# Cell type expression phenotyping from deconvolution of bulk RNAseq data

**Christopher K Tuggle** 

**Department of Animal Science** 

**Iowa State University** 

**AG2PI Field Day** 

March 12, 2025

NIFA Project # 2022-67015-37055



United States National Institute Department of of Food and Agriculture Agriculture



## Value of single cell transcriptomics for biology and genetics

Heterogeneous cell populations in tissues- analyses of "bulk" tissue misses a lot....





## ENCODE project showed enrichment of traitassociated SNPs is *biological context specific*.

## Disease associated variants often overlap with regulatory elements specific for target cell/tissue pathology

### **Cell-selective enrichment of trait-associated variants**



We need to test using context-specific information!

## Linking to Breeding Goals: Finding cell-type specific eQTL in tissues/cell mixes

#### **RESEARCH ARTICLE SUMMARY**

#### **PSYCHENCODE2**

# Single-cell genomics and regulatory networks for 388 human brains

Emani et al., Science 384, 862 (2024) 24 May 2024

Prashant S. Emani<sup>1,2</sup><sup>+</sup>, Jason J. Liu<sup>1,2</sup><sup>+</sup>, Declan Clarke<sup>1,2</sup><sup>+</sup>, Matthew Jensen<sup>1,2</sup><sup>+</sup>, Jonathan Warrell<sup>1,2</sup><sup>+</sup>, Chirag Gupta<sup>3,4</sup><sup>+</sup>, Ran Meng<sup>1,2</sup><sup>+</sup>, Che Yu Lee<sup>5</sup><sup>+</sup>, Siwei Xu<sup>5</sup><sup>+</sup>, Cagatay Dursun<sup>1,2</sup><sup>+</sup>, Shaoke Lou<sup>1,2</sup><sup>+</sup>, Yuhang Chen<sup>1,2</sup>, Zhiyuan Chu<sup>1</sup>, Timur Galeev<sup>1,2</sup>, Ahyeon Hwang<sup>5,6</sup>, Yunyang Li<sup>2,7</sup>, Pengyu Ni<sup>1,2</sup>, Xiao Zhou<sup>1,2</sup>, PsychENCODE Consortium<sup>1</sup>, Trygve E. Bakken<sup>8</sup>, Jaroslav Bendl<sup>9,10,11,12</sup>, Lucy Bicks<sup>13</sup>, Tanima Chatteriee<sup>1,2</sup>. Lijun Cheng<sup>14</sup>. Yuyan Cheng<sup>13,15</sup>, Yi Dai<sup>5</sup>, Ziheng Duan<sup>5</sup>, Mary Flaherty<sup>14</sup>, John F. Fullard<sup>9,10,11,12</sup>, Michael Gancz<sup>1,2</sup>, Diego Garrido-Martín<sup>16</sup>, Sophia Gaynor-Gillett<sup>14,17</sup>, Jennifer Grundman<sup>13</sup>, Natalie Hawken<sup>13</sup>, Ella Henry<sup>1,2</sup>, Gabriel E. Hoffman<sup>9,10,11,12,18,19</sup>. Ao Huang<sup>1</sup>. Yunzhe Jiang<sup>1,2</sup>, Ting Jin<sup>3,4</sup>, Nikolas L. Jorstad<sup>8</sup>, Riki Kawaguchi<sup>13,20</sup>, Saniya Khullar<sup>3,4</sup>, Jianyin Liu<sup>13</sup>, Junhao Liu<sup>5</sup>, Shuang Liu<sup>4</sup>, Shaojie Ma<sup>21,22</sup>, Michael Margolis<sup>13</sup>, Samantha Mazariegos<sup>13</sup>, Jill Moore<sup>23</sup>, Jennifer R. Moran<sup>14</sup>, Eric Nguyen<sup>1,2</sup>, Nishigandha Phalke<sup>23</sup>, Milos Pjanic<sup>9,10,11,12</sup>, Henry Pratt<sup>23</sup>, Diana Quintero<sup>13</sup>, Ananya S. Rajagopalan<sup>7</sup>, Tiernon R. Riesenmy<sup>24</sup>, Nicole Shedd<sup>23</sup>, Manman Shi<sup>14</sup>, Megan Spector<sup>14</sup>, Rosemarie Terwilliger<sup>25</sup>, Kyle J. Travaglini<sup>8</sup>, Brie Wamsley<sup>13</sup>, Gaoyuan Wang<sup>1,2</sup>, Yan Xia<sup>1,2</sup>, Shaohua Xiao<sup>13</sup>, Andrew C. Yang<sup>1,2</sup>, Suchen Zheng<sup>1,2</sup>, Michael J. Gandal<sup>26,27,28,29,30</sup>, Donghoon Lee<sup>9,10,11,12</sup>, Ed S. Lein<sup>8,31,32</sup>, Panos Roussos<sup>9,10,11,12,18,19</sup>, Nenad Sestan<sup>21</sup>, Zhiping Weng<sup>23</sup>. Kevin P. White<sup>33</sup>, Hyejung Won<sup>34</sup>, Matthew J. Girgenti<sup>25,35,36</sup>\*, Jing Zhang<sup>5</sup>\*, Daifeng Wang<sup>3,4,37</sup>\*, Daniel Geschwind<sup>13,20,27,28,38</sup>\*. Mark Gerstein<sup>1,2,7,24,39</sup>\*



#### ACKNOWLEDGMENTS

The authors thank the founder of the Allen Institute, P. G. Allen, for his vision, encouragement, and support. R.T. and M. J. Girgenti thank Keck Microarray Shared Resource (KMSR) and Yale Center for Genome Analysis (YCGA) at Yale University for their assistance with 10x Genomics single-cell RNA-seq services. Funding: Data were generated as part of the PsychENCODE Consortium, supported by U01DA048279, U01MH103339, U01MH103340, U01MH103346, U01MH103365, U01MH103392, U01MH116438, U01MH116441, U01MH116442, U01MH116488, U01MH116489, U01MH116492, U01MH122590, U01MH122591, U01MH122592, U01MH122849, U01MH122678, U01MH122681, U01MH116487, U01MH122509, R01MH094714, R01MH105472, R01MH105898, R01MH109677, R01MH109715, R01MH110905, R01MH110920, R01MH110921, R01MH110926, R01MH110927, R01MH110928, R01MH111721, R01MH117291, R01MH117292, R01MH117293, R21MH102791, R21MH103877, R21MH105853, R21MH105881, R21MH109956, R56MH114899, R56MH114901, R56MH114911, R01MH125516, R01MH126459, R01MH129301, R01MH126393, R01MH121521, R01MH116529, R01MH129817, R01MH117406, and P50MH106934 awarded to authors and collaborators A. Abvzov.

-NCTP mode

000

0

ase risk prediction

ated perturbations

ind gene regulatory

s and cell types.

 $\cap$ 

# Comprehensive mining of the blood cell transcriptome for improved phenomics in swine

Trait measured during each phase a

Quarantine Nursery

Rationale: Blood sampling is very practical and useful phenotype- but needs improvement to use in predictive genetics- a mixture of genetic effects on phenotypes

## Whole blood RNA patterns have been associate disease phenotypes in pigs in Resilience project

aNurHS1 Multiplier Herd Experimental Facilities > 2,400 samples ~24 d of age ~43 d ~71 d qNurHS2 with RNAseq data! Birth Weaning Entry Entry Growth rate Quarantine Challenge Farrowing **Challenge Finisher** NURSERY NURSERY Biosecure Natural disease challenge Biosecure **Blood Transcriptome Resilience phenotypes** Complete blood count Resilience Project natural disease challenge model (NDCM) Quantitative analysis of the blood transcriptome of young healthy pigs and its relationship with subsequent disease resilience Kyu-Sang Lim, Jian Cheng, Austin Putz, Qian Dong, Xuechun Bai, Hamid Beiki, Christopher K. Tuggle, Michael K. Dyck, Pig Gen Canada, Frederic Fortin, John C. S. Harding, Graham S. Plastow & Jack C. M. Dekkers BMC Genomics 22, Article number: 614 (2021) Cite this article

Table 3 The number of genes with expression levels in blood of young healthy pigs that were significantly (q < 0.20) associated with observed phenotypes, with or without accounting for blood cell composition, the

Without

395 (1656) °

171 (982)

744 (3224)

Blood

Collection

Number of genes from expression residuals with or without adjustment for cell composition

With+Without b

CBC

Differential Count

Eosinophil

Neutrophil

Band cell

Red blood cell

Lymphocyte

Monocyte

395

173

856

With

29 (1106)

101 (0)

2-3 mins

Plasma

WBCs

RBCs

Platelets

830 (3357)

Blood Smear

Centrifuged

mins

### Improving blood phenotypes: towards cell type specific phenotypes

Hypothesis: tools for accurate deconvolution of porcine whole blood transcriptome data into <u>cell-type-specific transcriptome data</u> will substantially improve molecular blood phenotypes as direct selection tools or as markers for animal traits



### **Cell Type Deconvolution of Cell Mixtures**

Nucleic Acids Research, 2024, **52**, 4761–4783 https://doi.org/10.1093/nar/gkae267 Advance access publication date: 15 April 2024 **Critical Reviews and Perspectives** 



# Fourteen years of cellular deconvolution: methodology, applications, technical evaluation and outstanding challenges

Hung Nguyen <sup>1,†</sup>, Ha Nguyen <sup>1,†</sup>, Duc Tran<sup>2</sup>, Sorin Draghici <sup>3,4,\*</sup> and Tin Nguyen <sup>1,\*</sup>

#### **Graphical abstract**



### **Cell Type Deconvolution**

- Computational methods developed to infer cell type compositions and GEP within heterogenous samples
- Disease mechanism, immune response, transcriptional states, cell proportions in disease phenotypes





## **Cell Type Deconvolution of Cell Mixtures**

**Assumption:** Total GE in bulk is *linear combination* of gene expression in individual cell types, weighted by cell type proportions.

Model:  $M = S \times F$ 

M=N genes x m samples S= N genes x C cell types F= C cell types x m samples

For every gene in bulk:

 $M_{N1} = F_1 \times S_{N1, C1} + F_2 \times S_{N1, C2} \dots$  $M_{N2} = F_1 \times S_{N2, C1} + F_2 \times S_{N2, C2} \dots$ 

### Output: cell type composition, CTS-GEP

Cell type proportion

Expression profiles







Gene expression signature matrix:

- Rows are genes
- Columns are cell types

**Starting Data:** 

Using SC RNAseq for relevant samples, create "pseudobulk" samples to be deconvoluted



**Starting Data:** 

Collect data from separate cell populations and create artificial mixtures, run bulk RNAseq of individuals and mixtures with known proportions





**Starting Data:** 

| Use    | bulk     | RN  | Aseq | and   |  |
|--------|----------|-----|------|-------|--|
| disso  | ociated  | C   | from |       |  |
| replic | cate     | san | in   |       |  |
| scRN   | Aseq     | to  | dete | rmine |  |
| cell p | oroporti | ons |      |       |  |

**Starting Data:** 

Run flow cytometry for isolated cells from blood, or tissue for cell composition, run bulk RNAseq of original sample



## **Comprehensive testing of methods using tissue atlases: scoring**

### **Starting Data:**



# Many, many methods for Deconvolution...

Evaluated using the multi-tissue reference from Tabula Sapiens and CellxGene

| A                 |           |     | the of | noa   |       |        | -     |      | NAME OF A DATABASE |   |      | Te of the local |         |        |       |
|-------------------|-----------|-----|--------|-------|-------|--------|-------|------|--------------------|---|------|-----------------|---------|--------|-------|
|                   |           |     | 1      |       | 10000 | - Inpu | t     | 3.1  | Output*            |   |      |                 |         |        |       |
|                   |           |     | 10     | æ     |       | 8      | à     | 2/1  | Ø /                | á   | 23   | 1               | R       |        |       |
|                   | .of       | · / | de .   | Pres. | d.    | de la  | and a | 000  | and de la          | NY A  | ABCY | applies         | a state | antry. | Alley |
| Potoronoo bacad   | Plan      | 100 | 8      | 6     | a Mar | E ST   | 100   | - GR | A. C.              | Oren.   | ASCO | SCOL            | CORD    | State  | USB   |
| Music             | 100       |     | 1      |       |       |        |       |      | WLCLE              | _   | _    |                 | _       | _      |       |
| NILG DIMI S       | G         |     | 4      | •     |       |        |       | 1    | WECLS              | -   |      |                 |         | _      |       |
| InDeconSec        | (p)       |     | 1      | 0     |       |        |       | 3    | W-CLS              |   |      |                 |         |        |       |
| AdRoit            | Get       |     | 1      | 0     |       |        |       | 1    | B.CLS              | -   |      | -               |         |        |       |
| RNA-Sieve         |           |     | 1      |       |       |        |       | 1    | MIE                | -   |      |                 |         |        |       |
| Scaden            | 0         |     | 1      |       |       |        |       |      | DNN                | 1   |      |                 |         |        |       |
| enatialDWI S      |           |     | 1      | 0     | 1     |        |       | 1    | WCIS               | -   |      |                 |         | -      |       |
| AutoGeneS         | -         |     | 1      | S     |       |        |       | 3    | CLS/SVB            |   |      |                 |         | -      |       |
| DecOT             | -         |     | 1      | s     |       |        |       | 1    | Ensemble           |   |      |                 |         | -      |       |
| BouesPrism        | G         |     | 1      | 0     |       |        |       | 1    | Bayesian           |   |      |                 |         | -      |       |
| DigitalDI Sorter  | -         |     | 1      |       |       |        |       |      | DNN                | -   |      |                 |         | -      |       |
| BaylCE            | G         |     | 1      |       |       |        |       | 1    | Reuppion           | - 1   |      |                 |         | _      |       |
| Dayloc            | 1         |     | 1      | 0     |       |        |       | 1    | CLS                |   | =    |                 | _       | -      |       |
| CPM               | œ         |     | 1      | 9     |       |        |       |      | SVR                |   |      |                 |         |        |       |
| BisqueBet***      | R         |     | 1      |       |       |        |       |      | CLS                |   | H    |                 |         |        |       |
| SCDC              | G         |     | 1      | S     |       |        |       | 5    | Ensemble           | 11  |      |                 |         |        |       |
| DAISM-DNN         | 2         |     | 1      | 9     |       |        | 1     | 1    | DNN                |   |      |                 |         |        |       |
| MOME              | G         |     | 1      |       |       |        |       | 1    | NME                |   | H    | -               |         |        |       |
| DeMixT            | R         |     | 1      |       |       |        |       | 1    | MLE                |   |      |                 | -       | -      |       |
| deconvSeq         | R         |     | 1      | s     |       |        |       | 1    | MLE                |   | E .  |                 |         |        |       |
|                   |           |     |        |       |       |        |       |      |                    |   |      |                 |         |        |       |
| CIBERSORT         | RO        |     |        | S     |       | 1      |       |      | v-SVR              | 1   |      |                 |         |        |       |
| MethylResolver    | R         |     |        | S     |       | 4      |       |      | CLS                |   |      | 12 m            | 19      |        |       |
| MIXTURE           | Re        |     |        | S     |       |        |       |      | v-SVR              | C.  |      |                 |         | 2      |       |
| FARDEEP           | R         |     |        | S     |       |        |       |      | CLS                | 2   |      | 14 × 5          | 1. NO   |        |       |
| MySort            |           |     |        | S     |       | *      |       |      | v-SVR              |   |      | 11              |         |        |       |
| NITUMID           | R         |     |        | S     |       | 1      |       | 1    | NMF                | 6   |      | 15 K            | 5       |        |       |
| quanTiseq         | 2         |     |        | S     |       | 1      |       |      | CLS                | -   |      | 20 C            |         |        |       |
| DeconRNASeq       | R         |     |        | S     |       |        |       |      | CLS                |   |      |                 | 1       |        |       |
| Bseq-SC           | R         |     |        | S     |       | 1      |       | 1    | v-SVR              |   |      |                 | 1000    |        |       |
| DCQ               |           |     |        | S     |       | 1      |       |      | R-CLS              |   | 0    |                 | 1       | 1      |       |
| DESeq2's unmix    | R         |     |        | F.    |       |        |       |      | CLS                | 6   |      |                 | 2       |        |       |
| dtangle           | R         |     |        | F     | 1     |        |       |      | Scoring            | -   |      |                 | 1 mm    |        |       |
| ARIC              | e .       |     |        | F     |       |        |       |      | W-SVR              |   | -    |                 | 1       |        | 1     |
| PREDE             | R         |     |        | F     |       |        |       | 1    | NME                |   |      | 12 12           | 1 CT    |        |       |
| EMeth             | R         |     |        | F     |       |        |       |      | MLE                | 0   |      |                 | -       |        | 12    |
| ImmuCellAl        | RO        |     |        | F     | 1     | 1      | 1     |      | CLS                |   |      |                 | S       |        |       |
| EPIC              | RO        |     |        | F     | 1     |        |       |      | W-CLS              | 19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-<br>19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-10-19-1 |      |                 |         |        |       |
| DeCompress        | R         | 1   |        | F     |       |        |       | 1    | Ensemble           | 1   |      | -               | 1000    |        | 10 A  |
| TICPE***          | R         |     |        | F     | 1     |        |       |      | Scoring            | 1   | 0    | 1               | 1       | 1      |       |
| Reference-free    |           |     |        |       |       |        |       |      |                    |   |      |                 |         |        |       |
| Linseed           | (B        | 1   |        |       |       |        |       |      | Scoring            |   |      |                 |         |        |       |
| TOAST             | CR CR     | 1   |        |       |       |        |       | 1    | NME/DCA            |   | -    | -               |         | -      |       |
| CollDistinguisher | G         | 1   |        |       |       |        |       | 3    | NME                |   | -    |                 |         | -      |       |
| DeconICA          | G         | 1   |        |       |       |        | 1     | 1    | NME                |   |      |                 | _       |        |       |
| debCAM            | (m        | 1   |        |       |       |        |       | 1    | NME                |   | 0    |                 |         | -      |       |
| BayesCCF          | A MARLINE | 1   |        |       |       |        |       | 1    | Baupsian           |   |      |                 |         |        |       |
| deconf            | (P        | 1   |        |       |       |        |       | 1    | NME                |   |      |                 |         |        |       |
| ReFACTor          | (D        | 1   |        |       |       |        | 1     | 112  | PCA                |   |      |                 |         |        |       |
| BayCount          | GR        | 1   |        |       |       |        |       | 1    | Bayesian           |   | H    |                 |         |        |       |
| SMC               | MAILAR    | 1   |        |       |       |        |       | 1    | Bayesian           | 1   | 0    | 1               | I       |        | I     |
|                   |           | -5- |        |       |       |        |       | -    | ouycoldii          | ·   | - C. |                 | 4       | *      |       |
| Semi-reference-fr | ee        |     |        |       |       |        |       |      |                    |   |      |                 |         |        |       |
| MCP-counter       | R         |     |        |       | 1     | 1      | 1     |      | Scoring            | 6   |      | 10 A            |         |        |       |
| Deblender         | MAILAS    | 1   |        |       | 1     |        |       | 1    | NMF                | S   |      | 19 (A)          |         |        | 3     |
| BisqueMarker      | R         |     |        |       | 1     |        | 1     |      | PCA                |   |      |                 |         |        |       |
| DRA               | (P        |     |        |       | 1     |        |       |      | Scoring            | 1 million (1997)  |      | 1               | 1 N     |        |       |

### **Curated path to choose method for specific situations!**



### Application for deconvolution of whole blood RNAseq data: Create Cell-type-specific signatures



Kristen

Byrne

Kapoor

Carrie Meeks



Mehak Kapoor



time points- total of 45 samples

Collect eight-cell type

Flow cytometry based Cell Counts (TRUTH) n=45

Create new reference information on immune-stimulated pigs



# Study Design: scRNAseq of PBMC from Salmonella challenged pigs



### Cells passing QC for each sample

| Pig ID   | 0D     | 2D     | 8D     | Total   |  |  |  |  |
|--|--------|--------|--------|---------|--|--|--|--|
| 842  | 10,306 | 10,164 | 16,274 | 36,744  |  |  |  |  |
| 852  | 6,877  | 4,304  | 5,151  | 16,332  |  |  |  |  |
| 853  | 8,319  | 4,600  | 7,089  | 20,008  |  |  |  |  |
| 854  | 1,913  | 5,285  | 1,927  | 9,125   |  |  |  |  |
| 864  | 4,959  | 11,236 | 4,671  | 20,866  |  |  |  |  |
| Total  | 32,374 | 35,589 | 35,112 | 103,075 |  |  |  |  |
| Cells with >300 genes, >400 UMIs and <15% mito |        |        |        |         |  |  |  |  |

sequences were retained, duplets removed.

# Additional data available to project

- scRNAseq of PBMC samples of 7 <u>healthy</u> pigs (10X Genomics)
- Total of 28,810 cells.
- 36 clusters grouped to 13 major porcine cell types.

- RNAseq of nine sorted populations of blood cells from 2 healthy pigs
- Cell populations covering all nucleated cells in blood



Herrera-Uribe and Wiarda, et al. 2021; Herrera-Uribe et al. 2023



### Application for deconvolution of whole blood RNAseq data: Create Cell-type-specific signatures



Kristen Byrne

Kapoor



Meeks

Mehak Kapoor



### **Feature Selection Tool: AutoGeneS**

- Automatic feature selection, does not rely on pre-defined markers.
- Deals with multi-collinearity, a critical problem in deconvolution when dealing with closely related cell types
- Employs multi-objective optimization as a solution to select a set of noncollinear genes.
- Aims at minimizing correlation and maximizing Euclidean distance



# Validation methods we can use





# AutogeneS PBMC only

- Used 7 sample published scRNAseq data set as reference
- MOO to select features and create gene expression signature matrix
- Check which cell types are highly correlated (will be difficult to distinguish) and refine cell type markers if needed





cell types

# Heatmap of GE signature matrix (400 features)



# **AutogeneS PBMC deconvolution results**

- Use GE Sig matrix from 7 sample data to predict proportions of new 15 sample scRNAseq data set (Salmonella infection)
- Used pseudobulking of 15 datasets
- Two models (NNLS, NuSVR)



### Within sample across cell types



### Pearson correlation range 0.67-0.96 across 15 samples

# **AutogeneS PBMC deconvolution results**

### Across samples for each cell type



## Summary

- Good predictions for most cell types
- This matrix could be used to deconvolute PBMC RNAseq datasets for major cell types
- More work needed for low abundance cells
- For whole blood, we need to add a major cell type not present in PBMC: PMNs (primarily neutrophils)

## **Future Plans- Deconvolution**

1. Create new GE signature matrix (GESM) from sc PBMC+PMN to test methods against

### datasets (with TRUTH):

- 48 Sal and Flu samples with flow cytometry cell type data and Quantseq RNA data
- 42 RFI WB samples with flow cytometry cell type data and RNAseq data

### Also test other tools such as CIBERSORTx, etc.

Will also use purified cell populations as additional verification/GESM creation

- 2. Accurate deconvolution of NDCM dataset: 2400+ Quantseq samples
- 3. Develop "short-hand" version of deconvolution GE signature with Nanostring genes

 $\rightarrow$  practical tool for deconvolution for phenotyping at cell type level to eliminate RNA prep, Quantseq and bioinformatics costs

## **Acknowledgments**





Juber Herrera Christopher Uribe Tuggle

Ryan Corbett

Pengxin Yang

Muskan Kapoor



Kapoor

Carrie Meeks









Crystal Loving

Jayne Wiarda Kristen Byrne

Sathesh Sivasankaran



Jack Dekkers



**IOWA STATE UNIVERSITY** 

Application of Genomics to Improve Disease Resilience and Sustainability

NIFA Project #2021-67015-34562





Joan Lunney







UNIVERSITY OF ALBERTA

National Institute of Food and Agriculture

in Pork Production

NIFA 2017-67007-26144 + 2018-67015-2701



