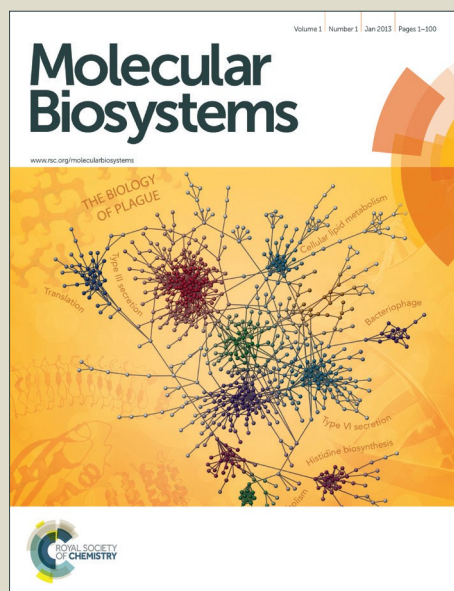


Molecular BioSystems

Accepted Manuscript



This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this *Accepted Manuscript* with the edited and formatted *Advance Article* as soon as it is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.



www.rsc.org/molecularbiosystems

1 **Reconstructing and analysing cellular states, space and time from**
 2 **gene expression profiles of many cells and single cells**

3 Mirko Francesconi^{a,b} & Ben Lehner^{a,b,c*}

4 **Affiliations**

5
 6 ^a EMBL-CRG Systems Biology Unit, Centre for Genomic Regulation (CRG), Dr. Aiguader
 7 88, 08003 Barcelona, Spain

8 ^b Universitat Pompeu Fabra (UPF), Dr. Aiguader 88, 08003 Barcelona, Spain

9 ^c Institució Catalana de Recerca i Estudis Avançats (ICREA), 08010 Barcelona, Spain

10 *Corresponding author

11 E-mail: ben.lehner@crg.eu

12
 13
 14 **Abstract**

15 Genome-wide gene expression profiling is a fast, cheap and standardised analysis that
 16 provides a high dimensional measurement of the state of a biological sample. In this
 17 review we describe some of the computational methods that can be applied to identify
 18 and interpret sources of variance in gene expression in whole organisms, organs,
 19 tissues or single cells. This allows the identification of constituent cell types and states
 20 in complex mixtures, the reconstruction of temporal trajectories of development,
 21 differentiation and progression, and the reconstruction of spatial patterning. When
 22 applied to genetically variable samples, these methods allow the efficient investigation
 23 of how genetic variation influences gene expression in space and time.
 24

1. Sources of variability in biological samples

Differences in gene expression measured by RNA sequencing or using DNA microarrays can be purely technical (this can be estimated¹) or due to biological differences between samples. These biological differences may be part of the experimental design or they may be due to uncontrolled experimental variation in the state of each sample²⁻⁷.

Uncontrolled experimental variation is usually regarded as a confounder but it might also be interesting when the sources of variation are correctly identified and understood (Fig1). In this review we provide an overview of the computational methods that can be used to identify and understand controlled and uncontrolled variance in gene expression datasets and highlight examples of how these have been used to make interesting discoveries.

In gene expression data from mixtures of cell types such as tissues, organs or tumours, an important source of variance derives from cell heterogeneity. This may come from the process of interest, for example a condition might change the proportion of different cell types within tissues⁸. Or different amounts of contaminating tissues could confound the analysis, for example with patient biopsies. In this case it is important to factor out this confounder to increase the power to detect differences in the tissue of interest.

Even when analysing populations of sorted pure cell types (or unicellular organisms) heterogeneity is still present. Cells are not static entities but they dynamically adjust their state in response to environmental stimuli. It might be hard to control all (micro-) environmental factors that can trigger some cell response. A typical example is growth rate: cells adjust their global gene expression according to growth rate - increasing for example ribosomes and translation-related genes and decreasing stress related genes - no matter what the growth-rate-limiting factor is^{9, 10}. Any perturbation that changes growth rate will consequently impact gene expression.

Moreover cells go through oscillations, for example cell cycles, metabolic cycles and circadian rhythms. This heterogeneity will be present in unsynchronized single cells but also potentially in bulk measures of asynchronous populations where oscillatory dynamics are convoluted in the average expression data; population measures reflect the average of gene expression of different cell cycle or other oscillatory stages

1 weighted by the fractions of cells that are at each stage. Conditions, treatments (or
2 confounders) that change the oscillatory dynamics will change these fractions and will
3 leave a trace on gene expression.

4
5 Even more complex dynamics are present when studying systems undergoing
6 development or differentiation. When studying fast developing model systems such as
7 *Drosophila melanogaster* or *Caenorhabditis elegans* that complete development, even a
8 few hours' difference in precise staging can introduce substantial variation in gene
9 expression^{5, 11, 12}. This might be related to the experimental conditions of interest (that
10 could cause a delay in development) or it could be a confounder if not properly
11 controlled (experimental batches might be at slightly different developmental stages).

12
13 In summary, even in a simple perturbation and expression profiling experiment it is
14 important to understand and correctly decompose gene expression variance into the
15 corresponding sources. This will help to: (1) better understand and interpret the global
16 effects of a treatment/condition/mutation, for example as developmental delays or
17 changes in growth rate or cell composition; (2) tease apart specific effects and direct
18 targets of treatment/condition/mutation beyond the developmental delays or changes
19 in growth rate; and (3) control for experimental confounders and increase the power to
20 detect the effects of interest. However, the real power of understanding and
21 decomposing expression profiles is the application to large datasets in which new cell
22 states, spatial patterning or temporal 'trajectories' of expression can be identified. In
23 addition, if this is performed in combination with genetically varying samples, the
24 impact of genetic variation on these states or spatial and temporal patterns can be
25 determined.

26
27 **2. Computational approaches for inferring cell states, sample compositions, time and**
28 **space from gene expression**

29
30 Global gene expression data is typically highly redundant because many genes (for
31 example those involved in the same biological processes) share correlated expression
32 profiles. Thus it is useful to represent (map) the high dimensional data listing the

expression of tens of thousands of genes onto a lower and more interpretable dimensional space, a task known as dimensionality reduction. Ideally after this step the data would be represented by a few dimensions that account for most of the variance in the data and where each dimension represents a distinct, interpretable biological process. At the same time it is desirable to filter out non-biological, uninteresting variance (technical noise).

***Ab initio* (unsupervised) methods for dimensionality reduction: principal components analysis**

Depending on the goal of the study and on the nature of the biological process of interest, different computational approaches for reducing dimensionality can be used. If the goal is to discover new cell types, states or trajectories *ab initio* from the data, one of most widely used technique is principal component analysis¹³ (PCA). PCA rotates the data into new orthogonal coordinate systems where the axes (components) are linear combinations of the original variables and represent the directions of maximal variance in the data. This means that the first component explains most of the variance, the second component explains most of the residual variance after subtracting the first one, and so on. Thus retaining only the first few principal components accounts for most of the variance present in the original variables and filters out noise. PCA can be performed by eigen decomposition^{14, 1514, 1514, 1514, 1514, 1514, 1514, 1514, 15} of the covariance matrix of the data or, more efficiently, by singular value decomposition (SVD) of the data matrix^{14, 15}.

For each component PCA outputs singular values (or eigenvalues) that indicate the variance explained, sample scores (sample coordinates on the component) and gene loadings (the coefficients of the genes in the linear combination, i.e. how much each genes contribute to the component).

Visualizing the sample scores on the first few principal components helps to provide an overview of the global structure of the data, for example by highlighting clusters or trajectories (Fig2). The biological meaning of each component can be deduced by analysing the gene loadings, for example by traditional gene set enrichment analysis^{16, 17}. To this purpose Chung et al. developed a method to systematically identify genes significantly associated to principal components avoiding over-fitting¹⁸. However

1 interpreting principal components is not always easy because the components might be
2 enriched in many distinct biological processes and, vice versa, the same biological
3 process can be enriched in different components. The reasons for this include that in
4 PCA the components are defined maximizing the variance explained as a criterion, and
5 they are also constrained to be orthogonal. These conditions and constraints mean that
6 the principal components will not necessarily correspond to separate biological sources
7 of variance. This is an important limitation of PCA when the aim is to clearly separate
8 and remove unobserved confounders from interesting sources of variance (signal).
9 Using PCA there is a risk is that interesting signal is also removed together with
10 confounders.

11

12 **Other unsupervised methods to decompose variance**

13 Additional methods that relax some or all PCA constraints have been developed for
14 better separating sources of variance and to increase interpretability, known as factor
15 analysis methods. Similarly to PCA, these methods search for linear combinations of the
16 data (factors) that best explain the correlations among the variables but improve
17 interpretability by allowing for further rotations that better capture the underlining
18 structure of the data. For example the varimax¹⁹ method further rotates the data after
19 PCA (preserving orthogonality) in a way such that the genes have high or low loadings
20 only in one factor. Promax²⁰ also allows oblique rotations thus relaxing the
21 orthogonality constraint. Other methods based on factor analysis have been recently
22 proposed to better estimate sources of variance with the aim of correcting gene
23 expression from hidden confounders, such as surrogate variable analysis (SVA)²¹,
24 probabilistic estimation of expression residuals (PEER)²² and remove unwanted
25 variation (RUV)²³. The last two methods also allow estimating hidden factors in a semi-
26 supervised manner only on selected gene sets (i.e. control genes) to minimize the risk of
27 explaining away the signal together with confounders.

28

29 Among unsupervised method to deconvolve sources of signals, Independent Component
30 Analysis²⁴ (ICA) is one of the most flexible. Its rationale stems from central limit
31 theorem, which states that mixtures (convolution) of independent signals tend to be
32 normally distributed. Thus an effective strategy to separate the hidden independent
33 source signals from the measured mixed signal is to find linear combinations

(components) that maximize non-gaussianity (rather than variance as in PCA). Several measure of non-gaussianity have been proposed such as kurtosis, negentropy or mutual information²⁴. Components obtained by ICA are linearly independent (a stronger condition than uncorrelated as in PCA) but they do not need to be orthogonal as in PCA. ICA has been applied to gene expression data^{25, 26} and it can outperform PCA in teasing apart independent biological processes underlying expression differences²⁶(Fig3).

Using a reference expression dataset

If one wants to match data to predetermined states or types then a good approach is to compare the data to existing reference expression profiles. A simple approach is to use a subset of relevant genes from a reference dataset to build a model that predicts the corresponding state in the dataset of interest. For example, the expression of many genes in yeast correlates linearly with growth rate under many different conditions⁹ and a simple linear model including these genes can infer the relative growth rates of new conditions from gene expression²⁷. Similarly, the proportion of cells in a sample in each stage of the cell cycle can be inferred by comparison to reference datasets defining sets of genes activated at different phases during the cell cycle. The expression level in each of these genes in an asynchronous population (vector A) can then be expressed as the weighted average of their expression at each cell cycle stage (matrix R) where the weights (matrix W) are the unknown fractions of cells at each cell cycle stage in the population ($A=WR$). These fractions can be determined by solving the system for W ²⁸. The same modelling framework can be used to deconvolve cell type fractions from gene expression data of whole tissues when cell type-specific expression signature are known⁸. Similarly, cell type-specific expression profiles can be inferred in complex tissues if the fraction of each different cell type in the tissue is known²⁹.

More powerful methods to match expression data to a reference dataset include partial least squares (PLS) and canonical correlation analysis (CCA) ⁵. These two related statistical techniques analyse the relationship between two datasets (covariance for PLS, correlation for CCA) of multiple dependent and multiple independent variables such as two gene expression datasets measuring the same genes in two different sets of conditions. They decompose the covariance (or correlation in CCA) between the two datasets by finding linear combinations of the reference dataset that best explain linear

1 combinations of the independent dataset, in a manner similarly to PCA except that in
2 this case only the variance shared between the two datasets is taken into account. Using
3 these methods as multiple advantages: first one can find multiple processes shared with
4 a reference and quantify how much variance is explained by each. Further, the approach
5 leaves the variance not explained by the reference untouched. This avoids that specific
6 signals of interest are explained away together with global confounders when the aim is
7 to correct gene expression before downstream analysis. This is in contrast to using
8 reference genes whose expression in the data of interest might reflect a combination of
9 underlining processes some of which might not be present in the reference.

10

11 **Tackling non-linearity**

12 Many biological processes such as the cell cycle, development and differentiation show
13 complex non-linear dynamics such as oscillations or bifurcations. In these cases, linear
14 methods (such as those described above) are a useful first step to reduce
15 dimensionality, visualize the data and filter out noise, but they cannot directly be used
16 to order the data along a non-linear dynamic process. Reconstructing non-linear
17 dynamics from the data might be challenging because classical distance measures are
18 not appropriate to define, for example, how close two data points are in a trajectory and
19 hence their ordering (Fig4a).

20

21 In some cases, simple transformations can be used to infer the correct dynamics and to
22 order the data. For example if PCA (or ICA) transformed data lie on (a portion of) a cycle
23 in a low dimensional space, a simple transformation in polar coordinates can recover
24 the correct order of the data points along the dynamics^{5, 14, 15}.

25

26 In the case of more complex dynamics, finding the geometry of the data and ordering or
27 clustering data points might be harder. When studying development and the data points
28 lie in a single trajectory, ordering them can be seen as an instance of the well known
29 travelling salesman problem to find the shortest path connecting all the points, for
30 which many algorithms have been developed³⁰.

31

32 However, often data points are arranged in a more complex way than one simple
33 trajectory as in the case of lineage bifurcations during differentiation³¹. General

1 approaches to this problem start by building a graph that connects data points only to
2 their nearest neighbours (with the aim to preserve only the local distances) and then
3 finding the minimum spanning tree (MST) that connects all the data points (Fig4B). In
4 the simplest case, the diameter of this graph represents the dynamic trajectory along
5 which data points can then be sorted³²(Fig4C). This strategy has been successfully
6 applied to uncover trajectories and bifurcations both for low and medium dimensional
7 data such as flow cytometry coupled with mass spectrometry (cyto-mass) or single cell
8 quantitative real time PCR (qRT-PCR) expression measurements^{33, 34} and for high
9 dimensional gene expression data after applying a linear dimensionality reduction step
10 such as PCA^{30, 32}, ICA³⁵ or a clustering step³⁶.

11
12 Alternatively, several methods have been developed that start by building a nearest
13 neighbour graph and use the shortest path (geodesic) distance between points instead
14 of the euclidean distance to perform non-linear dimensionality reduction and
15 clustering. Examples include Isomap³⁷, locally linear embedding (LLE)³⁸ and laplacian
16 eigenmaps³⁹. Other non-linear dimensionality reduction methods such as diffusion
17 maps⁴⁰ or the t-distributed stochastic linear embedding (t-SNE)⁴¹ are based on
18 alternatives to classical distance metrics but again with the same objective of preserving
19 local similarities rather than global ones. Isomaps³¹, diffusion maps⁴² and t-SNE⁴³ have
20 been used in a biological context to discover trajectories, bifurcations and cell
21 heterogeneity in medium or high dimensional data on differentiation, development and
22 disease.

23

24

25

26 3.Applications

27

28 Interpreting functional genomics data

29 One of the first applications of expression deconvolution was in the interpretation of
30 systematic functional genomic data such as analysing the consequences of gene
31 deletion. In an early study in yeast, Lu et al showed that it is possible to deconvolve the
32 fraction of cells in each cell cycle phase from bulk microarray expression data in

1 asynchronous populations by using reference genes that oscillate during the cell cycle²⁸.
 2 This deconvolution made it possible to evaluate the effects of various environmental
 3 and genetic perturbations on the cell cycle dynamics from bulk gene expression,
 4 characterizing both the specific phase of the cell cycle and the severity of defects (Fig5).
 5 For example, based on the changes in gene expression the authors inferred that about
 6 half out of the 300 tested gene deletions affect cell cycle progression.

7

8 More recently O'Duibhir et al performed a similar analysis on 1485 gene expression
 9 profiles of yeast gene deletion strains¹⁰. They first found that 25% of the 700 mutants
 10 that differ from wild type share a common expression signature that is very similar to a
 11 'slow growth' signature induced by nutrient limitation⁹ or environmental stress. In
 12 yeast, the growth rate, the stress response and metabolic activity are tightly
 13 coordinated with the cell cycle⁹, and the authors argued that a change in the fraction of
 14 cells at different cell cycle stages in a population can account for expression changes in
 15 many different experiments¹⁰.

16

17 **Analysis of complex tissues**

18 Gene expression deconvolution is also useful for interpreting physiological changes in
 19 complex samples such as tissues. Tissues are a mixture of cell types so differential
 20 expression can be driven by: (1) changes in the relative abundance of different types,
 21 (2) changes that occur only in a subset of cell types, (3) changes common to every cell
 22 type, or (4) a combination of these three. Expression deconvolution can help
 23 discriminate among these scenarios^{8, 29, 44-46}. Deconvolution of cell type fractions based
 24 on reference expression datasets showed that systemic lupus SLE patients have a
 25 specific increase in activated natural killer and T helper lymphocytes⁸. In contrast,
 26 deconvolution of cell-specific gene expression from whole blood samples in
 27 combination with cell-type frequency revealed that kidney transplant recipients
 28 experiencing rejection had hundreds of differentially expressed genes specifically in
 29 monocytes²⁹.

30

31 **Discovery of new cell types and states.**

32 Single cell RNA sequencing (RNA-seq) technologies⁴⁷⁻⁵⁰ are opening up new possibilities
 33 for the analysis of complex heterogeneous samples. Whole tissues can be dissociated

1 into single cells that can be separately profiled^{48, 51}. Whole genome single cell profiles
2 are inherently stochastic which makes the analysis of biological variance more
3 challenging⁵², nonetheless they dramatically improve the ability to discover and
4 characterize cell types and states. Cell types and states are classically defined by
5 measuring a combination of a few selected markers using flow cytometry. Even in the
6 most advanced configuration (flow cytometry coupled with mass spectrometry)⁵³, this
7 approach allows about 40 markers to be measured in parallel for each cell, thus
8 introducing selection bias. In contrast, RNA-seq allows the expression of all genes to be
9 profiled across hundreds or thousands of cells, so providing an unbiased *ab initio*
10 characterization of cell states and types including rare ones^{48, 51}. Jaitin et al., for
11 example, could decompose by hierarchical clustering the heterogeneous dendritic cell
12 group into four functionally distinct subclasses and showed how the relative abundance
13 of these cell types is remodelled after infection⁴⁸. In another recent study, Zeisel et al.
14 used single cell RNA-seq to identify 47 subclasses of cells in the mouse cortex and
15 hippocampus⁵¹.

17 **Development and differentiation**

18 Gene expression deconvolution is also useful for interpreting dynamical biological
19 processes from simple responses to stimuli⁵⁴ to more complex dynamics such as
20 development and differentiation in multicellular organisms. Cell state trajectories can
21 be reconstructed and data points can be ordered along these with little or no *a priori*
22 chronological information both from average^{5, 30, 31, 36} and single cell expression data^{35,}
23 ^{42, 54, 55}. Inferring the precise physiological time point of each sample from gene
24 expression can also be important in experiments where the exact chronological time
25 point at which each sample was collected was controlled and recorded, for example
26 because of heterogeneity in the rates of development between genotypes, individual
27 cells or experimental batches⁵.

29 For example, Shalek et al showed by using PCA that the response of single dendritic
30 cells to a pathogenic stimulus is variable in time and includes some precocious cells at
31 early time points that are more advanced in the dynamic response⁵⁴ and more similar to
32 cells at later time points. Studying early blood development in mouse embryos using
33 PCA and diffusion maps Moignard et al. revealed heterogeneity along the differentiation

1 dynamics not only within embryos collected at the same chronological time but also
2 between single cells within individual embryos⁴². These two examples highlight two
3 important advantages of using single cell data and data from single individuals to
4 analyse dynamical progressions: first, one can quantify how synchronous a process is;
5 and second, while in average data the differences in physiological time within each time
6 point are averaged out and decrease the signal, analysis of single cells or individuals
7 ordered along a time-series allows analysis of the full dynamic response⁴² even when
8 the system is asynchronous. Another example is provided by the analysis of blood
9 development⁵⁵, which provided more power to discover cascades of causal regulators of
10 differentiation^{35, 55}.

11

12 **Reconstructing 3D spatial gene expression**

13 In multicellular organisms gene expression not only varies in time but also in space.
14 Methods that retain full spatial information of genome-wide gene expression exist but
15 they are still limited in throughput and are laborious⁵⁶. Junker et al proposed a method
16 similar to tomography where a sample is cryo-sectioned in different directions, each
17 section is analysed by RNA-seq and spatial expression is mathematically reconstructed.
18 Applying this method, they constructed an atlas of 3D expression patterns for zebrafish
19 embryos. Although impressive, ambiguities remain in the atlas when genes are
20 expressed in more than one contiguous region because the system is
21 underdetermined⁵⁷. To overcome this fundamental limitation, RNA-seq would have to
22 be performed on slices at different angles across the sample, which would likely require
23 averaging across different embryos. An alternative approach for reconstructing 3D
24 expression patterns is to use the known spatial distributions of landmark genes, for
25 example mapped by *in situ* hybridisation. This idea has been applied to reconstruct
26 spatial gene expression in zebrafish embryos⁵⁸ and in the brains of annelid worms⁵⁹. In
27 these studies, samples were first dissociated into single cells that were RNA-seq profiled
28 and spatial gene expression was computationally reconstructed by measuring the
29 similarity of the expression of each gene to the marker genes with known spatial
30 expression patterns. A similar approach was also used for the spatial reconstruction of
31 a much smaller number of gene expression profiles in the mouse otocyst, the precursor
32 of the inner ear⁶⁰.

33

Interpreting the effects of genetic variation

A central goal of many fields of biology such as human genetics and plant or animal breeding is to understand how natural genetic variation amongst individuals alters their characteristics. Here too the decomposition of expression profiles can be useful, either to remove non-genetic variation and improve the power when asking how genetic variation influences gene expression or to infer additional phenotypic traits or hidden environmental perturbations from gene expression components.

Expression quantitative trait loci (eQTLs) are genetic variants that alter gene expression. eQTLs are identified by performing genome-wide expression profiling on genetically heterogeneous populations. Gene expression is influenced by many non-genetic factors that can obscure subtle genetic effects. These non-genetic factors can be known covariates such as sex or age but they are often hidden uncontrolled experimental variables. eQTL studies are particularly sensitive to these because they usually include many experimental batches. Controlling for both known and hidden confounders greatly increase the power to detect significant eQTLs⁷. The simplest and most widely used approach to improve *cis*-eQTL detection is to remove the first few principal components, because this only removes broad variance components preserving local genetic effects. This approach however is not well suited for improving the detection of genetic loci that cause large-scale gene expression changes in *trans* because their signal might be removed together with confounders.

Beyond increasing statistical power, the analysis of hidden confounders can be useful for discovering genotype-environment interactions, i.e. genetic variants that change gene expression differently under different conditions, for example cell type-⁶¹, tissue-⁶² or environmental-specific eQTLs⁶³. If a hidden confounder reflects a biological source of variance, it should be treated as a covariate instead of being corrected for when testing the genetic effects on gene expression. Parts et al., for example, used a sparse factor analysis model to infer different cellular states (defined by the activity of different molecular pathways) and showed that the effect of some genetic variants is highly dependent on the cellular state⁶.

1 As another example, Curtis et al. used an integrative clustering approach to discover
2 new molecular cancer subtypes from gene expression heterogeneity and to characterize
3 the impact of genetic variation on these cancer subtypes. They further showed that
4 integrative clustering combining gene expression and genomic information is predictive
5 of survival⁶⁴.

6
7 If a dynamic biological process underlies the hidden source of variance in an eQTL
8 dataset, then the data can be used to investigate how time-dependent processes such as
9 development are influenced by genetic variation. For example, we recently investigated
10 how natural genetic variation affects gene expression in both space and time during 12
11 hours of the development of *C. elegans* by identifying and exploiting small differences in
12 the exact physiological stages at which each sample was expression profiled ⁶⁵. The
13 physiological stage of each sample was inferred by comparing the expression to a
14 reference gene expression time series using CCA (Fig6A), and the tissue-specificity of
15 expression trends were inferred by comparison to expression profiles of sorted tissue
16 samples. This allowed us to examine how sequence variation in the genome alters how
17 genes are expressed in time (Fig6B), and also to ask whether these effects are tissue-
18 specific or not. In this way we were able to identify hundreds of examples where
19 genetic variation close to a gene increased the amplitude of oscillations, altered the rate
20 of induction, or completely altered the dynamics⁵.

22 Summary

23 We have highlighted in this review how both bulk and single cell gene expression data
24 can be decomposed into the constituent cell types and states and used to reconstruct
25 spatial and temporal patterns of expression. This allows biological processes to be
26 studied at multiple levels using a single expression dataset (that may actually have been
27 generated for a different purpose).

28
29 One important lesson to be learned is that ‘hidden’ confounders in expression data can
30 be more than artefacts to correct for. Rather, they can identify important biological
31 sources of variance that can be interpreted and used to make interesting discoveries.

32

1 Another important lesson is that, although it is tempting to apply a simplifying discrete
2 view of a process, this often results in a loss of information and power because of the
3 intrinsically continuous nature of many biological processes in both time and space.

4
5 As the cost of single cell RNA sequencing falls and the methods for sample preparation
6 become more routine, these kinds of analyses will become increasingly important and
7 widely used. In particular, we envisage that the application of single cell RNA
8 sequencing to complex samples such as human tissue samples will facilitate the analysis
9 of how genetic variation influences many different dynamic biological processes such as
10 disease progression, development and tissue composition.

11 **Acknowledgments**

12 Research in our laboratory is supported by a European Research Council Consolidator
13 grant IR-DC (616434), by the Spanish Ministry of Economy and Competitiveness
14 (BFU2011-26206 and 'Centro de Excelencia Severo Ochoa 2013-2017' SEV-2012-0208),
15 the AXA Research Fund, FP7 project 4DCellFate (277899), Agència de Gestió d'Ajuts
16 Universitaris i de Recerca (AGAUR), and the EMBL-CRG Systems Biology Program.

1
2

3
4
5
6
7
8

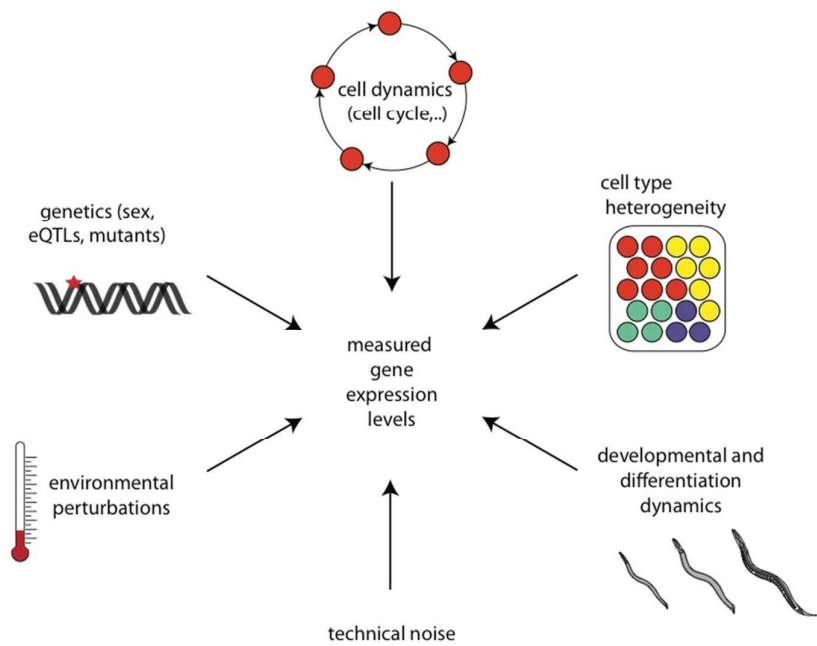


Fig 1 Sources of variance in gene expression data. Genome wide gene expression profiles can capture diverse intentional and unintentional influences.

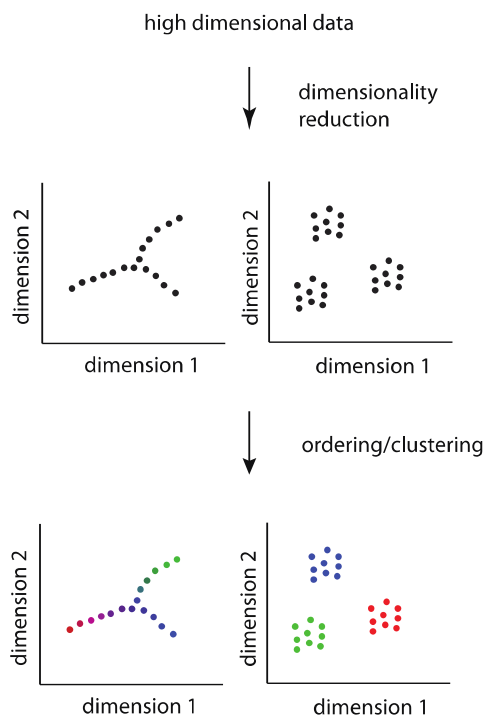


Fig 2. How to interpret high dimensional gene expression data. Projecting data in a low dimensional space (dimensionality reduction) is useful to filter out noise and helps to visualize the global structure of the data, for example highlighting clusters or trajectories. In the ideal case each dimension represents a distinct and interpretable biological process or function. Then clustering or methods to sort data points in trajectories can be applied.

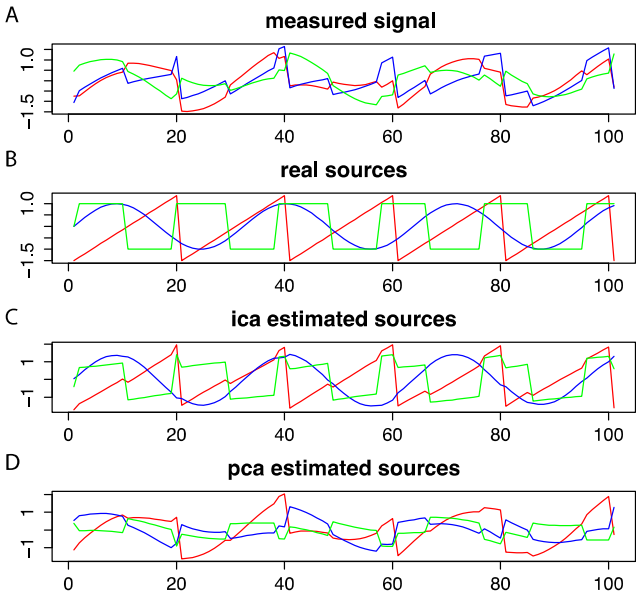


Fig 3. Finding the sources of variance in gene expression. Gene expression measures often include many mixed sources of signal. (A) Three measured signals constituted by a linear combination of (B) three original sources of signal along a time (or space) dimension. (C) ICA better estimates the original sources of signal than (D) PCA.

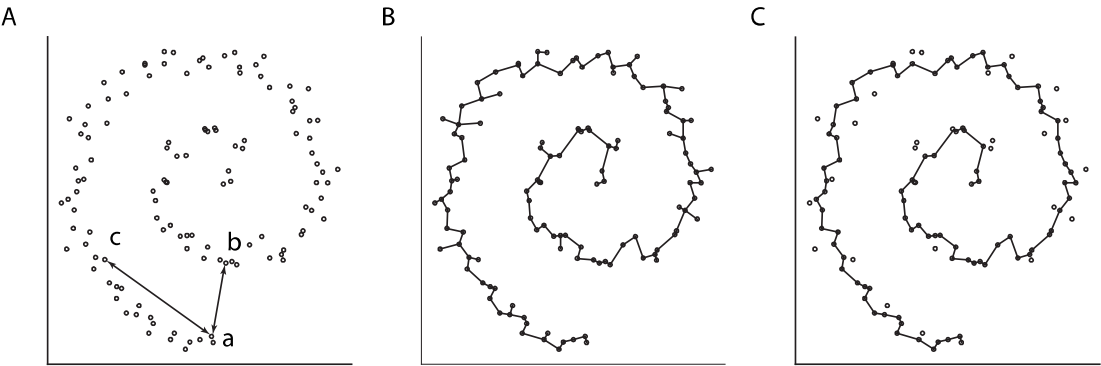


Fig 4. Identification of a non linear trajectory and the inference of the correct order of data points along this trajectory. A Synthetic dataset (jelly-roll) of points arranged in a spiral. Points a and b are closer than a and c in euclidean space despite being at opposite end of the spiral. B. The minimum spanning tree of the data – here

data points are connected according to their shortest path distance. C Data points are correctly sorted along a trajectory determined by the diameter of the minimum spanning tree (adapted from Magwene et al.³²).

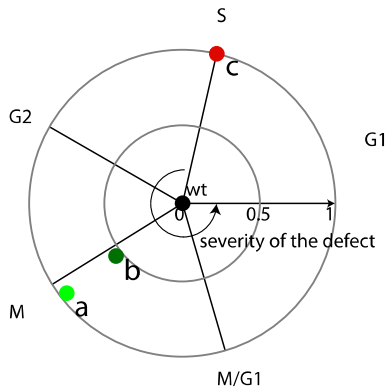


Fig 5. Interpreting the effect of genetic and environmental perturbations on cell cycle dynamics. The angular position indicates the phase of the cell cycle affected, while the distance from the centre (occupied by wild type) indicates the severity of the defect. Both mutants a and b affect mitosis but mutant a has a more severe effect. Mutant c affects S phase (adapted from Lu et al.²⁸).

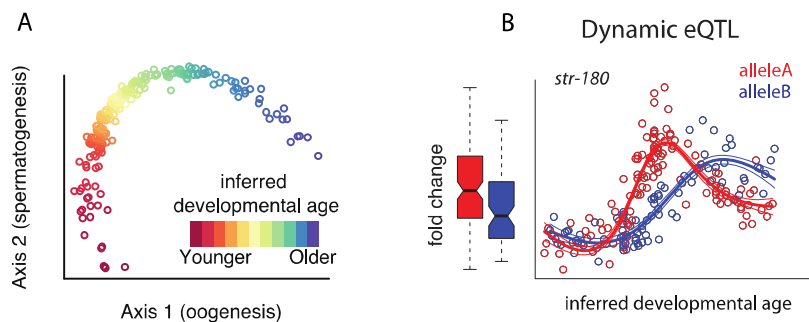


Fig 6. The impact of genetic variation on developmental dynamics. Worms switch from spermatogenesis to oogenesis when maturing. (A) the developmental age of each sample is inferred from the trajectory on the components related to spermatogenesis and oogenesis. (B) A dynamic eQTL analysis shows the complex effect of local genetic variation on the expression dynamics of the *str-180* gene (right) that could not be appreciated when developmental time is not included as a covariate in the analysis

(left). Analysis from Francesconi and Lehner⁵ of data from Rockman et al.⁶⁵.

References

1. P. Brennecke, S. Anders, J. K. Kim, A. A. Kolodziejczyk, X. Zhang, V. Proserpio, B. Baying, V. Benes, S. A. Teichmann, J. C. Marioni and M. G. Heisler, *Nat Methods*, 2013, **10**, 1093-1095.
2. N. Fusi, O. Stegle and N. D. Lawrence, *PLoS Comput Biol*, 2012, **8**, e1002330.
3. F. Buettner, K. N. Natarajan, F. P. Casale, V. Proserpio, A. Scialdone, F. J. Theis, S. A. Teichmann, J. C. Marioni and O. Stegle, *Nat Biotechnol*, 2015, **33**, 155-160.
4. N. Fusi, C. Lippert, K. Borgwardt, N. D. Lawrence and O. Stegle, *Bioinformatics*, 2013, **29**, 1382-1389.
5. M. Francesconi and B. Lehner, *Nature*, 2014, **505**, 208-211.
6. L. Parts, O. Stegle, J. Winn and R. Durbin, *PLoS Genet*, 2011, **7**, e1001276.
7. O. Stegle, L. Parts, R. Durbin and J. Winn, *PLoS Comput Biol*, 2010, **6**, e1000770.
8. A. R. Abbas, K. Wolslegel, D. Seshasayee, Z. Modrusan and H. F. Clark, *PLoS One*, 2009, **4**, e6098.
9. M. J. Brauer, C. Huttenhower, E. M. Airoidi, R. Rosenstein, J. C. Matese, D. Gresham, V. M. Boer, O. G. Troyanskaya and D. Botstein, *Mol Biol Cell*, 2008, **19**, 352-367.
10. E. O'Duibhir, P. Lijnzaad, J. J. Benschop, T. L. Lenstra, D. van Leenen, M. J. Groot Koerkamp, T. Margaritis, M. O. Brok, P. Kemmeren and F. C. Holstege, *Mol Syst Biol*, 2014, **10**, 732.
11. L. B. Snoek, M. G. Sterken, R. J. Volkers, M. Klatter, K. J. Bosman, R. P. Bevers, J. A. Riksen, G. Smant, A. R. Cossins and J. E. Kammenga, *Scientific reports*, 2014, **4**, 3912.
12. T. Hashimshony, M. Feder, M. Levin, B. K. Hall and I. Yanai, *Nature*, 2015, **519**, 219-222.
13. I. T. Jolliffe, *Principal Component Analysis*, Springer, Second edition edn., 2002.
14. O. Alter, P. O. Brown and D. Botstein, *Proc Natl Acad Sci U S A*, 2000, **97**, 10101-10106.
15. N. S. Holter, M. Mitra, A. Maritan, M. Cieplak, J. R. Banavar and N. V. Fedoroff, *Proc Natl Acad Sci U S A*, 2000, **97**, 8409-8414.
16. A. Subramanian, P. Tamayo, V. K. Mootha, S. Mukherjee, B. L. Ebert, M. A. Gillette, A. Paulovich, S. L. Pomeroy, T. R. Golub and E. S. Lander, *Proceedings of the National Academy of Sciences of the United States of America*, 2005, **102**, 15545-15550.
17. J. J. Goeman, S. A. Van De Geer, F. De Kort and H. C. Van Houwelingen, *Bioinformatics*, 2004, **20**, 93-99.
18. N. C. Chung and J. D. Storey, *Bioinformatics*, 2015, **31**, 545-554.
19. H. F. Kaiser, *Psychometrika*, 1958, **23**, 187-200.
20. A. E. Hendrickson and P. O. White, *British Journal of Statistical Psychology*, 1964, **17**, 65-70.
21. J. T. Leek and J. D. Storey, *PLoS Genet*, 2007, **3**, 1724-1735.
22. O. Stegle, L. Parts, M. Piipari, J. Winn and R. Durbin, *Nat Protoc*, 2012, **7**, 500-507.

- 1 23. J. A. Gagnon-Bartsch and T. P. Speed, *Biostatistics*, 2012, **13**, 539-552.
- 2 24. A. Hyvarinen and E. Oja, *Neural networks : the official journal of the International*
- 3 *Neural Network Society*, 2000, **13**, 411-430.
- 4 25. W. Liebermeister, *Bioinformatics*, 2002, **18**, 51-60.
- 5 26. S. I. Lee and S. Batzoglou, *Genome Biol*, 2003, **4**, R76.
- 6 27. E. M. Airoidi, C. Huttenhower, D. Gresham, C. Lu, A. A. Caudy, M. J. Dunham, J. R.
- 7 Broach, D. Botstein and O. G. Troyanskaya, *PLoS Comput Biol*, 2009, **5**, e1000257.
- 8 28. P. Lu, A. Nakorchevskiy and E. M. Marcotte, *Proc Natl Acad Sci U S A*, 2003, **100**,
- 9 10370-10375.
- 10 29. S. S. Shen-Orr, R. Tibshirani, P. Khatrri, D. L. Bodian, F. Staedtler, N. M. Perry, T.
- 11 Hastie, M. M. Sarwal, M. M. Davis and A. J. Butte, *Nat Methods*, 2010, **7**, 287-289.
- 12 30. L. Anavy, M. Levin, S. Khair, N. Nakanishi, S. L. Fernandez-Valverde, B. M. Degnan
- 13 and I. Yanai, *Development*, 2014, **141**, 1161-1166.
- 14 31. S. Huang, Y. P. Guo, G. May and T. Enver, *Dev Biol*, 2007, **305**, 695-713.
- 15 32. P. M. Magwene, P. Lizardi and J. Kim, *Bioinformatics*, 2003, **19**, 842-850.
- 16 33. P. Qiu, E. F. Simonds, S. C. Bendall, K. D. Gibbs, Jr., R. V. Bruggner, M. D. Linderman,
- 17 K. Sachs, G. P. Nolan and S. K. Plevritis, *Nat Biotechnol*, 2011, **29**, 886-891.
- 18 34. S. C. Bendall, K. L. Davis, A. D. Amir el, M. D. Tadmor, E. F. Simonds, T. J. Chen, D. K.
- 19 Shenfeld, G. P. Nolan and D. Pe'er, *Cell*, 2014, **157**, 714-725.
- 20 35. C. Trapnell, D. Cacchiarelli, J. Grimsby, P. Pokharel, S. Li, M. Morse, N. J. Lennon, K. J.
- 21 Livak, T. S. Mikkelsen and J. L. Rinn, *Nat Biotechnol*, 2014, **32**, 381-386.
- 22 36. P. Qiu, A. J. Gentles and S. K. Plevritis, *PLoS Comput Biol*, 2011, **7**, e1001123.
- 23 37. J. B. Tenenbaum, V. de Silva and J. C. Langford, *Science*, 2000, **290**, 2319-2323.
- 24 38. S. T. Roweis and L. K. Saul, *Science*, 2000, **290**, 2323-2326.
- 25 39. M. Belkin and P. Niyogi, *Neural computation*, 2003, **15**, 1373-1396.
- 26 40. R. R. Coifman, S. Lafon, A. B. Lee, M. Maggioni, B. Nadler, F. Warner and S. W.
- 27 Zucker, *Proc Natl Acad Sci U S A*, 2005, **102**, 7426-7431.
- 28 41. L. Van der Maaten and G. Hinton, *Journal of Machine Learning Research*, 2008, **9**,
- 29 85.
- 30 42. V. Moignard, S. Woodhouse, L. Haghverdi, A. J. Lilly, Y. Tanaka, A. C. Wilkinson, F.
- 31 Buettner, I. C. Macaulay, W. Jawaaid, E. Diamanti, S. Nishikawa, N. Piterman, V.
- 32 Kouskoff, F. J. Theis, J. Fisher and B. Gottgens, *Nat Biotechnol*, 2015, **33**, 269-276.
- 33 43. A. D. Amir el, K. L. Davis, M. D. Tadmor, E. F. Simonds, J. H. Levine, S. C. Bendall, D.
- 34 K. Shenfeld, S. Krishnaswamy, G. P. Nolan and D. Pe'er, *Nat Biotechnol*, 2013, **31**,
- 35 545-552.
- 36 44. H. Lahdesmaki, L. Shmulevich, V. Dunmire, O. Yli-Harja and W. Zhang, *BMC*
- 37 *Bioinformatics*, 2005, **6**, 54.
- 38 45. J. Clarke, P. Seo and B. Clarke, *Bioinformatics*, 2010, **26**, 1043-1049.
- 39 46. J. Ahn, Y. Yuan, G. Parmigiani, M. B. Suraokar, L. Diao, Wistuba, II and W. Wang,
- 40 *Bioinformatics*, 2013, **29**, 1865-1871.
- 41 47. S. Islam, A. Zeisel, S. Joost, G. La Manno, P. Zajac, M. Kasper, P. Lonnerberg and S.
- 42 Linnarsson, *Nat Methods*, 2014, **11**, 163-166.
- 43 48. D. A. Jaitin, E. Kenigsberg, H. Keren-Shaul, N. Elefant, F. Paul, I. Zaretsky, A. Mildner,
- 44 N. Cohen, S. Jung, A. Tanay and I. Amit, *Science*, 2014, **343**, 776-779.
- 45 49. S. Picelli, A. K. Bjorklund, O. R. Faridani, S. Sagasser, G. Winberg and R. Sandberg,
- 46 *Nat Methods*, 2013, **10**, 1096-1098.
- 47 50. T. Hashimshony, F. Wagner, N. Sher and I. Yanai, *Cell Rep*, 2012, **2**, 666-673.
- 48 51. A. Zeisel, A. B. M. Machado, S. Codeluppi, P. Lönnerberg, G. L. Manno, A. Juréus, S.
- 49 Marques, H. Munguba, L. He, C. Betsholtz, C. Rolny, G. Castelo-Branco, J. Hjerling-

1 Leffler and S. Linnarsson, *Science*, 2015.

2 52. O. Stegle, S. A. Teichmann and J. C. Marioni, *Nat Rev Genet*, 2015, **16**, 133-145.

3 53. S. C. Bendall, E. F. Simonds, P. Qiu, E. ad D Amir, P. O. Krutzik, R. Finck, R. V.

4 Bruggner, R. Melamed, A. Trejo, O. I. Ornatsky, R. S. Balderas, S. K. Plevritis, K.

5 Sachs, D. Pe'er, S. D. Tanner and G. P. Nolan, *Science*, 2011, **332**, 687-696.

6 54. A. K. Shalek, R. Satija, J. Shuga, J. J. Trombetta, D. Gennert, D. Lu, P. Chen, R. S.

7 Gertner, J. T. Gaublonne, N. Yosef, S. Schwartz, B. Fowler, S. Weaver, J. Wang, X.

8 Wang, R. Ding, R. Raychowdhury, N. Friedman, N. Hacohen, H. Park, A. P. May and

9 A. Regev, *Nature*, 2014, **510**, 363-369.

10 55. S. C. Bendall, K. L. Davis, E.-A. D. Amir, M. D. Tadmor, E. F. Simonds, T. J. Chen, D. K.

11 Shenfeld, G. P. Nolan and D. Pe'er, *Cell*, 2014, **157**, 714-725.

12 56. J. H. Lee, E. R. Daugharthy, J. Scheiman, R. Kalhor, J. L. Yang, T. C. Ferrante, R. Terry,

13 S. S. Jeanty, C. Li, R. Amamoto, D. T. Peters, B. M. Turczyk, A. H. Marblestone, S. A.

14 Inverso, A. Bernard, P. Mali, X. Rios, J. Aach and G. M. Church, *Science*, 2014, **343**,

15 1360-1363.

16 57. J. P. Junker, E. S. Noel, V. Guryev, K. A. Peterson, G. Shah, J. Huisken, A. P. McMahon,

17 E. Berezhikov, J. Bakkers and A. van Oudenaarden, *Cell*, 2014, **159**, 662-675.

18 58. R. Satija, J. A. Farrell, D. Gennert, A. F. Schier and A. Regev, *Nat Biotechnol*, 2015,

19 DOI: 10.1038/nbt.3192.

20 59. K. Achim, J. B. Pettit, L. R. Saraiva, D. Gavriouchkina, T. Larsson, D. Arendt and J. C.

21 Marioni, *Nat Biotechnol*, 2015, DOI: 10.1038/nbt.3209.

22 60. R. Durruthy-Durruthy, A. Gottlieb, B. H. Hartman, J. Waldhaus, R. D. Laske, R.

23 Altman and S. Heller, *Cell*, 2014, **157**, 964-978.

24 61. M. Ackermann, W. Sikora-Wohlfeld and A. Beyer, *PLoS Genet*, 2013, **9**, e1003514.

25 62. T. Flutre, X. Wen, J. Pritchard and M. Stephens, *PLoS Genet*, 2013, **9**, e1003486.

26 63. J. Gagneur, O. Stegle, C. Zhu, P. Jakob, M. M. Tekkedil, R. S. Aiyar, A.-K. Schuon, D.

27 Pe'er and L. M. Steinmetz, *PLoS Genet*, 2013, **9**, e1003803.

28 64. C. Curtis, S. P. Shah, S. F. Chin, G. Turashvili, O. M. Rueda, M. J. Dunning, D. Speed, A.

29 G. Lynch, S. Samarajiwa, Y. Yuan, S. Graf, G. Ha, G. Haffari, A. Bashashati, R.

30 Russell, S. McKinney, A. Langerod, A. Green, E. Provenzano, G. Wishart, S. Pinder,

31 P. Watson, F. Markowetz, L. Murphy, I. Ellis, A. Purushotham, A. L. Borresen-Dale,

32 J. D. Brenton, S. Tavaré, C. Caldas and S. Aparicio, *Nature*, 2012, **486**, 346-352.

33 65. M. V. Rockman, S. S. Skrovanek and L. Kruglyak, *Science*, 2010, **330**, 372-376.