

# Supplementary Material for “JIND: Joint Integration and Discrimination for Automated Single-Cell Annotation”

Mohit Goyal, Guillermo Serrano, Ilan Shomorony, Mikel Hernaez and Idoia Ochoa

## Contents

### List of Tables

S1	Comparing performances of JIND and JIND+ when the source and target batches are integrated with Seurat (Integrated) versus when Seurat integration is not performed (Batched). <i>raw</i> is the initial accuracy of the classifier, <i>rej</i> is the percentage of cells rejected by the classifier and <i>eff</i> is the effective accuracy after rejecting unconfident predictions. Best raw accuracy rates among the two cases (batched or integrated) are <b>boldfaced</b> . . . . .	4
S2	Impact of choosing different number of genes for training the prediction model. NN is the neural network based prediction model used by JIND and doesn't perform any batch alignment. JIND performs adversarial training to minimize distributional mismatch between source and target batches. JIND+ further fine-tunes assuming the confident predictions are correct on the target batch. <i>raw</i> is the initial accuracy of the classifier, <i>rej</i> is the percentage of cells rejected by the classifier and <i>eff</i> is the effective accuracy after rejecting unconfident predictions. . . . .	4
S3	Comparison position of ReLU activation w.r.t encoder. Case 1: When ReLU activation is a part of the classifier subnetwork and is applied on the latent code produced by the encoder subnetwork. Case 2: When ReLU activation is a part of encoder subnetwork applied before at the end of encoder to produce the latent code. JIND performs adversarial training to minimize distributional mismatch between source and target batches. JIND+ further fine-tunes assuming the confident predictions are correct on the target batch. <i>raw</i> is the initial accuracy of the classifier, <i>rej</i> is the percentage of cells rejected by the classifier and <i>eff</i> is the effective accuracy after rejecting unconfident predictions. . . . .	5

## List of Figures

S1	Heatmap for all differentially expressed genes between two groups on <i>PBMC</i> 10x_v5 dataset. <i>Monocytes FCGR3A</i> cells predicted by JIND+ as: <i>Monocytes FCGR3A</i> cells (G1) or <i>Monocytes CD14</i> cells (G2). The clustering of the cells is calculated with the entire set of genes. The results of this DE analysis can be found on the Supplementary Excel File. . . . .	6
S2	Heatmap for all differentially expressed genes between two groups on <i>PBMC</i> 10x_v5 dataset. <i>Monocytes FCGR3A</i> cells predicted by JIND+ as: <i>Monocytes FCGR3A</i> cells (G1) or <i>Monocytes CD14</i> cells (G2). The results of this DE analysis can be found on the Supplementary Excel File. . . . .	7
S3	Heatmap for all differentially expressed genes among two randomly chosen groups of <i>Acinar</i> cells present in <i>Pancreas</i> Mur16 dataset. The results of this DE analysis can be found on the Supplementary Excel File. . . . .	8
S4	Heatmap for all differentially expressed genes among two randomly chosen <i>Monocyte FCGR3A</i> cells present in <i>PBMC</i> 10x_v5 dataset. The results of this DE analysis can be found on the Supplementary Excel File. . . . .	9
S5	NN-based prediction model employed by JIND for cell identification. The network consists of two subnetworks, an encoder and a classifier, which are jointly trained. The input to the model is a vector (of dimension 5000 by default) containing the gene expression data for a cell, and the output consists of a probability vector indicating the likelihood of the cell belonging to each of the $K$ classes. . . . .	10

## Additional results

- Figure S1 shows different gene expression patterns across the two groups, G1: *Ductal* cells predicted as *Ductal* and , G2: *Ductal* cells predicted as *Acinar*, allowing differentiation between the two populations of cells. On the contrary, when DE analysis was performed on the groups, G1: randomly chosen subset of *Ductal* cells and , G2: remaining *Ductal* cells, (Figure S3) we neither observe meaningful clustering nor descriptive gene expression patterns necessary for differentiation. We can see the same results on the *PBMC* dataset on Figure S2 and Figure S4
- Table S1 shows the comparison between two cases, Batched: when JIND and JIND+ are evaluated

on datasets containing batch effects versus, Integrated: when JIND and JIND+ are run after Seurat integration on the same datasets.

Datasets		PBMC 10x_v3-10x_v5		Pancreas Bar16-Mu16		Pancreas Bar16-Seg16	
		Batched	Integrated	Batched	Integrated	Batched	Integrated
<i>JIND</i>	<i>raw</i>	<b>0.971</b>	0.968	<b>0.958</b>	0.946	<b>0.987</b>	0.946
	<i>rej</i>	0.07	0.06	0.05	0.10	0.05	0.08
	<i>eff</i>	0.986	0.985	0.974	0.979	0.997	0.979
<i>JIND+</i>	<i>raw</i>	<b>0.974</b>	0.971	0.959	<b>0.961</b>	<b>0.992</b>	0.961
	<i>rej</i>	0.03	0.03	0.03	0.09	0.02	0.05
	<i>eff</i>	0.985	0.978	0.971	0.980	0.997	0.980

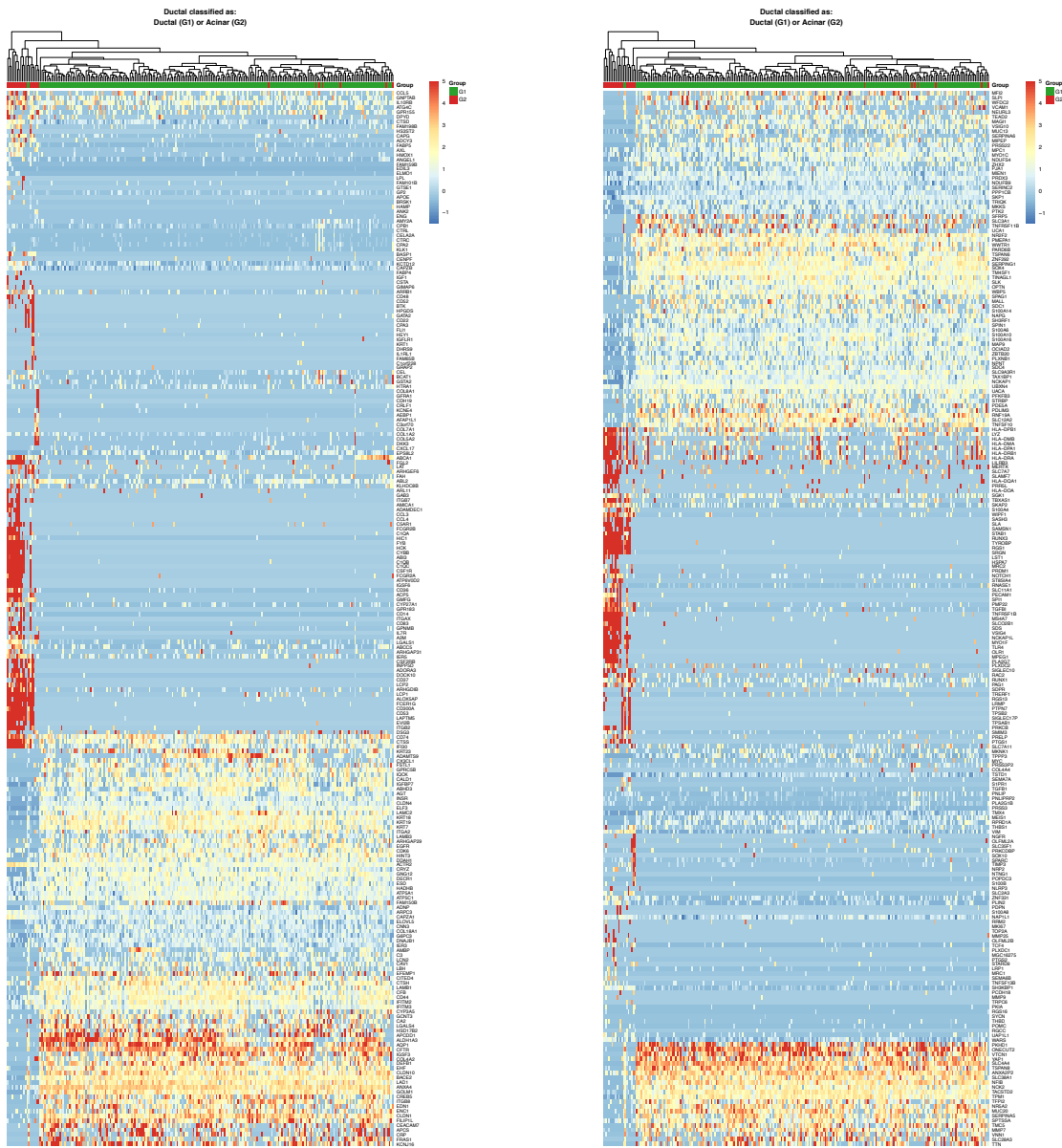
**Table S1:** Comparing performances of JIND and JIND+ when the source and target batches are integrated with Seurat (Integrated) versus when Seurat integration is not performed (Batched). *raw* is the initial accuracy of the classifier, *rej* is the percentage of cells rejected by the classifier and *eff* is the effective accuracy after rejecting unconfident predictions. Best raw accuracy rates among the two cases (batched or integrated) are **boldfaced**.

Methods	Datasets	Human-Hemato				Mouse Cortex			
	Metrics/#Genes	1000	3000	5000	7000	1000	3000	5000	7000
NN	<i>raw</i>	0.923	0.938	0.928	0.933	0.975	0.976	0.982	0.981
	<i>rej</i>	0.06	0.06	0.06	0.06	0.05	0.05	0.05	0.06
	<i>eff</i>	0.945	0.957	0.950	0.954	0.988	0.990	0.991	0.997
JIND	<i>raw</i>	0.922	0.939	0.927	0.931	0.976	0.970	0.982	0.981
	<i>rej</i>	0.07	0.06	0.06	0.06	0.06	0.03	0.05	0.06
	<i>eff</i>	0.944	0.957	0.948	0.953	0.989	0.982	0.991	0.996
JIND+	<i>raw</i>	0.919	0.931	0.927	0.932	0.976	0.970	0.976	0.978
	<i>rej</i>	0.04	0.04	0.04	0.04	0.02	0.03	0.04	0.04
	<i>eff</i>	0.933	0.945	0.943	0.948	0.985	0.950	0.984	0.990

**Table S2:** Impact of choosing different number of genes for training the prediction model. NN is the neural network based prediction model used by JIND and doesn't perform any batch alignment. JIND performs adversarial training to minimize distributional mismatch between source and target batches. JIND+ further fine-tunes assuming the confident predictions are correct on the target batch. *raw* is the initial accuracy of the classifier, *rej* is the percentage of cells rejected by the classifier and *eff* is the effective accuracy after rejecting unconfident predictions.

Methods	Datasets	PBMC 10x_v3-10x_v5		Pancreas Bar16-Mur16		Pancreas Bar16-Seg16	
		Case 1	Case 2	Case 1	Case 2	Case 1	Case 2
JIND	<i>raw</i>	0.971	0.972	0.958	0.881	0.987	0.86
	<i>rej</i>	0.07	0.07	0.05	0.39	0.05	0.44
	<i>eff</i>	0.986	0.988	0.974	0.987	0.997	0.9948
JIND+	<i>raw</i>	0.974	0.974	0.959	0.8805	0.992	0.901
	<i>rej</i>	0.03	0.03	0.03	0.27	0.02	0.23
	<i>eff</i>	0.985	0.984	0.971	0.979	0.997	0.984

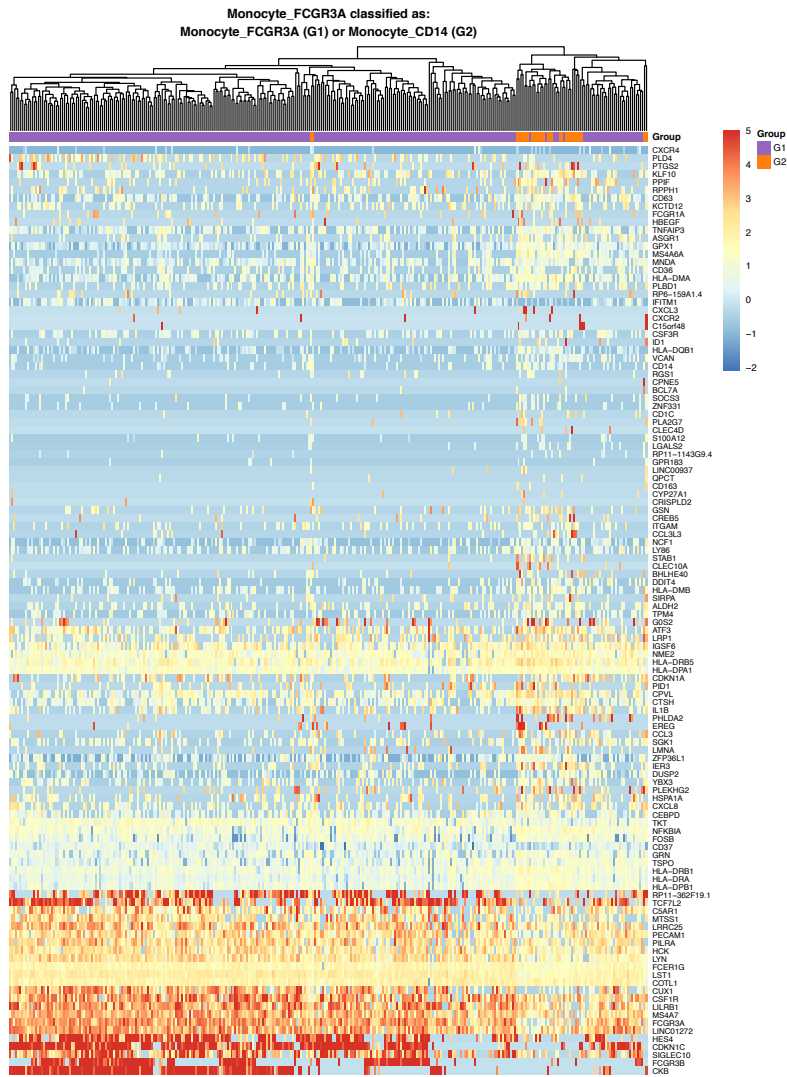
**Table S3:** Comparison position of ReLU activation w.r.t encoder. Case 1: When ReLU activation is a part of the classifier subnetwork and is applied on the latent code produced by the encoder subnetwork. Case 2: When ReLU activation is a part of encoder subnetwork applied before at the end of encoder to produce the latent code. JIND performs adversarial training to minimize distributional mismatch between source and target batches. JIND+ further fine-tunes assuming the confident predictions are correct on the target batch. *raw* is the initial accuracy of the classifier, *rej* is the percentage of cells rejected by the classifier and *eff* is the effective accuracy after rejecting unconfident predictions.



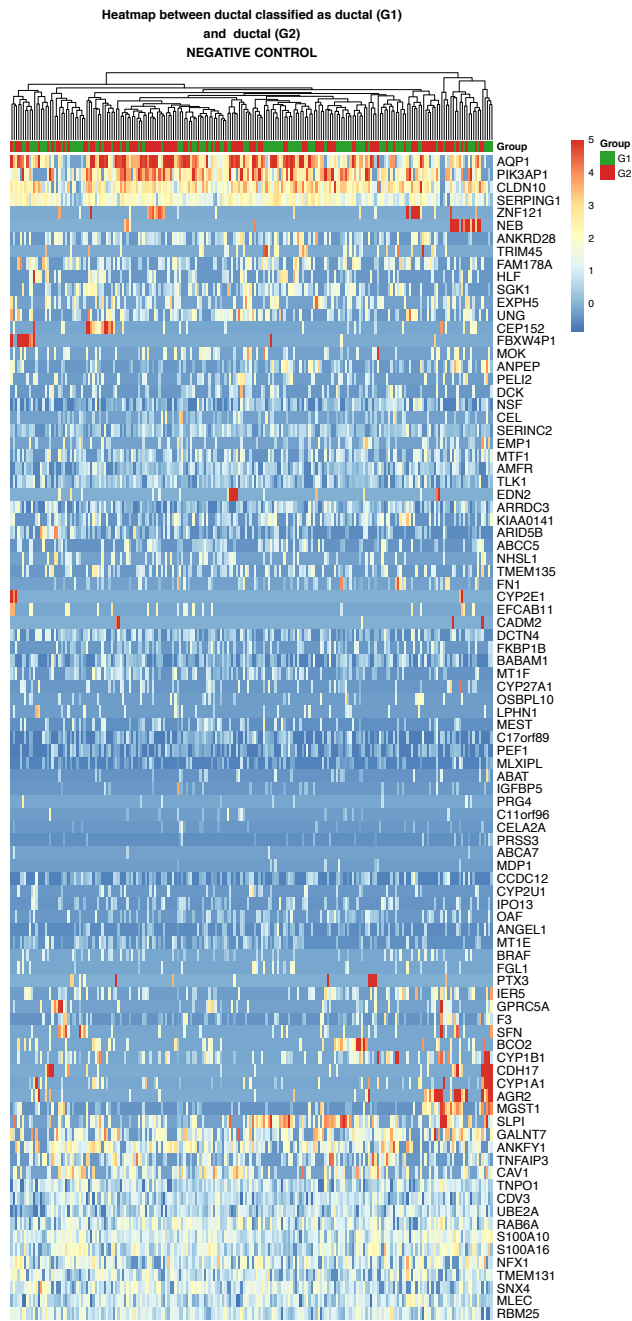
(a) First half of the DEG

(b) Second half of the DEG

**Figure S1:** Heatmap for all differentially expressed genes between two groups on *PBMC* 10x\_v5 dataset. *Monocytes FCGR3A* cells predicted by JIND+ as: *Monocytes FCGR3A* cells (G1) or *Monocytes CD14* cells (G2). The clustering of the cells is calculated with the entire set of genes. The results of this DE analysis can be found on the Supplementary Excel File.



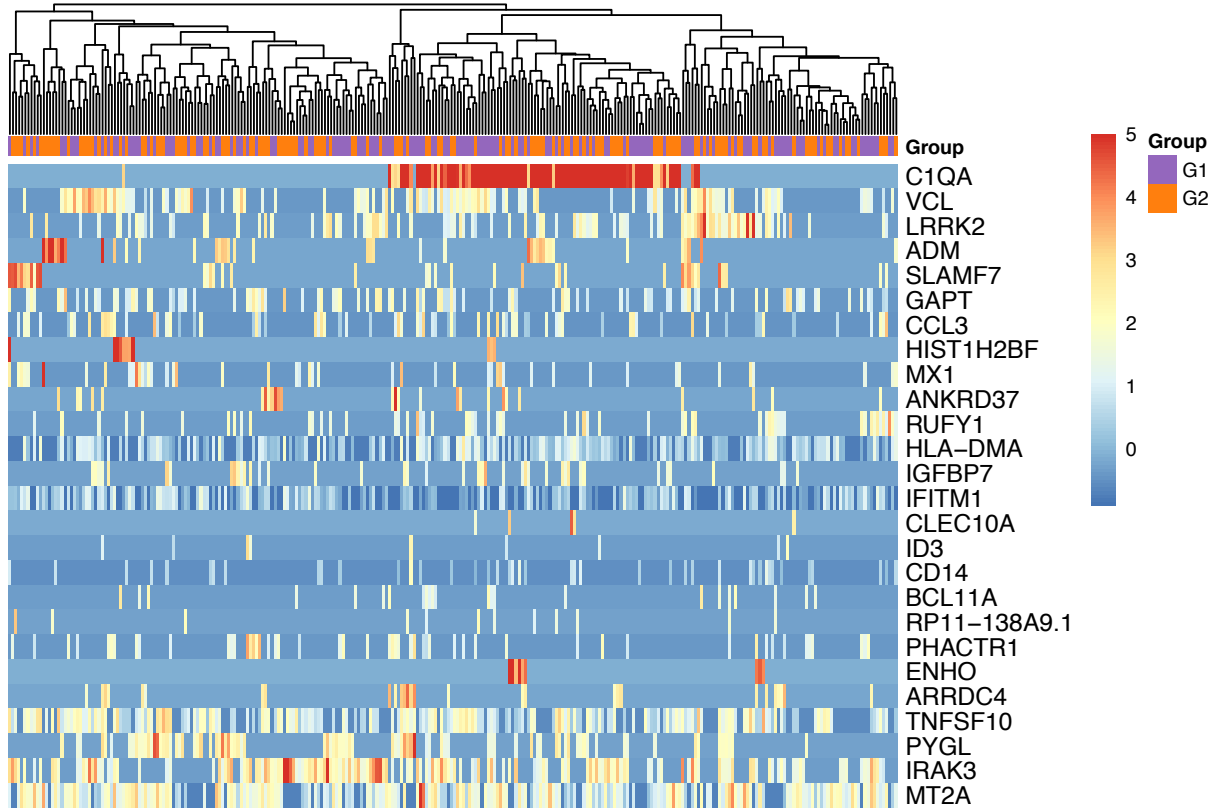
**Figure S2:** Heatmap for all differentially expressed genes between two groups on *PBMC* 10x\_v5 dataset. *Monocytes FCGR3A* cells predicted by JIND+ as: *Monocytes FCGR3A* cells (G1) or *Monocytes CD14* cells (G2). The results of this DE analysis can be found on the Supplementary Excel File.



**Figure S3:** Heatmap for all differentially expressed genes among two randomly chosen groups of *Acinar* cells present in *Pancreas* Mur16 dataset. The results of this DE analysis can be found on the Supplementary Excel File.

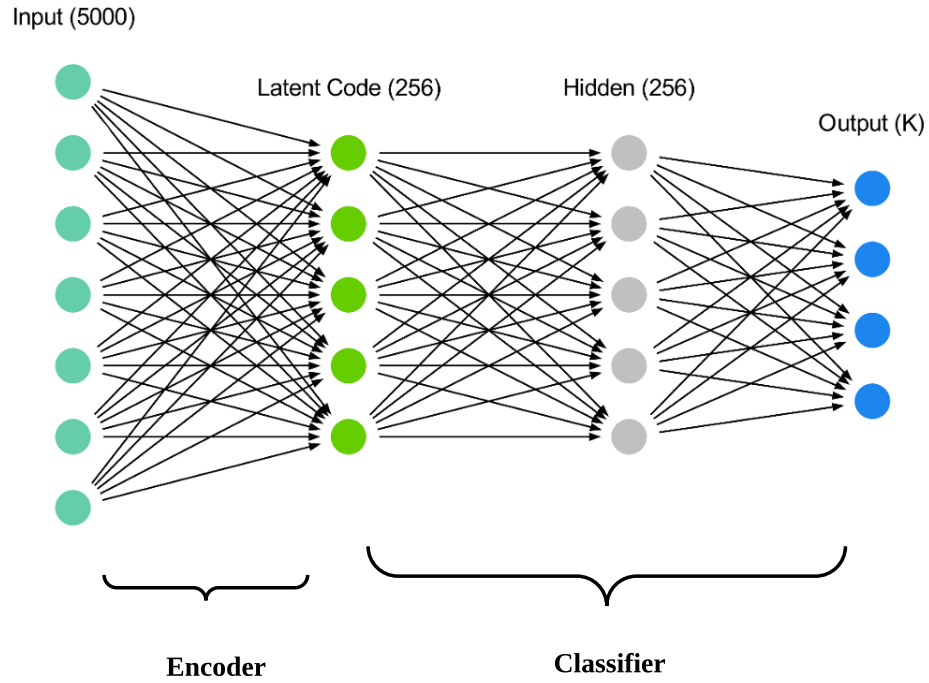


Heatmap between ductal classified as Monocyte\_FCGR3A (G1)  
 Monocyte\_FCGR3A (G2)  
 NEGATIVE CONTROL



**Figure S4:** Heatmap for all differentially expressed genes among two randomly chosen *Monocyte FCGR3A* cells present in *PBMC 10x\_v5* dataset. The results of this DE analysis can be found on the Supplementary Excel File.

## Neural Network based Prediction Model



**Figure S5:** NN-based prediction model employed by JIND for cell identification. The network consists of two subnetworks, an encoder and a classifier, which are jointly trained. The input to the model is a vector (of dimension 5000 by default) containing the gene expression data for a cell, and the output consists of a probability vector indicating the likelihood of the cell belonging to each of the  $K$  classes.