

# Search for and transformation of human cells and cell types with latent space representations

This manuscript ([permalink](#)) was automatically generated from [greenelab/czi-seed-rfa@f4557d2](#) on November 3, 2018.

## Authors

---

- **Loyal A. Goff**

 [0000-0003-2875-451X](#) ·  [loyale](#) ·  [loyalgoff](#)

Solomon H. Snyder Department of Neuroscience, Johns Hopkins University School of Medicine; Kavli Neurodiscovery Institute, Johns Hopkins University; McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University School of Medicine

- **Casey S. Greene**

 [0000-0001-8713-9213](#) ·  [cgreene](#) ·  [greenescientist](#)

Department of Systems Pharmacology and Translational Therapeutics, Perelman School of Medicine, University of Pennsylvania

- **Stephanie C. Hicks**

 [0000-0002-7858-0231](#) ·  [stephaniehicks](#) ·  [stephaniechicks](#)

Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health

- **Rob Patro**

 [0000-0001-8463-1675](#) ·  [rob-p](#)

Department of Computer Science, Stony Brook University

- **Elana J. Fertig**

 [0000-0003-3204-342X](#) ·  [ejfertig](#) ·  [FertigLab](#)

Department of Oncology, Sidney Kimmel Comprehensive Cancer Center, School of Medicine, Johns Hopkins University; Department of Applied Mathematics and Statistics, Whiting School of Engineering, Johns Hopkins University

- **Michael I. Love**

 [0000-0001-8401-0545](#) ·  [mikelove](#) ·  [mikelove](#)

Department of Biostatistics, University of North Carolina at Chapel Hill; Department of Genetics, University of North Carolina at Chapel Hill

- **Jane Roe**

 [XXXX-XXXX-XXXX-XXXX](#) ·  [janeroe](#)

Department of Something, University of Whatever; Department of Whatever, University of Something

# Abstract

---

**Instructions:** Describe your collaborative project, highlighting key achievements of the project; limited to 250 words.

## Five Key References

---

- Hicks refs: [1]
- projectR & scCoGAPS: [2]
- Alevin: [3]

## Project Team

---

### PI information

1. Loyal Goff (Submitter)
  - Title: Assistant Professor
  - Degrees: PhD
  - Type of organization: Academic
  - Tax ID: 52-0595110 (JHU)
  - Email: loyalgoff@jhmi.edu
2. Stephanie Hicks
  - Title: Assistant Professor
  - Degrees: PhD
  - Type of organization: Academic
  - Tax ID: 52-0595110 (JHU)
  - Email: shicks19@jhu.edu
3. Elana Fertig
  - Title: Associate Professor
  - Degrees: PhD
  - Type of organization: Academic
  - Tax ID: 52-0595110 (JHU)
  - Email: ejfertig@jhmi.edu
4. Casey Greene
5. Tom Hampton

## 6. Michael Love

- Title: Assistant Professor
- Degrees: Dr. rer. nat.
- Type of organization: Academic
- Tax ID: 56-6001393 (UNC)
- Email: milove@email.unc.edu

## 7. Rob Patro

- Title: Assistant Professor
- Degrees: PhD
- Type of Organization: Academic
- Tax ID: 16-1514621 (Stony Brook)
- Email: rob.patro@cs.stonybrook.edu

## Description (750 words TOTAL)

### 1. Loyal Goff

2. Stephanie C. Hicks is an Assistant Professor of Biostatistics at the Johns Hopkins Bloomberg School of Public Health. She is an expert in statistical methodology with a strong track record in processing and analyzing single-cell genomics data, including extensive experience developing fast, memory-efficient R/Bioconductor software to remove systematic and technical biases from scRNA-seq data [1]. Dr. Hicks will work together with Co-PIs to implement fast search algorithms in latent spaces (Aim 1) and to implement the methods developed into fast, scalable, and memory-efficient R/Bioconductor software packages (Aim 3).

3. Elana Fertig is an Associate Professor of Oncology and Applied Mathematics and Statistics at Johns Hopkins University. She developed the Bayesian non-negative matrix factorization algorithm CoGAPS [4] for latent space analysis. In collaboration with co-PI Goff, she adapted this tool to scRNA-seq data and developed a new transfer learning framework to relate the low-dimensional features in scRNA-seq data across data modalities, biological conditions, and organisms [2]. Dr. Fertig will work with the co-PIs to incorporate the error models from Aim 1 into the latent space representations, dimensionality estimation, and biological assessment metrics in Aim 2. She is developing standardized language for latent space representation in collaboration with co-PIs Goff and Greene [5] that will provide a strong foundation for standardization of these approaches across different unsupervised learning tools.

### 4. Casey Greene

### 5. Tom Hampton

6. Michael Love is an Assistant Professor of Biostatistics and Genetics at the University of North Carolina at Chapel Hill. He is a leading developer of statistical software for RNA-seq analysis in the Bioconductor Project, maintaining the widely used DESeq2 [6] and tximport [7] packages. He is a close collaborator with Dr. Rob Patro on bias-aware estimation of transcript

abundance from RNA-seq and estimation of uncertainty during transcript quantification [8].

Dr. Love will work with co-PIs to disseminate versioned reference cell type catalogs through widely used frameworks for genomic data analysis including R/Bioconductor and Python.

7. Rob Patro

## Proposal Body (2000 words)

---

High dimensional data can often be compressed into fewer dimensions without a substantial loss of information. For transcriptomic data, compressing on the gene dimension is most attractive: it can be done a single sample at a time, and genes are often co-regulated and thus provide information about each other. In the best case, the fewer dimensions both capture the biological sources of variability while ignoring noise and each dimension or small combination of dimensions aligns to interpretable biological processes.

For the Human Cell Atlas (HCA), low-dimensional representations enable efficient search and transformation, with the benefits becoming particularly pronounced as the number of cells and tissues becomes particularly large. Our **central hypothesis** is that these approaches will also enable search at the tissue and cell level of biologically-meaningful features and transformation in biologically meaningful ways, e.g., from a healthy to a disease context.

We propose to advance **base enabling technologies** for low-dimensional representations and to develop solutions for search and transformation that can be practically applied across the entirety of the HCA. We will continue to develop techniques that seek to learn interpretable biologically-aligned representations.

Latent space determination relies on computational techniques that learn low dimensional structure from high dimensional data. A wide variety of computational methods for dimensionality reduction have been applied to single cell data. These techniques fall largely into two classes: (1) linear techniques based upon matrix factorization and (2) non-linear techniques based upon manifold learning. Each specific technique selected for analysis will reveal distinct structures in the data. Much of the current research emphasizes demonstration of optimality of a single method. In contrast to this prevailing view, we posit that each method reveals a unique set of features associated with a set of *a priori* unknown biological processes. Together, the distribution of the ensemble of low dimensional representations from these methods will uncover the multiple, concurrent biological processes from high dimensional datasets. Previously, we have developed a common language for interpretation of matrix factorization methods to facilitate a unification of latent space methods [5]. The latent space team in the previous round of funding (containing PIs Fertig, Goff, Greene, and Patro) is also developing standardized output formats for these distinct methods. Together, this previous work will facilitate practical unification and biological assessment of latent spaces in this seed network.

Numerous factors impact the biological accuracy of latent spaces learned from dimension reduction techniques. Two notable examples include the preprocessing of the input data (Aim 1) and dimensionality of the low dimensional representation (Aim 2). The impact of each of these components on inference must be quantified and isolated before latent methods can be compared or unified. To that end, we select a single dimension reduction technique as the foundation for the technologies developed in our seed network. We use the Bayesian, non-negative matrix factorization method scCoGAPS [9] (PI Fertig). Previous work has demonstrated the robustness of this approach in inferring biologically relevant low dimensional representations of perturbation [10] and time course [12] from bulk datasets, leading to the winning solution in the HPN DREAM8 challenge [14]. Moreover, recent extension of this method to single cell data in the developing mouse retina distinguished simultaneously cellular identity, dynamic trajectories, and cell state within a single analysis [9]. Notably for the work proposed in this seed network, this method includes an uncertainty estimate that can be readily modified to account for measurement-specific technical variation [15]. Moreover, this same study demonstrated that low dimensional factorizations of tumor populations separate tumor samples relative to normals whereas higher dimensionalities separate tumor subtypes. Thus, CoGAPS has been shown to yield different factorizations reflective of biological hierarchies at different dimensionalities, rather than having a single dimensionality with a “correct” factorization. While technology development will focus on scCoGAPS, the tools developed as part of the seed network will be generalizable and compared to additional linear and non-linear methods developed in the literature **ADD RELEVANT CITATIONS** and by members of the broader CZI consortium. This single-method focus with subsequent expansion to additional methods will also facilitate educational initiatives to advance broader understanding of the interpretation of latent space solutions from high dimensional data (Aim 3).

*About fast and improved quant* Existing approaches for quantification from scRNA-seq data using tagged-end protocols (e.g. 10x Chromium, drop-Seq, inDrop, etc.) have no mechanism for accounting for reads mapping between multiple genes in the resulting quantification estimates, and therefore such reads are typically removed **CORRECT?**. This reduces quantification accuracy, and leads to systematic biases in gene expression estimates that correlate problematically with the size of gene families and gene function [3]. We have recently developed a novel approach for quantification from tagged-end data that accounts for reads mapping to multiple genomic loci (~15-25% of the reads in a typical experiment) in a principled and consistent way. We propose to expand on this work in a number of ways, building these capabilities into a production quality tool for the processing of scRNA-seq data. (1) Exploration of alternative models for UMI resolution, including a maximum likelihood (as opposed to maximum parsimony) model, (2) Development of new approaches for quality-control and filtering using features from the UMI-resolution graph, (3) Creation of a compressed and indexable data structure for the UMI-resolution graph to allow direct access, query, and, eventually, search.

The technologies to improve quantification will have a critical impact on the outcomes of latent spaces. However, there are currently no standardized, quantitative metrics to determine relative

uncovering of biology from low dimensional representations. We have developed new transfer learning methods to quantify the extent to which latent space representations from one set of training data are represented in another [2]. These tools provide a strong foundation to enable biological quantification of latent space representations by quantifying the extent to which those spaces transfer across datasets of related biological contexts.

## Aim 1

*Rationale:* need to add

*Fast search algorithms* (in low dimensions, quantifying differences between case and reference maps, twist: shared latent spaces / k-mers):

*Finding differences between maps:* The fast search algorithms will then be used quantify differences between a reference transcriptome map (the Human Cell Atlas) and non-reference transcriptome maps from other samples of interest, similar to the idea of using a reference genome to identify genomic differences in between a reference and non-reference genome. Globally quantifying differences between transcriptome maps is important because it allows for quantification of differences at the population or individual level between, for example, ten transcriptome maps from individuals with a particular phenotype to be compared to the Human Cell Atlas reference map (or other control transcriptome maps if available). Our metric to quantify differences will depend on the distributions of cell expression within and between individuals, which Dr. Hicks has extensive experience with [17]. We will leverage not only the cell-to-cell correlation structure within one transcriptome map (or human individual), but also the correlation structure across transcriptome maps (or multiple human individuals), which will share common latent spaces across individuals for a particular phenotype. Our initial approach will be to use linear mixed models to account for the correlation structure within and between transcriptome maps. The statistical method will be fast, memory-efficient and will scale to billions of cells because we work in the latent space with a significantly reduced number of dimensions (instead of billions, just hundreds or thousands).

*New models for UMI deduplication accounting for transcript-level information:* (Rob)

*Parsimony & likelihood based, integrated with gene-level uncertainty:*

- \* Everything FAST! API for search against HCA reference? (Rob)
- \* k-mer / quantified latent spaces (Casey / Rob)

## Aim 2

*Rationale:* Biological systems are comprised of diverse cell types with overlapping molecular phenotypes, and biological processes are often reused with modifications across cellular contexts. The functional output of these systems is determined by the interactions between these complex

components, rather than a single gene or cell. This suggests that fundamental biological mechanisms may broadly contribute to an observed state, with context-specific modifiers conferring selective susceptibility to disease. Latent space techniques are poised to reveal these fundamental mechanisms in the broad survey of single cell data across model systems and cellular contexts in the Human Cell Atlas. We hypothesize that the features learned from these techniques will define constitutive basis vectors that reflect discrete biological processes or features. Thus, these basis vectors will be shared across different biological systems, with context-specific perturbations indicating pathogenic differences in disease. *We propose a central suite of statistics for assessment and interpretation of latent space tools to define the identity and dimensionality of biological systems.*

*Quantifying latent space estimation with transfer learning:* A critical challenge to latent space methods is the quantification of methods performance. Numerous computational metrics have been developed to assess convergence of the low dimensional estimation. However, these metrics do not quantify whether the features in a low dimensional representation of scRNA-seq data represent biological processes in the measured system. The performance of these methods can be quantified directly in datasets for which cell types and states are known (e.g., perturbation experiments, controlled admixture experiments, etc). However, these annotations are lacking in most biological datasets limiting any such quantification. Transfer learning methods have been developed in machine learning to relate features learned in a source dataset to those in a new, target dataset in order to transfer annotations from one context to another. In this project, we will adapt these methods to quantify the performance of latent space methods by the extent to which learned low dimensional features from a source dataset transfer to a target dataset in a related biological context. We will benchmark the performance of the resulting metric on simulated datasets, cross-validation in scRNA-seq datasets with known cell types and states, and cross-study validation of systems in related biological contexts with known cell types and states. Gene set enrichment methods will also be used to explore the relevant biological processes described by individual basis vectors, and related bases will be identified through clustering and exploratory approaches in these benchmark datasets. Our transfer learning based metric will be piloted on low dimensional representations learned with scCoGAPS and then applied to a broader suite of latent space tools. We will release software for this transfer learning quantification of latent space representations in R and Python using standard latent space file formats developed by our team in the first year of HCA funding.

*Dimensionality estimation:* Dimensionality reduction methods are sensitive to the number of low features learned in each dataset. Many computational techniques optimize dimensionality by creating a cost function which penalizes models with higher number of features. Similar to the quantification metrics, these penalty terms do not reflect the extent to which features learned at a given dimensionality reflect biology. Moreover, many systems may have more than one biologically accurate low dimensional representation. Such multiple truths in data would be particular prominent in systems that can be subdivided into hierarchical classifications. For example, in the



case of cancer we observed that a low dimensional representation of bulk data learned from CoGAPS distinguished cancers from normals whereas a higher dimension distinguished tumor subtypes [15]. Both of these low dimensional representations are equally valid, and each reflects different biological features in the data. To find these multiple truths, we will develop a parallel framework to run scCoGAPS for multiple dimensionalities and quantify performance with our transfer-learning based metric on random subsets of the data. The dimensions with greatest cross-validated feature robustness will be retained as the optimal dimensionalities for each dataset. We will develop software to enable this cross-validation dimensionality estimation across multiple latent space methods. We note that this same software will provide a robust tool to define ensembles of low dimensional representations that reflect underlying biology learned across multiple latent space methods. **Rob: I'm not sure if you want to fill in some of your ideas re persistent homology instead. Very open to that idea and think it may be a nice, more efficient methodology than what's proposed here.**

*Search tool for latent spaces and reference cell types:* **Loyal, Casey – what are the datasets that will be used for this – I would think all healthy cells in a single system to enable quantification of context-specific in the next part of this aim.** Comprehensive identification of basis vectors across conditions is an area of active research for our group in the previous funding period. We will use scCoGAPS and other tools developed within our collaborative network to establish a compendium of basis vectors across our single cell catalog. Ensembles of the low dimensional features that represent robust biological features across methods using methods described above will be preserved as the 'biological basis' of the Human Cell Atlas. The weights of these bases will be correlated across all available metadata attributes for each cell to identify basis vectors that are associated with specific cellular contexts, disease states, technical parameters, or other phenotypic features. A reference catalog of gene weights for specific cell types will be defined by the set of basis vectors associated with cellular identity in datasets with known ground truth. We will adapt the software we developed for transfer learning of features from bulk data recount [16] to facilitate querying of signatures in new user-defined datasets (delivery of which is described in the next aim). As datasets accumulate and methods are refined, the biological basis and reference catalog of gene weights will evolve over time. To enable reproducible research leveraging HCA, we will implement a content-based versioning system, which identifies versions of the reference cell type catalog by the gene weights and transcript nucleotide sequences using a hash function. Such a hash-based versioning and provenance identification and detection framework has proven successful in the bulk RNA-seq context to support reproducible computational analyses [18].

*Differentiating context-specific latent spaces from latent spaces that are universal across biological contexts:* The search tool to define reference cell types based upon latent spaces was defined for healthy tissues from XXX (some control). Deviations of common cell types or states from the healthy baseline in other populations will indicate context-specific alterations, which may be associated with disease. To identify potentially pathogenic responses in target datasets, we will implement a random forest classifier into our transfer learning method to segregate cells based on

their usage of disease-associated basis vectors after projection. In other cases, disease may arise from changes in variation reflective of inter-cellular heterogeneity. Therefore, we will also develop methods to quantify variation from latent space vectors. Both methods will be incorporated in our latent space search tool. **Loyal: I'm not sure if this is what you had in mind. It may also be that these are reflected in the hierarchy of dimensionality – may want to incorporate here.**

### Aim 3

*Rationale:* Low dimensional representations provide a powerful means to analyze scRNA-seq and HCA data to make tasks faster, perform more biologically grounded analyses, and provide interpretable summaries of complex high-dimensional data. However, using these capabilities to the fullest extent requires an investment in robustness and evaluation, integration with widely used toolkits, and in a scalable education effort that can reach students from the undergraduate level and beyond. *We propose to enhance software usability and develop open instructional materials that we will use to deliver short-course training that includes the topics of single cell profiling, machine learning methods, low-dimensional representations, reference cell type catalogs, and tools developed by our group in response to this RFA.*

#### Aim 3.1: Training for next-generation of single cell data scientists (Tom)

Training / teaching (scRNAseq, low-dimensional representations, reference cell type catalog, RFA-developed tools)

Although the HCA data set will greatly increase the rate of biological discovery across many biomedical fields, a background in bioinformatics will be required to these data effectively. Our mission is to make this obstacle easy for most biologists and physician scientists to overcome, and to enable scientists who already possess a background in bioinformatics to use software created in this project effectively. We will address these goals by offering short courses and creating distributed communities.

1. Short Courses. Short courses (3-5 days) will introduce all features of bioinformatics required to access HCA and use the tools developed in this application successfully, such as:
  - Visualization and Exploration of High Dimensional Data
  - R Statistical Programming Environment
  - UNIX
  - Statistical Approaches for High Dimension Biomedical Data
  - Gene Set and Pathway Analysis
  - Low dimensional representations of high dimensional data
  - R and UNIX Tools Specific to HCA Data Analysis
  - scRNAseq tools
  - RFA-developed tools

2. Course Project. The course project will provide a hands-on experience of retrieving data from HCA, performing quality control, exploratory data analysis, identification of differentially abundant clusters, identification of marker genes associated with clusters, differential gene expression, and interpretation of results in the context of other data such as Gene Ontology and KEGG pathways. Groups of participants share their findings in a final presentation.
3. Distributed Communities. Where short courses jump start participants' ability to perform basic tasks and learn the tools relevant to scRNA-seq, distributed communities provide the opportunity of mastery of these techniques, which includes the ability to teach others. Our support of these communities will include:
  - Workshops and Seminars For Group Leaders
  - Access to Community-Developed Curricula and Documentation
  - Scholarships for Short Courses

### **Aim 3.2: Software Engineering and Integration into Bioconductor (Stephanie and Casey)**

We will implement the proposed methods from Aims 1 and 2 into robust software, which will be integrated into the R/Bioconductor and Python frameworks. The software will be fast, scalable, and memory-efficient because will leverage the computational tools previously developed by Bioconductor for single-cell data access to the HCA, data representation (`SingleCellExperiment`, `beachmat`, `DelayedArray`, `HDF5Array` and `rhdf5`) and data assessment and ameliorization of data quality (`scater`, `scraper`, `DropletUtils`).

### **Aim 3.3: Data Integration into Bioconductor (Stephanie and Mike)**

We will integrate the catalogs of reference cell types (Aim 2) into the R/Bioconductor and Python frameworks. The use of consortia data summaries as annotation and interpretive scaffold for outside datasets has proven widely successful for the ENCODE, Roadmap Epigenome Mapping, and GTEx projects. We will use HCA data to define similar summaries as annotation and scaffolding for local genomic datasets. Additional information about biological variability that can be transferred from individual or ensemble latent spaces, can also be used landmarks to ordinate differences among populations of cells observed in a local datasets, whether bulk or single-cell experiments. We will package and version reference cell types, including measures of technical and biological variability transferred from the latent spaces defined in earlier aims, and deliver these as structured data objects in Bioconductor and Python. We will leverage our expertise in working with Bioconductor both as core package developers and power users to enable on-the-fly downloading of reference cell types and variability via the *AnnotationHub* framework, which includes rich specification of metadata including the provenance and versioning of the catalog. Python workflows leveraging existing frameworks for genomic data analysis (Biopython, BioRanges, and bx-python) will also be supported via cross-platform packages such as Feather

that facilitate transfer of data objects between languages. We will develop *F1000Research* workflows for R and python demonstrating how HCA-defined reference cell types and tools developed in this RFA can be used within a typical genomic data analysis.

# References

---

**1. Missing data and technical variability in single-cell RNA-sequencing experiments**

Stephanie C Hicks, F William Townes, Mingxiang Teng, Rafael A Irizarry

*Biostatistics* (2017-11-06) <https://doi.org/gfb8g4>

DOI: [10.1093/biostatistics/kxx053](https://doi.org/10.1093/biostatistics/kxx053) · PMID: [29121214](https://pubmed.ncbi.nlm.nih.gov/29121214/)

**2. Decomposing cell identity for transfer learning across cellular measurements, platforms, tissues, and species.**

Genevieve L Stein-O'Brien, Brian S. Clark, Thomas Sherman, Christina Zibetti, Qiwen Hu, Rachel Sealfon, Sheng Liu, Jiang Qian, Carlo Colantuoni, Seth Blackshaw, ... Elana J. Fertig

*Cold Spring Harbor Laboratory* (2018-08-20) <https://doi.org/gd2xpn>

DOI: [10.1101/395004](https://doi.org/10.1101/395004)

**3. Alevin efficiently estimates accurate gene abundances from dscRNA-seq data**

Avi Srivastava, Laraib Malik, Tom Sean Smith, Ian Sudbery, Rob Patro

*Cold Spring Harbor Laboratory* (2018-06-01) <https://doi.org/gffk42>

DOI: [10.1101/335000](https://doi.org/10.1101/335000)

**4. CoGAPS: an R/C++ package to identify patterns and biological process activity in transcriptomic data**

Elana J. Fertig, Jie Ding, Alexander V. Favorov, Giovanni Parmigiani, Michael F. Ochs

*Bioinformatics* (2010-09-01) <https://doi.org/cwqsv4>

DOI: [10.1093/bioinformatics/btq503](https://doi.org/10.1093/bioinformatics/btq503) · PMID: [20810601](https://pubmed.ncbi.nlm.nih.gov/20810601/) · PMCID: [PMC3025742](https://pubmed.ncbi.nlm.nih.gov/PMC3025742/)

**5. Enter the Matrix: Factorization Uncovers Knowledge from Omics**

Genevieve L. Stein-O'Brien, Raman Arora, Aedin C. Culhane, Alexander V. Favorov, Lana X. Garmire, Casey S. Greene, Loyal A. Goff, Yifeng Li, Aloune Ngom, Michael F. Ochs, ... Elana J. Fertig

*Trends in Genetics* (2018-10) <https://doi.org/gd93tk>

DOI: [10.1016/j.tig.2018.07.003](https://doi.org/10.1016/j.tig.2018.07.003) · PMID: [30143323](https://pubmed.ncbi.nlm.nih.gov/30143323/)

**6. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2**

Michael I Love, Wolfgang Huber, Simon Anders

*Genome Biology* (2014-12) <https://doi.org/gd3zvn>

DOI: [10.1186/s13059-014-0550-8](https://doi.org/10.1186/s13059-014-0550-8) · PMID: [25516281](https://pubmed.ncbi.nlm.nih.gov/25516281/) · PMCID: [PMC4302049](https://pubmed.ncbi.nlm.nih.gov/PMC4302049/)

## **7. Differential analyses for RNA-seq: transcript-level estimates improve gene-level inferences**

Charlotte Soneson, Michael I. Love, Mark D. Robinson

*F1000Research* (2015-12-30) <https://doi.org/gdtgw8>

DOI: [10.12688/f1000research.7563.1](https://doi.org/10.12688/f1000research.7563.1) · PMID: [26925227](https://pubmed.ncbi.nlm.nih.gov/26925227/) · PMCID: [PMC4712774](https://pubmed.ncbi.nlm.nih.gov/PMC4712774/)

## **8. Salmon provides fast and bias-aware quantification of transcript expression**

Rob Patro, Geet Duggal, Michael I Love, Rafael A Irizarry, Carl Kingsford

*Nature Methods* (2017-03-06) <https://doi.org/gcw9f5>

DOI: [10.1038/nmeth.4197](https://doi.org/10.1038/nmeth.4197) · PMID: [28263959](https://pubmed.ncbi.nlm.nih.gov/28263959/) · PMCID: [PMC5600148](https://pubmed.ncbi.nlm.nih.gov/PMC5600148/)

## **9. Comprehensive analysis of retinal development at single cell resolution identifies NFI factors as essential for mitotic exit and specification of late-born cells**

Brian Clark, Genevieve Stein-O'Brien, Fion Shiau, Gabrielle Cannon, Emily Davis, Thomas Sherman, Fatemeh Rajaii, Rebecca James-Esposito, Richard Gronostajski, Elana Fertig, ... Seth Blackshaw

*Cold Spring Harbor Laboratory* (2018-07-27) <https://doi.org/gdwrzh>

DOI: [10.1101/378950](https://doi.org/10.1101/378950)

## **10. Gene expression signatures modulated by epidermal growth factor receptor activation and their relationship to cetuximab resistance in head and neck squamous cell carcinoma**

Elana J Fertig, Qing Ren, Haixia Cheng, Hiromitsu Hatakeyama, Adam P Dicker, Ulrich Rodeck, Michael Considine, Michael F Ochs, Christine H Chung

*BMC Genomics* (2012) <https://doi.org/gb3fgp>

DOI: [10.1186/1471-2164-13-160](https://doi.org/10.1186/1471-2164-13-160) · PMID: [22549044](https://pubmed.ncbi.nlm.nih.gov/22549044/) · PMCID: [PMC3460736](https://pubmed.ncbi.nlm.nih.gov/PMC3460736/)

## **11. CoGAPS matrix factorization algorithm identifies transcriptional changes in AP-2alpha target genes in feedback from therapeutic inhibition of the EGFR network**

Elana J. Fertig, Hiroyuki Ozawa, Manjusha Thakar, Jason D. Howard, Luciane T. Kagohara, Gabriel Krigsfeld, Ruchira S. Ranaweera, Robert M. Hughes, Jimena Perez, Siân Jones, ... Christine H. Chung

*Oncotarget* (2016-09-16) <https://doi.org/f9k8d8>

DOI: [10.18632/oncotarget.12075](https://doi.org/10.18632/oncotarget.12075) · PMID: [27650546](https://pubmed.ncbi.nlm.nih.gov/27650546/) · PMCID: [PMC5342018](https://pubmed.ncbi.nlm.nih.gov/PMC5342018/)

**12. Pattern Identification in Time-Course Gene Expression Data with the CoGAPS Matrix Factorization**

Elana J. Fertig, Genevieve Stein-O'Brien, Andrew Jaffe, Carlo Colantuoni

*Gene Function Analysis* (2013-10-24) <https://doi.org/f5j7xj>

DOI: [10.1007/978-1-62703-721-1\\_6](https://doi.org/10.1007/978-1-62703-721-1_6) · PMID: [24233779](https://pubmed.ncbi.nlm.nih.gov/24233779/)

**13. Integrated time course omics analysis distinguishes immediate therapeutic response from acquired resistance**

Genevieve Stein-O'Brien, Luciane T. Kagohara, Sijia Li, Manjusha Thakar, Ruchira Ranaweera, Hiroyuki Ozawa, Haixia Cheng, Michael Considine, Sandra Schmitz, Alexander V. Favorov, ... Elana J. Fertig

*Genome Medicine* (2018-05-23) <https://doi.org/gfc4dq>

DOI: [10.1186/s13073-018-0545-2](https://doi.org/10.1186/s13073-018-0545-2) · PMID: [29792227](https://pubmed.ncbi.nlm.nih.gov/29792227/) · PMCID: [PMC5966898](https://pubmed.ncbi.nlm.nih.gov/PMC5966898/)

**14. Inferring causal molecular networks: empirical assessment through a community-based effort**

Steven M Hill, Laura M Heiser, Thomas Cokelaer, Michael Unger, Nicole K Nesser, Daniel E Carlin, Yang Zhang, Artem Sokolov, Evan O Paull, ... Sach Mukherjee

*Nature Methods* (2016-02-22) <https://doi.org/f3t7t4>

DOI: [10.1038/nmeth.3773](https://doi.org/10.1038/nmeth.3773) · PMID: [26901648](https://pubmed.ncbi.nlm.nih.gov/26901648/) · PMCID: [PMC4854847](https://pubmed.ncbi.nlm.nih.gov/PMC4854847/)

**15. Preferential Activation of the Hedgehog Pathway by Epigenetic Modulations in HPV Negative HNSCC Identified with Meta-Pathway Analysis**

Elana J. Fertig, Ana Markovic, Ludmila V. Danilova, Daria A. Gaykalova, Leslie Cope, Christine H. Chung, Michael F. Ochs, Joseph A. Califano

*PLoS ONE* (2013-11-04) <https://doi.org/gcpgc6>

DOI: [10.1371/journal.pone.0078127](https://doi.org/10.1371/journal.pone.0078127) · PMID: [24223768](https://pubmed.ncbi.nlm.nih.gov/24223768/) · PMCID: [PMC3817178](https://pubmed.ncbi.nlm.nih.gov/PMC3817178/)

**16. MultiPLIER: a transfer learning framework reveals systemic features of rare autoimmune disease**

Jaclyn N Taroni, Peter C Grayson, Qiwen Hu, Sean Eddy, Matthias Kretzler, Peter A Merkel, Casey S Greene

*Cold Spring Harbor Laboratory* (2018-08-20) <https://doi.org/gfc9bb>

DOI: [10.1101/395947](https://doi.org/10.1101/395947)

**17. quantro: a data-driven approach to guide the choice of an appropriate normalization method.**

Stephanie C Hicks, Rafael A Irizarry

*Genome biology* (2015-06-04) <https://www.ncbi.nlm.nih.gov/pubmed/26040460>

DOI: [10.1186/s13059-015-0679-0](https://doi.org/10.1186/s13059-015-0679-0) · PMID: [26040460](https://pubmed.ncbi.nlm.nih.gov/26040460/) · PMCID: [PMC4495646](https://pubmed.ncbi.nlm.nih.gov/PMC4495646/)

#### 18. **tximeta**

Rob Patro Michael Love

*Bioconductor* (2018) <https://doi.org/gfddxw>

DOI: [10.18129/b9.bioc.tximeta](https://doi.org/10.18129/b9.bioc.tximeta)