



UCLA

Statistical Methods for Bulk and Single-cell RNA Sequencing Data

Jingyi Jessica Li

Associate Professor

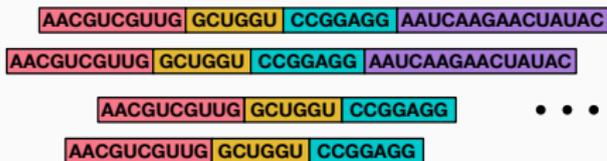
Department of Statistics

University of California, Los Angeles

<http://jsb.ucla.edu>

RNA sequencing (RNA-seq) technology

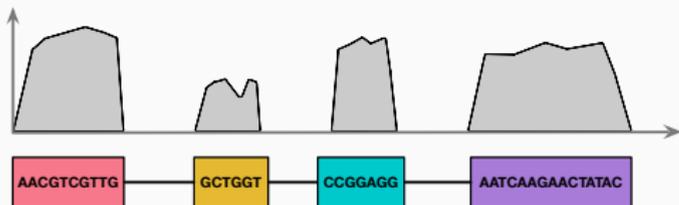
full length RNA isoforms
(unknown)



RNA-seq
experiments

statistical
inference

RNA-seq data
(observed)



RNA sequencing (RNA-seq) experiment

full length RNA isoforms
(1712 bp on average)

AACGUCGUUG GCUGGU CCGGAGG AAUCAAGAACUAUAC ...
AACGUCGUUG GCUGGU CCGGAGG AAUCAAGAACUAUAC

↓ fragmentation

RNA fragments
(< 600 bp)

AACGUCG UUG GCUGGU CCGG AGG AAUCAAGAACUAUAC ...
AACGUCGUUG GCUGGU CCGGAGG AAUC AAGAACUAUAC

RNA sequencing (RNA-seq) experiment

full length RNA isoforms
(1712 bp on average)

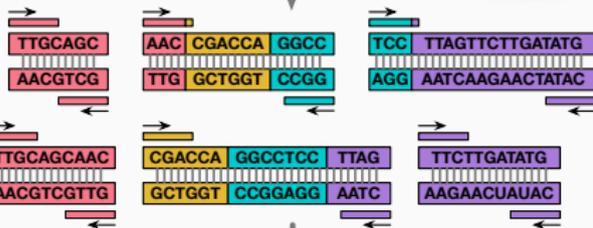
AACGUCGUUG GCUGGU CCGGAGG AAUCAAGAACUAUAC ...
AACGUCGUUG GCUGGU CCGGAGG AAUCAAGAACUAUAC

fragmentation

RNA fragments
(< 600 bp)

AACGUCG UUG GCUGGU CCGG AGG AAUCAAGAACUAUAC ...
AACGUCGUUG GCUGGU CCGGAGG AAUC AAGAACUAUAC

processing sequencing



RNA-seq reads
(< 300 bp)

AACG --- CAGC TTTG --- GGCC AGGA --- TATG ...
AACG --- CAAC GCTG --- TTAG AAGA --- TATG

of RNA-seq reads \propto isoform abundance \times isoform length

Mapping RNA-seq reads to the reference genome

full length RNA isoforms
(1712 bp on average)

AACGUCGUUG GCUUGU CCGGAGG AAUCAAGAACUAUAC ...
AACGUCGUUG GCUUGU CCGGAGG AAUCAAGAACUAUAC

processing

sequencing

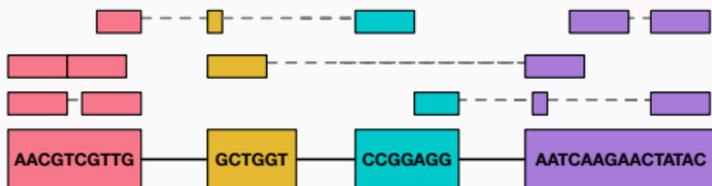


RNA-seq reads
(< 300 bp)

AACG --- CAGC --- TTG G --- GGCC --- AGG A --- TATG ...
AACG --- CAAC --- GCTG --- TTAG --- AAGA --- TATG

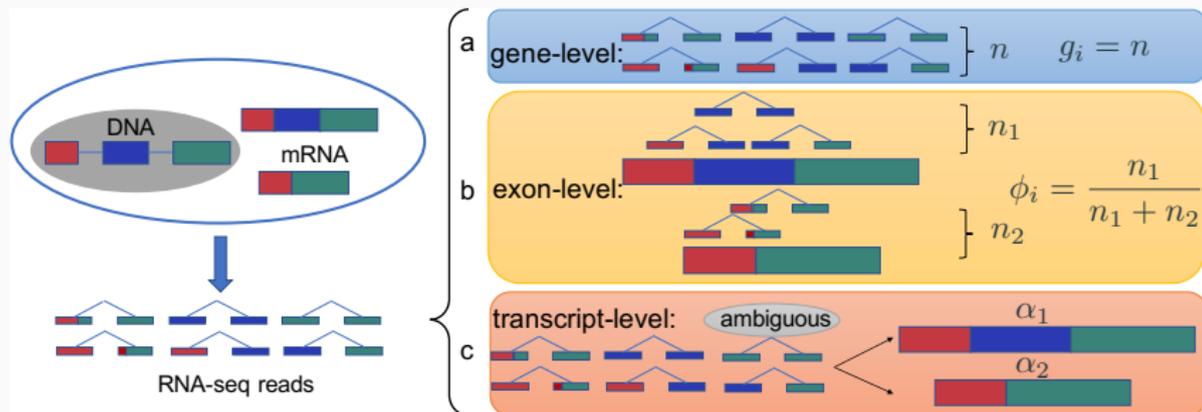
mapping (alignment)

RNA-seq reads
aligned to genome

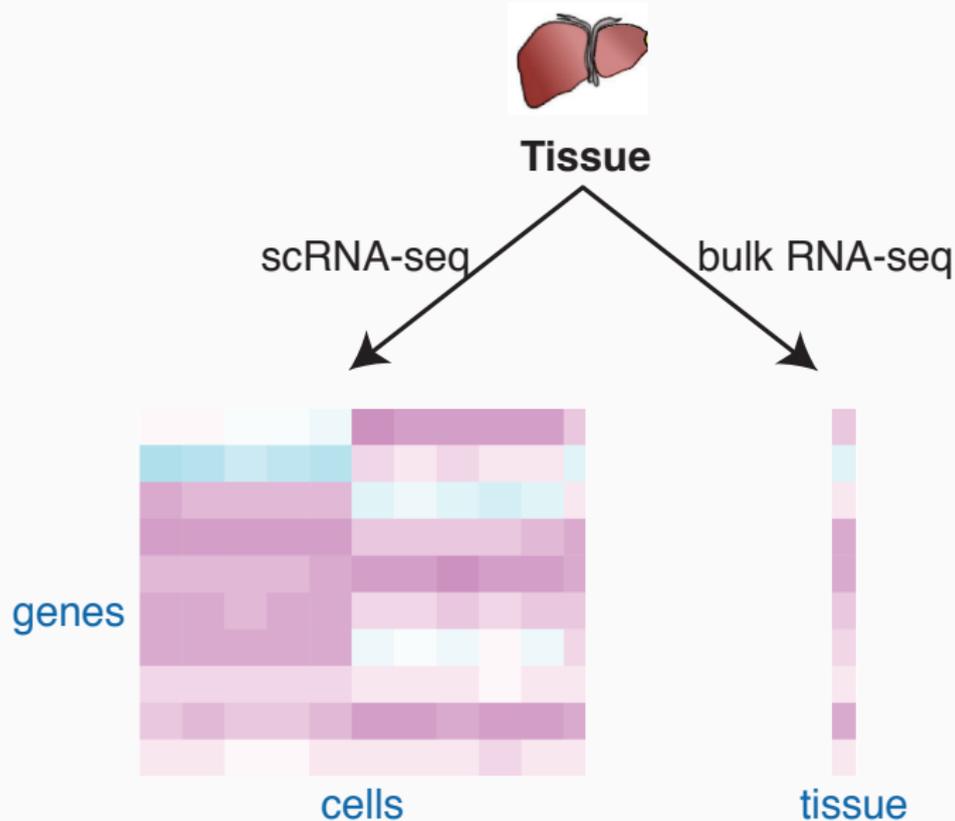


Reference-based RNA-seq data analysis

1. Align RNA-seq reads to a reference genome
2. Analyze aligned reads at three levels

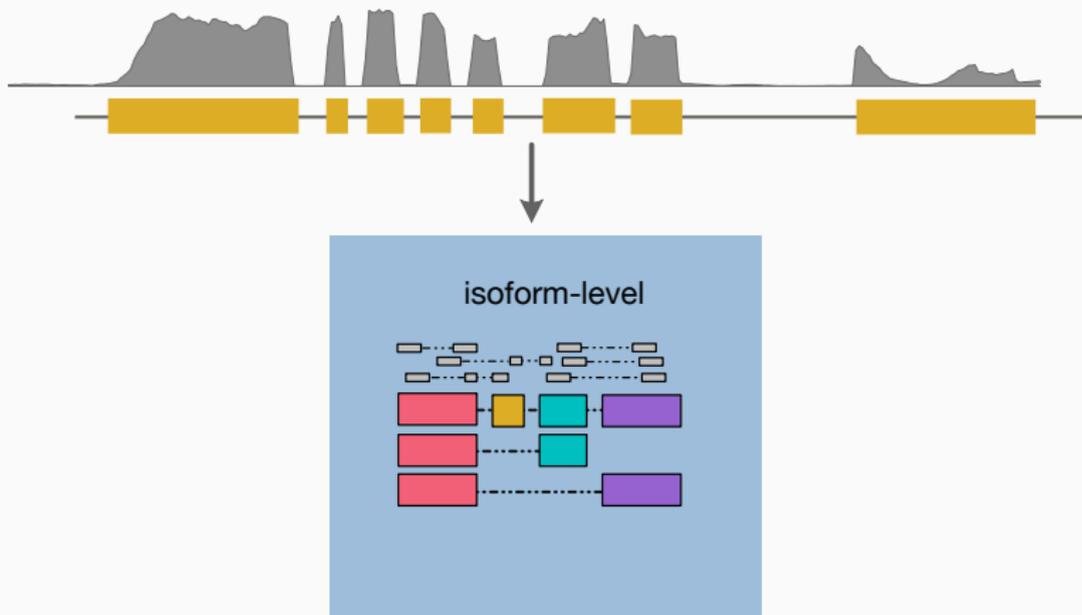


Single-cell (sc) vs. bulk RNA-seq at the gene level



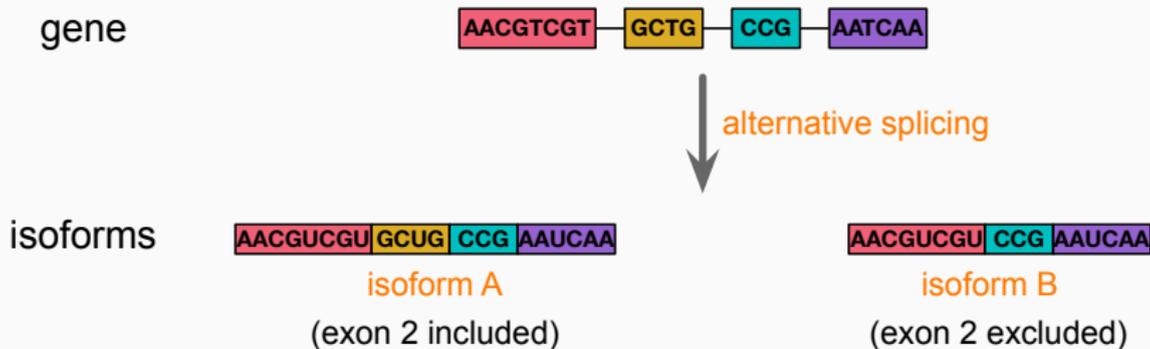
Bulk RNA-seq: transcript/isoform discovery & quantification

AIDE: annotation-assisted isoform discovery



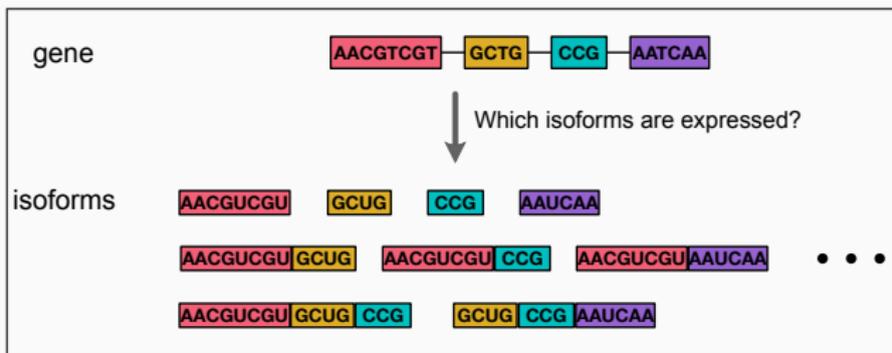
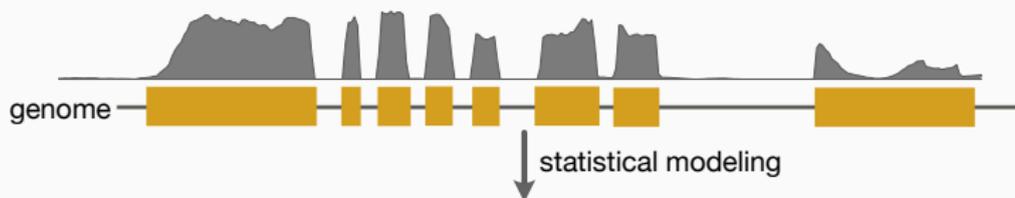
Isoform discovery: which isoforms are expressed?

- More than 90% genes undergo alternative splicing in mammals [Hooper, *Human Genomics*, 2014].
- At least 35% genetic diseases involve abnormal splicing [Manning et al., *Nature Reviews Mol. Cell Biol.* 2017].



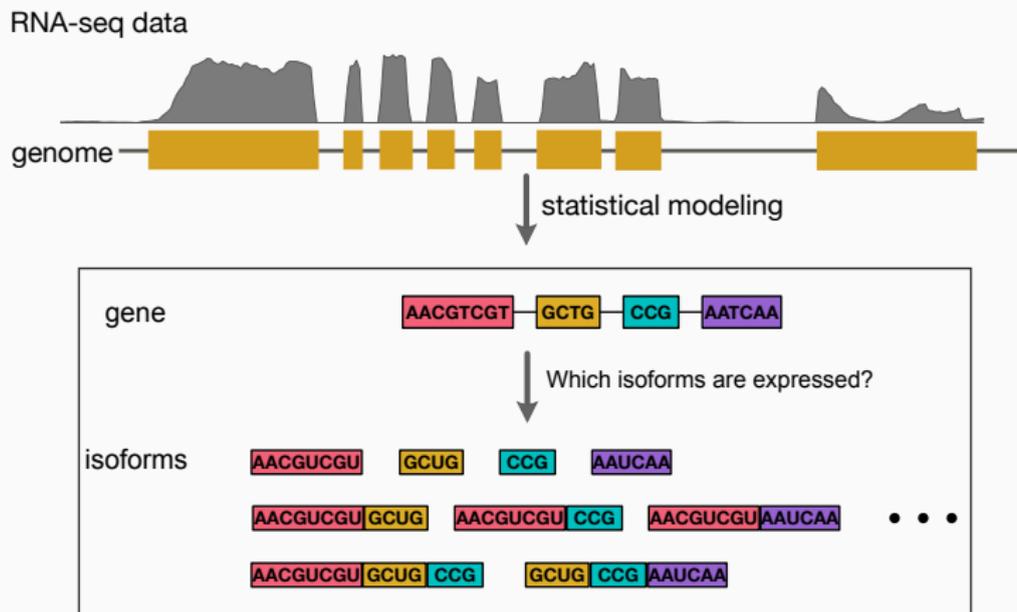
Isoform discovery: which isoforms are expressed?

RNA-seq data



Challenge 1: large number of candidate isoforms

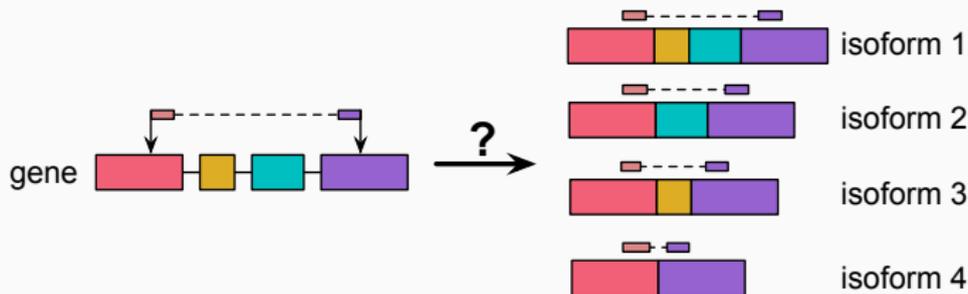
Variable size (# of candidate isoforms) = $2^{\# \text{ of exons}} - 1$



For this 4-exon gene, $2^4 - 1 = 15$ candidate isoforms

Challenge 2: great information loss

- RNA-seq reads are very short compared with full-length isoforms.
- Most RNA-seq reads do not uniquely map to a single isoform.



- Technical biases introduced into RNA-seq experiments.

Existing isoform discovery methods

State-of-the-art methods for isoform discovery:

- SLIER [Jiang et al., *Bioinformatics*, 2009]
- Cufflinks [Trapnell et al., *Nature Biotechnology*, 2010]
- SLIDE [Li et al., *Proc. Natl. Acad. Sci.* 2011]
- StringTie [Pertea et al., *Nature Biotechnology*, 2015]
- ...

Limitations:

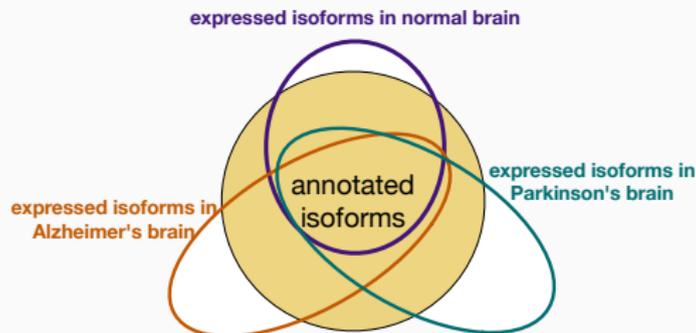
1. Low accuracy for genes with complex splicing structures.
2. Difficult to improve isoform-level performance.
[Kanitz et al., *Genome Biology*, 2015]
3. Usage of annotations results in false positives.

Usage of annotations results in false positives

Annotated isoforms are experimentally validated:



- *Ensembl* database: 203,903 isoforms
[Zerbino et al., *Nucleic Acids Research*, 2017]



False positives → false discoveries

Essay

Why Most Published Research Findings Are False

John P.A. Ioannidis

68,836

Save

3,644

Citation

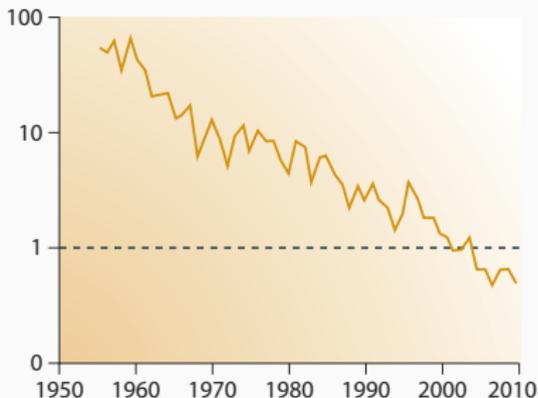
2,622,757

View

10,479

Share

Number of drugs per billion US\$ R&D spending



[Scannell et al., *Nat. Rev. Drug Discov.* 2012]

Highlights of the AIDE method

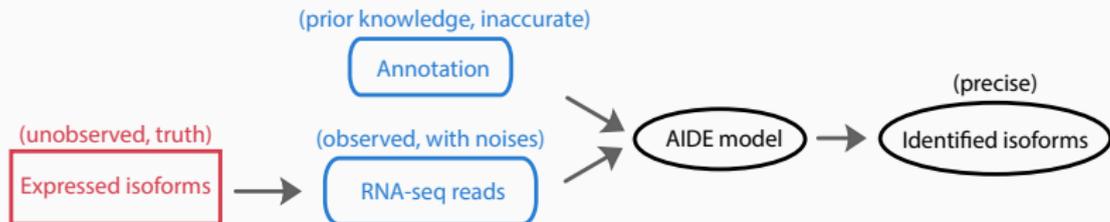
1. **Selectively leverage annotation information** to increase the precision and robustness of isoform discovery.

Highlights of the AIDE method

1. **Selectively leverage annotation information** to increase the precision and robustness of isoform discovery.
2. Practical probabilistic model to **account for technical biases**.
3. **Conservatively** identify isoforms that make statistically significant contributions to explaining the observed RNA-seq reads.

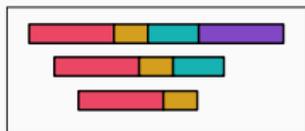
Highlights of the AIDE method

1. **Selectively leverage annotation information** to increase the precision and robustness of isoform discovery.
2. Practical probabilistic model to **account for technical biases**.
3. **Conservatively** identify isoforms that make statistically significant contributions to explaining the observed RNA-seq reads.
4. **First method to control false discoveries** by employing a statistical testing procedure.

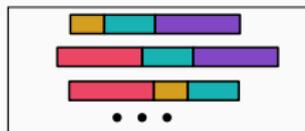


The stepwise selection in AIDE: two stages

annotated isoforms:



non-annotated isoforms:

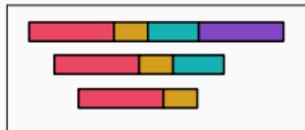


Stage 1: candidates are annotated isoforms only

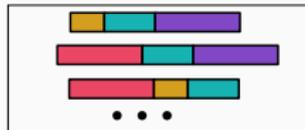
Initialization \longrightarrow Forward step \rightleftharpoons Backward step

The stepwise selection in AIDE: two stages

annotated isoforms:



non-annotated isoforms:



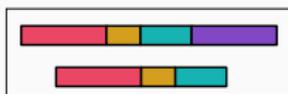
Stage 1: candidates are annotated isoforms only

Initialization



Forward step

Backward step

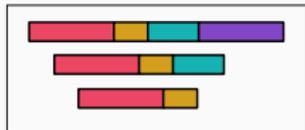


vs.

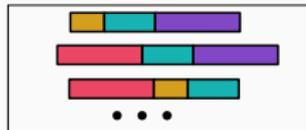


The stepwise selection in AIDE: two stages

annotated isoforms:



non-annotated isoforms:



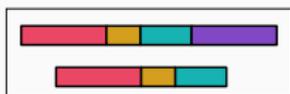
Stage 1: candidates are annotated isoforms only

Initialization



Forward step

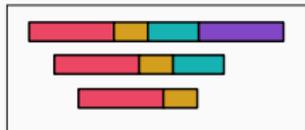
Backward step



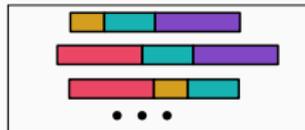
selected based on MLE

The stepwise selection in AIDE: two stages

annotated isoforms:



non-annotated isoforms:



Stage 1: candidates are annotated isoforms only

Initialization

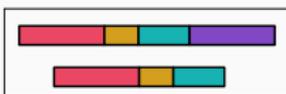


Forward step

Backward step

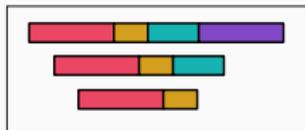


vs.
LRT

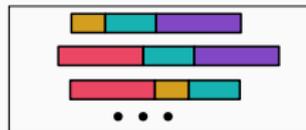


The stepwise selection in AIDE: two stages

annotated isoforms:



non-annotated isoforms:

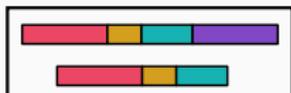


Stage 1: candidates are annotated isoforms only

Initialization



Forward step



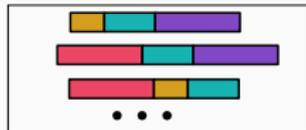
Backward step

The stepwise selection in AIDE: two stages

annotated isoforms:



non-annotated isoforms:



Stage 1: candidates are annotated isoforms only

Initialization



Forward step



Backward step

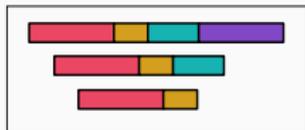


vs.

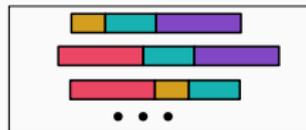


The stepwise selection in AIDE: two stages

annotated isoforms:



non-annotated isoforms:



Stage 1: candidates are annotated isoforms only

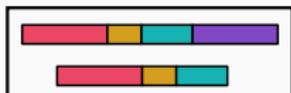
Initialization



Forward step



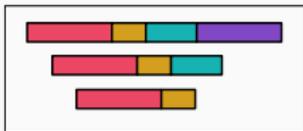
Backward step



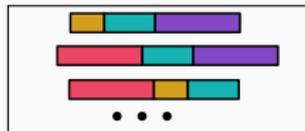
selected based on MLE

The stepwise selection in AIDE: two stages

annotated isoforms:



non-annotated isoforms:



Stage 1: candidates are annotated isoforms only

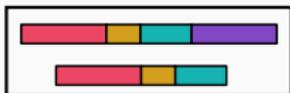
Initialization



Forward step



Backward step

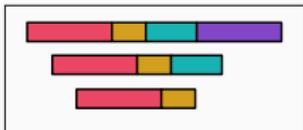


vs.
LRT

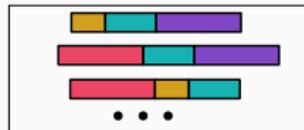


The stepwise selection in AIDE: two stages

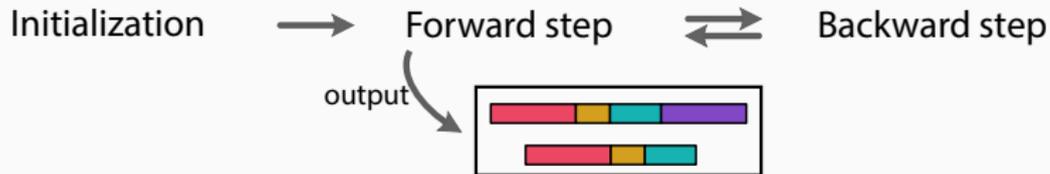
annotated isoforms:



non-annotated isoforms:

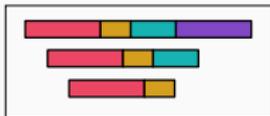


Stage 1: candidates are annotated isoforms only

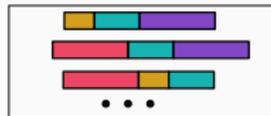


The stepwise selection in AIDE: two stages

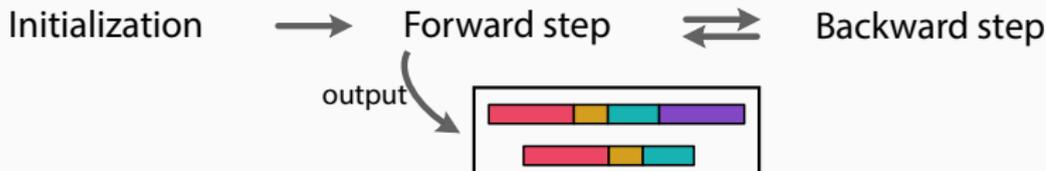
annotated isoforms:



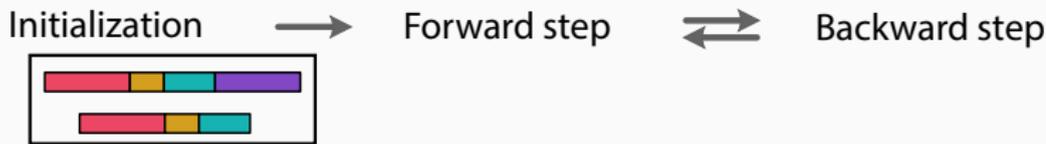
non-annotated isoforms:



Stage 1: candidates are annotated isoforms only

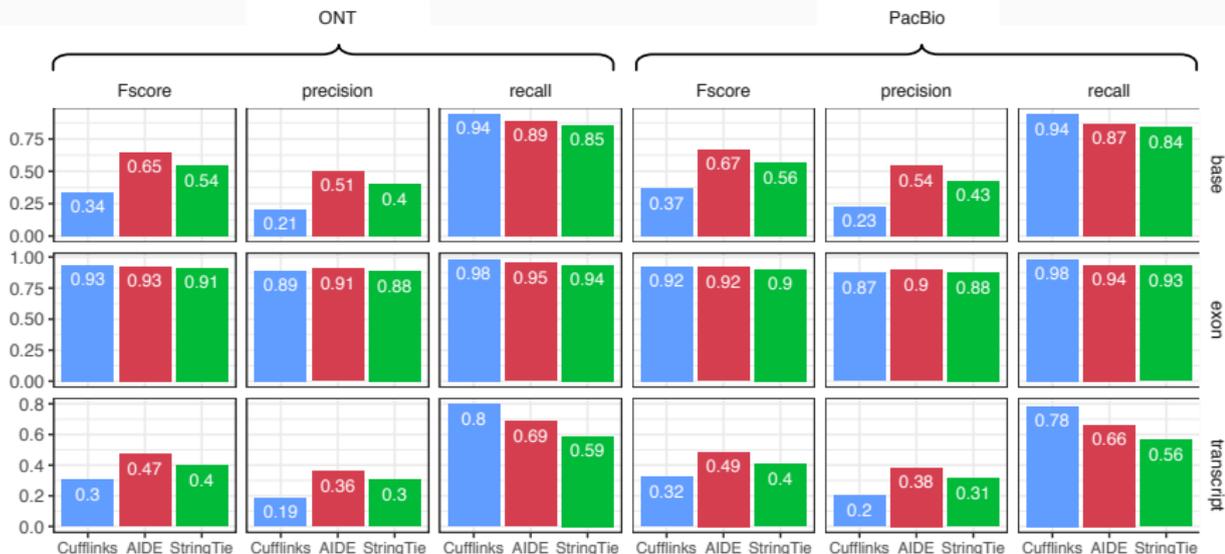


Stage 2: candidates are all possible isoforms



AIDE outperforms state-of-the-art methods

- Human embryonic stem cells
- Input: Illumina RNA-seq data
- Evaluation: PacBio and Nanopore ONT RNA-seq data



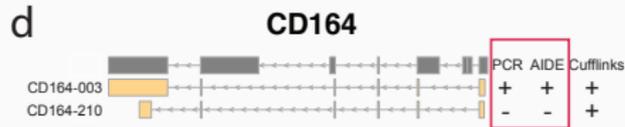
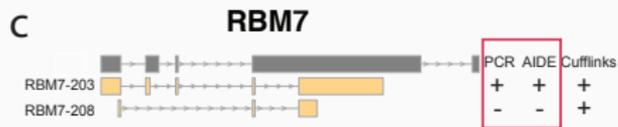
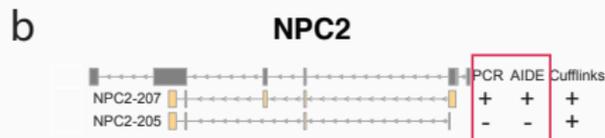
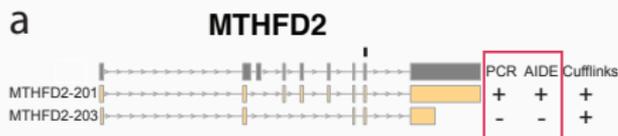
AIDE effectively reduces false discoveries in real data

- Data: breast cancer RNA-seq samples
- Six genes:
 - isoforms identified only by Cufflinks but not by AIDE
 - experimental validation (PCR)

AIDE effectively reduces false discoveries in real data

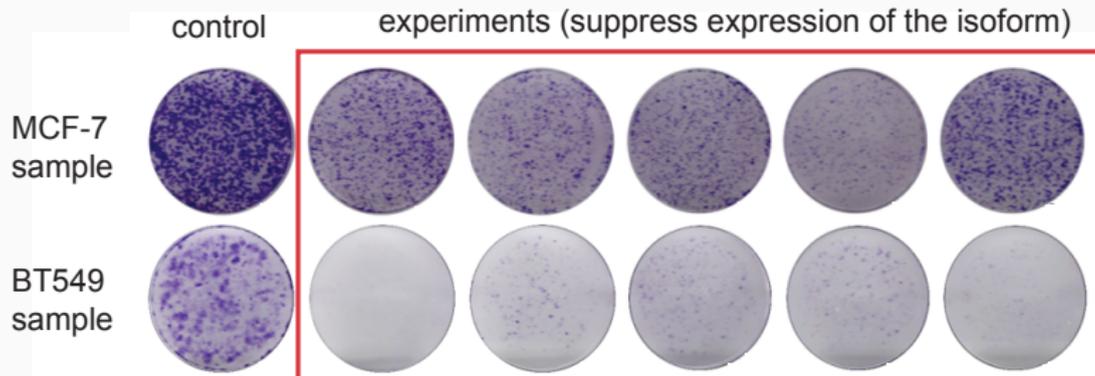
- Data: breast cancer RNA-seq samples
- Six genes:
 - isoforms identified only by Cufflinks but not by AIDE
 - experimental validation (PCR)
- Four genes:

the isoforms uniquely predicted by Cufflinks were false positives



AIDE discovers isoforms with biological significance

FGFR1



Summary of the AIDE method

- The first isoform discovery method that **directly controls false discoveries** by implementing the statistical model selection principle.
- Software: <https://github.com/Vivianstats/AIDE>
- Manuscript:



bioRxiv

THE PREPRINT SERVER FOR BIOLOGY

AIDE: annotation-assisted isoform discovery and abundance estimation from RNA-seq data

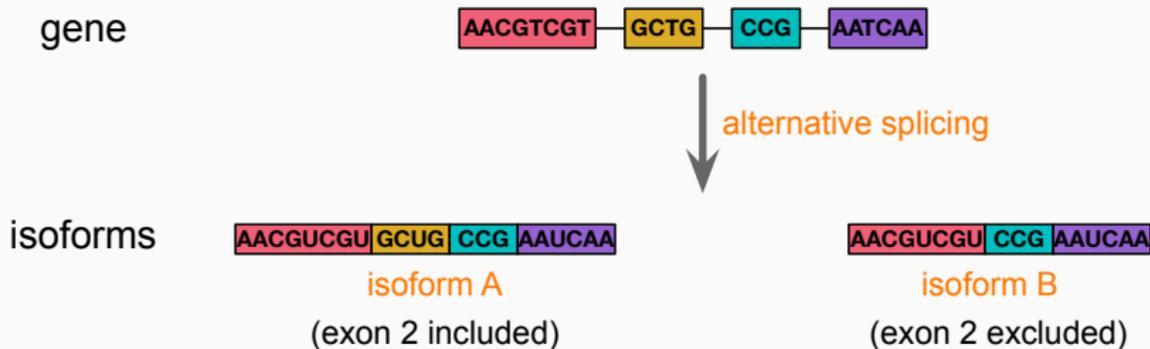
Wei Vivian Li, Shan Li,  Xin Tong, Ling Deng,  Hubing Shi, Jingyi Jessica Li

doi: <https://doi.org/10.1101/437350>

In press at *Genome Research*.

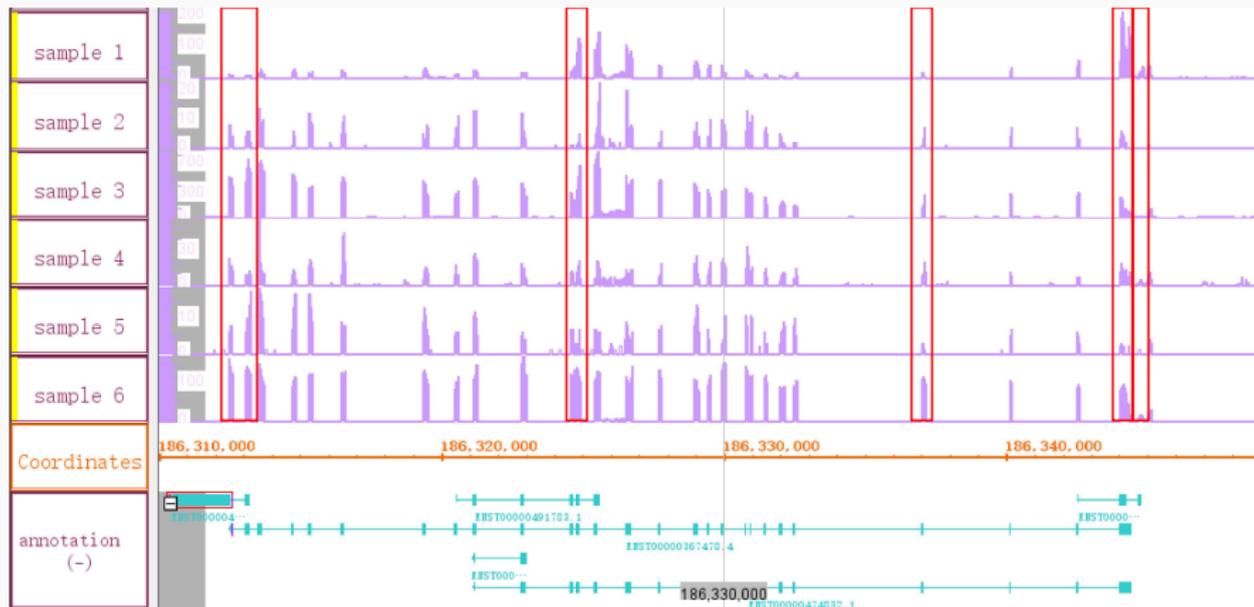
Isoform quantification: what are the isoform expression levels?

- More than 90% genes undergo alternative splicing in mammals [Hooper, *Human Genomics*, 2014].
- At least 35% genetic diseases involve abnormal splicing [Manning et al., *Nature Reviews Mol. Cell Biol.* 2017].



Motivation: multiple human ESC RNA-seq samples

chr1; gene: *TPR*



How to combine multiple RNA-seq samples?

Given D RNA-Seq (technical or biological) replicate samples and gene annotations, how to estimate the abundance of each annotated isoform for every gene?

How to combine multiple RNA-seq samples?

Given D RNA-Seq (technical or biological) replicate samples and gene annotations, how to estimate the abundance of each annotated isoform for every gene?

- Apply a single-sample method to **each sample separately** and then average the estimated isoform abundance across multiple samples?

How to combine multiple RNA-seq samples?

Given D RNA-Seq (technical or biological) replicate samples and gene annotations, how to estimate the abundance of each annotated isoform for every gene?

- Apply a single-sample method to **each sample separately** and then average the estimated isoform abundance across multiple samples?
 - This does not fully use the multi-sample information to reduce the variance in estimating isoform abundance

How to combine multiple RNA-seq samples?

Given D RNA-Seq (technical or biological) replicate samples and gene annotations, how to estimate the abundance of each annotated isoform for every gene?

- Apply a single-sample method to **each sample separately** and then average the estimated isoform abundance across multiple samples?
 - This does not fully use the multi-sample information to reduce the variance in estimating isoform abundance
- Apply a single-sample method to **a pooled sample from the D samples?**

How to combine multiple RNA-seq samples?

Given D RNA-Seq (technical or biological) replicate samples and gene annotations, how to estimate the abundance of each annotated isoform for every gene?

- Apply a single-sample method to **each sample separately** and then average the estimated isoform abundance across multiple samples?
 - This does not fully use the multi-sample information to reduce the variance in estimating isoform abundance
- Apply a single-sample method to **a pooled sample from the D samples?**
 - The estimated isoform abundance may be biased by outlier samples

Joint Modeling of **M**ultiple RNA-seq **S**amples for Accurate **I**soform **Q**uantification

Summary

- It is necessary to consider the heterogeneity of different samples to make robust isoform quantification

Summary

- It is necessary to consider the heterogeneity of different samples to make robust isoform quantification
- MSIQ is able to identify a consistent group of samples that are most representative of the biological condition

Summary

- It is necessary to consider the heterogeneity of different samples to make robust isoform quantification
- MSIQ is able to identify a consistent group of samples that are most representative of the biological condition
- MSIQ increases the accuracy of isoform quantification by incorporating the information from multiple samples

Summary

- It is necessary to consider the heterogeneity of different samples to make robust isoform quantification
- MSIQ is able to identify a consistent group of samples that are most representative of the biological condition
- MSIQ increases the accuracy of isoform quantification by incorporating the information from multiple samples
- Our proposed hierarchical model is an umbrella framework that are generalizable to incorporate more delicate consideration of read generating mechanisms

MSIQ: joint modeling of multiple RNA-seq samples for accurate isoform quantification

by Wei Vivian Li, Anqi Zhao, Shihua Zhang, and Jingyi Jessica Li

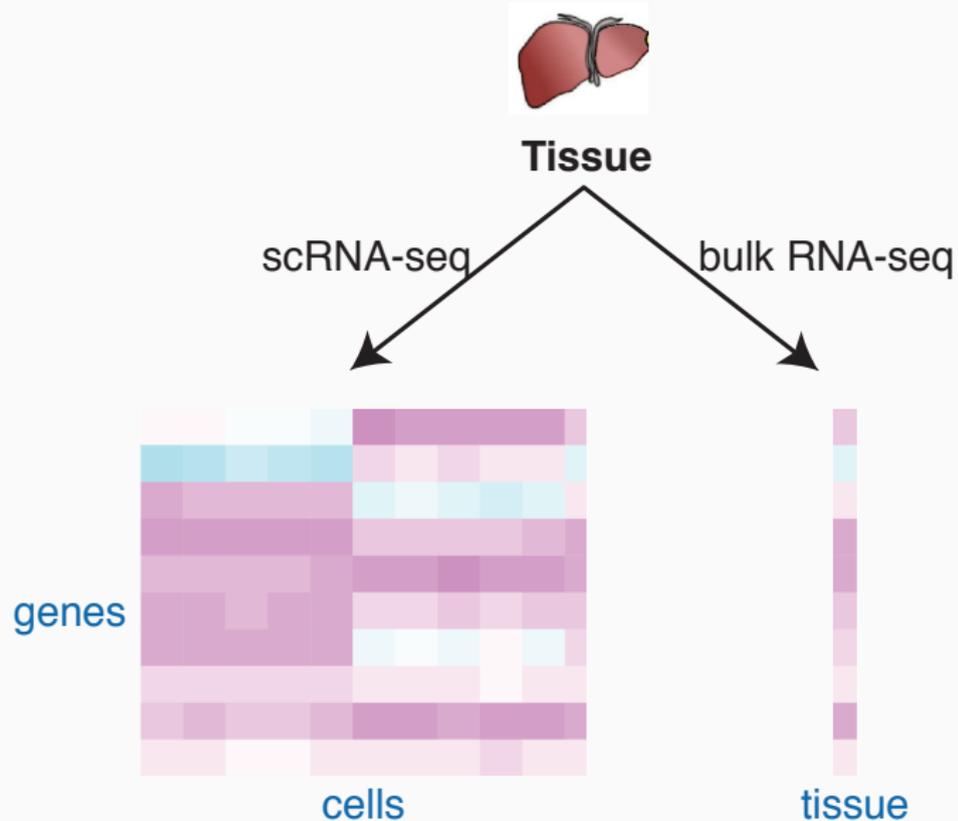
Annals of Applied Statistics 12(1):510–539

R package MSIQ

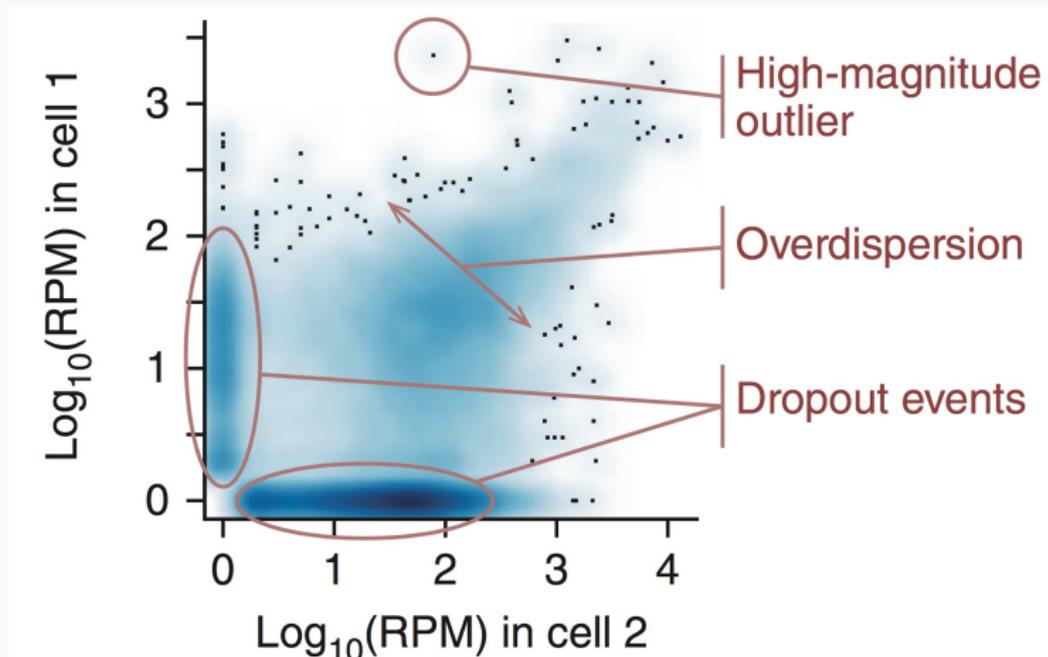
<http://github.com/Vivianstats/MSIQ>

Single-cell RNA-seq: dropout imputation & experimental design

scRNA-seq vs. bulk RNA-seq at the gene level



Dropout events in scRNA-seq



Dropout events in scRNA-seq

- A **dropout** event occurs when a transcript is expressed in a cell but is entirely undetected in its mRNA profile
- Dropout events occur due to low amounts of mRNA in individual cells
- The frequency of dropout events depends on scRNA-seq protocols
 - Fluidigm C1 platform: ~ 100 cells, ~ 1 million reads per cell
 - Droplet microfluidics: $\sim 10,000$ cells, $\sim 100K$ reads per cell
[Zilionis et al., *Nature Protocols*, 2017]
- **Trade-off**: given the same budget, more cells, more dropouts

Why do we need genome-wide explicit imputation methods?

Downstream analyses relying on the accuracy of gene expression measurements:

- differential gene expression analysis
- identification of cell-type-specific genes
- reconstruction of differentiation trajectory

It is important to adjust/correct the false zero expression values due to dropouts

Genome-wide explicit imputation for dropouts

Key points to consider:

- It is not ideal to impute all gene expressions
 - imputing expressions unaffected by dropout would introduce new bias
 - could also eliminate meaningful biological variation
- It is inappropriate to treat all zero expressions as missing values
 - some zero expressions may reflect true biological non-expression
 - zero expressions can be resulted from gene expression stochasticity

Genome-wide explicit imputation for dropouts

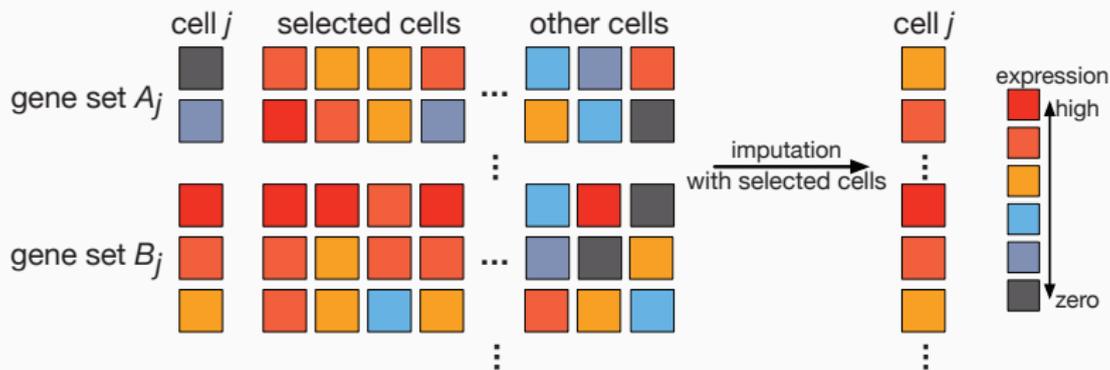
Key points to consider:

- It is not ideal to impute all gene expressions
 - imputing expressions unaffected by dropout would introduce new bias
 - could also eliminate meaningful biological variation
- It is inappropriate to treat all zero expressions as missing values
 - some zero expressions may reflect true biological non-expression
 - zero expressions can be resulted from gene expression stochasticity

How to determine which values are affected by the dropout events?

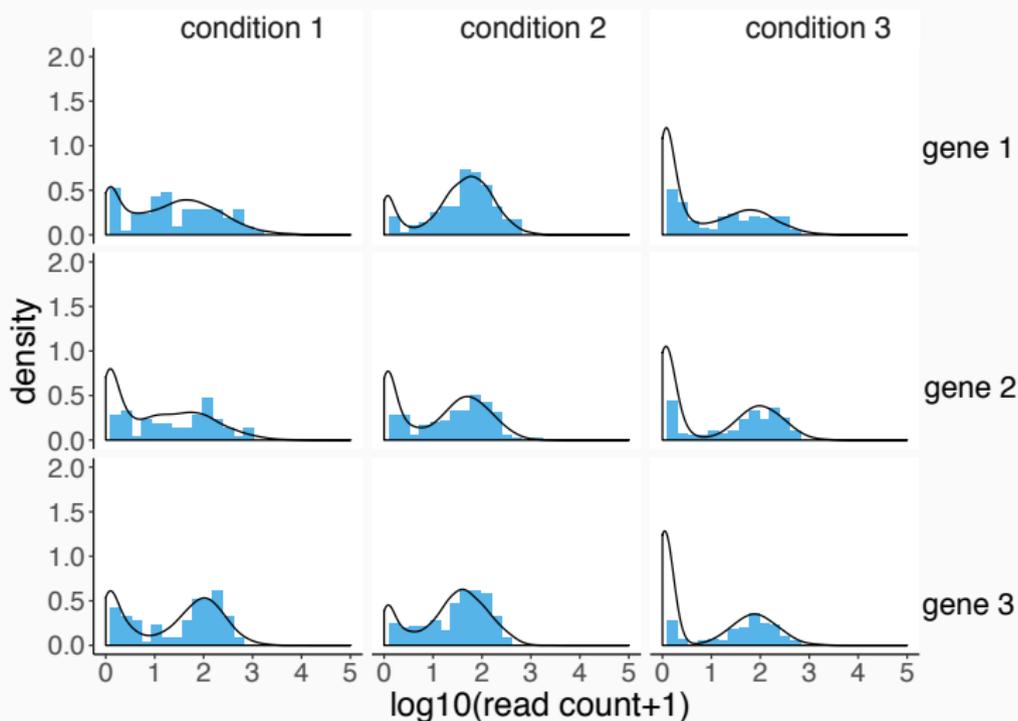
Our method: scImpute

1. For each gene, to determine which expression values are most likely affected by dropout events
2. For each cell, to impute the highly likely dropout values by borrowing information from the same genes' expression in similar cells

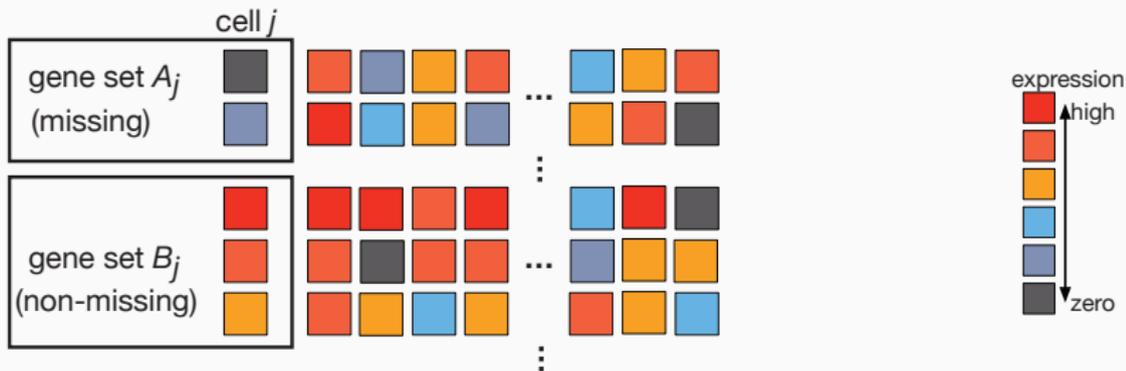


scImpute steps

1. Detection of cell subpopulations and outlier cells
2. Identification of dropout values



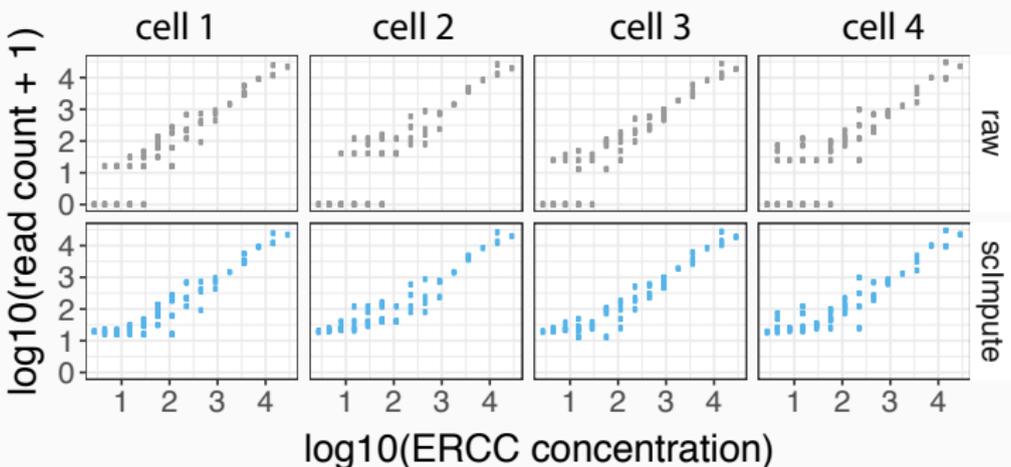
3. Imputation of gene expression cell by cell



Example 1: ERCC spike-ins

scImpute recovers the true expression of the ERCC spike-in transcripts, especially low abundance transcripts that are impacted by dropout events

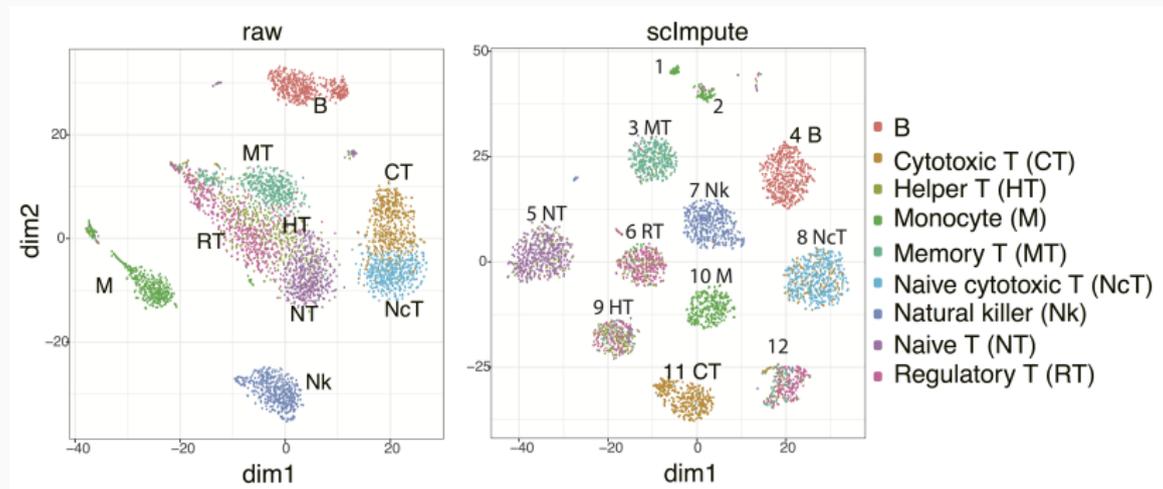
- 3,005 cells from the mouse somatosensory cortex region
- 57 ERCC transcripts



Example 2: cell clustering

4,500 peripheral blood mononuclear cells (PBMCs) from high-throughput droplet-based system 10x genomics [Zheng et al., *Nature communications*, 2017]

Proportion of zero expression is 92.6%



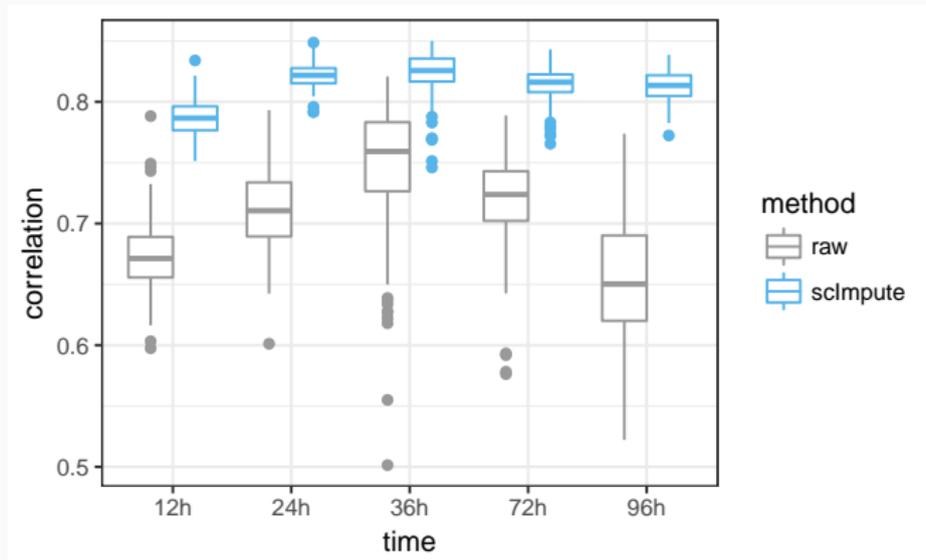
Example 3: gene expression dynamics

Bulk and single-cell time-course RNA-seq data profiled at 0, 12, 24, 36, 72, and 96 h of the differentiation of embryonic stem cells into definitive endoderm cells [Chu et al., *Genome biology*, 2016]

time point	00h	12h	24h	36h	72h	96h	total
scRNA-seq (cells)	92	102	66	172	138	188	758
bulk RNA-seq (replicates)	0	3	3	3	3	3	15

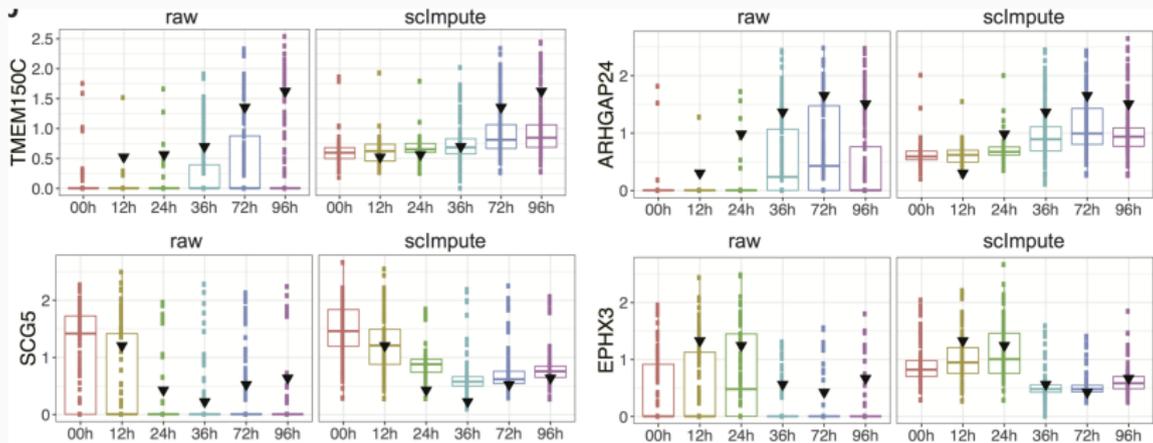
Example 3: gene expression dynamics

Correlation between gene expression in single-cell and bulk data



Example 3: gene expression dynamics

Imputed read counts reflect more accurate gene expression dynamics along the time course



Conclusions

- `sclImpute` is a flexible and easily interpretable statistical method that addresses the dropout events prevalent in scRNA-seq data
- `sclImpute` focuses on imputing the missing expression values of dropout genes, while retaining the expression levels of genes that are largely unaffected by dropout events
- `sclImpute` is compatible with existing pipelines or downstream analysis of scRNA-seq data, such as normalization, differential expression analysis, clustering and classification
- `sclImpute` scales up well when the number of cells increases

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

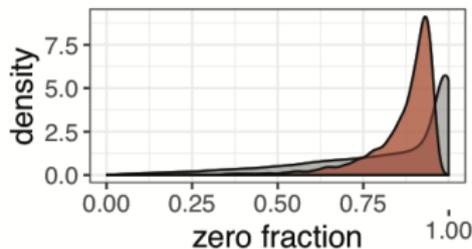
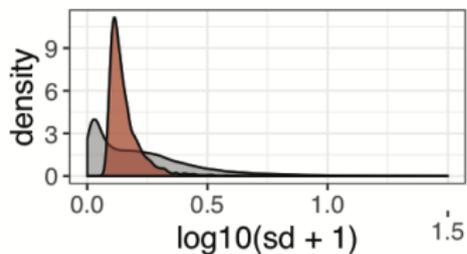
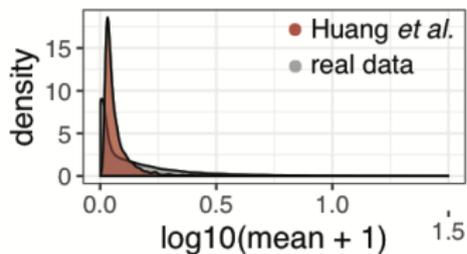
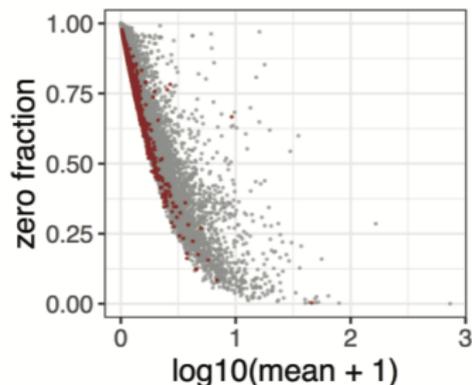
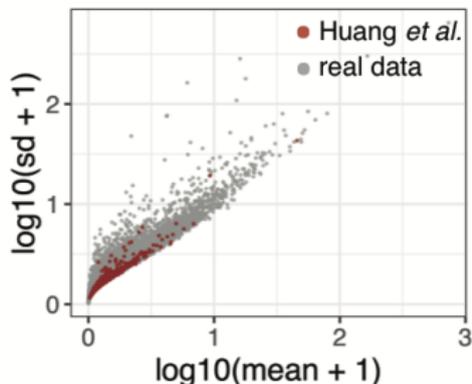
by Wei Vivian Li and Jingyi Jessica Li

Nature Communications 9:997

R package scImpute

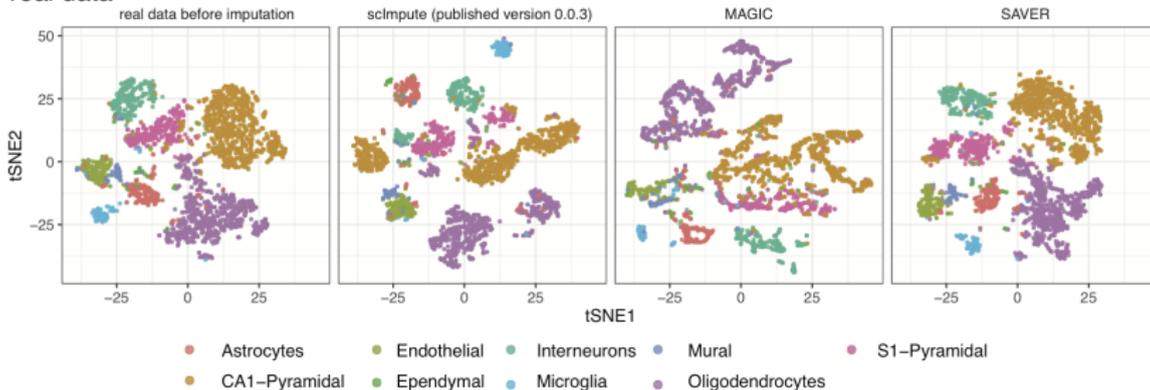
<https://github.com/Vivianstats/scImpute>

Real vs. semi-synthetic data

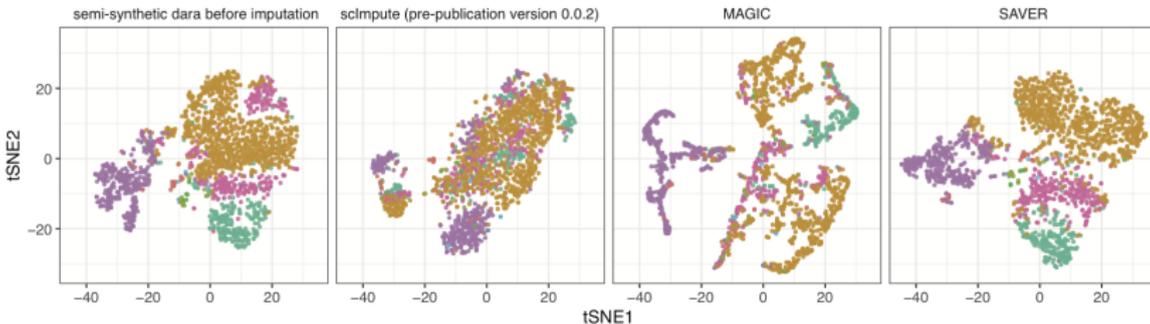


Real vs. semi-synthetic data

real data



Huang *et al.* semi-synthetic data



Benchmark standard

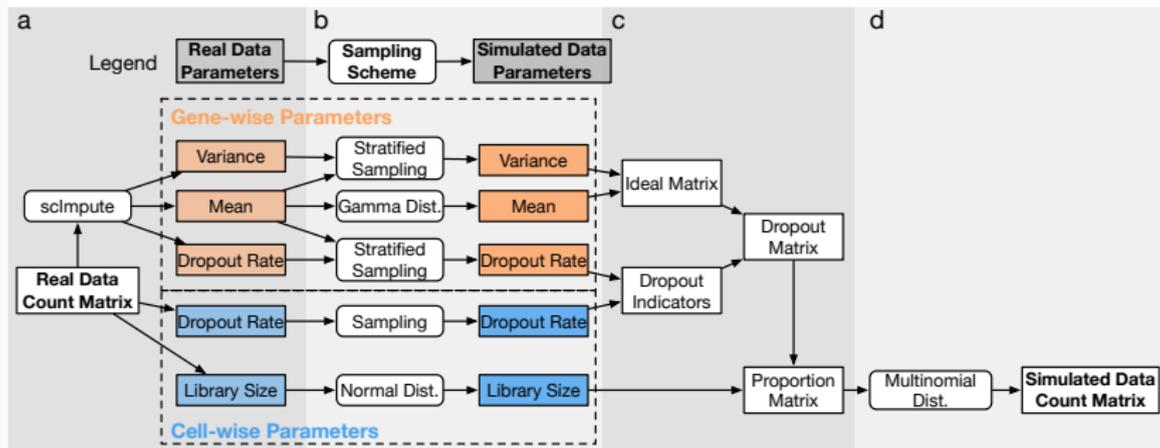
		labels used in Huang <i>et al.</i>						
		0	1	2	3	4	5	6
labels reported in Zeisel <i>et al.</i>	CA1-Pyramidal	442	20	289	1	4	42	40
	S1-Pyramidal	2	273	1	1	0	32	11
	Oligodendrocytes	0	0	0	282	0	62	2
	Interneurons	5	7	2	0	220	6	1
	Endothelial	0	0	0	0	1	0	14
	Microglia	0	0	0	0	0	0	6
	Mural	0	1	0	0	0	0	0
	Ependymal	0	0	0	0	0	0	7
	Astrocytes	0	1	0	2	0	1	20

Simulation-based scRNA-seq experimental design

Advantages of scDesign:

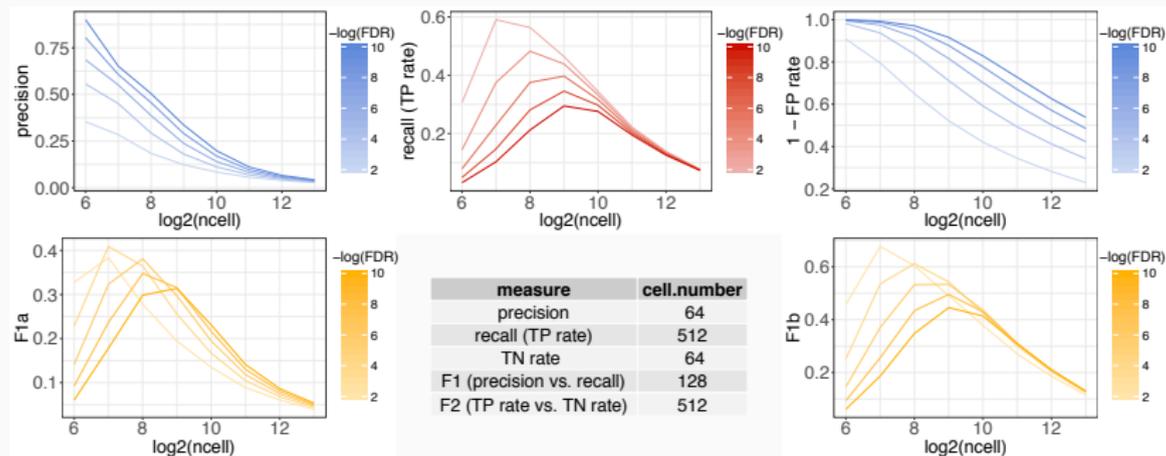
- Protocol-adaptive and data-adaptive: learn from
 - Public scRNA-seq data
 - Pilot-study data
- Generate synthetic data that well mimic real data under a pre-specified experimental setting
 - Assist experimental design & method development
- Flexible in accommodating user-specific analysis needs
- No experimental cost

Generative framework of scDesign



Experimental design for gene differential expression

Astrocytes vs. Oligodendrocytes (Fluidigm C1)



Bioinformatics, 35, 2019, i41–i50

doi: 10.1093/bioinformatics/btz321

ISMB/ECCB 2019

OXFORD

A statistical simulator scDesign for rational scRNA-seq experimental design

Wei Vivian Li¹ and Jingyi Jessica Li ^{1,2,*}

¹Department of Statistics, University of California, Los Angeles, CA 90095-1554, USA and ²Department of Human Genetics, University of California, Los Angeles, CA 90095-7088, USA

*To whom correspondence should be addressed.

Bulk RNA-seq

- Isoform identification: **AIDE**
- Isoform quantification: **MSIQ**

Single-cell RNA-seq

- Dropout imputation: **scImpute**
- Simulator & experimental design: **scDesign**

Acknowledgements

Dr. Wei Vivian Li

former PhD student @UCLA

currently Assistant Professor @Rutgers

Collaborators:

Dr. Hubing Shi (Sichuan University)

Dr. Xin Tong (USC)

Dr. Shihua Zhang (CAS)

Dr. Anqi Zhao (NUS)

