

# SCIENTIFIC DATA

## OPEN Data Descriptor: Single cell RNA sequencing of stem cell-derived retinal ganglion cells

Maciej Daniszewski<sup>1,2,\*</sup>, Anne Senabouth<sup>3,\*</sup>, Quan H. Nguyen<sup>3</sup>, Duncan E. Crombie<sup>1,2</sup>, Samuel W. Lukowski<sup>3</sup>, Tejal Kulkarni<sup>1,2</sup>, Valentin M. Sluch<sup>4</sup>, Jafar S. Jabbari<sup>5</sup>, Xitiz Chamling<sup>4</sup>, Donald J. Zack<sup>4,6</sup>, Alice Pébay<sup>1,2,†</sup>, Joseph E. Powell<sup>3,7,†</sup> & Alex W. Hewitt<sup>1,2,8,†</sup>

Received: 9 October 2017

Accepted: 22 December 2017

Published: 13 February 2018

We used single cell sequencing technology to characterize the transcriptomes of 1,174 human embryonic stem cell-derived retinal ganglion cells (RGCs) at the single cell level. The human embryonic stem cell line BRN3B-mCherry (A81-H7), was differentiated to RGCs using a guided differentiation approach. Cells were harvested at day 36 and prepared for single cell RNA sequencing. Our data indicates the presence of three distinct subpopulations of cells, with various degrees of maturity. One cluster of 288 cells showed increased expression of genes involved in axon guidance together with semaphorin interactions, cell-extracellular matrix interactions and ECM proteoglycans, suggestive of a more mature RGC phenotype.

Design Type(s)	transcription profiling by high throughput sequencing design
Measurement Type(s)	transcription profiling assay
Technology Type(s)	RNA sequencing
Factor Type(s)	Cell Surface Antigen
Sample Characteristic(s)	Homo sapiens • embryonic stem cell • H7-hESC

<sup>1</sup>Centre for Eye Research Australia, Royal Victorian Eye and Ear Hospital, East Melbourne 3002, Australia. <sup>2</sup>Ophthalmology, Department of Surgery, the University of Melbourne, Melbourne 3002, Australia. <sup>3</sup>Institute for Molecular Bioscience, University of Queensland, St Lucia, Brisbane 4067, Australia. <sup>4</sup>Wilmer Eye Institute, Johns Hopkins University School of Medicine, Baltimore, MD 21287, USA. <sup>5</sup>Australian Genome Research Facility, Melbourne 3051, Australia. <sup>6</sup>Departments of Neuroscience, Molecular Biology and Genetics, and Institute of Genetic Medicine, Johns Hopkins University School of Medicine, Baltimore, MD 21287, USA. <sup>7</sup>Queensland Brain Institute, University of Queensland, St Lucia, Brisbane 4072, Australia. <sup>8</sup>School of Medicine, Menzies Institute for Medical Research, University of Tasmania, Hobart 7000, Australia. \*These authors contributed equally to this work. †These authors jointly supervised this work. Correspondence and requests for materials should be addressed to A.P. (email: apebay@unimelb.edu.au) or to J.E.P. (email: j.powell@imb.uq.edu.au) or to A.W.H. (email: hewitt.alex@gmail.com).

## Background & Summary

Since the isolation of embryonic stem cells (ESCs)<sup>1–4</sup> and generation of induced pluripotent stem cells (iPSCs)<sup>5,6</sup>, pluripotent stem cells (PSCs) have made a tremendous contribution towards improving our understanding of mechanisms involved in development and disease. PSCs have the ability to self-renew and differentiate into all cell types of the body, thereby providing great potential for regenerative medicine and cell replacement therapies. Further, PSC-derived progeny allow the investigation of disease-affected cell types that are not readily accessible due to their anatomical location, such as retinal cells<sup>7–9</sup>. Utilising such disease-affected cells will also significantly improve the drug development pipeline through efficacy profiling and side effect or toxicity assessment<sup>10</sup>.

The development of RNA sequencing (RNA-seq) technology has allowed for the rapid quantification of individual gene transcripts. Integrating this high-throughput data with computational and statistical methods provides a toolbox to study the molecular functions of human tissues. Critically, to date, the majority of RNA-seq studies have been conducted on ‘bulk’ samples, consisting of millions of individual cells—the result of which is that transcript quantification represents the average signal across the cell population being studied. Recent developments to isolate single cells, and barcode their expressed transcripts has enabled the transcriptomes of single cells to be sequenced (scRNA-seq) in a high-throughput manner. By sequencing large number of single cells from an individual ‘sample’ it is now possible to dissect the cellular composition of apparently homogenous tissues or cell cultures<sup>11–13</sup>. sc-RNAseq also opens the possibility of examining rare cell populations that could not otherwise be resolved using bulk RNA-seq, and further characterising well-known cell types, for example oligodendrocytes<sup>14</sup> or sensory neurons<sup>15</sup>. Moreover, scRNA-seq may also be used for tracking cell lineage during differentiation, as movement between different cell types is associated with changes in gene expression. Thus, stages across a cell lineage can be distinguished by their unique transcriptional signature<sup>16</sup>. This technology has also been used in cell culture, in particular with PSCs, their differentiated progeny and organoids, including of the nervous system<sup>17,18</sup>, as a way to distinctively characterize cellular subpopulations. Results of such analyses can discern determinants of cell fates, and this information can then applied to *in vitro* differentiation experiments to increase efficiency of generating the tissue of interest<sup>19</sup>.

Retinal ganglion cells (RGCs) transmit visual information from the retina to the midbrain through the optic nerve. Many diseases, such as primary open angle glaucoma, Leber hereditary optic neuropathy and autosomal dominant optic atrophy manifest by degeneration or loss of RGCs and culminate in irreversible loss of sight. It is estimated that there are more than 30 subtypes of RGCs in the mammalian retina<sup>20</sup>, however, the molecular profiling of RGCs in human normal and disease tissue has proven difficult. Currently, studying optic neuropathies is hindered by the lack of non-invasive means for obtaining RGCs from living donors. This can now be circumvented by use of PSCs as a source of RGCs<sup>8</sup>. We recently described a protocol for the differentiation of human PSCs into functional RGCs<sup>7</sup>. RGCs generated through this method are functional, as exemplified by the presence of sodium and potassium currents, mature axon potentials and the expression of RGC-specific markers, including *BRN3B*, *ISL1* and *PRPH*<sup>7</sup>. Moreover, whole transcriptome analysis through bulk RNA-seq of our PSC-derived RGCs demonstrated close resemblance to sensory neurons, and cells from the ganglion cell layer<sup>7</sup>. Herein we present a dataset of scRNA-seq to characterize the transcriptome of RGCs derived from human ESCs (hESCs) at a single cell level.

## Methods

### Ethical approval

All experimental work performed in this study was approved by the Human Research Ethics committees of the University of Melbourne (0605017) with the requirements of the National Health & Medical Research Council of Australia (NHMRC) and conformed with the Declarations of Helsinki.

### Cell culture and retinal differentiation

The hESC *BRN3B*-mCherry reporter line (A81-H7)<sup>9</sup> was maintained on vitronectin-coated 6-well plates using StemFlex (Gibco). Culture medium was changed every second day. Cells were differentiated into RGCs as previously described<sup>7</sup>. Briefly, undifferentiated hESCs cultured in monolayer on vitronectin-coated plates were differentiated using RGC differentiation medium 2 (DMEM F12 with GlutaMAX, 10% KnockOut Serum Replacement (Invitrogen), SM1 (Stem Cell Tech), 10 ng/ml noggin (Sapphire Biosciences), 10 ng/ml Dickkopf-related protein 1 (DKK1, Peprotech), 10 ng/ml Insulin Growth Factor 1 (IGF1, Peprotech) and 5 ng/ml basic Fibroblast Growth Factor (bFGF, Merck). Medium was changed every 2–3 days. RGC differentiation was monitored by the appearance of mCherry-positive cells, reflective of *BRN3B* expression.

### Fluorescence-activated cell sorting (FACS)

On day 36 of differentiation, cells were washed with phosphate-buffered saline (PBS) and incubated with Accutase (Sigma, 37 °C, 5 min). Cells were then incubated in RGC differentiation medium supplemented with the ROCK inhibitor Y27632 (10 μM, Selleckchem, RGC+RI) and gently dissociated using a P1000 pipette, filtered using a 100 μm nylon strainer (BD Falcon) and centrifuged (300 g, 10 min). The cell pellet was resuspended in RGC+RI medium and incubated with THY1 antibody (Human THY1 FITC conjugated, Miltenyi, 130-095-403, 4 °C, 15 min). Cells were washed in RGC+RI medium, and

centrifuged (300 g, 3 min). Two modifications to our original protocol were performed. Firstly, selection of RGCs using THY1 was performed by FACS instead of the magnetic sorting we originally reported. Secondly, cells were prepared for sequencing immediately following THY1 selection and were not allowed to rest prior to being further processed. A cell pellet was resuspended in 500  $\mu$ l of RGC+RI prior to sorting with a BD FACSAria III cell sorter (Becton, Dickinson). Both THY1-positive (+ve) and THY1-negative (-ve) fractions were collected in 5 ml conical tubes (BD Falcon).

### Single-cell preparation

Both THY1-positive (+ve) and THY1-negative (-ve) fractions were subjected to library preparation using the Single Cell 3' Reagent Kit (10X Genomics) as per the manufacturer's instruction. This step was performed within 60 min of the FACS. Briefly, cell suspension was mixed using a wide-bore tip to determine cell concentration using a Countess Automated Cell Counter (Life Technologies). Cells were centrifuged for 5 min at 300 g and the cell pellet was resuspended in PBS with 0.04% BSA. The cell suspension was passed through a cell strainer to remove any remaining cell debris and large clumps and the cell concentration was determined again.

### Generation of single cell GEMs and sequencing libraries

Single cell suspensions were loaded onto 10X Genomics Single Cell 3' Chips along with the reverse transcription (RT) master mix as per the manufacturer's protocol for the Chromium Single Cell 3' v2 Library (10X Genomics; PN-120233), to generate single cell gel beads in emulsion (GEMs). Sequencing libraries were generated with unique sample indices (SI) for each sample. The resulting libraries were assessed by gel electrophoresis (Agilent D1000 ScreenTape Assay) and quantified with qPCR (Illumina KAPA Library Quantification Kit). Following normalization to 2 nM, libraries were denatured and diluted to 17pM of cluster generation using the Illumina cBot (HiSeq PE Cluster Kit v4). Libraries for the two samples were multiplexed respectively, and sequenced on an Illumina HiSeq 2500 (control software v2.2.68/ Real Time Analysis v1.18.66.3) using a HiSeq SBS Kit v4 (Illumina, FC-401-4003) in high-output mode as follows: 126 bp (Read 1), 8 bp (i7 Index), 8 bp (i5 Index), and 126 bp (Read 2).

### Mapping of reads to transcripts and cells

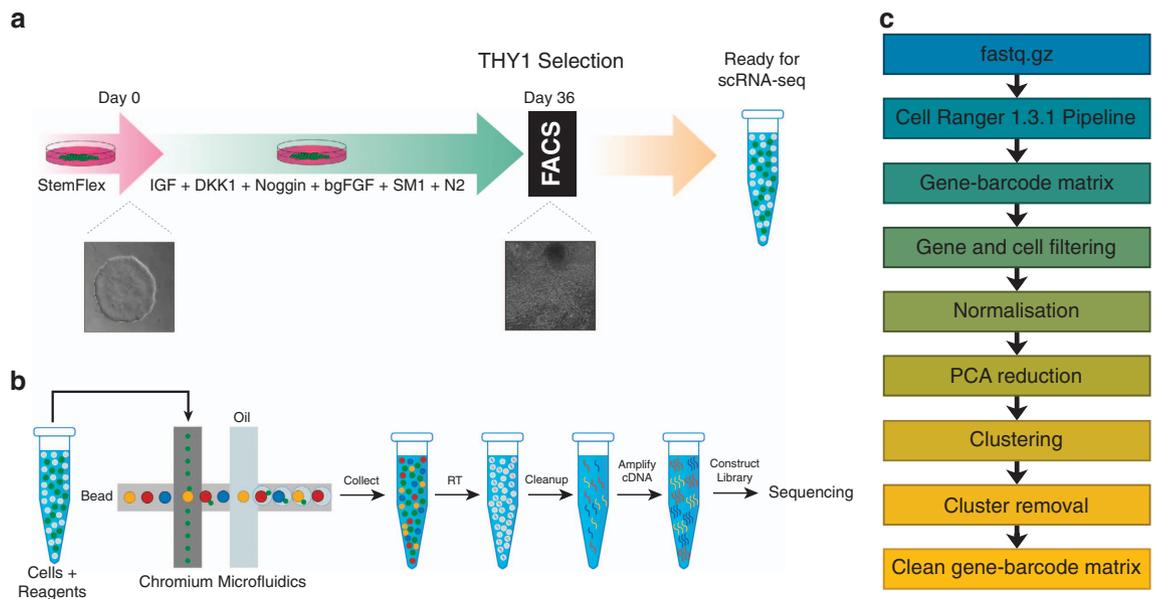
The sequencing data was processed into transcript count tables with the Cell Ranger Single Cell Software Suite 1.3.1 by 10X Genomics (<http://10xgenomics.com/>). Raw base call files from the HiSeq2500 sequencer were demultiplexed with the *cellranger mkfastq* pipeline into library-specific FASTQ files. As the libraries were sequenced using non-standard settings, *cellranger mkfastq* was run with the following parameters: `--use-bases-mask = "Y26n*,I8n*,n*,Y98n*" --ignore-dual-index`. The FASTQ files for each library were then processed independently with the *cellranger count* pipeline. This pipeline used STAR<sup>21</sup> to align cDNA reads to the Homo sapiens transcriptome (Sequence: GRCh38, Annotation: Gencode v25). Once aligned, barcodes associated with these reads – cell identifiers and Unique Molecular Identifiers (UMIs), underwent filtering and correction. Reads associated with retained barcodes were quantified and used to build a transcript count table. Resulting data for each sample were then aggregated using the *cellranger aggr* pipeline, which performed a between-sample normalization step and concatenated the two transcript count tables. Post-aggregation, the mapped data was processed and analyzed as described below.

### Preprocessing

To preprocess the mapped data, we constructed a cell quality matrix based on the following data types: library size (total mapped reads), total number of genes detected, percent of reads mapped to mitochondrial genes, and percent of reads mapped to ribosomal genes. Cells that had any of the 4 parameter measurements lower than 3x median absolute deviation (MAD) of all cells were considered outliers and removed from subsequent analysis (Table 1)<sup>22</sup>. In addition, we applied two thresholds to remove cells with mitochondrial reads above 20% or ribosomal reads above 50% (Table 1). To exclude genes that were potentially detected from random noise, we removed genes that were detected in fewer than 1% of all cells. The expression data was normalised on two levels to reduce possible systematic bias between samples and between cells. The first level of normalisation - between samples, was performed prior to data aggregation using the *cellranger aggr* depth equalisation method<sup>23</sup>. This method reduces potential confounding effects caused by differences in sequencing depths between samples by subsampling mapped reads from higher-depth libraries until the number of mapped reads per library were equal. The second level of normalisation - between cells, was performed after filtering using the deconvolution approach by Lun *et al.*<sup>24</sup>. This level of normalisation reduces bias possibly caused by technical variation such as cDNA synthesis, PCR amplification efficiency and sequencing depth for each cell. The deconvolution approach was chosen as it accounts for the sparse nature of expression data by pooling expression counts from groups of cells. As the sizes of the groups were linear (40, 60, 80, 100), the group-specific normalised size factors could be deconvolved into cell-specific size factors that were then used to scale the counts of individual cells. After normalisation, abundantly expressed ribosomal protein genes and mitochondrial genes were discarded. We have made available both the raw and normalised data (Data Citation 1).

Sample	Number of cells	Median reads per cell	Median genes per cell	Total genes detected	Median UMIs per cell	Total number of reads	Percent mapped reads	Remaining cells post filtering
1	1,090	124,127	3,528	21,317	13,575	135,299,096	61.10	993
2	194	493,659	5,188	19,812	26,218	95,769,921	62.80	181

**Table 1.** Summary statistics for sequencing and mapping data of two samples.

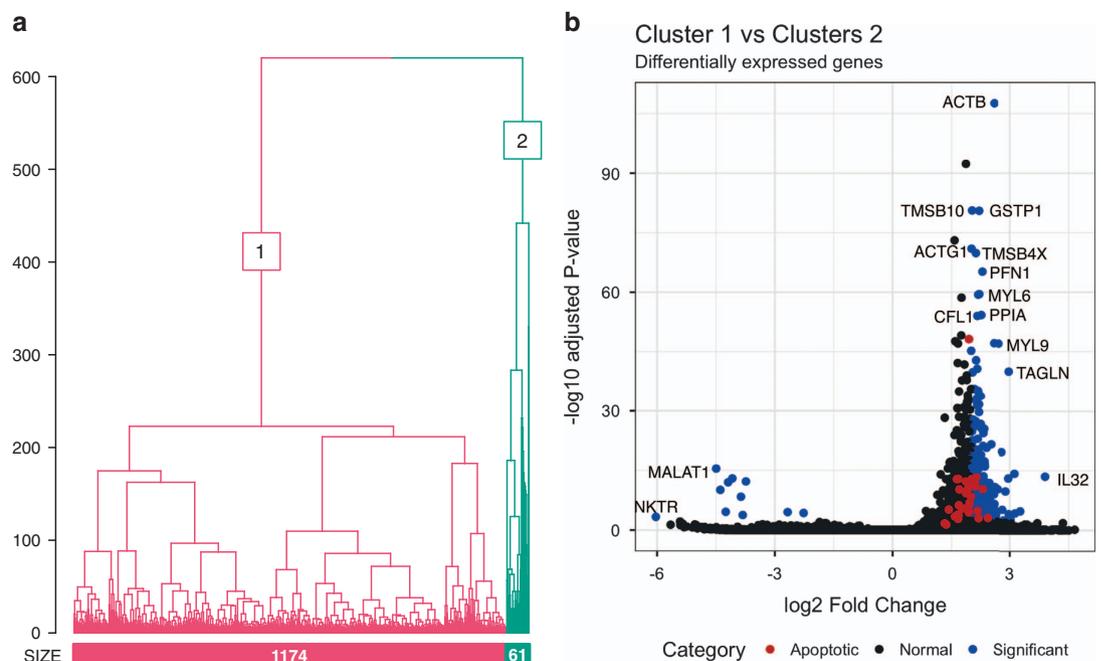


**Figure 1.** Schematic representation of the experimental workflow. **(a)** Guided differentiation of the reporter line BRN3B-mCherry A81-H7 hESCs into RGCs using IGF1, DKK1, Noggin, bFGF in a neural medium containing SM1 and N2, as described in Ref. 7. On day (d) 36, cells were sorted based on the expression of the marker THY1. Cells from both positive and negative THY fractions were then processed for scRNA-seq. Brightfield images describe cell morphology of undifferentiated hESCs prior to differentiation (d0) and post differentiation (d36) at time of sorting. **(b)** Single cell suspensions are prepared and libraries generated using the Chromium V2 chemistry. Libraries were sequenced on an Illumina HiSeq2500. **(c)** Sequence data is processed using bioinformatic pipelines, and analysis conducted on the resulting expression matrix.

### Identification of residual low-quality cells via clustering

We identified and removed a small group of cells with low-quality sequence data. These cells were not detected by initial filtering; instead, they were identified via clustering and enrichment of differentially expressed genes. The transcript count table underwent dimensionality reduction using Principal Component Analysis (PCA). This procedure was applied to the top 1,500 most variable genes using the *prcomp()* function in R<sup>25</sup>. The first 20 PCs were retained and a cell-PCA eigenvector matrix was used for clustering.

We applied an unsupervised clustering method that does not take into account any predetermined parameters to objectively identify single cell subpopulations<sup>26</sup>. This method is less biased compared to top-down clustering approaches, such as k-means. Briefly, to achieve high-resolution clustering capable of detecting small subpopulations and outliers, we applied bottom-up agglomerative hierarchical clustering to construct a dendrogram tree, where the highest resolution is one cell = one bottom branch. We used the reduced dataset containing the top 20 PCs described above to calculate an Euclidean distance matrix between cells, and organized cells into the dendrogram using the Ward's minimum distance so that similar cells are joint into larger groups of branches. To identify subpopulations, we applied an unsupervised, objective approach to merge the branches into a high-resolution and stable clustering result. The approach divided the dendrogram tree into 40 height-windows, ranging from 0.025 (from the bottom of the tree) to 1 (from the top). By iteratively and dynamically merging cells in each of the 40 height-windows, we generated 40 independent clustering results with varying resolutions. Clustering results were then compared quantitatively using adjusted Rand indices, which score pairs of cells that are the same or different between two clustering results<sup>27</sup>. The optimal clustering result was the most stable result across a range of consecutive tree-height values.



**Figure 2. Identification of residual low-quality cells via clustering.** (a) Unsupervised clustering of all cells into two subpopulations. The dendrogram tree displays distance and agglomerative clustering of the cells. Each branch represents one subpopulation. The clustering is based on the most stable clustering result across 40 tree cut heights. The branches are labelled with their subpopulation identification. The number of cells in each of the two populations are given below the branches. (b) The top significantly expressed genes at the initial filtering step between Cluster 1 and 2 shows upregulation of genes associated with apoptosis in Cluster 2.

To characterise the identified clusters, we performed pairwise differential expression analysis by fitting a general linear model and using a negative binomial test as described in the DESeq package<sup>28</sup>. Network analysis was then performed on significant differentially-expressed genes using Reactome functional interaction analysis<sup>28,29</sup>.

### Code availability

All code and usage notes are available at: <https://github.com/IMB-Computational-Genomics-Lab/RetinaGanglionCells>. This includes: computational bioinformatic pipelines that process sequence data in BCL format through to a mapped UMI expression matrix; scripts for quality-control, normalisation, clustering, differential expression and visualization.

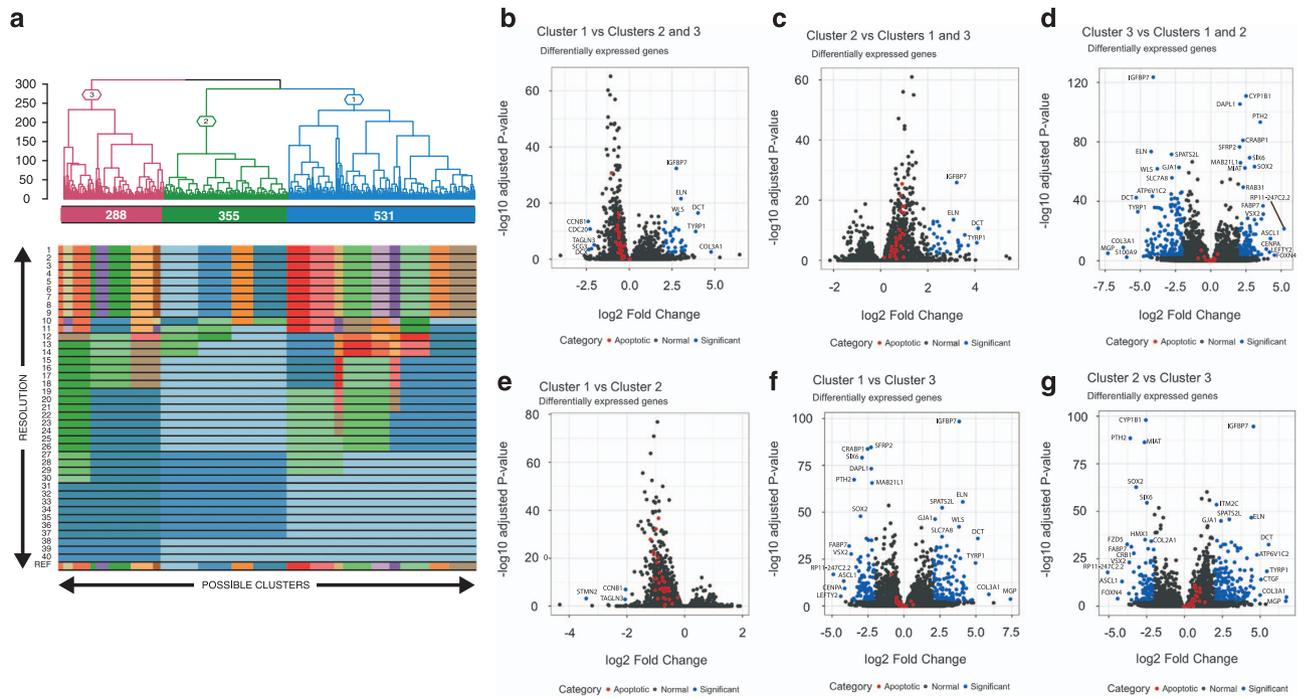
### Data Records

Data is available at ArrayExpress under accession number: E-MTAB-6108. Files consist of raw FASTQ files as well as a tab separated matrix of Transcripts Per Million for each cell passing quality control filtering. BAM files can be generated by using the supplied repository to process the FASTQ files via Cell Ranger.

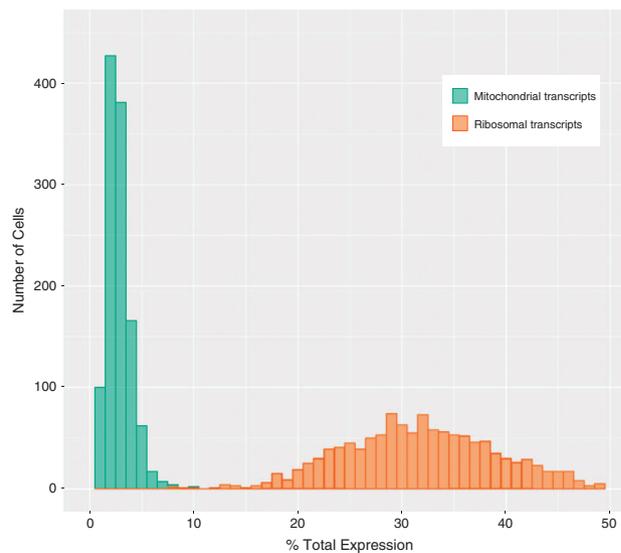
### Technical Validation

The hESC reporter line BRN3B-mCherry A81-H7 was differentiated to RGCs following our established protocol<sup>7</sup>, changing culture medium every second day. After 36 days, selection of RGCs using THY1 was performed by FACS. Both positive and negative THY fractions were harvested, and single cells harvested for library preparation and scRNA-Seq as outlined in Fig. 1. Processing our initial analysis identified a group of 61 cells whose expression levels indicated degradation and apoptosis (Fig. 2a and b). These 61 cells were removed from the data and the expression data from the remaining 1,174 healthy cells was re-normalised and analysed. Clustering of these 1,174 cells identified three distinct subpopulations consisting of 531, 355 and 288 cells (Fig. 3a). We performed differential expression analysis and subsequently pathway enrichment to characterise the molecular functions of these subpopulations (Fig. 3b-g). The proportion of reads mapped to mitochondrial and ribosomal genes are displayed in Fig. 4.

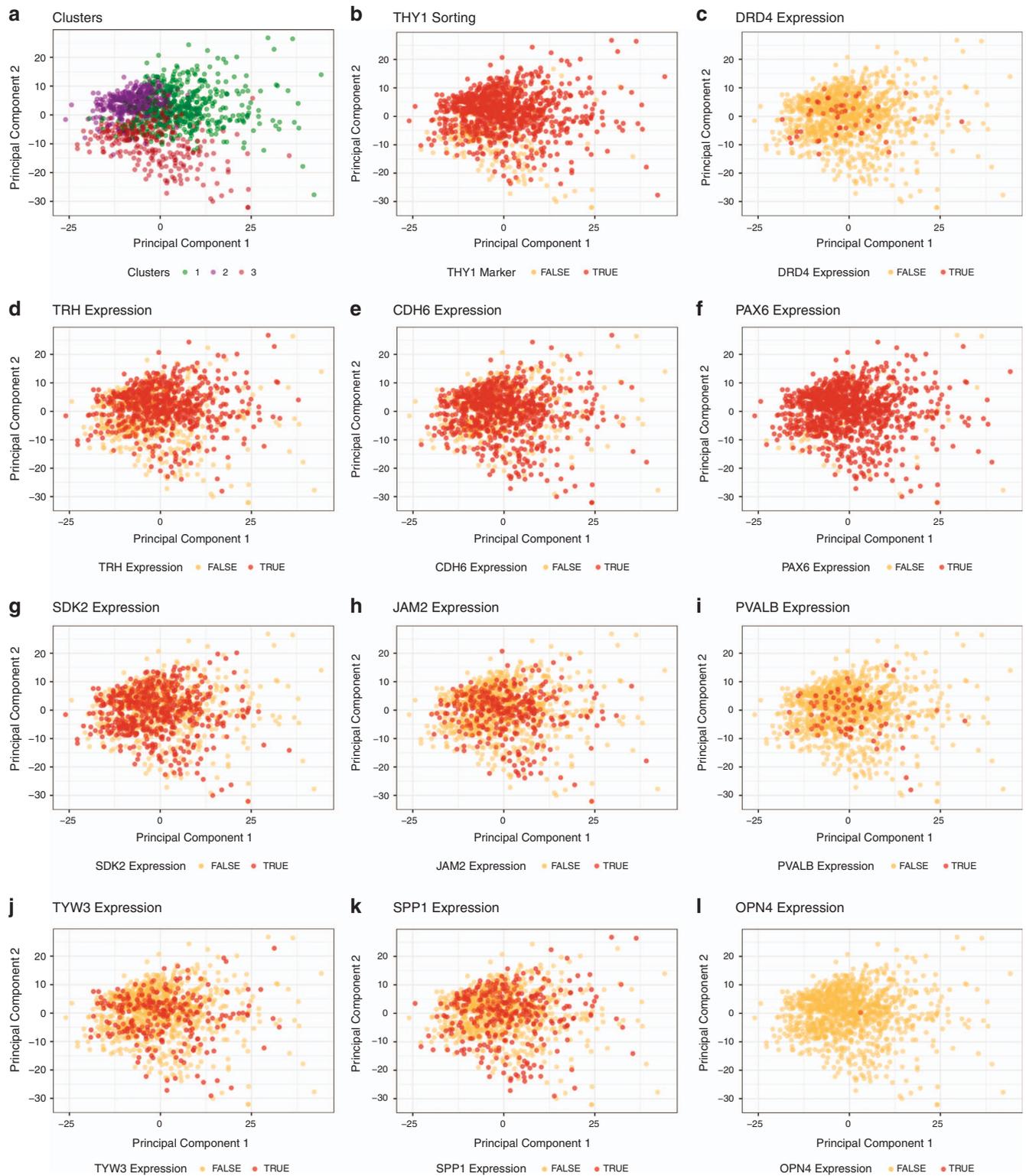
The 531 cells from subpopulation one were upregulated for genes associated with neural cell adhesion molecule signalling for neuronal outgrowth and Hedgehog pathway, which plays various roles in patterning of the central nervous system. Interestingly, genes implicated in collagen biosynthesis,



**Figure 3. Characterisation of filtered cells via differential expression.** (a) Unsupervised clustering of cells after filtering into three subpopulations. The dendrogram tree displays distance and agglomerative clustering of the cells. Each branch represents one subpopulation. The clustering is based on the most stable clustering result across 40 tree cut heights. The branches are labelled with their subpopulation identification. The numbers of cells in each of the three populations are given below the branches. The top significantly expressed genes of cells of each cluster vs other clusters; (b) one vs two and three, (c) two vs one and three, (d) three vs one and two. The top significantly expressed genes of cells in subpopulation one vs two (e), one vs three (f), two vs three (g). Genes represented in blue and black points are those in the top 0.5% highest  $-\log(P\text{-value})$ . Genes represented by red points are related to apoptotic pathways.



**Figure 4. Percentage of reads that mapped to mitochondrial and ribosomal transcripts.**



**Figure 5.** Cell heterogeneity within the RGC population using the PCA analysis. Shown are PC1 (x axis) vs PC2 (y axis) scores of cells from three clusters. (a) Cells from cluster 1 are plotted in red, cluster 2 in blue and cluster 3 in green. Relative expression of genes characteristic for the RGC subtypes: (b) THY1 (thymocyte antigen 1, CD90), (c) DRD4 (Dopamine receptor D4), (d) TRH (Thyrotropin-releasing hormone), (e) CDH6 (Cadherin-6), (f) PAX6 (Paired box 6), (g) SDK2 (Protein sidekick-2), (h) JAM2 (Junctional adhesion molecule B), (i) PVALB (Parvalbumin), (j) TYW3 (TRNA-YW synthesizing protein 3), (k) SPP1 (Secreted phosphoprotein 1, Osteopontin), (l) OPN4 (Melanopsin).

extracellular matrix proteoglycans and axon guidance were downregulated (Supplementary Table S1, Data Citation 2). This pattern of gene expression suggests a progenitor or an early differentiation state. Cells from subpopulation two contained upregulated genes associated with control of the Notch protein expression, implicated in the neuronal function and development, and DNA repair (Supplementary Table S2, Data Citation 2). Collectively, this pattern of gene expression is indicative of a more differentiated RGC phenotype than the cells in cluster one. The 288 cells identified as subpopulation three contained upregulated genes involved in axon guidance, together with semaphorin interactions, cell-extracellular matrix interactions and extracellular matrix proteoglycans. Furthermore, we observed significant downregulation of multiple genes associated with cell cycle (Supplementary Table S3, Data Citation 2). Taken together this indicates that this subpopulation three represents a more mature neuronal phenotype compared to cells in the other two subpopulations. Of note, one cell within subpopulation one was found to express *OPN4* a gene known to be expressed in intrinsically-photosensitive RGCs. These data indicate different levels of maturity of ESC derived RGC, with this conclusion supported by observed pathway enrichment (Supplementary Table S4-6, Data Citation 2). We have also conducted differential expression and pathway enrichment analysis to explore expression of genes associated with different RGC subtypes. In our analysis, we identified a number of genes that were previously shown to be expressed in various RGC subpopulations (Fig. 5). These genes could mark at least 9 separate RGC subtypes that grouped into three clusters<sup>30</sup>. This discrepancy may be explained by the fact that the number of RGC subtypes was estimated based on their morphological features but also expression of the molecular markers, including cell surface protein THY1, transcription factors from the Brn (Pou4F) family and RNA-binding protein RBPMS<sup>30</sup>. In our experimental design, we included the RGC enrichment step by sorting differentiated cells for THY1<sup>7</sup>. Limitation of this approach is exclusion of cells that could have the RGC identity without THY1 expression; however, we chose this marker as it is the cell surface protein and thus allows the maintenance of live cells in culture prior to sequencing or further characterisation.

### Usage Notes

Our experiment was designed to assess the different subpopulations of RGCs post differentiation from hESCs. hESC-derived RGCs obtained in a 36-day guided differentiation clustered into distinct subpopulations of neurons. Our initial analysis identified a group of 61 cells that showed a strong enrichment of stress and apoptosis pathways. This is possibly due to the FACS procedure itself, which can be stressful on cells. All post quality-control cells express genes relevant to RGC structure and functions. Altogether, our data provides strong support of an RGC identity of the cells in all clusters.

### References

- Martin, G. R. Isolation of a pluripotent cell line from early mouse embryos cultured in medium conditioned by teratocarcinoma stem cells. *Proc. Natl. Acad. Sci. USA* **78**, 7634–7638 (1981).
- Evans, M. J. & Kaufman, M. H. Establishment in culture of pluripotential cells from mouse embryos. *Nature* **292**, 154–156 (1981).
- Reubinoff, B. E., Pera, M. F., Fong, C. Y., Trounson, A. & Bongso, A. Embryonic stem cell lines from human blastocysts: somatic differentiation in vitro. *Nat. Biotechnol.* **18**, 399–404 (2000).
- Thomson, J. A. Embryonic Stem Cell Lines Derived from Human Blastocysts. *Science* **282**, 1145–1147 (1998).
- Takahashi, K. & Yamanaka, S. Induction of pluripotent stem cells from mouse embryonic and adult fibroblast cultures by defined factors. *Cell* **126**, 663–676 (2006).
- Takahashi, K. *et al.* Induction of pluripotent stem cells from adult human fibroblasts by defined factors. *Cell* **131**, 861–872 (2007).
- Gill, K. P. *et al.* Enriched retinal ganglion cells derived from human embryonic stem cells. *Sci. Rep* **6**, 30552 (2016).
- Gill, K. P., Hewitt, A. W., Davidson, K. C., Pébay, A. & Wong, R. C. B. Methods of Retinal Ganglion Cell Differentiation From Pluripotent Stem Cells. *Transl. Vis. Sci. Technol* **3**, 7 (2014).
- Sluch, V. M. *et al.* Differentiation of human ESCs to retinal ganglion cells using a CRISPR engineered reporter cell line. *Sci. Rep* **5**, 16595 (2015).
- Cook, D. *et al.* Lessons learned from the fate of AstraZeneca's drug pipeline: a five-dimensional framework. *Nat. Rev. Drug Discov.* **13**, 419–431 (2014).
- Shalek, A. K. *et al.* Single-cell transcriptomics reveals bimodality in expression and splicing in immune cells. *Nature* **498**, 236–240 (2013).
- Wills, Q. F. *et al.* Single-cell gene expression analysis reveals genetic associations masked in whole-tissue experiments. *Nat. Biotechnol.* **31**, 748–752 (2013).
- Kolodziejczyk, A. A. *et al.* Single Cell RNA-Sequencing of Pluripotent States Unlocks Modular Transcriptional Variation. *Cell Stem Cell* **17**, 471–485 (2015).
- Marques, S. *et al.* Oligodendrocyte heterogeneity in the mouse juvenile and adult central nervous system. *Science* **352**, 1326–1329 (2016).
- Usoskin, D. *et al.* Unbiased classification of sensory neuron types by large-scale single-cell RNA sequencing. *Nat. Neurosci.* **18**, 145–153 (2015).
- Etzrodt, M., Ende, M. & Schroeder, T. Quantitative single-cell approaches to stem cell research. *Cell Stem Cell* **15**, 546–558 (2014).
- Quadrato, G. *et al.* Cell diversity and network dynamics in photosensitive human brain organoids. *Nature* **545**, 48–53 (2017).
- Birey, F. *et al.* Assembly of functionally integrated human forebrain spheroids. *Nature* **545**, 54–59 (2017).
- Tang, F. *et al.* Tracing the derivation of embryonic stem cells from the inner cell mass by single-cell RNA-Seq analysis. *Cell Stem Cell* **6**, 468–478 (2010).
- Rouso, D. L. *et al.* Two Pairs of ON and OFF Retinal Ganglion Cells Are Defined by Intersectional Patterns of Transcription Factor Expression. *Cell Rep* **15**, 1930–1944 (2016).
- Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
- Leys, C., Ley, C., Klein, O., Bernard, P. & Licata, L. Detecting outliers: Do not use standard deviation around the mean, use absolute deviation around the median. *J. Exp. Soc. Psychol.* **49**, 764–766 (2013).
- Zheng, G. X. Y. *et al.* Massively parallel digital transcriptional profiling of single cells. *Nat. Commun* **8**, 14049 (2017).



24. Lun, A. T. L., Bach, K. & Marioni, J. C. Pooling across cells to normalize single-cell RNA sequencing data with many zero counts. *Genome Biol.* **17**, 75 (2016).
25. McCarthy, D. J., Campbell, K. R., Lun, A. T. L. & Wills, Q. F. Scater: pre-processing, quality control, normalization and visualization of single-cell RNA-seq data in R. *Bioinformatics* **33**, 1179–1186 (2017).
26. Nguyen, Q. *et al.* Single-Cell Transcriptome Sequencing Of 18,787 Human Induced Pluripotent Stem Cells Identifies Differentially Primed Subpopulations *bioRxiv*, doi:10.1101/119255 (2017).
27. Rand, W. R. Objective Criteria for the Evaluation of Clustering Methods. *J. Am. Stat. Assoc.* **66**, 846–850 (1971).
28. Anders, S. & Huber, W. Differential expression analysis for sequence count data. *Genome Biol.* **11**, R106 (2010).
29. Wu, G., Feng, X. & Stein, L. A human functional protein interaction network and its application to cancer data analysis. *Genome Biol.* **11**, R53 (2010).
30. Sanes, J. R. & Masland, R. H. The types of retinal ganglion cells: current status and implications for neuronal classification. *Annu. Rev. Neurosci.* **38**, 221–246 (2015).

## Data Citations

1. *ArrayExpress* E-MTAB-6108 (2017).
2. Daniszewski, M., Senabouth, A., Pébay, A., Powell, J. E. & Hewitt, A. W. *figshare* <https://doi.org/10.6084/m9.figshare.5715148> (2017).

## Acknowledgements

This work was supported by grants from the Sir Edward Dunlop Medical Research Foundation, the Ophthalmic Research Institute of Australia, Retina Australia, the Joan and Peter Clemenger Foundation, and a Philip Neal bequest. AWH & JEP are supported by National Health and Medical Research Council Fellowships, AP by an Australian Research Council Future Fellowship (FT140100047), and MD by an International Postgraduate Award Scholarship from the University of Melbourne. CERA receives Operational Infrastructure Support from the Victorian Government.

## Author Contributions

M.D.: concept and design, experimental work, interpretation of data, writing of manuscript, final approval of manuscript. A.S.: concept and design, experimental work, interpretation of data, writing of manuscript, final approval of manuscript. Q.N.: experimental work, interpretation of data, writing of manuscript, final approval of manuscript. D.E.C.: interpretation of data, final approval of manuscript. S.W.L.: experimental work, final approval of manuscript. T.J.: interpretation of data, final approval of manuscript. V.M.S.: provision of reagents, final approval of manuscript. J.S.J.: experimental work, final approval of manuscript. X.C.: provision of reagents, final approval of manuscript. D.J.Z.: provision of reagents, final approval of manuscript. A.P.: concept and design, financial support, interpretation of data, writing of manuscript, final approval of manuscript. J.E.P.: concept and design, financial support, interpretation of data, writing of manuscript, final approval of manuscript. A.W.H.: concept and design, financial support, interpretation of data, writing of manuscript, final approval of manuscript.

## Additional information

Supplementary Information accompanies this paper at <http://www.nature.com/sdata>

**Competing interests:** The authors declare no competing financial interests.

**How to cite this article:** Daniszewski, M. *et al.* Single cell RNA sequencing of stem cell-derived retinal ganglion cells. *Sci. Data* 5:180013 doi:10.1038/sdata.2018.13 (2018).

**Publisher's note:** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>

The Creative Commons Public Domain Dedication waiver <http://creativecommons.org/publicdomain/zero/1.0/> applies to the metadata files made available in this article.

© The Author(s) 2018