

1 **Supplementary material**

2 **0.1 Missing values and sample sizes per model**

3 Despite the aim of the study to run all seven models in this study on the same data set, we were faced with the challenge
 4 that differing requirements for each model with respect to missing values made an adjustment of sample sizes per model
 5 necessary. Combat [3, 2] accepts missing values, and could be thus run on the full data set. In contrast, ComBat Gam
 6 [7] does not. Thus, for ComBat Gam all subjects with a missing value in any of the 35 regions had to be excluded,
 7 which lead to a sample size reduction from 391 to 370 individuals for the training set, and from 168 to 156 individuals
 8 for the healthy test set. The normative modeling process is performed region wise and independently, thus only the
 9 subjects that contained missing subjects for that particular region were deleted.

10 **0.2 Model convergence, effective sample size and \hat{R}**

11 For the present project, each model run entailed a Monte-Carlo sampling process of 4000 iterations in Stan, of which
 12 2000 were disregarded as warm up. Stan allows for the computation of a number of diagnostics on the quality of the
 13 Markov Chain Monte Carlo (MCMC) sampling process, which are reported in the following.

14 **0.3 Model convergence**

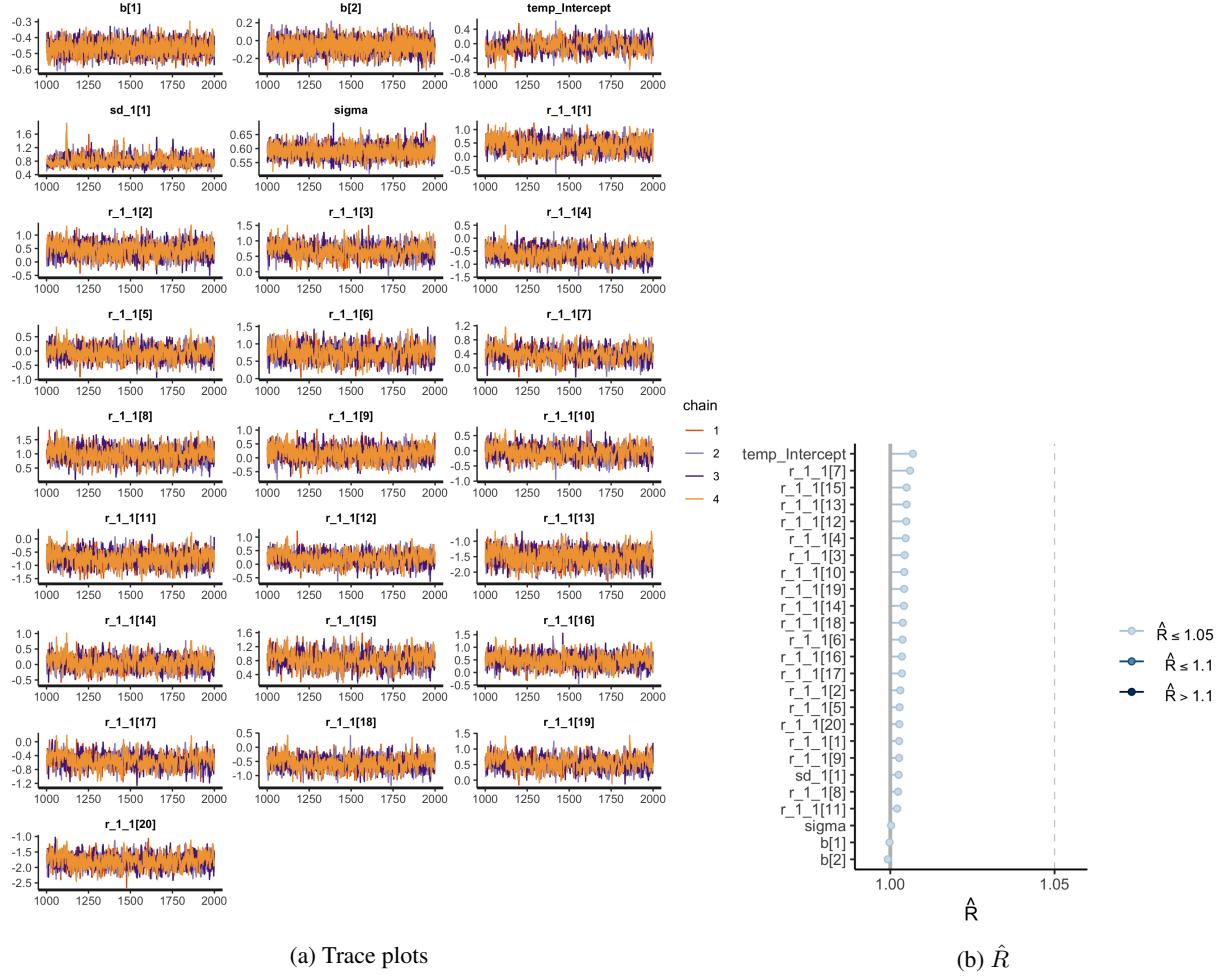
15 Markov chains are defined to only generate samples from the target distribution after the distribution has converged to
 16 an equilibrium, thus, when the distribution is considered to be the target density. In theory, this equilibrium can only
 17 asymptotically be reached, as the number of draws is theoretically infinite. In practice, the number of draws has to be
 18 a-priori set to a finite amount. As a consequence, the actual and convergence to the target density has to be monitored
 19 [1, 4, 8]. One way to monitor the convergence of a chain to equilibrium is to compare the chain to other randomly
 20 initialized chains. This can either be done via visual inspection, or using the scale reduction statistic, \hat{R} [5]. For visual
 21 inspection, the trace plots over 4 chains for all parameters can be found in Figs. 1a, 2a , 3a, 4a, 5a, 6a, 7a. The \hat{R} values,
 22 indicating the convergence of chains, can be found in Figs. 1b, 2b , 3b, 4b, 5b, 6b, 7b for each model, respectively.
 23 All \hat{R} values are <1.05, which provides good evidence that all chains have reached convergence and can therefore be
 24 considered to provide unbiased samples from the target density.

25 **0.4 Statistical comparison of measures of model performance**

26 All comparisons regarding measures of model performance were performed using two-way ANOVAs including the
 27 factors model (HBLM, HBGPM, Combat Gam, ComBat, ComBat without covariates, residuals, raw data) and set
 28 (train, test). Post hoc tests were performed using Tukey tests and corrected for multiple comparisons. Parametric tests
 29 such as ANOVA were deliberately chosen over their non-parametric equivalents, since deviations from gaussianity
 30 were negligible in the present data set and in the authors' opinion, the substantial loss of power with the choice of

31 non-parametric tests does not scale with the potential threat of violated modeling assumptions such as homoscedasticity
 32 and gaussianity.

Post-warmup iteration - Hierarchical Bayesian Linear Model



(a) Trace plots

(b) \hat{R}

Figure 1: Hierarchical Bayesian Linear Model

Post-warmup iteration - Hierarchical Bayesian Gaussian Process I

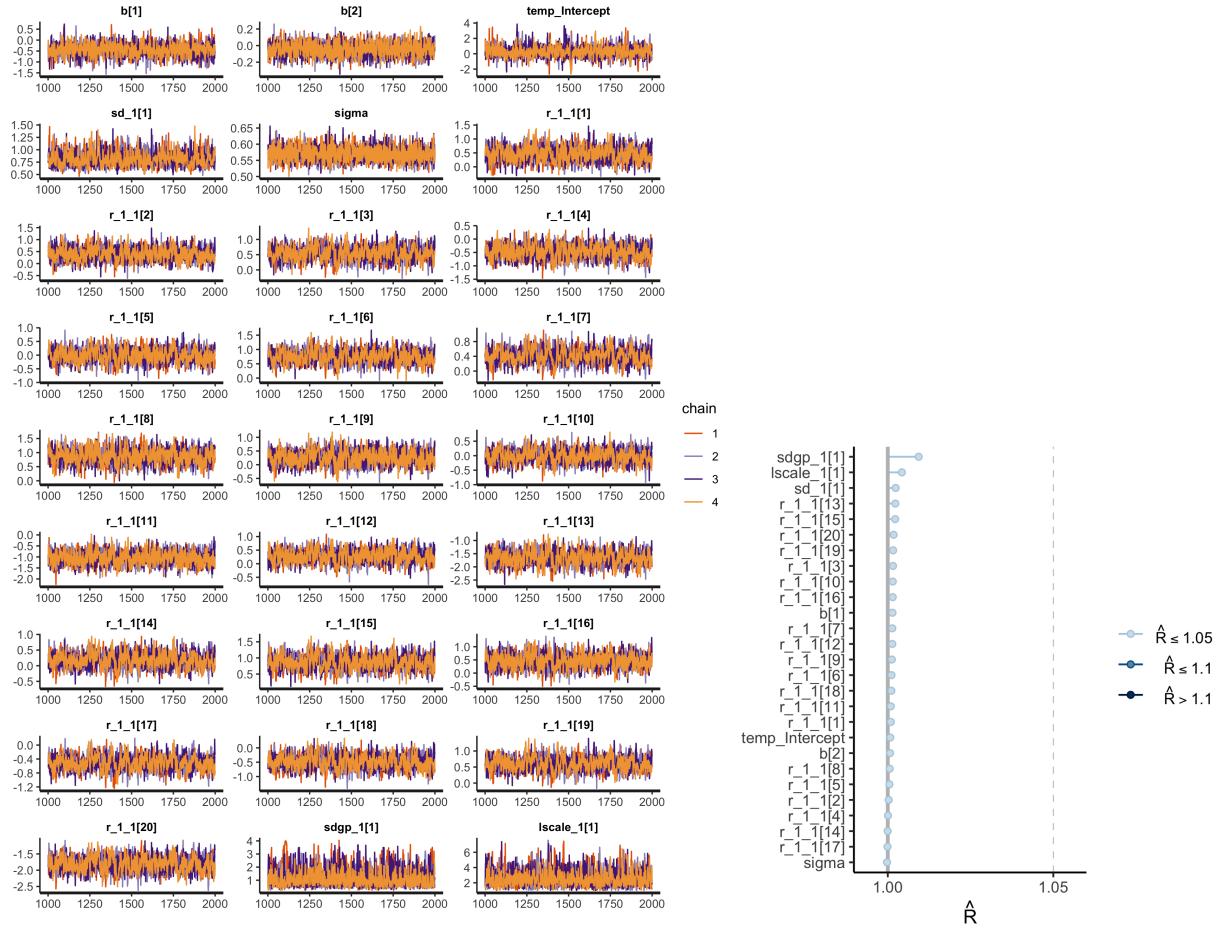


Figure 2: Hierarchical Bayesian Gaussian Process Model

Post-warmup iteration - ComBat Gam Model

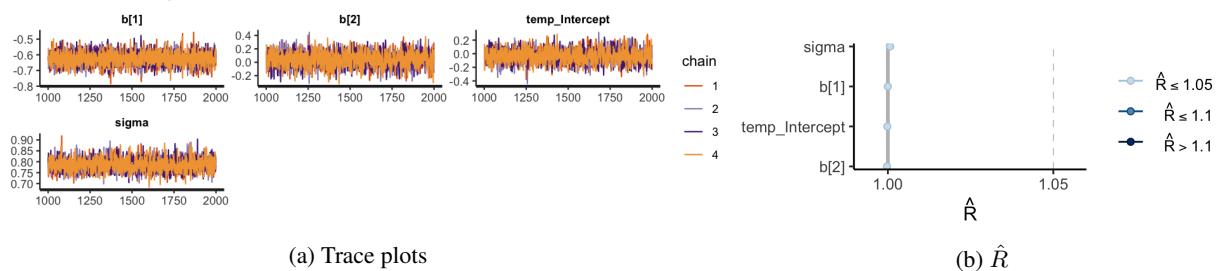


Figure 3: ComBat Gam Model

Post-warmup iteration - ComBat Model

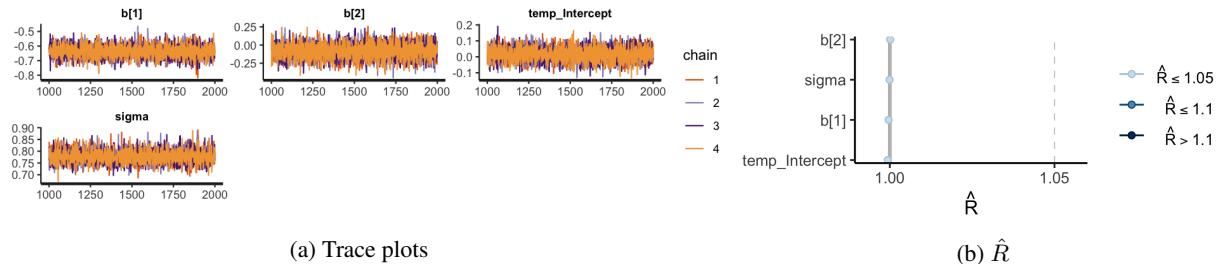


Figure 4: ComBat Model

Post-warmup iteration - ComBat w/o covariates Model

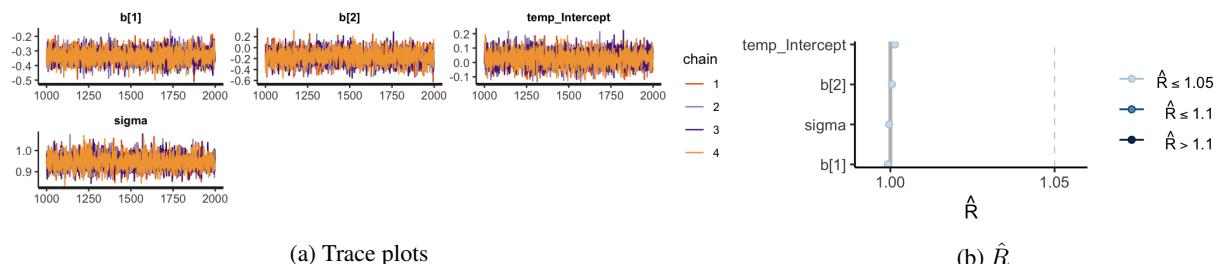


Figure 5: ComBat w/o covariates Model

Post-warmup iteration - Residuals Model

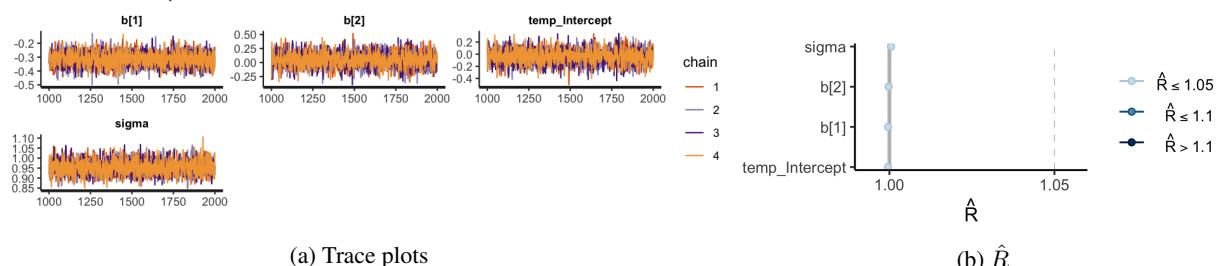


Figure 6: Residuals Model

Post-warmup iteration - Raw Data Model

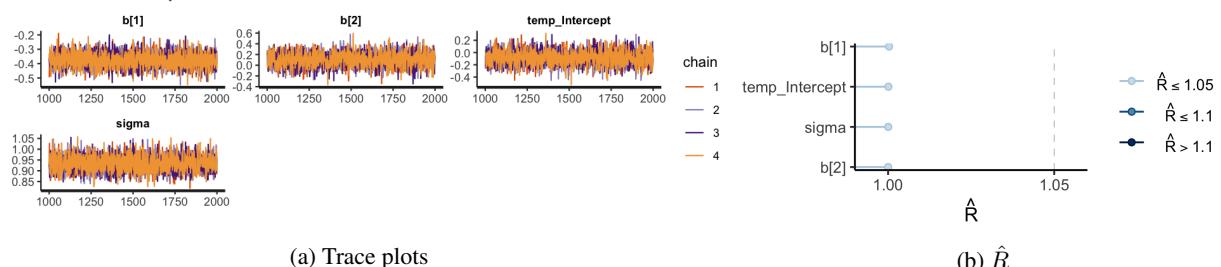


Figure 7: Raw Data Model

33 **0.4.1 Effective sample size**

34 One characteristic of MCMC methods is that samples will be auto- or anti-correlated within a chain, leading to a
 35 reduction of precision in the estimates of posterior quantities [6]. Stan uses the auto correlation ρ_t between samples n
 36 and $n+t$ with lag t to estimate the effective samples size N_{eff} of independent samples in the chain. N_{eff} is considered
 37 to have the same estimation power as N correlated samples and is then used, rather than N , to estimate precision and
 38 error measures [8]. N_{eff} for all models can be found in Figs. 8a, 8b, 8c, 8d, 8e, 8f, 8g.

39 **0.5 Results**

40 **0.6 Correlation between true and predicted value**

41 Correlations between true and predicted values for the HBLM and HBGPM are expressed in terms of the correlation
 42 coefficient ρ , calculated separately for each region. ρ ranged from 0.60 to 0.84 in the training and 0.55 to 0.80 in the
 43 test set. Overall, for just the Bayesian models, correlations were higher in the training set and dropped in the test set
 44 (training set, across all regions: $\bar{\rho}_{HBLM} = 0.73$, $SE = 0.05$; $\bar{\rho}_{HBGPM} = 0.75$, $SE = 0.06$; test set: $\bar{\rho}_{HBLM} = 0.69$, $SE =$
 45 0.06 ; $\bar{\rho}_{HBGPM} = 0.69$, $SE = 0.06$, $F[1, 136] = 18.82$, $p < 0.0001$). Correlations did not differ significantly between the
 46 HBLM and HBGPM. ($F[1, 136] = 2.16$, $p = 0.14$).

47 Comparisons to the other, non-Bayesian models showed that ρ was significantly higher for our models that
 48 included *site* as a predictor for with all other models in which site was harmonized for prior to running the normative
 49 models, both for training and test set (t-test *HBLM* and any other model $p < 0.001$, t-test *HBGPM* and any other model
 50 $p < 0.001$, both for training and test set.) The full distribution of the correlation coefficient ρ for all 35 regions per
 51 model can be found in Fig 3a, main text. In addition, a test comparing the performance of all models for all regions
 52 (Bayesian and non-Bayesian) showed that predictions made from training data were not overall more accurate than
 53 predictions from the test data (main effect *set*, $F[1, 476] = 0.30$, $p = \text{ns}$, interaction *set* \times *model*, $F[1, 476] = 3.50$, $p =$
 54 0.002). Further inspection showed that this might have been caused by the *residuals*, the *ComBat w/o covariates* and
 55 the *ComBat* model, where the test set performed *better* than the training set, canceling out performance benefits of the
 56 training data in the HBLM and HBGPM (see also Fig. 3a, main text).

57 **0.7 Standardized Root Mean Squared Errors**

58 We further evaluated the fit of the models by calculating the Standardized Root Mean Squared Error (SRMSE) between
 59 true values and predicted values per model per region. As expected, the SRMSE was larger for the test set ($M = 0.083$)
 60 and smaller for the training set ($M = 0.080$, [$F(1, 134278) = 59.28$, $p < 0.001$]). For both the training and the test
 61 set, the Bayesian models showed smaller SRMSEs than all other models across all regions ($p < 0.001$; training set:
 62 $SRMSE_{HBGPM} = 0.06$, $SE = 0.005$; $SRMSE_{HBLM} = 0.06$, $SE = 0.005$; $SRMSE_{ComBatGam} = 0.11$, $SE = 0.01$;
 63 $SRMSE_{residuals} = 0.09$, $SE = 0.002$; $SRMSE_{ComBat} = 0.08$, $SE = 0.007$, $SRMSE_{ComBat-w/o-covariates} = 0.09$,
 64 $SE = 0.002$; $SRMSE_{rawdata} = 0.08$, $SE = 0.005$; test set: $SRMSE_{HBGPM} = 0.06$, $SE = 0.005$; $SRMSE_{HBLM} =$

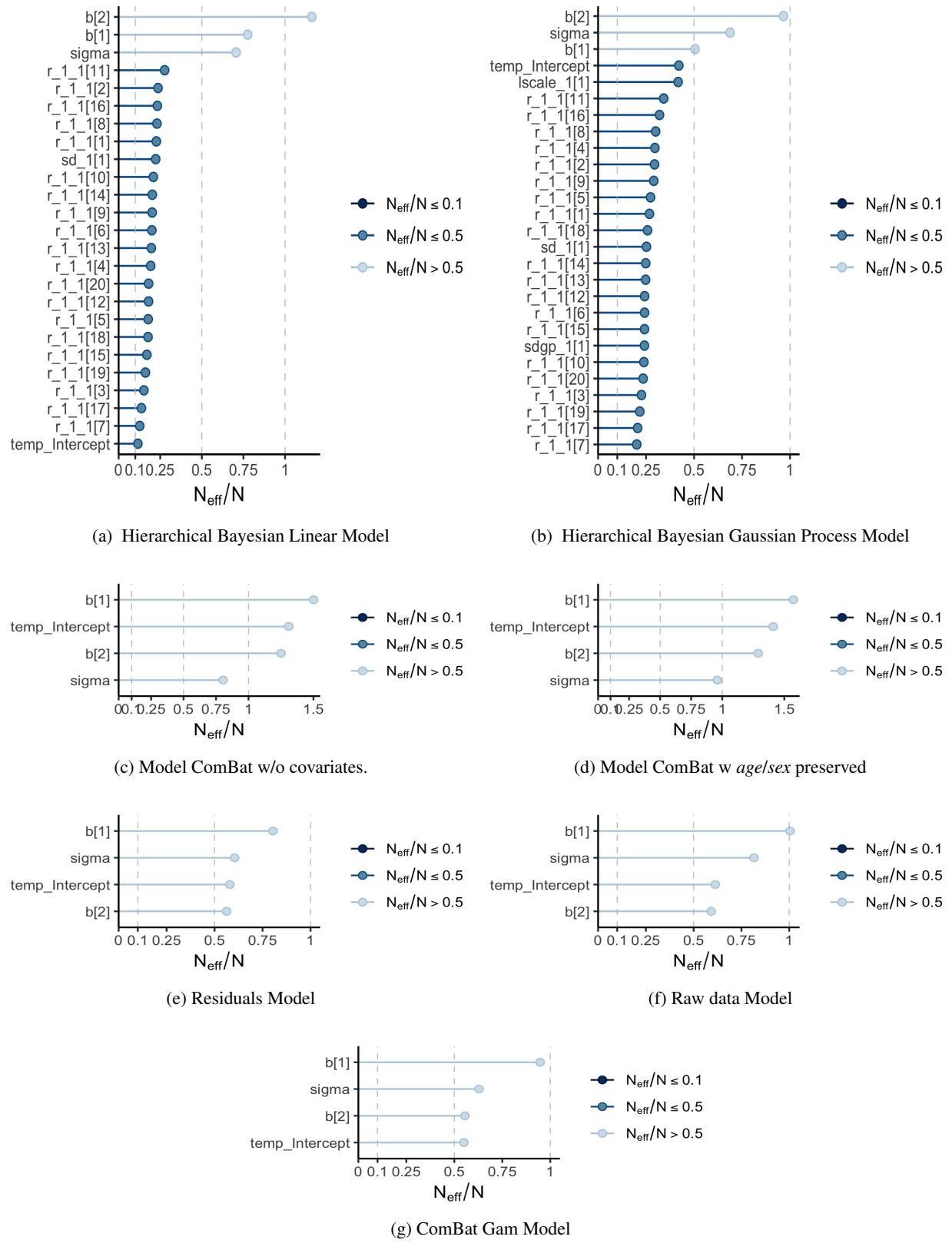
65 0.07, $SE = 0.006$; $SRMSE_{ComBatGam} = 0.12, SE = 0.01$; $SRMSE_{residuals} = 0.09, SE = 0.005$; $SRMSE_{ComBat}$
 66 $= 0.08, SE = 0.009$; $SRMSE_{ComBat-w/o-covariates} = 0.09, SE = 0.005$; $SRMSE_{rawdata} = 0.085, SE = 0.007$).
 67 Neither in the training nor the test set did the Bayesian models differ from each other (training set: contrast *HBLMR*
 68 - *HBLM*, $t = 2.33$, $p = ns.$; test set: contrast *HBLMR* - *HBLM*, $t = 1.14, p = ns.$). We also observed that both in the
 69 training and test set, the SRMSE of *ComBat w/o covariates* did not differ from the SRMSE of the *residuals* (training
 70 set: contrast *ComBat w/o covariates* - *residuals*, $t = 0.69, p = ns.$; test set: contrast *ComBat w/o covariates* - *residuals*, t
 71 $= -1.70, p = ns.$. The full distribution of SRMSE for all 35 regions per model can be found in Fig 3b, main text.

72 0.8 Explained variance

73 Analysis of the proportion of variance explained $EV = \frac{\sigma_{y-y}^2}{\sigma_y^2}$ per model per region were in line with the results reflected
 74 in ρ and SRMSE. EV was higher for the *HBLM* and *HBGPM*, with an average of 0.56 (*HBGPM*, range: 0.35-0.70) and
 75 0.53 (*HBLM*, range 0.35 - 0.67) for the training set and 0.50 (*HBGPM*, range 0.31 - 0.63) and 0.48 (*HBLM*, range 0.28 -
 76 0.60) for the test set across all cortical regions. The proportion of explained variance was substantially lower for the
 77 comparison models, with the ComBat and the ComBat Gam model performing best out of the comparison models, with
 78 an average of 0.31 for ComBat model for the training set (range: 0.00 - 0.51) and 0.33 for the test set (range -0.02 -
 79 0.58) across cortical regions, and an average of 0.31 for ComBat Gam for the training set (range: -0.01 - 0.51) and 0.22
 80 for the test set (range 0.03 - 0.46) but showing lower EV than the Bayesian models. Predictions derived from *residuals*
 81 and *ComBat w/o covariates* showed even lower EV, with *residuals* explaining an average of 0.07 for the training set
 82 (range: 0.00 - 0.15) and 0.11 for the test set (range 0.00 - 0.20) across cortical regions, and *ComBat* explaining an
 83 average of 0.09 for the training set (range: 0.00 - 0.17) and 0.12 for the test set (range: 0.00 - 0.25) across cortical
 84 regions. Thus, the *ComBat w/o covariates*, *ComBat* and *residuals* model performed even worse than predictions derived
 85 from *raw data*, which showed an average EV of 0.21 in the training set (range: 0.00 -0.46) and 0.20 in the test set
 86 (range 0.00 - 0.44) across cortical regions. These results include the interesting finding that the test set shows slightly
 87 higher EVs than the training set for all comparison models. An overview over the distribution of explained variance for
 88 training and the test set for all 35 regions for all models can be found in Fig. 3c, main text.

89 0.9 Log likelihood

90 The point-wise log likelihoods (LL) between the true and predicted were calculated for each data point, summed up per
 91 model across regions and averaged by the number of individuals in training and test set per model, respectively. The
 92 averaged summed LL across regions was closer to zero for the nonlinear Bayesian model than for the linear Bayesian
 93 model, both for the training and the test set ($\sum \frac{1}{n_{test}} LL_{HBGPM}$, test set: -1.109, $\sum \frac{1}{n_{test}} LL_{HBLM}$ test set: -1.121;
 94 $\sum \frac{1}{n_{train}} LL_{HBGPM}$, training set: -1.020, $\sum \frac{1}{n_{train}} LL_{HBLM}$, training set: -1.05. LL values were less close to zero
 95 for all comparison models, with the *Combat* model performing best for among those models, followed by the *raw data*
 96 model, the *residuals* model and the *ComBat w/o covariates* model (an overview of the log likelihood for all models is
 97 given in Tab. 4, main text) The distribution of the log likelihood for all regions is given in Fig. 3d, main text.

Figure 8: Effective sample sizes N_{eff} for all parameters

98 References

- 99 [1] Bob Carpenter et al. “Stan: A probabilistic programming language”. In: *Journal of statistical software* 76.1 (2017).
- 100 [2] Jean Philippe Fortin et al. “Harmonization of cortical thickness measurements across scanners and sites”. In: *NeuroImage* 167.June 2017 (2018), pp. 104–120. ISSN: 10959572. DOI: 10.1016/j.neuroimage.2017.11.024.
- 102 URL: <https://doi.org/10.1016/j.neuroimage.2017.11.024>.
- 103 [3] Jean Philippe Fortin et al. “Harmonization of multi-site diffusion tensor imaging data”. In: *NeuroImage* 161 (2017),
- 104 pp. 149–170. ISSN: 10959572. DOI: 10.1016/j.neuroimage.2017.08.047.
- 105 [4] Andrew Gelman, Daniel Lee, and Jiqiang Guo. “Stan: A Probabilistic Programming Language for Bayesian
- 106 Inference and Optimization”. In: *Journal of Educational and Behavioral Statistics* 40.5 (2015), pp. 530–543. ISSN:
- 107 10769986. DOI: 10.3102/1076998615606113.
- 108 [5] Andrew Gelman, Donald B Rubin, et al. “Inference from iterative simulation using multiple sequences”. In:
- 109 *Statistical science* 7.4 (1992), pp. 457–472.
- 110 [6] Charles Geyer. “Introduction to markov chain monte carlo”. In: *Handbook of markov chain monte carlo* 20116022
- 111 (2011), p. 45.
- 112 [7] Raymond Pomponio et al. “Harmonization of large MRI datasets for the analysis of brain imaging patterns through-
- 113 out the lifespan”. In: *NeuroImage* 208 (2020), p. 116450. ISSN: 1053-8119. DOI: 10.1016/j.neuroimage.
- 114 2019.116450.
- 115 [8] Stan Development Team. *Stan Modeling Language Users Guide and Reference Manual, Version 2.25*. 2020. URL:
- 116 <http://mc-stan.org/>.