**Supplemental Figure Legends**

**Supplemental Figure 1. Evaluation and benchmarking of the Capybara pipeline. Related to Figure 1 (A)** Distributions of quadratic programming (QP) metrics: Error, Lagrangian multiplier, and Deviance. Black line: distribution of these metrics in the reference. Other lines: ideal distributions for discrete, hybrid, and unknown classifications, which are modeled based on the reference metrics. Assuming each cell in the reference has a discrete identity, the ideal distribution of deviance for discrete cells is modeled based on the reference identity scores. The ideal distribution of deviance for hybrid exhibits a two standard deviation shift left of the discrete distribution; the deviance distribution for unknown exhibits a two standard deviation shift left of the hybrid distribution. **(B)** Evaluation using five human pancreatic datasets and the Allen Mouse Brain (AMB) atlas against ten other classifiers, using an established benchmarking pipeline (Abdelaal et al., 2019). The classification performance was evaluated by the area under the receiver operating characteristic (AUROC) and the mean total time in seconds. The color of each square labels the assessment of each aspect. **(C)** Cross-validation using *Tabula Muris* processed with 10x droplet-based or Smart-seq2 technologies. In this evaluation, 90 cells were sampled from each cell type of each tissue to construct the high-resolution reference. Sample cells used for reference construction were used as training sets for scMap and SingleCellNet. The remaining cells were used as test samples. The performance was evaluated by AUROC scores. **(D)** Demonstration of inappropriate reference selection: We applied Capybara to a publicly available Mouse PBMC dataset (https://www.10xgenomics.com/resources/datasets) with two references from the Mouse Cell Atlas (MCA; Han et al., 2018): peripheral blood (Left) and brain (Right). 86.3% of the PBMC cells are classified as unknown using the brain reference. **(E)** Classification of the simulated dataset (described in Figure 1) using scMap with cell or cluster mapping. The color of each square labels the percentage of the actual annotation mapped to each scMap annotation. **(F)** Increased hybrid cell percentages with an increasing number of cell types in the reference, using different annotation complexity in the AMB Atlas. The AMB contains three levels of annotations: 1) Three cell types breaking down into GABAergic, Glutamatergic, and non-neuronal categories; 2) 20 cell types, placing cells into detailed cell categories, such as astrocytes and neural progenitor (NP) cells; 3) 108 cell types, breaking the individual cell types into their subtypes, such as astrocytes expressing *Aqp4* and NP cells expressing *Cpne7* or *Met*.

**Supplemental Figure 2. Application of Capybara to classify hematopoietic cell identity. Related to Figure 2 and Table S2 (A)** Comparison of log-normalized marker gene expression between classified cell types and other cells. The following genes are used as key markers for

different populations: HSPCs (Hematopoietic Stem and Progenitor Cells): *Cd34*; Megakaryocytes: *Itga2b*; Neutrophils: *Cebpe*; Monocytes: *Csf1r*; Erythrocytes: *Car2*. We group the cells into two categories. One group represents all the cells classified as a specific cell type, and the other group represents the remaining cells that are not classified as the same identity. The comparison between the two groups was performed and tested with a two-sample Wilcoxon Rank Sum Test, using the wilcox.test function in R (****: P <=0.0001). **(B)** *Top:* Projection of hybrid cell populations on the PAGA-guided clustering, along with their corresponding discrete identities. *Bottom*: Comparison of pseudotime between hybrids and their discrete identity counterparts. A Wilcoxon test was used for significance testing. **(C)** PAGA-guided clustering of the Paul et al. dataset, consisting of a total of 2,730 myeloid progenitors enriched from mouse bone marrow, revealing a total of 24 clusters. As expected, PAGA connectivity links the clusters containing hybrid cells but does not pinpoint the hybrid state. **(D)** Correlation between log (PAGA connectivity scores) and log (transition scores).

**Supplemental Figure 3. Evaluation of hybrid cells using ground-truth lineage tracing dataset. Related to Figure 3 and Table S3. (A)** Manual annotation of the major differentiated cell populations identified from the Weinreb et al., 2020 hematopoiesis lineage tracing dataset, projected onto the SPRING embedding. **(B)** Comparison between Capybara classification and the manual annotation. We selected the major differentiated myeloid cell types, including basophils, eosinophils, mast cells, monocytes, and neutrophils, as a reference. Cell types not included in the reference are correctly identified as 'Unknown.' 95.1% of cells with unknown identities are labeled as undifferentiated in the original Weinreb et al. annotation. The color of each square denotes the percentage agreement between the manual and Capybara-based annotations. **(C)** Integrative analysis of monocyte-neutrophil hybrids, monocytes, neutrophils (identified by Capybara), with IG2 (Irf8-GMP2 bistable intermediates), and GMP (Granulocyte Monocyte Progenitor) cells (identified in Olsson et al., 2016). IG2 cells have the potential to differentiate into monocytes or granulocytes. To evaluate monocyte-neutrophil hybrids, we used Seurat V4 to integrate monocytes, neutrophils, and monocyte-neutrophil hybrids with the IG2 and GMP populations, showing the overlap between monocyte-neutrophil hybrids and the IG2 population. **(D)** The color of each square denotes the percentage of each population in each cluster. **(E)** Cosine similarity of percentage cluster representations across all populations, showing the highest similarity between monocyte-neutrophil hybrids and the IG2 population. **(F)** Comparison between RNA velocity and Capybara transition scores in cardiac reprogramming (Stone et al., 2019). *Left*: RNA velocity vectors projected onto the UMAP embedding (Square:

area containing a small number of moving vectors; Circle: area containing a large number of moving vectors). *Middle*: Transition scores projected onto the UMAP embedding (Square: area corresponding to the RNA velocity plot showing low transition scores; Circle: area corresponding to the RNA velocity plot showing high transition scores). *Right*: correlation between log (velocity vector counts) and log (transition scores).

**Supplemental Figure 4. Capybara analysis of direct cardiac reprogramming. Related to Figure 4 and Table S4. (A)** Initial tissue-level classification of the Stone et al., 2019 dataset reveals four major tissues, which the high-resolution atlas restricts to two major tissues with the higher-resolution reference. **(B)** Normalized gene expression of cardiac markers labeling atrial vs. ventricular cardiomyocytes (****: P <=0.0001, Wilcoxon test). **(C)** Breakdown of cell types listed as "Other" in Figure 4. **(D)** Detailed hybrid cell-type breakdown for each time point of the Stone et al., cardiac reprogramming time course. **(E)** Integration of 10x dataset generated in this study with days 7 and 14 from the Stone et al., 2019 dataset, Cosine similarity = 0.804 between the independent studies. *Bottom right:* Projection of the two major cardiomyocyte populations, atrial and ventricular cardiomyocytes, onto the integrated UMAP. **(F)** RNA FISH of cells expressing only *Myh6* or *Myh7*. Scale bars = 50μm **(G)** Negative staining controls (MEFs) for RNA FISH and immunostaining. Scale bars = 10μm **(H)** Normalized gene expression of *Actc1* and *Tnnc1* across atrial, ventricular, and AV hybrids in the scRNA-seq data**. (I)** Immunofluorescence for MYL7 (atrial) and MYL2 (ventricular) proteins showing cells expressing a single protein (discrete) or co-expressing both proteins (hybrids) Scale bars = 10μm. **(J)** Quantification of discrete cells identified from immunostaining and scRNA-seq. MYL7-expressing cells are enriched relative to MYL2-expressing cells, confirming the atrial/ventricular bias observed from scRNA-seq data.
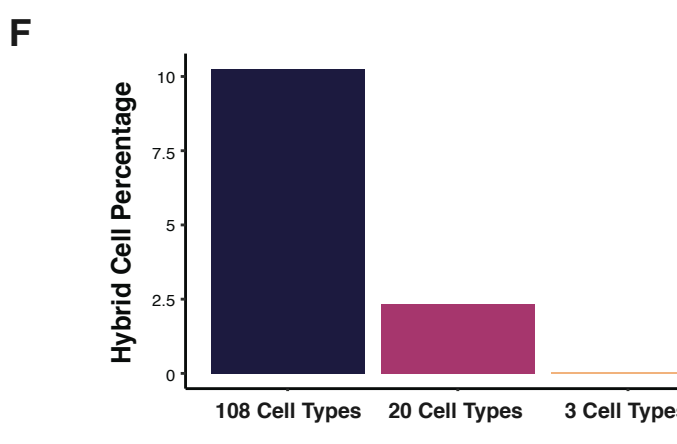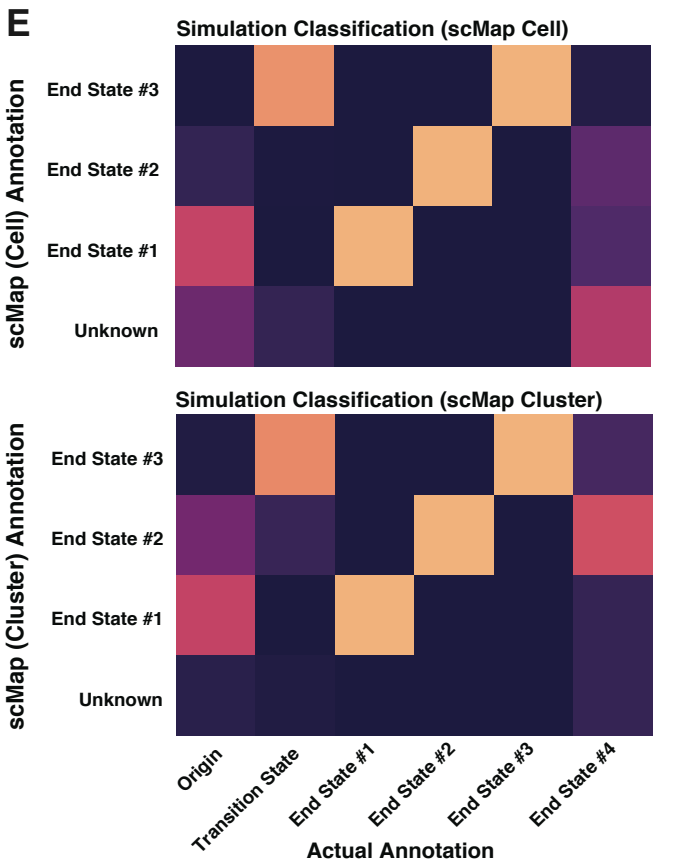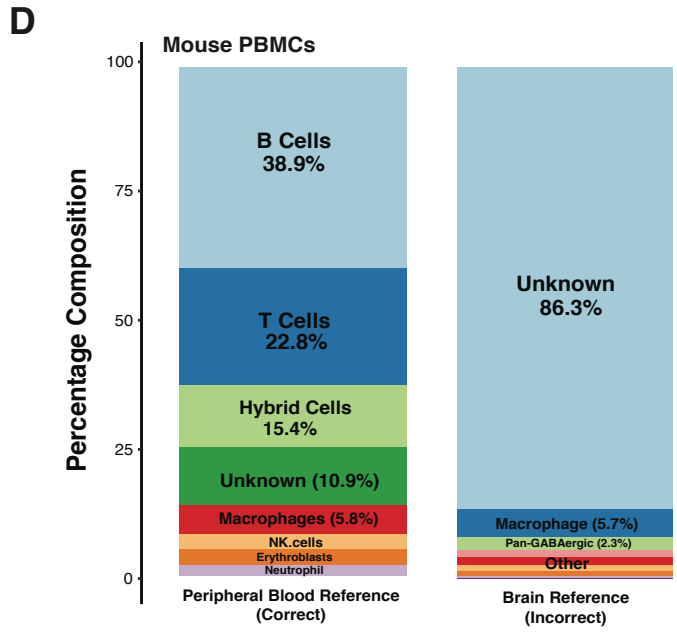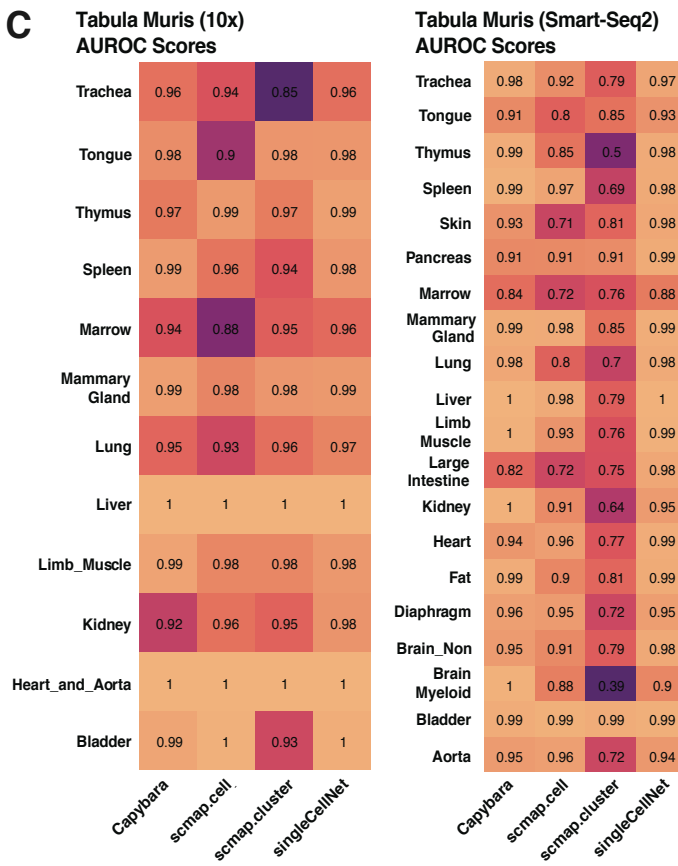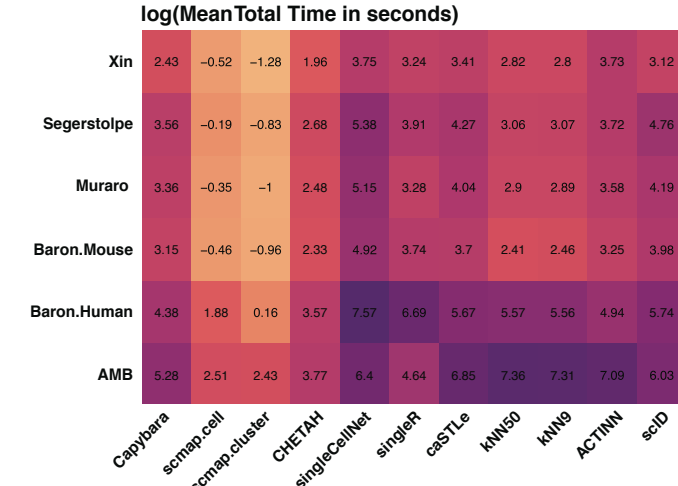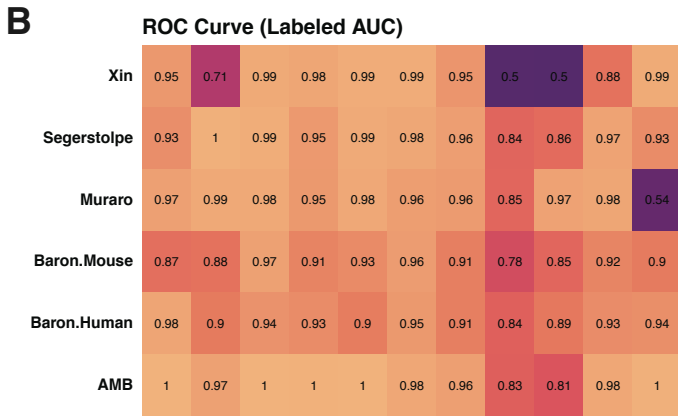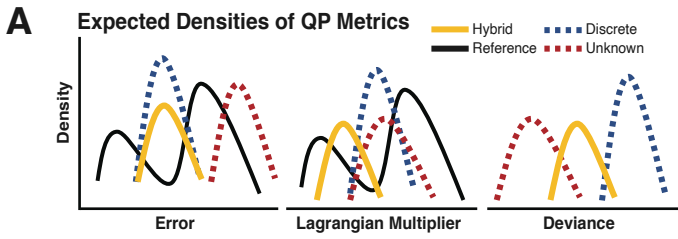
**Supplemental Figure 5. Capybara analysis of motor neuron differentiation and programming. Related to Figure 5 and Table S5. (A)** Capybara classification steps: with prior knowledge that these protocols aim to generate spinal motor neurons, we selected a single-cell spinal cord development atlas (Delile et al., 2019) as the high-resolution reference, omitting the general tissue selection step. We identified the major embryonic development stages corresponding to each protocol from Briggs et al., 2017. **(B)** UMAP embedding of this dataset, divided by protocol (Direct Programming: DP; Directed Differentiation: DD). *Top*: Projection of time points onto the UMAP embedding. *Bottom*: Projection of major discrete cell types, as identified via Capybara analysis, onto the UMAP embedding. **(C)** Transition scores for each protocol across experimental time points (****: P <=0.0001, n.s. = not significant, Wilcoxon test).

**(D-E)** Major hybrid populations in the direct programming and differentiation protocols. For each time point, we show the percentage of each hybrid type. **(F)** UMAP plot of the integrated datasets generated in this study, including four samples with different treatment groups. Expression of the motor neuron (MN) marker, *Mnx1* (left panel), and dorsal neuron marker, *Pou4f1* (right panel), are shown. **(G)** The major Capybara classification is shown for each cluster. **(H)** All discrete and hybrid cell type compositions for each treatment group.
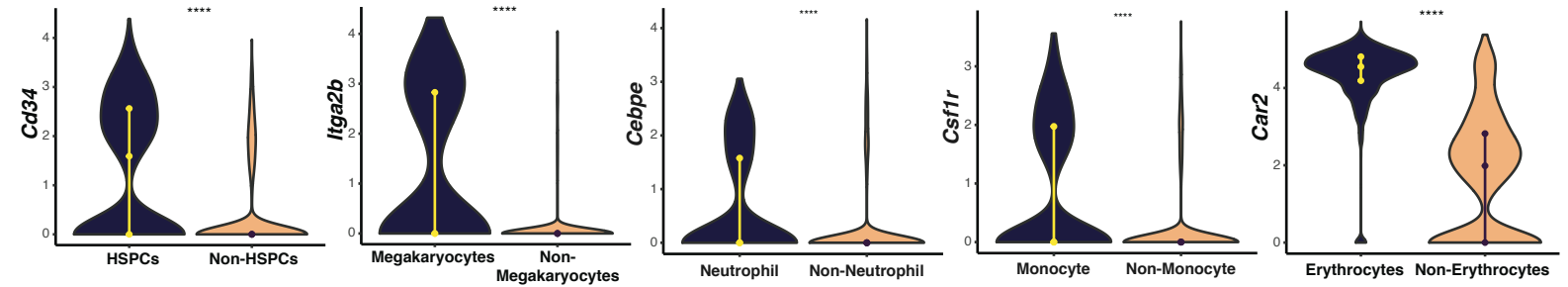
**Supplemental Figure 6. Capybara analysis of fibroblast to induced Endoderm Progenitor (iEP) Reprogramming. Related to Figure 6 and Table S6. (A)** Tissue-level classification of iEP reprogramming (Biddy et al., 2018) reveals seven major tissues. With the high-resolution reference, three relevant tissues are selected. **(B)** Cell-type classification of long-term cultured iEPs using an expanded reference with embryonic populations. We selected a foregut organogenesis atlas (Han et al., 2020) and a gut tube development atlas (Nowotschin et al., 2018). We combined the two references and performed Capybara analysis. 99.9% of cells were classified as 'unknown.' We then combined this endoderm development atlas with the regenerative liver atlas and MCA: the iEPs primarily classified as injured BECs. **(C)** Hybrid populations in day 28 and long-term iEPs. **(D)** 3D-rendering of a microscopy z-stack for 3D-cultured iEPs stained for DAPI, EpCAM, and CK19, demonstrating branching. Scale bars = 50-100 $\mu$m. **(E)** Expression of other BEC markers on the UMAP embedding, including *Krt19, Sox9,* and *Cftr*. **(F)** Module scores comparing identified injured BECs, MEF/Stromal cells (LT-iEP Dataset), and primary BECs (Pepe-Mooney et al., 2019). BEC markers used for module scores are identified from Verhulst et al., 2019 and are listed in **Table S6**.
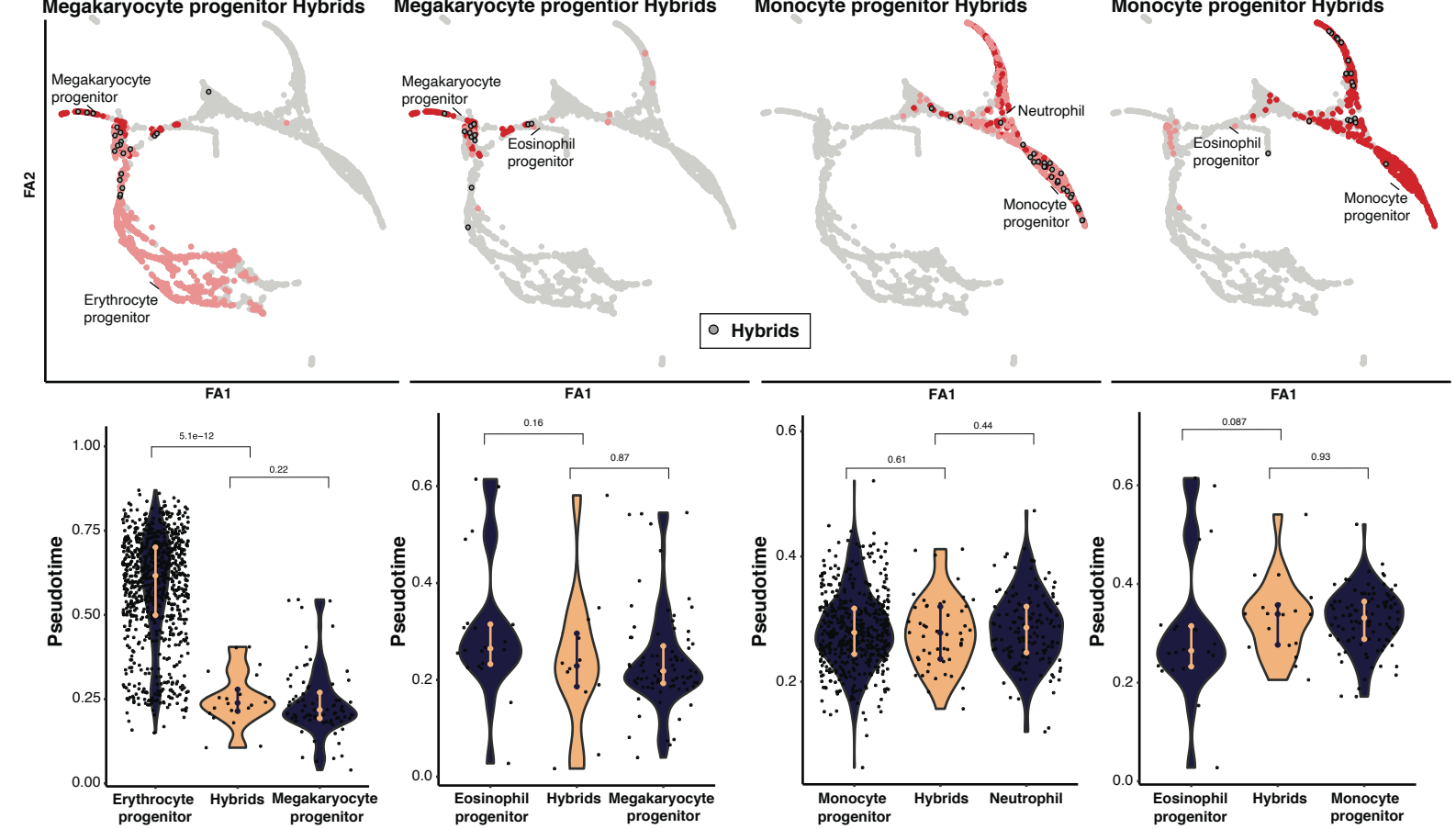
# Supplemental Figure 1

# Supplemental Figure 2

## A



## B

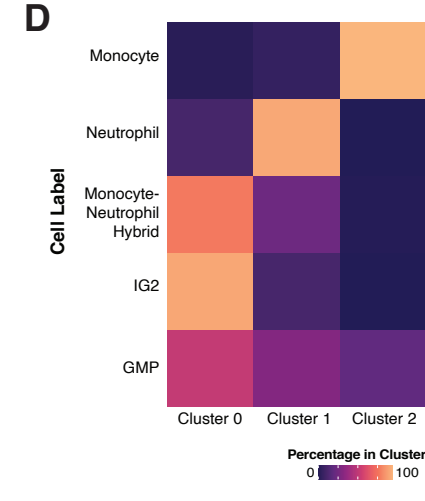Erythrocyte progenitor-Megakaryocyte progenitor Hybrids

Eosinophil progenitor-Megakaryocyte progenitor Hybrids

Neutrophil-Monocyte progenitor Hybrids

Eosinophil progenitor-Monocyte progenitor Hybrids



## C

PAGA Connectivity



## D

### Pearson Correlation = 0.84

# Supplemental Figure 3

**A** Weinreb differentiated populations



**B**



**C** Monocyte-Neutrophil/Hybrids (Weinreb et al., 2020) with IG2/GMP (Olsson et al., 2016)



**D**



**E**

|  | GMP | IG2 | Mono-Neutro Hybrid | Monocyte | Neutrophil |
|---|---|---|---|---|---|
| **GMP** | 1 | 0.87 | 0.93 | 0.35 | 0.58 |
| **IG2** | 0.87 | 1 | 0.99 | 0.02 | 0.24 |
| **Mono-Neutro Hybrid** | 0.93 | 0.99 | 1 | 0.04 | 0.40 |
| **Monocyte** | 0.35 | 0.02 | 0.04 | 1 | 0.05 |
| **Neutrophil** | 0.58 | 0.24 | 0.40 | 0.05 | 1 |

**F**

RNA Velocity Vectors          Transition Scores          Pearson Correlation = 0.77
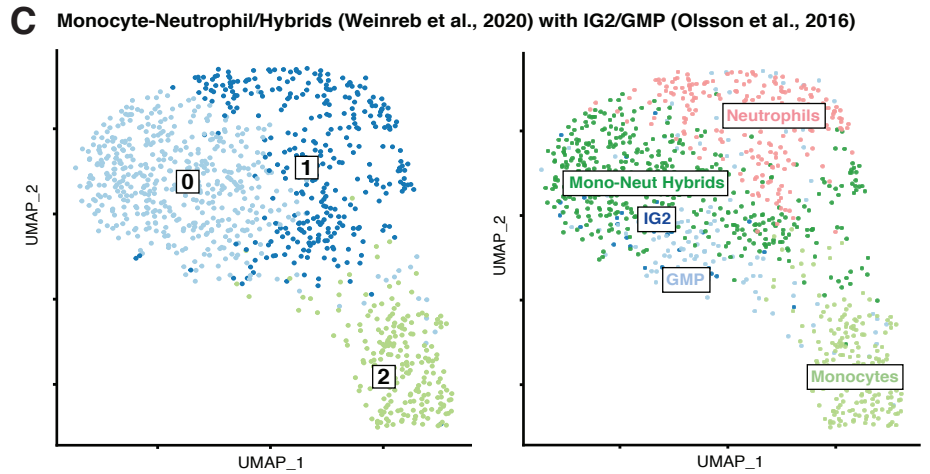
# Supplemental Figure 4

**A**

### Step 1: Tissue Level Classification

*Bulk Classification to MCA Tissues*



Stacked bar (Percentage Composition): Neonatal Skin, Fetal Stomach, Neonatal Heart, Fetal Lung, Other — *Cardiac Reprogramming*

### Step 2: Continuous Identity Measurement

*Single-cell resolution map*



Stacked bar (Percentage Composition): Neonatal Skin, Neonatal Heart, Fetal Lung, Other — *Cardiac Reprogramming*

**B**



Violin plots of Normalized Expression (Atrial vs Ventricular): *Myl12a* ****, *Acta2* ****, *Pfn1* ****, *Csrp2* ****, *Eef1a1* ****, *Arhgap31* ****

**C**



Percentage Composition, Other Cell Types: Brown Adipose Tissue, Endothelial Cell, Osteoblast, Epithelial Cell, Keratinocyte, Neuron, Adipocyte/Melanocyte/Acinar

**D** Stone et al. Hybrid Identities



Percentage of hybrid population across Day -1, Day 1, Day 2, Day 3, Day 7, Day 14 for hybrid identities: Muscle-Stromal 1, Muscle-Stromal 2, Muscle-Smooth Muscle, Macrophage1-Stromal 1, Macrophage1-Smooth Muscle, Macrophage1-Muscle, Erythroblast-Stromal, Macrophage1-Stromal 2, Smooth Muscle-Stromal, Artial CM-Muscle, Stromal3-Stromal2, DC-Macrophage1, Macrophage1-Macrophage2, DC-Macrophage2, BAT-DC, BAT-Cardiac Muscle, BAT-Ventricular CM, Artial CM-BAT, Artial CM-Ventricular CM

**E**



UMAP plots: Stone et al., Day 7; Stone et al., Day 14; iCM with SM (This Study); Target Cell Identity (Ventricular, Atrial). Axes UMAP_1, UMAP_2.

**F**



DAPI / *Myh6* (50 μm); DAPI / *Myh7* (50 μm)

**G**



DAPI, *Myl4*, *Actc1*, *Tnnc1*, *Myh6*, *Myh7*, DAPI, MYL7, MYL2 (10 μm each)

**H**



Violin plots: *Actc1* (Atrial, Ventricular, AV Hybrid); *Tnnc1* (Atrial, Ventricular, AV Hybrid)

**I**



Atrial, Ventricular, AV hybrid: DAPI, MYL7, MYL2, Merge (10 μm each)

**J**



MYL7 & MYL2 IF: Percent positive cells for IF (*) and scRNA-seq (n.s.); legend MYL7, MYL2

# Supplemental Figure 5

**A** Step 1: Tissue Level Classification

Step 2: Continuous Identity Measurement

**B** (Briggs et al., 2017) Direct Programming (DP) — Direct Differentiation (DD)

**C** DP Transition Scores — DD Transition Scores

**D** DP Hybrids

**E** DD Hybrids

**F** Mnx1 (MN) — Pou4f1 (Dorsal)

**G** This Study (Ngn2+Isl1+Lhx3 +/- RA/SAG)

**H** Discrete Identities — Hybrid Identities

# Supplemental Figure 6

## A

### Step 1: Tissue Level Classification

**Bulk Classification to MCA Tissues**



Percentage Composition

- Embryonic Mesenchyme
- Fetal Lung
- Embryonic Stem Cells
- Trophoblast Stem Cells
- Neonatal Skin
- Fetal Intestine
- Fetal Stomach
- Other

*iEP Reprogramming*

### Step 2: Continuous Identity Measurement

**Single-cell resolution map**



Percentage Composition

- MEF
- Neonatal Skin
- Embryonic Mesenchyme
- Fetal Stomach
- Fetal Liver
- Neonatal Pancreas
- Other

*iEP Reprogramming*

## B Expanded references



Han et al, 2020 & Nowotschin et al, 2019 Endoderm Atlas

Percentage Composition

Unknown (99.7%)

Discrete (0.3%)

*Long-term iEPs*

Endoderm Atlas + Regenerative Liver + MCA

Percentage Composition

Discrete (Injured BEC) (95.7%)

Hybrid (4%)

Other Discrete (0.3%)

*Long-term iEPs*

## C



Percentage Composition

**Day 28**

**Long-term iEPs**

- Injured BEC-MEF
- Injured BEC-Stromal Cell
- Normal BEC-Stromal Cell
- Injured BEC-Injured Hepatocyte
- Injured BEC-Normal BEC

## D



DAPI | EpCAM | CK19 | Merge

## E



*Krt19* | *Sox9* | *Ctfr*

UMAP_2 / UMAP_1

Expression

## F



BEC Module score

p < 2.22e−16

p < 2.22e−16

- MEF/Stromal
- BEC−Like iEPs
- Primary BECs