding of M/EEG-derived brain age,
701    more systematic comparisons against MRI-derived brain age (Xifra-Porxas et al. 2021) and
702    other measures of mental health and cognitive function are important objectives (Anatürk et
703    al. 2021; Dadi et al. 2021).

704    In the following we wish to point out a few imminent opportunities for turning the
705    limitations of the present work into future research projects, potentially, enabled by the results
706    and tools brought by the current benchmarks.


**Opportunities and suggestions for follow-up research using the benchmark tools**


707    *Model averaging.* In many instances, combining prediction models using model averaging
708    approaches can improve the prediction performance (O'Connor et al. 2021; Dadi et al. 2021;
709    Varoquaux et al. 2017). This could also be a practical way of combining the benchmarks into
710    a single model for subsequent generalization testing. Future studies could use this benchmark
711    to investigate model averaging approaches.


712    *The impact of deeper architectures.* An important design decision in deep neural networks is
713    the total depth of the neural network. Here we used previously published architectures
714    designed for EEG-based pathology decoding (Schirrmeister et al. 2017). Future studies could
715    build on top of this benchmark to explore the importance of deep architectures for brain age
716    modeling. Specifically, it would be possible to use methods from neural architecture search,
717    e.g., AutoPyTorch (Zimmer, Lindauer, and Hutter 2021), to design better-performing
718    architectures. Since this benchmark does not only provide access to diverse datasets in an
719    identical file format, but also enables direct comparison to others, it is the optimal starting point
720    for such an optimization while at the same time avoiding overfitting the architecture to a single
721    dataset.


722    *The role of preprocessing.* While data cleaning is of major importance for extracting
723    physiologically interpretable biomarkers, predictions from machine learning models tend to be

724    far less affected by noise (Sabbagh et al. 2020). On the other hand, artifacts and noise may

725    inform predictions, potentially reducing biological specificity. Future studies could benefit from

726    this benchmark to quantify the role of artifact signals for brain age predictions and develop

727    de-confounding strategies (Du et al. 2021; Mehrabi et al. 2021; Lu, Schölkopf, and Hernández-

728    Lobato 2018; Bica, Alaa, and Van Der Schaar 2020).

729

730    *Eyes-open versus eyes-closed.* Some of the datasets analyzed in this benchmark contain

731    resting-state signals under different conditions. In the lack of strong a-priori hypotheses, here

732    we simply pooled both conditions. It is currently unclear whether the relationship between

733    eyes-closed versus eyes-open resting-state may contain valuable information about brain

734    aging. It is imaginable, however, that signals induced by transient visual deprivation may

735    reveal levels of vigilance (Wong, DeYoung, and Liu 2016), which in turn may be altered by

736    neuropsychiatric conditions (Hegerl et al. 2012). Future work could benefit from the

737    benchmark to investigate the importance of eyes-closed versus eyes-open resting-state for

738    brain age modeling.

739    *Model inspection.* The interpretability of machine learning models is essential for clinical

740    impact (Rudin 2019; Ghassemi, Oakden-Rayner, and Beam 2021). This benchmark did not

741    cover methods for explaining the role of variable importance for model predictions. Future

742    work could validate the relative importance of M/EEG signals or features for brain age

743    modeling.

744    *Exploring the link with MRI and cognitive scores.* This study established the tools and methods

745    for basic benchmarks on prediction performance. However, to build useful brain age models,

746    it is essential to validate brain-age predictions to cognitive function, measures of health or

747    clinical endpoints (Dadi et al. 2021; Cole et al. 2018; Liem et al. 2017). To further establish

748    the relative merit of M/EEG over MRI, comparisons between the modalities are essential

749    (Engemann et al. 2020). Most of the datasets covered in this benchmark include MRI data,

750    social details and psychometric scores next to the M/EEG data, providing a wealth of

751    opportunities for deep cross-dataset validation of brain age measures.

## Conclusion

752    Computational benchmarks across M/EEG datasets and machine learning methods bear the

753    potential to enhance applications of machine learning in clinical neuroscience in several ways.

754    Standardization of data formats, software and analysis pipelines are important factors for the

755 scalability of predictive modeling of M/EEG. For stimulating the development of more
756 generalizable machine learning models it is crucial that a critical mass of M/EEG datasets be
757 analyzed by the international community. As the diversity of the datasets increases,
758 generalization gaps will manifest themselves, calling for computation methods for closing
759 these gaps. The implied learning process may eventually lead to developing more widely
760 applicable M/EEG-based biomarkers that are clinically robust across a wide range of
761 sociocultural contexts, clinical populations, recording sites and measurement techniques. We
762 hope that benchmarks, tools and resources resulting from this study will facilitate investigating
763 open scientific questions related to learning biomarkers of brain health on an ever-growing
764 number of M/EEG datasets from increasingly diverse real-world contexts.

## Acknowledgements

## Author contributions

authors in alphabetical order

**Conceptualization**: A.G., D.E.
**Data curation**: A.G., A.M., D.E., H.B., L.G.
**Software**: A.G., A.M., D.E., D.S., H.B., L.G., R.H.
**Formal analysis**: A.G., D.E.
**Supervision**: A.G., D.E.
**Funding acquisition**: A.G., D.E.
**Validation**: A.G., D.E.

**Investigation**: A.G., A.M., D.E., H.B., L.G.

**Visualization**: A.P., D.E.

**Methodology**: A.G., D.E.

**Project administration**: D.E.

**Writing—original draft**: D.E.

**Writing—review and editing**: A.G., A.M., D.E., D.S., H.B., L.G., T.B.

## Declaration of conflicts of interest

D.E. is a full-time employee of F. Hoffmann-La Roche Ltd.
H.B. receives graduate funding support from InteraXon Inc.

## References

Alexander, Lindsay M., Jasmine Escalera, Lei Ai, Charissa Andreotti, Karina Febre, Alexander Mangone, Natan Vega-Potler, et al. 2017. "An Open Resource for Transdiagnostic Research in Pediatric Mental Health and Learning Disorders." *Scientific Data* 4 (December): 170181.

Al Zoubi, Obada, Chung Ki Wong, Rayus T. Kuplicki, Hung-Wen Yeh, Ahmad Mayeli, Hazem Refai, Martin Paulus, and Jerzy Bodurka. 2018. "Predicting Age From Brain EEG Signals—A Machine Learning Approach." *Frontiers in Aging Neuroscience* 10: 184.

Anatürk, Melis, Tobias Kaufmann, James H. Cole, Sana Suri, Ludovica Griffanti, Enikő Zsoldos, Nicola Filippini, et al. 2021. "Prediction of Brain Age and Cognitive Age: Quantifying Brain and Cognitive Maintenance in Aging." *Human Brain Mapping* 42 (6): 1626–40.

Ang, Kai Keng, Zheng Yang Chin, Haihong Zhang, and Cuntai Guan. 2008. "Filter Bank Common Spatial Pattern (FBCSP) in Brain-Computer Interface." In *2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence)*, 2390–97.

Appelhoff, Stefan, Matthew Sanderson, Teon L. Brooks, Marijn van Vliet, Romain Quentin, Chris Holdgraf, Maximilien Chaumon, et al. 2019. "MNE-BIDS: Organizing Electrophysiological Data into the BIDS Format and Facilitating Their Analysis." *The Journal of Open Source Software* 4 (44). https://pure.mpg.de/rest/items/item_3192645/component/file_3192646/content.

Arnold, Jeffrey B. 2017. "Ggthemes: Extra Themes, Scales and Geoms for 'ggplot2.'" *R Package Version* 3 (0).

Babayan, Anahit, Miray Erbey, Deniz Kumral, Janis D. Reinelt, Andrea M. F. Reiter, Josefin Röbbig, H. Lina Schaare, et al. 2019. "A Mind-Brain-Body Dataset of MRI, EEG, Cognition, Emotion, and Peripheral Physiology in Young and Old Adults." *Scientific Data* 6 (February): 180308.

Babiloni, Claudio, Giuliano Binetti, Andrea Cassarino, Gloria Dal Forno, Claudio Del Percio, Florinda Ferreri, Raffaele Ferri, et al. 2006. "Sources of Cortical Rhythms in Adults during Physiological Aging: A Multicentric EEG Study." *Human Brain Mapping* 27 (2): 162–72.

Banville, Hubert, Omar Chehab, Aapo Hyvarinen, Denis Engemann, and Alexandre Gramfort. 2020. "Uncovering the Structure of Clinical EEG Signals with Self-Supervised

Learning." *Journal of Neural Engineering*, November. https://doi.org/10.1088/1741-2552/abca18.

Banville, Hubert, Sean U. N. Wood, Chris Aimone, Denis-Alexander Engemann, and Alexandre Gramfort. 2021. "Robust Learning from Corrupted EEG with Dynamic Spatial Filtering." *arXiv [cs.LG]*. arXiv. http://arxiv.org/abs/2105.12916.

Bao, Pinglei, Liang She, Mason McGill, and Doris Y. Tsao. 2020. "A Map of Object Space in Primate Inferotemporal Cortex." *Nature* 583 (7814): 103–8.

Barachant, Alexandre, Stéphane Bonnet, Marco Congedo, and Christian Jutten. 2012. "Multiclass Brain-Computer Interface Classification by Riemannian Geometry." *IEEE Transactions on Bio-Medical Engineering* 59 (4): 920–28.

Bashyam, Vishnu M., Guray Erus, Jimit Doshi, Mohamad Habes, Ilya Nasrallah, Monica Truelove-Hill, Dhivya Srinivasan, et al. 2020. "MRI Signatures of Brain Age and Disease over the Lifespan Based on a Deep Brain Network and 14 468 Individuals Worldwide." *Brain: A Journal of Neurology* 143 (7): 2312–24.

Bica, Ioana, Ahmed Alaa, and Mihaela Van Der Schaar. 2020. "Time Series Deconfounder: Estimating Treatment Effects over Time in the Presence of Hidden Confounders." In *Proceedings of the 37th International Conference on Machine Learning*, edited by Hal Daumé Iii and Aarti Singh, 119:884–95. Proceedings of Machine Learning Research. PMLR.

Bosch-Bayard, Jorge, Lidice Galan, Eduardo Aubert Vazquez, Trinidad Virues Alba, and Pedro A. Valdes-Sosa. 2020. "Resting State Healthy EEG: The First Wave of the Cuban Normative Database." *Frontiers in Neuroscience* 14 (December): 555119.

Breiman, Leo. 2001. "Random Forests." *Machine Learning* 45 (1): 5–32.

Buitinck, Lars, Gilles Louppe, Mathieu Blondel, Fabian Pedregosa, Andreas Mueller, Olivier Grisel, Vlad Niculae, et al. 2013. "API Design for Machine Learning Software: Experiences from the Scikit-Learn Project." *arXiv [cs.LG]*. arXiv. http://arxiv.org/abs/1309.0238.

Bycroft, Clare, Colin Freeman, Desislava Petkova, Gavin Band, Lloyd T. Elliott, Kevin Sharp, Allan Motyer, et al. 2018. "The UK Biobank Resource with Deep Phenotyping and Genomic Data." *Nature* 562 (7726): 203–9.

Cabeza, Roberto, Nicole D. Anderson, Jill K. Locantore, and Anthony R. McIntosh. 2002. "Aging Gracefully: Compensatory Brain Activity in High-Performing Older Adults." *NeuroImage* 17 (3): 1394–1402.

Chambon, Stanislas, Mathieu N. Galtier, Pierrick J. Arnal, Gilles Wainrib, and Alexandre Gramfort. 2018. "A Deep Learning Architecture for Temporal Sleep Stage Classification Using Multivariate and Multimodal Time Series." *IEEE Transactions on Neural Systems and Rehabilitation Engineering: A Publication of the IEEE Engineering in Medicine and Biology Society* 26 (4): 758–69.

Cheveigné, Alain de, and Lucas C. Parra. 2014. "Joint Decorrelation, a Versatile Tool for Multichannel Data Analysis." *NeuroImage* 98 (September): 487–505.

Choy, Tricia, Elizabeth Baker, and Katherine Stavropoulos. 2021. "Systemic Racism in EEG Research: Considerations and Potential Solutions." *Affective Science*, May. https://doi.org/10.1007/s42761-021-00050-0.

Cole, James H. 2020. "Multimodality Neuroimaging Brain-Age in UK Biobank: Relationship to Biomedical, Lifestyle, and Cognitive Factors." *Neurobiology of Aging* 92 (August): 34–42.

Cole, James H., and Katja Franke. 2017. "Predicting Age Using Neuroimaging: Innovative Brain Ageing Biomarkers." *Trends in Neurosciences* 40 (12): 681–90.

Cole, James H., Katja Franke, and Nicolas Cherbuin. 2019. "Quantification of the Biological Age of the Brain Using Neuroimaging." In *Biomarkers of Human Aging*, edited by Alexey Moskalev, 293–328. Cham: Springer International Publishing.

Cole, James H., Rudra P. K. Poudel, Dimosthenis Tsagkrasoulis, Matthan W. A. Caan, Claire Steves, Tim D. Spector, and Giovanni Montana. 2017. "Predicting Brain Age with Deep Learning from Raw Imaging Data Results in a Reliable and Heritable Biomarker." *NeuroImage* 163 (December): 115–24.

Cole, J. H., S. J. Ritchie, M. E. Bastin, M. C. Valdés Hernández, S. Muñoz Maniega, N. Royle, J. Corley, et al. 2018. "Brain Age Predicts Mortality." *Molecular Psychiatry* 23 (5): 1385–92.

Dadi, Kamalaker, Gaël Varoquaux, Josselin Houenou, Danilo Bzdok, Bertrand Thirion, and Denis Engemann. 2021. "Population Modeling with Machine Learning Can Enhance Measures of Mental Health." *GigaScience* 10 (10). https://doi.org/10.1093/gigascience/giab071.

Dähne, Sven, Frank C. Meinecke, Stefan Haufe, Johannes Höhne, Michael Tangermann, Klaus-Robert Müller, and Vadim V. Nikulin. 2014. "SPoC: A Novel Framework for Relating the Amplitude of Neuronal Oscillations to Behaviorally Relevant Parameters." *NeuroImage* 86 (February): 111–22.

Damoiseaux, J. S., C. F. Beckmann, E. J. Sanz Arigita, F. Barkhof, Ph Scheltens, C. J. Stam, S. M. Smith, and S. A. R. B. Rombouts. 2008. "Reduced Resting-State Brain Activity in the 'Default Network' in Normal Aging." *Cerebral Cortex* 18 (8): 1856–64.

Denissen, Stijn, Denis Alexander Engemann, Alexander De Cock, Lars Costers, Johan Baijot, Jorne Laton, Iris-Katharina Penner, et al. 2021. "Brain Age as a Surrogate Marker for Information Processing Speed in Multiple Sclerosis." *medRxiv*.

Devarajan, Kavya, S. Bagyaraj, Vinitha Balasampath, Jyostna E., and Jayasri K. 2014. "EEG-Based Epilepsy Detection and Prediction." *IACSIT International Journal of Engineering and Technology* 6 (3): 212–16.

Dosenbach, Nico U. F., Binyam Nardos, Alexander L. Cohen, Damien A. Fair, Jonathan D. Power, Jessica A. Church, Steven M. Nelson, et al. 2010. "Prediction of Individual Brain Maturity Using fMRI." *Science* 329 (5997): 1358–61.

Driscoll, I., C. Davatzikos, Y. An, X. Wu, D. Shen, M. Kraut, and S. M. Resnick. 2009. "Longitudinal Pattern of Regional Brain Volume Change Differentiates Normal Aging from MCI." *Neurology* 72 (22): 1906–13.

Du, Mengnan, Fan Yang, Na Zou, and Xia Hu. 2021. "Fairness in Deep Learning: A Computational Perspective." *IEEE Intelligent Systems* 36 (4): 25–34.

Duncan, L., H. Shen, B. Gelaye, J. Meijsen, K. Ressler, M. Feldman, R. Peterson, and B. Domingue. 2019. "Analysis of Polygenic Risk Score Usage and Performance in Diverse Human Populations." *Nature Communications* 10 (1): 3328.

Engemann, Denis A., Oleh Kozynets, David Sabbagh, Guillaume Lemaître, Gael Varoquaux, Franziskus Liem, and Alexandre Gramfort. 2020. "Combining Magnetoencephalography with Magnetic Resonance Imaging Enhances Learning of Surrogate-Biomarkers." *eLife* 9 (May). https://doi.org/10.7554/eLife.54055.

Engemann, Denis A., Federico Raimondo, Jean-Rémi King, Benjamin Rohaut, Gilles Louppe, Frédéric Faugeras, Jitka Annen, et al. 2018. "Robust EEG-Based Cross-Site and Cross-Protocol Classification of States of Consciousness." *Brain: A Journal of Neurology* 141 (11): 3179–92.

Esteller, R., J. Echauz, T. Tcheng, B. Litt, and B. Pless. 2001. "Line Length: An Efficient Feature for Seizure Onset Detection." In *2001 Conference Proceedings of the 23rd Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, 2:1707–10 vol.2.

Esteller, R., G. Vachtsevanos, J. Echauz, and B. Litt. 2001. "A Comparison of Waveform Fractal Dimension Algorithms." *IEEE Transactions on Circuits and Systems I: Fundamental Theory and Applications* 48 (2): 177–83.

Ewers, Michael, Reisa A. Sperling, William E. Klunk, Michael W. Weiner, and Harald Hampel. 2011. "Neuroimaging Markers for the Prediction and Early Diagnosis of Alzheimer's Disease Dementia." *Trends in Neurosciences* 34 (8): 430–42.

Ferrucci, Luigi, Marta Gonzalez-Freire, Elisa Fabbri, Eleanor Simonsick, Toshiko Tanaka, Zenobia Moore, Shabnam Salimi, Felipe Sierra, and Rafael de Cabo. 2020. "Measuring Biological Aging in Humans: A Quest." *Aging Cell* 19 (2): e13080.

Fischl, Bruce. 2012. "FreeSurfer." *NeuroImage* 62 (2): 774–81.

Fry, Anna, Thomas J. Littlejohns, Cathie Sudlow, Nicola Doherty, Ligia Adamska, Tim Sprosen, Rory Collins, and Naomi E. Allen. 2017. "Comparison of Sociodemographic

and Health-Related Characteristics of UK Biobank Participants with Those of the General Population." *American Journal of Epidemiology* 186 (9): 1026–34.

Garcés, Pilar, David López-Sanz, Fernando Maestú, and Ernesto Pereda. 2017. "Choice of Magnetometers and Gradiometers after Signal Space Separation." *Sensors* 17 (12). https://doi.org/10.3390/s17122926.

Gaubert, Sinead, Federico Raimondo, Marion Houot, Marie-Constance Corsi, Lionel Naccache, Jacobo Diego Sitt, Bertrand Hermann, et al. 2019. "EEG Evidence of Compensatory Mechanisms in Preclinical Alzheimer's Disease." *Brain: A Journal of Neurology* 142 (7): 2096–2112.

Gemein, Lukas A. W., Robin T. Schirrmeister, Patryk Chrabąszcz, Daniel Wilson, Joschka Boedecker, Andreas Schulze-Bonhage, Frank Hutter, and Tonio Ball. 2020. "Machine-Learning-Based Diagnostics of EEG Pathology." *NeuroImage* 220 (October): 117021.

Ghassemi, Marzyeh, Luke Oakden-Rayner, and Andrew L. Beam. 2021. "The False Hope of Current Approaches to Explainable Artificial Intelligence in Health Care." *The Lancet. Digital Health* 3 (11): e745–50.

"Global Brain Consortium Homepage." n.d. Accessed November 30, 2021. https://globalbrainconsortium.org/.

Gonneaud, Julie, Alex T. Baria, Alexa Pichet Binette, Brian A. Gordon, Jasmeer P. Chhatwal, Carlos Cruchaga, Mathias Jucker, et al. 2021. "Accelerated Functional Brain Aging in Pre-Clinical Familial Alzheimer's Disease." *Nature Communications* 12 (1): 1–17.

Gorgolewski, Krzysztof J., Tibor Auer, Vince D. Calhoun, R. Cameron Craddock, Samir Das, Eugene P. Duff, Guillaume Flandin, et al. 2016. "The Brain Imaging Data Structure, a Format for Organizing and Describing Outputs of Neuroimaging Experiments." *Scientific Data* 3 (June): 160044.

Gramfort, Alexandre, Martin Luessi, Eric Larson, Denis A. Engemann, Daniel Strohmeier, Christian Brodbeck, Roman Goj, et al. 2013. "MEG and EEG Data Analysis with MNE-Python." *Frontiers in Neuroscience* 7 (December): 267.

Gramfort, Alexandre, Martin Luessi, Eric Larson, Denis A. Engemann, Daniel Strohmeier, Christian Brodbeck, Lauri Parkkonen, and Matti S. Hämäläinen. 2014. "MNE Software for Processing MEG and EEG Data." *NeuroImage* 86 (February): 446–60.

Güntekin, Bahar, Tuba Aktürk, Xianghong Arakaki, Laura Bonanni, Claudio Del Percio, Rebecca Edelmayer, Francesca Farina, et al. 2021. "Are There Consistent Abnormalities in Event-related EEG Oscillations in Patients with Alzheimer's Disease Compared to Other Diseases Belonging to Dementia?" *Psychophysiology*, August. https://doi.org/10.1111/psyp.13934.

Harati, A., S. López, I. Obeid, J. Picone, M. P. Jacobson, and S. Tobochnik. 2014. "The TUH EEG CORPUS: A Big Data Resource for Automated EEG Interpretation." In *2014 IEEE Signal Processing in Medicine and Biology Symposium (SPMB)*, 1–5.

Harris, Charles R., K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, et al. 2020. "Array Programming with NumPy." *Nature* 585 (7825): 357–62.

Hastie, Trevor, Robert Tibshirani, Jerome Friedman, and James Franklin. 2005. "The Elements of Statistical Learning: Data Mining, Inference and Prediction." *The Mathematical Intelligencer* 27 (2): 83–85.

Hegerl, Ulrich, Kathrin Wilk, Sebastian Olbrich, Peter Schoenknecht, and Christian Sander. 2012. "Hyperstable Regulation of Vigilance in Patients with Major Depressive Disorder." *The World Journal of Biological Psychiatry: The Official Journal of the World Federation of Societies of Biological Psychiatry* 13 (6): 436–46.

Henrich, J., and S. Heine. 2010. "The Weirdest People in the World?" *Behavioral and Brain Sciences 33 (2-3): 61-83*.

Hernandez-Gonzalez, Gertrudis, Maria L. Bringas-Vega, Lidice Galán-Garcia, Jorge Bosch-Bayard, Yenisleidy Lorenzo-Ceballos, Lester Melie-Garcia, Lourdes Valdes-Urrutia, Marcia Cobas-Ruiz, Pedro A. Valdes-Sosa, and Cuban Human Brain Mapping Project (CHBMP). 2011. "Multimodal Quantitative Neuroimaging Databases and Methods: The

Cuban Human Brain Mapping Project." *Clinical EEG and Neuroscience: Official Journal of the EEG and Clinical Neuroscience Society (ENCS)* 42 (3): 149–59.

He, Tong, Ru Kong, Avram J. Holmes, Minh Nguyen, Mert R. Sabuncu, Simon B. Eickhoff, Danilo Bzdok, Jiashi Feng, and B. T. Thomas Yeo. 2020. "Deep Neural Networks and Kernel Regression Achieve Comparable Accuracies for Functional Connectivity Prediction of Behavior and Demographics." *NeuroImage* 206 (February): 116276.

Ibanez, Agustin, Mario A. Parra, Christopher Butler, and Latin America and the Caribbean Consortium on Dementia (LAC-CD). 2021. "The Latin America and the Caribbean Consortium on Dementia (LAC-CD): From Networking to Research to Implementation Science." *Journal of Alzheimer's Disease: JAD* 82 (s1): S379–94.

Inouye, T., K. Shinosaki, H. Sakamoto, S. Toi, S. Ukai, A. Iyama, Y. Katsuda, and M. Hirano. 1991. "Quantification of EEG Irregularity by Use of the Entropy of the Power Spectrum." *Electroencephalography and Clinical Neurophysiology* 79 (3): 204–10.

Jas, Mainak, Denis A. Engemann, Yousra Bekhti, Federico Raimondo, and Alexandre Gramfort. 2017. "Autoreject: Automated Artifact Rejection for MEG and EEG Data." *NeuroImage* 159 (October): 417–29.

Jas, Mainak, Eric Larson, Denis A. Engemann, Jaakko Leppäkangas, Samu Taulu, Matti Hämäläinen, and Alexandre Gramfort. 2018. "A Reproducible MEG/EEG Group Study With the MNE Software: Recommendations, Quality Assessments, and Good Practices." *Frontiers in Neuroscience* 12: 530.

Jayaram, Vinay, and Alexandre Barachant. 2018. "MOABB: Trustworthy Algorithm Benchmarking for BCIs." *Journal of Neural Engineering* 15 (6): 066011.

Jonsson, B. A., G. Bjornsdottir, T. E. Thorgeirsson, L. M. Ellingsen, G. Bragi Walters, D. F. Gudbjartsson, H. Stefansson, K. Stefansson, and M. O. Ulfarsson. 2019. "Brain Age Prediction Using Deep Learning Uncovers Associated Sequence Variants." *Nature Communications* 10 (1): 5409.

Kernbach, Julius M., B. T. Thomas Yeo, Jonathan Smallwood, Daniel S. Margulies, Michel Thiebaut de Schotten, Henrik Walter, Mert R. Sabuncu, et al. 2018. "Subspecialization within Default Mode Nodes Characterized in 10,000 UK Biobank Participants." *Proceedings of the National Academy of Sciences of the United States of America* 115 (48): 12295–300.

Khan, Sheraz, Javeria A. Hashmi, Fahimeh Mamashli, Konstantinos Michmizos, Manfred G. Kitzbichler, Hari Bharadwaj, Yousra Bekhti, et al. 2018. "Maturation Trajectories of Cortical Resting-State Networks Depend on the Mediating Frequency Band." *NeuroImage* 174 (July): 57–68.

Kietzmann, Tim C., Patrick McClure, and Nikolaus Kriegeskorte. 2019. "Deep Neural Networks in Computational Neuroscience." In *Oxford Research Encyclopedia of Neuroscience*.

King, J-R, and S. Dehaene. 2014. "Characterizing the Dynamics of Mental Representations: The Temporal Generalization Method." *Trends in Cognitive Sciences* 18 (4): 203–10.

King, J-R, L. Gwilliams, C. Holdgraf, and J. Sassenhagen. 2018. "Encoding and Decoding Neuronal Dynamics: Methodological Framework to Uncover the Algorithms of Cognition." https://hal.archives-ouvertes.fr/hal-01848442/.

King, Kevin S., Julia Kozlitina, Roger N. Rosenberg, Ronald M. Peshock, Roderick W. McColl, and Christine K. Garcia. 2014. "Effect of Leukocyte Telomere Length on Total and Regional Brain Volumes in a Large Population-Based Cohort." *JAMA Neurology* 71 (10): 1247–54.

Kostas, Demetres, and Frank Rudzicz. 2020. "Thinker Invariance: Enabling Deep Neural Networks for BCI across More People." *Journal of Neural Engineering* 17 (5): 056008.

Larson-Prior, L. J., R. Oostenveld, S. Della Penna, G. Michalareas, F. Prior, A. Babajani-Feremi, J-M Schoffelen, et al. 2013. "Adding Dynamics to the Human Connectome Project with MEG." *NeuroImage* 80 (October): 190–201.

LeCun, Yann, Patrick Haffner, Léon Bottou, and Yoshua Bengio. 1999. "Object Recognition with Gradient-Based Learning." In *Shape, Contour and Grouping in Computer Vision*, edited by David A. Forsyth, Joseph L. Mundy, Vito di Gesú, and Roberto Cipolla, 319–

45. Berlin, Heidelberg: Springer Berlin Heidelberg.

Leeuwen, K. G. van, H. Sun, M. Tabaeizadeh, A. F. Struck, M. J. A. M. van Putten, and M. B. Westover. 2019. "Detecting Abnormal Electroencephalograms Using Deep Convolutional Networks." *Clinical Neurophysiology: Official Journal of the International Federation of Clinical Neurophysiology* 130 (1): 77–84.

Leonelli, S. 2016. "Data-Centric Biology." https://www.degruyter.com/document/doi/10.7208/9780226416502/html.

Liem, Franziskus, Gaël Varoquaux, Jana Kynast, Frauke Beyer, Shahrzad Kharabian Masouleh, Julia M. Huntenburg, Leonie Lampe, et al. 2017. "Predicting Brain-Age from Multimodal Imaging Data Captures Cognitive Impairment." *NeuroImage* 148 (March): 179–88.

Loeffler, Markus, Christoph Engel, Peter Ahnert, Dorothee Alfermann, Katrin Arelin, Ronny Baber, Frank Beutner, et al. 2015. "The LIFE-Adult-Study: Objectives and Design of a Population-Based Cohort Study with 10,000 Deeply Phenotyped Adults in Germany." *BMC Public Health* 15 (July): 691.

Lu, Chaochao, Bernhard Schölkopf, and José Miguel Hernández-Lobato. 2018. "Deconfounding Reinforcement Learning in Observational Settings." *arXiv [cs.LG]*. arXiv. http://arxiv.org/abs/1812.10576.

Mather, Karen Anne, Anthony Francis Jorm, Ruth Adeline Parslow, and Helen Christensen. 2011. "Is Telomere Length a Biomarker of Aging? A Review." *The Journals of Gerontology. Series A, Biological Sciences and Medical Sciences* 66 (2): 202–13.

McKinney, Wes, and Others. 2011. "Pandas: A Foundational Python Library for Data Analysis and Statistics." *Python for High Performance and Scientific Computing* 14 (9): 1–9.

Mehrabi, Ninareh, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. "A Survey on Bias and Fairness in Machine Learning." *ACM Comput. Surv.*, 115, 54 (6): 1–35.

Niso, Guiomar, Krzysztof J. Gorgolewski, Elizabeth Bock, Teon L. Brooks, Guillaume Flandin, Alexandre Gramfort, Richard N. Henson, et al. 2018. "MEG-BIDS, the Brain Imaging Data Structure Extended to Magnetoencephalography." *Scientific Data* 5 (June): 180110.

Niso, Guiomar, Christine Rogers, Jeremy T. Moreau, Li-Yuan Chen, Cecile Madjar, Samir Das, Elizabeth Bock, et al. 2016. "OMEGA: The Open MEG Archive." *NeuroImage* 124 (Pt B): 1182–87.

Obeid, Iyad, and Joseph Picone. 2016. "The Temple University Hospital EEG Data Corpus." *Frontiers in Neuroscience* 10 (May): 196.

O'Connor, David, Evelyn M. R. Lake, Dustin Scheinost, and R. Todd Constable. 2021. "Resample Aggregating Improves the Generalizability of Connectome Predictive Modeling." *NeuroImage* 236 (August): 118044.

Oostenveld, R., and P. Praamstra. 2001. "The Five Percent Electrode System for High-Resolution EEG and ERP Measurements." *Clinical Neurophysiology: Official Journal of the International Federation of Clinical Neurophysiology* 112 (4): 713–19.

Päivinen, Niina, Seppo Lammi, Asla Pitkänen, Jari Nissinen, Markku Penttonen, and Tapio Grönfors. 2005. "Epileptic Seizure Detection: A Nonlinear Viewpoint." *Computer Methods and Programs in Biomedicine* 79 (2): 151–59.

Paszke, Adam, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, et al. 2019. "PyTorch: An Imperative Style, High-Performance Deep Learning Library." In *Advances in Neural Information Processing Systems*, edited by H. Wallach, H. Larochelle, A. Beygelzimer, F. d\textquotesingle Alché-Buc, E. Fox, and R. Garnett. Vol. 32. Curran Associates, Inc. https://proceedings.neurips.cc/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf.

Pedersen, Thomas Lin. 2019. "Patchwork: The Composer of Plots." *R Package Version* 1 (0).

Pedregosa, F., G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, et al.

2011. "Scikit-Learn: Machine Learning in Python." *Journal of Machine Learning Research: JMLR* 12: 2825–30.

Pernet, Cyril R., Stefan Appelhoff, Krzysztof J. Gorgolewski, Guillaume Flandin, Christophe Phillips, Arnaud Delorme, and Robert Oostenveld. 2019. "EEG-BIDS, an Extension to the Brain Imaging Data Structure for Electroencephalography." *Scientific Data* 6 (1): 103.

Perslev, Mathias, Sune Darkner, Lykke Kempfner, Miki Nikolic, Poul Jørgen Jennum, and Christian Igel. 2021. "U-Sleep: Resilient High-Frequency Sleep Staging." *NPJ Digital Medicine* 4 (1): 72.

Poldrack, Russell A., Grace Huckins, and Gael Varoquaux. 2020. "Establishment of Best Practices for Evidence for Prediction: A Review." *JAMA Psychiatry (Chicago, Ill.)* 77 (5): 534–40.

Raffel, Joel, James Cole, Chris Record, Sujata Sridharan, David Sharp, and Richard Nicholas. 2017. "Brain Age: A Novel Approach to Quantify the Impact of Multiple Sclerosis on the Brain (P1.371)." *Neurology* 88 (16 Supplement). https://n.neurology.org/content/88/16_Supplement/P1.371.short.

Richman, J. S., and J. R. Moorman. 2000. "Physiological Time-Series Analysis Using Approximate Entropy and Sample Entropy." *American Journal of Physiology. Heart and Circulatory Physiology* 278 (6): H2039–49.

Roberts, S. J., W. Penny, and I. Rezek. 1999. "Temporal and Spatial Complexity Measures for Electroencephalogram Based Brain-Computer Interfacing." *Medical & Biological Engineering & Computing* 37 (1): 93–98.

Rodrigues, Pedro Luiz Coelho, Christian Jutten, and Marco Congedo. 2019. "Riemannian Procrustes Analysis: Transfer Learning for Brain-Computer Interfaces." *IEEE Transactions on Bio-Medical Engineering* 66 (8): 2390–2401.

Roy, Yannick, Hubert Banville, Isabela Albuquerque, Alexandre Gramfort, Tiago H. Falk, and Jocelyn Faubert. 2019. "Deep Learning-Based Electroencephalography Analysis: A Systematic Review." *Journal of Neural Engineering* 16 (5): 051001.

Rudin, Cynthia. 2019. "Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and Use Interpretable Models Instead." *Nature Machine Intelligence* 1 (5): 206–15.

Sabbagh, D., P. Ablin, G. Varoquaux, and A. Gramfort. 2019. "Manifold-Regression to Predict from MEG/EEG Brain Signals without Source Modeling." *arXiv Preprint arXiv*. https://arxiv.org/abs/1906.02687.

Sabbagh, David, Pierre Ablin, Gaël Varoquaux, Alexandre Gramfort, and Denis A. Engemann. 2020. "Predictive Regression Modeling with MEG/EEG: From Source Power to Signals and Cognitive States." *NeuroImage* 222 (November): 116893.

Scahill, Rachael I., Chris Frost, Rhian Jenkins, Jennifer L. Whitwell, Martin N. Rossor, and Nick C. Fox. 2003. "A Longitudinal Study of Brain Volume Changes in Normal Aging Using Serial Registered Magnetic Resonance Imaging." *Archives of Neurology* 60 (7): 989–94.

Schiratti, J-B, Jean-Eudes Le Douget, Michel Le Van Quyen, Slim Essid, and Alexandre Gramfort. 2018. "An Ensemble Learning Approach to Detect Epileptic Seizures from Long Intracranial EEG Recordings." In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 856–60.

Schiratti, J-B, Jean-Eudes Le Douget, Michel Le Van Quyen, Slim Essid, and Alexandre Gramfort. 2018. "An Ensemble Learning Approach to Detect Epileptic Seizures from Long Intracranial EEG Recordings." *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. https://doi.org/10.1109/icassp.2018.8461489.

Schirrmeister, Robin Tibor, Jost Tobias Springenberg, Lukas Dominique Josef Fiederer, Martin Glasstetter, Katharina Eggensperger, Michael Tangermann, Frank Hutter, Wolfram Burgard, and Tonio Ball. 2017. "Deep Learning with Convolutional Neural Networks for EEG Decoding and Visualization." *Human Brain Mapping* 38 (11): 5391–5420.

Schulz, Marc-Andre, B. T. Thomas Yeo, Joshua T. Vogelstein, Janaina Mourao-Miranada,

Jakob N. Kather, Konrad Kording, Blake Richards, and Danilo Bzdok. 2020. "Different Scaling of Linear Models and Deep Learning in UKBiobank Brain Images versus Machine-Learning Datasets." *Nature Communications*. https://doi.org/10.1038/s41467-020-18037-z.

Schumacher, Julia, Nicola J. Ray, Calum A. Hamilton, Paul C. Donaghy, Michael Firbank, Gemma Roberts, Louise Allan, et al. 2021. "Cholinergic White Matter Pathways in Dementia with Lewy Bodies and Alzheimer's Disease." *Brain: A Journal of Neurology*, October. https://doi.org/10.1093/brain/awab372.

Shafto, Meredith A., Lorraine K. Tyler, Marie Dixon, Jason R. Taylor, James B. Rowe, Rhodri Cusack, Andrew J. Calder, et al. 2014. "The Cambridge Centre for Ageing and Neuroscience (Cam-CAN) Study Protocol: A Cross-Sectional, Lifespan, Multidisciplinary Examination of Healthy Cognitive Ageing." *BMC Neurology* 14 (October): 204.

Shekh Ibrahim, Sharifah Aida, Nurfaten Hamzah, Athirah Raihanah Abdul Wahab, Jafri Malin Abdullah, Nurul Hashimah Ahamed Hassain Malim, Putra Sumari, Zamzuri Idris, et al. 2020. "Big Brain Data Initiative in Universiti Sains Malaysia: Challenges in Brain Mapping for Malaysia." *The Malaysian Journal of Medical Sciences: MJMS* 27 (4): 1–8.

Sitt, Jacobo Diego, Jean Remi King, Imen El Karoui, Benjamin Rohaut, Frederic Faugeras, Alexandre Gramfort, Laurent Cohen, Mariano Sigman, Stanislas Dehaene, and Lionel Naccache. 2014. "Large Scale Screening of Neural Signatures of Consciousness in Patients in a Vegetative or Minimally Conscious State." *Brain: A Journal of Neurology* 137 (8): 2258–70.

Smith, Stephen M., Thomas E. Nichols, Diego Vidaurre, Anderson M. Winkler, Timothy E. J. Behrens, Matthew F. Glasser, Kamil Ugurbil, Deanna M. Barch, David C. Van Essen, and Karla L. Miller. 2015. "A Positive-Negative Mode of Population Covariation Links Brain Connectivity, Demographics and Behavior." *Nature Neuroscience* 18 (11): 1565–67.

Spiegelhalter, David. 2016. "How Old Are You, Really? Communicating Chronic Risk through 'effective Age' of Your Body and Organs." *BMC Medical Informatics and Decision Making* 16 (1). https://doi.org/10.1186/s12911-016-0342-z.

Stokes, Mark G., Michael J. Wolff, and Eelke Spaak. 2015. "Decoding Rich Spatial Information with High Temporal Resolution." *Trends in Cognitive Sciences* 19 (11): 636–38.

Sun, Haoqi, Luis Paixao, Jefferson T. Oliva, Balaji Goparaju, Diego Z. Carvalho, Kicky G. van Leeuwen, Oluwaseun Akeju, et al. 2019. "Brain Age from the Electroencephalogram of Sleep." *Neurobiology of Aging* 74 (February): 112–20.

Taulu, Samu, Juha Simola, and Matti Kajola. 2005. "Applications of the Signal Space Separation Method." *Signal Processing, IEEE Transactions on* 53 (9): 3359–72.

Taylor, Jason R., Nitin Williams, Rhodri Cusack, Tibor Auer, Meredith A. Shafto, Marie Dixon, Lorraine K. Tyler, Richard N. Henson, and Others. 2017. "The Cambridge Centre for Ageing and Neuroscience (Cam-CAN) Data Repository: Structural and Functional MRI, MEG, and Cognitive Data from a Cross-Sectional Adult Lifespan Sample." *NeuroImage* 144: 262–69.

Teixeira, C. A., B. Direito, H. Feldwisch-Drentrup, M. Valderrama, R. P. Costa, C. Alvarado-Rojas, S. Nikolopoulos, et al. 2011. "EPILAB: A Software Package for Studies on the Prediction of Epileptic Seizures." *Journal of Neuroscience Methods* 200 (2): 257–71.

Tibor Schirrmeister, Robin, Lukas Gemein, Katharina Eggensperger, Frank Hutter, and Tonio Ball. 2017. "Deep Learning with Convolutional Neural Networks for Decoding and Visualization of EEG Pathology." *arXiv E-Prints*, August, arXiv:1708.08012.

Tietz, Marian, T. J. Fan, D. Nouri, and Others. 2017. "Skorch: A Scikit-Learn Compatible Neural Network Library That Wraps PyTorch." July.

Valdes-Sosa, Pedro A., Lidice Galan-Garcia, Jorge Bosch-Bayard, Maria L. Bringas-Vega, Eduardo Aubert-Vazquez, Iris Rodriguez-Gil, Samir Das, et al. 2021. "The Cuban Human Brain Mapping Project, a Young and Middle Age Population-Based EEG, MRI, and Cognition Dataset." *Scientific Data* 8 (1): 45.

Van Essen, David C., Stephen M. Smith, Deanna M. Barch, Timothy E. J. Behrens, Essa

Yacoub, Kamil Ugurbil, and WU-Minn HCP Consortium. 2013. "The WU-Minn Human Connectome Project: An Overview." *NeuroImage* 80 (October): 62–79.

Varoquaux, Gaël, Pradeep Reddy Raamana, Denis A. Engemann, Andrés Hoyos-Idrobo, Yannick Schwartz, and Bertrand Thirion. 2017. "Assessing and Tuning Brain Decoders: Cross-Validation, Caveats, and Guidelines." *NeuroImage* 145 (Pt B): 166–79.

Virtanen, Pauli, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, et al. 2020. "SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python." *Nature Methods* 17 (3): 261–72.

Völker, Martin, Robin T. Schirrmeister, Lukas D. J. Fiederer, Wolfram Burgard, and Tonio Ball. 2018. "Deep Transfer Learning for Error Decoding from Non-Invasive EEG." In *2018 6th International Conference on Brain-Computer Interface (BCI)*, 1–6.

Walhovd, K. B., A. M. Fjell, J. Brewer, L. K. McEvoy, C. Fennema-Notestine, D. J. Hagler Jr, R. G. Jennings, D. Karow, A. M. Dale, and Alzheimer's Disease Neuroimaging Initiative. 2010. "Combining MR Imaging, Positron-Emission Tomography, and CSF Biomarkers in the Diagnosis and Prognosis of Alzheimer Disease." *AJNR. American Journal of Neuroradiology* 31 (2): 347–54.

Wickham, Hadley. 2011. "Ggplot2." *Wiley Interdisciplinary Reviews. Computational Statistics* 3 (2): 180–85.

Wong, Chi Wah, Pamela N. DeYoung, and Thomas T. Liu. 2016. "Differences in the Resting-State fMRI Global Signal Amplitude between the Eyes Open and Eyes Closed States Are Related to Changes in EEG Vigilance." *NeuroImage* 124 (Pt A): 24–31.

Wrigglesworth, Jo, Nurathifah Yaacob, Phillip Ward, Robyn L. Woods, John McNeil, Elsdon Storey, Gary Egan, et al. 2021. "Brain-Predicted Age Difference Is Associated with Cognitive Processing in Later-Life." *Neurobiology of Aging*, October. https://doi.org/10.1016/j.neurobiolaging.2021.10.007.

Xifra-Porxas, Alba, Arna Ghosh, Georgios D. Mitsis, and Marie-Hélène Boudrias. 2021. "Estimating Brain Age from Structural MRI and MEG Data: Insights from Dimensionality Reduction Techniques." *NeuroImage* 231 (117822): 117822.

Yamins, Daniel L. K., and James J. DiCarlo. 2016. "Using Goal-Driven Deep Learning Models to Understand Sensory Cortex." *Nature Neuroscience* 19 (3): 356–65.

Ye, Elissa, Haoqi Sun, Michael J. Leone, Luis Paixao, Robert J. Thomas, Alice D. Lam, and M. Brandon Westover. 2020. "Association of Sleep Electroencephalography-Based Brain Age Index With Dementia." *JAMA Network Open* 3 (9): e2017357.

Yger, Florian, Maxime Berar, and Fabien Lotte. 2017. "Riemannian Approaches in Brain-Computer Interfaces: A Review." *IEEE Transactions on Neural Systems and Rehabilitation Engineering: A Publication of the IEEE Engineering in Medicine and Biology Society* 25 (10): 1753–62.

Zimmer, Lucas, Marius Lindauer, and Frank Hutter. 2021. "Auto-Pytorch: Multi-Fidelity MetaLearning for Efficient and Robust AutoDL." *IEEE Transactions on Pattern Analysis and Machine Intelligence* 43 (9): 3079–90.