1  **Title**

2  Brain age as an estimator of neurodevelopmental outcome: A deep learning approach for neonatal

3  cot-side monitoring

4

5

6  **Author names**

7  Amir Ansari[1a], Kirubin Pillay[2a], Luke Baxter[2], Emad Arasteh [1,3], Anneleen Dereymaeker[4], Gabriela

8  Schmidt Mellado[2], Katrien Jansen[4,5], Gunnar Naulaers[4], Aomesh Bhatt[2], Sabine Van Huffel[1], Caroline

9  Hartley[2], Maarten De Vos[1,5], Rebeccah Slater[2]*

10

11  [a]These authors contributed equally

12

13

14  **Author affiliations**

15  [1] Department of Electrical Engineering (ESAT), STADIUS Center for Dynamical Systems, Signal

16  Processing and Data Analytics, KU Leuven, Leuven, Belgium

17  [2] Department of Paediatrics, University of Oxford, Oxford, UK

18  [3] Department of Neonatology, Wilhelmina Children's Hospital, University Medical Center Utrecht,

19  Utrecht, the Netherlands

20  [4] Department of Development and Regeneration, University Hospitals Leuven, Neonatal Intensive

21  Care Unit, KU Leuven, Leuven, Belgium

22  [5] Department of Development and Regeneration, University Hospitals Leuven, Child Neurology, KU

23  Leuven, Leuven, Belgium

24

25

26  **Contact information**

27  * Corresponding author

28  Paediatric Neuroimaging Group, Department of Paediatrics, Level 2 Children's Hospital, John

29  Radcliffe Hospital, University of Oxford, Oxford, United Kingdom

30  rebeccah.slater@paediatrics.ox.ac.uk

31

32

**Highlights**

- Preterm stress exposure leads to long-term neurodevelopmental deficits
- Deficits are quantifiable using EEG-based brain age prediction errors
- Our deep-learning solution for brain age prediction outperforms previous approaches
- Predictions are achieved with only 20 mins EEG and a single bipolar channel
- Prediction errors correlate with long-term Bayley scale neurodevelopmental outcomes

**Abstract**

The preterm neonate can experience stressors that affect the rate of brain maturation and lead to long-term neurodevelopmental deficits. However, some neonates who are born early follow normal developmental trajectories. Extraction of data from electroencephalography (EEG) signals can be used to calculate the neonate's brain age which can be compared to their true age. Discrepancies between true age and brain age (the brain age delta) can then be used to quantify maturational deviation, which has been shown to correlate with long-term abnormal neurodevelopmental outcomes. Nevertheless, current brain age models that are based on traditional analytical techniques are less suited to clinical cot-side monitoring due to their dependency on long-duration EEG recordings, the need to record activity across multiple EEG channels, and the manual calculation of predefined EEG features which is time-consuming and may not fully capture the wealth of information in the EEG signal. In this study, we propose an alternative deep-learning approach to determine brain age, which operates directly on the EEG, using a Convolutional Neural Network (CNN) block based on the Inception architecture (called Sinc). Using this deep-learning approach on a dataset of preterm infants with normal neurodevelopmental outcomes (where we assume brain age = postmenstrual age), we can calculate infant brain age with a Mean Absolute Error (MAE) of 0.78 weeks (equivalent to a brain age estimation error for the infant within +/- 5.5 days of their true age). Importantly, this level of accuracy can be achieved by recording only 20 minutes of EEG activity from a single channel. This compares favourably to the degree of accuracy that can be achieved using traditional methods that require long duration recordings (typically >2 hours of EEG activity) recorded from a higher density 8-electrode montage (MAE = 0.73 weeks). Importantly, the deep learning model's brain age deltas also distinguish between neonates with normal and severely abnormal outcomes (Normal MAE = 0.71 weeks, severely abnormal MAE = 1.27 weeks, p=0.02, one-way ANOVA), making it highly suited for potential clinical applications. Lastly, in an independent dataset collected at an independent site, we demonstrate the model's generalisability in age prediction, as accurate age predictions were also observed (MAE of 0.97 weeks).

**Keywords**

Preterm, Electroencephalography, Machine Learning, Artificial Intelligence, Convolutional Neural Network, Bayley Scales

## 1. Introduction

The newborn infant's brain is undergoing rapid developmental change, influenced by both genetic and environmental factors (Colonnese et al., 2010; Milh et al., 2007; Wess et al., 2017). Relative to their term-born counterparts, infants born prematurely are at increased risk of poorer long-term neurodevelopmental outcomes (Blencowe et al., 2013; Wallois et al., 2020). This risk of impairment increases with the degree of prematurity at birth and the presence of gross morphological lesions, but can also be brought about by subtler environmental stressors (Scher, 2008), excessive exposure to painful stimuli (Grunau, 2013; Moultrie et al., 2017), and pharmacological interventions (Duerden et al., 2016; Malk et al., 2014).

The early identification of abnormal neurodevelopment is essential to identify infants at greatest risk who might benefit most from developmental care interventions (Burke, 2018). To date, neurological assessment of the newborn has remained predominantly subjective (Dempsey et al., 2018). For example, trained neonatologists and clinical neurophysiologists visually inspect infant's brain activity using electroencephalography (EEG) to determine if brain function is developmentally age-appropriate or dysmature (Scher, 1997), based on developmentally changing EEG features characteristic of maturational status (André et al., 2010). While these trained individuals can estimate age with an error of two weeks for preterm babies and one week for term babies, these estimates can be highly variable across reviewers (Stevenson et al., 2020b). Subjectivity, inter-rater variability, and requirement of specialist EEG interpretation are central issues that severely limit the reliability and generalisability of many current neurological assessment methods. There is an urgent need for objective and automated neuromonitoring that can be used cot-side to identify infants at increased risk of abnormal neurodevelopmental outcomes.

To this end, a variety of metrics have been developed to capture key maturational characteristics from the preterm EEG (De Wel et al., 2017; Dereymaeker et al., 2016; Lavanga et al., 2017; Pillay et al., 2018; Tolonen et al., 2007), and these measures have been combined using machine learning algorithms to successfully predict infants' brain age (O'Toole et al., 2016; Stevenson et al., 2017). An infant's brain age is their predicted age from a model that has been trained using brain-based features (structural or functional) as predictors and true age as the response. In adults, the difference between the brain age and the true age, termed the brain age delta, has been demonstrated to be more than random noise prediction error, but in fact is of biological and clinical value (Smith et al., 2019; Vidal-Pineiro et al., 2021).

In infants, analogous findings have been observed. Recently, we trained a Random Forest (RF) regression model using a data-driven approach that combined 226 EEG features and demonstrated a significant correlation between the infants' brain age delta and the severity of their abnormal neurodevelopmental outcome, where the neurodevelopmental outcomes were assessed behaviourally using the Bayley Scales of Infant Development (BSID-II) at a 9-month follow-up test occasion (Pillay et al., 2020). Additionally, an independent research group showed a similar correlation when training a multivariate regression model for brain age estimation (Stevenson et al., 2020a). These studies established the proof-of-concept in infant populations that the inter-individual variability in automatically and objectively generated brain age deltas could be used to

116 risk-stratify infants in the first few weeks of postnatal life according to neurodevelopmental
117 outcomes.

118

119 However, a major limitation to these studies is their lack of clinical utility. A large number of features
120 are required to summarize the EEG data, which are computationally time-consuming to calculate.
121 These approaches rely on pre-staging the EEG recording into sleep states (i.e. sleep-staging) or burst
122 periods which require additional algorithms (Dereymaeker et al., 2017b; Palmu et al., 2010).
123 Furthermore, multiple EEG channels are required as well as at least 1 hour EEG recording duration.
124 These data-heavy requirements severely limit the ease with which these methods can be
125 incorporated into the busy clinical environment.

126

127 Here, we directly address these barriers to clinical utility by adopting a deep learning approach.
128 Deep learning has demonstrated superior performance over traditional machine learning methods,
129 has excellent performance on a reduced number of EEG channels, and tends to perform predictions
130 faster once trained (Ansari et al., 2018). Furthermore, deep learning models are gaining popularity
131 in preterm EEG analysis for classifying seizures (Ansari et al., 2019; O'Shea et al., 2021) and for
132 automated sleep-staging (Ansari et al., 2020). Together, these observations suggest deep learning
133 could offer a promising approach for cot-side monitoring and assessment of neurological function.

134

135 In the current study, we implement a novel Convolutional Neural Network (CNN)-based
136 architecture, inspired by Google's Inception model (and its variants), to generate infant brain age
137 predictions using dramatically reduced EEG data requirements compared to previous proof-of-
138 concept studies. We use our established RF model as a "gold standard" benchmark of performance,
139 a model which requires eight EEG channels, at least 1 hour EEG recording duration, and EEG data
140 sleep-staging. We train the RF and deep learning models on a training dataset, and subsequently
141 test the models' performance on two independent datasets, demonstrating robust external
142 validation. Using our deep learning approach, we achieve performance comparable to our RF model
143 benchmark, while requiring only a single EEG channel (1-channel bipolar montage), 20 mins EEG
144 recording duration, and no EEG data sleep-staging. Our deep learning model is able to accurately
145 predict infant age within the first few weeks of postnatal life, and generates brain age deltas with
146 magnitudes that significantly differ between infants with normal and severely abnormal
147 neurodevelopmental outcomes assessed using BSID-II at 9-month follow-up. This study thus
148 demonstrates potential clinical utility for an objective and automated deep learning-based
149 approach to cot-side assessment of infants' neurological function and neurodevelopmental
150 outcomes.

151

152

153 **2. Methods**
154 **2.1. Participants**
155 *2.1.1. Study Design*
156 Data were collected in three independent cohorts. The first cohort, referred to as dataset $\mathcal{D}1$, was
157 used to train the models and compare the relative performances among models e.g. models with
158 different architectures, different channel montages, and different recording durations. The second

159  cohort, referred to as dataset $\mathcal{D}2$, was used to independently test the trained RF and deep learning
160  models in their brain age prediction performances, and to assess the association between brain age
161  deltas and 9-month BSID-II follow-up outcomes. The third cohort, referred to as dataset $\mathcal{D}3$, was
162  used to further test the generalisability of the deep learning model to predict brain age in this
163  dataset collected at an independent site by an independent research team.

165  *2.1.2. Recruitment*
166  EEG data for datasets $\mathcal{D}1$ and $\mathcal{D}2$ were recorded from the Neonatal Intensive Care Unit (NICU) at
167  UZ Leuven Hospitals, Leuven, Belgium. Infants were recruited and data recorded with informed
168  consent from the parents and in accordance with the guidelines approved by the ethics committee
169  of the University Hospitals, Leuven. All infants had a gestational age (GA) at birth less than 32 weeks,
170  and between two and four recordings were obtained during their stay in the NICU.

172  Infants in dataset $\mathcal{D}3$ were selected from a database of previously recorded data collected at the
173  Newborn Care Unit and Maternity wards of the John Radcliffe Hospital (Oxford University Hospitals
174  NHS Foundation Trust, Oxford, United Kingdom). Ethical approval was obtained from the UK
175  National Research Ethics Service (reference: 12/SC/0447) and parental written informed consent
176  was obtained before each participant was studied.

178  All participant recruitment was conducted in accordance with the standards set by the Declaration
179  of Helsinki and Good Clinical Practice guidelines.

181  *2.1.3. Datasets*
182  Datasets $\mathcal{D}1$ and $\mathcal{D}2$ were collected as previously described (Pillay et al., 2020). Dataset $\mathcal{D}1$ consists
183  of n=40 infants (111 recordings) with postmenstrual age range (PMA) at time of recording of 27.3–
184  43.1 weeks, with mean recording duration of 8h 07m (standard deviation: 5h 55m) and mean
185  number of recordings per infant of 2.8 (standard deviation: 1.6). All infants in dataset $\mathcal{D}1$ were
186  selected for normal neurodevelopmental outcome at 24-months follow-up age based on
187  behavioural assessment using BSID-II.

189  Dataset $\mathcal{D}2$ consists of n=43 infants (142 recordings). One infant with a single recording was
190  excluded as our objective with this dataset was to assess longitudinal multi-recording trajectories.
191  The analysed dataset $\mathcal{D}2$ thus consists of n=42 infants (141 recordings) with a PMA at recording
192  range of 27.3–42.0 weeks, mean recording duration of 7h 05m (standard deviation: 5h 43m), and
193  mean number of recordings per infants of 3.3 (standard deviation: 1.4). Unlike dataset $\mathcal{D}1$, dataset
194  $\mathcal{D}2$ includes infants with a range of both normal and abnormal outcomes, grouped by BSID-II scores
195  at 9-month follow-up (Pillay et al., 2020). N=22 infants (71 recordings) had normal outcome i.e. no
196  neurodevelopmental impairment (NDI); n=10 infants (36 recordings) had mild abnormal outcome
197  (mild NDI); and n=10 infants (34 recordings) had moderate-to-severe abnormal outcome (mild-to-
198  severe NDI) or died (Pascal et al., 2020).

200  Dataset $\mathcal{D}3$ consists of n=73 infants, each recorded on a single test occasion (thus 73 recordings).
201  Infants were included in this dataset for the current study if they had at least 20 minutes of EEG

202  data recorded and if the EEG was assessed as normal for age by a trained clinical neurophysiologist
203  (author GSM). The infants had a median PMA at recording of 35.3 weeks (interquartile range: 33.3
204  – 36.9, range: 28.0 – 42.6) and postnatal age of 14 days (interquartile range: 5 – 41, range: 0 – 112).
205  The mean recording duration was 50 minutes (standard deviation: 18 minutes).
206
207  **2.2. EEG data**
208  *2.2.1. Setup*
209  For dataset $\mathcal{D}1$ and $\mathcal{D}2$, data were recorded using a sampling frequency of 250 Hz using Brain RT
210  OSG Equipment (Mechelen, Belgium). In a few cases, the EEG was sampled at 256 Hz due to some
211  setup variations on the Brain RT device used. All recordings were performed with nine electrodes in
212  a referential montage: Fp1, Fp2, C3, C4, T3, T4, O1, O2, and Cz reference (Figure 1).
213
214  For dataset $\mathcal{D}3$, EEG recordings were acquired from DC to 800 Hz using a SynAmps RT 64-channel
215  headbox and amplifiers (Compumedics Neuroscan). Activity was recorded using CURRY scan7
216  neuroimaging suite (Compumedics Neuroscan), with a sampling rate of 2000 Hz. Between 8 and 25
217  electrodes were used for recording, positioned according to the modified international 10-20
218  system, including C3 and C4 (those used in the analysis here), with reference at Fz and ground at
219  Fpz. The scalp was cleaned with preparation gel (Nuprep gel, D.O. Weaver and Co.) and disposable
220  Ag/AgCl cup electrodes (Ambu Neuroline) were placed with conductive paste (Elefix EEG paste,
221  Nihon Kohden).
222
223  *2.2.2. Preprocessing*
224  For the deep learning approaches in datasets $\mathcal{D}1$ and $\mathcal{D}2$, each recording was downsampled to 64
225  Hz to reduce the number of parameters required to train the model. The downsampling routine
226  included pre-filtering to prevent aliasing using a low-pass filter with cut-off frequency 32 Hz.
227  Filtering and downsampling was performed using the scipy.signal.resample_poly function.
228  Recordings were then split into 30-second segments and the amplitudes standardized such that the
229  mean and standard deviation of the amplitudes were zero and one, respectively. The mean and
230  standard deviation were obtained by standardizing the data (across all channels) in the training set,
231  with these values carried forward to standardize the test sets (see below). Finally, any segments
232  where the absolute differences (compared to the mean) at any point exceeded 600 μV were
233  rejected as artefact.
234
235  For the RF approach in datasets $\mathcal{D}1$ and $\mathcal{D}2$, which relied on an explicit pre-calculation of many
236  established features, the pre-processing approach (resampling and standardization) was different
237  and specific to each calculated feature, as described previously (Pillay et al., 2018).
238
239  For dataset $\mathcal{D}3$, pre-processing was matched to the $\mathcal{D}1$ and $\mathcal{D}2$ deep learning approach. We applied
240  a low-pass 32 Hz anti-aliasing filter followed by downsampling to 64 Hz. For standardization of
241  dataset $\mathcal{D}3$, the mean and standard deviation of $\mathcal{D}1$ were used.
242
243  **2.3. Brain age prediction architectures**
244  *2.3.1. Sinc architecture*

6

245 Figure 2a shows the block-diagram of the proposed deep neural network for brain age prediction.
246 As input, the network processes a 30 s multi-channel EEG segment. Each input segment has
247 dimension $C$ x 1920 where $C$ is the number of EEG channels and 1920 is the total number of
248 timepoints in the 30 s segment (30 s duration x 64 Hz sampling frequency). Each segment has a
249 single output label that is a continuous PMA value.
250
251 The model includes a series of convolutional layers with exponential linear unit (ELU) activations,
252 maximum and average pooling layers to downsample the data, normalization layers for faster
253 training convergence, and a dense layer with linear activation to perform the final regression and
254 produce a brain age estimate. As each convolutional layer is designed to extract specific
255 characteristics from the EEG, these are analogous to a (trainable, data-driven) feature extraction
256 layer. More generally, the proposed architecture can be grouped into a more traditional, sequential
257 CNN block that can be described as an initial feature extraction stage, followed by the two
258 successive Sinc (i.e. Shared Inception) blocks that form a second feature extraction stage.
259
260 We previously introduced Sinc as a powerful CNN-based block for extracting multi-scale temporal
261 information from infant EEG, namely sleep state classification (Ansari et al., 2021). Sinc is an
262 extension of Google's Inception block, where the original independent and parallel convolutional
263 branches are now boosted via parameter sharing. As shown in Figure 2b, the output from each
264 preceding branch is additionally fed into the subsequent one, with the overall output of Sinc
265 comprising the concatenation of all multi-scale convolutions in the block (see also Figure 2biv). This
266 increases the number of temporal scales achievable (by allowing a wider range of receptive fields),
267 when compared to an Inception layer, while avoiding the need to scale up the number of trainable
268 parameters as a result. Only two hyper-parameters are required for a Sinc block: $M$ (the number of
269 convolutional branches), and $N$ (the number of convolutional filters used in each branch). When
270 using a single-channel EEG segment as input ($C$ = 1), the total number of trainable parameters in
271 the complete model is 620K.
272
273 *2.3.2. Alternative deep learning architectures*
274 Four different deep learning architectures were also considered, based on recent key developments
275 in the CNN domain (Figure 2b), and the Sinc model was compared against these architectures: No
276 FEII (following the same design as the Sinc architecture but without the entire Feature Extraction II
277 portion), CNN (replaces the Sinc blocks with the same convolutional neural network layer used
278 elsewhere in the model), Residual (similar to CNN but including the additional residual shortcut),
279 and Inc (replacing the Sinc blocks with traditional Inception blocks). These architectures are
280 described in Supplementary Information S.1.
281
282 *2.6.2. Random Forest (RF) architecture*
283 In addition to comparing the Sinc architecture to alternative deep neural network architectures, it
284 was also compared against our established and previously published RF approach (Pillay et al.,
285 2020). The RF model was developed using a large set of pre-calculated features, derived after the
286 EEG is classified into different sleep stages using an additional, unsupervised algorithm known as
287 Cluster-based Adaptive Sleep Staging (CLASS) that we have also previously developed (Dereymaeker

7

288  et al., 2017b; Pillay et al., 2020). The pre-calculated features were derived from an EEG literature
289  review covering the amplitude domains, Fourier transforms, Wavelet transforms, Empirical Mode
290  Decompositions (EMDs) and other complexity measures (such as entropy and fractal analysis).
291  These features were calculated across all channels and a final median taken across channels as input
292  into the RF model. The RF is an ensemble method that uses a large set (or 'forest') of trained decision
293  trees to provide an averaged final prediction. Each tree is trained on a bootstrapped sample of the
294  dataset and a random selection of the features which is shown to provide better prediction accuracy
295  than from an individual tree and can also provide an implicit measure of the important features
296  used in the algorithm. The RF model uses 1500 trees and utilizes all features for each tree split. All
297  steps and hyperparameter choices used here are the same as our previous published RF approach
298  (Pillay et al., 2020).

299

300  **2.4. Model training and relative performance assessments using dataset $\mathcal{D}1$**

301  *2.4.1. Splitting dataset $\mathcal{D}1$ into training set and test set*

302  Dataset $\mathcal{D}1$ was used to train and test all models. By using a cohort of only normal outcome data, it
303  is assumed that predicted brain age equates to true PMA. This allows training of a normative model
304  to predict the PMA, and therefore brain age, for a normally developing baby (Pillay et al., 2020;
305  Stevenson et al., 2020a, 2017). Further data can then be assessed against this trained model to
306  identify deviations. Dataset $\mathcal{D}1$ was divided by recording into age-stratified training and test sets of
307  size 50 and 47 recordings, respectively (Supplementary Information S.2.; Supplementary Figure 1).

308

309  *2.4.2. Model training in dataset $\mathcal{D}1$*

310  For the training of both deep learning models and the RF model, the mean squared error (MSE) loss
311  was used. For the RF model, the model was re-trained using the original training procedure as
312  previously outlined (Pillay et al., 2020). For the deep learning models, model training included early
313  stopping, Gaussian noise addition, recording segmentation into 30 s segments, and ensemble
314  learning. These four components are described in detail in Supplementary Information S.3, with
315  early stopping, Gaussian noise addition, and ensemble learning included to increase robustness of
316  the model.

317

318  *2.4.3. Model testing in dataset $\mathcal{D}1$*

319  *2.4.3.1. Assessing model performance*

320  The ultimate goal of each brain age prediction model is to generate a single brain age prediction
321  estimate per EEG recording. For the deep learning models, each deep learning model generates ten
322  brain age prediction estimates per 30 s segment of an EEG recording (as a 10-learner ensemble
323  method was used, see Supplementary Information S.3.). During testing, all contiguous 30 s segments
324  across each recording are used with the number of 30 s segments therefore dependent on the
325  overall EEG recording duration. To aggregate a deep learning model's predictions to a single value
326  per recording, the median across the ten ensemble predictions per 30 s segment is determined,
327  then a further median across all 30 s segments in the recording is taken resulting in the final
328  prediction estimate.  This is different to the RF model strategy, where a single brain age prediction
329  estimate is generated per recording by manually calculating features in the 30 s segments and taking
330  the medians across all segments before brain age prediction is performed.  Across all recordings in

331    the test set in $\mathcal{D}1$, there were a total of 30K segments used. For both the deep learning models and

332    RF model, the final prediction estimate for a recording is used to generate the prediction error (or

333    absolute prediction error) for that recording.

334

335    *2.4.3.2. Reducing EEG channel requirements*

336    The established RF model uses eight channels in a referential montage (Figure 1) to predict infant

337    brain age. The performance of the RF model, the Sinc model, and the other deep learning models

338    were assessed and compared using this initial setup. Subsequently, the deep learning models were

339    re-trained and performances compared as the number of EEG channels were systematically

340    reduced: 4-channel referential (C3, C4, T3 and T4), 2-channel referential (C3 and C4), and finally a

341    1-channel bipolar (C3-C4) montage (Figure 1). Channels were selected to ensure good symmetry

342    across the midline of the scalp and ample coverage. The 1-channel bipolar montage was selected

343    for its similarity to setups used in clinical amplitude-integrated EEG (aEEG) monitors. EEG pre-

344    processing was independently repeated each time, with the amplitude standardisation step

345    recalculated on the reduced channel configurations. After (re)training using the training set, each

346    model generated a brain age prediction per recording in the test set. This set of predictions was

347    used to generate a set of absolute errors per model. Using one-sample paired t-tests ($p < 0.05$

348    significance level), we assessed the model performances by comparing t-statistic magnitudes and

349    tested for statistically significant differences between the Sinc model's mean absolute error and the

350    mean absolute error of each of the alternative models (RF model and other deep learning models).

351

352    It is worth noting that the 1-channel bipolar montage used for our analyses was achieved by ignoring

353    the additional channels unnecessary for this montage. This approach is distinct to a true clinical

354    scenario when only a 1-channel bipolar montage would be used during recording. Our assumption,

355    which we believe to be reasonable, is that both approaches to 1-channel bipolar montage setup are

356    closely matched for this specific use case. However, this assumption should be tested in future

357    external validations of the deep learning model using clinical grade bipolar montage data.

358

359    *2.4.3.3. Reducing EEG recording duration requirements*

360    Having demonstrated the high performance of the deep learning Sinc model using the full-length

361    EEG recording duration with only a 1-channel bipolar setup (section 2.4.3.2.), we next assessed the

362    Sinc model performance using the 1-channel bipolar setup (Figure 1) as the EEG recording duration

363    was systematically varied. We examined a range of recording lengths from 0.5–120 min and

364    compared Sinc model performance on these reduced recording durations relative to the Sinc model

365    performance with the full-length EEG recording duration to identify an appropriate reduced

366    recording duration. To get a reduced recording from a single full recording, we randomly sampled

367    each reduced duration segment from the full recording, generating an absolute error value per

368    reduced duration segment. Due to the arbitrary nature of selecting a reduced recording segment

369    from a full recording, we repeated the procedure using 1000 bootstrapped samples from which a

370    mean absolute error was derived per recording per reduced recording duration. A minimum

371    reduced recording duration was identified as the duration at which the prediction performance,

372    measured using the mean absolute error, noticeably drops below that of the full duration Sinc

373    model. Finally, the mean absolute error of this reduced duration (20 mins – see Section 3.1.) 1-

374 channel Sinc model was compared to the mean absolute error of the 8-channel full duration RF
375 model using a one-sample paired t-tests (p<0.05 significance level). For the reduced duration Sinc
376 data, the initial 20 mins of each recording was selected for t-test analysis.

377

378 **2.5. Interpreting Sinc model performance using dataset $\mathcal{D}1$**
379 Deep neural networks are notorious for being black-box machines, limiting interpretability when
380 compared to machine learning approaches and traditional visual assessment approaches.
381 Nevertheless, methods are improving to visualize these networks to understand how they were
382 trained and their potential to generalise well on new data. In this study, two visualization techniques
383 were applied to further understand Sinc model performance: input-loss minimisation and uniform
384 manifold approximation and projection, see Supplementary Information S.4.

385

386 **2.6. External validation of Sinc model performance using dataset $\mathcal{D}2$**
387 The final Sinc model was trained on the entire dataset $\mathcal{D}1$ using the 1-channel bipolar setup (Figure
388 1). Similarly, the final "gold standard" RF model was trained on the entire dataset $\mathcal{D}1$ using the 8-
389 channel referential setup. When applying the final Sinc model to the independent hold-out dataset
390 $\mathcal{D}2$, the 1-channel bipolar setup was used and a 20 mins recording duration was randomly sampled
391 from the full duration EEG recording. When applying the final RF model to dataset $\mathcal{D}2$, the 8-channel
392 referential setup and full EEG recording were used.

393

394 *2.6.1. PMA prediction in independent hold-out dataset*
395 To assess the generalisability of the Sinc model to predict infants' PMA on independent data, the
396 normal BSID-II outcome data from the independent hold-out dataset $\mathcal{D}2$ was used. To assess the
397 association between true PMA and predicted PMA, a linear mixed effects regression model was
398 used (p<0.05 significance level). Random intercepts were introduced to group repeated recordings
399 from the same infant. Associations between true PMA and predicted PMA were also assessed for
400 the mild abnormal and severe abnormal groups. Similarly, the RF model was used to generate PMA
401 predictions for the normal, mild abnormal, and severe abnormal groups in dataset $\mathcal{D}2$, and
402 associations between true age and predicted age were assessed in an identical manner to the Sinc
403 model.

404

405 The two models' PMA prediction performances were compared using the linear mixed effects
406 regression models' z-statistic magnitudes per BSID-II outcome cohort. Additionally, Bland-Altman
407 analysis (Bland and Altman, 1999, 1986) was used to assess Sinc-RF model agreement in absolute
408 error magnitude, pooled across all data in dataset $\mathcal{D}2$. Bland-Altman analysis was implemented in
409 R v4.1.1 (R Core Team, 2018) using a publicly available package
410 (https://rdrr.io/cran/BlandAltmanLeh) to estimate the bias (Sinc minus RF) and limits of agreement,
411 along with 95% confidence intervals. Model agreement was assessed using individual recording
412 absolute errors and using within-infant multi-recording mean absolute errors, to assess the
413 influence of within-infant multi-recording averaging on model agreement (Bland-Altman plot y-axis,
414 limits of agreement width) and average prediction error magnitude (Bland-Altman plot x-axis,
415 range).

416

417 *2.6.2. Associating brain age delta magnitude to 9-month BSID-II follow-up outcomes*

418 The association between an infant's brain age delta magnitude and 9-month BSID-II follow-up

419 outcomes (normal, mild abnormal, severe abnormal) was assessed for all infants in dataset $\mathcal{D}2$. For

420 each infant, a brain age delta (absolute error) was determined per recording, and the mean absolute

421 error (MAE) across an infant's multiple recordings was used as an estimate of that infant's brain age

422 delta i.e. the deviation between their brain age and their true age. This per-infant MAE thus

423 represents an infant's overall brain neurodevelopmental trajectory deviation, with a larger

424 trajectory deviation corresponding to greater deviations from the norm.

425

426 Trajectory deviations across all infants in dataset $\mathcal{D}2$ were then grouped by neurodevelopmental

427 outcome (as defined in section 2.1.1.) and significant differences between groups assessed using

428 one-way ANOVA ($p<0.05$ significance level). Tukey's post-hoc test, which corrects for multiple

429 comparisons ($p<0.05$ significance level), was used to identify significant pair-wise comparisons.

430 Additionally, the two models' BSID-II outcome group separation performances were compared

431 using the pairwise standardised effect size (Cohen's D, estimated using MATLAB's meanEffectSize

432 function) magnitudes per contrast: mild minus normal, severe minus mild, and severe minus normal.

433

434 Finally, group-wise (normal, mild abnormal, severe abnormal) differences in GA, PMA and the

435 number of recordings in each infant's trajectory were checked using one-way ANOVA to assess their

436 potential influence as confounding factors.

437

438 **2.7 External validation of Sinc model performance using dataset $\mathcal{D}3$**

439 The final Sinc model was applied to the independent dataset $\mathcal{D}3$ collected at an independent centre

440 (Oxford, UK). The 1-channel bipolar montage (C3-C4) and the first 20-minutes of each recording

441 were used in the analysis.

442

443 The association between true PMA and predicted PMA was assessed using Pearson correlation (z-

444 statistic calculated using the Fisher r-to-z transform, $p<0.05$ significance level). Z-statistics are

445 reported for the results of both datasets $\mathcal{D}2$ and $\mathcal{D}3$. Each infant in dataset $\mathcal{D}3$ was recorded on a

446 single test occasion; the group-level MAE was calculated as the mean across all recordings of each

447 subject-level brain age delta i.e. each infant's error in predicted versus true age.

448

449 The brain age delta estimate can have a dependency with age – an age association bias that is known

450 to occur for several distinct reasons such as regression dilution (Smith et al., 2019). To correct for

451 this age association bias, we adjusted the predicted brain age using the linear regression between

452 the brain age delta and the true age (Smith et al., 2019). To assess the generalisability of this

453 correction to new data we adjusted the predicted brain age using leave-one-subject-out cross-

454 validation, calculating the MAE of the held-out subject compared with its true age.

455

456

457 **3. Results**

458 **3.1. The Sinc model outperforms alternative model architectures in predicting infant age (dataset**

459 $\mathcal{D}1$**)**

460 A comparison of model performance across the Sinc model and four alternative candidate deep
461 learning models, with reduced channel setups is summarised in Figure 3a and Supplementary Table
462 1. Using the 8-channel setup and the full recording duration data of dataset $\mathcal{D}1$, the Sinc model out-
463 performed both the established benchmark RF model (Sinc error = 0.73 weeks, RF error = 1.01
464 weeks, n = 47 recordings, t-statistic = 1.44, p = 0.078) as well as the candidate deep learning models.
465 When the number of recording channels was reduced from eight to one (bipolar channel, C3-C4),
466 the Sinc model had consistently lower MAE values compared with alternative models and exhibited
467 a total drop in performance of only 0.05 weeks (Sinc: 8-channel MAE = 0.73 weeks, 1-channel MAE
468 = 0.78 weeks). Furthermore, the 1-channel bipolar Sinc model outperformed the 8-channel
469 referential RF model (Sinc error = 0.78 weeks, RF error = 1.01 weeks, n = 47 recordings, t-statistic =
470 1.13, p = 0.13).

471

472 The Sinc model prediction error recorded from a single channel with full recording duration
473 (duration: median = 4h 25m, IQR = 4h 4m–7h 10m) was compared to Sinc model prediction error
474 using a single channel and reduced recording durations ranging from 0.5–120 mins (Figure 3b).
475 Using only 20 mins of EEG recording, the mean Sinc model prediction error was equivalent to using
476 the full recording duration. Using the established RF method as a benchmark, which relied on an 8-
477 channel setup and full-length recordings, the proposed Sinc model outperformed this benchmark
478 while having practical setup requirements that are far more achievable and practical for use in a
479 clinical environment (Sinc error = 0.79 weeks, RF error = 1.01 weeks, n = 47 recordings, t-statistic =
480 1.07, p = 0.14). While the three Sinc models' performances (8-channel full duration, 1-channel full
481 duration, 1-channel 20 min duration) did not statistically significantly differ to the benchmark RF
482 model performance, the Sinc model's performances were marginally but consistently improved (t-
483 statistics = 1.44, 1.13, 1.07, respectively, with positive t-statistics indicating larger MAE for RF).

484

485 **3.2. Sinc model may determine age using degree of EEG continuity (dataset $\mathcal{D}1$)**
486 To shed light on the specific EEG features that the deep learning Sinc model is likely utilising for the
487 brain age prediction, a method called input-loss minimisation was used to generate synthetic EEG
488 data that would force the model to make a brain age prediction of 30 weeks, 35 weeks, and 40
489 weeks PMA, respectively (Figure 4). Visually examining the synthetic EEG data shows that EEG
490 continuity and bursting were qualitatively distinguishing features and are therefore likely features
491 that the Sinc model used to characterise age-dependent activity. The 30-week synthetic data
492 reflects aspects of high discontinuity with short, high amplitude bursts and long-duration inter-burst
493 intervals (approximately 5-20 s) (Figure 4a). With increasing PMA, the inter-burst interval durations
494 decreased and burst periods widened, and by term-age, the signal was almost fully continuous with
495 no clear burst or inter-burst interval patterns (Figure 4c).

496

497 Using UMAP to visualise the data inputs to the three Sinc blocks (FEI, FEII, and Regression), a clear
498 separation of features occurs, beginning with a low-level followed by high-level feature extraction
499 (Figure 5). At the stage of inputs to Regression, the data can visually be seen to separate such that
500 datapoints increase almost monotonically with PMA (Figure 5c). This clear progression is indicative
501 that the network weights are trained well in the intermediate layers, and this visualisation provides
502 further insight into the role of each block.

12

503

504 **3.3. Sinc model brain age prediction generalises accurately to an independent hold-out dataset**
505 **(dataset $\mathcal{D}2$)**

506 Having established the Sinc model in dataset $\mathcal{D}1$ (section 3.1), this model was applied to a healthy
507 cohort of infants' data from the independent hold-out dataset $\mathcal{D}2$. Using 1-channel bipolar EEG data
508 of 20 min recording duration, the Sinc model's predicted ages were statistically significantly
509 correlated with infants' true PMA (Normal: n = 22 infants, z-statistic = 33.32, p < 0.0001) (Figure 6ai),
510 demonstrating that the model successfully generalises to independent data. The Sinc model also
511 generated predicted ages that were statistically significantly correlated with infants' true PMA for
512 the infants in dataset $\mathcal{D}2$ that had abnormal BSID-II follow-up outcomes (Mild abnormal: n = 10
513 infants, z-statistic = 18.03, p < 0.0001; Severe abnormal: n = 10 infants, z-statistic = 15.54, p < 0.0001)
514 (Figure 6aii). Infants with abnormal BSID-II follow-up outcomes were not used in training the Sinc
515 model, and so age predictions for these cohorts were, as expected, less accurate than those of the
516 healthy outcome cohort and thus exhibited weaker correlations (although still very strong) between
517 brain age and true age.

518

519 Using the 8-channel EEG setup and the entire recording duration, the RF model generated age
520 predictions that were statistically significantly correlated with infants' true PMA for both the normal
521 outcome and abnormal outcome cohorts (Normal: z-statistic = 22.89, p < 0.0001; Mild abnormal: z-
522 statistic = 12.51, p < 0.0001; Severe abnormal: z-statistic = 10.76, p < 0.0001) (Figure 6b). While the
523 brain age prediction correlation results for both the novel Sinc model and the established RF model
524 were very strong and highly significant for all three infant cohorts, the Sinc model consistently
525 outperformed the RF model per cohort (consistently larger z-statistics). Importantly, Sinc's
526 improved prediction accuracy was achieved while using dramatically lower EEG data requirements.

527

528 To quantitatively assess the level of agreement in PMA prediction performance between the RF and
529 Sinc models, we generated Bland-Altman plots of absolute prediction errors for the entirety of
530 dataset $\mathcal{D}2$ (pooled normal, mild abnormal, and severe abnormal outcome data) based on both
531 individual recordings (n = 141 recordings in total) (Figure 6ci) and individual infants (n = 42 infants
532 in total) (Figure 6cii). In both instances, there was a statistically significant negative bias reflecting
533 the reduced prediction error using the Sinc model (per-recording: mean bias = -0.202, 95% CI = [-
534 0.387, -0.016]; per-infant: mean bias = -0.231, 95% CI = [-0.444, -0.017]). Assessing the individual
535 recordings data, the limits of agreement were -2.435 and 2.032 with 95% CI = [-2.756, -2.115] and
536 [1.712, 2.353], respectively (Figure 6ci). Assessing the individual infants' data (multi-recording
537 average per infant), the limits of agreement were -1.573 and 1.112 with 95% CI = [-1.943, -1.204]
538 and [0.742, 1.481], respectively (Figure 6cii). The narrower limits of agreement width using the
539 infant-level assessment highlights a noticeable increase in Sinc-RF model agreement when using
540 multi-recording average prediction errors per infant rather than prediction errors based on
541 individual recordings, due to the reduced random noise variance as a consequence of the multi-
542 recording averaging. Using multi-recording average prediction errors per infant, we can expect 95%
543 of absolute prediction error differences between the RF and Sinc models to be approximately ±1.5
544 weeks, and the Sinc model to have a smaller prediction error of approximately 0.23 weeks on
545 average.

13

**3.4. Sinc model brain age deltas are associated with 9-month follow-up neurodevelopmental outcomes (dataset $\mathcal{D}2$)**

The variability in brain age delta magnitudes between infants with normal and abnormal BSID-II follow-up outcomes forms the foundation of the possibility of using brain age prediction to risk-stratify infants in the first few weeks of postnatal life according to neurodevelopmental outcomes. Here, using the Sinc model, the average brain age deltas for the normal, mild abnormal, and severe abnormal outcomes groups assessed using the BSID-II at nine months postnatal age were found to significantly differ (Normal: mean MAE = 0.71, n = 22 infants; Mild abnormal: mean MAE = 0.79, n = 10 infants; Severe abnormal: mean MAE = 1.27, n = 10 infants; one-way ANOVA: f-statistic = 4.24, p = 0.02) (Figure 7a). Significant differences between the mean deltas for the normal and severe abnormal groups were observed using post-hoc analysis adjusted for multiple comparisons (Tukey test: q-statistic = 4.20, p = 0.02) (Figure 7a). Taken together, these results indicate that Sinc model brain age delta magnitudes, generated using a single channel and 20 mins recording duration, scale with clinically informative BSID-II outcomes that are assessed several months later.

As reported previously, the RF model's brain age deltas also significantly differed between the three BSID-II outcome cohorts (Normal: mean MAE = 0.83, Mild abnormal: mean MAE = 1.13, Severe abnormal: mean MAE = 1.63, one-way ANOVA: f-statistic = 4.96, p = 0.01) (Figure 7b), with significant differences observed between the mean prediction errors for the normal and severe abnormal groups (Tukey test: q-statistic = 4.36, p = 0.01) (Figure 7b).

Quantitatively assessing the magnitude of the group average MAE separation between BSID-II outcome cohorts, a similar trend was observed for both the Sinc and RF models (Figure 7c). Both models exhibited poorest separation between the normal and mild abnormal outcome cohorts (group separation effect size: Sinc Cohen's D = 0.186; RF Cohen's D = 0.585), an intermediate degree of separation between the mild abnormal and severe abnormal outcome cohorts (group separation effect size: Sinc Cohen's D = 0.71; RF Cohen's D = 0.557), and greatest separation between the normal and severe abnormal outcome cohorts (group separation effect size: Sinc Cohen's D = 1.104; RF Cohen's D = 1.146) (Figure 7c).

No significant differences were identified between outcome groups for the potential confounding variables. Sinc model MAEs one-way ANOVA results (n = 42): GA: f-statistic = 0.93, p = 0.40; PMA: f-statistic = 0.51, p = 0.60; trajectory recording number: f-statistic = 0.28, p = 0.76).

**3.5 Sinc model accurately predicts brain age in data collected at an independent site (dataset $\mathcal{D}3$)**

The Sinc model was applied to an independent dataset collected at an independent centre (Oxford, UK; dataset $\mathcal{D}3$). The Sinc model's predicted ages were significantly correlated with the infant's true PMA (n = 73 infants, Pearson correlation coefficient r=0.91, z-statistic=1.52, p < 0.0001, Figure 8a), with good prediction accuracy (MAE = 0.97 weeks). This highlights that the Sinc model can generate age predictions using single recordings per infant for accurate group-level analysis at an independent site.

Unlike dataset $D2$, a noticeable bias in age prediction was visible in dataset $D3$ (Figure 8a). The magnitude of the brain age delta was significantly negatively correlated with the infant's true PMA (r =-0.24, p<0.01, Figure 8b). To generate unbiased brain age delta values, this age association should be minimised (Smith et al., 2019). A simple linear regression model trained on dataset $D3$, and validated using leave-one-out cross-validation, reduces this bias (Figure 8c). This additional linear model could be used in novel single-subject data collected at this site to produce brain age deltas with minimal age association bias. However, the biological value of the brain age deltas in dataset $D3$ has yet to be established. This dataset currently does not have follow-up BSID-II outcomes, so the association between brain age deltas and follow-up outcomes could not be assessed.

## 4. Discussion

This study presents the first deep learning architecture for the prediction of brain age from infant EEG activity. The model is based on a deep CNN structure incorporating the new Sinc block for enhanced multi-scale decompositions, with prediction likely utilising between-infant differences in their EEG continuity and bursting characteristics. Relative to previous proof-of-concept studies (Pillay et al., 2020; Stevenson et al., 2020a), the current deep learning approach was able to predict infant brain age with comparable accuracy and generate brain age delta magnitudes that were significantly associated with neurodevelopmental outcome at a 9-month follow-up using BSID-II assessment. Importantly, the current approach achieved this using dramatically reduced EEG data utilisation requirements, relying on only a single channel bipolar montage and 20 mins recording duration. This is important as it suggests that future systems utilising this method may only require single-channel capabilities which is simpler to set up and makes EEG data acquisition easier. This streamlined model, which can be applied in an objective and automated manner, thus demonstrates potential clinical utility for cot-side monitoring assessment of neurological well-being.

The chosen development strategy for the Sinc model involves training and testing the model first on a normal development dataset $D1$ and then additionally assessing performance in two independent datasets ($D2$ and $D3$, the latter collected at an independent site). Although we performed a single split on $D1$ for initial training and testing and could have used alternative techniques (such as cross validation), the goal was to assess relative performance with this dataset when comparing models, channel numbers, and recording durations. We kept the training and test splits in $D1$ consistent across these comparisons ensuring that relative differences in performance were meaningfully comparable. Furthermore, by showing high performance in the brain age prediction in the independent datasets, which was comparable to the held-out test set performance in $D1$, we can justify with confidence that the training strategies and choices made have still resulted in a robust generalisable model.

The model performed well on data collected at an independent site, despite differences in data collection such as EEG recording equipment and research personnel. This importantly suggests that the model is generalisable and could easily be employed for clinical use across multiple hospitals. Interestingly, an age association bias in model estimates could be observed between the predicted

632 age and true age when the model was applied to dataset $\mathcal{D}3$ (Oxford dataset), with the model likely
633 to overestimate age in the youngest infants and underestimate age in the oldest infants. The bias
634 was not observed in dataset $\mathcal{D}2$ (Leuven dataset). Bias in brain age predictions can arise from a
635 number of factors (Smith et al., 2019): for example, "regression dilution" due to errors in
636 measurement of the predictors (dataset $\mathcal{D}3$ used single recordings per infant, while dataset $\mathcal{D}2$
637 used multiple recordings per infant affording reduced measurement error). Using leave-one-
638 subject-out cross-validation, we demonstrated that it was possible to minimise this bias in dataset
639 $\mathcal{D}3$, suggesting that this correction would be generalisable for future infants collected at this centre.

641 Throughout our analyses, we used our previously published (Pillay et al., 2020) RF model as a "gold
642 standard" benchmark against which our novel Sinc model's performance was assessed. The RF
643 model used an 8-channel referential montage, over an hour of EEG recording, required sleep-staging
644 and an explicit pre-calculation of over 200 established features, while the Sinc model required only
645 a 1-channel bipolar montage and a 20 min recording duration, no sleep-staging, and included an
646 implicit feature extraction step. In all analyses, the Sinc model either performed comparably to or
647 out-performed the RF model. Additionally, in work published by an independent group (Stevenson
648 et al., 2020a), brain age deltas exhibited greatest separation between infants with normal and
649 severely abnormal BSID-II follow-up outcomes – an observation that is consistent with the current
650 study's findings, further supporting the results of the Sinc model.

652 Although a quantitative analysis of model speed was beyond the scope of this study, it is clear from
653 previous studies (Pillay et al., 2020; Stevenson et al., 2020a) that the requirement to extract multiple
654 features (some highly complex and non-linear), as well as the need to pre-stage the EEG based on
655 sleep state or states of discontinuity would slow performance, and this is suggested in a related
656 study on neonatal sleep-staging (Ansari et al., 2018). With the right accelerated hardware, however,
657 the proposed model (once trained) performs brain age predictions very quickly. This simplified
658 analysis pipeline lends itself well for hospital use if fast feedback is required in high-intensity
659 contexts, for instance, while the infant is in critical or post-operative care.

661 A further advantage of the Sinc model over the other deep learning architectures tested here is the
662 introduction of the Sinc block which, with a reasonable number of parameters, achieves a highly
663 non-linear architecture for performing multi-scale analysis (Ansari et al., 2021). The streamlined
664 preprocessing and feature extraction as well as the highly non-linear nature of the Sinc model are
665 invaluable attributes that provide flexibility for extraction of key signal characteristics and result in
666 a more focused feature set. The deep learning Sinc model is thus a flexible and efficient approach
667 for use with neonatal EEG data, which are data that typically exhibits highly variable and diverse
668 signal patterns.

670 Using the trained Sinc model to generate synthetic EEG data (Figure 4), our results suggest the
671 model's predictive performance may rely on identifying signal characteristics related to changes in
672 the EEG discontinuity with age (related to bursts and inter-burst intervals). This finding relates
673 sensibly to other findings in the current paper as well as established understanding of infant EEG
674 maturation. Regarding our present findings, the Sinc model's performance did not drop

16

675  substantially going from eight channels to one, or full recording duration to 20 mins. This might
676  suggest that the feature extraction stages of the architecture may be more tuned to global channel-
677  independent characteristics (such as bursting and continuity), as opposed to spatially-dependent
678  characteristics (such as inter-channel synchrony). Further, if the model relies on identifying changes
679  in burst/inter-burst cycling and encodes this in a highly multi-scale manner, this may indicate that
680  information on an infant's burst/inter-burst cycling may be sufficiently discernible from a 20-minute
681  EEG recording, with additional data providing diminished returns in discriminatory power.

682

683  Regarding infant EEG maturation, the progression of burst/inter-burst activity to continuous activity
684  is the expected characteristic developmental trajectory from preterm to term age (André et al.,
685  2010). Interestingly, these discontinuity patterns are also key for human experts when performing
686  visual age prediction (Dereymaeker et al., 2017a; Husain, 2005). Observing this link between the
687  synthetic inputs generated by the trained model and expected maturational trends strongly
688  suggests the Sinc model is relying on biophysiologically sensible signal features, which is important
689  for the generalisability of a model to novel data. We can tentatively suggest further similarities
690  between the Sinc model's generated synthetic EEG data and prominent features in the RF model. In
691  agreement with our previous work (Pillay et al., 2020), prominent features chosen by the
692  comparison RF model retrained in this study were based on the Line Length Burst %, a measure of
693  the percentage of burst periods in the EEG (Koolen et al., 2014), as well as measures of skewness of
694  the EEG amplitudes, which measure the asymmetry of a distribution compared to a Gaussian
695  distribution. Line Length Burst % would be expected to change with PMA as the burst periods
696  decrease with age and the EEG transitions to a more continuous pattern. Similarly, during this
697  transition, the distribution shifts away from a symmetrical Gaussian distribution as the number of
698  high positive bursts or spike amplitudes decreases. When comparing to the simulated results of Sinc
699  in Figure 4, we see similar behaviour is also identified by this trained neural network emphasising
700  the importance of this EEG characteristic across age.

701

702  We also note potentially interesting amplitude effects that are visible when looking at the model's
703  synthetic data across eight channels. For example, channels C3 and C4 have larger signal amplitudes
704  relative to other channels. While amplitude is a feature that changes with maturation (André et al.,
705  2010) making inter-subject variability in amplitude of potential value for brain age prediction, one
706  must be cautious when interpreting this subtler cross-channel amplitude effect in the synthetic
707  data. These amplitude effects may reflect a biophysiologically interesting phenomenon or may be
708  an artefactual consequence of proximity to the Cz reference electrode. Future work on the Sinc
709  model may help shed light on the potential role of amplitude effects.

710

711  Additionally, the role of motion artefacts, potentially related to sleep state and general motor
712  activity levels, could influence prediction performance. We applied a very simple amplitude-
713  threshold approach for artefact removal, and while this eliminates any major baseline drifts, periods
714  of recording drop-off or high-amplitude motion artefacts, some subtler artefacts likely remain. It is
715  unclear whether any residual motion effects influence prediction performance (either beneficially
716  or detrimentally). However, the lack of motion-like signals in the model-generated synthetic EEG
717  data suggests motion is unlikely to be playing a major role.

718

719 The ultimate interest in studying brain age delta magnitude is that neurological dysfunction can

720 manifest in infants' EEG as both accelerated or slowed maturation relative to a normative trajectory

721 (Scher, 1997; Watanabe et al., 1999), and these functional maturational deviations have prognostic

722 value (Iyer et al., 2015; Tokariev et al., 2019). The present study focused on the prognostic value of

723 preterm and term age resting-state brain function as a basis for risk-stratification using 9-month

724 BSID-II follow-up as the relevant outcome. However, as with any scale, there are limitations to BSID-

725 II predictive validity (Hack et al., 2005). Clinical decision making regarding the provision of

726 developmental care interventions (Burke, 2018) using deep learning-based predictions of infant

727 brain age would benefit from advancing the prognostic validity of the brain age delta metric. For

728 example, demonstrating associations between the metric and additional follow-up outcome

729 metrics, such as executive function (Dai et al., 2021), would improve validity. Additionally,

730 understanding the association between the metric and contemporaneous structural (e.g. body

731 weight, brain structural MRI) and functional (e.g. sensory-evoked neural and behavioural responses,

732 brain functional MRI) indices of development would be beneficial. We note that in the severe

733 outcome group of dataset $\mathcal{D}2$, a particularly large deviation was identified at 27.3 weeks PMA (see

734 Figure 6aii,bii). When investigating this infant's recording further (by AD), it was confirmed that the

735 baby was indeed very clinically unstable, with a history of seizure activity, generally suppressed

736 baseline EEG and alternating, abnormal rhythmic activity. Further investigations into associations

737 between the brain age delta magnitude and these contemporaneous and follow-up assessments

738 will be highly valuable in advancing model validity and appreciating the potential clinical value of

739 the Sinc brain age prediction model.

740

741 It is important to note that the focus of this manuscript was to provide an efficient diagnostic

742 approach for identifying abnormal brain maturation and to additionally show that this metric

743 correlates strongly with long term neurodevelopmental outcome. We do not, however, suggest a

744 cause for deviations between true age and brain age (i.e. brain age deltas) in this study nor that this

745 is directly associated to specific environmental or genetic causes. There is increasing evidence that

746 large brain age deltas may be a symptom of pre-existing conditions from birth (such as genetic

747 factors or low birth weight) which has a lasting impact on the infant's development presented

748 through alterations in brain age trajectories (Vidal-Pineiro et al., 2021). Regardless of the specific

749 causes of brain age deltas, it is clear that the magnitudes of these deviations are of biological and

750 clinical interest, and the ability to track and estimate brain age deviations with a model such as Sinc

751 provide a means to identify effects as soon as they manifest potentially allowing for rapid clinical

752 responses.

753

754

755 **5. Conclusions**

756 We outline a deep learning approach for infant brain age prediction and follow-up BSID-II outcome

757 risk-stratification with dramatically reduced EEG data requirements relative to previous proof-of-

758 concept studies. In an independent hold-out dataset, our Sinc model accurately predicts infant brain

759 age and significantly distinguishes infants with normal outcome from those with severely abnormal

760 outcome using a 1-channel bipolar montage setup and 20 min recording duration. The model also

18

761 accurately predicts infant brain age when applied to data collected at an independent site. This
762 objective and automated deep learning approach thus displays potential clinical utility for cot-side
763 monitoring and use in neurological function assessment. A major next objective will be the efficient
764 deployment of this model into the hospital setting using clinical grade bipolar montage data.

765

766

767 **Data availability statement**
768 Due to ethical restrictions and the sensitive nature of these data, it is not possible to publicly share
769 the supporting data.

770

771

772 **Code availability statement**
773 The underlying code for the deep learning models, including the training, validation, and testing
774 processes are openly available for download using the following GitHub link:
775 https://github.com/amirans65/brainagemodel.

776

777

778 **CRediT authorship contribution statement**

792

793

794 **Acknowledgements**

816
817

818   **Declaration of competing interests**
819   The authors declare no conflicts of interest.

820
821

822   **References**

823   André, M., Lamblin, M.-D., dAllest, A.M., Curzi-Dascalova, L., Moussalli-Salefranque, F.,
824      NguyenTheTich, S., Vecchierini-Blineau, M.-F., Wallois, F., Walls-Esquivel, E., Plouin, P., 2010.
825      Electroencephalography in premature and full-term infants. Developmental features and
826      glossary. Neurophysiologie Clinique/Clinical Neurophysiology 40, 59–124.
827      https://doi.org/10.1016/j.neucli.2010.02.002

828   Ansari, A.H., Cherian, P.J., Caicedo, A., Naulaers, G., De Vos, M., Van Huffel, S., 2019. Neonatal
829      Seizure Detection Using Deep Convolutional Neural Networks. Int J Neural Syst 29, 1850011.
830      https://doi.org/10.1142/S0129065718500119

831   Ansari, A.H., De Wel, O., Lavanga, M., Caicedo, A., Dereymaeker, A., Jansen, K., Vervisch, J., De Vos,
832      M., Naulaers, G., Van Huffel, S., 2018. Quiet sleep detection in preterm infants using deep
833      convolutional neural networks. J Neural Eng 15, 066006. https://doi.org/10.1088/1741-
834      2552/aadc1f

835   Ansari, A.H., De Wel, O., Pillay, K., Dereymaeker, A., Jansen, K., Van Huffel, S., Naulaers, G., De Vos,
836      M., 2020. A convolutional neural network outperforming state-of-the-art sleep staging
837      algorithms for both preterm and term infants. J Neural Eng 17, 016028.
838      https://doi.org/10.1088/1741-2552/ab5469

839   Ansari, A.H., Pillay, K., Dereymaeker, A., Jansen, K., Van Huffel, S., Naulaers, G., De Vos, M., 2021. A
840      Deep Shared Multi-Scale Inception Network Enables Accurate Neonatal Quiet Sleep
841      Detection with Limited EEG Channels. IEEE J Biomed Health Inform PP.
842      https://doi.org/10.1109/JBHI.2021.3101117

843   Bland, J.M., Altman, D.G., 1999. Measuring agreement in method comparison studies. Stat Methods
844      Med Res 8, 135–160. https://doi.org/10.1177/096228029900800204

845   Bland, J.M., Altman, D.G., 1986. Statistical methods for assessing agreement between two methods
846      of clinical measurement. Lancet 1, 307–310.

847   Blencowe, H., Lee, A.C.C., Cousens, S., Bahalim, A., Narwal, R., Zhong, N., Chou, D., Say, L., Modi, N.,
848      Katz, J., Vos, T., Marlow, N., Lawn, J.E., 2013. Preterm birth-associated neurodevelopmental

849       impairment estimates at regional and global levels for 2010. Pediatr Res 74 Suppl 1, 17–34.
850       https://doi.org/10.1038/pr.2013.204

851 Burke, S., 2018. Systematic review of developmental care interventions in the neonatal intensive
852       care unit since 2006. J Child Health Care 22, 269–286.
853       https://doi.org/10.1177/1367493517753085

854 Colonnese, M.T., Kaminska, A., Minlebaev, M., Milh, M., Bloem, B., Lescure, S., Moriette, G., Chiron,
855       C., Ben-Ari, Y., Khazipov, R., 2010. A Conserved Switch in Sensory Processing Prepares
856       Developing Neocortex for Vision. Neuron 67, 480–498.
857       https://doi.org/10.1016/j.neuron.2010.07.015

858 Dai, D.W.T., Franke, N., Wouldes, T.A., Brown, G.T.L., Tottman, A.C., Harding, J.E., PIANO Study
859       Group, 2021. The contributions of intelligence and executive function to behaviour problems
860       in school-age children born very preterm. Acta Paediatr 110, 1827–1834.
861       https://doi.org/10.1111/apa.15763

862 De Wel, O., Lavanga, M., Dorado, A.C., Jansen, K., Dereymaeker, A., Naulaers, G., Van Huffel, S.,
863       2017. Complexity Analysis of Neonatal EEG Using Multiscale Entropy: Applications in Brain
864       Maturation and Sleep Stage Classification. Entropy 19, 516.
865       https://doi.org/10.3390/e19100516

866 Dempsey, E.M., Kooi, E.M.W., Boylan, G., 2018. It's All About the Brain—Neuromonitoring During
867       Newborn Transition. Seminars in Pediatric Neurology, Fetal Neurology 28, 48–59.
868       https://doi.org/10.1016/j.spen.2018.05.006

869 Dereymaeker, A., Koolen, N., Jansen, K., Vervisch, J., Ortibus, E., De Vos, M., Van Huffel, S., Naulaers,
870       G., 2016. The suppression curve as a quantitative approach for measuring brain maturation
871       in preterm infants. Clin Neurophysiol 127, 2760–2765.
872       https://doi.org/10.1016/j.clinph.2016.05.362

873 Dereymaeker, A., Pillay, K., Vervisch, J., De Vos, M., Van Huffel, S., Jansen, K., Naulaers, G., 2017a.
874       Review of sleep-EEG in preterm and term neonates. Early Hum Dev 113, 87–103.
875       https://doi.org/10.1016/j.earlhumdev.2017.07.003

876 Dereymaeker, A., Pillay, K., Vervisch, J., Van Huffel, S., Naulaers, G., Jansen, K., De Vos, M., 2017b.
877       An Automated Quiet Sleep Detection Approach in Preterm Infants as a Gateway to Assess
878       Brain Maturation. Int J Neural Syst 27, 1750023.
879       https://doi.org/10.1142/S012906571750023X

880 Duerden, E.G., Guo, T., Dodbiba, L., Chakravarty, M.M., Chau, V., Poskitt, K.J., Synnes, A., Grunau,
881       R.E., Miller, S.P., 2016. Midazolam dose correlates with abnormal hippocampal growth and
882       neurodevelopmental outcome in preterm infants. Ann Neurol 79, 548–559.
883       https://doi.org/10.1002/ana.24601

884 Efron, B., Tibshirani, R.J., 1994. An Introduction to the Bootstrap. Chapman and Hall/CRC, New York.
885       https://doi.org/10.1201/9780429246593

886 Grunau, R.E., 2013. Neonatal pain in very preterm infants: long-term effects on brain,
887       neurodevelopment and pain reactivity. Rambam Maimonides Med J 4, e0025.
888       https://doi.org/10.5041/RMMJ.10132

889 Hack, M., Taylor, H.G., Drotar, D., Schluchter, M., Cartar, L., Wilson-Costello, D., Klein, N., Friedman,
890       H., Mercuri-Minich, N., Morrow, M., 2005. Poor predictive validity of the Bayley Scales of
891       Infant Development for cognitive function of extremely low birth weight children at school
892       age. Pediatrics 116, 333–341. https://doi.org/10.1542/peds.2005-0173

893 Husain, A.M., 2005. Review of neonatal EEG. Am J Electroneurodiagnostic Technol 45, 12–35.

894 Iyer, K.K., Roberts, J.A., Hellström-Westas, L., Wikström, S., Hansen Pupp, I., Ley, D., Vanhatalo, S.,
895       Breakspear, M., 2015. Cortical burst dynamics predict clinical outcome early in extremely
896       preterm infants. Brain 138, 2206–2218. https://doi.org/10.1093/brain/awv129

897  Koolen, N., Jansen, K., Vervisch, J., Matic, V., De Vos, M., Naulaers, G., Van Huffel, S., 2014. Line
898       length as a robust method to detect high-activity events: automated burst detection in
899       premature EEG recordings. Clin Neurophysiol 125, 1985–1994.
900       https://doi.org/10.1016/j.clinph.2014.02.015

901  Lavanga, M., De Wel, O., Caicedo, A., Jansen, K., Dereymaeker, A., Naulaers, G., Van Huffel, S., 2017.
902       Monitoring Effective Connectivity in the Preterm Brain: A Graph Approach to Study
903       Maturation. Complexity 2017, e9078541. https://doi.org/10.1155/2017/9078541

904  Malk, K., Metsäranta, M., Vanhatalo, S., 2014. Drug effects on endogenous brain activity in preterm
905       babies. Brain Dev 36, 116–123. https://doi.org/10.1016/j.braindev.2013.01.009

906  Milh, M., Kaminska, A., Huon, C., Lapillonne, A., Ben-Ari, Y., Khazipov, R., 2007. Rapid cortical
907       oscillations and early motor activity in premature human neonate. Cereb. Cortex 17, 1582–
908       1594. https://doi.org/10.1093/cercor/bhl069

909  Moultrie, F., Slater, R., Hartley, C., 2017. Improving the treatment of infant pain. Current Opinion in
910       Supportive and Palliative Care 11, 112–117.
911       https://doi.org/10.1097/SPC.0000000000000270

912  O'Shea, A., Ahmed, R., Lightbody, G., Pavlidis, E., Lloyd, R., Pisani, F., Marnane, W., Mathieson, S.,
913       Boylan, G., Temko, A., 2021. Deep Learning for EEG Seizure Detection in Preterm Infants. Int
914       J Neural Syst 31, 2150008. https://doi.org/10.1142/S0129065721500088

915  O'Toole, J.M., Boylan, G.B., Vanhatalo, S., Stevenson, N.J., 2016. Estimating functional brain
916       maturity in very and extremely preterm neonates using automated analysis of the
917       electroencephalogram. Clin Neurophysiol 127, 2910–2918.
918       https://doi.org/10.1016/j.clinph.2016.02.024

919  Palmu, K., Stevenson, N., Wikström, S., Hellström-Westas, L., Vanhatalo, S., Palva, J.M., 2010.
920       Optimization of an NLEO-based algorithm for automated detection of spontaneous activity
921       transients in early preterm EEG. Physiol Meas 31, N85-93. https://doi.org/10.1088/0967-
922       3334/31/11/N02

923  Pascal, A., Naulaers, G., Ortibus, E., Oostra, A., De Coen, K., Michel, S., Cloet, E., Casaer, A., D'haese,
924       J., Laroche, S., Jonckheere, A., Plaskie, K., Van Mol, C., Delanghe, G., Bruneel, E., Van
925       Hoestenberghe, M.-R., Samijn, B., Govaert, P., Van den Broeck, C., 2020.
926       Neurodevelopmental outcomes of very preterm and very-low-birthweight infants in a
927       population-based clinical cohort with a definite perinatal treatment policy. Eur J Paediatr
928       Neurol 28, 133–141. https://doi.org/10.1016/j.ejpn.2020.06.007

929  Pillay, K., Dereymaeker, A., Jansen, K., Naulaers, G., De Vos, M., 2020. Applying a data-driven
930       approach to quantify EEG maturational deviations in preterms with normal and abnormal
931       neurodevelopmental outcomes. Sci Rep 10, 7288. https://doi.org/10.1038/s41598-020-
932       64211-0

933  Pillay, K., Dereymaeker, A., Jansen, K., Naulaers, G., Van Huffel, S., De Vos, M., 2018. Automated
934       EEG sleep staging in the term-age baby using a generative modelling approach. J Neural Eng
935       15, 036004. https://doi.org/10.1088/1741-2552/aaab73
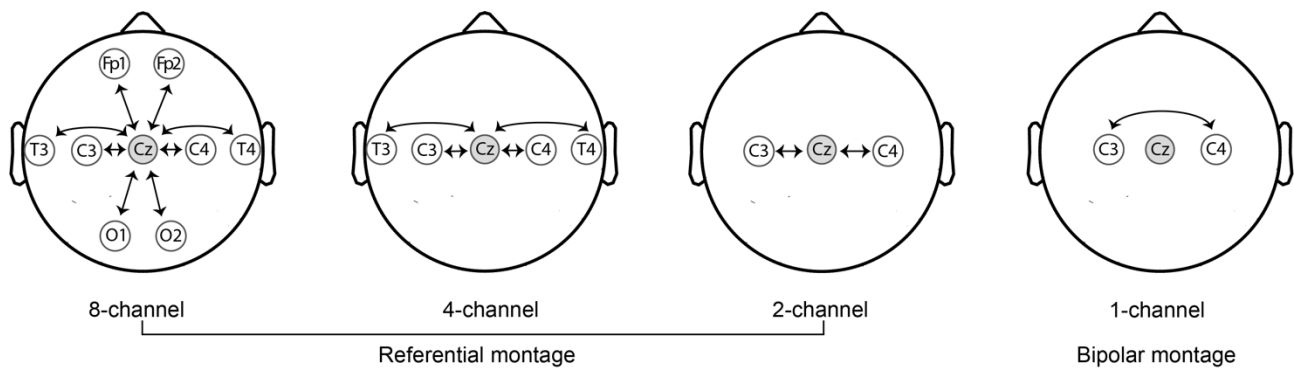
936  R Core Team, 2018. R: A language and environment for statistical computing.

937  Scher, M.S., 2008. Ontogeny of EEG-sleep from neonatal through infancy periods. Sleep Med 9, 615–
938       636. https://doi.org/10.1016/j.sleep.2007.08.014

939  Scher, M.S., 1997. Neurophysiological assessment of brain function and maturation. II. A measure
940       of brain dysmaturity in healthy preterm neonates. Pediatr Neurol 16, 287–295.
941       https://doi.org/10.1016/s0887-8994(96)00009-4

942  Smith, S.M., Vidaurre, D., Alfaro-Almagro, F., Nichols, T.E., Miller, K.L., 2019. Estimation of brain age
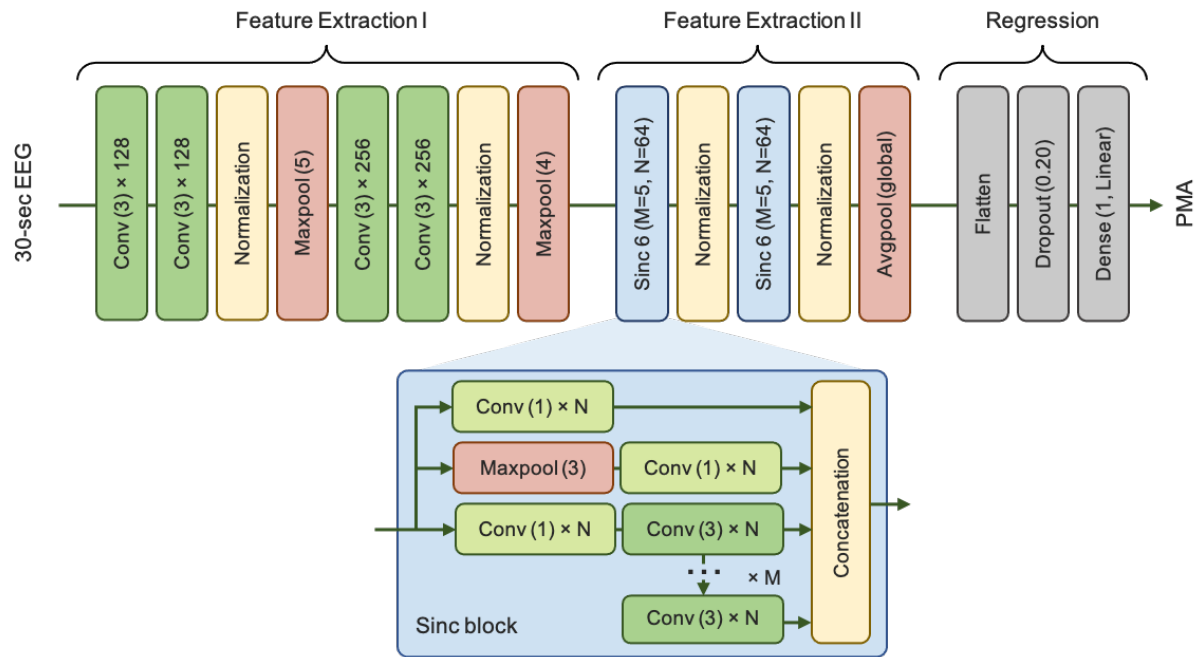943       delta from brain imaging. NeuroImage 200, 528–539.
944       https://doi.org/10.1016/j.neuroimage.2019.06.017

945 Stevenson, N.J., Oberdorfer, L., Koolen, N., O'Toole, J.M., Werther, T., Klebermass-Schrehof, K.,
946   Vanhatalo, S., 2017. Functional maturation in preterm infants measured by serial recording
947   of cortical activity. Sci Rep 7, 12969. https://doi.org/10.1038/s41598-017-13537-3

948 Stevenson, N.J., Oberdorfer, L., Tataranno, M.-L., Breakspear, M., Colditz, P.B., Vries, L.S. de,
949   Benders, M.J.N.L., Klebermass-Schrehof, K., Vanhatalo, S., Roberts, J.A., 2020a. Automated
950   cot-side tracking of functional brain age in preterm infants. Annals of Clinical and
951   Translational Neurology 7, 891–902. https://doi.org/10.1002/acn3.51043

952 Stevenson, N.J., Tataranno, M.-L., Kaminska, A., Pavlidis, E., Clancy, R.R., Griesmaier, E., Roberts,
953   J.A., Klebermass-Schrehof, K., Vanhatalo, S., 2020b. Reliability and accuracy of EEG
954   interpretation for estimating age in preterm infants. Ann Clin Transl Neurol 7, 1564–1573.
955   https://doi.org/10.1002/acn3.51132

956 Tokariev, A., Roberts, J.A., Zalesky, A., Zhao, X., Vanhatalo, S., Breakspear, M., Cocchi, L., 2019.
957   Large-scale brain modes reorganize between infant sleep states and carry prognostic
958   information for preterms. Nat Commun 10, 2619. https://doi.org/10.1038/s41467-019-
959   10467-8

960 Tolonen, M., Palva, J.M., Andersson, S., Vanhatalo, S., 2007. Development of the spontaneous
961   activity transients and ongoing cortical activity in human preterm babies. Neuroscience 145,
962   997–1006. https://doi.org/10.1016/j.neuroscience.2006.12.070

963 Vidal-Pineiro, D., Wang, Y., Krogsrud, S.K., Amlien, I.K., Baaré, W.F., Bartres-Faz, D., Bertram, L.,
964   Brandmaier, A.M., Drevon, C.A., Düzel, S., Ebmeier, K., Henson, R.N., Junqué, C., Kievit, R.A.,
965   Kühn, S., Leonardsen, E., Lindenberger, U., Madsen, K.S., Magnussen, F., Mowinckel, A.M.,
966   Nyberg, L., Roe, J.M., Segura, B., Smith, S.M., Sørensen, Ø., Suri, S., Westerhausen, R.,
967   Zalesky, A., Zsoldos, E., Walhovd, K.B., Fjell, A., 2021. Individual variations in 'brain age' relate
968   to early-life factors more than to longitudinal brain change. eLife 10, e69995.
969   https://doi.org/10.7554/eLife.69995

970 Wallois, F., Routier, L., Bourel-Ponchel, E., 2020. Impact of prematurity on neurodevelopment, in:
971   Gallagher, A., Bulteau, C., Cohen, D., Michaud, J.L. (Eds.), Handbook of Clinical Neurology,
972   Neurocognitive Development: Normative Development. Elsevier, pp. 341–375.
973   https://doi.org/10.1016/B978-0-444-64150-2.00026-5

974 Watanabe, K., Hayakawa, F., Okumura, A., 1999. Neonatal EEG: a powerful tool in the assessment
975   of brain damage in preterm infants. Brain Dev 21, 361–372. https://doi.org/10.1016/s0387-
976   7604(99)00034-0

977 Wess, J.M., Isaiah, A., Watkins, P.V., Kanold, P.O., 2017. Subplate neurons are the first cortical
978   neurons to respond to sensory stimuli. Proc Natl Acad Sci U S A 114, 12602–12607.
979   https://doi.org/10.1073/pnas.1710793114

980

981

982

**Figures**



8-channel    4-channel    2-channel    1-channel

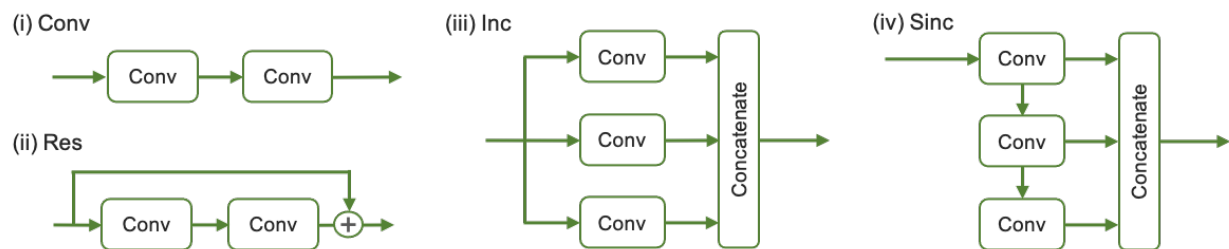Referential montage                    Bipolar montage

***Figure 1: EEG montages used during analysis.*** *All recordings in datasets $\mathcal{D}1$ and $\mathcal{D}2$ were acquired with eight recording EEG electrodes in positions: Fp1, Fp2, C3, C4, T3, T4, O1, O2, and a reference electrode placed at Cz (shaded in grey). The arrows represent the specific channels used during analysis. For dataset $\mathcal{D}3$, analysis was conducted using the 1-channel bipolar montage. Recordings were initially acquired with electrode positions Cz, CPz, C3, C4, Oz, FCz, T3 and T4, and a reference electrode at Fz.*

**Figure 2: Deep learning architectures. a.** *Block-diagram of the proposed Sinc network architecture, including the typical structure of the Sinc block.* **b.** *Illustrative block diagrams of different blocks in the deep architectures: (i) Sequential Convolutional layers, (ii) Residual block, (iii) Inception block, (iv) Shared Inception (Sinc) block.*
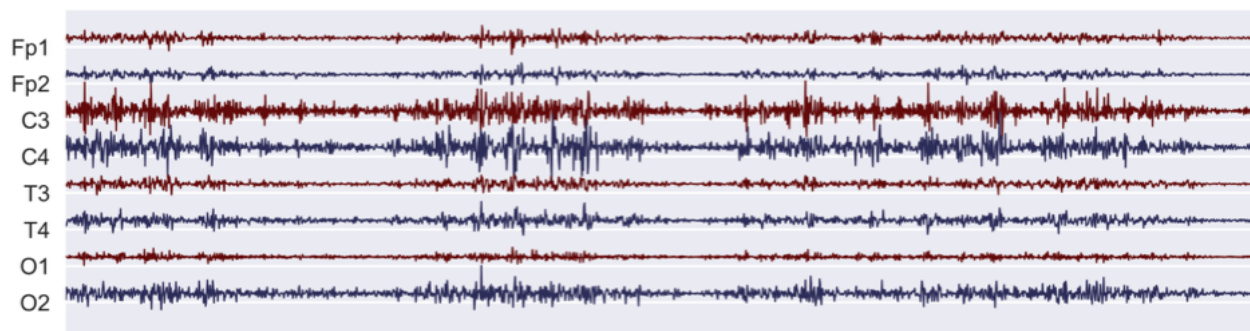
*Figure 3: The Sinc model outperforms alternative architectures in predicting infant brain age.* Brain age prediction performance (MAE) using dataset $\mathcal{D}1$ test set. *a. Each line represents a different model, and each model uses the entire recording duration. See Supplementary Table 1 for plotted values. The RF model is the established benchmark, which uses eight channels. The Sinc model consistently outperforms both the RF model and the alternative deep learning models, with a lower prediction error using a single channel (MAE = 0.78 weeks) than the RF model using eight channels (MAE = 1.01 weeks). b. The Sinc model's performance using a single channel and the full recording duration (MAE = 0.78 weeks, dotted line) was used as a benchmark to assess Sinc model performance with a single channel and systematically reduced recording durations (solid line). Performance using the reduced recording durations are matched to the full recording duration when recordings of 20 mins or longer are used; using less than 20 mins recording duration exhibits a gradual drop in prediction performance. Shaded intervals denote the standard deviation for the reduced recording durations. Note, MAE performance suggests a drop below the full signal performance beyond 20 min duration. This is due to the bootstrap sampling error (Efron and Tibshirani, 1994), and this inherent bias is a fluctuation about the full recording MAE with standard deviation <1. We can assume that the MAE beyond 20 mins is equivalent to the MAE when the full recording duration is used. As it is too computationally intensive to show performance beyond 2 hour signal durations the random variation cannot be fully shown here. Abbreviations: FE = feature extraction; CNN = convolutional neural network; Inc = inception; Sinc = shared inception; RF = random forest; ch = channel; MAE = mean absolute error.*
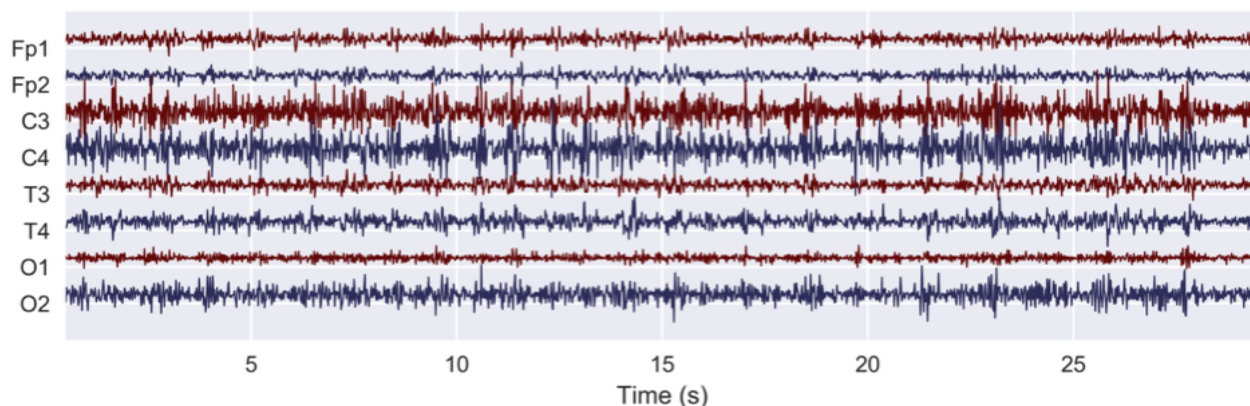
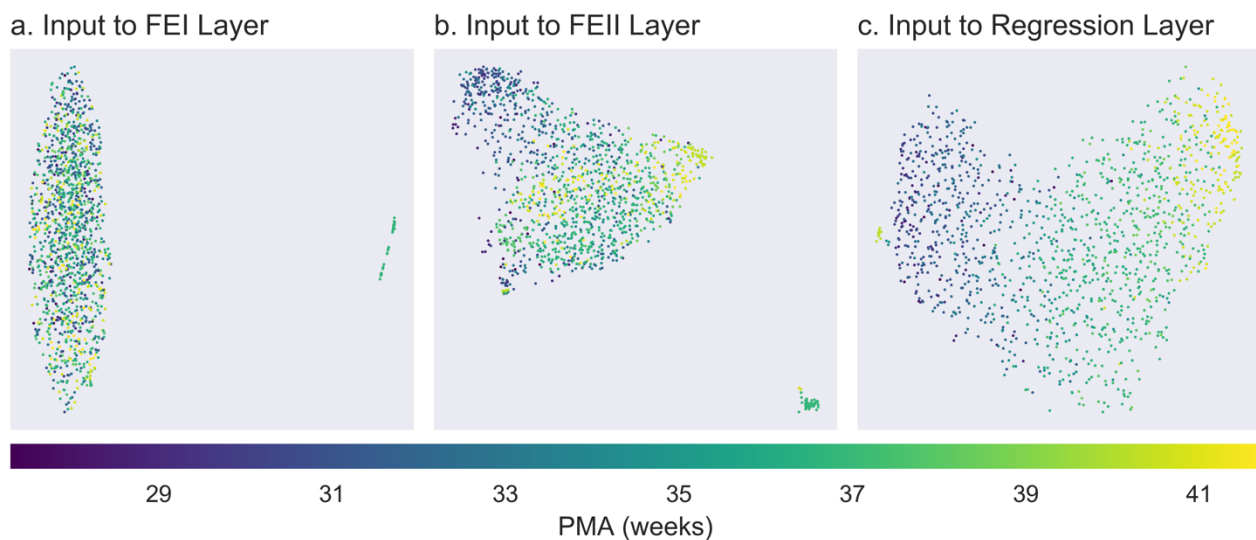**Figure 4: Synthetic EEG data generated using the Sinc model highlight changes in discontinuity characteristics with PMA, reminiscent of maturational trends seen in real EEG data.** *Results are generated using the input-loss minimization technique for three target PMAs (30, 35, and 40 weeks) spanning the early preterm to term age range. This is performed for the 8-channel full recording duration case. The degree of continuity in activity can be seen to increase with PMA.*
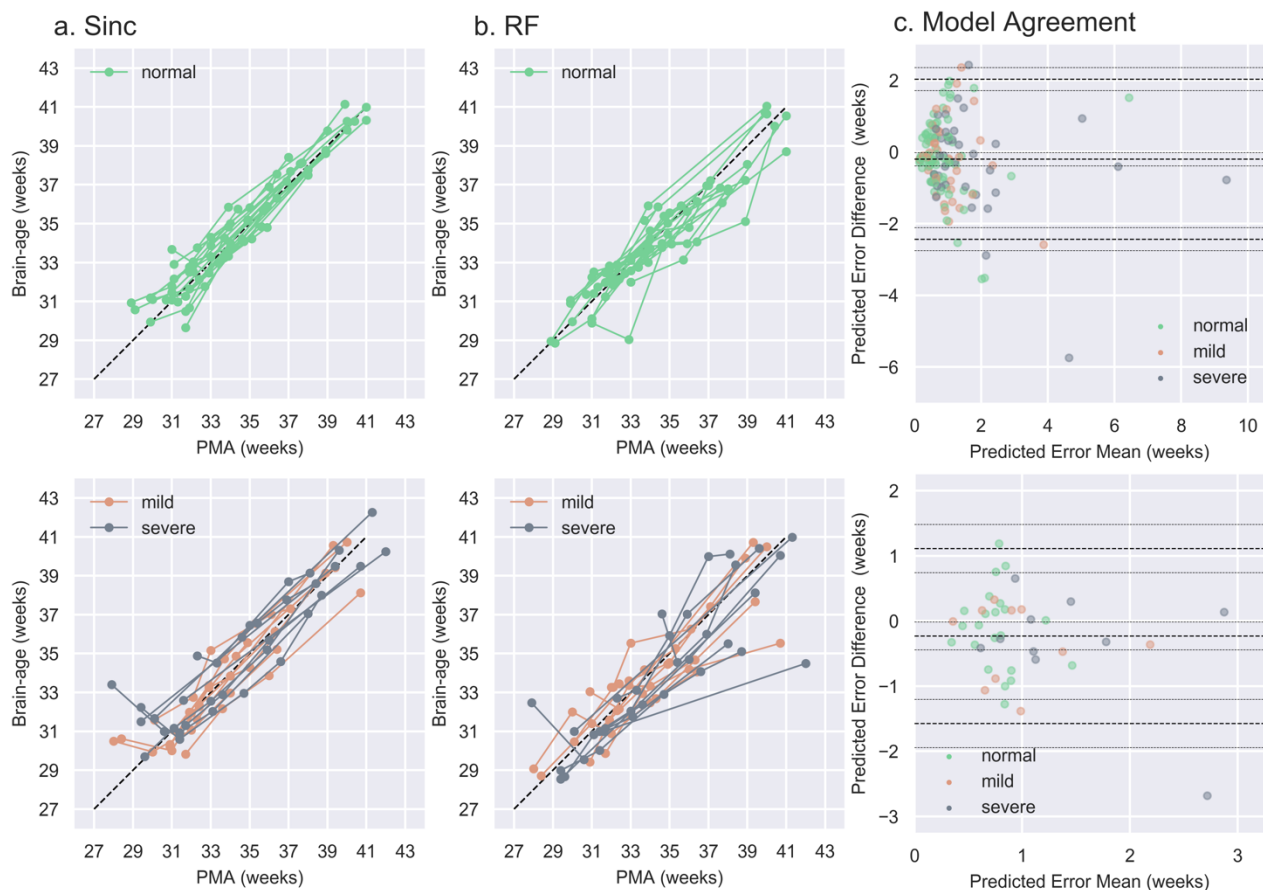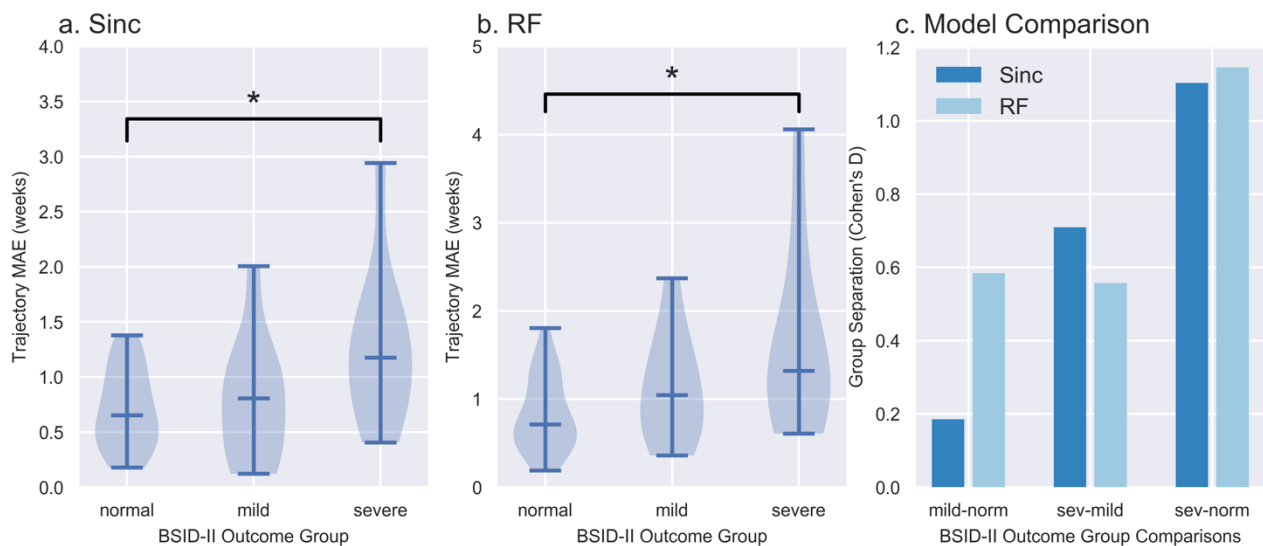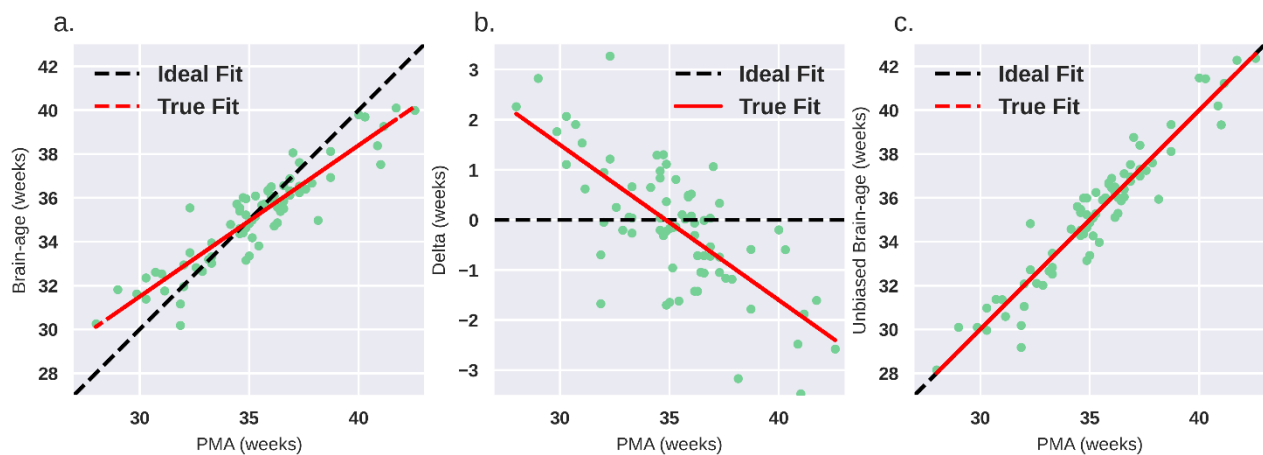
**Figure 5: Visualising Sinc model performance using UMAPs.** *Visualization of the inputs at various blocks in the proposed model: Feature Extraction I (FEI), Feature Extraction II (FEII), and Regression (see Figure 1b). Results are shown for the 1-channel, full recording duration case in Dataset $\mathcal{D}$1. An increasing separation of the features with respect to PMA is seen on moving from a-c. This clear progression is indicative that the network weights are trained well in the intermediate layers, and this visualisation provides further insight into the role of each CNN block. **a.** Input to FEI has not yet been processed, so there is no separation of inputs to FEI. **b.** The input to FEII is the output from FEI. It is evident that the role of FEI is to perform a low-level 'feature extraction' that performs an initial separation between the very preterm (blue dots) and preterm and term age groups (green and yellow dots) i.e. a general separation between strong discontinuity and continuity in the EEG. **c.** The input to Regression is the output of FEII. The FEII stage performs a higher-level feature extraction providing further discriminatory power, allowing better separation of these mid-age (31-37 weeks) and term age groups. Furthermore, at the stage of input to Regression, we observe that the PMAs of the datapoints from left to right increase almost monotonically such that the very left and right datapoints correspond to the extremely young and old neonates, respectively, while the middle ages are almost uniformly distributed in-between. Abbreviations: UMAP = uniform manifold approximation and projection.*

28

**Figure 6: Brain age prediction models generalise to independent dataset $\mathcal{D}2$. a.** *Sinc model brain age predictions for infants with (i) normal and (ii) abnormal BSID-II follow-up outcomes (Dataset $\mathcal{D}2$). Each string of connected points is a single infant's longitudinally-assessed multi-recording trajectory, and the dashed black line is the y=x line along which perfect predictions would lie.* **b.** *RF model brain age predictions for infants with (i) normal and (ii) abnormal BSID-II follow-up outcomes.* **c.** *Bland-Altman plots to assess agreement between Sinc and RF models' PMA prediction performances, quantified using absolute prediction errors. In both plots, the x-axis is the mean prediction error of the two models, and the y-axis is the difference in prediction errors (Sinc minus RF). The heavy grey lines are the mean bias and limits of agreement, while the light grey lines indicate the 95% CI for the bias and limits of agreement. (i) Per-recording model agreement assessment. (ii) Per-infant model agreement assessment i.e. multi-recording average per infant. Note the greater model agreement (narrower limits of agreement along y-axis) and reduced average prediction error (shorter range along x-axis) when using the multi-recording average prediction error in (ii) compared to the single recording prediction error in (i). Abbreviations: Sinc = shared inception; RF = random forest; PMA = postmenstrual age.*

29

L074

L075 ***Figure 7: Brain age delta magnitudes scale with 9-month follow-up neurodevelopmental***
L076 ***outcomes. a.*** *Sinc model absolute prediction error magnitudes (brain age deltas) for each of the*
L077 *three BSID-II outcome cohorts: normal, mild abnormal, and severe abnormal. The average prediction*
L078 *error is larger for poorer 9-month follow-up BSID-II neurodevelopmental outcomes, and the mean*
L079 *prediction error for the severe abnormal group is significantly larger than that of the normal group.*
L080 ***b.*** *RF model absolute prediction error magnitudes for each of the three BSID-II outcome cohorts. The*
L081 *average prediction error is larger for poorer 9-month follow-up BSID-II neurodevelopmental*
L082 *outcomes, and the mean prediction error for the severe abnormal group is significantly larger than*
L083 *that of the normal group.* ***c.*** *The x-axis displays each of the three combinations of pairwise*
L084 *comparisons for the three BSID-II outcome cohorts: mild minus normal, severe minus mild, and*
L085 *severe minus normal. For each model, the y-axis displays the standardised effect size (Cohen's D)*
L086 *separating each pair of BSID-II outcome cohort. Sinc = shared inception; RF = random forest; MAE =*
L087 *mean absolute error; BSID-II = Bayley scale of infant development; * = statistically significant.*

L088

L089

L090

**Figure 8: Sinc model brain age prediction generalises to dataset $\mathcal{D}3$.** *In each panel (a-c), each point indicates a single infant (n=73); the dashed black line is the ideal fit line; and the red solid line is the true fit line (least squares). **a.** Sinc model brain age predictions for dataset $\mathcal{D}3$. The ideal fit line is the y=x line of perfect prediction. The misalignment between the ideal fit line and the true fit line indicates an age association bias. **b.** Correlation between the brain age delta (predicted age minus true age) and the infant's true age. The ideal fit line is the y=0 line of zero age association bias. The slope of the true fit line indicates the magnitude and direction of the age association bias. **c.** The predicted brain age after adjusting for the delta age association bias using leave-one-out cross validation. The ideal fit line is the y=x line of perfect prediction.*