



Published in final edited form as:

*Nat Methods*. 2010 April ; 7(4): 311–317. doi:10.1038/nmeth.1442.

## Identifying single-cell molecular programs by stochastic profiling

Kevin A. Janes<sup>1,2,\*</sup>, Chun-Chao Wang<sup>2</sup>, Karin J. Holmberg<sup>2</sup>, Kristin Cabral<sup>3</sup>, and Joan S. Brugge<sup>1,\*</sup>

<sup>1</sup> Department of Cell Biology, Harvard Medical School, Boston, MA 02115, USA

<sup>2</sup> Department of Biomedical Engineering, University of Virginia, Charlottesville, VA 22908, USA

<sup>3</sup> Molecular Genetics Core Facility, Children's Hospital Boston and Harvard Medical School, Boston, MA 02115, USA

### Abstract

Cells within tissues can be morphologically indistinguishable yet show molecular expression patterns that are remarkably heterogeneous. Here, we describe an approach for comprehensively identifying coregulated, heterogeneously expressed genes among cells that otherwise appear identical. The technique, called “stochastic profiling”, involves the repeated, random selection of very-small cell populations via laser-capture microdissection, followed by a customized single-cell amplification procedure and transcriptional profiling. Fluctuations in the resulting gene-expression measurements are then analyzed statistically to identify transcripts that are heterogeneously co-expressed. We stochastically profiled matrix-attached human epithelial cells in a three-dimensional culture model of mammary-acinar morphogenesis. Of 4,557 transcripts, we identified 547 genes with strong cell-to-cell expression differences. Clustering of this heterogeneous subset revealed several molecular “programs” implicated in protein biosynthesis, oxidative-stress responses, and nuclear factor- $\kappa$ B signaling, which were independently confirmed by RNA fluorescence *in situ* hybridization. Thus, stochastic profiling can reveal single-cell heterogeneities without measuring individual cells explicitly.

### INTRODUCTION

Cell-to-cell variations in gene and protein expression play an important role in the development and function of many tissues<sup>1, 2</sup>. Fluctuations at the single-cell level can be

Users may view, print, copy, download and text and data- mine the content in such documents, for the purposes of academic research, subject always to the full Conditions of use: [http://www.nature.com/authors/editorial\\_policies/license.html#terms](http://www.nature.com/authors/editorial_policies/license.html#terms)

\* Correspondence: [kjanes@virginia.edu](mailto:kjanes@virginia.edu) and [joan\\_brugge@hms.harvard.edu](mailto:joan_brugge@hms.harvard.edu).

#### AUTHOR CONTRIBUTIONS

K.A.J. conceived of the study, performed the computational simulations, designed and optimized the experimental protocols, performed the stochastic profiling experiments, analyzed the data, and wrote the initial draft of the manuscript. C.C.W. validated the RNA FISH riboprobes, performed the RNA FISH experiments, and edited the manuscript. K.J.H. cloned and prepared the RNA FISH riboprobes, segmented the images for quantitation, and edited the manuscript. K.C. optimized the microarray hybridization protocols and performed the microarray hybridization experiments. J.S.B. supervised the overall research progress and contributed to the initial draft of the manuscript.

#### COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

masked or completely misrepresented when analyzed at the population level<sup>3</sup>. This makes heterogeneities problematic for interpreting bulk measurements from large numbers of cells, such as from tumors or developing organs. Yet, such non-uniformities often uncover interesting molecular patterns that can reveal important mechanisms for the regulation of cell fate<sup>4, 5</sup>. Identifying heterogeneities is thus key for gaining a deeper understanding of tissue physiology.

The challenge in discovering heterogeneities is that cells of the same type may appear phenotypically indistinguishable. Heterogeneities at the molecular level can be uncovered by immunochemistry, but the markers must be selected *a priori* and analyzed in small groups. While more parameters can be screened simultaneously with flow cytometry<sup>3</sup>, this involves substantial tissue processing to isolate single cells from solid tissues. Extraction of individual cells is possible *in situ* using laser-capture microdissection<sup>6</sup>, but aside from large cells such as neurons and cardiomyocytes<sup>7, 8</sup>, there is usually not enough biological material to measure the expression of all but the most-abundant transcripts.

Last and most importantly, there is the conceptual hurdle of interpreting measurements from a single cell. Regulated cell-to-cell heterogeneities will appear as fluctuations in one-cell measurements. However, fluctuations will also be observed because of random biological variation, which may be functionally inconsequential<sup>9</sup>, and measurement error, which can be enormous<sup>10</sup>. The inability to separate contributions from these different sources has precluded using single-cell approaches to study the coordination of pathways that are heterogeneously activated.

We sought to address these challenges by developing an approach, called stochastic profiling, which is based on small-population averaging of randomly chosen cells. As a first application, we examined single-cell gene expression in a three-dimensional culture model of mammary acinar morphogenesis<sup>11</sup>. The sensitivity, precision, and quantitative accuracy of stochastic profiling make it an attractive technique for studying endogenous transcriptional heterogeneities in development and cancer.

## RESULTS

To reveal the dichotomous expression of a gene (“Gene B”), which is expressed at high levels in one population and at low levels in another (Fig. 1a), we repeatedly select very-small cell populations at random and measure the average gene expression from each random sampling (Step 1). Then, we construct a reference histogram from homogeneously expressed genes (“Gene A”), which estimates the sampling fluctuations when no dichotomy is present (Step 2). Last, we compare the estimated reference distribution to fluctuations of candidate genes measured from the same stochastic samplings (Step 3). The Gene-B distribution will deviate from the Gene-A reference because of differences in the proportion of subpopulations that were collected at each sampling (Fig. 1a). In addition, dichotomously expressed genes that are coregulated at the single-cell level (“Gene B” and “Gene C”) will have deviations that correlate across repeated samplings. Therefore, we can in principle reveal heterogeneous expression programs composed of multiple genes by clustering patterns of sampling fluctuations.

## Theoretical validation of stochastic sampling

We used computer simulations to help define the required sampling conditions and characterize the expression heterogeneities that stochastic sampling detects. Cells transcribe genes in exponential “bursts”<sup>12</sup>, which give rise to log-normal distributions of mRNA species in the population<sup>13</sup> (see below). Single-cell gene-expression levels were thus modeled as log-normal probability distributions with coefficients of variation (CVs) proportional to the log-standard deviation (Fig. 1b). Together, the model described the reference and dichotomous distributions with four parameters: the CV of the reference distribution ( $CV_a$ ), the CV of the distributions in the gene that is dichotomously expressed ( $CV_b$ ), the magnitude of the expression difference between the dichotomous subpopulations ( $D$ ), and the fraction of cells with high expression for the gene that is dichotomously expressed ( $F$ ) (Fig. 1b).

After selecting values for  $CV_a$ ,  $CV_b$ ,  $D$ , and  $F$ , we simulated the experiments and centered the sampling fluctuations of each gene on their respective log mean (Fig. 1a). Next, the sampling fluctuations of the dichotomously expressed gene were compared against a log-normal distribution using the log-standard deviation calculated from the reference distribution. The discrepancy between the log-normal reference and the sampling fluctuations of the dichotomously expressed gene was then assessed for statistical significance by a  $\chi^2$  goodness-of-fit test (see **Online Methods**).

As a control for the modeled stochastic samplings, we simulated a parallel set of control samplings, where all the parameters were the same but  $F$  was set to zero (i.e., no dichotomy). These control samplings identified false positives, which were scored as different from the reference simply because the model CVs were poorly matched ( $CV_a \ll CV_b$ ; Fig. 1c,d). When the reference and dichotomy CVs were poorly matched in the opposite direction ( $CV_a \gg CV_b$ ), there was the danger of false negatives, because a dichotomous sampling distribution could be misinterpreted as a log-normal distribution with a larger CV (Fig. 1c,e). Effective stochastic sampling occurred when the reference and dichotomy CVs were roughly comparable, so that significant deviations from the reference were observed only when  $F > 0$  ( $P < 0.05$ , Fig. 1c,f).

We first sought to determine the maximum number of cells that, when averaged, could confidently identify heterogeneities across a wide range of  $CV_b$ . Direct estimates of transcriptional noise are not available, but studies in yeast have found that protein levels can fluctuate with CVs ranging from ~12-38% (Ref. <sup>14</sup>). We independently varied  $CV_a$  and  $CV_b$  over this range for different numbers of cells sampled and then identified the CV combinations that gave false positives, false negatives, and effective stochastic sampling. When  $CV_a$  was very low (< 20%), we found that there was a substantial likelihood of false positives, which was independent of the number of cells sampled (Fig. 1g). Conversely, when  $CV_a$  was very high (> 30%), there was a danger of false negatives, which increased dramatically when more than 10 cells were sampled (Fig. 1g). With 10-cell samplings and an intermediate reference distribution ( $CV_b \sim 25\text{--}30\%$ ), we were able to achieve effective stochastic sampling across nearly all  $CV_b$  values (Fig. 1g). Furthermore, using these parameters, stochastic sampling could successfully identify dichotomies as small as 5–6 fold

(Fig. 1h), with relatively little dependence on the dichotomy fraction above ~5% (Fig. 1i). When  $F < 0.05$ , the dichotomy is too rare to detect reliably in 10-cell samplings, and we observed a sharp increase in false negatives (Supplementary Fig. 1). We conclude that stochastic sampling of up to 10 cells is sufficient to detect many dichotomies when given a reference for the “average” non-dichotomous sampling fluctuations.

### Optimization of small-cell PCR for stochastic profiling

Based on the simulation-derived estimates, we then developed a poly(A)-PCR amplification procedure for accurately profiling gene expression in 10 microdissected cells. Poly(A) PCR can amplify large quantities of polyadenylated transcripts from minute samples<sup>15</sup>. This technique has previously been modified to improve either single-cell representation of genes or detection sensitivity for low-abundance transcripts<sup>10, 16</sup>. To optimize the technique for stochastic sampling, we designed a “small-cell” poly(A) PCR that maximizes both the reproducibility between measurement replicates and the quantitative accuracy of genes measured from 10 cells (Supplementary Fig. 2 and **Online Methods**). Accuracy and precision were validated by serially diluting microdissected cells before the amplification and then quantifying high- to low-abundance genes post-amplification by real-time quantitative PCR (RT-qPCR). The dilutions are critical to ensure that quantitative differences in transcript levels are not artificially increased or decreased during the procedure. To date, this quantitative accuracy has only been shown when amplification is omitted entirely<sup>17</sup>, which substantially limits the number of transcripts that can analyzed from the same sample.

For a large panel of genes with varying abundances, we found that small-cell poly(A) PCR was highly accurate and reproducible for 3–100 cells (Fig. 2a–h and Supplementary Fig. 3). The median amplification efficiency ( $E$ ) across all genes measured was 99.5%, and for individual genes, the efficiency was comparable to that of the RT-qPCR primers themselves ( $E_p$ ). This suggested that the poly(A)-PCR procedure was not skewing changes in the abundance of individual genes. Overall 10-cell reproducibility as measured by RT-qPCR was 0.36 cycle thresholds ( $C_T$ ), which corresponds to an amplification precision of ~28% if  $E = 100\%$  ( $2^{0.36} - 1 = 28\%$ , Fig. 2i). Importantly, we found for many genes that the accuracy and precision of poly(A) PCR decreased substantially when single-cell equivalents of RNA were used (Fig. 2a–h and Supplementary Fig. 3). Several genes were not reproducibly detectable (e.g., Fig. 2e–h), whereas others deviated from the log-linear standard predicted from the 3–100-cell dilution series (e.g., Fig. 2a,c,d). These results were obtained from microdissected breast-epithelial cells with an average diameter of ~10  $\mu\text{m}$ . Therefore, many more cell types should be quantifiable using a small-cell (rather than single-cell) approach together with stochastic profiling.

### Adapting small-cell PCR to oligonucleotide microarrays

A key step toward accurate 10-cell quantification was limiting the number of amplification cycles in small-cell poly(A) PCR to no more than 30 (Supplementary Fig. 2). With 10 microdissected cells, a 30-cycle amplification typically yielded ~10 ng of unlabeled cDNA, which was insufficient for oligonucleotide microarrays. We therefore reamplified a fraction of the poly(A) cDNA and added aminoallyl-dUTP for subsequent fluorophore labeling,

yielding ~1.5 µg of labeled cDNA per 10-cell sample. The conditions for reamplification differed from small-cell poly(A) PCR (see **Online Methods**) and were carefully monitored with real-time pilot experiments to identify the maximum number of cycles that kept all samples in the exponential phase of amplification. We found that doing so maintained the quantitative accuracy toward high- and low-abundance transcripts (Supplementary Fig. 4). Furthermore, repeat reamplifications using the same starting cDNA pool confirmed that reamplification added little measurement error to the final microarray measurements (Supplementary Fig. 5). Hybridization of reamplified samples to Illumina HumanRef-8 microarrays consistently detected 7,000–8,000 transcripts (median detection  $P < 0.1$ ). This result compares favorably with an earlier study from our group<sup>18</sup>, in which ~8,700 transcripts were detected by standard profiling approaches using RNA extracted from large populations of the cells used here. We conclude that our experimental platform is sufficiently accurate and sensitive to quantify much of the transcriptome for stochastic profiling.

### Stochastic profiling of epithelial acinar morphogenesis

As a proof of principle, we tested the feasibility of stochastic profiling in a three-dimensional (3D) culture model of mammary-epithelial acinar morphogenesis<sup>11</sup>. For this culture model, individual MCF10A mammary-epithelial cells are seeded in reconstituted basement membrane and develop to form proliferation-arrested, hollow acinar structures comprised of 50–100 cells when fully mature. Each acinus is clonal and thus isogenic, but many signaling and cell-fate dichotomies nonetheless emerge during morphogenesis. For example, matrix-attached cells of the outer acinus appear grossly similar but show variable expression of phospho-Akt<sup>19</sup>, phospho-myosin light chain<sup>20</sup>, and the CDK inhibitor p27 (Ref. 21). The overall extent of such cell-to-cell heterogeneities and their role in morphogenesis has not been defined.

We focused the stochastic profiling on matrix-attached cells in developing 3D cultures, because these cells comprise the final acinar structure that resembles the lobular unit of the breast *in vivo*<sup>11</sup>. Matrix-attached cells are also readily identified in cryosections of 3D structures and can be microdissected as single cells with high accuracy (Supplementary Fig. 6 and **Online Methods**). We obtained transcriptional profiles for 16 independent 10-cell samplings of matrix-attached cells along with 16 measurement controls. The control samples consisted of independent amplifications from a common starting pool of 160 microdissected cells. These amplification replicates were used to gauge the measurement error associated with profiling gene expression from 10 “average” cells.

Our analysis focused on the 4,557 transcripts that were clearly identified in all 32 microarrays (16 samplings plus 16 controls,  $P < 0.1$ ). First, we identified the subset of transcripts whose sampling CV was significantly higher than the corresponding control CV in amplification replicates (see **Online Methods**). We reasoned that the independent samplings of such transcripts would provide a good estimate for normal biological variation with only a minor contribution from measurement noise. Many transcripts were eliminated from the subset because their independent 10-cell sampling measurements were highly reproducible. For example, 1,332 genes had a sampling CV  $< 20\%$ , meaning that the

corresponding control CV would have needed to be  $< \sim 10\%$  to be included in the analysis. We presume that the majority of these transcripts are homogeneously expressed or show heterogeneities too small or infrequent to be detected experimentally.

Next, we clustered the independent measurements of the 1,003 genes in the subset, using Euclidean distance as a metric to sort transcripts roughly by sampling CV (Fig. 3a). We observed a plateau of low and consistent sampling CVs, followed by an abrupt increase where sampling fluctuations seemed to become more irregular and less random. We defined the transcripts in the early plateau as the reference-gene set (Fig. 3a) and found that the median sampling CV in this set was 19% with an interquartile range of 14–26% (Supplementary Fig. 7a). We fed these empirically derived parameters into our earlier model and found that stochastic profiling should be effective across the entire interquartile range of CVs (Supplementary Fig. 7b). Last, we compared sampling fluctuations of individual transcripts against a log-normal reference distribution with  $CV_a = 0.19$  at a false-discovery rate of 0.05 (Supplementary Fig. 7c). Overall, stochastic profiling identified 547 genes whose expression was predicted to be strongly heterogeneous (12% of all transcripts consistently detected).

### Discovery of heterogeneous single-cell programs

We standardized and reclustered the sampling data for the candidate heterogeneities to organize genes by their pattern of sampling fluctuations (Fig. 3b). The analysis identified multiple clusters that had strong links to recognized biological processes. The first cluster contained many genes involved in protein synthesis, including ribosomal subunits (*RPS6*, *RPL38*, etc.), initiation-elongation factors (*EIF3M*, *EEF2*), and chaperones (*SEC61G*, *TBCA*). This cluster also contained the basal-progenitor markers, *KRT5* (Ref. <sup>22</sup>) and an *ALDH* isoform<sup>23</sup>, and the *JUND* transcription factor. The second cluster was comprised of several transcripts connected with oxidative-stress responses and proliferative suppression, such as *PRDX4*, *FAM120A*<sup>24</sup>, *SERP1* (Ref. <sup>25</sup>), and *FOXO1*<sup>26</sup>. The third cluster was the smallest but contained a large proportion of genes known to be initiators (*ILIR1*), effectors (*NFKBIA*), or markers (*BIRC3*, *SOD2*) of nuclear factor- $\kappa$ B (NF- $\kappa$ B) signaling<sup>27</sup>. NF- $\kappa$ B signaling heterogeneity was also observed posttranslationally by localization of the p65 subunit of NF- $\kappa$ B and expression of I $\kappa$ B $\alpha$ , an upstream inhibitor of NF- $\kappa$ B (Supplementary Fig. 8). Taken together, the correlated sampling fluctuations and shared biological function within clusters suggested these were molecular programs that were induced heterogeneously in single cells.

We next sought to validate the stochastic-profiling predictions by an independent method. We developed an RNA fluorescence *in situ* hybridization (FISH) procedure for dual tracking of gene-expression variation in individual cells (see **Online Methods**). Our two-color RNA FISH protocol was optimized for specificity (**Supplementary Figs. 9 and 10**) and for reliably detecting single-cell coregulatory patterns between selected transcripts (Fig. 4a-c). Using RNA FISH, we observed pronounced cell-to-cell expression heterogeneities for nearly all transcripts identified by stochastic profiling that were examined (Fig. 4a-c, **Supplementary Fig. 11**, and **Supplementary Note 1**). Conversely, we observed more-uniform expression for two genes, *GAPDH* and *HINT1*, whose stochastic-sampling

fluctuations were not different than the reference distribution (Supplementary Fig. 12). Thus, stochastic profiling can separate acute single-cell heterogeneities from transcripts with normal expression variability.

In addition, for gene pairs in the same cluster, we found highly concordant patterns of strong and weak expression among individual cells (Fig. 4a-c, **Supplementary Fig. 11**, and **Supplementary Note 2**). Analysis of cell-to-cell fluorescence intensities revealed that matrix-attached cells were almost exclusively “double negative” (weakly expressing both genes) or “double positive” (strongly expressing both genes) (Fig. 4d). Cells that strongly expressed one gene but not the other (“single positive”) were too rare to constitute a meaningful subpopulation and were likely filtered out by stochastic profiling (Fig. 4d and Supplementary Fig. 1). Together, this indicates that clusters of genes with similar stochastic-sampling fluctuations are heterogeneously coexpressed with high probability.

As a final validation, we checked whether genes in separate stochastic-profiling clusters were distinguishable on the single-cell level by RNA FISH. The observed concordance between clusters ranged from no discernable correlation (Fig. 5a-c) to pairs with stronger covariation (Fig. 5d-i). Nevertheless, for each gene pairing, we repeatedly identified single-positive cells at frequencies that should be detected by stochastic profiling (> 9–10%, Figs. 1i and 5b,e,h). Inclusion of these single-positive cells during stochastic sampling would be sufficient to perturb any correlated fluctuations, providing an explanation for the distinct clusters shown in Figure 3b. Indeed, using the RNA FISH measurements as the basis for simulated stochastic samplings, we estimated probability distributions that largely captured the stochastic-profiling measurements (Figs. 3b and 5c,f,i).

## DISCUSSION

Transcriptional heterogeneities can emerge from purely stochastic cell-fate decisions<sup>1, 2, 28</sup>, but they can also be instructed by differences in the microenvironment<sup>29</sup>. Stochastic profiling does not make a distinction between these heterogeneities but provides a means for identifying them so that the underlying mechanisms can be studied thereafter. The biggest advantage of stochastic profiling is its improved accuracy and reproducibility, which becomes possible when 10 cells are measured instead of one (Fig. 2). Although measurements are not explicitly single cell, the entire procedure requires only a few hundred cells, meaning that stochastic profiling should be amenable to most *ex vivo* tissue specimens in the future.

Our first application of stochastic profiling uncovered many genes not previously suspected to show heterogeneous regulation during morphogenesis. MCF10A cells have a basal-progenitor expression profile<sup>30</sup>, suggesting that some heterogeneities could be due to partial differentiation of single cells in 3D. The existence of a heterogeneous stress-response program is particularly intriguing, because it raises the possibility that individual cells might occupy stressful niches caused by local cell-cell interactions and basement-membrane composition.

Another interesting question raised by the study here is whether the single-cell programs identified by stochastic profiling are coordinated during morphogenesis. For the gene clusters imaged simultaneously by RNA FISH, we found that the single-positive populations were not equally populated. For example, high *JUND* expression could be found in cells with low *IL1R1* or *FOXO1* expression, but cells with the opposite pattern were extremely rare (Fig. 5e,h). Future work will focus in greater depth on these dependencies and their possible role during morphogenesis.

The extent to which heterogeneously activated pathways *in vivo* might obscure phenotypes or create patterns in tissues is only beginning to be studied<sup>1</sup>. The bottleneck is not in studying the role of heterogeneities, but rather in identifying them in the first place. Stochastic profiling provides a valuable tool for analyzing the coordination of such pathways quantitatively and systematically.

## ONLINE METHODS

### Monte Carlo simulations

Stochastic sampling simulations (Fig. 1) were performed in MATLAB (Mathworks) with the statistics toolbox. For each simulation, the model assumed a binomial distribution for the cellular dichotomy and log-normal distribution of measured transcripts<sup>13, 31</sup>.  $CV_a$  and  $CV_b$  were varied between 12–38% to approximate biologically plausible values<sup>14</sup> and then re-run with empirically derived values (Supplementary Fig. 7a,b). The distribution of 48 population-averaged samplings was log-mean centered and compared to a log-normal distribution with a standard deviation estimated from 48 reference samplings. The  $\chi^2$  goodness of fit between the dichotomous and reference distributions was done using the `chi2gof` function with 10 bins. The  $\chi^2$  test directly evaluates the relative differences between observed and expected values on the sampling histogram and is a robust, conservative test for this application<sup>32</sup>. Bins were pooled if the observed or expected value in a bin was less than five. Each  $CV_a$ ,  $CV_b$ ,  $D$ , and  $F$  parameter set was run 50 times to measure the median  $P$  values and the associated nonparametric confidence intervals. Stochastic sampling was deemed effective when the median  $P$  value for  $F = 0$  was less than 0.05 and the median  $P$  value for  $F = 0$  was greater than 0.05. The source code for the simulations is available in the Supplementary Software.

For the simulation of probability density functions (Fig. 5c,f,i), the single-cell FISH intensities from Fig. 5b,e,h were randomly combined as 10-cell averages for each gene pair for 5,000 iterations. These bootstrapped estimates were standardized and then compiled as two-dimensional histograms by using the `hist2` function with 20 bins.

### Cell lines

The MCF10A-5E clone was isolated by limiting dilution of the parental MCF10A line (ATCC) and selected for its homogeneous behavior in 3D. MCF10A-5E cells were maintained as described previously for MCF10A cells<sup>33</sup>.



### Frozen sectioning of 3D cultures

To allow embedding of 3D cultures, a plastic coverslip was cut to size and placed at the base of an 8-well chamber slide (BD Biosciences) before starting. Coverslipped chamber slides were then coated with Matrigel (BD Biosciences), and 3D culture of MCF10A-5E cells was performed as described previously<sup>33</sup>. For fresh frozen sections (used for laser capture microdissection), coverslips were washed with PBS and then embedded directly in NEG 50 on a dry ice-isopentane bath. For fixed frozen sections (used for RNA FISH), coverslips were washed in PBS and fixed in 3.7% paraformaldehyde for 15 min. After three 5 min washes in PBS, samples were cryopreserved in 15% sucrose for 15 min, 30% sucrose for 15 min, and then embedded in NEG 50 as described above. Sectioning was performed at  $-24^{\circ}\text{C}$  on a cryostat (Leica). Embedded specimens and cryosections were stored at  $-80^{\circ}\text{C}$  until further use.

### Laser capture microdissection

8  $\mu\text{m}$  sections were cut on plain glass slides and kept at  $-24^{\circ}\text{C}$  during sectioning and  $-80^{\circ}\text{C}$  during storage. After removing from  $-80^{\circ}\text{C}$ , slides were fixed immediately in 75% ethanol for 30 sec, followed by distilled water for 30 sec. Fixed slides were stained for 30 sec with nuclear fast red (Vector Laboratories) containing 1 U  $\text{ml}^{-1}$  RNAsin Plus (Promega), then washed twice in distilled water for 15 sec. Stained slides were dehydrated with an ethanol series (30 sec each of 70%, 95%, and 100% ethanol) and cleared with xylene for 2 min. After air drying for 5–10 min, slides were stored in a dessicator and used immediately.

Before microdissection, slides were cleaned with a PrepStrip (Arcturus) to remove loosely adherent material. Microdissection was performed on a Pixcell II instrument (Arcturus) using Capsure HS LCM caps (Arcturus). 750  $\mu\text{s}$  laser shots were used at 50–65 mW power to achieve single-cell resolution (Supplementary Fig. 6). For this study, matrix-attached cells were sampled at 3–4 random positions across  $\sim 3$  acini to focus on matrix-dependent (rather than acinus-dependent) heterogeneities. After microdissection, LCM caps were cleaned with an adhesive note to remove biological material adjacent to the dissected cells.

### Small-cell quantitative mRNA amplification

Samples were eluted from the microdissection caps by adding 4  $\mu\text{l}$  digestion buffer (1.25 $\times$  MMLV RT buffer [Invitrogen], 100  $\mu\text{M}$  dNTPs [Roche], 0.08 OD  $\text{ml}^{-1}$  oligo(dT)<sub>24</sub>, and 250  $\mu\text{g}$   $\text{ml}^{-1}$  proteinase K [Sigma]) and incubating at  $42^{\circ}\text{C}$  for 1 hr. Digested samples were spun into PCR tubes and quenched with 1  $\mu\text{l}$  of digestion stop buffer (1.5 U  $\text{ml}^{-1}$  Prime RNase inhibitor [Eppendorf], 1.5 U  $\text{ml}^{-1}$  RNAGuard [Amersham], and 5 mM freshly prepared PMSF). The quenched samples were then processed by using poly(A) PCR<sup>15</sup> that was heavily modified to allow quantitative amplification of high- and low-abundance transcripts.

4.5  $\mu\text{l}$  of the quenched samples were transferred into thin-walled 0.2-ml PCR tubes, and 0.5  $\mu\text{l}$  of Superscript III (Invitrogen) was added. The first-strand synthesis reaction was incubated at  $50^{\circ}\text{C}$  for 15 min and then heat-inactivated at  $70^{\circ}\text{C}$  for 15 min. The samples were placed on ice and spun for 2 min at 14,000 rpm on a benchtop centrifuge at  $4^{\circ}\text{C}$ . Next, 1  $\mu\text{l}$  of RNase H solution (2.5 U  $\text{ml}^{-1}$  RNase H [USB Corporation], 12.5 mM  $\text{MgCl}_2$ ) was

added, and the reaction was incubated at 37°C for 15 min. After RNase H treatment, the reaction was poly(A) tailed with 3.5 µl of 2.6× tailing solution (80 U terminal transferase [Roche], 2.6× terminal transferase buffer [Invitrogen], 1.9 mM dATP) for 15 min at 37°C and then heat-inactivated at 65°C for 10 min. The samples were placed on ice and spun for 2 min at 14,000 rpm on a benchtop centrifuge at 4°C. To each sample, 90 µl of ThermoPol PCR buffer was added to a final concentration of 1× ThermoPol buffer (New England Biolabs), 2.5 mM MgSO<sub>4</sub>, 1 mM dNTPs (Roche), 100 µg ml<sup>-1</sup> BSA (Roche), 10 U AmpliTaq (Applied Biosystems), and 5 µg AL1 primer<sup>15</sup>. Each reaction was split into three thin-walled 0.2-ml PCR tubes and amplified according to the following thermal cycling scheme: four cycles of 1 min at 94°C (denaturation), 2 min at 32°C (annealing), and 6 min plus 10 sec per cycle at 72°C (extension); 21 cycles of 1 min at 94°C (denaturation), 2 min at 42°C (annealing), and 6 min 40 sec plus 10 sec per cycle at 72°C (extension). The reaction was cooled, placed on ice, and the three tubes from each sample were pooled and amplified according to the following thermal cycling scheme: five cycles of 1 min at 94°C (denaturation), 2 min at 42°C (annealing), and 6 min at 72°C (extension). Further thermal cycling led to overamplification and loss of quantitative accuracy (K.A.J. and J.S.B., unpublished observations). Samples were stored at -20°C until use.

### Real-time quantitative PCR (RT-qPCR)

RT-qPCR of amplified material from stochastic sampling was measured as described previously<sup>34</sup>, except that tenfold less of each amplified sample was used as the starting cDNA template. Primer sequences and concentrations are shown in Supplementary Table 1.

### Small-cell reamplification and microarray hybridization

Amplified small-cell samples were reamplified and aminoallyl labeled in a 100 µl reaction containing 1× High-Fidelity buffer (Roche), 3.5 mM MgCl<sub>2</sub>, 200 µM dATP, dCTP, and dGTP, 40 µM dTTP (Roche), 160 µM aminoallyl-dUTP (Ambion), 100 µg ml<sup>-1</sup> BSA (Roche), 5 µg AL1 primer, and 1 µl amplified cDNA. Each reaction was amplified according to the following thermal cycling scheme: 1 min at 94°C (denaturation), 2 min at 42°C (annealing), and 3 min at 72°C (extension). In pilot experiments, 20 µl of this reaction for each stochastic sampling was monitored in the presence of 0.25× SYBR Green on a LightCycler II real-time PCR instrument (Roche). The number of amplification cycles (~20) was selected to ensure that all samples remained in the exponential phase during amplification<sup>35</sup>. Samples were purified on a PureLink column (Invitrogen), ethanol precipitated, and labeled with Alexa 555 amine-reactive dye (Invitrogen) according to the manufacturer's recommendation. Labeling efficiency was ~2 dye molecules per 100 bases.

For microarray hybridization, 1 µg Alexa 555-labeled cDNA (total volume: 5 µl) was mixed with 10 µl GEX hybridization buffer (Illumina). Samples were denatured at 94°C for 4 min and then added directly to HumanRef-8 Expression BeadChips (Illumina) prewarmed at 58°C. Slides were incubated at 58°C for 20 hr and washed according to the manufacturer's recommendations. After drying, slides were scanned on a BeadArray reader (Illumina) with a scan setting of "Direct hybridization 1". Samples were normalized to their mean overall fluorescence intensity relative to the overall dataset and then to the median fluorescence intensity of all transcripts detected ( $P < 0.1$ ) on each sample for subsequent analysis.

### Riboprobe synthesis

A 175–225 bp fragment of each gene was cloned by PCR into pcDNA3 (Invitrogen) from an MCF10A cDNA library generated by first-strand synthesis with Superscript III (Invitrogen) and an oligo(dT)<sub>24</sub> primer. Plasmids were linearized with the appropriate restriction enzymes and purified by phenol-chloroform extraction and ethanol precipitation. Riboprobes were synthesized from the linearized template by using the MAXIscript Sp6/T7 kit (Ambion) as recommended, except that *in vitro* transcriptions were incubated for 2 hr and Sp6 *in vitro* transcriptions were performed at 40°C to increase yield. Digoxigenin (DIG)- and dinitrophenyl (DNP)-labeled riboprobes were synthesized with 35% DIG-UTP (Roche) or DNP-UTP (Perkin Elmer) and 65% unlabeled UTP. After DNase digestion, riboprobes were ethanol precipitated, resuspended in RNase-free water to 0.2 µg ml<sup>-1</sup>, and stored at -80°C.

### Multicolor RNA fluorescence *in situ* hybridization (RNA FISH)

5 µm frozen sections of day 10 structures were cut on Superfrost Plus slides (Fisher), air dried, and stored at -80°C until further use. Slides were thawed at room temperature until completely dry, treated with 0.2 N HCl for 10 min, and washed in PBS for 5 min. Slides were then postfixed in 3.7% paraformaldehyde for 15 min, washed 2 × 10 min in PBS, and once in freshly prepared 0.1 M triethanolamine (pH 8.0) for 10 min. Samples were next acetylated with 0.25% acetic anhydride in freshly prepared 0.1 M triethanolamine (pH 8.0) for 5 min and washed in 2× SSC for 10 min. Slides were dehydrated with an ethanol series (2 min each of 70%, 95%, and 100% ethanol), and sections were covered with hybridization solution (1 mg ml<sup>-1</sup> yeast tRNA, 10% dextran sulfate in 2× SSC, 50% formamide) containing 50–500 ng ml<sup>-1</sup> of each riboprobe. Sections were covered with Parafilm, sealed with rubber cement, and incubated at 42°C in a humidified chamber for 14–16 hr.

After hybridization, slides were soaked in 2× SSC at 37°C for 5 min, the Parafilm was removed, and slides were washed in 2× SSC, 50% formamide for 30 min at 55°C, followed by 0.1× SSC for 30 min at 55°C. Slides were equilibrated in PBS for 10 min and then blocked for 1 hr at room temperature with 1× Western Blocking Reagent (Roche) in PBS + 0.3% Tween-20. After blocking, slides were incubated 1 hr at room temperature with 1× Western Blocking Reagent (Roche) in PBS + 0.3% Tween-20 containing anti-digoxin (1:500, Jackson ImmunoResearch) and anti-DNP (1:1,000, Invitrogen). Slides were washed 3 × 5 min in PBS and incubated for 1 hr at room temperature with 1× Western Blocking Reagent (Roche) in PBS + 0.3% Tween-20 containing Alexa 488-conjugated goat anti-rabbit (1:200, Invitrogen) and Alexa 555-conjugated goat anti-mouse (1:200, Invitrogen). Slides were washed 3 × 5 min in PBS and cell membranes were labeled with 20 µg ml<sup>-1</sup> Alexa 350-conjugated wheat-germ agglutinin for 5 min at room temperature. After two 5 min washes in PBS, autofluorescence was quenched with 10 mM CuSO<sub>4</sub> in 50 mM NH<sub>4</sub>Ac (pH 5.0) for 10 min<sup>36</sup>. Slides were washed with PBS for 5 min and mounted with 0.5% n-propyl gallate in PBS + 90% glycerol<sup>37</sup>.

### Immunofluorescence

5 µm sections of day 10 structures were cut on Superfrost Plus slides (Fisher), air dried, and stored at -80°C until further use. Slides were thawed at room temperature until completely

dry, hydrated  $3 \times 5$  min in PBS, and then blocked for 1 hr at room temperature with  $1 \times$  Western Blocking Reagent (Roche) in PBS + 0.3% Tween-20. After blocking, slides were incubated overnight at room temperature with  $1 \times$  Western Blocking Reagent (Roche) in PBS + 0.3% Tween-20 containing anti-p65 (A) (1:100, Santa Cruz) or anti-I $\kappa$ B $\alpha$  (C-21) (1:500, Santa Cruz). Slides were washed  $3 \times 5$  min in PBS and incubated for 1 hr at room temperature with  $1 \times$  Western Blocking Reagent (Roche) in PBS + 0.3% Tween-20 containing Alexa 555-conjugated goat anti-rabbit (1:200, Invitrogen). Slides were washed  $3 \times 5$  min in PBS and counterstained with  $0.5 \mu\text{g ml}^{-1}$  DAPI (Sigma) for 5 min at room temperature. After two 5 min washes in PBS, autofluorescence was quenched with 10 mM CuSO<sub>4</sub> in 50 mM NH<sub>4</sub>Ac (pH 5.0) for 10 min<sup>36</sup>. Slides were washed with PBS for 5 min and mounted with 0.5% n-propyl gallate in PBS + 90% glycerol<sup>37</sup>.

## Microscopy

Frozen sections and coverslips were imaged with a  $40 \times 1.3$  NA oil objective on an BX51 upright fluorescence microscope (Olympus) with the following filter sets: ET-DAPI (excitation: 325–375 nm, dichroic: 400 nm, emission: 435–485 nm), ET-FITC (excitation: 450–490 nm, dichroic: 495 nm, emission: 500–550 nm), ET-CY3 (excitation: 520–570 nm, dichroic: 565 nm, emission: 570–640 nm), and ET-CY5 (excitation: 590–650 nm, dichroic: 660 nm, emission: 665–735 nm). Images were captured with an Orca R2 CCD camera (Hamamatsu) at  $2 \times 2$  binning and exposure times that filled 90% of the camera bit depth, with the exception of the RNA FISH sense controls (**Supplementary Figs. 9 and 10**) where the exposure time was matched to the antisense image. Displayed images were rainbow pseudocolored with a linear lookup table that covered the full range of the data for each fluorescence channel.

## Image segmentation and quantification

Single cells from RNA FISH images were segmented by hand based on wheat-germ agglutinin staining (DAPI channel), and traced image segments were then applied to the DIG- and DNP-labeled riboprobe stainings (FITC and Cy3 channels). Median fluorescence intensities per cell for each riboprobe were calculated, and individual images were normalized to the maximum observed intensity in each channel for comparison across multiple images.

## Statistical analysis

Statistical analyses of RT-qPCR measurements were performed on the cycle thresholds of the measured genes. This is equivalent to a  $\log_2$  transformation, which allows log-normal distributions to be treated as normal distributions<sup>13, 31</sup>. Estimation of the coefficient of variation for amplification replicates (Fig. 2i) was done in Igor Pro (WaveMetrics) by nonlinear least-squares curve fitting of the mean-centered cycle thresholds to a normal distribution with a mean of zero. Confidence intervals on CVs were calculated with McKay's transformation<sup>38</sup>, and non-overlapping 90% confidence intervals were considered significantly different.  $\chi^2$  goodness-of-fit tests for sampling fluctuations were performed in MATLAB with the `chi2gof` function, a mean of zero, and a standard deviation equal to the

reference distribution (false-discovery rate = 0.05). Nonparametric confidence intervals for the RNA FISH subpopulations were based on a binomial distribution.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

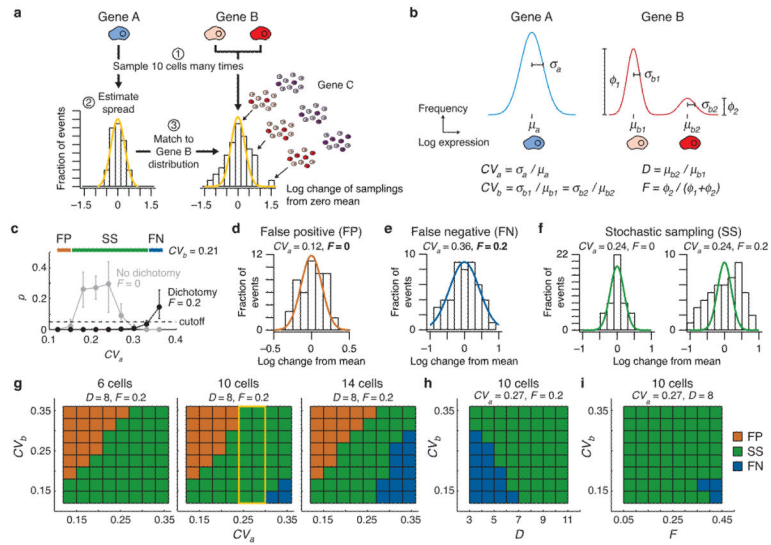
## ACKNOWLEDGMENTS

We thank Tim McDaniel (Illumina) for generously providing the microarrays used in this study, Greg Cox (Molecular Probes) for advice during development of the RNA FISH protocol, and Christian Reinhardt (MIT) for critically reading the manuscript. This work was supported by the National Institutes of Health (5-R01-CA105134-07 to J.S.B.), the National Institutes of Health Director's New Innovator Award Program (1-DP2-OD006464-01 to K.A.J.), the Mary Kay Ash Charitable Foundation (to K.A.J.), and the Pew Scholars Program in the Biomedical Sciences (to K.A.J.).

## REFERENCES

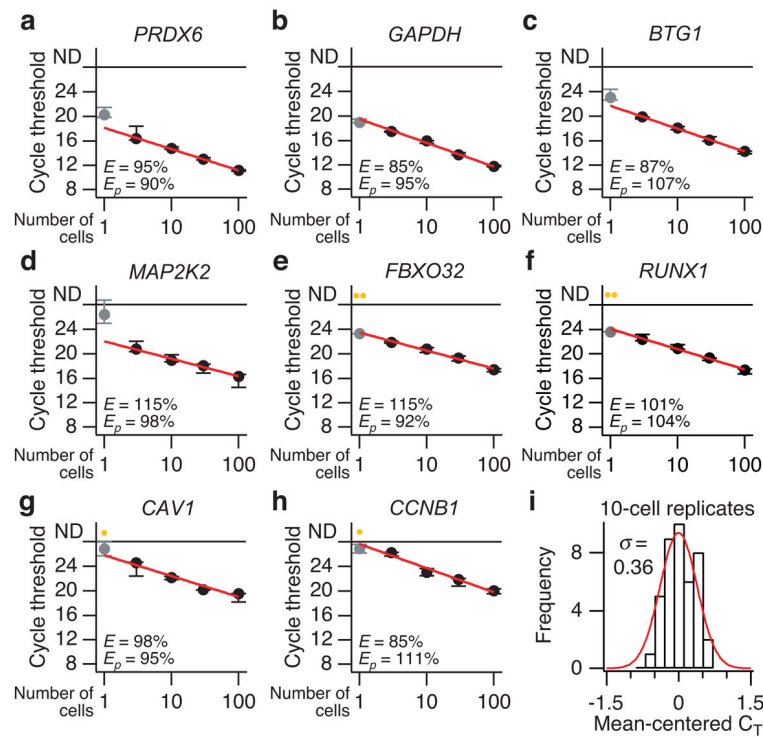
1. Wernet MF, et al. Stochastic spineless expression creates the retinal mosaic for colour vision. *Nature*. 2006; 440:174–180. [PubMed: 16525464]
2. Chang HH, Hemberg M, Barahona M, Ingber DE, Huang S. Transcriptome-wide noise controls lineage choice in mammalian progenitor cells. *Nature*. 2008; 453:544–547. [PubMed: 18497826]
3. Irish JM, Kotecha N, Nolan GP. Mapping normal and cancer cell signalling networks: towards single-cell proteomics. *Nat Rev Cancer*. 2006; 6:146–155. [PubMed: 16491074]
4. Ferrell JE Jr, Machleder EM. The biochemical basis of an all-or-none cell fate switch in *Xenopus* oocytes. *Science*. 1998; 280:895–898. [PubMed: 9572732]
5. Altan-Bonnet G, Germain RN. Modeling T cell antigen discrimination based on feedback control of digital ERK responses. *PLoS Biol*. 2005; 3:e356. [PubMed: 16231973]
6. Emmert-Buck MR, et al. Laser capture microdissection. *Science*. 1996; 274:998–1001. [PubMed: 8875945]
7. Tietjen I, et al. Single-cell transcriptional analysis of neuronal progenitors. *Neuron*. 2003; 38:161–175. [PubMed: 12718852]
8. Bahar R, et al. Increased cell-to-cell variation in gene expression in ageing mouse heart. *Nature*. 2006; 441:1011–1014. [PubMed: 16791200]
9. Fraser HB, Hirsh AE, Giaever G, Kumm J, Eisen MB. Noise minimization in eukaryotic gene expression. *PLoS Biol*. 2004; 2:e137. [PubMed: 15124029]
10. Kurimoto K, et al. An improved single-cell cDNA amplification method for efficient high-density oligonucleotide microarray analysis. *Nucleic Acids Res*. 2006; 34:e42. [PubMed: 16547197]
11. Debnath J, Brugge JS. Modelling glandular epithelial cancers in three-dimensional cultures. *Nat Rev Cancer*. 2005; 5:675–688. [PubMed: 16148884]
12. Golding I, Paulsson J, Zawilski SM, Cox EC. Real-time kinetics of gene activity in individual bacteria. *Cell*. 2005; 123:1025–1036. [PubMed: 16360033]
13. Bengtsson M, Stahlberg A, Rorsman P, Kubista M. Gene expression profiling in single cells from the pancreatic islets of Langerhans reveals lognormal distribution of mRNA levels. *Genome Res*. 2005; 15:1388–1392. [PubMed: 16204192]
14. Newman JR, et al. Single-cell proteomic analysis of *S. cerevisiae* reveals the architecture of biological noise. *Nature*. 2006; 441:840–846. [PubMed: 16699522]
15. Brady G, Iscove NN. Construction of cDNA libraries from single cells. *Methods Enzymol*. 1993; 225:611–623. [PubMed: 8231874]
16. Hartmann CH, Klein CA. Gene expression profiling of single cells on large-scale oligonucleotide arrays. *Nucleic Acids Res*. 2006; 34:e143. [PubMed: 17071717]
17. Taniguchi K, Kajiya T, Kambara H. Quantitative analysis of gene expression in a single cell by qPCR. *Nat Methods*. 2009; 6:503–506. [PubMed: 19525960]

18. Schmelzle T, et al. Functional role and oncogene-regulated expression of the BH3-only factor Bmf in mammary epithelial anoikis and morphogenesis. *Proc Natl Acad Sci U S A*. 2007; 104:3787–3792. [PubMed: 17360431]
19. Debnath J, Walker SJ, Brugge JS. Akt activation disrupts mammary acinar architecture and enhances proliferation in an mTOR-dependent manner. *J Cell Biol*. 2003; 163:315–326. [PubMed: 14568991]
20. Pearson GW, Hunter T. Real-time imaging reveals that noninvasive mammary epithelial acini can contain motile cells. *J Cell Biol*. 2007; 179:1555–1567. [PubMed: 18166657]
21. Pearson GW, Hunter T. PI-3 kinase activity is necessary for ERK1/2-induced disruption of mammary epithelial architecture. *Breast Cancer Res*. 2009; 11:R29. [PubMed: 19457236]
22. Rakha EA, Reis-Filho JS, Ellis IO. Basal-like breast cancer: a critical review. *J Clin Oncol*. 2008; 26:2568–2581. [PubMed: 18487574]
23. Ginestier C, et al. ALDH1 is a marker of normal and malignant human mammary stem cells and a predictor of poor clinical outcome. *Cell Stem Cell*. 2007; 1:555–567. [PubMed: 18371393]
24. Tanaka M, et al. A novel RNA-binding protein, Ossa/C9orf10, regulates activity of Src kinases to protect cells from oxidative stress-induced apoptosis. *Mol Cell Biol*. 2009; 29:402–413. [PubMed: 19015244]
25. Yamaguchi A, et al. Stress-associated endoplasmic reticulum protein 1 (SERP1)/Ribosome-associated membrane protein 4 (RAMP4) stabilizes membrane proteins during stress and facilitates subsequent glycosylation. *J Cell Biol*. 1999; 147:1195–1204. [PubMed: 10601334]
26. Gross DN, van den Heuvel AP, Birnbaum MJ. The role of FoxO in the regulation of metabolism. *Oncogene*. 2008; 27:2320–2336. [PubMed: 18391974]
27. Karin M, Ben-Neriah Y. Phosphorylation meets ubiquitination: the control of NF- $\kappa$ B activity. *Annu Rev Immunol*. 2000; 18:621–663. [PubMed: 10837071]
28. Laslo P, et al. Multilineage transcriptional priming and determination of alternate hematopoietic cell fates. *Cell*. 2006; 126:755–766. [PubMed: 16923394]
29. Yakoby N, et al. A combinatorial code for pattern formation in *Drosophila* oogenesis. *Dev Cell*. 2008; 15:725–737. [PubMed: 19000837]
30. Neve RM, et al. A collection of breast cancer cell lines for the study of functionally distinct cancer subtypes. *Cancer Cell*. 2006; 10:515–527. [PubMed: 17157791]
31. Warren L, Bryder D, Weissman IL, Quake SR. Transcription factor profiling in individual hematopoietic progenitors by digital RT-PCR. *Proc Natl Acad Sci U S A*. 2006; 103:17807–17812. [PubMed: 17098862]
32. Sheskin, DJ. *Handbook of Parametric and Nonparametric Statistical Procedures*. Edn. 4th.. Chapman & Hall; New York: 2007.
33. Debnath J, Muthuswamy SK, Brugge JS. Morphogenesis and oncogenesis of MCF-10A mammary epithelial acini grown in three-dimensional basement membrane cultures. *Methods*. 2003; 30:256–268. [PubMed: 12798140]
34. Miller-Jensen K, Janes KA, Brugge JS, Lauffenburger DA. Common effector processing mediates cell-specific responses to stimuli. *Nature*. 2007; 448:604–608. [PubMed: 17637676]
35. Nagy ZB, et al. Real-time polymerase chain reaction-based exponential sample amplification for microarray gene expression profiling. *Anal Biochem*. 2005; 337:76–83. [PubMed: 15649378]
36. Schnell SA, Staines WA, Wessendorf MW. Reduction of lipofuscin-like autofluorescence in fluorescently labeled tissue. *J Histochem Cytochem*. 1999; 47:719–730. [PubMed: 10330448]
37. Giloh H, Sedat JW. Fluorescence microscopy: reduced photobleaching of rhodamine and fluorescein protein conjugates by n-propyl gallate. *Science*. 1982; 217:1252–1255. [PubMed: 7112126]
38. McKay AT. Distribution of the Coefficient of Variation and the Extended 't' Distribution. *J Roy Stat Soc*. 1932; 95:695–698.



**Figure 1. Small-cell profiling by stochastic sampling can distinguish transcriptional heterogeneities from normal biological variation**

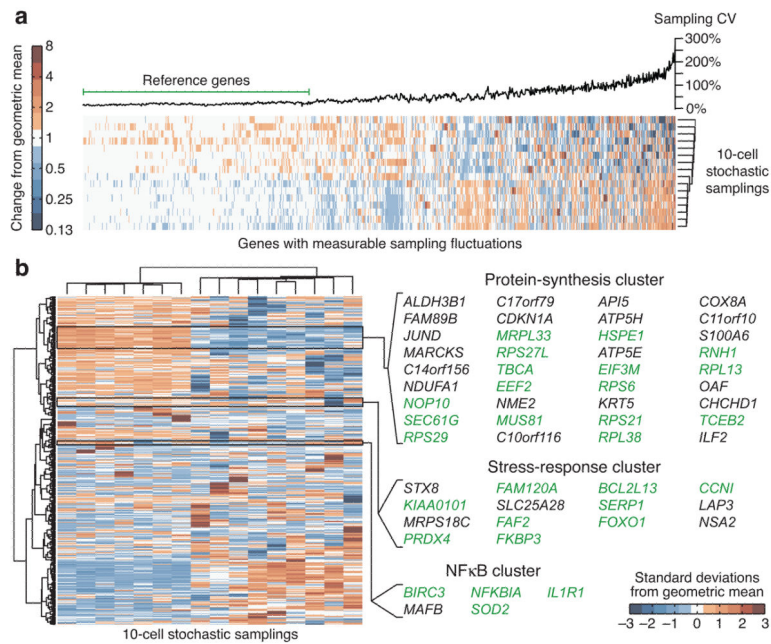
(a) The statistical and empirical steps of stochastic sampling, as described in the main text. The distributions shown were based on 48 simulated samplings with the following model parameters:  $CV_a = 25\%$ ,  $CV_b = 25\%$ ,  $D = 8$ ,  $F = 0.2$  (as defined in the main text). (b) Theoretical population distributions of a constitutively expressed gene (Gene A) and a dichotomy with two subpopulations (Gene B). (c) Identifying false positives (FP, brown), false negatives (FN, blue), and effective stochastic sampling (SS, green) through Monte Carlo simulations. Stochastic-sampling experiments were simulated as described in the **Online Methods** with the indicated parameters and  $D = 8$ . Data are shown as the median  $p$  value for the  $\chi^2$  goodness of fit between the test and reference distributions  $\pm 90\%$  nonparametric confidence intervals from 50 simulations of 48 samplings. (d-f) Examples of false positives (brown), false negatives (green), and effective stochastic sampling (blue) for  $D = 8$  and  $CV_b = 0.21$ . (g) Effective stochastic sampling with up to 10 averaged cells. Note that when 10 cells are averaged and  $CV_a = 25\text{--}30\%$  (yellow box), stochastic sampling is effective for all values of  $CV_b$ . (h) Effective stochastic sampling for dichotomies with expression differences greater than fivefold. (i) Stochastic sampling is not strongly dependent on the relative proportion of subpopulations in a dichotomy.



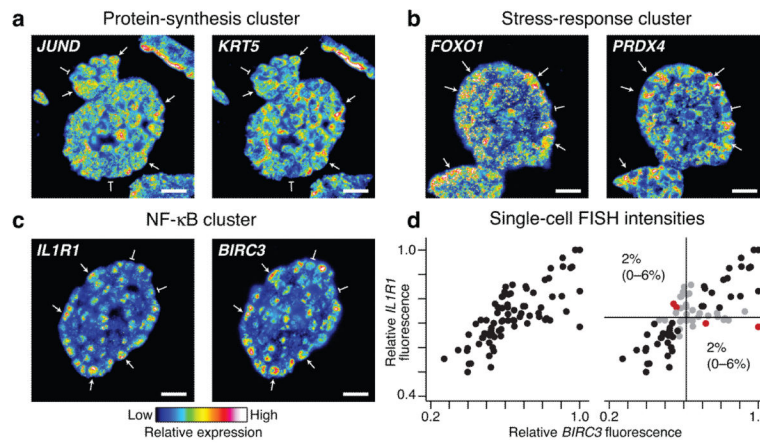
**Figure 2. Quantitative and reproducible small-cell amplification of high- to low-abundance transcripts from 3–100 cells**

(a–h) The RT-qPCR cycle threshold for each gene is plotted as a function of starting cellular material and is shown as the median  $\pm$  range of three replicate small-cell amplifications. Amplification efficiencies ( $E$ ) based on a log-linear fit of the 3–100-cell dilutions (red line) are listed along with primer efficiencies ( $E_p$ ) calculated by serially diluting the template before RT-qPCR. Genes are ordered **a** through **h** in the order of increasing median cycle threshold from the 10-cell replicates, which was used as an approximation of relative abundance (lower cycle thresholds suggest increased relative abundance). Note that the one-cell amplifications (gray) of higher-abundance transcripts (**a–d**) often deviate from the log-linear fit, and the one-cell amplification of lower-abundance transcripts (**e–h**) are frequently not detectable (yellow, ND). (i) Reproducible small-cell amplification of 10 cells. The cycle thresholds from 10-cell amplification replicates of all genes were mean centered, grouped, and fit to a normal distribution. The standard deviation ( $\sigma$ ) of the mean-centered cycle thresholds ( $C_T$ ) was 0.36, corresponding to a coefficient of variation of 28%, assuming that amplicons double after each cycle (i.e., 100% efficiency,  $2^{0.36} - 1 = 0.28$ ).



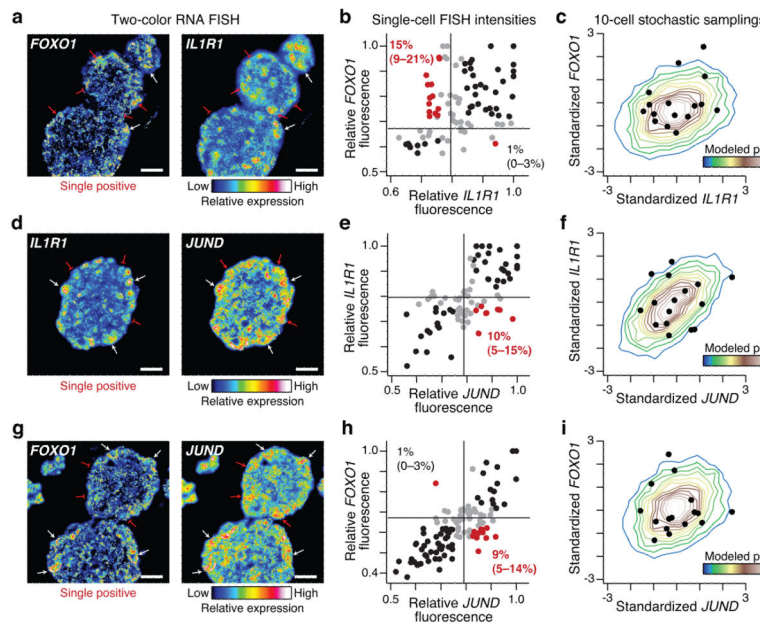


**Figure 3. Stochastic profiling of matrix-attached cells at day 10 of MCF10A morphogenesis**  
**(a)** Hierarchical clustering of unscaled sampling fluctuations for transcripts with measurable biological variation. Genes with sampling variations greater than measurement error were clustered using a Euclidean distance metric and average linkage. The genes with consistent CV values (left) were used as the reference subset for calculating an appropriate reference distribution used to test for heterogeneous expression. **(b)** Hierarchical clustering of scaled sampling fluctuations for transcripts predicted to be heterogeneously expressed by stochastic profiling. Candidate heterogeneities were scaled to unit variance and clustered using a Euclidean distance metric and Ward's linkage. Selected clusters were examined for enriched biological functions. Genes consistent with the assigned functions are highlighted in green.



**Figure 4. Stochastic profiling identifies clusters of heterogeneously coexpressed transcripts**

Two-color RNA FISH images were collected at day 10 of MCF10A morphogenesis for (a) *JUND* and *KRT5* in the protein-synthesis cluster, (b) *FOXO1* and *PRDX4* in the stress-response cluster, and (c) *IL1R1* and *BIRC3* in the NF- $\kappa$ B cluster. Images are pseudocolored to highlight quantitative differences in fluorescence intensity, and single cells showing strong coexpression are highlighted with arrows (high expression) or flat markers (low expression). Two-color images for 3–4 additional gene pairs within each cluster are shown in Supplementary Figure 11. (d) *BIRC3*–*IL1R1* images were segmented to quantify average fluorescence intensities in single cells as described in the **Online Methods**. Data are shown from cells in four independent acini after normalization to the maximum observed cellular fluorescence signal in each image. Gates were defined as the 25<sup>th</sup> percentile centered on the median fluorescence intensity (black lines) for each gene. Observations that were within the range of the gates were scored as neither positive nor negative (gray). Single positive cells (red) are shown as the percentage of the overall cell population with 90% confidence intervals in parentheses. For a–c, scale bar is 20  $\mu$ m.



**Figure 5. Stochastic profiling distinguishes heterogeneous expression patterns that are not exclusively coexpressed**

Two-color RNA FISH images were collected at day 10 of MCF10A morphogenesis and compared to the stochastic-profiling data for (a-c) *FOXO1* and *IL1R1*, (d-f) *IL1R1* and *JUND*, and (g-i) *FOXO1* and *JUND*. (a,d,g) Representative pseudocolored images, containing single-positive cells highlighted with red arrows (high expression) or flat markers (low expression) Scale bar is 20  $\mu$ m. (b,e,h) Fluorescence intensities are shown from cells in four independent acini after normalization to the maximum observed cellular fluorescence signal in each image. Gates were defined as the 25<sup>th</sup> percentile centered on the median fluorescence intensity (black lines) for each gene. Observations that were within the range of the gates were scored as neither positive nor negative (gray). The percentages of single-positive cells (red) in the overall cell population are shown with 90% confidence intervals in parentheses. (c,f,i) RNA FISH distributions of b,e,h were resampled as 10-cell averages, standardized, and the resulting probability density function (pdf) was compared to the standardized sampling fluctuations from stochastic profiling (black circles) shown in Figure 3b.