# Supplemental Information
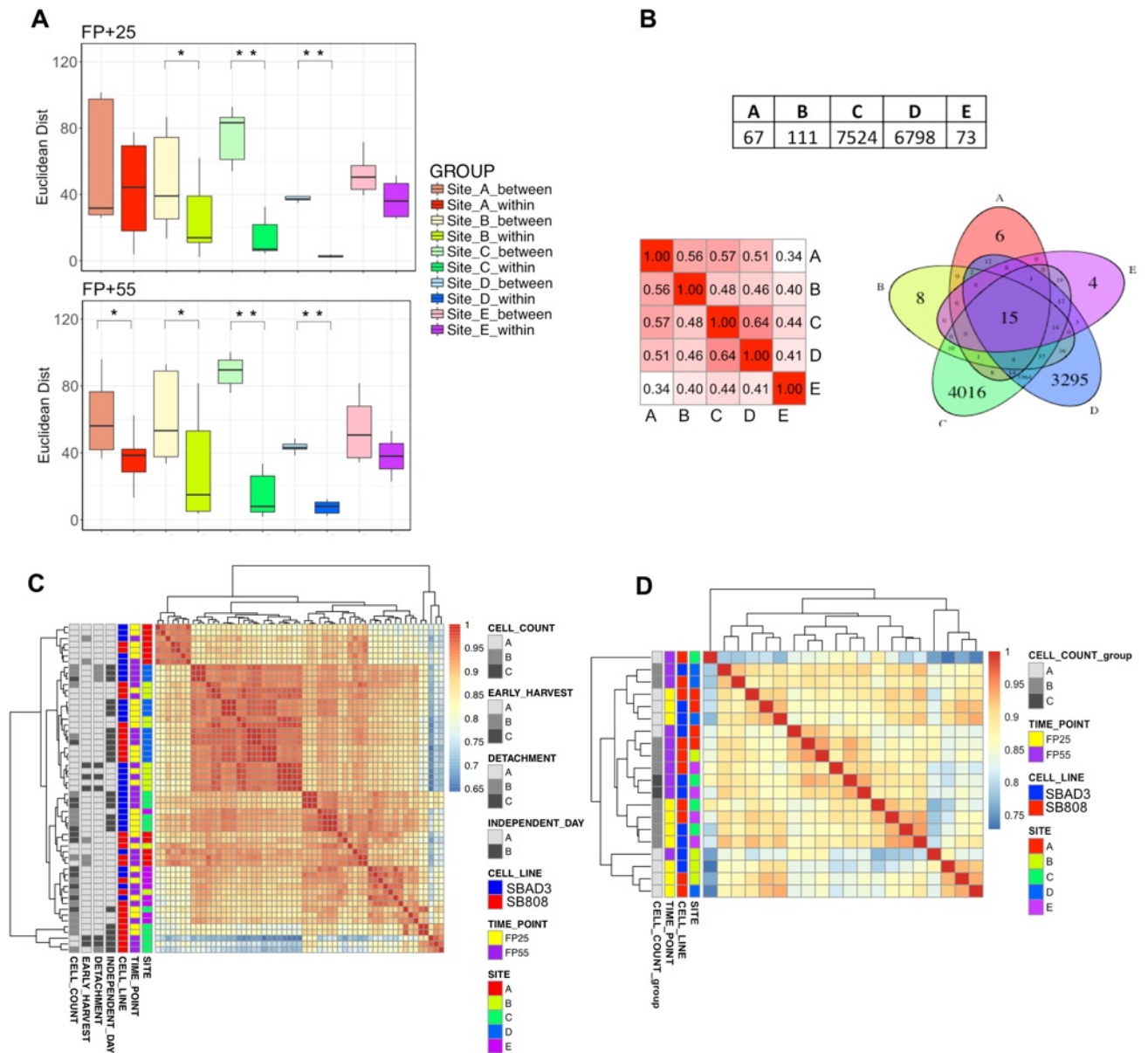
# Reproducibility of Molecular Phenotypes after Long-Term Differentiation to Human iPSC-Derived Neurons: A Multi-Site Omics Study

Viola Volpato, James Smith, Cynthia Sandor, Janina S. Ried, Anna Baud, Adam Handel, Sarah E. Newey, Frank Wessely, Moustafa Attar, Emma Whiteley, Satyan Chintawar, An Verheyen, Thomas Barta, Majlinda Lako, Lyle Armstrong, Caroline Muschet, Anna Artati, Carlo Cusulin, Klaus Christensen, Christoph Patsch, Eshita Sharma, Jerome Nicod, Philip Brownjohn, Victoria Stubbs, Wendy E. Heywood, Paul Gissen, Roberta De Filippis, Katharina Janssen, Peter Reinhardt, Jerzy Adamski, Ines Royaux, Pieter J. Peeters, Georg C. Terstappen, Martin Graf, Frederick J. Livesey, Colin J. Akerman, Kevin Mills, Rory Bowden, George Nicholson, Caleb Webber, M. Zameel Cader, and Viktor Lakics
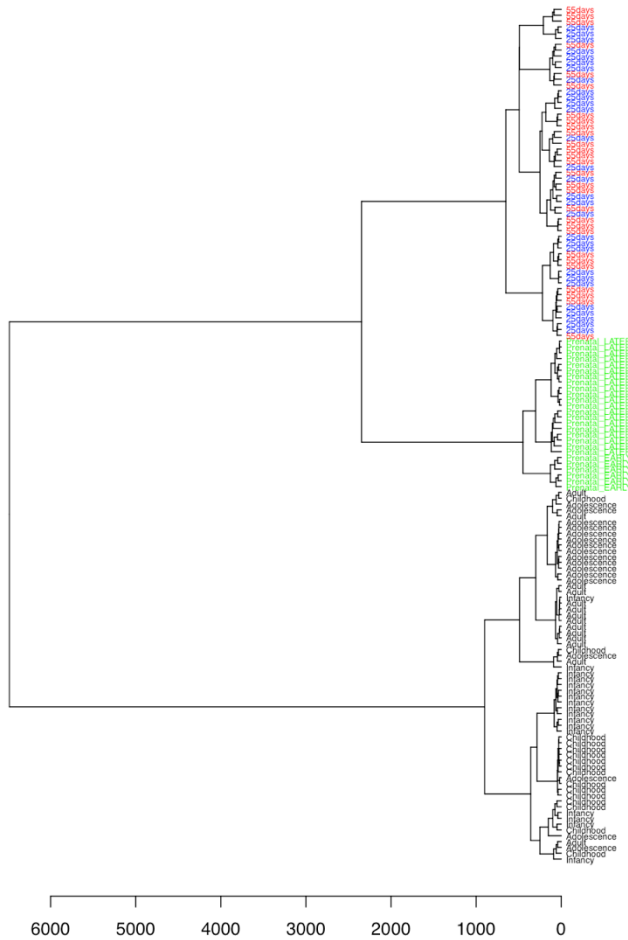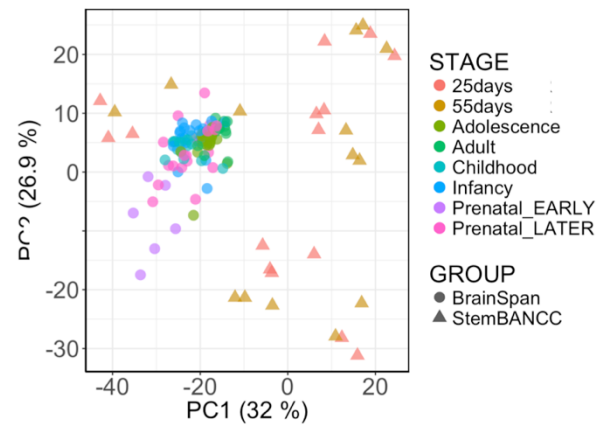
**Supplemental Figures**



**Figure S1, related to Figure 2. Variation in omics readouts reproducibility within and between laboratories.**

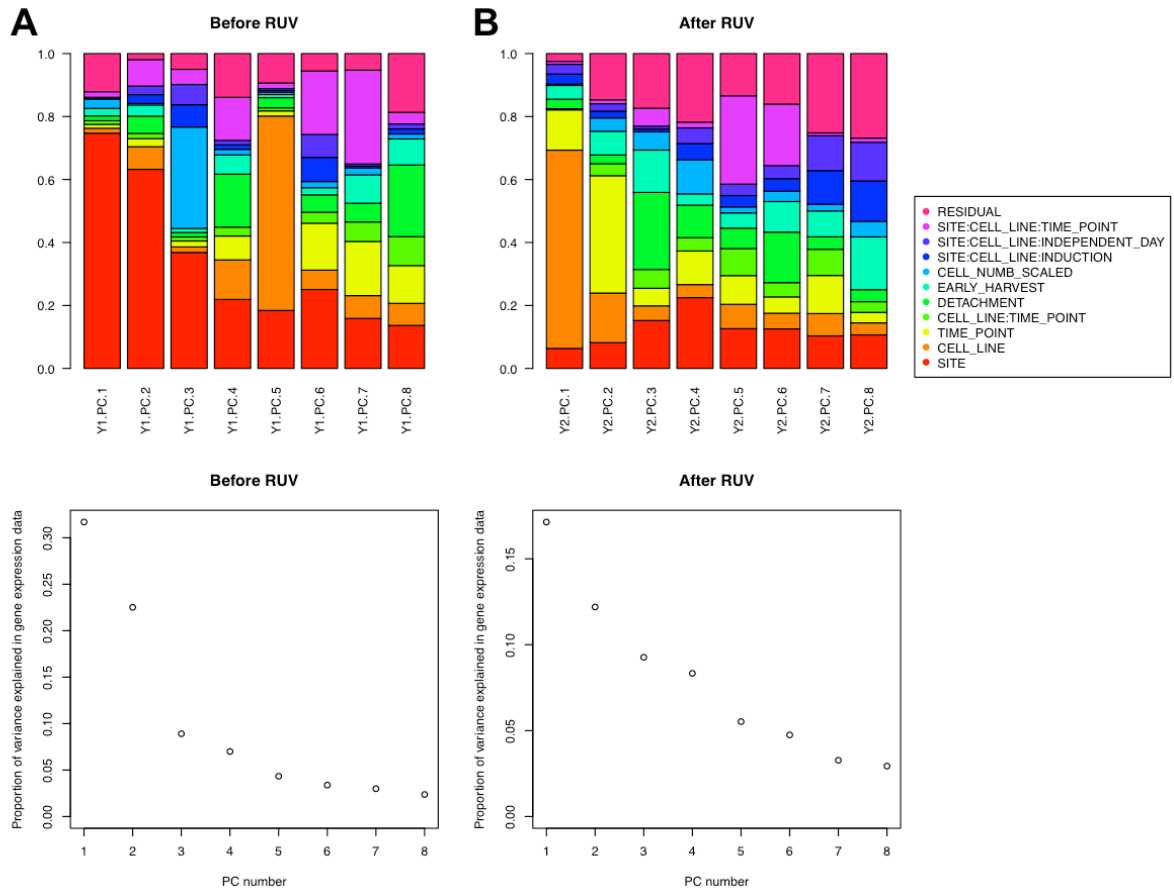**A) Euclidean distances within cell line replicates and between all replicates of different cell lines.** Euclidean distances are calculated between the gene expression profiles of each sample within each laboratory and between replicates of different lines in any laboratory by time point (FP+25 top, FP+55 bottom). Within each laboratory the expression profiles derived from replicates of the same line are significantly closer to each other than those between replicates of different lines for 4 out of 5 laboratories (** p-value<0.005, * p-value<0.05). Box-and-whisker graphs represent distributions, where the span of the box is the interquartile range (IQR) and includes the median (bold line). The ends of the upper and lower whiskers represent the data point with the maximum distance from the third and first quartiles, respectively, but no further than 1.5 * IQR.

Data beyond the end of the whiskers are outliers. **B) Results of differential expression analysis between genotype within each laboratory.** Large difference between laboratories in the number of DE genes between genotypes controlling for time point variation before RUVs correction (top). Only 15 DE genes are found in common between all laboratories indicating a remarkably low degree of cross-laboratory reproducibility (bottom right). Semantic similarity scores between the top 40 enriched GO BP terms in DE genes between genotypes controlling for time point variation before RUVs correction within laboratory (bottom left). **C) Heatmap of pair-wise Spearman's rank based correlations of gene expression between all samples.** Correlations are computed on 13,373 genes expressed across all samples (coloured based on all known covariates both biological and technical, see Methods). **D) Heatmap of Spearman's rank based correlations of protein abundances normalised for total protein amount.** 1,037 proteins observed across all samples were used. As observed for the transcriptomics data, the heatmap did not show clustering of samples by genotype.
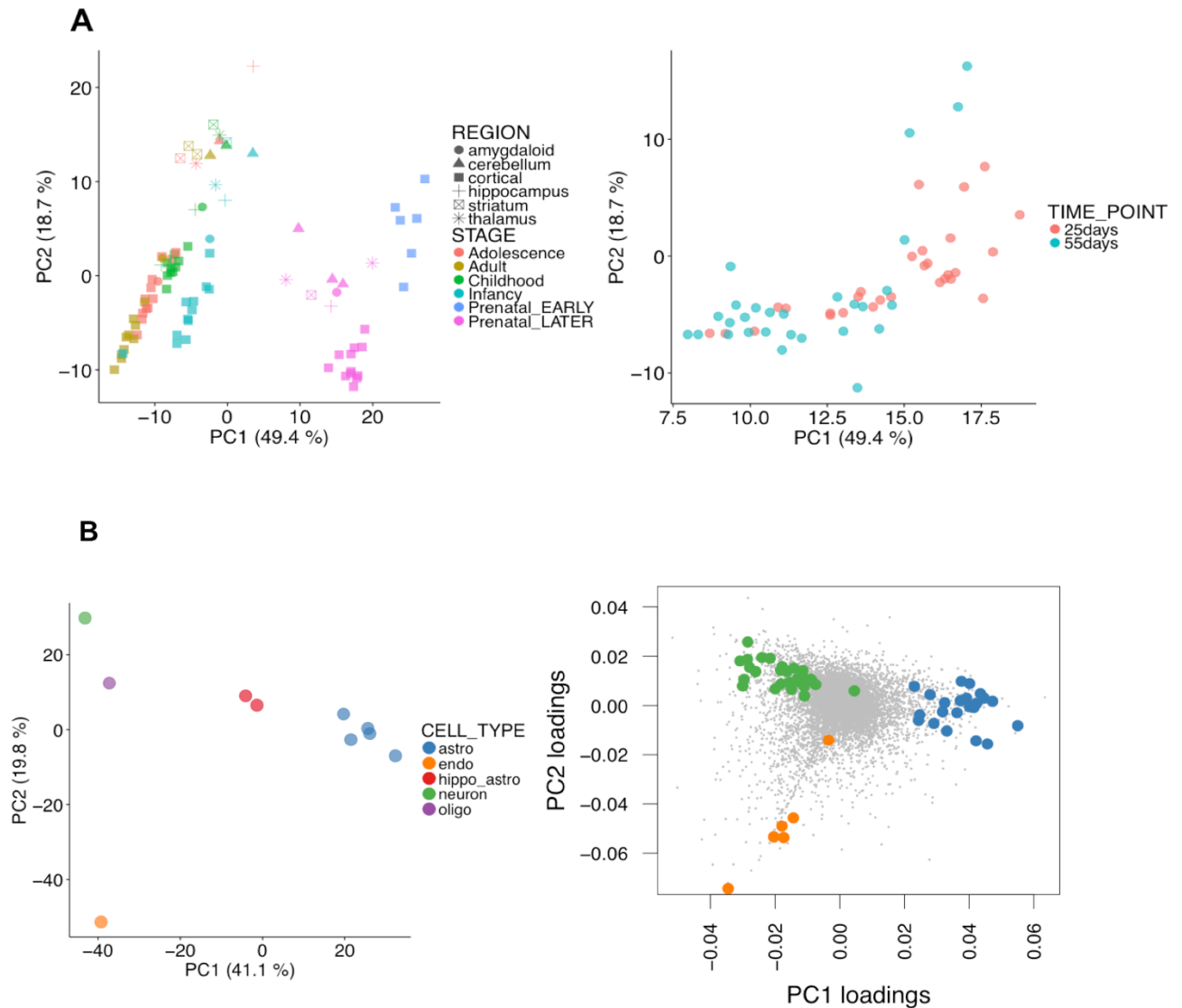
**Figure S2, related to Figure 3. Comparison of temporal gene expression profiles of StemBANCC samples with post-mortem brain samples from the BrainSpan Atlas of the Developing Human Brain.**

**A) Hierarchical clustering of the 57 StemBANCC samples and post-mortem brain samples from the BrainSpan Atlas.** StemBANCC samples cluster together with fetal post-mortem brain samples (BrainSpan fetal samples in green, BrainSpan postnatal samples in black, StemBANCC FP+25 samples in blue and StemBANCC FP+55 samples in red). **B) PCA plot of StemBANCC samples from the SBAD3 line before RUVs correction and the BrainSpan Atlas**. BrainSpan samples were projected on principal component axes of StemBANCC CTR samples before RUVs correction. The direction of human neuronal maturation was not recapitulated by non-normalized gene expression profiles.

**Figure S3, related to Figure 5. Variance component analysis of gene data before and after RUVs correction.**

Variance captured by first 8 principal components of gene counts before (A) and after RUVs correction on 5 RUVs factors (B). Proportions of variance explained by known covariates (see Methods) for any principal component (top) and proportions of explained variance in gene expression data by any principal component (bottom) are reported.

**Figure S4, related to Figure 5. A neuronal maturation axis and a neuron-glia axis are derived from external gene expression data and used to explain laboratory variability in StemBANCC samples. A) Identification of a transcriptional axis of maturation.**

PCA plot of BrainSpan samples (A) on a set of 787 'cortical marker genes' identified using GTEx data (see Methods) shows that BrainSpan data clearly cluster by sample age. StemBANCC samples were projected (right) on the transcriptional maturation axis (first component) identified in PCA plot (left). Except for the three FP+55 SB808 outlier samples from laboratory C that cluster with the FP+25 samples, the identified transcriptional maturation axis clearly separates samples by time point. The position of StemBANCC samples projected along this axis was used as covariate named 'MATURITY' in subsequent variance component analysis (see Methods). **B) Identification of extended lists of cell-type specific genes.** PCA plot (left) of RNAseq data from purified human brain cell types including neurons, astrocytes, oligodendrocytes and endothelial cells (Zhang et al., 2016). Extended lists of astrocyte, endothelial and neuron markers are identified based on the gene loadings (right) on first and second principal components from PCA plot (left). Reported in bold are mouse cell type specific markers identified as described in Methods (Single Cells section) and the colour code is the same for both figures.
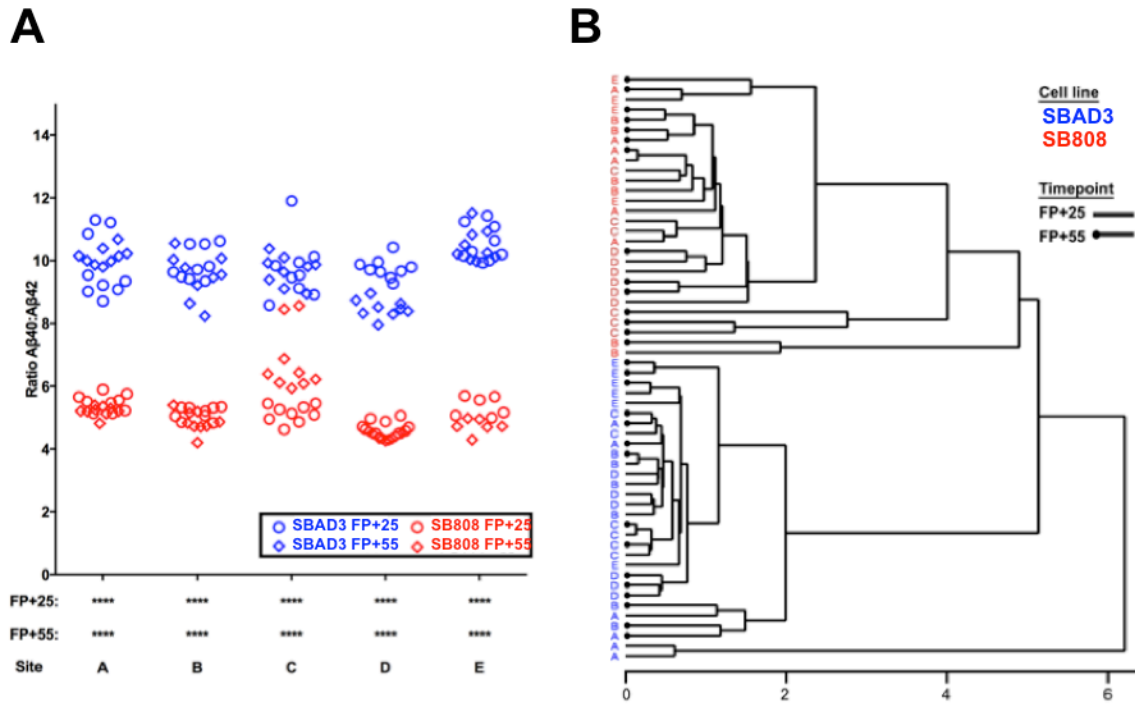
**Figure S5, related to Figure 6. Identification of subpopulation in Single Cell data.**

**A) Cell identity of cell sub-populations in the SB808 (left) and SBAD3 line (right).** To classify which brain cell type is resembled by the different subpopulations, we used a transcriptional purified cell mouse cortex catalogue, containing data for neurons, astrocytes, microglia, endothelial cells, pericytes, and various maturation states of oligodendrocytes (Zhang et al., 2014). From mouse gene expression values, we ranked the genes according to their fold change in each cell type. To evaluate the specificity of each population for a specific cell type, we compared the sum of expression values for top50 of cell specific marker gene with those of random sampling of 50 genes. In this figure, we report the log10 of the empiric p-value associated with an enrichment in markers genes of different cell type in the mouse cortex, in the cell sub-populations of the SB808 and SBAD3

line. **B) Proportion of false negative in term of DE genes due to the cell heterogeneity.** We randomly sampled populations of 100 cells from each line to simulated conditions of cell heterogeneity and performed differential expression analyses for each simulation. By considering 396 SB808 cells and 375 SBAD3 cells, we identified 192 DE genes between SB808 vs SBAD3 line with FDR adjusted p-value less than 1%. We then randomly sampled 100 cells from each line to simulated conditions of cell heterogeneity and performed differential expression analyses for each simulation. This graph represents for each 192 DE genes the fraction of simulations where a positive DE gene is not detected as DE gene (y-axis) according the log10 of p-value associated with differential expression level between SB808 vs SBAD3. We found 10/192 DE genes with proportion of false negative less than 5% that means the cell heterogeneity will affect the sensitivity to detect DE of 80% of positive DE genes. **C) Overlap between the differentially regulated genes in the single cell SB808-vs- SBAD3 subpopulations with the Bulk DE genes**. Down-regulated genes (top) and up-regulated genes (bottom) Red star indicates significance. **D) Differentially-expressed genes and pathways between genotypes vary between iPSC sub-populations.** By considering the SB808and SBAD3 cell populations together, we identified six subpopulations of cells by using unsupervised hierarchical clustering approach on the expression profiles. We detected genes DE between SB808 and SBAD3 line within each subpopulation, excepted for the cluster 5 that did not included no SBAD3 cells. **Correlation of the log fold-change of the differentially-expressed genes between subpopulations (top-left)**. For five sets of DE between SB808 and SBAD3 line in each subpopulation (row), we compared their fold-change in four others populations (column) by using a correlation test based on their fold change in both subpopulation. The value in each cell corresponds to Pearson correlation coefficient between fold-change estimation of set of genes in two subpopulations. The single and double stars indicate when the p-value and q-value are less than 5% respectively. **Overlap between down (top-right) and up (bottom-left) regulated DE genes (SB808 vs SBAD3) between different subpopulations.** For each sub-population, we identified down (lower expression level in SB808 than SBAD3 line) and up (higher expression level in SB808 than SBAD3 line) DE genes. We examined what is fraction (number in each cell) of these down/up (row) DE was detected in others sub-populations (column). The right number associated with each subpopulation gives the number of down and up regulated detected in a given subpopulation. We did not find up-regulated genes in C6. **There is little overlap in the GO Biological Processes associated with differential expressed genes between the two cells in each iPSC sub-populations (bottom-right).** We identified the top10 of Biological Process Gene Ontology pathway associated with DE genes of each iPSC subpopulation. For each topGO pathwa, we examined if they were also enriched in DE in others sub-population. The heatmap plot represents the -log10 of p-value associated with enrichment analyses by considering DE in each of five sub-population for these different topGO pathwayset. **E) Principal component analysis on normalized expression matrix (771 SC libraries x 12835 genes).**

**Figure S6, related to Figure 1: Robust reproducibility of a published disease phenotype in neuronally differentiated iPSC-based model system in multiple laboratories**

**A)** Compared to SBAD3, SB808 neurons display a decreased β-amyloid 40/42 ratio at both time points across all test centers. **B)** Hierarchical clustering based on the normalized proportion of each measured β-amyloid species in the supernatants of cortical iPSC-derived neurons. With very few outliers, samples from multiple test centers cluster by cell line rather than test center, indicating the robust reproducibility of this particular disease signal (Duff et al., 1996; Sproul et al., 2014).

**Supplemental Methods**

**Generation and maintenance of STEMBANCC iPSC lines**

Fibroblasts were transduced with the reprogramming vectors KOS (a polycistronic vector encoding *KLF4, OCT4, SOX2*), h*c-Myc* and h*Klf4* following manufacturer's instructions. One week after transduction, fibroblasts were disaggregated and plated onto feeder layers of mitotically inactivated mouse embryonic fibroblasts in hESC culture medium (KO-DMEM, 20% Knockout™ Serum Replacement, 0.1 mM nonessential amino acids, 2 mM L-glutamine, 100 units/mL penicillin and 8ng/mL human recombinant bFGF, all from Thermo Fisher Scientific) at a density of 8,000 cells per well of a six well plate. The cultures undergoing reprogramming were maintained at $37^{o}$C and 5% $CO_2$ in hESC medium for 2-3 weeks or until colonies with typical hESC morphology appeared. Individual colonies were mechanically dissected and plated onto fresh feeder plates for up to 8 passages before being adapted to feeder free conditions which involved plating on Matrigel-coated plates (Corning, NY) with mTeSR1 media (Stem Cell Technologies, Vancouver, British Columbia, Canada) which was also the medium used in further experiments. A detailed quality control analysis has undertaken in the lines used in the study, which included G-banding, SNP-array-based karyotyping and whole exome sequencing. Mycoplasma testing is routinely performed at each of the participating sites, at the time when they received the cells from the central repository and after passaging. To differentiate iPSCs into cortical neurons, a detailed SOP (see Supplemental Methods) has been used by all 5 partners, based on a previously published method (Shi et al., 2012b, 2012a) Briefly, confluent monolayer iPSCs were induced by dual-SMAD inhibition for 12 days followed by three weeks of progenitor expansion and differentiation. Independent inductions were seeded into 12-well culture dishes at a final plating density of 8.5 x $10^4$ cells /$cm^2$. Samples were collected at 25 and 55 days after final plating (FP+25, FP+55s), and the morphology of the iPSC-derived neurons at 10 days after final plating is depicted in **Figure 1B**.

**Assessment of β-amyloid species in culture supernatants**

48-hour conditioned media collected from triplicate wells was spun at 1000 x g to remove cellular debris and supernatant stored at -80 °C until use. Samples were assayed for soluble APPβ (MesoScale Diagnostics) and soluble β-amyloid-1-38, β-amyloid-1-40 and β-amyloid-1-42 by multiplexed immunoassay (MesoScale Diagnostics).

**<u>Standard Operating Procedure for cortical differentiation of human IPS cells</u>**

Readapted from Shi et al 2012a, 2012b.

| Product name | Supplier | Catalogue number |
|---|---|---|
| BD Matrigel hESC-qualified Matrix | BD Biosciences | 354277 |

| | | |
|---|---|---|
| mTESR 1 | Stem Cell Technlgs | 05850 |
| Geltrex (ready to use) | Life Technologies | A1569601 |
| Y-27632 Rock inhibitor | Cell Guidance Systems | SM02-10 |
| Essential-8 | Life Technologies | A1517001 |
| Ultrapure 0.5M EDTA | Life Technologies | 15575020 |
| Laminin | Sigma | L2020 |
| DMEM:F12 +glutamax | Life Technologies | 31331 |
| Insulin (10mg/ml) | Sigma | I9278 |
| 2-mercaptoethanol (50mM) | Life Technologies | 31350 |
| Non essential amino acids (100x) | Life Technologies | 11140 |
| Sodium Pyruvate (100mM) | Sigma | S8636 |
| Pens/Strep (10000 U/ul) | Life Technologies | 15140 |
| N2 | Life Technologies | 17502048 |
| B27 | Life Technologies | 17504044 |
| L-Glutamine (200mM) | Life Technologies | 25030024 |
| Neurobasal | Life Technologies | 12348 |
| SB431542 | Tocris | 1614 |
| Dorsomorphin | Tocris | 3093 |

**Neural Maintenance media (1L)**

500ml DMEM:F12 +glutamax

0.25ml Insulin

1ml 2-mercaptoethanol

5ml Non essential amino acids

5ml Sodium Pyruvate*

2.5ml Pens/Strep

5ml N2

10ml B27

5ml L-Glutamine* (or glutamax)

500ml Neurobasal

Store at 4°C and use within 3 weeks (*corrected from Shi et al)

**SB431542**

Supplied as a powder, 10mg, MWt (pure compound) = 384.39

Resuspend to 10mM in DMSO, prepare 50ul aliquots and freeze at -20°C

Use at 1:1000 (10uM)

**Dorsomorphin**

Supplied as a powder, 10mg, MWt (pure compound) = 472.41

Resuspend to 10mM in DMSO, dilute 10mM stock to 1mM with ddH2O

Prepare 50ul aliquots and freeze at -20°C

use at 1:1000 (1uM)

**Neural induction media (10ml)**

10ml Neural maintenance media

10ul SB431542

10ul Dorsomorphin

Store at 4°C and use within 5 days

**Procedure**

(**Step numbers** refer to steps in Shi et al 2012a)

- 1.1 **Steps 1-22:** Routine maintenance of hIPSCs. Essential-8 media & geltrex can be substituted for mTESR & Matrigel if preferred.
- 1.2 **Step 23:** Passage cells 2:1 in the presence of Rock Inhibitor. Starting material: 2 wells of a 6 well plate hIPSCs in mTESR (or E8) at 70-90% confluency
  - Pre-coat 1 well of a 6 well plate with 1ml matrigel
  - Pre-warm 0.5mM EDTA to 37°C
  - Aspirate media from two nearly confluent wells of iPSc, wash each with 1ml each of PBS/well (room temp)
  - Aspirate PBS and add 1ml pre-warmed EDTA/well then immediately remove
  - Add 1ml EDTA, incubate 37℃ 4 – 6 min
  - Remove matrigel from coated well (NB don't allow to dry)
  - Check cells have start to detach from each other but not from the plastic
  - Carefully aspirate EDTA and flush the loosened cells with 1ml mTeSR/10uM ROCKi, using a p1000 tip, moving around the well to ensure even flushing. Don't pipette up and down multiple times, as this will result in the patches disintegrating too much.
  - Transfer all 1ml to the new coated well. Repeat with second well of iPSCs and transfer to the new well giving a total vol of 2ml.

- Transfer carefully to incubator, swirling in figure-of-8 to ensure even dispersal of cells

1.3 **Step 24:** Neural induction. (NOTE: inductions have been optimised in 35mm dish or single well of six-well plate). 24hrs after plating, check the cells have reached 100% confluence, wash the cells once with PBS and add 2ml of neural induction medium per well. This is day 0. If the cells are not 100% confluent continue to incubate in mTESR for 1 more day before switching to induction medium. Any gaps in the sheet of cells at this stage will contribute to non-specific differentiation. Refresh induction media daily.

1.4 **Steps 26-31:** On day 12 after induction, the cells should have formed a dense neuro-epithelial sheet (may well appear 'yellow' and 'lumpy'). Passage the cells with dispase as follows:
- Pre-coat 2 wells of a 6 well plate with laminin (1ml per well, 10ug/ml laminin in PBS. Coat at 37°C for 4hrs- overnight)
- Add 200ul dispase stock directly to the 2ml media in the well of the 6-well plate.
- Incubate at 37°C for 3 mins. NOTE: Dispase can be left on as long as 30 mins if sheet is not easy to detach
- Remove cells, keeping sheet as intact as possible by pipetting carefully two or three times from the edge. Clumps should be clearly visible by eye.
- Add 10 ml fresh neural induction medium to a 15ml tube and transfer the clumps into this tube. Allow the clumps to settle in the bottom, then discard the supernatant carefully. Repeat this wash.
- Remove the laminin from the wells, gently resuspend the cells, again without breaking clumps up, in 4ml of neural induction medium and transfer 2ml to each of the two pre-coated laminin wells.
- Incubate the cells overnight to allow the cells to reattach, and change the medium to neural maintenance medium +20ng/ml FGF2 the next day. If the clumps are not attached the following day, they can be transferred to a fresh laminin coated well. Media can now be refreshed at 48hr intervals.

1.5 **Step 34:** After 4 days of FGF treatment, withdraw FGF. Cells can be split 1:2 with dispase when rosettes start to meet, or if neural crest cells begin to appear. Careful dispase passaging should leave non-specific cells attached, and lift off neural rosettes.
Note: There is a possibility that control and disease-specific lines behave differently in terms of differentiation speed, so for the faster line less passages will be needed to reach a stage were the cells are ready for final plating (indicated by the appearance of a critical number of neurons). Careful observation of cell morphology (looking for the appearance of neurons) is crucial to determine how long the first period until final plating should last and this will be likely to be different for the two lines (also reflected by different numbers of dispase and accutase passages for the two lines). This also means that the final plating point (estimated to be D35 but maybe a lot less for a faster growing line) will be considered to be zero, D60 and D90 will be calculated from this time point (e.g +25 and +55 days after this time point). If you prefer to do the final plating at the same time for both lines, the faster growing line can be cryopreserved and thawed later on when the slower growing line is also ready for final plating.

1.6 **Step 42-49:** Passaging to single cells. On day 25 after induction (±1 day), cells can be dissociated with accutase at a ratio of 1:1.
- Precoat well with laminin as above. Remove the medium and wash cells once with 2ml PBS (MgCl2 and CaCl2 free)
- Add 0.5ml Accutase per 35mm well. Incubate the cells in Accutase at 37°C for 5mins.
- Pipette up and down to detach the cells and dilute into 10ml neural maintenance medium. Centrifuge cells at 400g for 5min, repeat the wash and spin, then resuspend in 2ml neural maintenance medium and transfer to laminin coated well.
- Replace the media the day after plating, and every 48hrs subsequently. Cells can be expanded 1:2 when the well reaches 90%-100% confluency (approx. every two to three days)

1.7 **Freezing/Thawing** (optimum stage for freezing is between d28 and d31)
- Following dissociation of the culture with Accutase as described, resuspend cortical stem cells in 1ml neural freezing medium (10% DMSO in neural maintenance media + 20 ng/ml FGF2) per 35mm dish of cells.
- Aliquot 1 ml of the cell suspension into each cryovial.

- Freeze in a CoolCell freezing container at − 80 °C overnight. Transfer the cryovials to liquid nitrogen for long-term storage.
- Thawing NSCs.
- Partially thaw the cells in a 37 °C water bath.
- Transfer the partially thawed NSCs to 10 volumes of room-temperature neural maintenance medium.
- Centrifuge the cells once at 400g for 3 min
- Gently resuspend the cells in 2 ml of neural maintenance medium, and plate into poly-ornithine/laminin-coated 35-mm dishes at 50,000 cells per cm2 (or 1 vial/well). Addition of 20 ng/ml FGF2 to media for the first 12-24hrs after thawing can greatly improve survival.
- Withdraw FGF the following day, and resume culturing of cells as per protocol.

1.8     **Step50, Final plating:** (As neurons are fragile, survival rate after passage is low. For this reason we routinely passage for the final time around day 35). Ideal starting material before final plating is >=60cm2 almost confluent NPs.


Remove the medium and wash cells once with 2ml PBS (MgCl2 and CaCl2 free).
- Add 0.5ml Accutase per 35mm well. Incubate the cells in Accutase at 37°C for 5mins.
- Pipette up and down gently to detach the cells and dilute into 10ml neural maintenance medium.
- Centrifuge cells at 400g for 5min, repeat the wash and spin, then re-suspend in neural maintenance medium.
- Count cells and dilute in neural maintenance media to 3x10^5 cells per ml.
- Remove laminin, and pipette 1ml re-suspended cells per well into 9 wells of each 12w plate. Thus final plating at 300K per well = 85k/cm2
- Leave enough cells to enable quality control with immunochemistry (seed cells on poly/ornithine coated dishes or coverslips to confirm cortical identity, according to your own protocol)

1.9     Replace the media the day after plating, and every 48hrs subsequently, performing full media changes each time, using 1ml media per 12w.


1.10    Add laminin (1/100 in maintenance media, final conc 0.01 mg/ml) every 10 days at d44, d54, d64, d74, d84


1.11    **FP+25:** Sample all media 48hrs after last change (total wells =18, pool into 6 tubes). FP+25 omics samples: row A collect for RNA, row B harvest for protein .


1.12    **FP+55:** Sample all media 48hrs after last change (total wells =9, pool into 3 tubes). FP+55 omics samples: row A collect for RNA, row B harvest for protein .


1.13    **Sample collection.**
- **To sample supernatants for beta-amyloid measurements** (supernatants of the "RNA" samples will be used for MSD measurements)**:**
  - 48 hrs after last media change, remove 1ml media from wells A1, A2 and A3 into three separate collection tubes, spin at 1000g for 5mins to remove cells. Collect the supernatant from each collection tube into lo-bind eppendorf tubes. Store at -80C. Unlike for the supernatants of proteomics and metabolomics wells, for the MSD measurements we keep the three "A" wells separate!
- **To sample supernatants for proteomics**:
  - 48 hrs after last media change, pool 3x1ml media from rows B into two collection tubes, spin them at 1000g for 5 mins to remove cells. Collect 1ml supernatant into lo-bind eppendorf tubes. Store at -80C. From these conditioned media samples.

- **To sample cells for transcriptomics (totalRNA):** (estimate 3ug/well)

- Pre-warm PBS wash at 37C, pre cool PBS collection buffer on ice.
- Wash well with warm PBS. Aspirate off.
- Add 1ml ice cold PBS to well and pipette to detach cells (with scraping if needed, but this is not usually necessary).
- Collect and pool PBS/cells from all three 'A' samples. Spin at >=1000g for 5mins, discard supernatant. Lyse cells in 600ul RLT buffer plus Beta-mercaptoethanol (as per Qiagen RNeasy Mini) and pipet to mix. Homogenize the lysate using QIAshredder spin column, and store homogenized lysate at -80C until needed.
- Continue to purify RNA as per RNeasy protocol, eluting in 30ul RNase free water. Quantify using Qubit RNA broad range assay.

- **To sample cells for proteomics:** (d90 estimate >500ug protein per well)
  - Prepare 37C PBS, (wet)ice cold PBS, and ethanol/dry ice bath
  - Wash well with warm PBS. Aspirate off.
  - Add 1ml ice cold PBS to 'B' wells and pipette to remove cells
  - Aspirate cells and collect in lo-bind tubes. Spin at 1600g for 5mins, discard supernatant and snap cool pellet on ethanol/dry ice bath. Store at -80C

1.14 **QC:**
**For QC purposes, collect phase contrast images of each induction frequently, in particular at neuroepithelial sheet stage, on appearance of rosettes, after dissociation to single cells, and pre/post final plating.**

Immunostaining for TBR1 and CTIP2 after day 45 (FP+10) to confirm cortical identity.

Example immuno protocol:

Fix with 4% PFA 20min RoomTemp

Wash with TBST (1x TBS plus 0.3% Triton) 5min RT

Repeat wash 3x

Block 4% Goat serum/TBST 1hr RT

Primary antibodies overnight 4C in blocking buffer

(Final concentration: ab31940 at 1ug/ml, ab18465 at 1ug/ml)

Wash 3x with TBST 5min RT

Secondary antibodies in TBST 1hr RT (1:500)

Wash 3x with TBST and counterstain with DAPI if required.

1.15 MSD kits used:
K15200G-1 - Aβ Peptide Panel1 (6E10) V-PLEX Plus Kit: 25ul sample

K15120E-1 - sAPPalpha/sAPPbeta Kit: 25ul sample

One neural induction refers to a single well of a 6-well plate at d0. Multiple wells may be derived from one induction but replicates should be always maintained separately after this stage.

## Bulk RNAseq Analysis

### Sample Collection

Pooled cells from three wells were detached in ice-cold PBS and total RNA was extracted using the RNeasy Mini Kit (Qiagen) following the manufacturer's instructions.

### Sequencing, Mapping and gene count estimation:

All sequencing was carried out in single partner centre on an Illumina Hiseq4000 obtaining 75PE reads. Basic quality control screenings on unmapped reads and sequence mapping were performed through CGAT pipeline pipeline_readqc.py. The quality of the sequencing was assessed by FASTQC software (version 0.9.3), (http://www.bioinformatics.babraham.ac.uk/projects/fastqc/). RNA-seq data were mapped to the hg19 assembly via STAR version 2.2.0c (Dobin et al., 2012). Read alignments were merged in single BAM file output per sample (57 in total). Reads were filtered to remove those not uniquely mapped (mapping quality equal to 255) and all ribosomal and mitochondrial RNA reads. Gene-level read counts were obtained using FeatureCount (Liao et al., 2014). Cuffquant and Cuffnorm tools from Cufflinks program (Trapnell et al., 2010) were used to calculate fragments per kilobase per million reads (FPKM) from the merged BAM files. Only for a direct comparison between pooled proteomic data (20 samples) and transcriptomic data (57 samples), BAM files from replicate samples were merged and gene quantification process was repeated as before.

### Data normalisation

Upper quartile (UQ) normalisation has been applied (in-house code) on raw gene counts to correct for library size differences. UQ normalised counts have been used throughout the paper when considering not RUV corrected gene data. UQ normalised counts were also used as input for RUVs correction and differential gene expression analysis.

### Analysis of sources of variation

To minimize the impact of unwanted sources of variation we used methods implemented in the R-Project packages EDASeq(Muschet et al., 2016) and RUVseq(Heywood et al., 2015). The following steps were applied: first raw count data were normalized by upper quartile (UQ). Second, RUVs (remove unwanted variation method) was used to infer factors explaining transcriptome-wide variance components. Information about replicates' structure was given as input to RUVs in order to retain variation coming from the covariates of interest (cell line and time point). Third, these factors were regressed out from the UQ-normalised gene counts (normCounts function used from EDASeq package) and RUVs corrected gene data were used to perform further analyses. The same approach was applied to samples from the control cell line only to expose a "clean" time point signal regardless of any cell line effect.

To explain the variance captured by PCA principal components before and after RUV correction and the RUVs factors in terms of meta-data variables, an over-dispersed Gaussian response model was used. We fitted a Bayesian generalized linear multilevel model using the MCMCglmm R package (Hadfield, 2010). Eleven meta-data variables, or covariates, (SITE, CELL_LINE, TIME_POINT, DETACHMENT, EARLY_HARVEST, CELL_COUNT, MATURITY, CELL_LINE:TIME_POINT, SITE:CELL_LINE:TIME_POINT, SITE:CELL_LINE:INDUCTION, SITE:CELL_LINE:INDEPENDENT) were modelled. Posterior samples of variance proportions were obtained by standardizing the sum of posterior variances across covariates to sum to one. By fitting regression models between the first twenty RUVs factors and the covariates, we were able to estimate proportions of variance captured by RUVs factors and explained in terms of known covariates. Notably, since the factors were estimated based on genotype and time point replicates, the variation coming from these two biological covariates of interest was marginal as expected from our application of RUVs. The residual variance

was generally low in almost all RUVs factors, suggesting that most of the variation captured by RUVs is attributable to the confounding factors we modelled.

The same approach was used to explain variance attributable to each gene (before RUV correction) in terms of meta-data variables by using an over-dispersed Poisson response model. Genes were ranked based on proportion of variance explained by each meta-data variable and the top 100 were used for functional enrichment analysis (GO pathways). To dissect variation captured by any RUVs factor that we found to be explained by SITE-origin, SITE-specific means of proportions of variance were correlated to SITE-specific meta-data variables by fitting linear regression models.

**Differential gene expression analysis**

Within-laboratory differentially expressed (DE) genes were estimated before and after RUVs correction using Limma moderated t-statistic (Limma R package, Ritchie et al., 2015) for either cell line or time point effect correcting for the experimental structure (i.e. design formula defined as "~CELL_LINE + TIME_POINT + CELL_LINE:TIME_POINT") at FDR <= 0.01. Multi-laboratory DE genes were estimated as before but also correcting for laboratory origin before and after RUVs correction. In both cases, UQ normalized data was used as input for Limma.

**Heterogeneity analysis**

To quantify the effect of heterogeneity across laboratories we used the $I^2$ quantity(Higgins et al., 2003) that describes the percentage of total variation across experiments. $I^2$ is calculated based on Cochran's Q(Higgins and Thompson, 2002) test for the null hypothesis that all experiments identify the same effect (cell line and time point effects in the present study). FDR adjusted p-values (Benjamini–Hochberg) are calculated for the test. $I^2$ values range between 0% and 100%, where 0%, 50% and 75% correspond to no, moderate and high heterogeneity, respectively(Higgins et al., 2003). We used voom, lmFit and eBayes functions from Limma R package(Ritchie et al., 2015) to estimate the effects of interest and extract standard deviations for any effect and in any gene. Q and $I^2$ measures were then calculated for either cell line or time point effect correcting for the experimental structure, as described previously (Higgins and Thompson, 2002; Higgins et al., 2003).

**Comparison with public data**

Gene expression profiles were compared to public RNA-Seq data sets (FPKM) from the Genotype-Tissue Expression Portal (GTEx consortium, 2015) [downloaded on September 2015] and the BrainSpan Atlas of the Developing Human Brain (Miller et al., 2014). GTEx data were averaged by tissue and BrainSpan data were averaged by tissue and age to obtain six main age groups (early-prenatal, late-prenatal, infancy, childhood, adolescence and adult). Both data sets were corrected for batch effect, using ComBat function from sva R package(Johnson et al., 2007) , and log-transformed for Principal Component Analysis (PCA) (scaled and centered) and Hierarchical Clustering Analysis (scaled, Euclidean distance, Ward method). CORTECON(van de Leemput et al., 2014) gene clusters for temporal cortex development were downloaded from the database available at http://cortecon.neuralsci.org. Genes assigned uniquely to any stage-specific cluster were used in our enrichment analysis.

**Identification of cortical markers and a transcriptional axis of neuronal maturation from GTEx and BrainSpan datasets**

A set of 787 'cortical marker genes' was identified using GTEx data as genes showing at least five-fold higher RPKM level in three GTEx brain cortical tissues as compared to a group of at least 40 (~80%) "non-cortical" GTEx tissues. A transcriptional maturation axis was then identified on this set of cortical marker genes from PCA of BrainSpan data that clearly cluster by sample age (**Figure S4A left**). The position of STEMBANCC samples projected along this axis (first principal components) was used as covariate named 'MATURITY' in subsequent variance component analysis (**Figure S4A right**).

**Identification of cell type specific markers**

We used RNA-Seq data(Zhang et al., 2016) from purified human brain cell types including neurons, astrocytes, oligodendrocytes and endothelial cells available at http://www.brainrnaseq.org. In a PCA of these data (**Figure S4B left**), principal components 1 and 2 distinguish between three main groups namely neurons and oligodendrocytes, astrocytes and endothelial cells. Given that positive control cell type markers defined as in section '**Brain cell class of different iPSC subpopulations**' (Single Cell Methods) lie correctly along the identified axes of PCA gene loadings (**Figure S4B right**), we derived extended lists of cell type-specific genes contributing to PCA coordinates of the respective cell type groups. We defined different sets of cell type specific genes using different stringency thresholds on the gene loadings. These sets were compared to DE genes from our samples using a hypergeometric test for significant overlap.

**GO pathway enrichment analysis**

We performed a classical enrichment analysis by testing the over-representation of gene ontology biological processes (GO BP) terms within the group of differentially expressed genes using a Fisher test. Semantic similarity between lists of enriched GO terms was calculated using GOSemSim R package(Yu et al., 2010) (Wang method(Wang et al., 2007)).

**Implementing RUVs**

RUVs (Risso et al., 2014) assumes that the biological covariates of interest are constant across replicates. Using a set of negative control samples, or replicates, to column-center the counts RUVs estimates sources of unwanted variation on a set of control genes. In our analysis we used all genes that are expressed (count>=1) across all samples. Information about replicates' structure is given as input to RUVs. A further tuning parameter for RUVs is the number of k estimable factors whose choice should be driven by sample size, extent of technical effects and of differential expression. We estimated a maximum of 20 RUVs factors that we analysed in terms of variance decomposition, improvement of site reproducibility and increase in number of DE genes between covariates of interest. However, when referring to RUV corrected gene counts in all other analyses we intend gene counts normalised on the first five RUVs factors. This corresponds to a first clear clustering of samples on PCA plot based on the two covariates of interest.

**Variance component analysis**

To explain the variance captured by any RUVs factor in terms of meta-data variables we fitted a Bayesian generalized linear multilevel model using the MCMCglmm R package (Hadfield, 2010). Quantitative explanatory covariates were "CELL_NUMB_SCALED" and "Maturity_SCALED" (scaled to zero mean and unit variance).

The "Maturity" covariate consisted of scores on the first principal component of the 57 samples on a set of cortical marker genes to reflect the neuronal maturation axis (**Figure S4A, see Methods**).

Categorical covariates were:

- SITE, modelling variation across sites,
- CELL_LINE, modelling site-homogeneous variation across cell lines,
- TIME_POINT modelling site-homogeneous variation across time points,
- CELL_LINE:TIME_POINT modelling site-homogeneous interaction between time point and cell line
- DETACHMENT,
- EARLY_HARVEST,
- SITE:CELL_LINE:INDUCTION, modelling inter-induction variation, having one level for each of 29 inductions across the whole experiment
- SITE:CELL_LINE:INDEPENDENT_DAY, modelling day-specific inter-induction variation, i.e. attributable to inductions being performed on different days (as was the case at three of the five labs), having one level for each of the 21 (site, cell line, day of induction) triples in the experiment.
- SITE:CELL_LINE:TIME_POINT, modelling site-heterogeneous variation across cell lines and time points, one level for each of the 20 (site, cell line, time point) triples in the experiment.

Each of the covariates SITE:CELL_LINE:INDUCTION, SITE:CELL_LINE:INDEPENDENT_DAY, SITE:CELL_LINE:TIME_POINT was modelled hierarchically (as a "random" effect) with its own variance component. Variance components were allocated non-informative Inverse-Gamma priors (shape = 0.01, rate = 0.01). Other covariates were treated as "fixed" effects, with parameters given non-informative priors of Gaussian distributions with zero mean and standard deviation set at 100 times the SD of the model's dependent variable. The models were fitted using Markov Chain Monte Carlo (MCMC), with samples collected for 500,000 iterations (with a thinning interval of 100) after a burn-in phase of 50,000.

At each thinned MCMC iteration, samples were saved and extracted from the posterior distributions of fixed and random effects and used to estimate the posterior distribution of variance proportions. Variances were extracted as follows: for any particular covariate in the linear model (fixed or random) encoded by the design matrix X and with parameters $\beta$ (with $\beta$ modelled hierarchically in the case of a random covariate), the variance attributable to the effect at MCMC iteration j was quantified as the sample variance of the fitted values at that iteration, i.e. $\text{var}(X\beta^{(j)})$. Posterior samples of variance proportions were obtained by standardizing the sum of posterior variances across covariates to sum to one

**Marker Genes**

The neuron-specific stage markers used to illustrate the relative cell culture heterogeneity before and after RUV (**Figure 4B**) were as follows:

Cortex: FOXG1, FOXP2, LHX2, OTX1, EMX1, OTX2, LHX9, EMX2

Layers: TBR1, OTX1, CTIP2, FEZF2, ETV1/ER81, SATB2, CUX1, RORB, BRN2, FOXG1

Mature: SYP, SLC17A7, DLG4

**Other analyses**

Principal component analysis is performed on log-transformed gene counts through prcomp R function (center=T,scale=T). Hierarchical clustering is performed through hclust R function (on scaled data, Euclidean distance, Ward method). SVD single value decomposition analysis is performed on scaled gene counts through svd R function. Heatmaps are created using pheatmap R function.

<u>**Single Cell RNAseq Analysis**</u>

**Quality control of single cell RNA sequence data**

We checked the quality of RNA sequencing data using the FASTQC software (version 0.9.3) (http://www.bioinformatics.babraham.ac.uk/projects/fastqc/) via the CGAT pipeline pipeline_readqc.py and reported the summary results of FASTQC by plate and by library**.**

**Exons genes annotations file**

We generated annotations within the ENSEMBL gene set after reconciliation with the UCSC genome assembly from human genome (hg19) by using the CGAT pipeline pipeline_annotations.py. The generated gtf file provided the information regarding exon parts of transcripts. This set includes both coding and non-coding transcripts. Coding transcripts span both the UTR and the CDS. We merged this file with ERCC-spike-ins annotations.

**Alignment and quality control of alignment of single cell**

We aligned our RNA sequences to the human genome (hg19) using STAR (version 2.3.0). STAR is a mapper developed for RNA-seq data and is able to ignore adapters by clipping. We generated the index required by STAR using the following options:

--runMode genomeGenerate

--genomeFastaFiles genome softmasked fasta file (hg19)

--sjdbGTFfile gtf containing all known gene models (generated with CGAT pipeline pipeline_annotations.py)

--outFilterType BySJout

 We aligned reads with the CGAT pipeline pipeline_mapping.py (option: make mapping) using STAR default options and:

-- runMode               alignReads

-- genomeLoad               LoaDaNsdRemove

-- outStd          SAM

-- outSAMstrandField          intronMotif

-- outSAMunmapped          Within

-- outFilterType          BySJout

For the batch1, the two bam files of each of the eight libraries coming from two lanes were then merged by using samtools (version 1.8) by running the CGAT pipeline pipeline_mapping.py with the option make mergeBAMFiles. We compiled the statistic regarding the quality of mapping by using make buildBAMStats of CGAT pipeline pipeline_mapping.py.

**Count read overlapping exon annotations and basic Metrics**

Uniquely mapped read pairs were counted using featureCounts  subread-1.5.0 by using the exons annotations generated by approach described above. To evaluate the sequencing output and the amount of usable data, we used some metrics reported by featuresCounts including sequencing depth (Number of sequenced read pairs - Count), percentage of mapped reads, numbers of mapped reads aligning to various annotated genomic features, namely non-exonic coverage (No_feature), exonic-coverage (Ambiguous_mapping, Multiple_mapping and uniquely mapped reads- ENSEMBL_Genes).

We observed a low average proportion of uniquely mapped read (38%). The 5' to 3' coverage plots show mostly uniform coverage of all samples with only a few potentially failed samples (or blanks) that show spiky coverage. The coverage plots do not suggest degradation for most of the samples.

**Quality control and cells filtering**

All QC metrics and plot diagnostic were computed by using the R package scater (1.0.4)  (doi: http://dx.doi.org/10.1101/069633). Firstly we noted that for the batch1 including four plates, the libraries were associated with low libraries complexity measure. We determined whether RNA in each captured cell was degraded by studying the total % of mapped reads compared to the proportion of reads mapped to spike-in molecules.

We excluded the libraries with:

(1) less than 2000 expressed genes
(2) low complexity, where the % of 200 most expressed features (genes and ERCC-spike in) represented more than 50% of total number of counts.
(3) low % of endogenous RNA, for which the % of  ERCC spike-in > 14%
(4) low number of mapped read for which the total counts < $10^6$
(5) bulk libraries that was used as control libraries.

 By using these filters we removed 669 libraries with aberrant patterns from downstream analyses. Please Note that all libraries from batch 1 (batch test) prepared with higher concentrations in ERCC spike-inns and TSO than others batch were discarded.

We performed QC diagnostic at feature level. We observed 1197 features was observed with detectable expression in 50% of libraries and that the top 20 expressed features (including 10 ERRC spike-ins) consumed ~ 25% of reads. We removed the features for which the means of counts where less than one and thus considered 12825 features. We then computed the normalized log-expression values with R package scran (1.0.4) by adding one to each count, dividing by the size factor for that cell, and log transforming.

**Important explanatory variables**

From normalized expression values matrix (771 single cell libraries x 12835 genes), we then identified variables that drive variation in expression data across cells by using a linear model for each cell feature and by plotting the distribution of their marginal $R^2$ values. As expected we found that the read depth (total features) and the library complexity (% counts top 100 features) were the two most important explanatory variables. We observed that the site origin was also an important explanatory variable (> 20% of variance) suggesting that two sites produced different cell types (**Figures S5E**). Finally, we noted that the variables batch (2,3,4) and genotype explained a low fraction of the total variance. Since we wanted to identify the causes explaining the cell population differences between sites, we did not correct for these effects.

**Identification of iPSC subpopulations: clustering analyses**

We performed hierarchical clustering on the Euclidean distances between cells, using Ward's criterion to minimize the total variance within each cluster. This yields a dendrogram that groups together cells with similar expression patterns across the chosen genes. Clusters are explicitly defined by applying dynamic tree cut to the dendrogram(Langfelder et al., 2008) (R package dynamicTreeCut). This approach exploits the shape of the branches in the dendrogram to refine the cluster definitions. We identified four, five and six cells subpopulations in SB808, SBAD3 and SB808 and SBAD3 together respectively (**Figure 6A**)

**Gene markers between iPSC subpopulations**

Markers of specific subpopulations were identified by looking at genes that are consistently differentially expressed in the largest subpopulation compared to the others (cluster 1 vs others clusters). DE analyses were performed with the R package edgeR (3.14.0)(Robinson et al., 2010) that uses negative binomial (NB) distributions to model the read counts for each sample. We estimated the NB dispersion parameter that quantifies the biological variability in expression across cells in the same cluster. Large dispersion estimates (> 0.5) were observed due to technical noise with single cell RNA-seq data (in contrast to bulk data where values of 0.05–0.2 are more typical). We then used the design matrix to fit a NB General Linear Model to the counts for each gene. Finally, the top ten sets of DE genes from each pairwise comparison were considered to be as potential marker genes separating subpopulations (**Table S5**).

**GO pathway associated with markers genes between iPSC subpopulations**

We performed gene set enrichment analysis with topGO R package (2.24.0) by considering top 50 marker genes identified above. We performed a classical enrichment analysis by testing the over-representation of GO BP terms within the group of differentially expressed genes by Fisher test. We then listed the top 20 significant GO BP terms identified (**Table S6**).

**Brain cell class of different iPSC subpopulations**

To classify which brain cell type is resembled by the different subpopulations, we used a public database containing transcriptional datasets of purified cortical cells, such as neurons, astrocytes, microglia, endothelial cells, pericytes, and various maturation states of oligodendrocytes(Zhang et al., 2014). Mouse data, instead of human data as used in bulk analysis, has been chosen for its larger number of samples available for each cell type. From Fragments Per Kilobase of transcript per Million (FPKM) mouse gene expression values (http://web.stanford.edu/group/barres_lab/brain_rnaseq.html) we ranked the genes according to their fold change in each cell type. To evaluate the specificity of each population for a specific cell type, we compared the sum of expression values for top 50 cell specific marker genes with those of random sampling of 50 genes (**Figure S5A**).To simplify the heatmap, we display only the top 20 marker genes for each cell type that were expressed in 25% single cell libraries (**Figure 6A**).

**Comparison of differentially-expressed genes and pathways between genotypes between iPSC subpopulations**

We performed differential expression analyses between the two line within each subpopulation. As cluster 5 included no control iPSC cells, we excluded this cluster from downstream analyses. We performed differential expression analyses with edge R with the same approach as described above to identify marker genes for the subpopulations by considering DE genes based on empiric p-value < 0.001. We then examined the overlap between up- or down-regulated genes among different iPSC subpopulations (**Figure S5D**). Further we calculated the correlation coefficients and their significance values for log fold change between DE of different iPSC subpopulations (**Figure S5D**). We reported DE genes by iPSC subpopulations in **Table S3**.

To ensure that there was no overlap between GO pathways associated with distinct sets of DE genes of each iPSC subpopulation, we compared the p-values of enrichment analyses for the top 10 BP GO terms associated to DE genes in each iPSC subpopulations. GO enrichment analysis was performed as described above using topGO R package.

**Proteomics**

**Materials and Reagents**

All materials were of analytical and mass-spectrometry grade. DL-dithiothreitol (DTT), iodoacetamide, ASB-14, Tris base and urea were all purchased from Sigma-Aldrich. UPLC-MS grade acetonitrile (ACN), Formic acid (FA) and water were obtained from Fluka, and sequencing-grade modified porcine trypsin from Promega. All buffers and solutions were prepared using ultra-pure 18 MΩ water (MilliQ) and UPLC solvents using UPLC-MS grade water.

**In-solution proteolytic digestion**

Lysate and supernatant were stored at -80 °C until use. The frozen cell pellets (about 2 million cells) were thawed, dissolved in lysis buffer containing 100mM Tris HCl pH 7.8, 6M Urea, 2M Thiourea, 2% ASB-14. Cells were then sonicated for 5 min to disrupt cell membranes followed by shaking for 1h in room temperature in order to solubilize proteins. Samples were reduced with the addition of 1.5 µl of DTT (30 mg/mL) for 1h at room temperature and then alkylated with 3 µL of iodoacetamide (36 mg/mL) for 30 min in the dark. To dilute the urea, 155 µL of ultra-pure water was added prior to addition of 10 µL of sequencing-grade trypsin (Promega) (0.1 µg/µL). Trypsin digestion was carried out for 12 hours at 37°C, followed by desalting and concentrated using C18 Isolute columns (Biotage).

**Label-free proteomic analyses (2D-LC-MSe)**

1 pmol of yeast enolase reference standard (Waters) was added to the each sample and 1 µg of the sample analysed using a 2D-LC-nanoESI-MSe using a nanoAcquity UPLC 2D-LC system and Synapt G2-Si mass spectrometer (Waters, Manchester, UK). Peptides were fractionated in 8 steps (8.7%, 11.8%, 13.6%, 15.3%, 17.1%, 19.3%, 22.5% and 50% of solvent B) under high pH conditions using XBridge BEH C18 Trap Column, 130Å, 5 µm, 180 µm x 50 mm (Waters). Solvent A was 20 mM ammonium formate in water (pH 9) and solvent B was 100% acetonitrile (ACN). Eluted peptides were trapped in the Symmetry C18, 100Å 5 µm, 300 µm x 20 mm (Waters) trap column and subsequently separated under low pH conditions on a nanoAcquity C18 column Peptide BEH C18, 130Å, 1.7 µm, 75 µm x 150 mm (Waters). Analytical chromatography was performed using a 60-min gradient starting at 97% solvent A, ramping to 40% solvent B in 40 min, then to 85% solvent B over 2 min (held for 3 min) and finally decreased to 97% solvent A in 15 min. Mobile phase A contained 95% $H_2O$, 5% DMSO, 0.1% FA and mobile phase B was 95% ACN, 5% DMSO and 0.1% FA.

**Principal Component Analysis**

Protein abundances of 1037 proteins, which are expressed across all samples, are corrected for batch effect using ComBat function from sva R package(Leek et al., 2012) (two batches done separately for the two time points). PCA is performed by prcomp R function (scale=T, center=T).

PCA of protein abundances and gene counts at FP+55 time point is performed on the same set of 1037 proteins/genes after correcting for differences between the two types of data using ComBat function.

**References**

Consortium, GTEx. (2015). The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue gene regulation in humans. Science *348*, 648–660.

Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* bts635 (2012). doi:10.1093/bioinformatics/bts635

Duff, K., Eckman, C., Zehr, C., Yu, X., Prada, C.M., Perez-tur, J., Hutton, M., Buee, L., Harigaya, Y., Yager, D., et al. (1996). Increased amyloid-beta42(43) in brains of mice expressing mutant presenilin 1. Nature *383*, 710–713.

Hadfield, J.D. (2010). MCMC Methods for Multi-Response Generalized Linear Mixed Models: The MCMCglmm *R* Package. J. Stat. Softw. *33*.

Heywood, W.E., Galimberti, D., Bliss, E., Sirka, E., Paterson, R.W., Magdalinou, N.K., Carecchio, M., Reid, E., Heslegrave, A., Fenoglio, C., et al. (2015). Identification of novel CSF biomarkers for neurodegeneration and their validation by a high-throughput multiplexed targeted proteomic assay. Mol. Neurodegener. *10*.

Higgins, J.P.T., and Thompson, S.G. (2002). Quantifying heterogeneity in a meta-analysis. Stat. Med. *21*, 1539–1558.

Higgins, J.P.T., Thompson, S.G., Deeks, J.J., and Altman, D.G. (2003). Measuring inconsistency in meta-analyses. BMJ *327*, 557–560.

Johnson, W.E., Li, C., and Rabinovic, A. (2007). Adjusting batch effects in microarray expression data using empirical Bayes methods. Biostat. Oxf. Engl. *8*, 118–127.

Langfelder, P., Zhang, B., and Horvath, S. (2008). Defining clusters from a hierarchical cluster tree: the Dynamic Tree Cut package for R. Bioinforma. Oxf. Engl. *24*, 719–720.

Leek, J.T., Johnson, W.E., Parker, H.S., Jaffe, A.E., and Storey, J.D. (2012). The sva package for removing batch effects and other unwanted variation in high-throughput experiments. Bioinformatics *28*, 882–883.

Liao, Y., Smyth, G. K. & Shi, W. (2014). featureCounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinforma. Oxf. Engl.* **30,** 923–930.

Miller, J.A., Ding, S-L., Sunkin, S.M., Smith, K.A., Ng, L., Szafer, A., Ebbert, A., Riley, Z., Royall, J., Aiona, K., et al. (2014). Transcriptional landscape of the prenatal human brain. Nature. *508*, 199-206

Muschet, C., Möller, G., Prehn, C., Hrabě de Angelis, M., Adamski, J., and Tokarz, J. (2016). Removing the bottlenecks of cell culture metabolomics: high-throughput normalization procedure, correlation of metabolites to cell number, and impact of the harvesting method. (Submitted).

Risso, D., Ngai, J., Speed, T.P., and Dudoit, S. (2014). Normalization of RNA-seq data using factor analysis of control genes or samples. Nat. Biotechnol. *32*, 896–902.

Ritchie, M.E., Phipson, B., Wu, D., Hu, Y., Law, C.W., Shi, W., and Smyth, G.K. (2015). limma powers differential expression analyses for RNA-sequencing and microarray studies. Nucleic Acids Res. *43*, e47.

Robinson, M.D., McCarthy, D.J., and Smyth, G.K. (2010). edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. Bioinformatics *26*, 139–140.

Shi, Y., Kirwan, P., Smith, J., Robinson, H.P.C., and Livesey, F.J. (2012a). Human cerebral cortex development from pluripotent stem cells to functional excitatory synapses. Nat. Neurosci. *15*, 477–86, S1.

Shi, Y., Kirwan, P., and Livesey, F.J. (2012b). Directed differentiation of human pluripotent stem cells to cerebral cortex neurons and neural networks. Nat. Protoc. *7*, 1836–46.

Sproul, A.A., Jacob, S., Pre, D., Kim, S.H., Nestor, M.W., Navarro-Sobrino, M., Santa-Maria, I., Zimmer, M., Aubry, S., Steele, J.W., et al. (2014). Characterization and Molecular Profiling of PSEN1 Familial Alzheimer's Disease iPSC-Derived Neural Progenitors. PLoS ONE *9*, e84547.

Trapnell, C. *et al.* Transcript assembly and abundance estimation from RNA-Seq reveals thousands of new transcripts and switching among isoforms. *Nat. Biotechnol.* **28,** 511–515 (2010).

van de Leemput, J., Boles, N.C., Kiehl, T.R., Corneo, B., Lederman, P., Menon, V., Lee, C., Martinez, R.A., Levi, B.P., Thompson, C.L., et al. (2014). CORTECON: A Temporal Transcriptome Analysis of In Vitro Human Cerebral Cortex Development from Human Embryonic Stem Cells. Neuron *83*, 51–68.

Wang, J.Z., Du, Z., Payattakool, R., Yu, P.S., and Chen, C.-F. (2007). A new method to measure the semantic similarity of GO terms. Bioinformatics *23*, 1274–1281.

Yu, G., Li, F., Qin, Y., Bo, X., Wu, Y., and Wang, S. (2010). GOSemSim: an R package for measuring semantic similarity among GO terms and gene products. Bioinforma. Oxf. Engl. *26*, 976–978.

Zhang, Y., Chen, K., Sloan, S.A., Bennett, M.L., Scholze, A.R., O'Keeffe, S., Phatnani, H.P., Guarnieri, P., Caneda, C., Ruderisch, N., et al. (2014). An RNA-sequencing transcriptome and splicing database of glia, neurons, and vascular cells of the cerebral cortex. J. Neurosci. Off. J. Soc. Neurosci. *34*, 11929–11947.

Zhang, Y., Sloan, S.A., Clarke, L.E., Caneda, C., Plaza, C.A., Blumenthal, P.D., Vogel, H., Steinberg, G.K., Edwards, M.S.B., Li, G., et al. (2016). Purification and Characterization of Progenitor and Mature Human Astrocytes Reveals Transcriptional and Functional Differences with Mouse. Neuron *89*, 37–53.