RESEARCH ARTICLE

WILEY

# Estimating cortical thickness trajectories in children across different scanners using transfer learning from normative models

C. Gaiser[1,2] | P. Berthet[3,4] | S. M. Kia[5,6,7] | M. A. Frens[1] |
C. F. Beckmann[5,8,9] | R. L. Muetzel[10,11] | Andre F. Marquand[5,8]

[1]Department of Neuroscience, Erasmus MC, University Medical Centre Rotterdam, Rotterdam, The Netherlands

[2]The Generation R Study Group, Erasmus MC—Sophia Children's Hospital, University Medical Centre Rotterdam, Rotterdam, The Netherlands

[3]Department of Psychology, University of Oslo, Oslo, Norway

[4]Norwegian Center for Mental Disorders Research (NORMENT), University of Oslo, and Oslo University Hospital, Oslo, Norway

[5]Donders Institute for Brain, Cognition, and Behavior, Radboud University, Nijmegen, The Netherlands

[6]Department of Psychiatry, Utrecht University Medical Center, Utrecht, The Netherlands

[7]Department of Cognitive Science and Artificial Intelligence, Tilburg University, Tilburg, The Netherlands

[8]Department of Cognitive Neuroscience, Radboud University Medical Center, Nijmegen, The Netherlands

[9]Centre for Functional MRI of the Brain, University of Oxford, Oxford, UK

[10]Department of Child and Adolescent Psychiatry, Erasmus MC—Sophia Children's Hospital, University Medical Centre Rotterdam, Rotterdam, The Netherlands

[11]Department of Radiology and Nuclear Medicine, Erasmus MC—Sophia Children's Hospital, University Medical Centre Rotterdam, Rotterdam, The Netherlands

**Correspondence**
Andre F. Marquand, Donders Centre for Cognitive Neuroimaging, Donders Institute for Brain, Cognition and Behaviour, Kapittelweg 29, 6525 EN Nijmegen, The Netherlands.
Email: andre.marquand@donders.ru.nl

**Funding information**
Erasmus MC Fellowship, Grant/Award Number: 2021.042; European Research Council, Grant/Award Number: 10100118; Netherlands Organization for Scientific Research Vici Grant, Grant/Award Number: 17854; NWO-CAS, Grant/Award Number: 012-200-013; Sophia Foundation, Grant/Award Number: S18-20; The Netherlands Organization for Scientific Research (NWO), and the Ministry of Health, Welfare, and Sport; Wellcome Trust, Grant/Award Number: 215698/Z/19/Z; Wellcome Trust Collaborative Award in Science, Grant/Award Number: 215573/Z/19/Z; ZonMw; Dutch Scientific Organization

## Abstract

This work illustrates the use of normative models in a longitudinal neuroimaging study of children aged 6–17 years and demonstrates how such models can be used to make meaningful comparisons in longitudinal studies, even when individuals are scanned with different scanners across successive study waves. More specifically, we first estimated a large-scale reference normative model using Hierarchical Bayesian Regression from $N = 42{,}993$ individuals across the lifespan and from dozens of sites. We then transfer these models to a longitudinal developmental cohort ($N = 6285$) with three measurement waves acquired on two different scanners that were unseen during estimation of the reference models. We show that the use of normative models provides individual deviation scores that are independent of scanner effects and efficiently accommodate inter-site variations. Moreover, we provide empirical evidence to guide the optimization of sample size for the transfer of prior knowledge about the distribution of regional cortical thicknesses. We show that a transfer set containing as few as 25 samples per site can lead to good performance metrics on the test set. Finally, we demonstrate the clinical utility of this approach by showing

C. Gaiser and P. Berthet contributed equally to this work.

that deviation scores obtained from the transferred normative models are able to detect and chart morphological heterogeneity in individuals born preterm.

**Practitioner Points**

- We show successful transfer learning from large-scale normative models to a new cohort.
- As few as 25 scans per site are needed to adapt prior knowledge from normative models to new sites.
- Resulting deviation scores from the normative model are free of site-effects and are able to uncover morphological heterogeneity in individuals born preterm.

## 1 | INTRODUCTION

Identifying structural or functional biomarkers of psychiatric and neurological illnesses across the lifespan has received increasing attention in recent years. Many of these disorders present symptoms that begin during childhood and adolescence (Bayer et al., 2021; Rogers & de Brito, 2016; Solmi et al., 2022; Whittle et al., 2020). There is, however, large interindividual heterogeneity in symptoms and underlying biology (DeLisi, 2008; Fuhrmann et al., 2022; Mills et al., 2021; Tamnes et al., 2017), making it challenging to pinpoint the precise underlying neurobiological substrates. Longitudinal datasets provide particularly valuable insights on the temporal evolution of brain development and offer considerable potential to understand the emergence of psychopathology and to parse this heterogeneity across individuals.

To detect and understand this heterogeneity and atypicality, there is a need to better characterize typical neurodevelopment (Insel, 2014; Volpe, 2009). In recent years, the availability of large datasets has greatly assisted efforts to understand interindividual variability in brain development (Bethlehem et al., 2022; Rutherford, Fraza, et al., 2022). For example, large scale studies using cortical volume, cortical thickness (CT) and surface area have identified a general decrease in these metrics with age, after adolescence (Bethlehem et al., 2022; Frangou et al., 2022; Rutherford, Fraza, et al., 2022; Tamnes et al., 2017; Thambisetty et al., 2010). CT has been shown to more accurately reflect underlying pathophysiological mechanisms than gray matter volume analysis (Clarkson et al., 2011; Hutton et al., 2009; Pereira et al., 2012; Zhao et al., 2022). However, these large data resources have expanded in scale via large, long-running longitudinal cohort studies. While the benefits of these large and unique cohorts are obvious, such studies also impose particular difficulties. For example, data must often be aggregated across multiple study centers, which necessitates dealing with site effects and, across developmental time scale, subjects are often scanned with different scanner hardware and/or software at successive timepoints. As a result, there is often little or no overlap in terms of age of participants and site effects in successive acquisition waves. Such nontrivial differences across sites, scanners, and timepoints have been difficult to account for statistically in analyses. Therefore, in addition to longitudinal data, novel methodological tools that map interindividual differences are needed to generate new insights.

Normative modeling approaches have recently emerged as a tool to better understand longitudinal developments with neuroimaging data (Marquand et al., 2019; Marquand, Rezek, et al., 2016). These approaches produce statistical inference at the individual level, without relying on strong assumptions about clustering of individuals or population structure (Antoniades et al., 2021; Cole, 2012; Marquand et al., 2019). Instead, symptoms in individual patients can be related to extreme deviation from the normative range (Fraza et al., 2021; Marquand, Wolfers, et al., 2016; Zabihi et al., 2019). This has shown the potential to detect morphological differences in patient populations which were not evident using standard techniques (Remiszewski et al., 2022). Additionally, a *Hierarchical Bayesian Regression* (HBR) approach to normative modeling has been shown to efficiently accommodate inter-site variation and to provide good computational scaling, which is important when using large studies, longitudinal studies, or combining smaller studies, that are acquired across multiple sites (Bayer et al., 2021; Kia et al., 2022; Rutherford, Fraza, et al., 2022). It also supports federated (i.e., decentralized) multisite normative modeling to transfer previously trained models onto unseen sites, while benefiting from the training on the large reference datasets (Kia et al., 2022; Rutherford, Kia, et al., 2022). This is especially interesting given that, in longitudinal studies running over several years, changes of scanner hardware, software and/or scan protocols are the norm rather than the exception, which generates a need to correct for the resulting scanner effects.

In this work, we provide a case study in using the transfer of prior knowledge about CT distributions from normative models derived from a large reference (e.g. lifespan) cohort to better estimate parameters on a smaller target (e.g. clinical) cohort. For this, we use longitudinal CT data from the Generation R study (Jaddoe et al., 2006; Kooijman et al., 2016; White et al., 2018), which contains data from children aged 6–17 years scanned in two different scanners, unseen by the reference models. The narrow age-range makes this study a good candidate for transfer learning in that it is beneficial to transfer information learned from a large lifespan cohort to obtain precise estimates of the slope or trajectory of developmental effects across a narrower age range. This method provides important advantages: one, it allows meaningful comparison of individuals scanned on different scanners, while taking advantage of previous knowledge, built from

large publicly available datasets to set informed hyperpriors: expected mean and variance of the distribution of samples for each region of interests. This, in turn, provides three benefits to the study, on providing more accurate predictions from the models thanks to the use of the mentioned informed priors; second this enables to reduce the ratio of training samples necessary to learn developmental trajectories for to the unseen sites, thereby enabling more participants to be allocated to the test set, and thus improving statistical power (Pan & Yang, 2010). Third, we will show that it provides a means to draw meaningful inferences within individuals across timepoints, even when follow-up scans are derived from a different scanner. This work also aims to offer some guidance on the methodology, for example, providing empirical estimates of the number of samples required for the transfer of knowledge from previous learnings and choices in transfer configurations, for example, factors included as batch effects. Finally, we provide a demonstration of the clinical utility of this approach by using it to understand interindividual differences in brain morphology resulting from preterm birth.

## 2 | METHODS

### 2.1 | Normative modeling

We estimated normative models using HBR to predict CT from age, sex and scanner site, for each *region of interest* (ROI) using the freely available PCNtoolkit python package, version 0.22 (Rutherford, Kia, et al., 2022).

### 2.1.1 | Reference models

We assembled a large reference cohort containing $n = 42,993$ (95%) healthy individuals to train the normative models before validating this model on $n = 2682$ controls and patients (5%, stratified by sites)

from a collection of mostly publicly available MRI datasets across 77 sites and 45,675 participants (see Tables 1 and 2 for details). The reference model is available on the PCNportal (https://pcnportal.dccn.nl/). CT measures were obtained from FreeSurfer processing (versions 5.3 or 6.0), as referred in the publications associated with the datasets (Dale et al., 1999; Fischl et al., 1999, 2002; Fischl & Dale, 2000). The Destrieux atlas was used to parcellate the brain into ROIs (Destrieux et al., 2010). One normative model was estimated per ROI. Linear HBR models were estimated using fixed effects of age and batch (i.e., random) effects for site and sex. In practice, this allows each site and sex to have different slopes, intercepts and variances. We included only data from the first visit when multiple visits were available (i.e., UKBB and ABCD). Only participants with complete data on ROIs were included.

Estimated reference models performed well according to accuracy metrics (explained variance: mean = 0.44, SD = 0.13, *standardized mean squared error* (SMSE): mean = 0.55, SD = 0.13, and *mean standardized log loss* (MSLL): mean = −0.37, SD = 0.14). Outputs include hyperparameters defining the mean and variance of the site-specific mean effects and variance, estimated during the training over the collection of datasets. This can be used as informed priors when adapting the normative models to unseen target sites. These hyperparameters are adapted to the unseen site using a holdout subset of the target dataset, that is, the adaptation set. This allows to reduce the number of samples used for adaptation while retaining a low variance of the estimations.

### 2.1.2 | Target cohort

As target cohort, 8523 $T_1$-weighted MRI scans from the population-based longitudinal Generation R study (Jaddoe et al., 2006; Kooijman et al., 2016) were used. In short, the Generation R study is a prospective cohort study from fetal life until adulthood that is designed to find early markers for typical and atypical development, growth, and

**TABLE 1** Overview and demographics of participants in the reference cohort.

| Datasets | N | Number of scanners | Age range | Gender ratio F/M |
|---|---|---|---|---|
| ABCD | 9605 | 29 | 9–11 | 48.8/51.2 |
| CAMCAN | 582 | 1 | 8–88 | 52.2/47.8 |
| CMI | 633 | 2 | 5–22 | 40.9/59.1 |
| FCON | 928 | 18 | 8–79 | 58.5/41.5 |
| HCP | 1049 | 1 | 22–37 | 54.0/46.0 |
| HCPAD | 1262 | 5 | 8–100 | 54.4/45.6 |
| HCPEP | 56 | 4 | 17–36 | 35.8/64.2 |
| IXI | 529 | 1 | 20–86 | 56.1/43.9 |
| NKI | 438 | 1 | 7–85 | 64.8/35.2 |
| OASIS | 1542 | 5 | 43–97 | 61.2/38.8 |
| OPN | 580 | 6 | 8–58 | 55.5/44.5 |
| PNC | 1344 | 1 | 8–23 | 53.3/46.7 |
| UKBB | 24,445 | 3 | 44–82 | 52.9/47.1 |
| Total | 42,993 | 77 | 5–100 | 52.5/47.5 |

| Datasets | N (N patients) | Number of scanners | Age range | Gender ratio F/M |
|---|---|---|---|---|
| ABCD | 499 (0) | 28 | 9–11 | 44.5/55.5 |
| CAMCAN | 32 (0) | 1 | 18–82 | 34.4/65.6 |
| CMI | 38 (0) | 2 | 6–22 | 31.6/68.4 |
| FCON | 67 (23) | 14 | 15–78 | 44.8/55.2 |
| HCP | 64 (0) | 1 | 22–36 | 60.9/39.1 |
| HCPAD | 53 (0) | 5 | 8–83 | 49.0/51.0 |
| HCPEP | 123 (122) | 4 | 17–35 | 39.0/61.0 |
| IXI | 24 (0) | 1 | 20–68 | 62.5/37.5 |
| NKI | 20 (0) | 1 | 11–84 | 60.0/40.0 |
| OASIS | 349 (274) | 5 | 46–96 | 50.1/49.9 |
| OPN | 27 (0) | 6 | 8–35 | 51.8/48.2 |
| PNC | 76 (0) | 1 | 9–21 | 47.4/52.6 |
| UKBB | 1310 (0) | 3 | 45–80 | 54.3/45.7 |
| Total | 2682 (419) | 72 | 6–96 | 50.4/49.6 |

**TABLE 2** Overview and demographics of participants in the validation set.

health. Almost 10,000 pregnant women living in Rotterdam, The Netherlands, were enrolled in the study between 2002 and 2006. Data from the children and caregivers was collected at several time points and written informed consent and/or assent was obtained from all participants. All study procedures were approved by the Medical Ethics Committee of the Erasmus Medical Center. The imaging protocol and quality assessment is extensively described by White et al. (2018). MRI scans were acquired in three waves using two different scanners, making the cohort an ideal validation set to investigate the transfer of hyperparameters from a reference dataset to an unseen target set. In longitudinal studies running over several years, changes of scanner hardware, software and/or scan protocols are often inevitable, which generates a need to correct for the resulting scanner effects. In the first wave, 1033 participants (484 females, age range: [6–10]) were imaged with a 3 T MR750 Discovery MRI scanner, while in the second ($n = 3920$, 1977 female, age range: [9–12]) and third wave ($n = 3570$, 1866 female, age range: [13–17]) a 3 T MR750w Discovery scanner (General Electric, Milwaukee, WI, USA) was used. After exclusion of scans with incidental findings ($n = 58$), braces ($n = 1067$), and low-quality visual inspection ratings of FreeSurfer reconstructions ($n = 2067$), a total of 6285 scans were included in the target dataset. Since several exclusion factors can be present in a single scan, Supplementary Figure 1 illustrates the extend of overlap. Figure 1 shows a histogram of age and scanner distributions in the target dataset.

## 2.2 | Transfer of hyperparameters from reference models to target cohort

By making use of the Generation R study cohort, we set out to show the advantage of transferring the hyperparameters to an unseen site by (1) determining the optimal number of samples needed for adaptation to the target cohort, (2) validating the recalibration of data to the
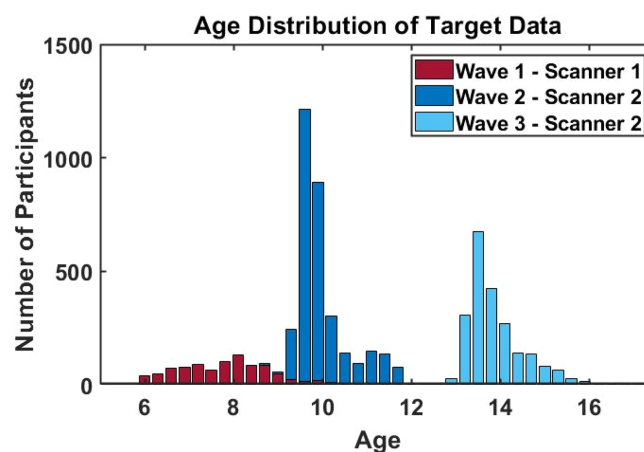


**FIGURE 1** Histogram of the scanning waves and age distributions in the Generation R target dataset.

target cohort and successful removal of site-effects by comparing raw and scanner corrected values, and (3) illustrating the utility of site-corrected deviations scores to uncover changes in morphology between groups and individuals. In the following, these three aims are described in more detail.

### 2.2.1 | Optimal sample size for parameter adaptation

In order to determine the optimal number of samples in the adaptation set, we leveraged the large amount of data available in the Generation R cohort. As described above, to prevent bias, held-out data should be used for adapting the parameters of the normative model to the target cohort (see Kia et al., 2020, 2022 for details). The number of scans in the adaptation set was varied ranging from 5 to

300 scans and model metrics (explained variance, SMSE, MSLL) of the subsequent models were calculated for each sample size. The resulting information is particularly useful for small imaging cohorts, since cohorts with smaller sample sizes can employ the current approach to boost power by making use of the hyperpriors inferred from large data. Yet, this is only viable if the samples needed to recalibrate the models can be kept to an optimal minimum. As a reference for model performance using the transfer approach, we estimated normative models using only the Generation R cohort. Models using only Generation R data were generated with a 50/50 split in training set ($n = 3143$) and test set ($n = 3142$) while ensuring equal age and sex distributions. Parameters were identical to the ones used in the reference model.

## 2.2.2 | Validation of adaptation

Additionally, two aspects of the Generation R study design make the cohort an ideal target set to validate the successful recalibration of the normative models to an unseen site. First, scans of participants that have repeated measurements over all three scanning time points are present. Uncorrected CT values of a participant with scans across all three measurement waves (and therefore across both scanners) show heterogeneity over time points that is partly due to biological changes over time and partly due to confounding site-effects. After successful recalibration of the normative model, we expect resulting z-scores, which are in principle free of site-effects, of the same participant to be in a similar range while raw values will differ. Second, there is an overlapping age range (8.6–10.7 years of age), in which scans from both scanners were obtained (Figure 1). Z-scores of participants from wave 1 (scanner 1), that fall in the overlapping age range of waves 1 and 2 should be distributed similarly after recalibration as z-scores in the same age range of wave 2 (scanner 2), while raw, uncorrected values differ due to scanner effects. Therefore, scans of participants with measurements at all three scanning time points ($n = 1317$) and scans from the first imaging wave that fall in the overlapping age range ($n = 211$) were withheld from the adaptation set that was used to recalibrate the reference normative model to the new unseen site. As outlined above, these scans hold valuable information that will be used to determine the successful calibration of the models by comparing raw CT before adaptation and corrected estimates after adaptation.

## 2.2.3 | Clinical application of normative estimates

Finally, we used the resulting site-effect free estimates to illustrate their potential to uncover morphological deviations in clinical cohorts by contrasting estimates in CT per ROI between participants in the Generation R cohort born preterm (gestational age < 37 weeks, $n = 339$) and children born at term ($n = 5646$). Pre-term birth interrupts a vulnerable period for brain development, as processes such as synaptogenesis, axonal growth, and neuronal migration, take place during the third semester (Volpe, 2009). Therefore, deviation scores from the normative

models can for instance be used to explore the variability in CT within children born preterm, but also to find ROIs that differ between children born preterm and at term. Notably, these deviation scores are free of site-effects and therefore especially suited for longitudinal MRI designs, as it is the case with the Generation R study.

# 3 | RESULTS

## 3.1 | Transfer results

### 3.1.1 | Optimal number of samples for parameter adaptation

We first determined the optimal number of subjects needed in the adaptation set. Figure 2 shows evaluation metrics for each ROI as the sample size of the adaptation set increases. Performance of the model reaches a plateau around 100 subjects. We thus adapted the initial reference models to the unseen sites of the Generation R study on $n = 300$ (4.8%) ($n = 100$ for scanner 1 in wave 1; $n = 200$ for scanner 2 in waves 2 and 3) and tested the models on the remaining participants ($n = 5985$; $n = 813$ for scanner 1 in wave 1, $n = 5172$ for scanner 2 in waves 2 and 3). For each adaptation set, subjects from wave 1 and 3 were sampled randomly, whereas subjects from wave 2 were sampled pseudorandomly to ensure a uniform cover of the full range of the narrow and highly peaked age distribution in this wave (Figure 1). Sampling was done using the *datasample* function (without replacement) in MATLAB (Mathworks, USA). While model performance reached a performance ceiling at approximately 100 scans per scanner/wave in the adaptation sample, only slight concessions in model performance are present as adaptation sample size decreases to only 25 scans. In Table 3 performance metrics are reported for transfer learning with 25 and 100 scans per adaptation sample, as well as for normative models trained on half of the Generation R cohort without the use of hyperparameters from the reference models (reference for baseline performance for models covering a narrower age range compared to the reference lifespan model).

As can be seen in Figure 2, model performance seems to differ in cerebral regions, with better evaluation metrics in occipital and frontal compared to the remaining ROIs. We therefore tested the influence of average ROI area (as reported by Destrieux et al., 2010) on model performance. We find that ROI area and evaluation metrics were correlated ($r_{\text{Explained Variance}} = .21$, $p_{\text{Explained Variance}} = .012$; $r_{\text{SMSE}} = -.21$, $p_{\text{SMSE}} = .011$; $r_{\text{MSLL}} = -.18$, $p_{\text{MSLL}} = .025$), with larger ROIs outperforming smaller ROIs (Supplementary Figure 2).

### 3.1.2 | Adaptation settings

We furthermore tested different adaptation settings. Scans in waves 2 and 3 of the target cohort were acquired on the same scanner, however at different time points. Therefore, we compared two different adaptation settings: First, treating waves 2 and 3 as the same site, and second, treating waves 2 and 3 as different sites. When treated as
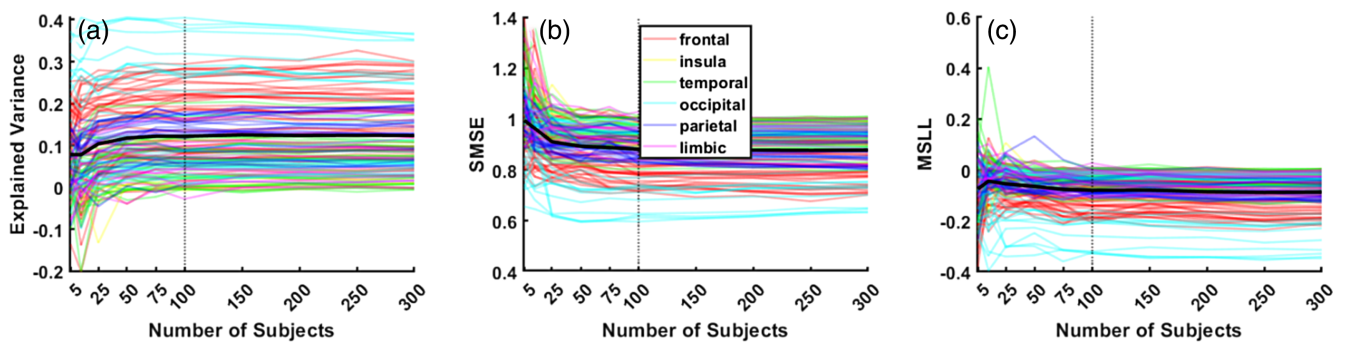
**FIGURE 2** Comparison of model performance as the number of subjects in adaptation set increases. Colored lines show evaluation metrics per region of interest (ROI), color coded according to cerebral area. The black line illustrates the mean across all ROIs. Model performance reaches a plateau at approximately 100 scans per wave in the adaptation sample (vertical dotted line).

**TABLE 3** Performance metrics of models that used 25 and 100 scans for adaptation to the target cohort. Additionally, performance metrics of models trained on half of the target cohort without the use of transfer learning are reported. These serve as a comparison for optimal performance, since performance metrics are generally lower in cohorts with a narrower age range, and therefore less variance to explain, compared to life-span models. While models trained on only Generation R data perform slightly better, they come at the cost of sacrificing a large percentage of data points to train cohort specific models.

| | Explained variance mean (SD) | SMSE mean (SD) | MSLL mean (SD) |
|---|---|---|---|
| Transfer learning with 25 scans | 0.10 (0.10) | 0.91 (0.10) | −0.05 (0.08) |
| Transfer learning with 100 scans | 0.12 (0.10) | 0.88 (0.10) | −0.08 (0.08) |
| Normative Models trained on Generation R data only | 0.17 (0.12) | 0.83 (0.12) | −0.12 (0.10) |

same sites, we found a slight bias for lower deviation scores (z-scores) when running the adaptation with only scans from wave 2, and higher deviation scores when running the adaptation with only scans from wave 3, in particular for frontal ROIs. The effect of the different adaptation settings on all ROIs is shown in Supplementary Figure 3 and is explicitly illustrated in an example ROI in Figure 3. Panel (a) shows that the model is more successful in reparametrizing the raw data to centiles when each time point of measurement is handled as a separate site-effect. Possible sources for such effects might stem from changes in scanner software, changes in image quality with age (i.e., motion artifacts), or sample variability. In our target cohort, scanner software was upgraded after the first 370 scans of wave 2 but was otherwise identical in waves 2 and 3. However, age-related improvements in images quality are frequently reported in the literature and quality assurance, measured as topological defects in the surface reconstruction for FreeSurfer processed MRI data (https://github.com/Deep-MI/qatools-python), does show improvements in image quality with age across the three waves ($Mean_{wave\ 1} = 229.06$, $SD_{wave\ 1} = 98.56$; $Mean_{wave\ 2} = 213.89$, $SD_{wave\ 2} = 67.15$; $Mean_{wave\ 3} = 166.97$, $SD_{wave\ 3} = 48.63$).

### 3.1.3 | Adaptation validation

After choosing for an adaptation setting treating the three measurement waves as different batch effects, we validated the success of the adaptation of the reference model to the target cohort by examining the differences between raw CT values and corrected deviations

(z-scores) after transfer of the subjects which were withheld from the adaptation sets (see 2.2.2 Validation of adaptation). Scans of participants with repeated measurements at all imaging waves (a random sample of 10 participants is depicted by colored lines) show a decline over time in raw CT (Figure 4a). As expected, thinning of the cortex can be observed with age, however, the raw CT values are confounded by noise stemming from site-effects of the different measurement waves. In the resulting z-scores of the withheld subjects, these site-effects are removed as demonstrated by stable deviations from the normative model within a participant (Figure 4b). The same holds true for the withheld subjects from measurement wave 1 that fall in the overlapping age range (8.6–10.7 years of age) of waves 1 (scanner 1) and 2 (scanner 2) (Figure 4c-e). While raw CT values in the overlapping age range vary vastly between the two measurement waves ($t(2874) = 13.4$, $p < .001$), with a tendency of higher values in measurement wave 1 compared to wave 2 (Figure 4c), this difference is slightly reduced when correcting for sex (Figure 4d) ($t(2874) = 11.4$, $p < .001$) and practically absent in the sex- and additionally site-effect corrected z-scores (Figure 4e) ($t(2874) = 1.0$, $p = .324$). Therefore, we can meaningfully compare individuals on the basis of z-scores, bearing in mind that the z-scores are defined with respect to a lifespan based normative model.

### 3.2 | Relating site-effect corrected z-scores to gestational age

To illustrate the usefulness of the resulting models, we compared extreme deviations, acquired at the level of individuals, between
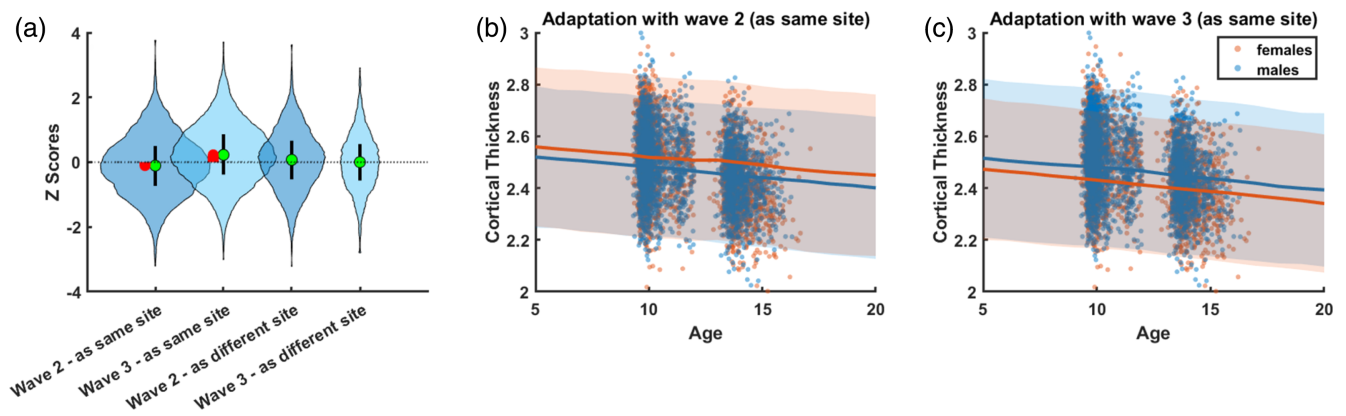
**FIGURE 3** Effects of different recalibration configurations on the target cohort illustrated in an example ROI (inferior frontal sulcus). Panel (a) shows z-score distributions when measurement waves 2 and 3 of the Generation R target cohort are treated as the same (same scanner) or different sites in a frontal example ROI. Median and interquartile range are represented by green dots and black bars, respectively. For each measurement wave, we would expect the median z-score to be around 0. However, this is not the case if measurement waves 2 and 3 are treated as the same site. The difference from 0 is indicated by red bars. By examining the cortical thickness (CT) trajectories in panel (b) and (c), we see that this might be due to a misestimation of mean and variance in females when both waves are treated as the same site.



**FIGURE 4** Validation of transferring the reference normative model to the target cohort using two groups of subjects that were withheld from the adaptation set: (1) subjects with repeated measurements at all three imaging waves (random sample of ten participants depicted by colored lines, panels a and b); (2) subjects from imaging waves 1 and 2 that fall in the overlapping age range of both scanners [8.6–10.7] (depicted by darker shaded red and blue dots and lines, panels a–e). Panels (a) and (c) show raw cortical thickness values. Panels (b), (d), and (e) show sex-effect (panel d) or sex- and site-effect corrected z-scores of the same participants (panels b and e). For consistency, the same region of interest (ROI) (inferior frontal sulcus) as in the previous figures is illustrated.

children born preterm and children born at term in the target cohort. Percentages of individuals with an extreme z-score (larger/smaller than 2) per ROI are shown in Figure 5. In the children born at term, we find approximately 2.5% of children with extreme negative and extreme positive z-scores respectively across ROIs. Exceptions are primarily smaller ROIs (sulcus intermedius primus [left and right], posterior ramus of the lateral sulcus [left], anterior transverse collateral sulcus [right], orbital sulcus [right]) where areas with thicker cortices than expected can be observed. Importantly, extreme deviations are much more prevalent with children born preterm with the most pronounced extreme positive deviations (thicker cortex than expected) found in the left pericallosal sulcus and lateral aspect of the
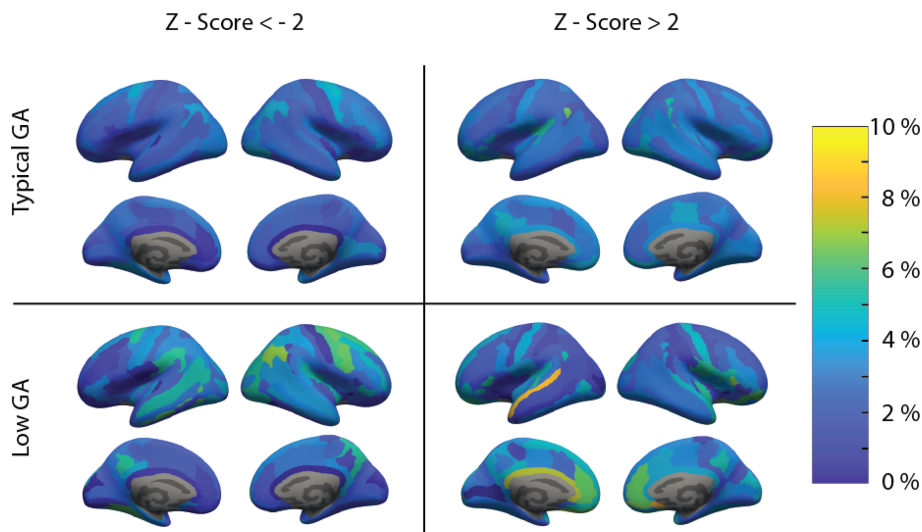
**Z - Score < - 2**     **Z - Score > 2**

Typical GA

Low GA

**FIGURE 5** Differences in site-effect corrected z-scores between children born preterm (low gestational age [GA]) and children born term (typical GA). On the left side, extreme negative deviations (cortex thinner than expected) are illustrated. On the right, extreme positive deviations (cortex thicker than expected) are shown.

superior temporal gyrus, as well as in the anterior part of cingulate gyrus and sulcus of both hemispheres. The most striking extreme negative deviations (thinner cortex than expected) can be seen on the left hemisphere in the superior and inferior temporal sulcus, lingual sulcus, superior part of the precentral sulcus, supramarginal gyrus, and on the right hemisphere in the superior and inferior part of the precentral sulcus, superior frontal sulcus, angular gyrus, precentral gyrus, and the precuneus.

These regions are consistent with previous findings on CT differences in adolescents born preterm. Pronounced cortical thinning has been found persistently in areas surrounding the central sulcus and temporal lobes (Martinussen et al., 2005; Nagy et al., 2011; Zubiaurre-Elorza et al., 2012) as well as thicker cortices in frontal regions surrounding the anterior cingulate cortex (Bjuland et al., 2013). The current approach has been shown to capture structural deviations better than case–control studies as they are more sensitive to individual heterogeneity (Remiszewski et al., 2022). It also offers improved insights in longitudinal cohorts, as these deviation scores are not cofounded by site-effects.

## 4 | DISCUSSION

In this study, we used information from normative models that were initially trained on a large number of samples, scanned over 77 sites, as prior knowledge for the parameters of the CT distributions when adapting these models to the two scanners of the longitudinal Generation R study.

We report three main findings: first, transfer learning is successful and allows for meaningful comparisons between individuals from different scanners, and sexes, as previously reported (Kia et al., 2022). Second, we quantified the number of samples in the transfer set needed to obtain good performance metrics on the test set and show that relatively few samples are sufficient for good performance (approximately $n = 25$). This provides the added benefit of improving the statistical power of statistical analyses on the resulting larger test

set. While we used 100 samples per measurement wave in the adaptation site, slightly smaller adaptation samples decreased the evaluation metrics only marginally. Third, we show that the deviations from these normative models are meaningful in that they are altered in a highly individualized manner in individuals born preterm.

Our results support the finding that normative models capture the general trend of decreasing CT with age, as reported in previous studies (Bethlehem et al., 2022; Frangou et al., 2022; Rutherford, Fraza, et al., 2022; Tamnes et al., 2017; Thambisetty et al., 2010). Interestingly, we found that the model performed better when each measurement wave of the transfer cohort was treated as a separate site-effect, even though two of three waves were acquired on the same scanner. This could be due to sample variability, a misestimation of parameters in the female cohort, or it might be linked to the fact that scan quality tends to improve with age. For future studies, it may be useful to treat distinct measurement intervals as separate batch-effects, resulting in a factorial design of sex × scanners × waves, even if the scanner setup has not changed, to produce more precise models. Our recommendations might differ for longer timescales, such as nonlinear or non-Gaussian lifespan trajectories, which usually requires more data (de Boer et al., 2022). They might also differ in more fine-grained parcellations, as we found ROI area to correlate positively with model performance. However, the methods we introduce can be used to determine the optimal number of subjects for such cases.

The successful validation of the use of transfer learning with normative models opens the door for further investigations exploring the relationship between deviation scores and various phenotypes. Individual-level deviations, as obtained through normative models, have been shown to provide stronger effects than typical case–control studies using uncorrected raw measurements (Rutherford, Fraza, et al., 2022) and are therefore particularly suitable for exploring and investigating individual differences within and across datasets. For longitudinal cohorts, an approach to quantitatively assess the significance of within-subject changes over time (as, e.g., illustrated by participants with repeated measurements in Figure 4b) was recently

introduced (Rehak Buckova et al., 2023). The used federated learning framework makes it possible to use the models presented in this work as informed priors (models are available online via PCNportal [https://pcnportal.dccn.nl/]) to investigate CT in smaller and/or clinical cohorts in such a way.

With the current study, we also illustrate how normative models can be used to study clinical phenotypes, by investigating the relation between extreme deviations scores and the gestational age at birth, that is between children born at-term and preterm. While children born at-term show an expected distribution of approximately 2.5% of z-scores higher than 2 or lower than −2 respectively, children born preterm are more likely to have extreme deviations in specific ROIs, which are consistent with previous literature showing pronounced differences in particular in frontal and temporal cortices. While we show a comparison between groups, the current approach does not require clustering of individuals into groups but instead can be used to make inferences about heterogeneity within clinical groups as well as about deviations on an individual level.

## 4.1 | Limitations and future directions

We demonstrate that evaluation metrics level off after 100 scans in the adaptation set, with as few as 25 scans leading to effective transfer of knowledge. Although this marks a considerable improvement in terms of scanning costs compared to developing models for each new research objective, it might prevent small cohorts from utilizing the current approach, given that neuroimaging studies typically include a median sample size of 25 participants (Marek et al., 2022). It is important to note, however, that out-of-sample testing is a critical step in assessing the generalizability of the model. Future work could investigate the possibility of utilizing the deviation scores from all available data point. This could be achieved by generating numerous randomly sampled adaptation sets, running the adaption to the target cohort for each set, and subsequently analyzing their convergence. Furthermore, our work estimates normative models on a single ROI, thereby neglecting any spatial interdependencies between brain regions.

## 5 | CONCLUSION

Using longitudinal CT data from the Generation R study of children aged 6–17 years, we present an application of transfer learning of large-scale normative models which produce good performance metrics even with an adaptation set containing as few as 25 scans. The resulting deviation scores allow for meaningful comparisons across scanner site and sex. Using these obtained deviation scores, we were able to show localized differences in CT between children born preterm and children born at-term.

## CONFLICT OF INTEREST STATEMENT

The authors declare no conflicts of interest.

## DATA AVAILABILITY STATEMENT

MRI datasets from the reference model are mostly publicly available: CAMCAN, PNC, HCP-Aging, HCP-Dev, HCP-EP, OASIS, OPN, IXI, NKI-RS, UKBB, ABCD, and CMI-HBN. Data from the Generation R cohort (target cohort) are not publicly available due to legal and ethical restrictions; however, access can be requested via the Generation R administration (secretariaat.genr@erasmusmc.nl). The reference model is available on the PCNportal (https://pcnportal.dccn.nl/) and code to generate normative models and transfer knowledge from existing models to new sites is freely available via the PCNtoolkit (https://github.com/amarquand/PCNtoolkit).

## ORCID

C. Gaiser https://orcid.org/0000-0003-3658-1946

## REFERENCES

Antoniades, M., Haas, S. S., Modabbernia, A., Bykowsky, O., Frangou, S., Borgwardt, S., & Schmidt, A. (2021). Personalized estimates of brain structural variability in individuals with early psychosis. *Schizophrenia Bulletin*, 47(4), 1029–1038. https://doi.org/10.1093/schbul/sbab005

Bayer, J. M. M., Dinga, R., Kia, S. M., Kottaram, A. R., Wolfers, T., Lv, J., Zalesky, A., Schmaal, L., & Marquand, A. (2021). Accommodating site variation in neuroimaging data using normative and hierarchical Bayesian models. *BioRxiv*, 2021.02.09.430363. Retrieved from https://www.biorxiv.org/content/10.1101/2021.02.09.430363v2.abstract

Bethlehem, R. A. I., Seidlitz, J., White, S. R., Vogel, J. W., Anderson, K. M., Adamson, C., Adler, S., Alexopoulos, G. S., Anagnostou, E., Areces-Gonzalez, A., Astle, D. E., Auyeung, B., Ayub, M., Bae, J., Ball, G., Baron-Cohen, S., Beare, R., Bedford, S. A., Benegal, V., ... Alexander-Bloch, A. F. (2022). Brain charts for the human lifespan. *Nature*, 604(7906), 525–533. https://doi.org/10.1038/s41586-022-04554-y

Bjuland, K. J., Løhaugen, G. C. C., Martinussen, M., & Skranes, J. (2013). Cortical thickness and cognition in very-low-birth-weight late teenagers. *Early Human Development*, 89(6), 371–380. https://doi.org/10.1016/j.earlhumdev.2012.12.003

Clarkson, M. J., Cardoso, M. J., Ridgway, G. R., Modat, M., Leung, K. K., Rohrer, J. D., Fox, N. C., & Ourselin, S. (2011). A comparison of voxel

and surface based cortical thickness estimation methods. *NeuroImage*, 57(3), 856–865. https://doi.org/10.1016/j.neuroimage.2011.05.053

Cole, T. J. (2012). The development of growth references and growth charts. *Annals of Human Biology*, 39(5), 382–394. https://doi.org/10.3109/03014460.2012.694475

Dale, A. M., Fischl, B., & Sereno, M. I. (1999). Cortical surface-based analysis. II. Inflation, flattening, and a surface-based coordinate system. *NeuroImage*, 9(2), 179–194. Retrieved from http://linkinghub.elsevier.com/retrieve/pii/S1053811998903950

de Boer, A. A. A., Kia, S. M., Rutherford, S., Zabihi, M., Fraza, C., Barkema, P., Westlye, L. T., Andreassen, O. A., Hinne, M., Beckmann, C. F., & Marquand, A. (2022). Non-gaussian normative modelling with hierarchical Bayesian regression. BioRxiv, 2022.10.05.510988. https://doi.org/10.1101/2022.10.05.510988

DeLisi, L. E. (2008). The concept of progressive brain change in schizophrenia: Implications for understanding schizophrenia. *Schizophrenia Bulletin*, 34(2), 312–321. https://doi.org/10.1093/schbul/sbm164

Destrieux, C., Fischl, B., Dale, A., & Halgren, E. (2010). Automatic parcellation of human cortical gyri and sulci using standard anatomical nomenclature. *NeuroImage*, 53(1), 1–15. https://doi.org/10.1016/j.neuroimage.2010.06.010

Fischl, B., & Dale, A. M. (2000). Measuring the thickness of the human cerebral cortex from magnetic resonance images. *Proceedings of the National Academy of Sciences of the United States of America*, 97(20), 11050–11055. https://doi.org/10.1073/pnas.200033797

Fischl, B., Salat, D. H., Busa, E., Albert, M., Dieterich, M., Haselgrove, C., van der Kouwe, A., Killiany, R., Kennedy, D., Klaveness, S., Montillo, A., Makris, N., Rosen, B., & Dale, A. M. (2002). Whole brain segmentation. *Neuron*, 33(3), 341–355. https://doi.org/10.1016/s0896-6273(02)00569-x

Fischl, B., Sereno, M. I., & Dale, A. M. (1999). Cortical surface-based analysis. I. Segmentation and surface reconstruction. *NeuroImage*, 9(2), 195–207. Retrieved from http://www.ncbi.nlm.nih.gov/pubmed/9931269

Frangou, S., Modabbernia, A., Williams, S. C. R., Papachristou, E., Doucet, G. E., Agartz, I., Aghajani, M., Akudjedu, T. N., Albajes-Eizagirre, A., Alnæs, D., Alpert, K. I., Andersson, M., Andreasen, N. C., Andreassen, O. A., Asherson, P., Banaschewski, T., Bargallo, N., Baumeister, S., Baur-Streubel, R., … Dima, D. (2022). Cortical thickness across the lifespan: Data from 17,075 healthy individuals aged 3–90 years. *Human Brain Mapping*, 43(1), 431–451. https://doi.org/10.1002/hbm.25364

Fraza, C. J., Dinga, R., Beckmann, C. F., & Marquand, A. F. (2021). Warped Bayesian linear regression for normative modelling of big data. *NeuroImage*, 245(May), 118715. https://doi.org/10.1016/j.neuroimage.2021.118715

Fuhrmann, D., Madsen, K. S., Johansen, L. B., Baaré, W. F. C., & Kievit, R. A. (2022). The midpoint of cortical thinning between late childhood and early adulthood differs across individuals and regions: Evidence from longitudinal modelling in a 12-wave sample. BioRxiv, 261 (March), 1–23. https://doi.org/10.1016/j.neuroimage.2022.119507

Hutton, C., Draganski, B., Ashburner, J., & Weiskopf, N. (2009). A comparison between voxel-based cortical thickness and voxel-based morphometry in normal aging. *NeuroImage*, 48(2), 371–380. https://doi.org/10.1016/j.neuroimage.2009.06.043

Insel, T. R. (2014). Mental disorders in childhood. *Jama*, 311(17), 1727. https://doi.org/10.1001/jama.2014.1193

Jaddoe, V. W. V., Mackenbach, J. P., Moll, H. A., Steegers, E. A. P., Tiemeier, H., Verhulst, F. C., Witteman, J. C. M., & Hofman, A. (2006). The Generation R Study: Design and cohort profile. *European Journal of Epidemiology*, 21(6), 475–484. https://doi.org/10.1007/s10654-006-9022-0

Kia, S. M., Huijsdens, H., Dinga, R., Wolfers, T., Mennes, M., Andreassen, O. A., Westlye, L. T., Beckmann, C. F., & Marquand, A. F. (2020). Hierarchical Bayesian regression for multi-site normative modeling of neuroimaging data. 1–12. Retrieved from http://arxiv.org/abs/2005.12055

Kia, S. M., Huijsdens, H., Rutherford, S., de Boer, A., Dinga, R., Wolfers, T., Berthet, P., Mennes, M., Andreassen, O. A., Westlye, L. T., Beckmann, C. F., & Marquand, A. F. (2022). Closing the life-cycle of normative modeling using federated hierarchical Bayesian regression. *PLoS One*, 17(12), e0278776. https://doi.org/10.1371/journal.pone.0278776

Kooijman, M. N., Kruithof, C. J., van Duijn, C. M., Duijts, L., Franco, O. H., van IJzendoorn, M. H., de Jongste, J. C., Klaver, C. C. W., van der Lugt, A., Mackenbach, J. P., Moll, H. A., Peeters, R. P., Raat, H., Rings, E. H. H. M., Rivadeneira, F., van der Schroeff, M. P., Steegers, E. A. P., Tiemeier, H., Uitterlinden, A. G., … Jaddoe, V. W. V. (2016). The Generation R Study: Design and cohort update 2017. *European Journal of Epidemiology*, 31(12), 1243–1264. https://doi.org/10.1007/s10654-016-0224-9

Marek, S., Tervo-Clemmens, B., Calabro, F. J., Montez, D. F., Kay, B. P., Hatoum, A. S., Donohue, M. R., Foran, W., Miller, R. L., Hendrickson, T. J., Malone, S. M., Kandala, S., Feczko, E., Miranda-Dominguez, O., Graham, A. M., Earl, E. A., Perrone, A. J., Cordova, M., Doyle, O., … Dosenbach, N. U. F. (2022). Reproducible brain-wide association studies require thousands of individuals. *Nature*, 603, 654–660. https://doi.org/10.1038/s41586-022-04492-9

Marquand, A. F., Kia, S. M., Zabihi, M., Wolfers, T., Buitelaar, J. K., & Beckmann, C. F. (2019). Conceptualizing mental disorders as deviations from normative functioning. *Molecular Psychiatry*, 24(10), 1415–1424. https://doi.org/10.1038/s41380-019-0441-1

Marquand, A. F., Rezek, I., Buitelaar, J., & Beckmann, C. F. (2016). Understanding heterogeneity in clinical cohorts using normative models: Beyond case-control studies. *Biological Psychiatry*, 80(7), 552–561. https://doi.org/10.1016/j.biopsych.2015.12.023

Marquand, A. F., Wolfers, T., Mennes, M., Buitelaar, J., & Beckmann, C. F. (2016). Beyond lumping and splitting: A review of computational approaches for stratifying psychiatric disorders. *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging*, 1(5), 433–447. https://doi.org/10.1016/j.bpsc.2016.04.002

Martinussen, M., Fischl, B., Larsson, H. B., Skranes, J., Kulseng, S., Vangberg, T. R., Vik, T., Brubakk, A.-M., Haraldseth, O., & Dale, A. M. (2005). Cerebral cortex thickness in 15-year-old adolescents with low birth weight measured by an automated MRI-based method. *Brain*, 128(11), 2588–2596. https://doi.org/10.1093/brain/awh610

Mills, K. L., Siegmund, K. D., Tamnes, C. K., Ferschmann, L., Wierenga, L. M., Bos, M. G. N., Luna, B., Li, C., & Herting, M. M. (2021). Inter-individual variability in structural brain development from late childhood to young adulthood. *NeuroImage*, 242(August), 118450. https://doi.org/10.1016/j.neuroimage.2021.118450

Nagy, Z., Lagercrantz, H., & Hutton, C. (2011). Effects of preterm birth on cortical thickness measured in adolescence. *Cerebral Cortex*, 21(2), 300–306. https://doi.org/10.1093/cercor/bhq095

Pan, S. J., & Yang, Q. (2010). A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(10), 1345–1359. https://doi.org/10.1109/TKDE.2009.191

Pereira, J. B., Ibarretxe-Bilbao, N., Marti, M. J., Compta, Y., Junqué, C., Bargallo, N., & Tolosa, E. (2012). Assessment of cortical degeneration in patients with Parkinson's disease by voxel-based morphometry, cortical folding, and cortical thickness. *Human Brain Mapping*, 33(11), 2521–2534. https://doi.org/10.1002/hbm.21378

Rehak Buckova, B., Fraza, C., Rehak, R., Kolenic, M., Beckmann, C., Spaniel, F., Marquand, A., & Hlinka, J. (2023). Using normative models pre-trained on cross-sectional data to evaluate longitudinal changes in neuroimaging data. bioRxiv, 2023.06.2009.544217.

Remiszewski, N., Bryant, J. E., Rutherford, S. E., Marquand, A. F., Nelson, E., Askar, I., Lahti, A. C., & Kraguljac, N. V. (2022). Contrasting case-control and normative reference approaches to capture clinically

relevant structural brain abnormalities in patients with first-episode psychosis who are antipsychotic naive. *JAMA Psychiatry*, *79*(11), 1133–1138. https://doi.org/10.1001/jamapsychiatry.2022.3010

Rogers, J. C., & de Brito, S. A. (2016). Cortical and subcortical gray matter volume in youths with conduct problems a meta-analysis. *JAMA Psychiatry*, *73*(1), 64–72. https://doi.org/10.1001/jamapsychiatry.2015.2423

Rutherford, S., Fraza, C., Dinga, R., Kia, S. M., Wolfers, T., Zabihi, M., Berthet, P., Worker, A., Verdi, S., Andrews, D., Han, L. K. M., Bayer, J. M. M., Dazzan, P., McGuire, P., Mocking, R. T., Schene, A., Sripada, C., Tso, I. F., Duval, E. R., … Marquand, A. F. (2022). Charting brain growth and aging at high spatial precision. *eLife*, *11*, 1–15. https://doi.org/10.7554/eLife.72904

Rutherford, S., Kia, S. M., Wolfers, T., Fraza, C., Zabihi, M., Dinga, R., Berthet, P., Worker, A., Verdi, S., Ruhe, H. G., Beckmann, C. F., & Marquand, A. F. (2022). The normative modeling framework for computational psychiatry. *Nature Protocols*, *17*(July), 1711–1734. https://doi.org/10.1038/s41596-022-00696-5

Solmi, M., Radua, J., Olivola, M., Croce, E., Soardo, L., Salazar de Pablo, G., il Shin, J., Kirkbride, J. B., Jones, P., Kim, J. H., Kim, J. Y., Carvalho, A. F., Seeman, M. v., Correll, C. U., & Fusar-Poli, P. (2022). Age at onset of mental disorders worldwide: Large-scale meta-analysis of 192 epidemiological studies. *Molecular Psychiatry*, *27*(1), 281–295. https://doi.org/10.1038/s41380-021-01161-7

Tamnes, C. K., Herting, M. M., Goddings, A. L., Meuwese, R., Blakemore, S. J., Dahl, R. E., Güroğlu, B., Raznahan, A., Sowell, E. R., Crone, E. A., & Mills, K. L. (2017). Development of the cerebral cortex across adolescence: A multisample study of inter-related longitudinal changes in cortical volume, surface area, and thickness. *Journal of Neuroscience*, *37*(12), 3402–3412. https://doi.org/10.1523/JNEUROSCI.3302-16.2017

Thambisetty, M., Wan, J., Carass, A., An, Y., Prince, J. L., & Resnick, S. M. (2010). Longitudinal changes in cortical thickness associated with normal aging. *NeuroImage*, *52*(4), 1215–1223. https://doi.org/10.1016/j.neuroimage.2010.04.258

Volpe, J. J. (2009). Brain injury in premature infants: A complex amalgam of destructive and developmental disturbances. *The Lancet Neurology*, *8*(1), 110–124. https://doi.org/10.1016/S1474-4422(08)70294-1

White, T., Muetzel, R. L., el Marroun, H., Blanken, L. M. E., Jansen, P., Bolhuis, K., Kocevska, D., Mous, S. E., Mulder, R., Jaddoe, V. W. V., van der Lugt, A., Verhulst, F. C., & Tiemeier, H. (2018). Paediatric population neuroimaging and the Generation R Study: The second wave. *European Journal of Epidemiology*, *33*(1), 99–125. https://doi.org/10.1007/s10654-017-0319-y

Whittle, S., Vijayakumar, N., Simmons, J. G., & Allen, N. B. (2020). Internalizing and externalizing symptoms are associated with different trajectories of cortical development during late childhood. *Journal of the American Academy of Child and Adolescent Psychiatry*, *59*(1), 177–185. https://doi.org/10.1016/j.jaac.2019.04.006

Zabihi, M., Oldehinkel, M., Wolfers, T., Frouin, V., Goyard, D., Loth, E., Charman, T., Tillmann, J., Banaschewski, T., Dumas, G., Holt, R., Baron-Cohen, S., Durston, S., Bölte, S., Murphy, D., Ecker, C., Buitelaar, J. K., Beckmann, C. F., & Marquand, A. F. (2019). Dissecting the heterogeneous cortical anatomy of autism spectrum disorder using normative models. *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging*, *4*(6), 567–578. https://doi.org/10.1016/j.bpsc.2018.11.013

Zhao, Y., Zhang, Q., Shah, C., Li, Q., Sweeney, J. A., Li, F., & Gong, Q. (2022). Cortical thickness abnormalities at different stages of the illness course in schizophrenia. *JAMA Psychiatry*, *79*(6), 560. https://doi.org/10.1001/jamapsychiatry.2022.0799

Zubiaurre-Elorza, L., Soria-Pastor, S., Junque, C., Sala-Llonch, R., Segarra, D., Bargallo, N., & Macaya, A. (2012). Cortical thickness and behavior abnormalities in children born preterm. *PLoS One*, *7*(7), e42148. https://doi.org/10.1371/journal.pone.0042148

## SUPPORTING INFORMATION

Additional supporting information can be found online in the Supporting Information section at the end of this article.