

SUPPLEMENTAL MATERIAL for:
Transcriptional and Cellular Diversity of the Human Heart

Nathan R. Tucker,^{1,2,3,#} Mark Chaffin,^{1,#} Stephen J. Fleming,^{1,4} Amelia W. Hall,^{1,2}
Victoria A. Parsons,² Kenneth Bedi,⁵ Amer-Denis Akkad,^{1,6} Caroline N. Herndon,¹ Alessandro Arduini,¹
Irinna Papangeli,^{1,6} Carolina Roselli,^{1,7} François Aguet,⁸ Seung Hoan Choi,¹ Kristin G. Ardlie,⁸ Mehrtaash
Babadi,^{1,4} Kenneth B. Margulies,⁵ Christian M. Stegmann,^{1,6} and Patrick T. Ellinor^{1,2,*}

1. Precision Cardiology Laboratory, The Broad Institute, Cambridge, MA, USA 02142
 2. Cardiovascular Research Center, Massachusetts General Hospital, Boston, MA, USA 02114
 3. Masonic Medical Research Institute, Utica, NY, USA 13501
 4. Data Sciences Platform, The Broad Institute of MIT and Harvard, Cambridge, MA, USA 02142
 5. Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, USA 19104
 6. Precision Cardiology Laboratory, Bayer US LLC, Cambridge, MA, 02142
 7. University Medical Center Groningen, University of Groningen, 9712 CP, Groningen, NL
 8. The Broad Institute of MIT and Harvard, Cambridge, MA, USA 02142
- # These authors contributed equally

Running Title: Single cell transcriptomics of the human heart

Word count:

Keywords: Heart, single cell sequencing, cardiovascular disease, genetics

Corresponding Author:

Patrick T. Ellinor, MD, PhD
The Broad Institute of MIT and Harvard
75 Ames Street
Cambridge, MA 02142
ellinor@mgh.harvard.edu

Table of Contents:

Description	Page
Supplemental Methods	3
Supplemental Tables II, IV-VIII: provided as .xls files	8
Table I. Resources and Reagents	8
Table II. Summary of quality control metrics for each sample processed.	
Table III. Number of cells observed in each cell cluster within the global map.	10
Table IV. Marker genes for clusters identified in the joint UMAP plot.	
Table V. Marker genes for subclusters identified following secondary clustering of major cell types.	
Table VI. Results of differential expression analysis for each chamber level comparison for cardiomyocytes, endothelium, fibroblasts, macrophages and pericytes.	
Table VII. Results of differential expression analysis between males and females within all chambers, and each chamber individually, for cardiomyocytes, endothelium, fibroblasts, macrophages, and pericytes.	
Table VIII. LD score regression analyses with various cell type selectivity thresholds for marker genes	
Figure I: Analytic workflow and quality control metrics.	11
Figure II: Establishment of quality control metrics.	12
Figure III: Subclustering analysis of cardiomyocytes and pericytes.	14
Figure IV: Histological analysis of donor tissue.	15
Figure V: Ontology analysis of left vs right specific genes	16
Figure VI: Intersection of snRNAseq data with clinical testing panels and eQTL data.	17

Expanded Methods:

Single nucleus RNA-sequencing

Given the relatively small amount of input tissue (<100mg) and a desire to process nuclear RNAs from thaw to reverse transcription as quickly as possible, we chose a nuclear isolation strategy that omits gradient centrifugation or FACS-based sorting. Cytoplasmic fragments or suspected doublets that may be retained by this strategy were removed informatically as described in the next section. Single nucleus suspensions were generated by a series of cellular membrane lysis, differential centrifugation and filtration steps. Approximately 100mg of tissue was cryosectioned at 100 μ m on a Leica CM1950 cryostat to enable liberation of nuclei from the tissue while minimizing mechanical manipulation. Tissue sections were homogenized in a dounce homogenizer after suspension in 4mL of ice cold lysis buffer containing propidium iodide for nuclear staining (250mM Sucrose, 25mM KCl, 0.05% IGEPAL-630, 3mM MgCl₂, 1 μ M DTT, 10mM Tris pH 8.0). After 5 minutes incubation, large debris was pelleted at 20g for 1 min in a Beckman Coulter Allegra X-15R swinging bucket centrifuge. Supernatant was brought to 8mL of total volume with nuclear wash buffer (PBS + 3mM MgCl₂ + 0.01% BSA) then filtered sequentially through a 100 μ m and 20 μ m filter (pluriSelect Life Science). Nuclei were pelleted at 400g for 5 minutes at 4C, washed in 4mL of nuclear wash buffer and repelleted. After removal of wash buffer, nuclei were resuspended in approximately 500 μ L of cold nuclear resuspension buffer (Nuclear wash buffer + 0.4U/ μ L of murine RNase inhibitor (New England Biolabs)) with gentle trituration then counted on a hemocytometer. Cells were loaded into the 10x Genomics microfluidic platform (Single cell 3' solution, v2) for an estimated recovery of 5000 cells per device. Processing of libraries was performed according to manufacturer's instructions with a few modifications. First, nuclei were incubated at 4C for 30 minutes after emulsion generation to promote nuclear lysis. Second, the reverse transcription protocol was modified for one of the two replicates to be 42C for 20 minutes then 53C for 120 minutes. This is noted as (_2) in the Sample ID column of the sample information table (Supplemental Table ST1). Libraries were multiplexed at an average of 4 libraries per flow cell on an Illumina Nextseq550 in the Broad Institute's Genomics Platform.

Sample selection and quality control

In total, 56 single nuclei RNA-seq experiments were performed from all four chambers of the human heart in 7 distinct biological individuals, processed in duplicate. Reads from single nuclei experiments were de-multiplexed and aligned to a GRCh38 human pre-mRNA reference using the 10x Genomics toolkit Cell Ranger 2.1.1 and default parameters with the exception of setting the --expect-cells flag to 5000 based on library preparation.

For each experiment, the distribution of the number of unique molecular identifiers (UMI) was visually inspected to identify experiment failures. Any experiment where the fraction of reads in cells was less than 33% (n=7) or the median UMI per cell was less than 600 (n=6) based on Cell Ranger were excluded from further analysis. In total 9 experiments failed on these criteria. These cutoffs corresponded to poor structure in the UMI decay curve.

Additionally, genetic concordance was checked between all experiments of the same biological individual using the Genome Analysis Toolkit (GATK) [46] method *CrosscheckFingerprints* on the single cell 10x aligned reads. A list of approximately 6,300 sites was provided and samples were considered concordant if their corresponding LOD score was greater than 10. Two experiments from the right

ventricle of individual **P1681** were discordant with the remaining experiments from the same individual and were subsequently removed.

Post-Sample Selection Processing: CellBender

All 45 experiments passing initial quality control were processed using the *remove-background* tool from CellBender v0.1 to determine which droplets contain a cell and to correct gene count matrices by removing ambient background RNA contamination. For complete details on the CellBender *remove-background* model, see <https://github.com/broadinstitute/CellBender>. Briefly, CellBender performs Bayesian inference in the context of a probabilistic model to remove ambient RNA by estimating the contribution of ambient background RNA captured in each droplet and adjusting the count matrix appropriately. The CellBender model does require that there are some unambiguously empty droplets containing only ambient, background RNA. One additional experiment was removed because of poor definition in the UMI decay curve which prevented CellBender from converging appropriately.

Specifically, CellBender *remove-background* was run on a Tesla K80 GPU using the following parameters: expected-cells 5000, total-droplets-included 20000, low-count-threshold 50, epochs 300, z-dim 200, z-layers 1000, empty-drop-training-fraction 0.3. After an examination of the output cell calls, it was determined that unusually high ambient RNA had led to a failure on sample RV_1666_2, which was subsequently rerun with the following parameters altered: total-droplets-included 15000, low-count-threshold 200.

Four chamber map aggregation

Prior to cell clustering, additional low quality cells were removed on a per-experiment basis to account for large variability in sequencing depth and complexity between experiments. These pre-processing steps were performed using Seurat 2.3.4. In brief, cells were removed from an experiment if: 1) the number of genes detected was less than 100 or greater than a predefined upper outlier cutoff, 2) the number of UMI for the cell was greater than a predefined upper outlier cutoff, or 3) the percent of mitochondrial gene content was greater than 5%. The upper outlier cutoff was calculated as the third quartile plus 1.5 times the interquartile range. Upper cutoffs were used to minimize the introduction of multiplets into downstream clustering. These criteria reduced the total number of putative cells from 373,243 to 287,269 for subsequent aggregation.

Highly variable genes were selected to perform cell clustering using Seurat 2.3.4. The aggregated cell count matrix was first normalized by dividing the number of UMI for each transcript by the total UMI for the cell, multiplying by 10,000, and taking the natural log of these results. Variable genes were found globally using the *FindVariableGenes* function in Seurat which bins genes by average expression and calculates a Z score for dispersion within each bin [1]. Normalized expression bounds were set between 0.03 and 5, and genes with dispersion Z score greater than 0.5 were selected. In total, 1,969 genes remained for clustering.

Because of large biological heterogeneity between samples, we applied the single-cell variational inference (scVI) framework to map cells from different samples onto a joint coordinate system.[41] In brief, scVI uses deep neural networks to learn the underlying distributions of cell-level expression, while accounting for batch variables. By treating the individual, rather than the experiment, as a batch indicator, this procedure aligns cells accounting for heterogeneity between individuals. Note that the success of scVI in removing biological batch effects implies that in our dataset, batch effect is dominated by biological inter-individual variability and is not technical in nature. The inferred latent space can then

be used for downstream clustering of cells. We applied scVI 0.3.0 on the previously identified 1,969 genes to estimate 50 latent variables. A neighborhood graph of cells was built based on these 50 latent variables using scanpy 1.4 [10] (*scanpy.pp.neighbors*) selecting a cosine distance metric and using 15 neighbors. Cells were subsequently placed into clusters using the Louvain algorithm (*scanpy.api.tl.louvain*) with default parameters and a resolution parameter of 1.0. To visualize cells in a high dimensional space, uniform manifold approximation and projection (UMAP) was applied to the same latent space using a cosine distance metric, and default parameters.[47]

Calculation of intronic and exonic reads

scR-Invex (Aaron Graubert, François Aguet; <https://github.com/broadinstitute/scrinvex>) was run in order to count the reads in each BAM file output by CellRanger 2.1.1 *count* that mapped to intronic, exonic, and junction regions of the transcriptome.

Sub-chamber visualization

Sub-chamber maps were generated using global latent variables and retaining cluster identities from the global cell map. Default UMAP parameters were used with the exception of setting *metric='cosine'*, *spread = 1*, *min_dist = .3*, and *n_neighbors = 25* for left ventricle and *metric='cosine'*, *spread = 1*, *min_dist = .1*, and *n_neighbors = 15* for left atrium, right ventricle, and right atrium.

Tissue staining and microscopy

RNA *in situ* hybridization was performed using the RNAscope 2.5 High Definition and RNAscope Multiplex Fluorescent v2 assays from Advanced Cell Diagnostics, Inc. (catalog numbers 322370 and 323100, respectively) following the manufacturer's protocols with the following modifications. During tissue preparation, fresh frozen tissue was sectioned at 15 μ m and mounted onto Superfrost Plus Slides (VWR). Fixation was performed in 4% PFA for 10 min at 4°C. Protease treatment was performed using RNAscope Protease III for 15 min at RT. Probes for *ADAMTS4*, *DCN* and *RYR2* were obtained from Advanced Cell Diagnostics. Hematoxylin and eosin or Masson's Trichrome stains were performed on paraffin embedded sections according to standard protocols.

Marker gene identification

Genes discriminating each cluster were identified by calculating the area under the receiver operating characteristic curve (AUC) for all genes comparing cells from the target cluster to all other cells not included in that cluster. The AUC will indicate how well a gene discriminates cells of a given cluster from those of all other clusters, with a value of 0.50 designating no discrimination and a values of 1 designating perfect discrimination. Two clusters determined to contain a high fraction of cytoplasmic material based on the proportion of exonic reads detected with scR-Invex were excluded from these calculations. Data were normalized by dividing the number of UMI for each transcript by the total UMI for the cell, multiplying by 10,000, and taking of the natural log of this result. The classifier used to calculate the AUC was built by taking the normalized expression values in each cell as predictions, and the cell cluster assignment as the class being predicted. Genes with an AUC greater than 0.70 and an average natural log fold-change > 0.6 were selected as markers of a cluster. These cutoffs were chosen to balance selecting genes with moderate discrimination for the cluster of interest, while not being overly inclusive so that a reasonably sized set of genes (hundreds for the larger clusters and a less than 10 for the smallest cluster) was available for each cluster. A less stringent AUC cutoff was required than

others have used in the past [48,49] as: 1) the reference group for a cluster of interest sometimes contains cells with similar expression profiles (e.g., multiple fibroblast clusters) and 2) single nuclei RNA-seq create a generally higher noise ratio than whole cell RNA-seq. In some clusters, these gene lists did not provide sufficient information to identify cell types. In those cases, genes expressed in > 5% of cluster cells that showed a standardized positive predictive value (PPV50) > 0.90 were also examined [50,51]. For each gene, the PPV50 between the target cluster of interest and all cells from other clusters quantifies the probability that a cell is of the target cluster when it expresses that gene of interest (UMI > 0), standardizing to an equal number of cells between clusters (prevalence=50%). Standardization was necessary given the highly variable numbers of cells within each cluster, which systematically reduce the number of genes found with an unstandardized PPV for smaller clusters.

Sub-clustering

To uncover potential sub-clusters of cells within the major clusters identified above, a simple sub-clustering procedure was performed for the follow cell types: cardiomyocytes (clusters 3, 4, 6, 15), fibroblasts (cluster 1, 2, 14), endothelial cells (cluster 9 and 10), pericytes (cluster 7), macrophages (cluster 8), adipocytes (cluster 11), vascular smooth muscle (cluster 13), neuronal (cluster 16), and lymphocytes (cluster 17). A new neighborhood graph was built for each of these groups using cosine distance based on the latent variables derived from the global scVI model. Louvain clustering was applied using varying resolution (ranging from 0.2 to 0.6) to establish new sub-groups. AUC was calculated for each sub-cluster compared to the remaining cells in the cluster. Genes with AUC greater than 0.65 and an average natural log fold-change > 0.5 were selected as markers of sub-clusters. A more liberal cutoff was employed here as cell sub-clusters look more similar to one another on average. Similarly to the global map, when necessary genes expressed in > 5% of target subcluster cells with a PPV50 > 0.90 were examined to help determine potential cell sub-types. Putative spurious sub-groups were identified based on an elevated proportion of exonic reads to all reads, which are often marked by increased mitochondrial genes.

Gene ontology analysis

Gene ontology analysis was performed using the R package GOstats version 2.46.0. Ensembl identifiers were mapped to Entrez gene identifiers when possible for compatibility with GOstats gene ontologies. The gene universe was set to all protein coding genes that were successfully mapped and only gene sets with a minimum size of 5 were considered for enrichment testing.[52] For each set of marker genes, a hypergeometric test was performed to test for enrichment of genes in each ontology, considering only ontologies with at least one gene overlapping the given marker list. A Bonferroni significance threshold was used for each set of markers correcting for the number of ontologies tested.

Genome-wide association study integration

For six cardiometabolic traits with genome-wide association studies (GWAS), we looked for enriched heritability around marker genes of given cell types using stratified linkage disequilibrium (LD) score regression.[53,54] We considered major cell types for this analysis, excluding low quality sub-clusters identified as described above. Only genes with a total of at least 10 counts across all cells were considered. Gene coordinates were used from the GRCh37 Ensembl reference to align with LD score regression methods. When genes from the GRCh38 Ensembl reference were not available in GRCh37 Ensembl reference, coordinates were lifted back using liftOver [55] when possible. In total, 25,968 genes

were considered. For each cluster, a new set of marker genes were identified based on having at least some discrimination for the cluster of interest over other cell types (AUC > 0.55). Single nucleotide polymorphisms (SNPs) within 100 KB of any gene identified this way were annotated for LD score regression based on 1000G European individuals. The LD score regression model was run including the baseline annotations generated in *Finucane et al 2015* [56] only considering high quality HapMap3 SNPs. The six GWAS traits used included atrial fibrillation,[33] PR Interval,[57] QT Interval,[58] coronary artery disease,[59] LDL,[60] and type 2 diabetes.[61] European ancestry-specific results were used when available to be most consistent with the LD reference panel.

Additionally, specific GWAS genes for atrial fibrillation, PR interval, QT interval, and coronary artery disease were highlighted based on colocalization of a genome-wide association signal and an expression quantitative trait loci (eQTL) from bulk sequence data of the relevant chamber of the heart. Colocalization was performed in a 1 MB region around the sentinel SNP of a GWAS locus using the `coloc.abf` function from the `coloc` package in R.[62] Allele frequency data was derived from the same European 1000 Genomes [63] samples used in the LD score regression analysis described above. Left ventricle eQTL data was taken from the Genotype-Tissue Expression (GTEx) project [32] based on 272 samples and left atrial eQTL data was taken from the MAGNet repository (<http://www.med.upenn.edu/magnet/>) based on 101 individuals.[33] Genes estimated to have a greater than 0.60 probability that the GWAS signal and eQTL signal share a causal variant we considered putative GWAS genes.

Table I: Resources and Reagents

Reagent	Source	Identifier
Single Cell 3' chip V2	10x Genomics	120236
Single Cell 3' v2 Mod 2	10x Genomics	120237
Mrine RNase inhibitor	New England Biolabs	M0314
AMPure XP	Beckman Coulter	A63881
α -ADAMTS4 RNAscope probe	Advanced Cell Diagnostics Inc. (ACD)	
α -DCN RNAscope probe	ACD	
α -RYS2 RNAscope probe	ACD	
RNAscope 2.5 High Definition assay	ACD	322370
RNAscope Multiplex Fluorescent v2 assay	ACD	323100
pluriStrainer	pluriSelect Life	43-50020-50
Software and Algorithms		
CellRanger 2.1.1		https://support.10xgenomics.com/
seurat 2.3.4	[44]	https://satijalab.org/seurat/
CellBender 0.1	[64]	https://github.com/broadinstitute/CellBender
scR-Invex		https://github.com/broadinstitute/scrinvex
scanpy 1.4	[10]	https://scanpy.readthedocs.io/en/stable/
scVI 0.3.0	[41]	https://scvi.readthedocs.io/en/master/
umap 0.3.7	[47]	https://umap-learn.readthedocs.io/en/latest/
R 3.5.0		https://www.r-project.org/
lme4 1.1-21		https://cran.r-project.org/web/packages/lme4/index.html
GOSTats 2.46.0	[65]	https://bioconductor.org/packages/release/bioc/html/GOstats.html
ldsc 1.0.0	[53,56]	https://github.com/bulik/ldsc

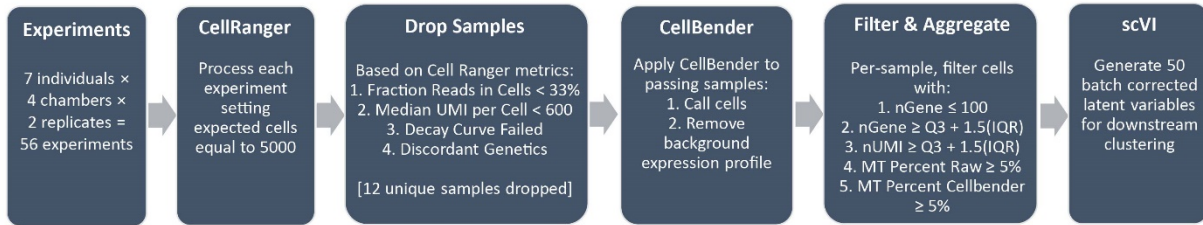
Coloc 3.2-1		https://cran.r-project.org/web/packages/coloc/index.html
-------------	--	---

Table III. Number of cells observed in each cell cluster within the global map.

	Number of Nuclei	Percent of Nuclei
1. Fibroblast I	39083	13.6%
2. Fibroblast II	38555	13.4%
3. Atrial Cardiomyocyte	34051	11.9%
4. Ventricular Cardiomyocyte I	28671	10.0%
5. Cytoplasmic Cardiomyocyte I	24695	8.6%
6. Ventricular Cardiomyocyte II	24528	8.5%
7. Pericyte	18467	6.4%
8. Leukocyte	17468	6.1%
9. Endothelium I	17142	6.0%
10. Endothelium II	10781	3.8%
11. Adipocyte	8658	3.0%
12. Cytoplasmic Cardiomyocyte I	6070	2.1%
13. Vascular Smooth Muscle	5740	2.0%
14. Fibroblast III	5582	1.9%
15. Ventricular Cardiomyocyte III	4707	1.6%
16. Neuronal	1568	0.5%
17. Lymphocyte	1503	0.5%

Figure I. Analytic workflow and quality control metrics. A: Analytic workflow for post-sequencing quality control through initial clustering. B: Table detailing the contribution of each sample to the final map by chamber and replicate number. UMI = unique molecular identifier; nGene = number of genes detected in a given cell; nUMI = the total UMI in a given cell; Q3 = 75th percentile; IQR = Interquartile Range; MT = Mitochondrial; RA = Right Atrium; LA = Left Atrium; RV = Right Ventricle; LV = Left Ventricle.

A



B

RA (4 Female, 2 Male)			LA (3 Female, 3 Male)		
RA_1221_1	RA_1681_1	RA_1708_1	LA_1600_1	LA_1681_1	LA_1708_1
	RA_1681_2	RA_1708_2	LA_1600_2	LA_1681_2	LA_1708_2
RA_1600_1	RA_1702_1	RA_1723_1	LA_1666_1	LA_1702_1	LA_1723_1
RA_1600_2	RA_1702_2	RA_1723_2	LA_1666_2	LA_1702_2	LA_1723_2
RV (3 Female, 2 Male)			LV (3 Female, 3 Male)		
RV_1221_1	RV_1702_1	RV_1723_1	LV_1221_1	LV_1681_1	LV_1708_1
RV_1221_2	RV_1702_2	RV_1723_2		LV_1681_2	LV_1708_2
RV_1666_1	RV_1708_1		LV_1666_1	LV_1702_1	LV_1723_1
RV_1666_2	RV_1708_2		LV_1666_2	LV_1702_2	LV_1723_2

Male:

- P1702
- P1666
- P1681

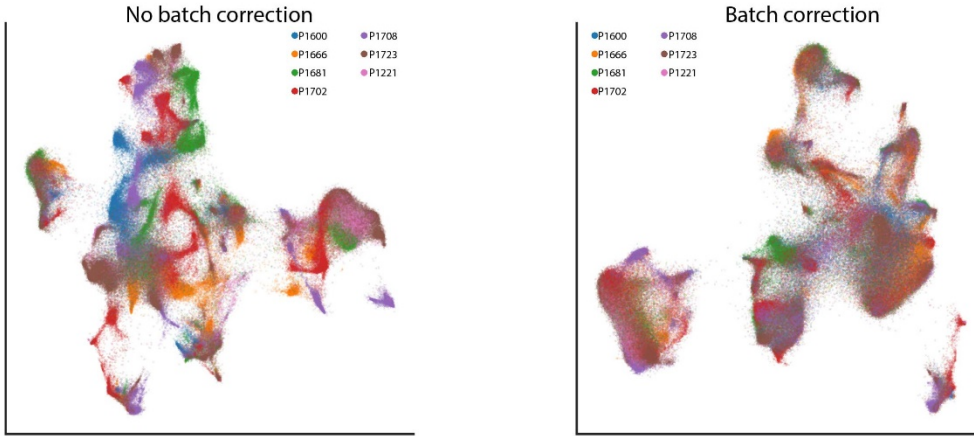
Female:

- P1221
- P1600
- P1723
- P1708

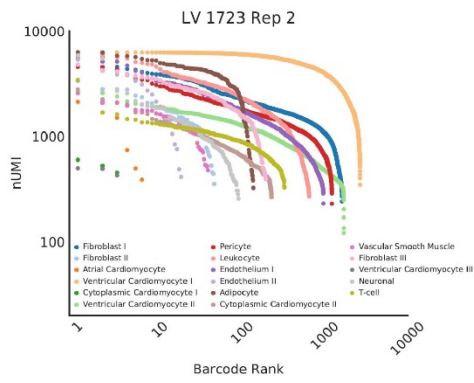
Figure II. Establishment of quality control metrics. A: UMAP plots generated both before and after batch correction by scVI. Colors of each dot correspond to the patient sample from which the cell arose. B: UMI decay curve from sample LV_1723_2 broken down by cell type after cell calling by CellBender. Colors correspond to the cell types as labeled. C: UMAP plot which displays the absolute percentage of mitochondrial (MT) reads in each cell. D: Density plot displaying the distribution of reads which map to exonic regions over total reads by cell type. Color corresponds to the cell clusters as called within the global Louvain clustering. Both clusters labeled as “cytoplasmic” indicate a higher than normal number of reads in the exonic regions.

Figure II:

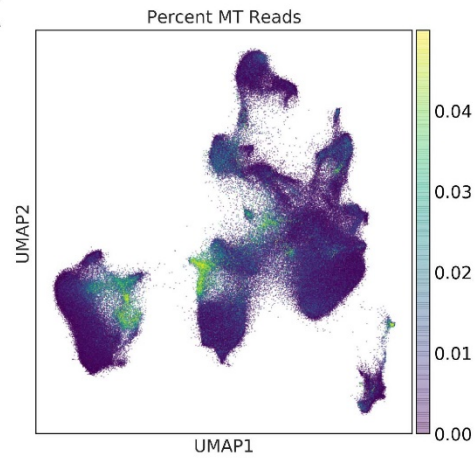
A



B



C



D

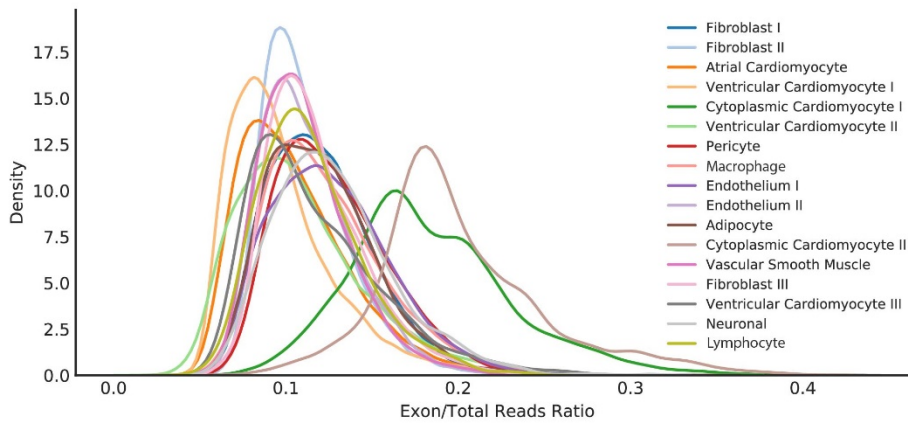


Figure III. Subclustering analysis of cardiomyocytes and pericytes. Results of major cell type subclustering for cardiomyocytes (A) and pericytes (B). Left panel displays the distribution of the identified subclusters within the global UMAP plot for all chambers. Each dot represents a cell, colored by its respective subcluster. Center panel is a density plot which displays the ratio of reads which lie in exonic regions as a percentage of total reads in each subcluster. Right panel is a dot plot displaying the genes with the highest AUC for each subcluster or those which are specifically mentioned within the text (*KCP*). The size of the dot represents the percentage of cells in the cluster in which each gene is detected and the color reflects the mean \log_2 expression. CM: Cardiomyocytes, PC: Pericytes.

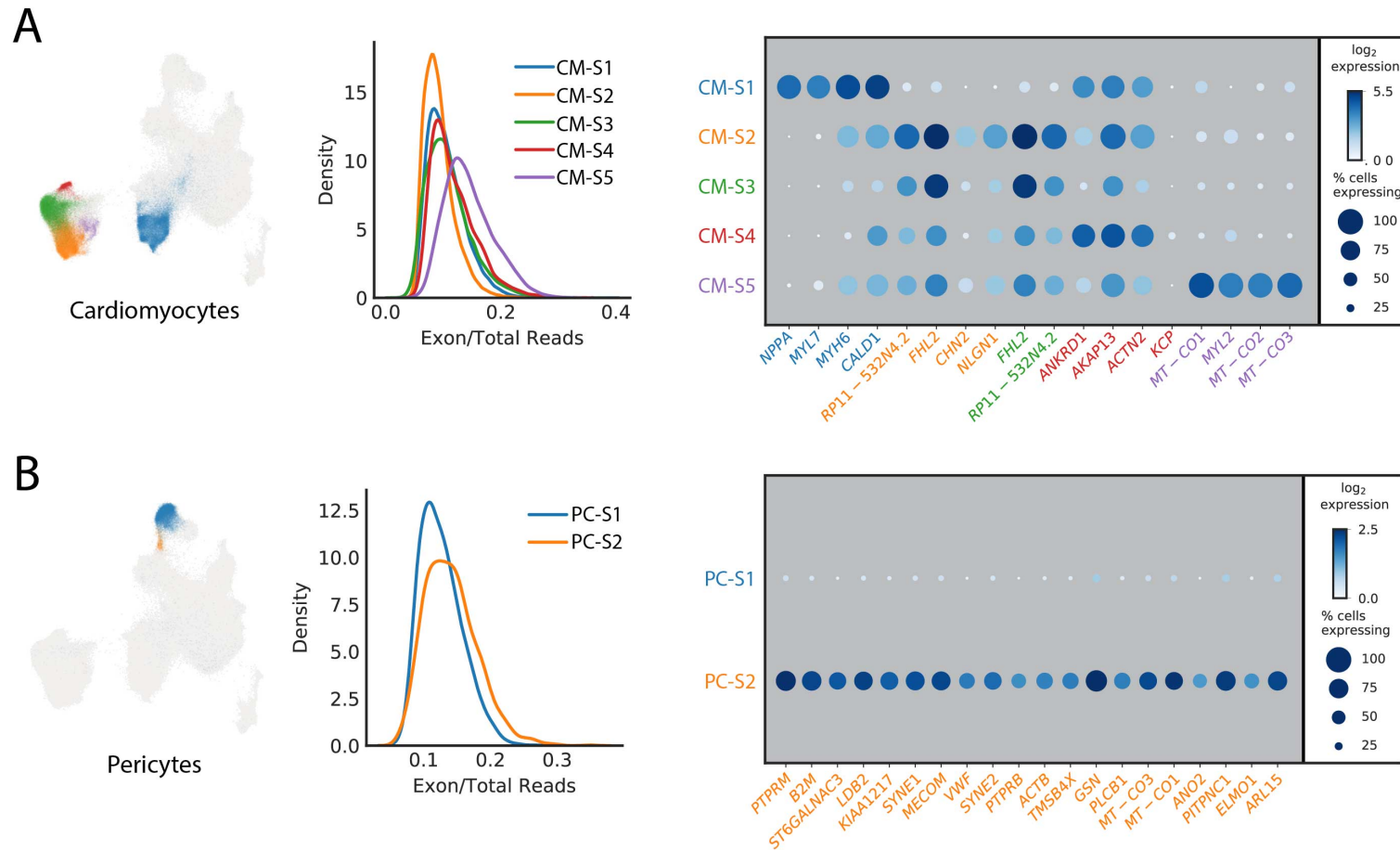
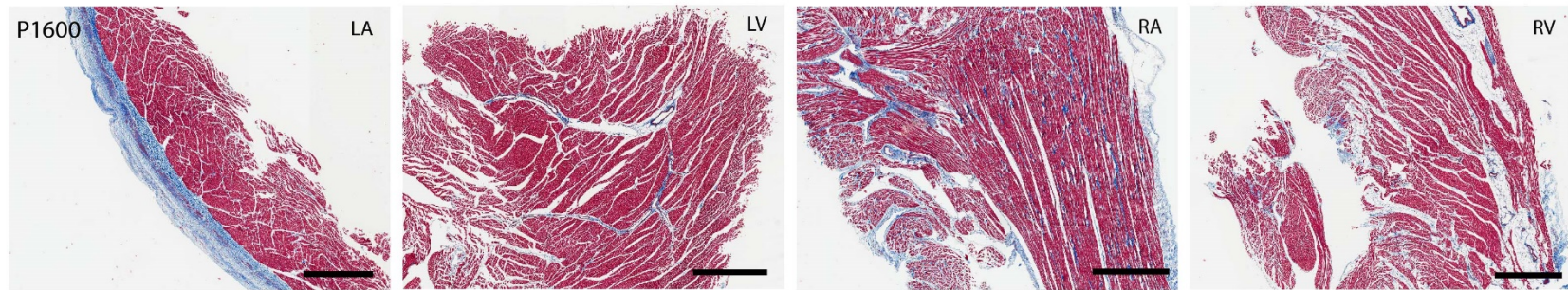


Figure IV. Histological analysis of donor tissue. A: Masson's Trichrome staining of representative regions from all four chambers in patient P1600. Scale bar indicates 1mm. B: Hematoxylin/eosin staining of samples from patient P1723. Scale bar represents 1mm. Arrows in right panel highlight regions of myocardial adiposity. RA: Right Atrium, LA: Left atrium, RV: Right ventricle, LV: Left ventricle.

A



B

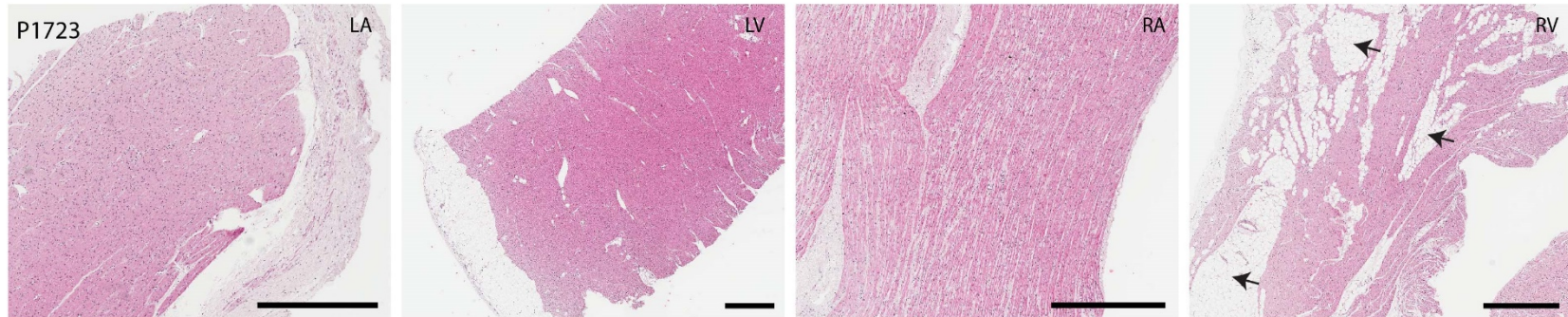


Figure V. Gene ontology analysis of left versus right specific genes. Enrichment analysis as performed by GOSTats using all genes which are specific for left or right sidedness at an FDR adjusted P-value of less than 0.01. Ngene: Number of genes which were significant and used for the ontology analysis.

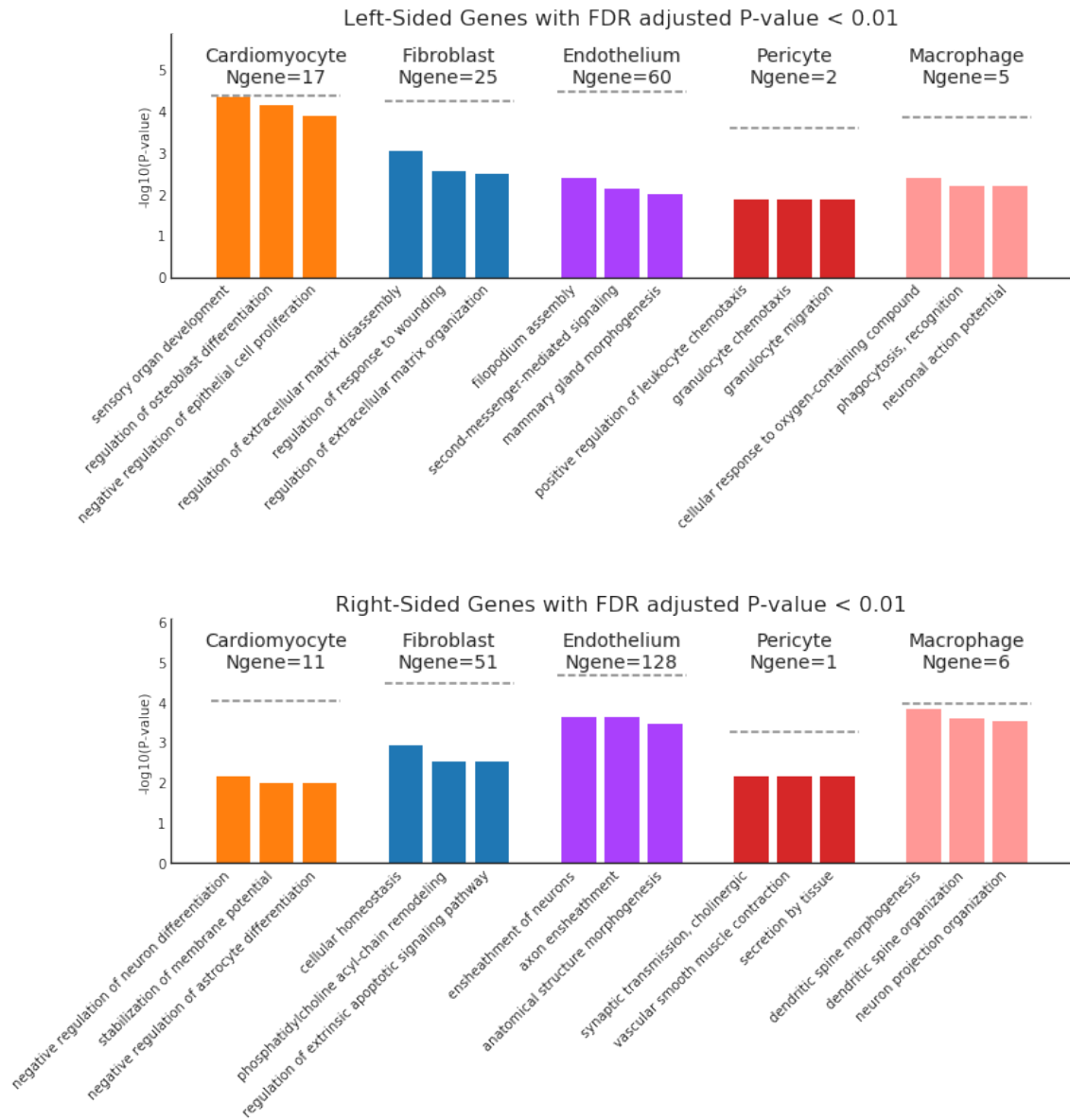


Figure VI. Intersection of snRNAseq data with clinical testing panels and eQTL data. A: Dot plot for genes currently on standard arrhythmia clinical testing panels. The size of each dot represents the percent of cells in which the gene of interest is detected and the shading represents the relative expression of the gene. Color of the genes correspond to the cell type for which the AUC reaches 0.70 or greater. Genes with black color indicate no cell type which reaches this threshold. Size and shade of the dot corresponds percentage of cells and relative expression, respectively. B: Dot plots for genes identified for eQTL analysis of left atrial and left ventricular tissue. SNPs used for eQTL analysis are those derived from the genome-wide association studies listed for cardiovascular traits and diseases.

Figure VI:

