# ADDITIONAL FILE 1: SUPPLEMENTARY FIGURES

## Screening for genes that accelerate the epigenetic ageing clock in humans reveals a role for the H3K36 methyltransferase NSD1

Daniel E. Martin-Herranz[1,2], Erfan Aref-Eshghi[3,4], Marc Jan Bonder[1,5], Thomas M. Stubbs[2], Sanaa Choufani[6], Rosanna Weksberg[6], Oliver Stegle[1,5,7], Bekim Sadikovic[3,4], Wolf Reik[8,9,10,*], Janet M. Thornton[1,*]

[1]European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge, UK; [2]Chronomics Ltd, Cambridge, UK; [3]Department of Pathology and Laboratory Medicine, Western University, London, Canada; [4]Molecular Genetics Laboratory, Molecular Diagnostics Division, London Health Sciences Centre, London, Canada; [5]European Molecular Biology Laboratory, Genome Biology Unit, Heidelberg, Germany; [6]Genetics and Genome Biology Program, Research Institute, The Hospital for Sick Children, Toronto, Canada; [7]Division of Computational Genomics and Systems Genetics, German Cancer Research Center (DKFZ), Heidelberg, Germany; [8]Epigenetics Programme, The Babraham Institute, Cambridge, UK; [9]Centre for Trophoblast Research, University of Cambridge, Cambridge, UK; [10]Wellcome Sanger Institute, Hinxton, Cambridge, UK; *These authors jointly supervised this work
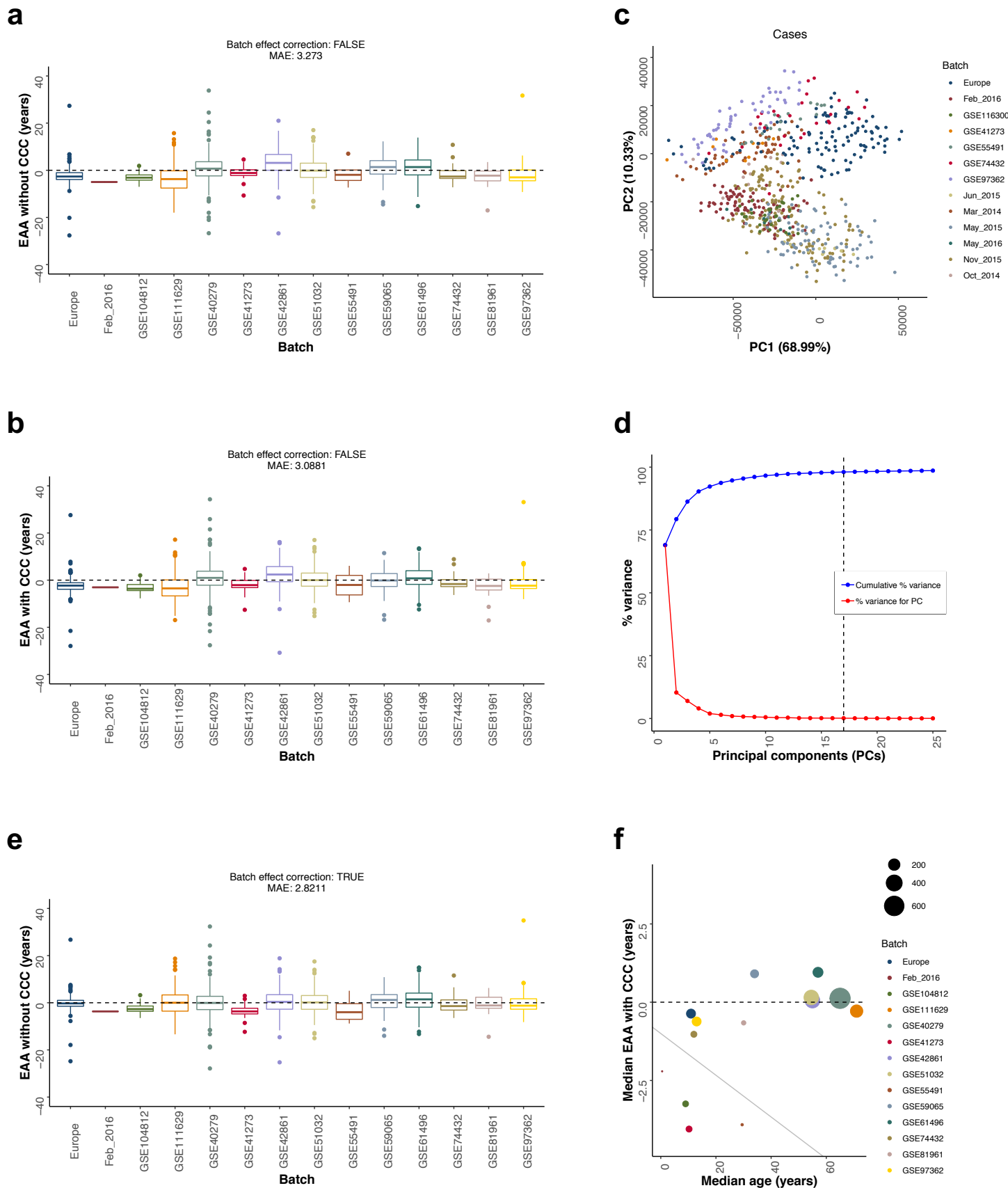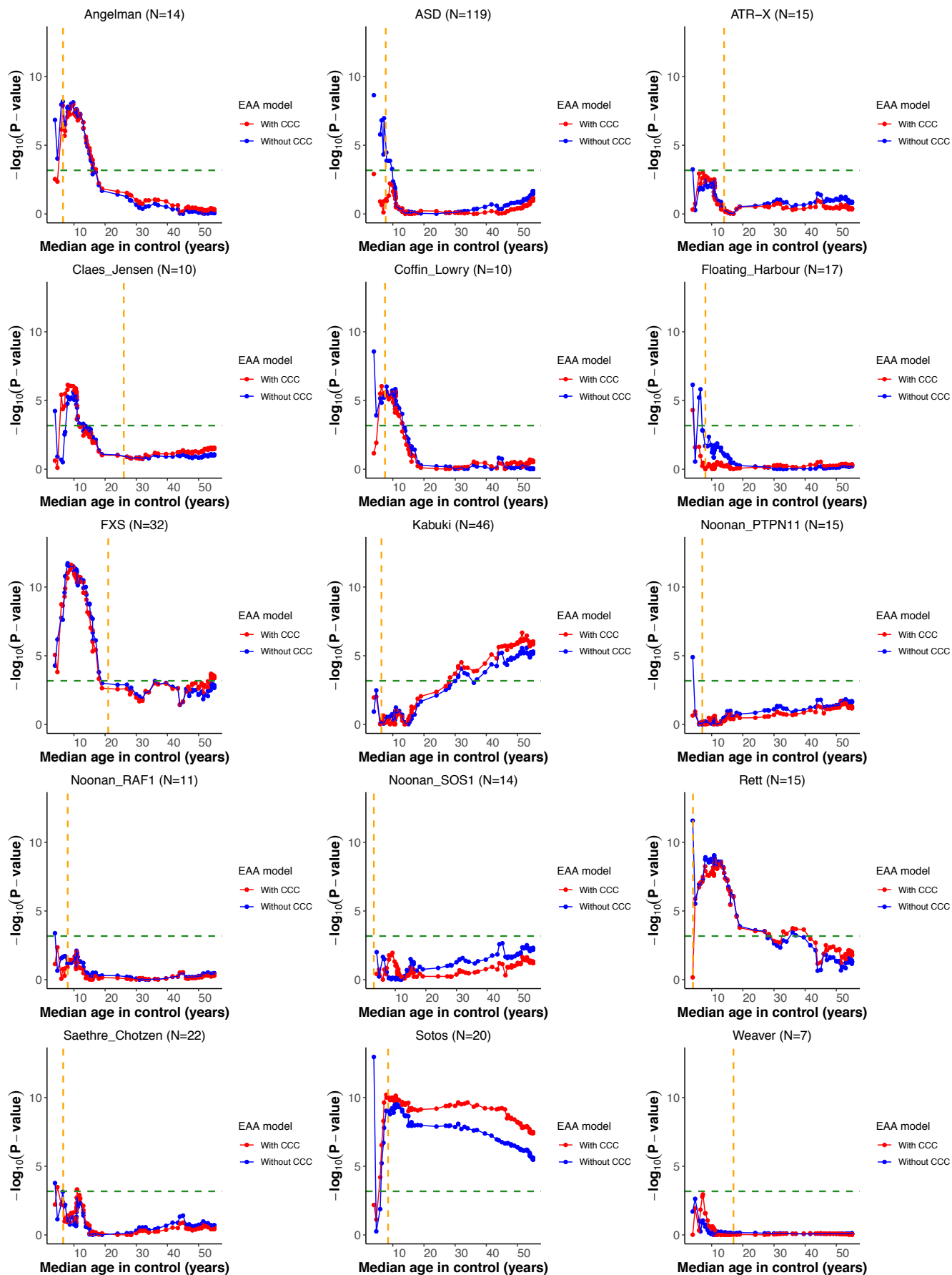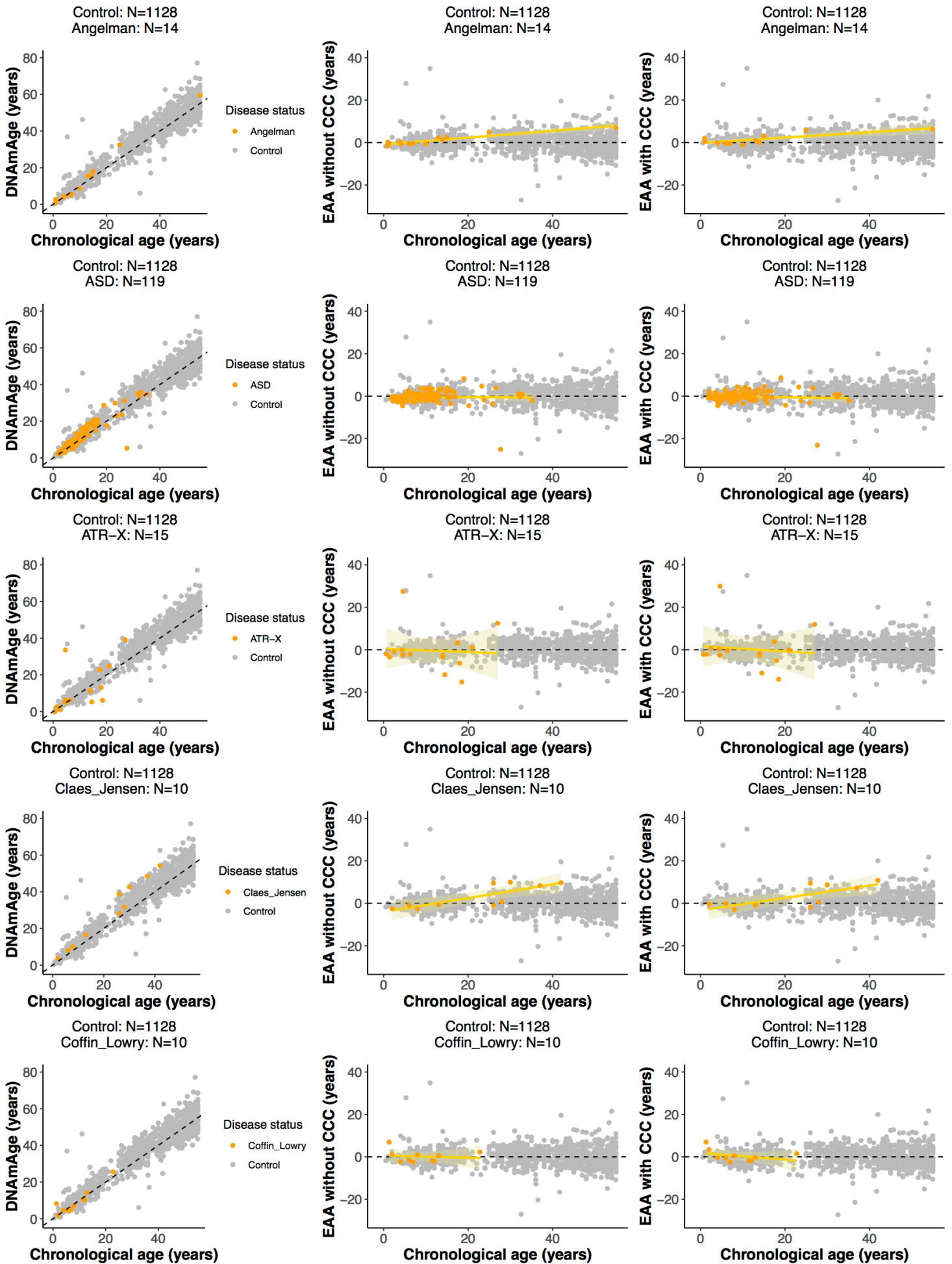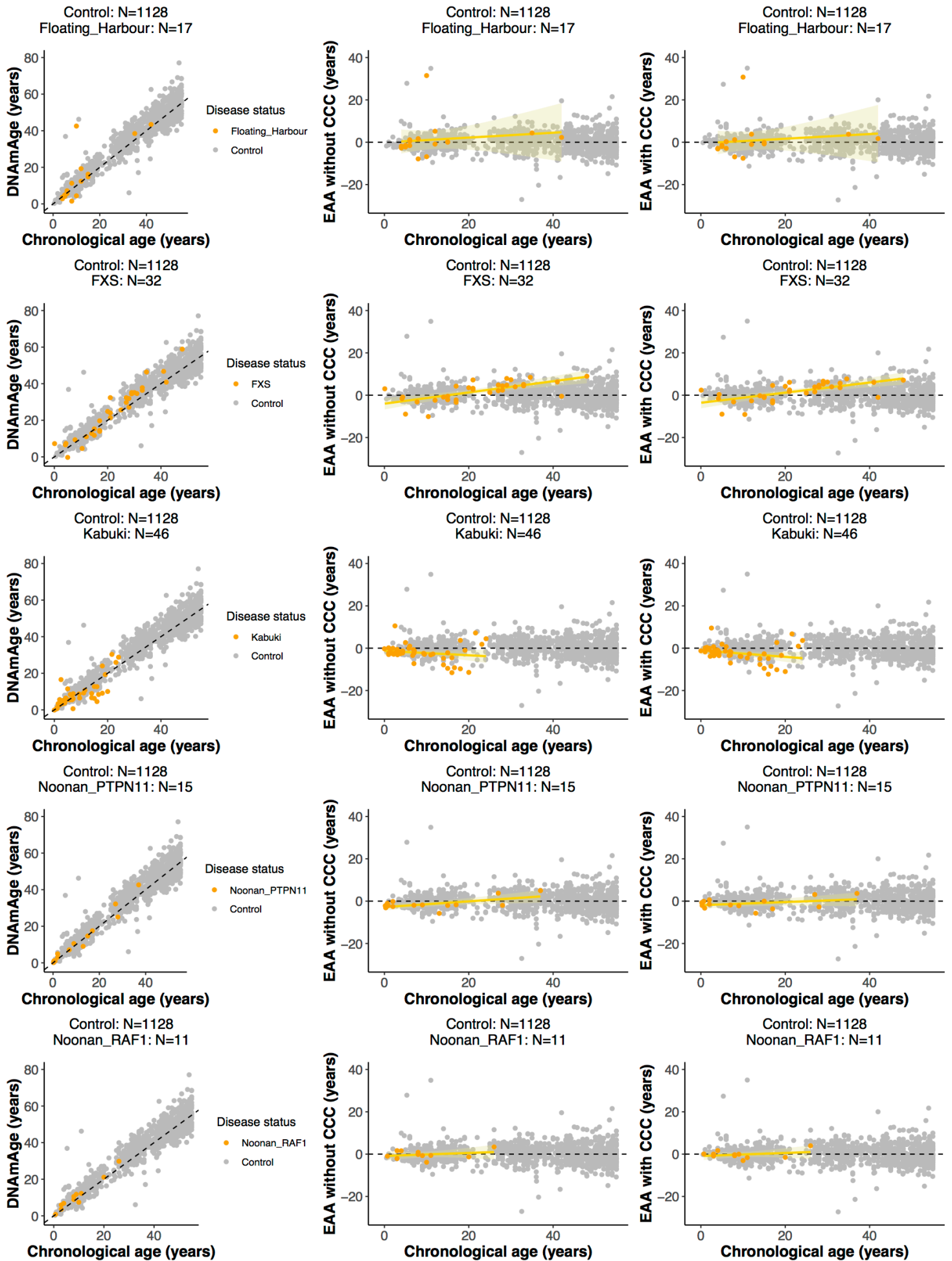
# Figure S1

**Figure S1.** Complementary to Fig. 1. **a.** Distribution of the epigenetic age acceleration (EAA) without cell composition correction (CCC) for the different control batches, before applying batch effect correction. MAE: median absolute error. **b.** Distribution of the EAA with cell composition correction (CCC) for the different control batches, before applying batch effect correction. **c.** Scatterplot showing the values of the first two principal components (PCs) for the cases (developmental disorder samples) after performing PCA on the control probes of the 450K arrays. Each point corresponds to a different case sample and the colours represent the different batches. The different batches cluster together in the PCA space, showing that the control probes indeed capture technical variation. Please note that all the PCA calculations were done with more samples from cases and controls than those that were included in the final screening since it was performed before the filtering step (see Methods). **d.** Plot showing the percentages of technical variance explained by the different PCs from the control probes. The dashed line represents the optimal number of PCs (17) that was finally used. **e.** Distribution of the EAA without cell composition correction (CCC) for the different control batches, after applying batch effect correction. **f.** After batch effect correction, deviations from a median EAA of zero (dotted black line) in some of the control batches can be explained by other causes. The grey line separates in the lower left corner those weird batches (Feb_2016, GSE104812, GSE41273, GSE55491), which have a small sample size and/or a low median age.

# Figure S2

**a**



Angelman (N=14) · ASD (N=119) · ATR−X (N=15) · Claes_Jensen (N=10) · Coffin_Lowry (N=10) · Floating_Harbour (N=17) · FXS (N=32) · Kabuki (N=46) · Noonan_PTPN11 (N=15) · Noonan_RAF1 (N=11) · Noonan_SOS1 (N=14) · Rett (N=15) · Saethre_Chotzen (N=22) · Sotos (N=20) · Weaver (N=7)

EAA model: With CCC, Without CCC

x-axis: Median age in control (years); y-axis: $-\log_{10}(\text{P}-\text{value})$

**b**

Control: N=1128
Angelman: N=14

Control: N=1128
Angelman: N=14

Control: N=1128
Angelman: N=14

Control: N=1128
ASD: N=119

Control: N=1128
ASD: N=119

Control: N=1128
ASD: N=119

Control: N=1128
ATR−X: N=15

Control: N=1128
ATR−X: N=15

Control: N=1128
ATR−X: N=15

Control: N=1128
Claes_Jensen: N=10

Control: N=1128
Claes_Jensen: N=10

Control: N=1128
Claes_Jensen: N=10

Control: N=1128
Coffin_Lowry: N=10

Control: N=1128
Coffin_Lowry: N=10

Control: N=1128
Coffin_Lowry: N=10

Disease status
- Angelman
- Control

Disease status
- ASD
- Control

Disease status
- ATR−X
- Control

Disease status
- Claes_Jensen
- Control

Disease status
- Coffin_Lowry
- Control

Control: N=1128
Floating_Harbour: N=17

Control: N=1128
Floating_Harbour: N=17

Control: N=1128
Floating_Harbour: N=17

Control: N=1128
FXS: N=32

Control: N=1128
FXS: N=32

Control: N=1128
FXS: N=32

Control: N=1128
Kabuki: N=46

Control: N=1128
Kabuki: N=46

Control: N=1128
Kabuki: N=46

Control: N=1128
Noonan_PTPN11: N=15

Control: N=1128
Noonan_PTPN11: N=15

Control: N=1128
Noonan_PTPN11: N=15

Control: N=1128
Noonan_RAF1: N=11

Control: N=1128
Noonan_RAF1: N=11

Control: N=1128
Noonan_RAF1: N=11

**Figure S2.** Complementary to Fig. 2. **a.** Effect of changing the median age of the controls when performing the screening for epigenetic age acceleration (EAA) in the different developmental disorders. The dashed green line displays the significance level of $\alpha$ = 0.01 after Bonferroni correction. The dashed orange line displays the median age for the samples in the developmental disorder considered. In blue: EAA model without cell composition correction (CCC). In red: EAA model with CCC. **b.** Left panel: scatterplot showing the relation between epigenetic age (*DNAmAge*) according to Horvath's model and chronological age of the samples for a given developmental disorder (orange) and control (grey). Each sample is represented by one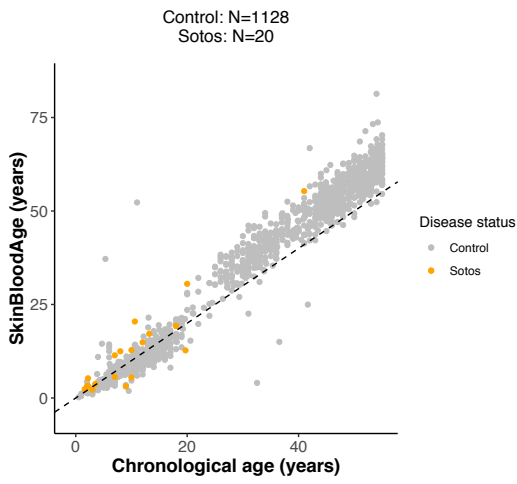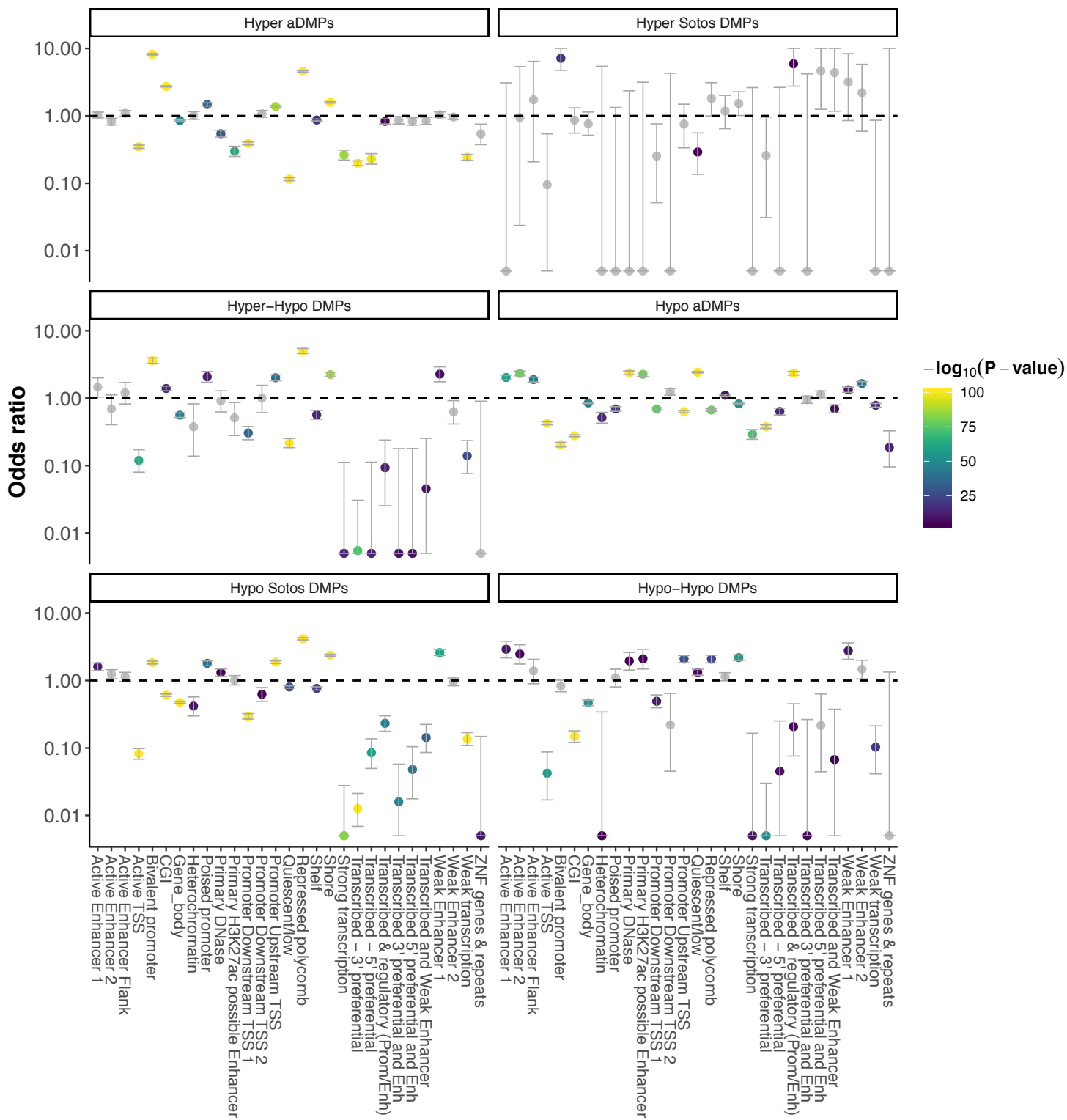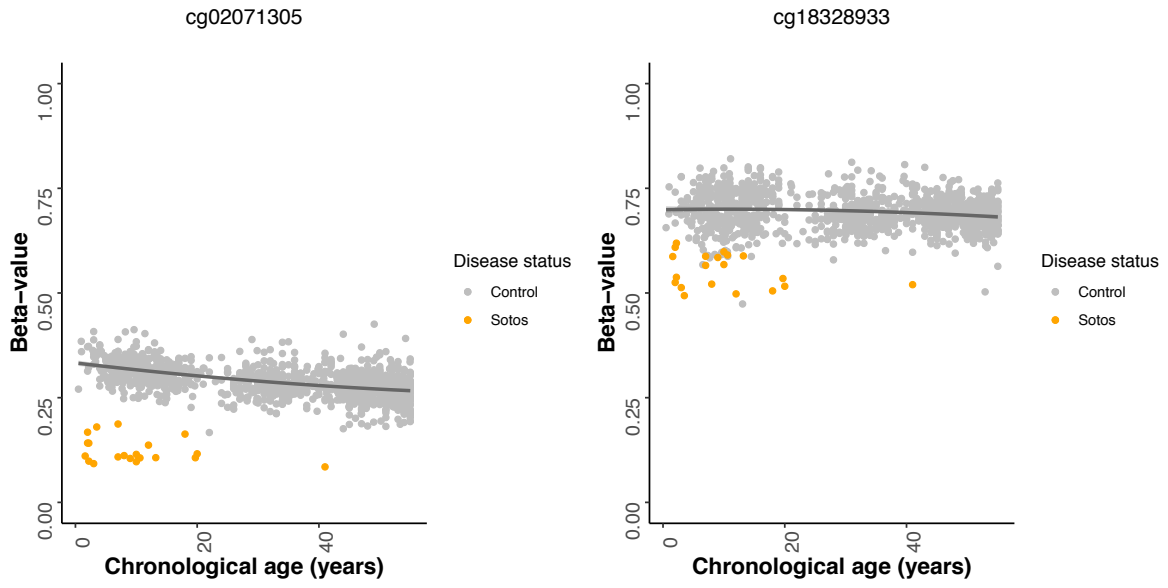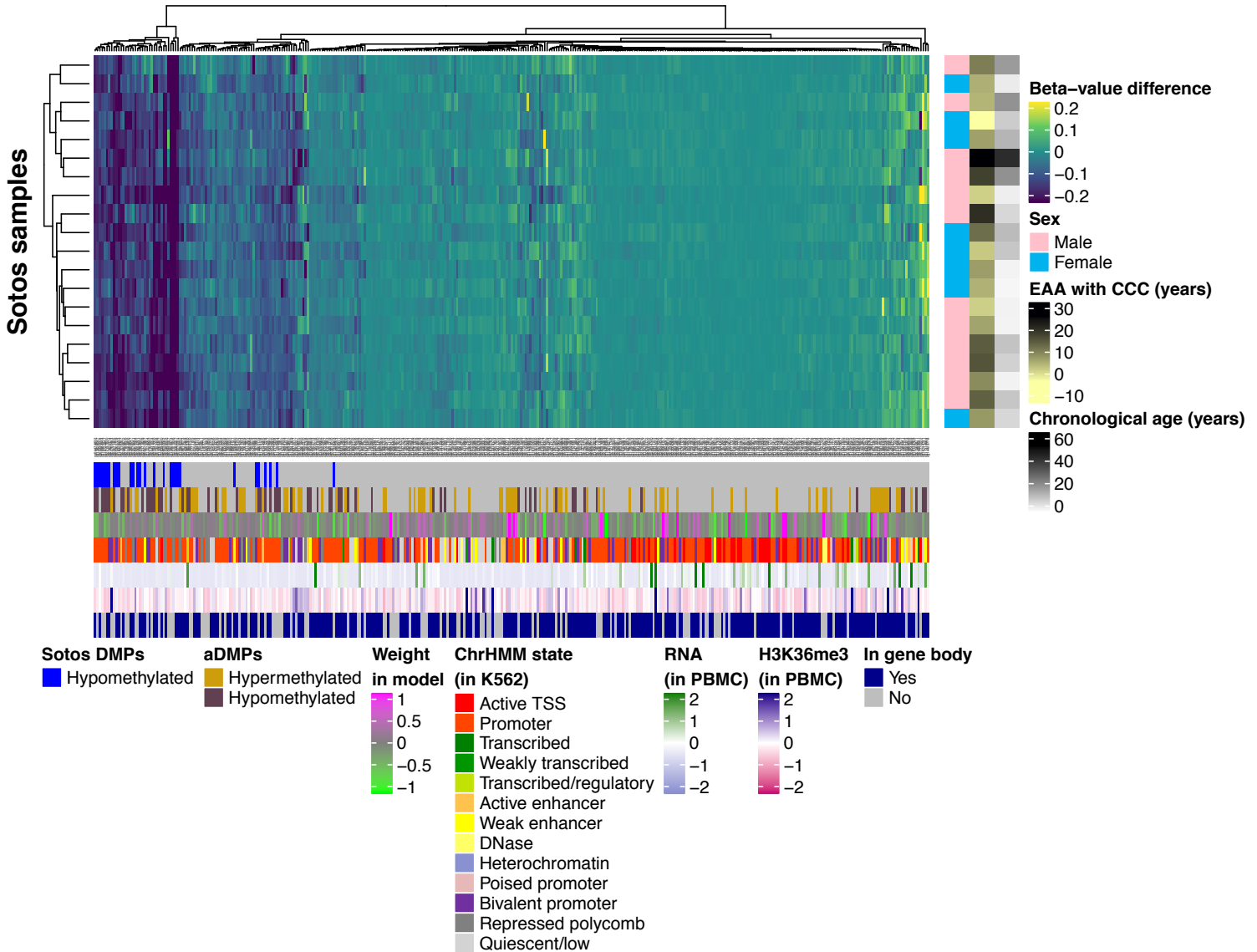 point. The black dashed line represents the diagonal to aid visualisation. Middle and right panels: scatterplots showing the relation between the epigenetic age acceleration (EAA) (without and with CCC respectively) and chronological age of the samples for a given developmental disorder (orange) and control (grey). Each sample is represented by one point. The yellow line represents the linear model *EAA ~ Age*, with the standard error shown in the light yellow shade. **c.** Left panel: scatterplot showing the relation between epigenetic age (*HannumAge*) according to Hannum's model and chronological age of the samples for Sotos (orange) and control (grey). Each sample is represented by one point. The black dashed line represents the diagonal to aid visualisation. Middle and right panels: boxplots showing the comparisons of the EAA distributions (according to Hannum's clock) for Sotos and control samples (without and with CCC respectively). The p-values (two-sided Wilcoxon's test, before multiple testing correction) are shown above the boxplots. **d.** As in c., but using Lin's epigenetic clock. **e.** As in c., but using the skin-blood epigenetic clock.

# Figure S3

**a**

**b**

**c**

cg02071305

cg18328933

Disease status
- Control
- Sotos

**d**

**Horvath's clock CpGs**

Sotos samples

Beta−value difference
- 0.2
- 0.1
- 0
- −0.1
- −0.2

**Sex**
- Male
- Female

**EAA with CCC (years)**
- 30
- 10
- 0
- −10

**Chronological age (years)**
- 60
- 40
- 20
- 0

**Sotos DMPs**
- Hypomethylated

**aDMPs**
- Hypermethylated
- Hypomethylated

**Weight in model**
- 1
- 0.5
- 0
- −0.5
- −1

**ChrHMM state (in K562)**
- Active TSS
- Promoter
- Transcribed
- Weakly transcribed
- Transcribed/regulatory
- Active enhancer
- Weak enhancer
- DNase
- Heterochromatin
- Poised promoter
- Bivalent promoter
- Repressed polycomb
- Quiescent/low

**RNA (in PBMC)**
- 2
- 1
- 0
- −1
- −2

**H3K36me3 (in PBMC)**
- 2
- 1
- 0
- −1
- −2

**In gene body**
- Yes
- No

**e**

**Horvath's clock CpGs**

Features:
H3K9me3_ENCFF713QZB
H3K9me3_ENCFF319EBK
H3K9me3_ENCFF033IPJ
H3K9me3_ENCFF171WZC
RNF2_ENCFF071CIY
RNF2_ENCFF847TGB
RNF2_ENCFF320VKN
RNF2_ENCFF857HEZ
H3K27me3_ENCFF412KUE
H3K27me3_ENCFF150RIG
H3K27me3_ENCFF265VZG
EZH2_ENCFF516PTT
H3K4me1_ENCFF457WMB
H3K36me3_ENCFF643USH
H3K36me3_ENCFF249WVX
H3K4me3_ENCFF573QMJ
H3K9ac_ENCFF455IGC
H3K4me3_ENCFF796FFT
H3K4me3_ENCFF303YKC
H3K27ac_ENCFF759GIZ
H3K27ac_ENCFF737OJY
H3K4me1_ENCFF100NYH
H3K9ac_ENCFF211ORP

**Z-score (in PBMC):** 4, 2, 0, −2, −4

**Cell type:** B cell, K562, PBMC

**Sotos DMPs:** Hypomethylated

**aDMPs:** Hypermethylated, Hypomethylated

**Weight in model:** 1, 0.5, 0, −0.5, −1

**ChrHMM state (in K562):** Active TSS, Promoter, Transcribed, Weakly transcribed, Transcribed/regulatory, Active enhancer, Weak enhancer, DNase, Heterochromatin, Poised promoter, Bivalent promoter, Repressed polycomb, Quiescent/low

**RNA (in PBMC):** 2, 1, 0, −1, −2

**In gene body:** Yes, No

**f**

**Odds ratio**

Panels: All Horvath, Hyper aDMPs, Hypo aDMPs, Hypo Sotos DMPs

**−log₁₀(P−value):** $-\log_{10}(P-\text{value})$ — 8.4, 8.2, 8.0

X-axis categories: Active Enhancer 1, Active Enhancer 2, Active Enhancer Flank, Active TSS, Bivalent promoter, CGI, Gene_body, Heterochromatin, Poised promoter, Primary DNase, Primary H3K27ac possible Enhancer, Promoter Downstream TSS 1, Promoter Downstream TSS 2, Promoter Upstream TSS, Quiescent/low, Repressed polycomb, Shelf, Shore, Strong transcription, Transcribed, Transcribed − 3' preferential, Transcribed − 5' preferential, Transcribed & regulatory (Prom/Enh), Transcribed 5' preferential and Enh, Transcribed and Weak Enhancer, Weak Enhancer 1, Weak Enhancer 2, Weak transcription, ZNF genes & repeats

**g**



Feature: H3K27ac, Feature: H3K4me3, Feature: H3K36me3, Feature: H3K27me3, Feature: H3K9ac, Feature: H3K4me1, Feature: H3K9me3, Feature: RNF2, Feature: EZH2, Feature: RNA, Feature: Replication_timing, Feature: LaminB1 — boxplots comparing Control vs In subset across All Horvath, Hyper aDMPs, Hypo aDMPs, Hypo Sotos DMPs.

**Figure S3.** Complementary to Fig. 3. **a.** Enrichment for the categorical (epi)genomic features considered when comparing the different genome-wide subsets of differentially methylated positions (DMPs) in ageing and Sotos against a control (see Methods). The y-axis represents the odds ratio (OR), the error bars show the 95% confidence interval for the OR estimate and the colour of the points codes for -$\log_{10}$(p-value) obtained after testing for enrichment using Fisher's exact test. An OR > 1 shows that the given feature is enriched in the subset of DMPs considered, whilst an OR < 1 shows that it is found less than expected. The 'Hyper-Hypo DMPs' subset results from the intersection between the hypermethylated DMPs in ageing and the hypomethylated DMPs in Sotos. The 'Hypo-Hypo DMPs' subset results from the intersection between the hypomethylated DMPs in ageing and Sotos. In grey: features that did not reach significance using a significance level of $\alpha$ = 0.01 after Bonferroni correction. **b.** Boxplots showing the distributions of scores (see Methods) for the continuous (epi)genomic features considered when comparing the different genome-wide subsets of differentially methylated positions (DMPs) in ageing and Sotos against a control (see Methods). The p-values (two-sided Wilcoxon's test, before multiple testing correction) are shown above the boxplots. The number of DMPs belonging to each subset (in green) and the median value of the feature score (in dark red) are shown below the boxplots. NFC: 'normalised fold change'; NRE: 'normalised RNA expression'; WTS: 'wavelet-transformed signals'; NRC: 'normalised read counts'. **c.** DNA methylation (beta-value) profiles for two of the clock CpG sites (cg02071305 and cg18328933). A linear model (displayed in dark grey) can be fixed to each CpG site to model the changes in beta-value with chronological age in the controls (grey). Afterwards, the difference of the Sotos samples beta-values (orange) with the controls can be estimated. **d.** Heatmap displaying the differential methylation patterns for Sotos samples (rows) when compared with controls in each one of the 353 clock CpGs (columns). Hierarchical clustering was performed in both rows and columns. RNA refers to the 'normalised RNA expression' (NRE, see Methods). H3K36me3 refers to the H3K36me3 histone modification 'normalised fold change' (NFC, see Methods). aDMPs: differentially methylated positions during ageing. EAA: epigenetic age acceleration. CCC: cell composition correction. PBMC: peripheral blood mononuclear cells. **e.** Heatmap displaying the scores for the different continuous (epi)genomic features (rows) in each one of the 353 clock CpGs (columns). The names of the features include the ENCODE ID. Hierarchical clustering was performed in both rows and columns. **f.** Same as a., but focused on the 353 clock CpG sites. **g.** Same as b., but focused on the 353 clock CpG sites.
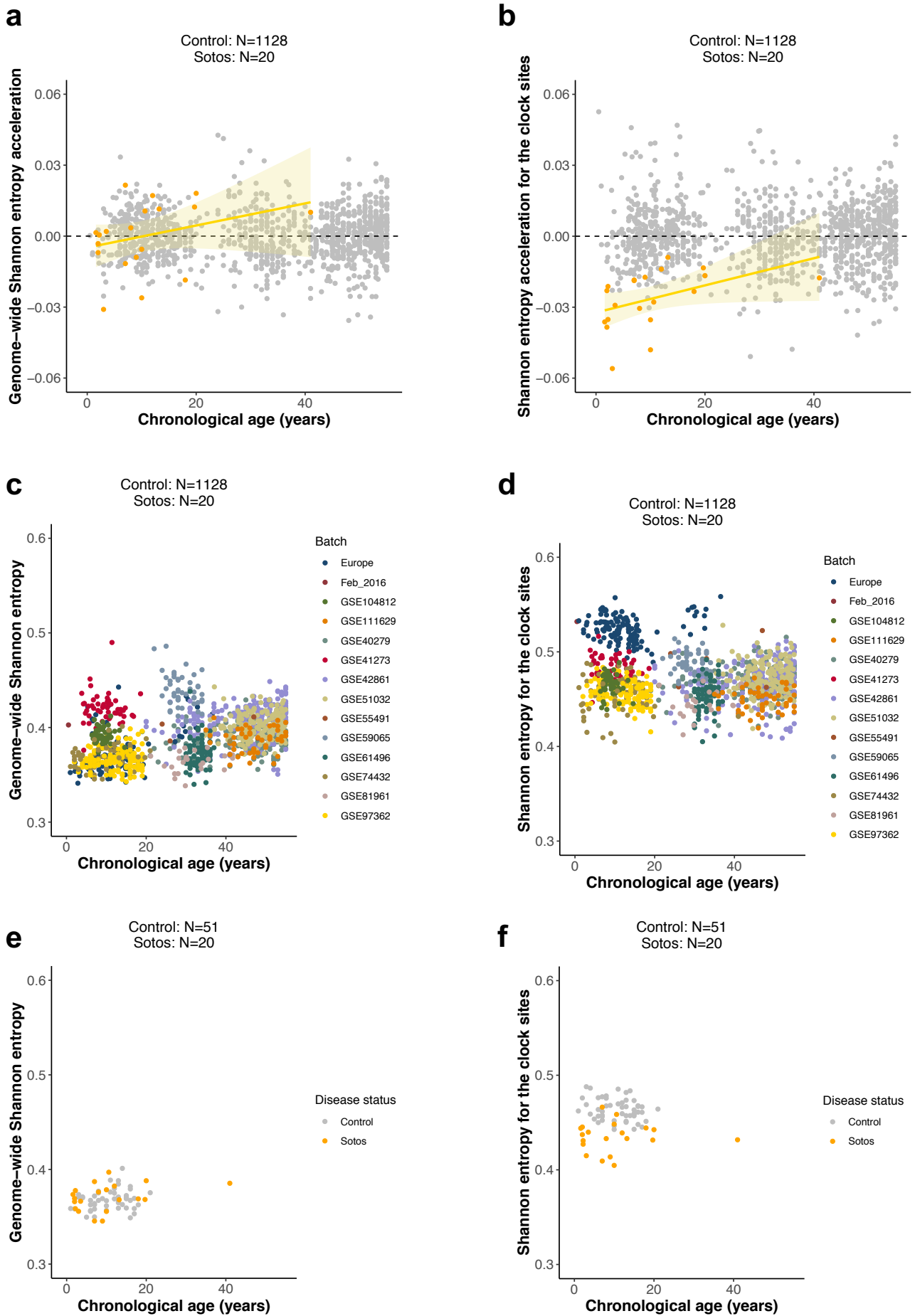
**Figure S4**

**Figure S4.** Complementary to Fig. 4. **a.** Scatterplot showing the relation between the genome-wide Shannon entropy acceleration (*gSEA*) and chronological age of the samples for Sotos (orange) and healthy controls (grey). Each sample is represented by one point. The yellow line represents the linear model *gSEA ~ Age*, with the standard error shown in the light yellow shade. **b.** Same as a., but using the Shannon entropy acceleration calculated only for the 353 CpG sites in the Horvath epigenetic clock (*cSEA*). **c.** Scatterplot showing the effects of the different batches on the genome-wide Shannon entropy calculations. Each sample is represented by one point and coloured according to the batch that they belong to. **d.** Same as c., but using the Shannon entropy calculated only for the 353 CpG sites in the Horvath epigenetic clock. **e.** Scatterplot showing the relation between genome-wide Shannon entropy and chronological age of the samples for Sotos (orange) and healthy controls (grey) in the GSE74432 batch (i.e. the batch where all the Sotos samples come from). Each sample is represented by one point. **f.** Same as e., but using the Shannon entropy calculated only for the 353 CpG sites in the Horvath epigenetic clock.
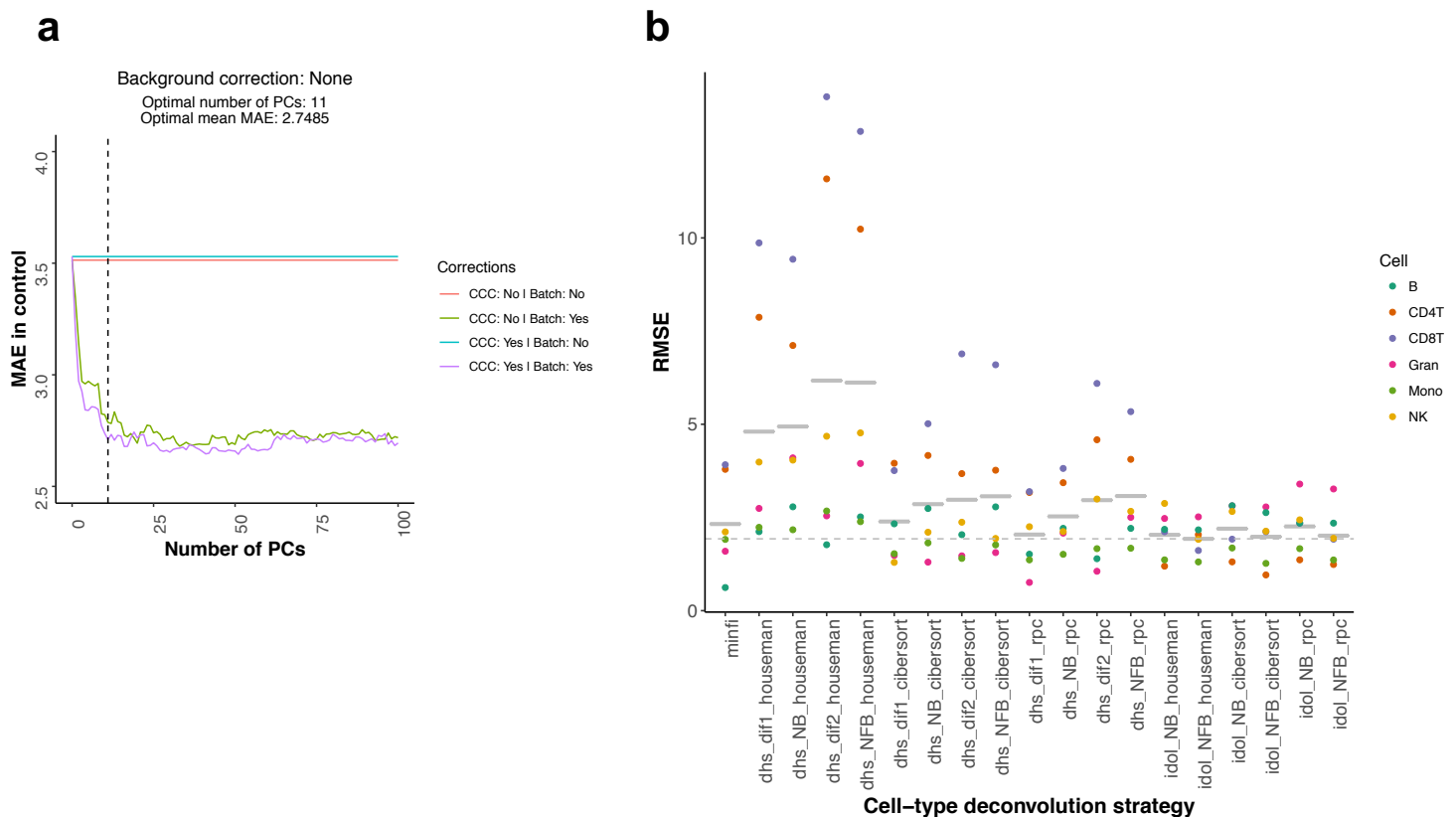
# Figure S5

**Figure S5.** Complementary to Methods. **a.** Plot showing how the median absolute error (MAE) of the prediction in the control samples, that should tend to zero, is reduced when the PCs capturing the technical variation are included as part of the modelling strategy (see Methods). In this case, no background correction was performed, as opposed to the results in Fig. 1c. The dashed line represents the optimal number of PCs (11) that was finally used. The optimal mean MAE is calculated as the average MAE between the green and purple lines. CCC: cell composition correction. **b.** Benchmarking of the different strategies for cell-type deconvolution in blood. The x-axis shows the different strategies that were tested (a more detailed description of these strategies can be found in Additional file 4). The y-axis shows the root mean square error (RMSE) obtained when comparing our predictions with the real proportions of cells in a gold-standard dataset (GSE77797). The grey horizontal solid lines represent the mean for the RMSE across cell types and the grey dashed line the minimum of these values.