



IEEE



# RIA: a novel Regression-based Imputation Approach for single-cell RNA sequencing

The 11th IEEE International Conference on Knowledge and Systems Engineering (KSE 2019)  
Da Nang, Vietnam

## Authors

**Bang Tran**

*Computer Science &  
Engineering  
University of Nevada, Reno  
Reno, USA*

**Duc Tran**

*Computer Science &  
Engineering  
University of Nevada, Reno  
Reno, USA*

**Hung Nguyen**

*Computer Science &  
Engineering  
University of Nevada, Reno  
Reno, USA*

**Tin Nguyen**

*Computer Science &  
Engineering  
University of Nevada, Reno  
Reno, USA*

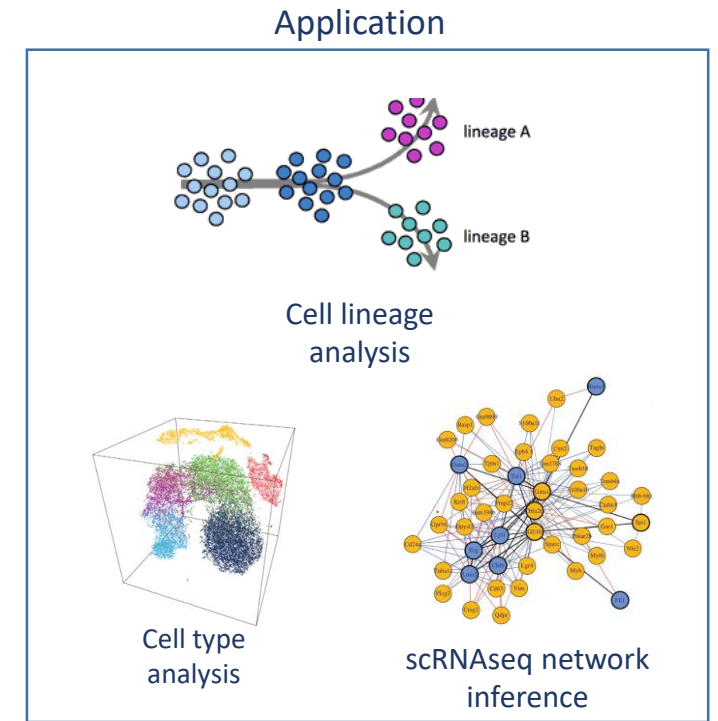
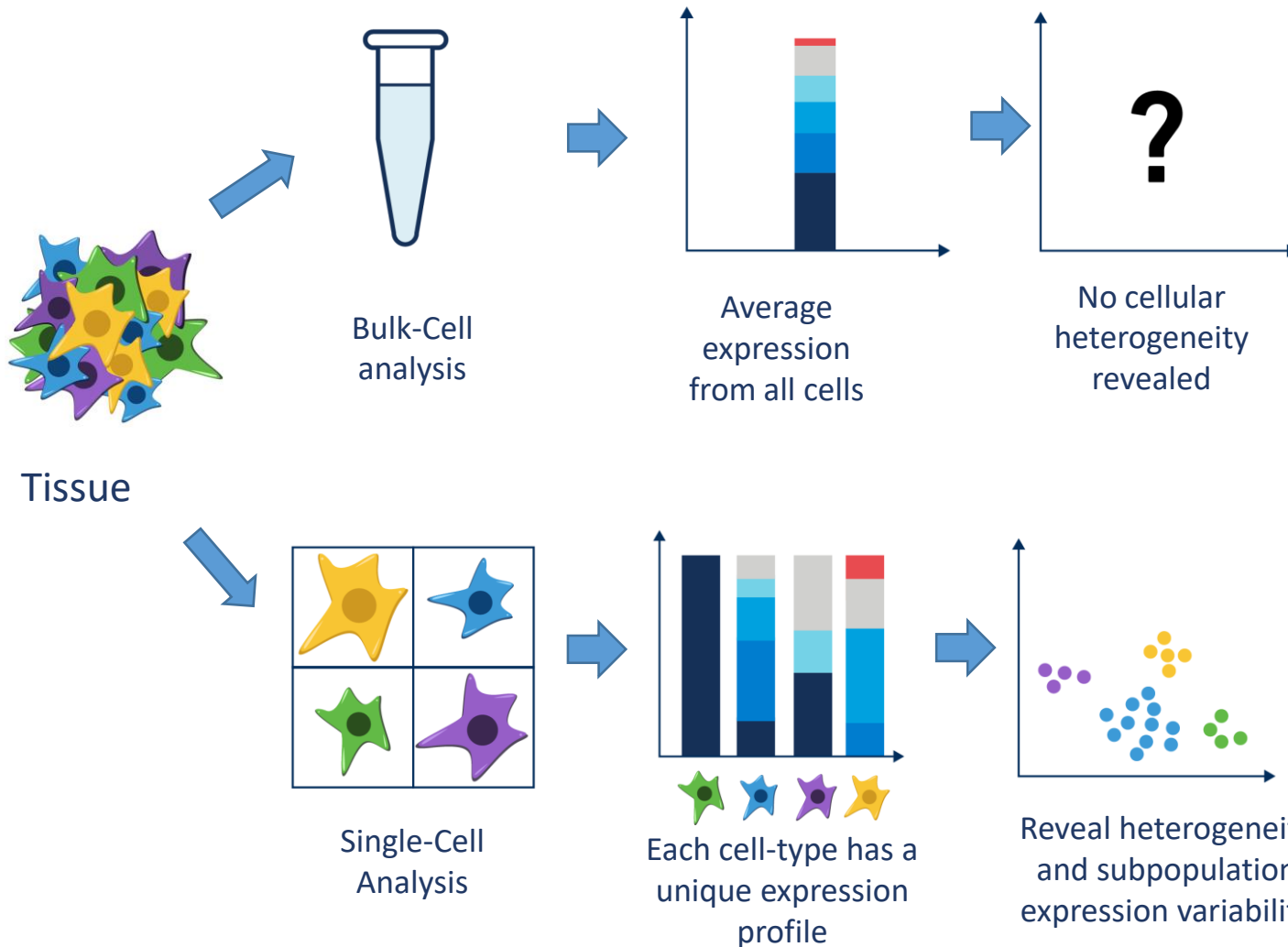
**Nam Sy Vo**

*Computational  
Biomedicine  
Vingroup Big Data  
Institute*

# Background

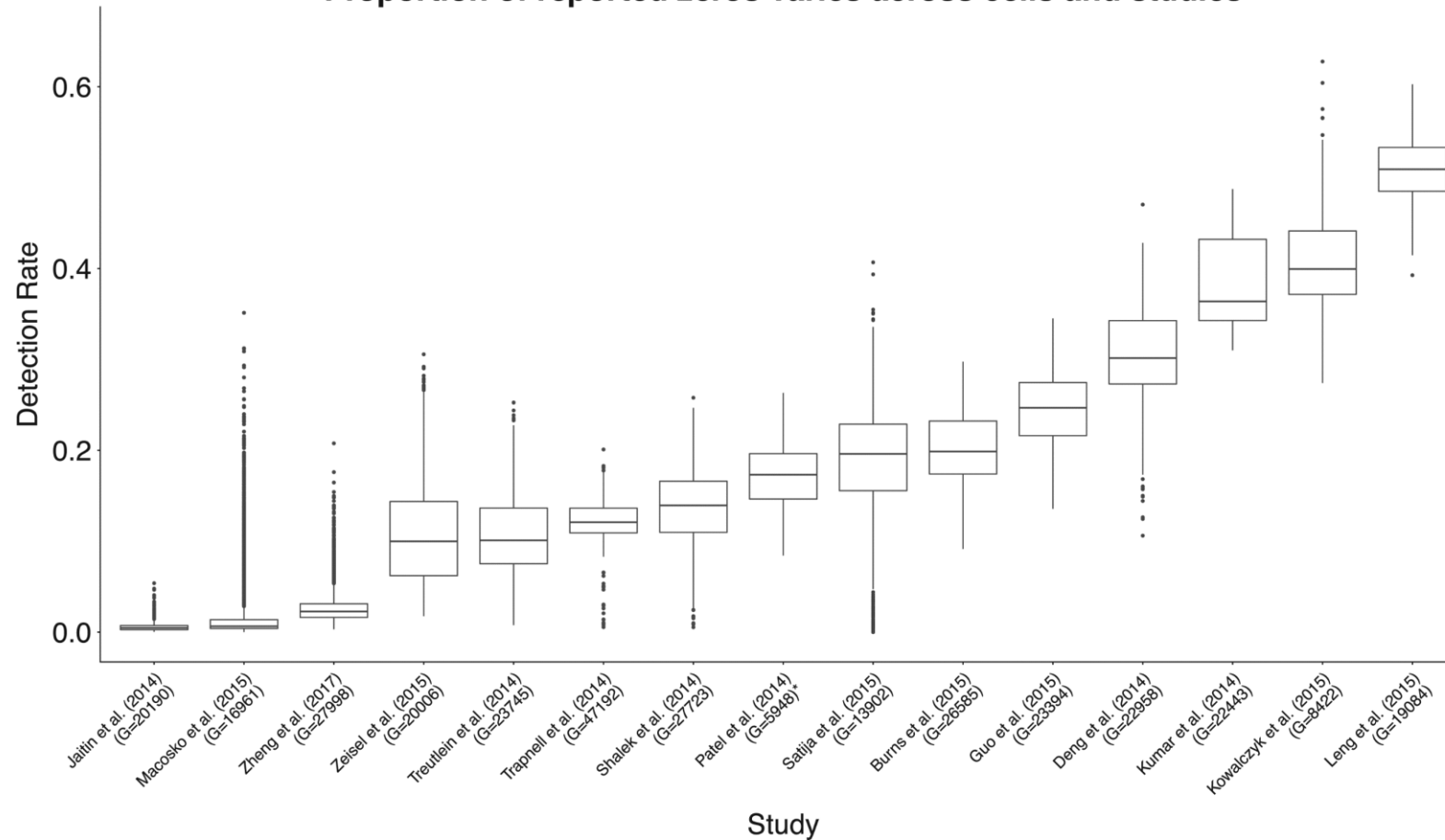
---

# Single – Cell RNA Sequencing



Source: 10xgenomics.com

Proportion of reported zeros varies across cells and studies

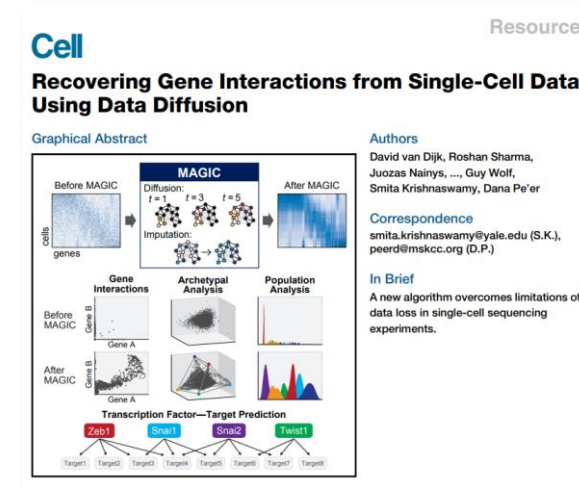


Dropout introduces zero expression value of genes in datasets.

# Single – Cell RNAseq Imputation State of The Art

## MAGIC [van Dijk et al., 2017]:

- **Citation:** 123
- **Lab:** Krishnaswamy Lab - Yale
- Using Markov Affinity-based Graph Imputation



## scImpute [Kwak et al., 2017]:

- **Citation:** 104
- **Lab:** UCLA
- Using statistical approach.

## An accurate and robust imputation method scImpute for single-cell RNA-seq data

Wei Vivian Li & Jingyi Jessica Li

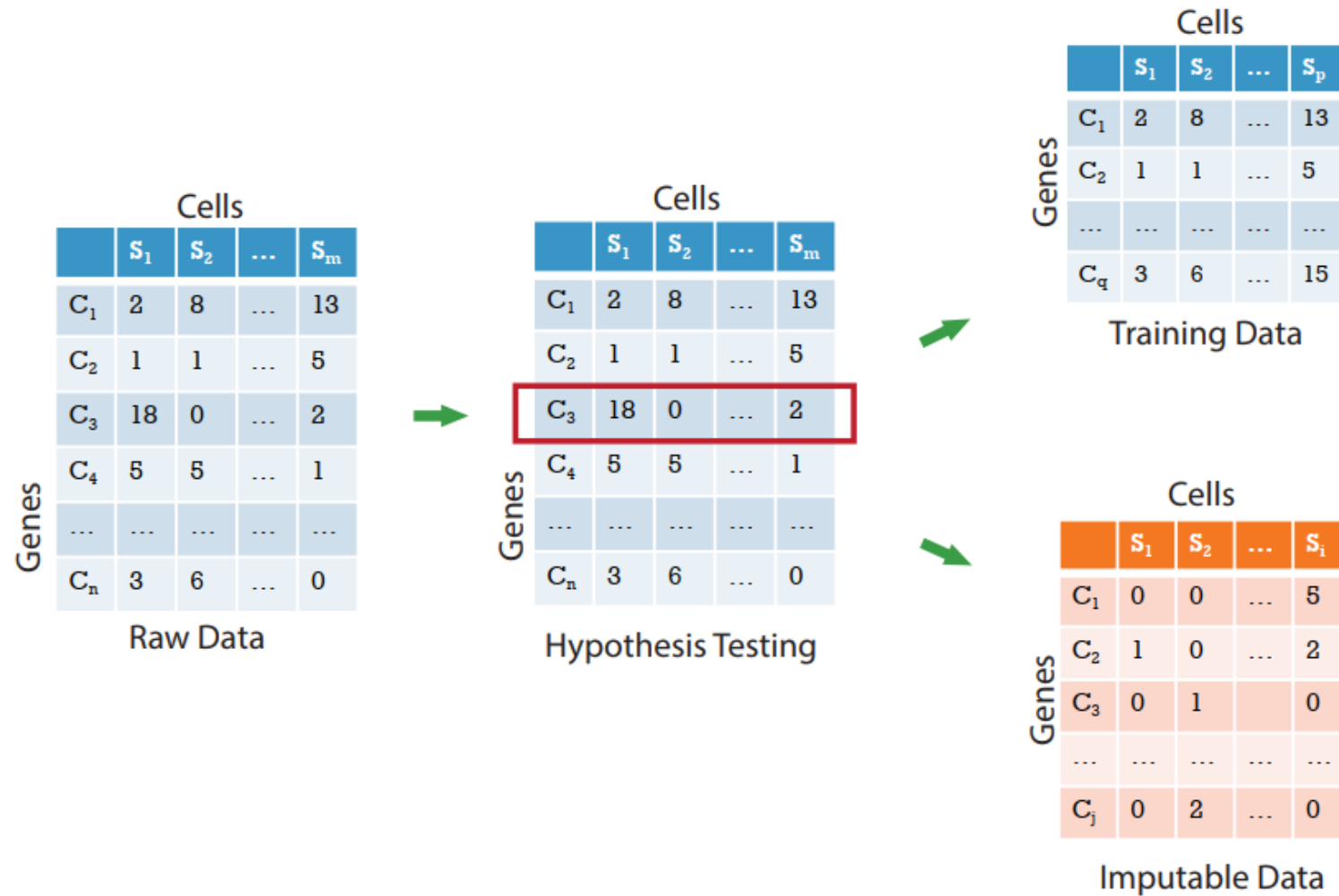
Nature Communications 9, Article number: 997 (2018) | Download Citation

## Method: RIA a novel Regression-based Imputation Approach for single-cell RNA sequencing

---

- Identifying genes that are impacted by dropouts
- Impute using genes with high confidence

# Hypothesis Testing and Dropout Identification





# Generalized Linear Regression Model

Cells

	$s_1$	$s_2$	...	$s_p$
$C_1$	2	8	...	13
$C_2$	1	1	...	5
...	...	...	...	...
$C_q$	3	6	...	15

Genes

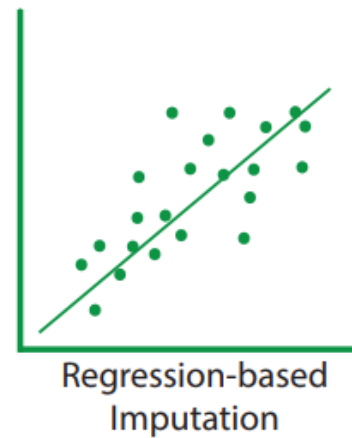
Training Data

Cells

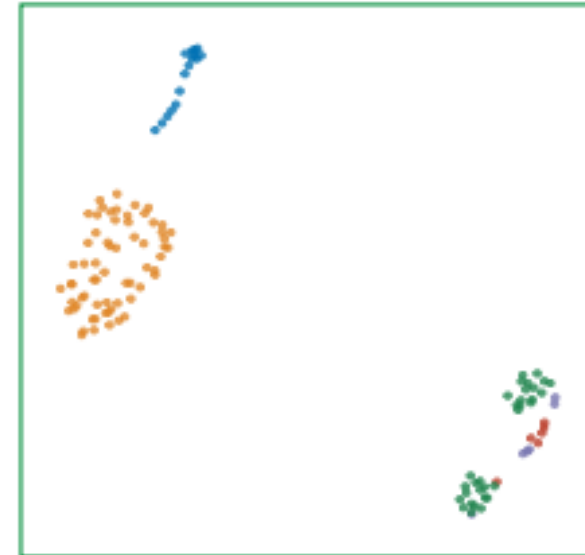
	$s_1$	$s_2$	...	$s_i$
$C_1$	0	0	...	5
$C_2$	1	0	...	2
$C_3$	0	1	...	0
...	...	...	...	...
$C_j$	0	2	...	0

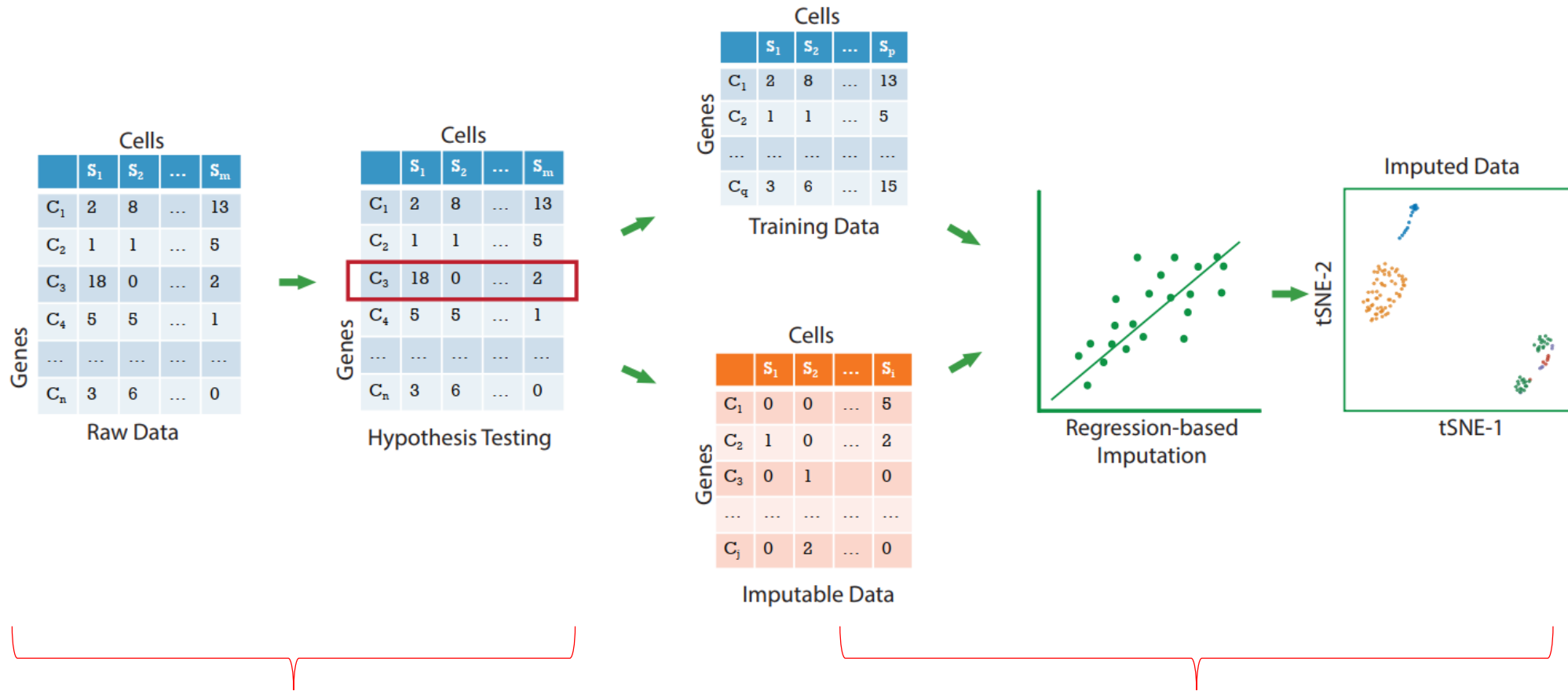
Genes

Imputable Data



Imputed data





Hypothesis testing approach

Generalized linear model

# Results

---

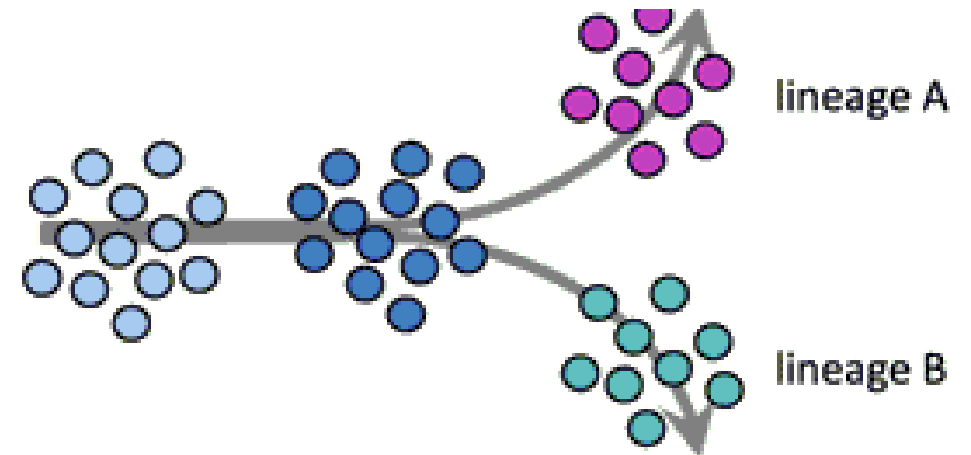
- MAGIC
- scImpute

TABLE I  
SINGLE-CELL DATA OBTAINED FROM NIH GEO

Dataset	Accession ID	Size	K	Organism	Protocol
Biase[27]	GSE57249	49	4	Mouse Embryo	SMARTer
Yan[28]	GSE36552	90	6	Human Embryo	Tang
Goolam[29]	E-MTAB-3321	124	5	Mouse Embryo	Smart-Seq2
Deng[30]	GSE45719	268	6	Mouse Embryo	Smart-Seq2
Zeisel[31]	GSE60361	3,005	9	Mouse Brain	STRT-Seq



Cell type clustering



Time-trajectory inference

TABLE III  
COMPARISON USING JACCARD INDEX

Dataset	Jaccard Index			
	Raw	RIA	scImpute	MAGIC
Biase	0.589	<b>0.708</b>	0.339	0.289
Yan	0.498	<b>0.498</b>	0.473	0.146
Goolam	0.496	<b>0.892</b>	0.375	0.312
Deng	0.524	<b>0.781</b>	0.395	0.518
Zeisel	0.651	<b>0.683</b>	0.605	0.285

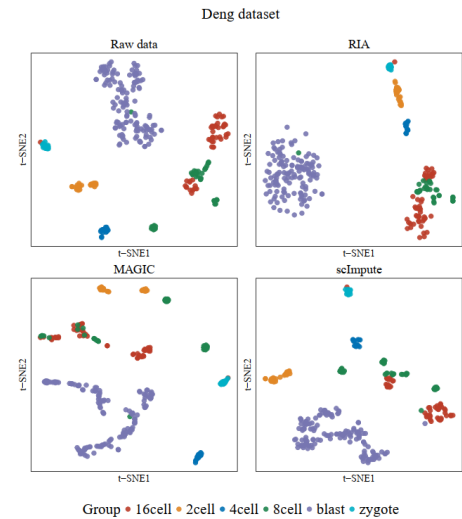
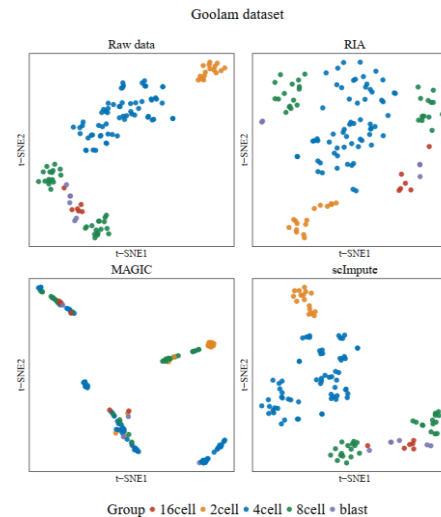
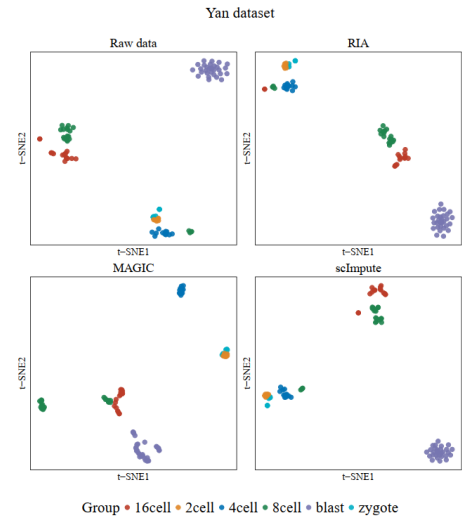
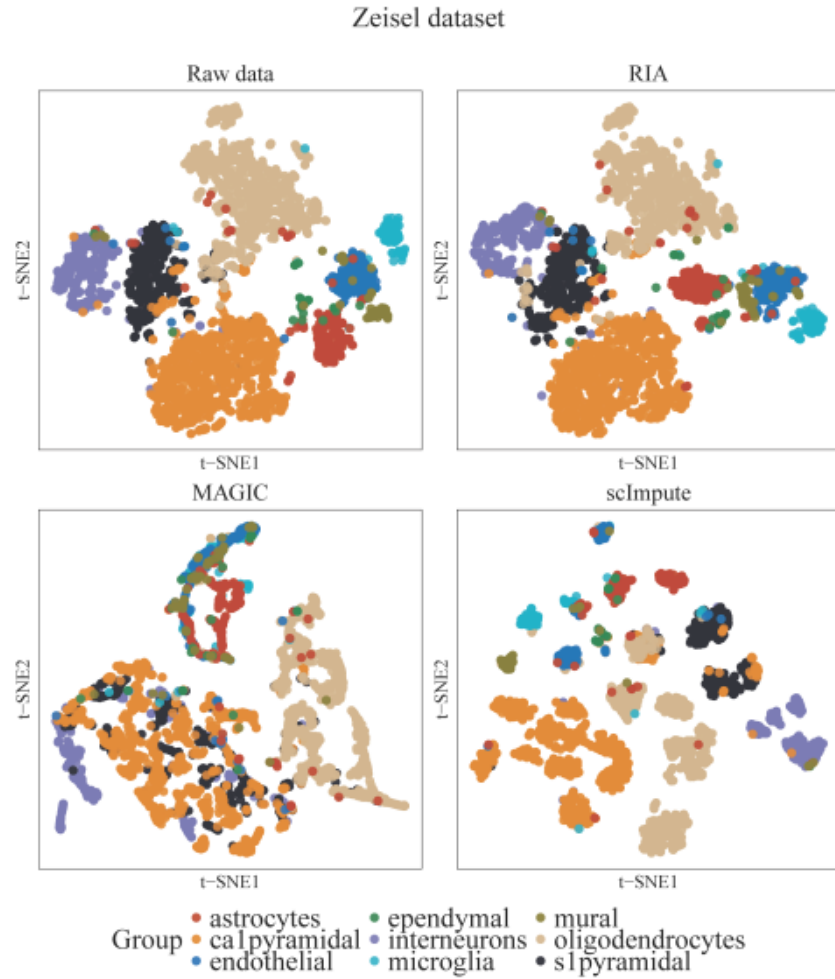
TABLE IV  
COMPARISON USING PURITY INDEX

Dataset	Purity Index			
	Raw	RIA	scImpute	MAGIC
Biase	0.795	<b>0.836</b>	0.449	0.612
Yan	0.711	<b>0.778</b>	0.733	0.467
Goolam	0.822	<b>0.952</b>	0.693	0.621
Deng	0.805	<b>0.839</b>	0.627	0.750
Zeisel	0.876	<b>0.893</b>	0.840	0.668

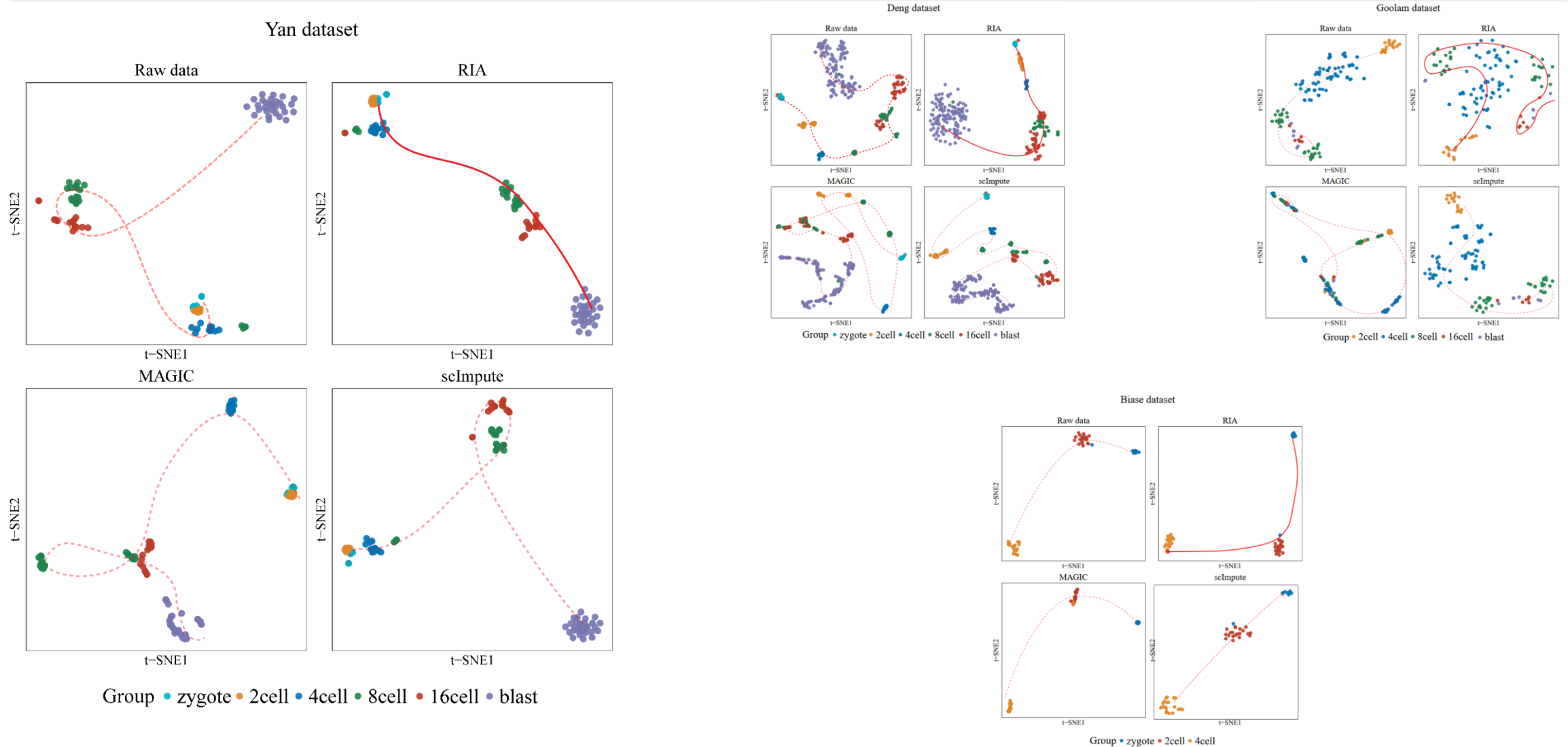
TABLE II  
COMPARISONS USING ADJUSTED RAND INDEX (ARI).

Dataset	Adjusted Rand Index			
	Raw	RIA	scImpute	MAGIC
Biase	0.558	<b>0.711</b>	-0.009	0.154
Yan	0.558	<b>0.573</b>	0.507	0.029
Goolam	0.501	<b>0.914</b>	0.321	0.197
Deng	0.549	<b>0.815</b>	0.229	0.483
Zeisel	0.738	<b>0.768</b>	0.689	0.289

# RIA Preserves Original Transcriptome Landscape



# RIA Improves Time Trajectory Inference for TSCAN



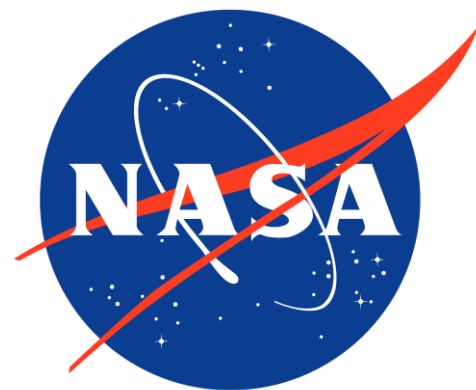
- RIA presents a new method to recover missing values caused by dropout events in scRNA-seq data.
- RIA uses a statistical hypothesis testing to identify the set of genes that are likely to be affected by dropouts.
- RIA imputes missing values by using highly correlated genes that share similar biological characteristics.
- RIA dramatically outperforms existing state-of-the-art approaches in improving the identification of cell populations.
- RIA is able to recover temporal trajectories in embryonic development stages.
- RIA is a fast method that can impute thousands of cells with tens of thousands of genes in minutes.



# Acknowledgements



People in Bioinformatics Lab at  
University of Nevada Reno



This work was partially supported by the  
National Aeronautics and Space  
Administration (NASA) under Grant  
Number 80NSSC19M0170

[1]. Kharchenko, P. V., Silberstein, L., & Scadden, D. T. (2014). Bayesian approach to single-cell differential expression analysis. *Nature methods*, 11(7), 740.

# Q & A

**Thank you!**