



**HAL**  
open science

## Exploiting Single-Cell Quantitative Data to Map Genetic Variants Having Probabilistic Effects

Florent Chuffart, Magali Richard, Daniel Jost, Claire Burny, H el ene Duplus-Bottin, Yoshikazu Ohya, Ga el Yvert

► **To cite this version:**

Florent Chuffart, Magali Richard, Daniel Jost, Claire Burny, H el ene Duplus-Bottin, et al.. Exploiting Single-Cell Quantitative Data to Map Genetic Variants Having Probabilistic Effects. *PLoS Genetics*, 2016, 12 (8), pp.e1006213. 10.1371/journal.pgen.1006213 . hal-01976590

**HAL Id: hal-01976590**

**<https://hal.science/hal-01976590v1>**

Submitted on 13 Jan 2025

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destin ee au d ep ot et  a la diffusion de documents scientifiques de niveau recherche, publi es ou non,  emanant des  tablissements d'enseignement et de recherche fran ais ou  trangers, des laboratoires publics ou priv es.

RESEARCH ARTICLE

# Exploiting Single-Cell Quantitative Data to Map Genetic Variants Having Probabilistic Effects

Florent Chuffart<sup>1</sup>, Magali Richard<sup>1</sup>, Daniel Jost<sup>2,3</sup>, Claire Burny<sup>1</sup>, H el ene Duplus-Bottin<sup>1</sup>, Yoshikazu Ohya<sup>4</sup>, Ga el Yvert<sup>1\*</sup>

**1** Laboratoire de Biologie et de Mod elisation de la Cellule, Ecole Normale Sup erieure de Lyon, CNRS, Universit  de Lyon, Lyon, France, **2** Laboratoire de Physique, Ecole Normale Sup erieure de Lyon, CNRS, Universit  de Lyon, Lyon, France, **3** University Grenoble Alpes, CNRS, TIMC-IMAG lab, UMR 5525, Grenoble, France, **4** Department of Integrated Biosciences, Graduate School of Frontier Sciences, University of Tokyo, Kashiwa, Chiba, Japan

\* [Gael.Yvert@ens-lyon.fr](mailto:Gael.Yvert@ens-lyon.fr)



CrossMark  
click for updates

 OPEN ACCESS

**Citation:** Chuffart F, Richard M, Jost D, Burny C, Duplus-Bottin H, Ohya Y, et al. (2016) Exploiting Single-Cell Quantitative Data to Map Genetic Variants Having Probabilistic Effects. *PLoS Genet* 12(8): e1006213. doi:10.1371/journal.pgen.1006213

**Editor:** Justin C. Fay, Washington University School of Medicine, UNITED STATES

**Received:** February 12, 2016

**Accepted:** July 2, 2016

**Published:** August 1, 2016

**Copyright:**   2016 Chuffart et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** The yeast morphological data corresponds to the experiments described in Nogami et al. *PLoS Genetics* 2007. For the present study, raw images were re-analyzed using CalMorph 1.0. The single-cell values and genotypes used are available in the supplementary material of this article. The flow cytometry data corresponding to yeast galactose response is made available from <http://flowrepository.org> under accession number FR-FCM-ZZPA.

**Funding:** This work was supported by the European Research Council under the European Union's Seventh Framework Programme FP7/2007-2013

## Abstract

Despite the recent progress in sequencing technologies, genome-wide association studies (GWAS) remain limited by a statistical-power issue: many polymorphisms contribute little to common trait variation and therefore escape detection. The small contribution sometimes corresponds to incomplete penetrance, which may result from probabilistic effects on molecular regulations. In such cases, genetic mapping may benefit from the wealth of data produced by single-cell technologies. We present here the development of a novel genetic mapping method that allows to scan genomes for single-cell Probabilistic Trait Loci that modify the statistical properties of cellular-level quantitative traits. Phenotypic values are acquired on thousands of individual cells, and genetic association is obtained from a multi-variate analysis of a matrix of Kantorovich distances. No prior assumption is required on the mode of action of the genetic loci involved and, by exploiting all single-cell values, the method can reveal non-deterministic effects. Using both simulations and yeast experimental datasets, we show that it can detect linkages that are missed by classical genetic mapping. A probabilistic effect of a single SNP on cell shape was detected and validated. The method also detected a novel locus associated with elevated gene expression noise of the yeast galactose regulon. Our results illustrate how single-cell technologies can be exploited to improve the genetic dissection of certain common traits. The method is available as an open source R package called *ptlmapper*.

## Author Summary

Genetic association studies are usually conducted on phenotypes measured at the scale of whole tissues or individuals, and not at the scale of individual cells. However, some common traits, such as cancer, can result from a minority of cells that adopted a special behavior. From one individual to another, DNA variants can modify the frequency of such

Grant Agreement n°281359 attributed to GY. DJ was supported by the Institut Rhône-Alpin des Systèmes Complexes and program AGIR of University Grenoble Alpes. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

**Competing Interests:** The authors have declared that no competing interests exist.

cellular behaviors. The body of one of the individuals then harbours more misbehaving cells and is therefore predisposed to a macroscopic phenotypic change, such as disease. Such genetic effects are probabilistic, they contribute little to trait variation at the macroscopic level and therefore largely escape detection in classical studies. We have developed a novel statistical method that uses single-cell measurements to detect variants of the genome that have non-deterministic effects on cellular traits. The approach is based on a comparison of distributions of single-cell traits. We applied it to colonies of yeast cells and showed that it can detect mutations that change cellular morphology or molecular regulations in a probabilistic manner. This opens the way to study multicellular organisms from a novel angle, by exploiting single-cell technologies to detect genetic variants that predispose to certain diseases or common traits.

## Introduction

Modern genetics aims to identify DNA variants contributing to common trait variation between individuals. A high motivation to map such variants is shared worldwide because many heritable traits relate to social and economical preoccupations, such as human health or agronomical and industrial yields. In addition to the molecular knowledge they provide, these variants fuel the development of personalized and predictive medicine as well as the improvement of economically-relevant plants, animal breeds or biotechnology materials. However, this high ambition is accompanied by a major challenge: common traits are under the control of numerous variants that each contribute little to phenotypic variation [1], and this modest contribution of each variant hampers the statistical power to detect them. Power is further limited by the multiplicity of linkage tests when scanning whole genomes. The consequence of this has been debated under the term "missing heritability": most of the genetic variants of interest remain to be identified. Currently, this issue is handled by modelling the effect of known or hidden factors, and by scaling up sample size up to tens of thousands of individuals [2–4]. Practically, however, cohort size cannot be infinitely increased, and relevant factors are difficult to choose. Studies would therefore greatly benefit from a better detection of small genetic effects, and from a reduction of the number of genomic loci to test.

Small-effect variants are typically associated with predisposition (or incomplete penetrance): carriers of a mutation display a phenotype at increased frequency, but not all of them do. In this probabilistic context, the statistical properties of cellular traits may sometimes become informative: a tissue may break because cells have an increased probability to detach, a tumor may emerge because a cell type has an increased probability of somatic mutations, a chemotherapy may fail if cancer cells have an increased probability to be in a persistent state. In other words, molecular events in one or few cells can have devastating consequences at the multicellular level. As discussed previously [5], cellular-scale probabilities are likely related to the genotype and this relation may sometimes underlie genetic predisposition [6]. Striking examples are genetic factors affecting the mutation rate of somatic divisions and thereby modifying cancer predisposition. These loci have a probabilistic effect on a cellular trait: the amount of *de novo* mutations in the cell's daughter. Other loci may modulate the heterogeneity between isogenic cancer cells that underlies tumour progression [7,8] and resistance to chemotherapy [9–11]. They would then change the fraction of problematic cells between individuals and thereby disease progression or treatment outcome.

Fortunately, the experimental throughput of single-cell measurements has recently exploded. Technological developments in high-throughput flow cytometry [12], multiplexed

mass-cytometry [13], image content analysis [14–16] and droplet-based single-cell transcriptome profiling [17,18] now offer the possibility to estimate empirically the statistical distribution of numerous molecular and cellular single-cell quantitative traits. We therefore propose to scan genomes for variants that modify single-cell traits in a probabilistic manner, which we call single-cell Probabilistic Trait Loci (scPTL). This requires to monitor not only the macroscopic trait of many individuals but also a relevant cellular trait in many cells of these individuals. After scPTL are found, they can constitute a set of candidate loci to be directly tested for a possible small effect on the macroscopic trait of interest.

Methods are needed to detect scPTL. With its fast generation time, high recombination rate and reduced genome size, the unicellular yeast *Saccharomyces cerevisiae* offers a powerful experimental framework for developing such methods. Using this model organism, scPTL were discovered by treating one statistical property of the single-cell trait, such as its variance in the population of cells, as a quantitative trait and by applying Quantitative Trait Locus (QTL) mapping to it [19,20]. However, this approach is limited because it is difficult to anticipate *a priori* which summary statistics must be used.

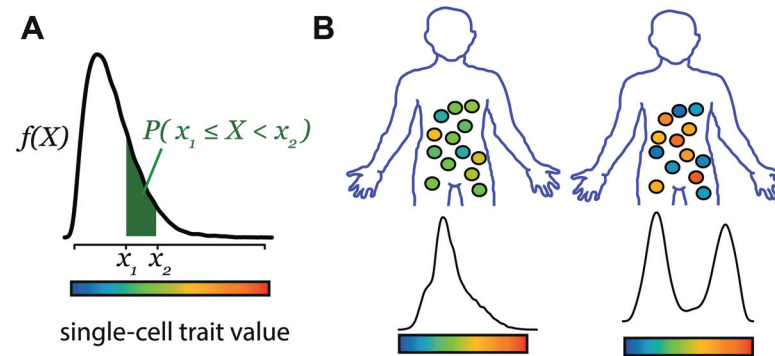
We present here the development of a genome-scan method that exploits all single-cell values with no prior simplification of the cell population phenotype. Using simulations and existing single-cell data from yeast, we show that it can detect genetic effects that were missed by conventional linkage analysis. When applied to a novel experimental dataset, the method detected a locus of the yeast genome where natural polymorphism modifies cell-to-cell variability of the activation of the GAL regulon. This work shows how single-cell quantitative data can be exploited to detect probabilistic effects of DNA variants. Our approach is conceptually and methodologically novel in quantitative genetics. Although we validated it using a unicellular organism, it opens alternative ways to apprehend the genetic predisposition of multicellular organisms to certain complex traits.

## Results

### Definitions

We specify here the concepts and definitions that are used in the present study. Let  $X$  be a quantitative trait that can be measured at the level of individual cells.  $X$  is affected by the genotype of the cells and by their environmental context. However, even for isogenic cells sharing a common, supposedly homogeneous environment,  $X$  may differ between the cells. To describe the values of  $X$  among cells sharing a common genotype and environment, we define a *single-cell quantitative trait density function*  $f$  [5] as the function underlying the probability that a cell expresses  $X$  at a given level (Fig 1A). Statistically speaking,  $f$  represents the probability density function of the random variable  $X$ . In the present study, this function  $f(X)$  constitutes the 'phenotype' of the individual from whom the cells are studied. As for any macroscopic phenotype, it can depend on the environmental context of the individual (diet, age, disease. . .) as well as on its genotype. Single-cell trait density functions also obviously depend on the properties of the cells that are studied, such as their differentiation state or proliferation rate.

We focus here on the effect of the genotype. Conceptually, cells from one individual may follow a density function of  $X$  that is different from the one followed by cells of another individual, because of genotypic differences between the two individuals (Fig 1B). The important concept is that the genetic difference has probabilistic consequences: it changes the probability that a cell expresses  $X$  at a given level, but it does not necessarily change  $X$  in most of the cells. Depending on the nature of trait  $X$  and how the two functions differ, such a genetic effect can have implications on macroscopic traits and predisposition to disease [5]. The term *single-cell Probabilistic*



**Fig 1. Concept and definitions.** **A)** A cellular trait is considered as a random variable  $X$  with density function  $f$ . The probability  $P(x_1 \leq X < x_2)$  ( $= \int_{x_1}^{x_2} f(X)dX$ ) that one cell expresses  $X$  at a value comprised between  $x_1$  and  $x_2$  is given by the shaded area. **B)**  $f$  differs between individuals because of environmental and genetic factors.

doi:10.1371/journal.pgen.1006213.g001

*Trait Locus* will refer here to a genetic locus modifying *any* characteristics of  $f$  (that is, changing allele A in allele B at the locus changes the density function  $f$  of  $X$ , i.e.  $f_B \neq f_A$ ).

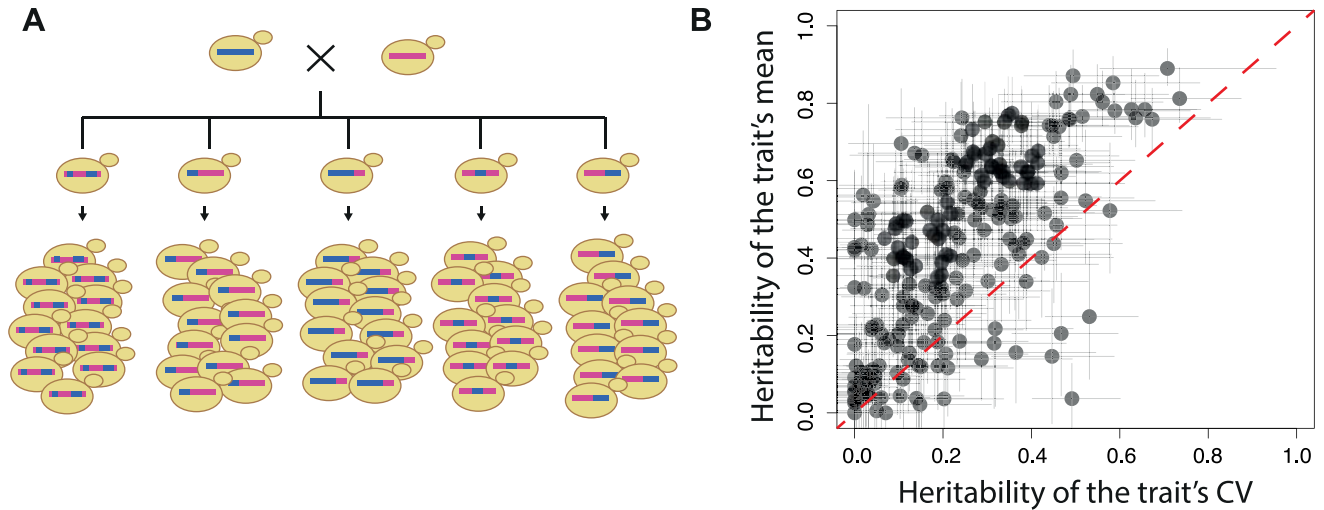
A quantitative trait locus (QTL) linked to  $X$  is a location on a chromosome where a genetic variant changes the mean or the median of  $X$  in the cell population. Similarly, a varQTL is a genetic locus changing the variance of  $X$  and a cvQTL is a genetic locus changing the coefficient of variation (standard deviation divided by the mean, abbreviated CV) of  $X$  in the cell population. All three types of loci (QTL, varQTL and cvQTL) assume a change in  $f$  and they are therefore special cases of scPTL. However, not all scPTL are QTL: many properties of  $f$  may change while preserving its mean, median, variance or CV. The purpose of the present study was to develop an approach that could identify scPTL without knowing *a priori* how it might change  $f$ .

### Mean and variability of cellular traits can have distinct genetic heritabilities

An important question before investing efforts in scPTL mapping is whether genotypes can modify  $f$  without affecting its expected value (the mean of  $X$ ). If not, then QTL mapping will capture the genetic modifiers of  $f$  and searching for more complex scPTL is not justified. In contrast, if other-than-mean genotypic changes of  $f$  are frequent, then scPTL can considerably complement QTL to control single-cell traits. In this case, scPTL mapping becomes important.

In multicellular organisms, cell types and intermediate differentiation states constitute the predominant source of cellular trait variation. Studying their single-cell statistical characteristics requires accounting for the developmental status of the cells. This constitutes a major challenge that can be avoided by studying unicellular organisms. The yeast *S. cerevisiae* provides the opportunity to study individual cells that all belong to a single cell type, in the context of a powerful genetic experimental system. By analysing specific gene expression traits in this organism, we and others identified loci that meet the definition of scPTL but not of QTL [20,21]. This illustrated that, for some traits, scPTL mapping could complement classic quantitative genetics to identify the genetic sources of cellular trait variation.

To estimate if non-QTL scPTL are frequent, we re-analysed an experimental dataset corresponding to the genetic segregation of many single-cell traits in a yeast cross (Fig 2A). After a round of meiosis involving two unrelated natural backgrounds of *S. cerevisiae*, individual segregants had been amplified by mitotic (clonal) divisions and traits of cellular morphology were acquired by semi-automated fluorescent microscopy and image analysis [22]. This way, for



**Fig 2. Genetic heritabilities for mean and CV of cellular yeast traits are not correlated.** **A)** Scheme of the experimental data used to compute genetic heritabilities. The dataset is from [22]. **B)** Broad-sense genetic heritability of mean (y-axis) and CV (coefficient of variation, x-axis) for 220 traits describing the morphology of individual cells (see [methods](#)). Each dot corresponds to one trait. Negative values were set to zero. Bars: 95% confidence intervals. Spearman correlation coefficient: 0.671.

doi:10.1371/journal.pgen.1006213.g002

each of 59 segregants, 220 single-cell traits were measured in about 200 isogenic cells, which enabled QTL mapping of these traits. We reasoned that if all scPTL of a trait are also QTL, then a high genetic heritability of any property of  $f$  should coincide with a high genetic heritability of the expected value of  $f$ . In particular, the coefficient of variation (CV) of a single-cell trait should then display high heritability only if the mean value of the trait also does. To see if this was the case, we computed for each trait the broad-sense genetic heritabilities of both the mean and CV of the trait. Note that the genetic heritability computed here is not the same as the mitotic heritability of cellular traits transmitted from mother to daughter cells. Here, a value (mean or CV) is computed on a population of cells, and its heritability corresponds to the proportion of its variation that can be attributed to genetic differences between the cell populations (see [methods](#)). Overall, heritability of mean was higher than heritability of CV, and the two types of heritabilities were correlated ([Fig 2B](#)). We also observed that several traits had high heritability of CV and low heritability of their mean value, or *vice versa*. This indicates that, for some traits, genetic factors exist that modify the trait CV but not the trait mean. This observation is in agreement with the complex CV-vs-mean dependency previously reported in this type of data [23,24]. We therefore sought to develop a method that can detect scPTL that do not necessarily correspond to QTL.

### Principle of scPTL mapping

One way to identify scPTL from experimental measures is to compute a summary statistic of the trait distribution, such as one of its moments, and then scan for QTL controlling this quantity. This approach is particularly appropriate when searching for specific genetic effects on  $f$ , such as a change in the level of cell-to-cell variability, and a few previous studies successfully used it to map varQTL and cvQTL [19,20,22,25,26]. However, it is less adapted when nothing is known on the way  $f$  may depend on genetic factors.

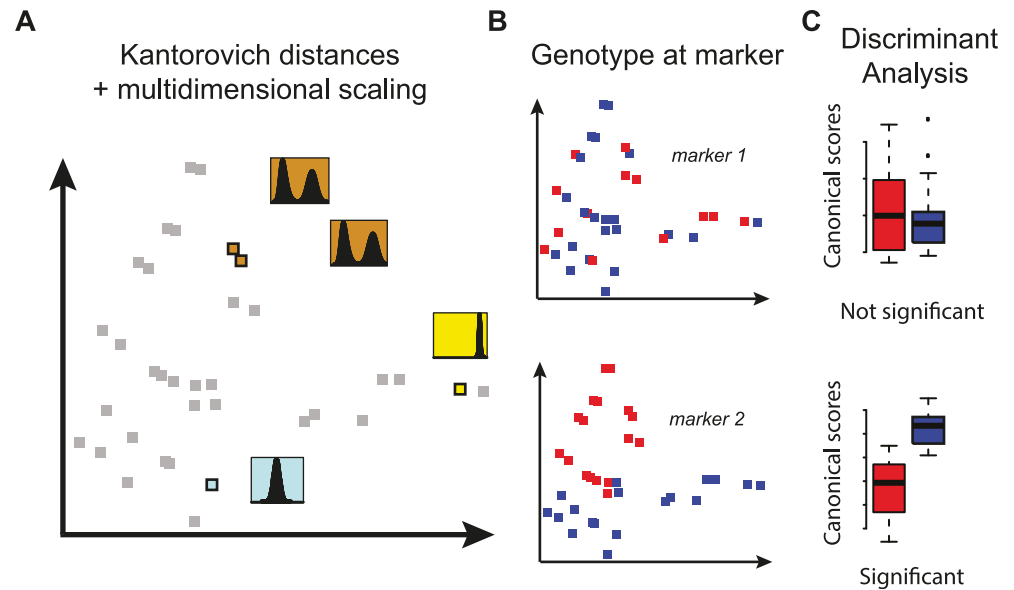
Scanning for scPTL considers the entire distribution of single-cell trait values as the phenotype of interest and searches the genome for a statistical association with *any* change in the distribution. We assume that for a set of genotypic categories (individuals for multicellular

organisms, or populations of cells for unicellular ones), a cellular trait has been quantified in many individual cells of the same type. This way, the observed distribution of the trait constitutes the phenotypic measure of individuals. We also assume that a genetic map is available and the individuals have been genotyped at marker positions on the map. The method we propose is based on three steps. First, a distance is computed for all pairs of individuals in order to quantify how much their phenotype differs. We chose the Kantorovich metric (also known as the Wasserstein distance or the earth-mover's distance) to measure this distance because, unlike the Kullback-Leibler divergence, it satisfies the conditions of non-negativity, symmetry and triangle inequality and, unlike the Hellinger distance, it does not converge to a finite upper limit when the overlap between distributions diminishes [27]. The Kantorovich metric can be viewed as the minimum energy required to redistribute one heap of earth (one  $f$ -function) into another heap (a second  $f$ -function). It has enabled developments in various fields, ranging from mathematics [28] to economy (the minimal transportation problem) [29,30] to the detection of states from molecular dynamics data [27]. The next two steps are inspired from methods used in ecology, where spatial distinctions between groups are often searched after determining distances between individuals [31,32]. In step 2 of our method, individuals are placed in a vectorial space while preserving as best as possible the distance between them (Fig 3A). This is achieved by multi-dimensional scaling, a dimension-reduction algorithm [33]. The third step is the genetic linkage test itself. At every genetic marker available, a linear discriminant analysis is performed to interrogate if individuals of different genotypic classes occupy distinct sectors of the phenotypic space (Fig 3B and 3C). The optimal choice of dimensionality is determined dynamically and a permutation test assesses statistical significance in the context of the corresponding degrees of freedom. Note that if the dimensions have been reduced to a single one, then canonical analysis is not needed: the phenotypic value of each individual has become a scalar and linkage can be performed by standard QTL mapping. Finally, scPTL linkage is scored using the Wilks' lambda statistics. Statistical inference is made using empirical  $p$ -values produced by permutations where the identities of individuals are re-sampled. The full procedure is described in details in the methods section.

## Detection of simulated scPTL

We first evaluated if our method could detect scPTL from simulated datasets. To do this, we considered a probabilistic single-cell trait governed by a positive feedback of molecular regulations. This is representative of the expression level of a gene with positive autoregulation. As depicted in Fig 4A, the employed model is based on three parameters. For each individual, a set of parameter values was chosen and single-cell values of expression were generated by stochastic simulations. We chose to simulate a scPTL that modified the expected values of the parameters so that the skewness of cellular trait distribution is affected. To do so, we considered a panel of individuals and their genotype at 200 markers evenly spaced every 5cM. Parameter values of each individual were drawn from Gaussians and the mean of these Gaussians depended on the genotype at the central marker. This defined two sets of phenotypes that are depicted by blue and red histograms in Fig 4B. A universal noise term  $\eta$  was added to introduce intra-genotype inter-individual variation which, in real datasets, could originate from limited precision of measurements or from non-genetic biological differences between individuals. For each of five increasing values of  $\eta$ , about 130 individuals were simulated.

We first scanned the generated dataset by QTL mapping, treating either the mean trait or its variance as the phenotype of interest. This way, the central scPTL locus was detected only when intra-genotype noise was null or very low (Fig 4C). This was anticipated because the mean and variance of the simulated trait values slightly differed between the two sets of



**Fig 3. Principle of scPTL mapping.** A cohort of multi-cellular individuals (or unicellular clones) with differing genotypes is used. For each individual (or clone), a cellular trait  $X$  is measured on a population of cells, and the observed distribution of  $X$  corresponds to the 'phenotype' of the corresponding individual. **A)** Kantorovich distances are computed for all pairs of individuals. The resulting distance matrix is used to place individuals in a multidimensional space. Proximity of individuals (grey and colored squares) in this space reflects comparable phenotypes (distributions in insets). **B)** Individuals are 'labeled' (blue vs. red) by their genotype at one genetic marker. **C)** A canonical discriminant analysis is performed to test if the genotype at the marker discriminates individuals in the phenotypic space. In the examples displayed, genetic linkage is significant at marker 2 but not at marker 1.

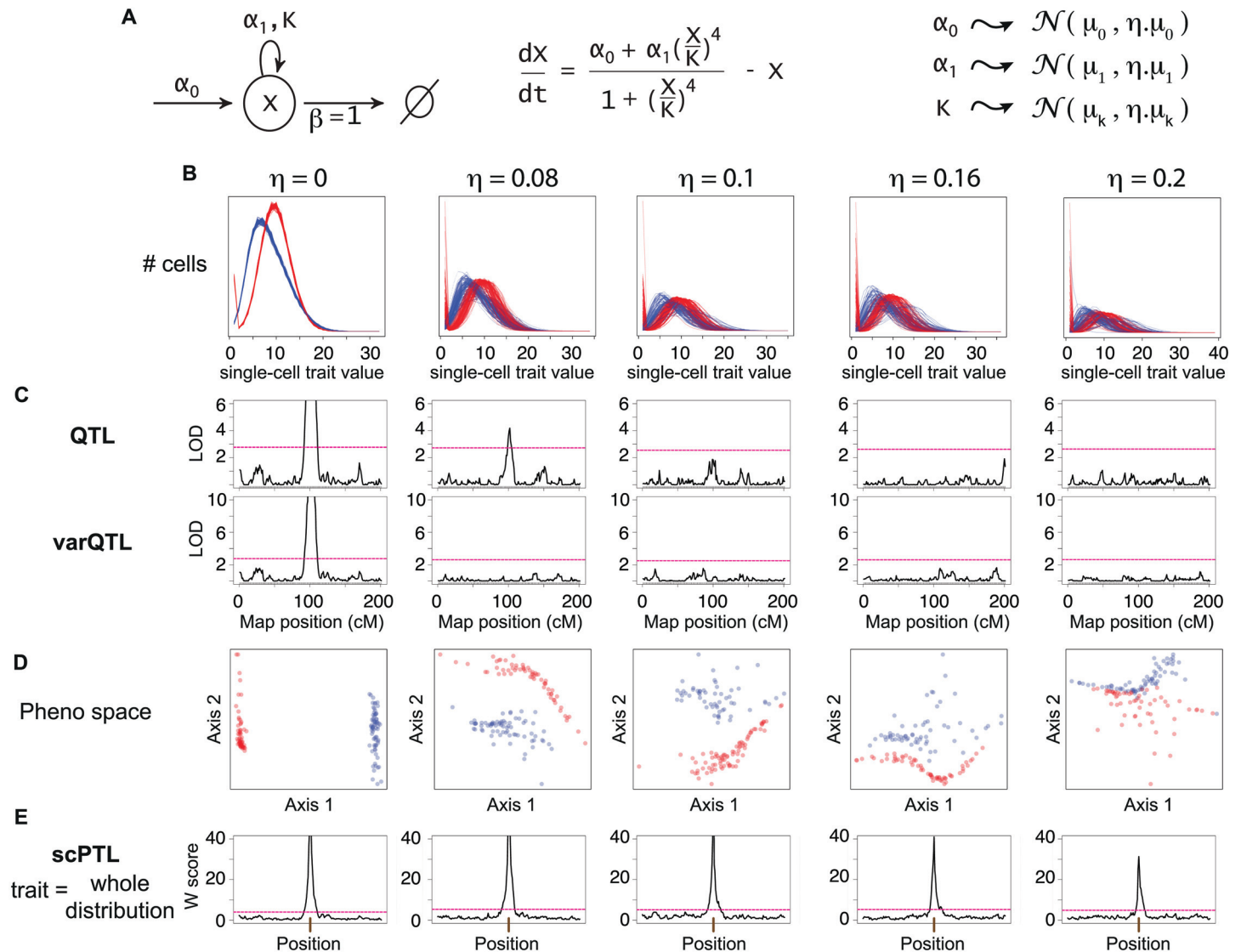
doi:10.1371/journal.pgen.1006213.g003

individuals. In contrast, our new method allowed to robustly detect the scPTL locus even in the presence of high (up to 20%) intra-genotype noise (Fig 4D and 4E).

### scPTL mapping of yeast cell morphology

The results described above using a simulated dataset suggest that the method can complement usual QTL mapping strategies. To explore if this was also the case when using real experimental data, we applied scPTL scans to the dataset of Nogami et al. [22] mentioned above (Fig 2A) where 220 single-cell traits were measured in about 200 cells from segregants of a yeast cross. We applied three genome x phenome scans, each one at FDR = 10%. Two consisted of QTL interval mapping and were done by considering either the mean cellular trait value of the population of cells or the coefficient of variation of the cellular trait as the population-level quantitative trait to be mapped. The third scan was done using the novel method described here to map scPTL. Significant linkages obtained from this scan are available in S1 Table. As shown in Fig 5, the three methods produced complementary results. We detected more linkages with the scPTL method than with the 2 QTL scans combined (71 vs. 61 traits mapped). This illustrates the efficiency of using the full data (whole distribution) of the cell population rather than using a summary statistic (mean or CV). In addition, we expected that a fraction of scPTL would match QTL, because QTL controlling the mean or CV of cellular traits are specific types of scPTL. This was indeed the case, with 67% of scPTL corresponding to loci that were detected by at least one of the two QTL scans. For 11 cellular traits, a locus was found by QTL or cvQTL mapping but it was missed by the scPTL scan. This illustrates that the methods have different power and sensitivity. Importantly, 22 cellular traits were associated to scPTL that were not



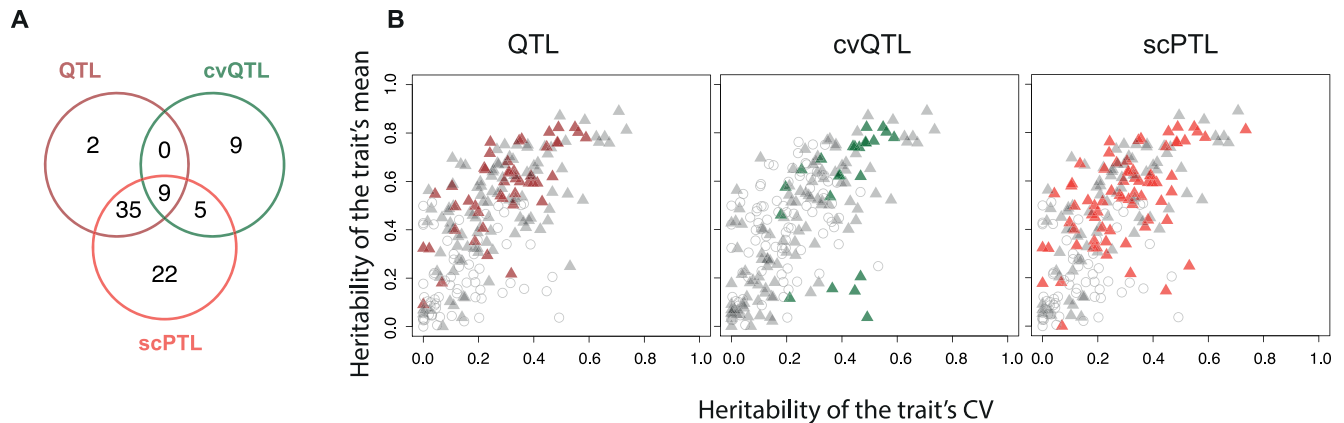


**Fig 4. Test on simulations.** **A**) A model of gene expression with positive feedback was used to simulate data (see [methods](#)). For each of ~130 distinct individuals, parameters  $\alpha_0$ ,  $\alpha_1$  and  $K$  were drawn from Gaussian distributions and then used to generate independent values of  $X$  in 10,000 cells of each individual. Mean values  $\mu_0$ ,  $\mu_1$  and  $\mu_k$  depended on the genotype of individuals at a locus located in the middle of a 200cM genetic map. Other sources of inter-individual variability were modeled by the extrinsic noise strength  $\eta$ . **B**) Distributions obtained (one per individual) at various values of  $\eta$ . Color: genotype at the locus controlling  $\mu_0$ ,  $\mu_1$  and  $\mu_k$ . **C**) QTL scans. For each individual, the mean (upper panels) or the variance (lower panels) of  $X$  were considered as quantitative traits and the map was scanned using interval mapping. Red dashed line: genome-wide significance threshold at 0.05. **D**) Coordinates of individuals (dots) in the phenotypic space obtained after computing Kantorovich distances and applying multi-dimensional scaling. Only the first two dimensions are shown. **E**) scPTL scan. At every marker position, linear discriminant analysis was performed. W score:  $-\log_{10}(\Lambda)$ , where  $\Lambda$  is the Wilks' lambda statistics of discrimination (see [methods](#)). Red dashed line: empirical genome-wide significance threshold at 0.05 (see [methods](#)).

doi:10.1371/journal.pgen.1006213.g004

detected by the QTL search, suggesting that some probabilistic effects may affect poorly the trait's mean or CV. Altogether, these observations highlight the complementarity of the different approaches and show that scPTL mapping can improve the detection of genetic variants governing the statistical properties of single-cell quantitative traits.

Examples of scPTL of yeast cellular morphology are shown in [Fig 6](#). One of the cellular traits measured was the distance between the center of the mother cell and the brightest point of DNA staining ([Fig 6A](#)). No QTL was found when searching genetic modifiers of the mean or

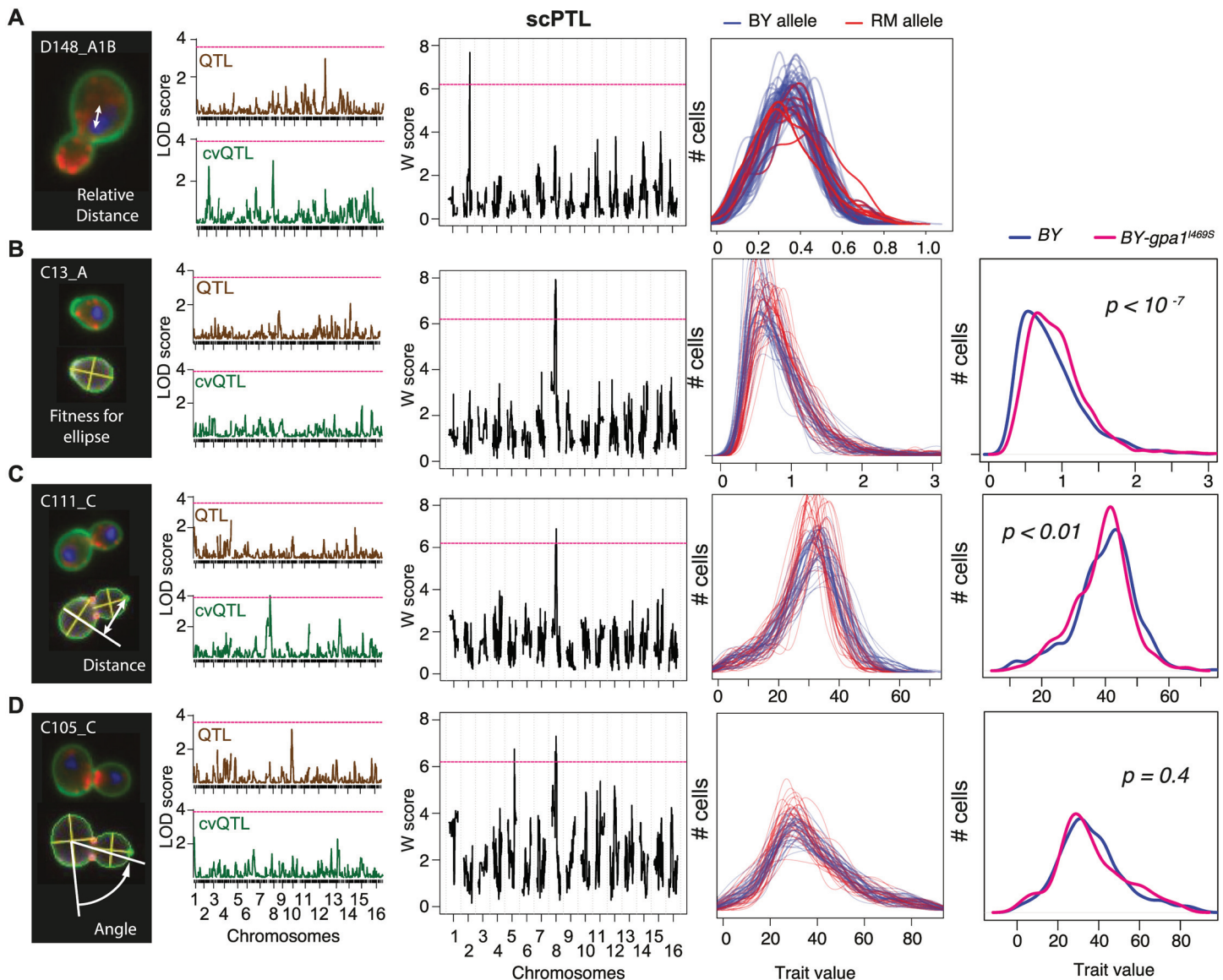


**Fig 5. Complementarity of scPTL and QTL mapping using experimental data.** The data of Nogami et al. [22] was used to perform genomic scans for QTL, cvQTL and scPTL. For scPTL mapping, we used both the first-axis only and multiple dimensions of the phenotypic space (see methods), and the results were pooled. **A**) Venn diagram showing the number of traits for which a significant locus was found in the genome by each method (each at FDR = 10%). **B**) Same representation as Fig 2B showing the traits successfully mapped by each method. The traits that passed the heritability filter (see methods) and were considered for mapping are shown as triangles, and colored if mapping was successful.

doi:10.1371/journal.pgen.1006213.g005

CV of this trait, but a significant scPTL was mapped on chromosome II. When displaying trait distributions, it was apparent that segregants carrying the BY genotype at the locus had reduced cell-cell variability of the trait as compared to segregants having the RM genotype (Fig 6A, right panel). Consistently, a small cvQTL peak was seen on chromosome II, although this peak did not reach genome-wide statistical significance. This trait, which relates to the statistical properties of DNA migration during the early phase of cell division, provided a biological example where scPTL scan identified a genetic modulator of cell-to-cell variability that was missed by the QTL approach.

Three other traits were of particular interest because they mapped to a position on chromosome VIII where a functional SNP was previously characterized in this cross. This SNP corresponds to a non-synonymous I->S mutation at position 469 in the  $G_{\alpha}$  protein Gpa1p. It targets a domain that is essential for physical interaction with pheromone receptors Ste2p and Ste3p [34,35]. In the presence of pheromone, Gpa1p is released from the receptor and triggers a signalling cascade of molecular response that causes cell-cycle arrest and cell elongation (a process called 'shmooing'). In the absence of pheromone, improper binding of Gpa1p<sup>I469S</sup> to the receptor causes residual activation of the pathway in the BY strain, as seen by transcriptomic profiling [36], which explains why BY cells are more elongated [24] and proliferate slower [37] than RM cells. Here we saw that this locus is a scPTL, but not a QTL, of the degree to which cells are elliptical (Fig 6B). Displaying the distributions of this trait in each segregant revealed a remarkable amount of variability between the segregants, and that the BY allele at the locus corresponded to a modest reduction of the trait value as compared to the RM allele (sharper mode at slightly lower value). To see if this was due to the GPA1<sup>I469S</sup> mutation, we examined the data from a BY strain where this mutation was cured [22]. Remarkably, the single amino-acid substitution caused a mild but statistically significant redistribution of the trait values (Fig 6B). This change was comparable to the difference seen among the segregants, demonstrating the causality of the GPA1<sup>I469S</sup> SNP. Another trait, corresponding to the distance between the bud tip and the short axis of the mother cell, also mapped to this locus, with the RM allele associated to greater cell-cell variability, and data from the GPA1<sup>I469S</sup> allele-replaced strain validated this SNP as the causal polymorphism (Fig 6C). These observations suggest that either the residual activation of the pathway in absence of pheromone is not uniform among



**Fig 6. Mapping single-cell probabilistic traits of cellular morphology.** For each of 59 recombinant BYxRM yeast strains, four quantitative traits were measured on ~200 individual cells [22]. From left to right: description of the trait; results from QTL scans applied to the mean (brown) or coefficient of variation (green) of the trait; results from scPTL scan (pink dashed line: genome-by-phenome significance threshold at FDR = 10%; single-cell trait density functions computed from the data, where each line corresponds to a recombinant strain (color: genotype at scPTL); when relevant, single-cell trait density functions of nearly-isogenic BY strains differing for one non-synonymous SNP in the *GPA1* gene are shown ( $p$ : statistical significance of the corresponding two-sample Kolmogorov-Smirnov test). **A)** Trait D148\_A1B is the distance between the nuclear brightest point and the mother center, relative to the mother size. An scPTL was detected on chromosome 2. **B)** Trait C13\_A is the fitness of the cell outline to the best adjusted ellipse. **C)** Trait C111\_C is the distance between the bud tip and the extension of the mother short axis. **D)** Trait C105\_C corresponds to the position of the budding site. It is the angle between the long axis of the mother cell and the line defined by the mother center and the middle point of neck.

doi:10.1371/journal.pgen.1006213.g006

BY cells, or the proper inactivation of the pathway is not complete in all RM cells. This, and the fact that the mutation does not prevent BY cells from proliferating (as compared to pheromone-arrested cells), indicate that the detachment of Gpa1p<sup>1469S</sup> from the receptor is a rare event that has probabilistic effects on the cellular phenotype. Further investigations based on biochemistry, dynamic recording of individual cells and stochastic modelling are needed to understand how variation in binding affinity accounts for this effect. The results described here

illustrate that scPTL scans can identify individual SNPs that modify single-cell trait distributions without necessarily affecting the trait mean.

Finally, another trait corresponding to the angle of bud site position mapped to two scPTL loci and no QTL. One of these loci contained the GPA1 gene on chromosome VIII. Although the phenotype of bud site selection is not related to 'shmooing', we examined if the GPA1<sup>1469S</sup> SNP was involved and found that it was not: the allele-replaced strain did not show a different trait distribution than its control (Fig 6D). Thus, other genetic polymorphisms at the locus should participate to the statistical properties of cellular morphology, by affecting the position of budding sites.

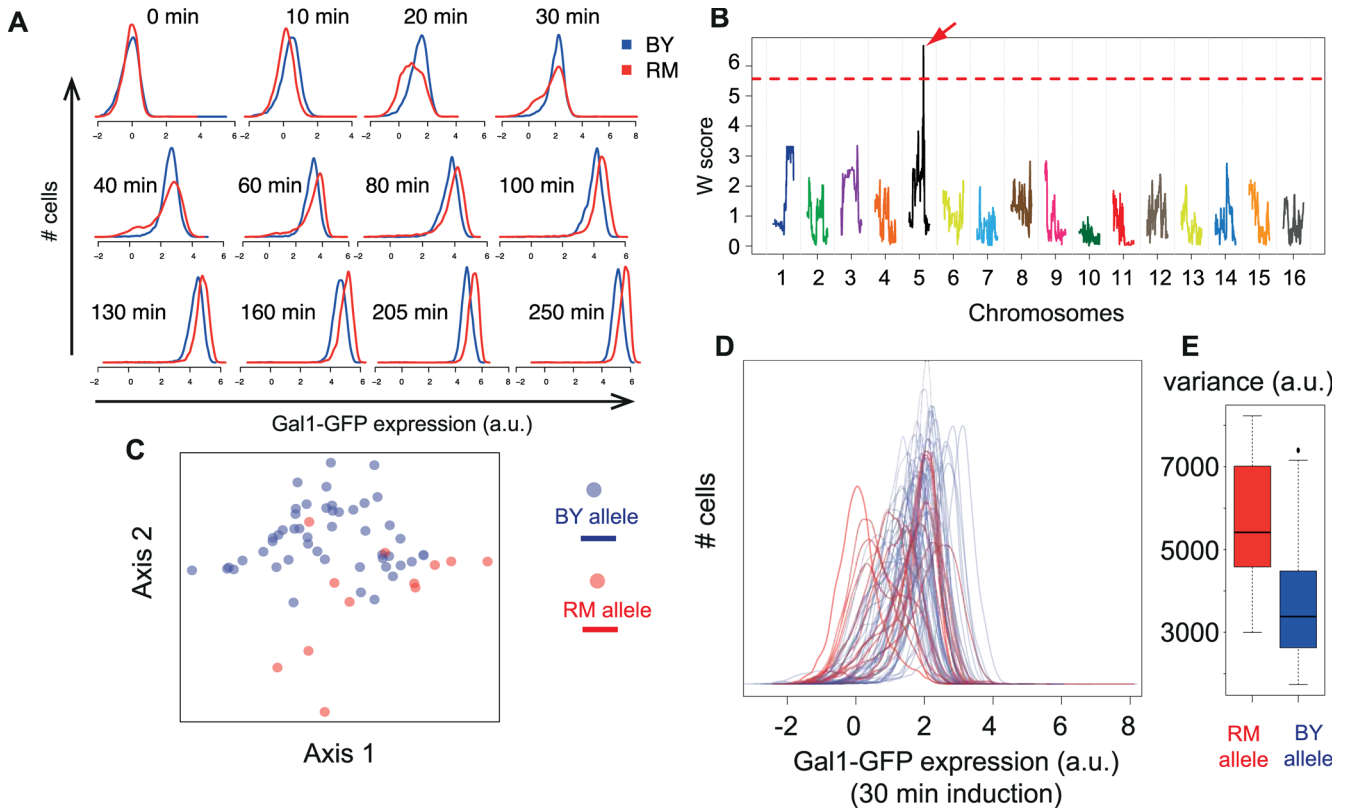
### scPTL scanning detects a new yeast locus that modulates cell-cell variability of the transcriptional response to galactose

We then explored if scPTL scanning could provide new results when applied to a molecular system that had been extensively characterized by classical genetics. The system we chose was the yeast GAL regulon which, in addition to be one of the best described regulatory network, presented several advantages. Natural strains of *S. cerevisiae* are known to display differences in its regulation [38,39] and the transcriptional response of cell populations can be tracked by flow cytometry. This provides data from large numbers of cells and therefore a good statistical power to compare single-cell trait distributions. In addition, acquisitions on many genotypes are possible using 96-well plates. We reasoned that if features of the cell population response segregate in the BY x RM cross (described above for morphology), then scPTL scanning might identify genetic variants having non-deterministic effects on the regulation of GAL genes.

We first compared the dynamics of transcriptional activation of the network in the two strains BY and RM. This was done by integrating a P<sub>Gall</sub>-GFP reporter system in the genome of the strains, stimulating them by addition of galactose in the medium, and recording the response by flow cytometry. As shown in Fig 7, both strains responded and full activation of the cell population was reached after ~2 hours of induction. Interestingly, remarkable differences were observed between the two strains regarding the distribution of the cellular response. The BY strain showed a gradual increase of expression through time that was relatively homogeneous among the cells (unimodal distribution with relatively low variance), whereas the RM strain showed elevated cell-cell heterogeneity at intermediate activation time points (higher variance, with fraction of non-induced cells). This suggested that genetic polymorphisms between the strains might control the level of heterogeneity of the cellular response at these intermediate time points.

We sought to map one or more of these genetic factors. To do so, we acquired the response of 60 meiotic segregants of the BY x RM cross. Using the data collected at each time point, we scanned the genome for scPTL of the reporter gene expression level using the novel genome-scan method described above. The procedure identified a locus on chromosome V position 350,744 that was highly significant (genome-wide  $p$ -value < 0.001) at 30 minutes post induction, the time at which heterogeneity markedly differed between the BY and RM strains (Fig 7B and 7C). The locus was also significant at times 20 min ( $p$  < 0.005) and 40 min ( $p$  < 0.005) post induction.

Visualizing the distributions of single-cell expression levels at 30 minutes revealed that the RM and BY genotypes at this locus corresponded to high and low cell-cell heterogeneity, respectively (Fig 7D and 7E). Thus, this locus explains, at least in part, the different levels of heterogeneity observed between the parental strains. It should therefore also be detected as a varQTL or cvQTL. This was indeed the case: the LOD score linking the locus to the variance of expression was 4.5 and reached statistical significance ( $P$  = 0.005). Importantly, the scPTL was



**Fig 7. Detection of a scPTL for the cellular response to galactose.** **A)** Time-course flow cytometry acquisitions of the response to galactose in strains BY and RM. Cells were cultivated in raffinose 2% and were shifted to a medium containing Raffinose 2% and Galactose 0.5%. After the indicated time, cultures were fixed with paraformaldehyde and analysed by flow cytometry. Histograms correspond to the fluorescent values obtained on cells gated for cell-size (see [methods](#)). **B)** Genome scan for scPTL affecting the response after 30 minutes induction. Data similar to panel A was generated for 60 segregants, and the histograms obtained at 30 min post-induction (shown in D), together with the genotypes from [66] were used for scPTL mapping using the multi-dimensions method. The linkage profile (W score) obtained when retaining the first two dimensions is shown, colored by chromosome. Dotted line: significance threshold at genome-wide  $p$ -value  $< 0.005$ . Arrow: significant scPTL on chromosome 5. **C)** Two-dimensional coordinates of the 60 segregants in the phenotypic space (30 min induction time). color: genotype at the scPTL locus. **D)** Phenotypes (histograms of single-cell expression value) of the 60 segregants after 30 min induction, colored by the genotype of the segregants at the scPTL locus. a.u.: arbitrary unit. **E)** Boxplot summarizing the variance of histograms from panel D grouped by the genotype at the scPTL locus.

doi:10.1371/journal.pgen.1006213.g007

not a QTL: the locus genotype did not correlate with the mean level of expression of the population of cells (LOD score  $< 2.8$ ).

When surveying the genomic annotations of the locus [40], we realized that it contained no obvious candidate gene that would explain an effect on the heterogeneity of the response (such as genes known to participate to the transcriptional response). One potentially causal gene was DOT6, which encodes a poorly characterized transcription factor that was shown to shuttle periodically between the cytoplasm and nucleus of the cells in standard growth conditions [41]. Given that i) the shuttling frequency of such factors can sometimes drive the response to environmental changes and ii) numerous non-synonymous BY/RM genetic polymorphisms were present in the gene, we constructed an allele-replacement strain for DOT6 and tested if the gene was responsible for the scPTL linkage. This was not the case. Strains BY and BY-DOT6<sup>RM</sup> (isogenic to BY except for the DOT6 gene which was replaced by the RM allele) displayed very similar transcriptional responses at intermediate times of induction (S1 Fig). Fine-mapping of the locus and a systematic gene-by-gene analysis are now needed to precisely identify the polymorphisms involved. By highlighting a novel genetic locus modulating cell-cell

variability of the transcriptional response to galactose, our results show that scPTL scanning can provide new knowledge on the fine structure of a well-studied system.

## Discussion

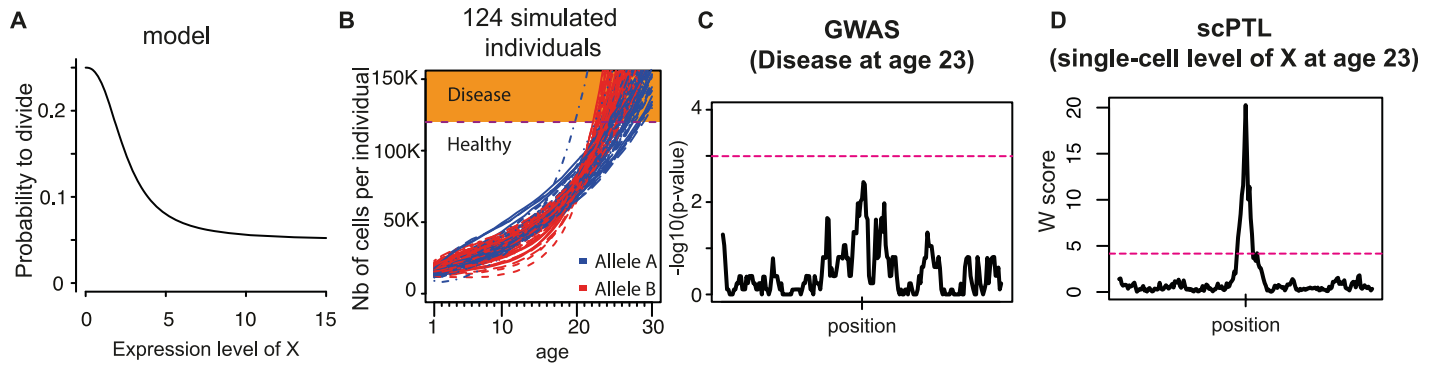
We have developed a novel method to scan genomes for genetic variants affecting the probabilistic properties of single-cell traits. We validated the method using data from colonies of a unicellular organism, which constitutes a first step before transferring the method to multicellular organisms. Our approach extends the usual genetic analysis of quantitative traits both conceptually and methodologically: by incorporating large samples of phenotypic values at the cellular scale, variants that have probabilistic effects can be detected and their possible contribution to trait heritability at the macroscopic (multicellular) scale can be investigated.

### scPTL and the genetic predisposition to disease

When considering macroscopic phenotypes, it is important to distinguish the situations where scPTL mapping is biologically relevant from those where it is not. The determinants of human height, for example, act via countless cells, of multiple types, and over a very long period of time (~ 16 years). In such cases, the macroscopic trait results from multiple effects that are cumulated and considering the probabilistic individual contribution of specific cells is inappropriate. Similarly, many tissular traits heavily rely on communications between cells and probabilistic changes in a few may not affect the collective output of the cell population. In contrast, a number of macroscopic traits can be affected by particular events happening in rare cells or at a very precise time (see below). In these cases, studying the probabilities of a biological outcome in the relevant cells or of a molecular event within the critical time interval can provide invaluable information on the emergence of the macroscopic phenotype, and scPTL mapping then becomes relevant.

A striking example of such traits is cancer. Genetic predisposition is conferred by variants affecting somatic mutation rates and these loci are special cases of scPTL: the cellular trait they modify is the amount of *de novo* mutations in the cell's daughter. These variants have classically been identified by genetic linkage of the macroscopic trait (disease frequency in families and cohorts), and their role on the maintenance of DNA integrity was deduced afterwards by molecular characterizations. For a review on the genetics of cancer syndrome predisposition, see [42,43]. scPTL mapping is also relevant to the non-genetic heterogeneity of cancer cells which was shown to be associated with tumour progression [7,8] and treatment efficiency [9–11]. Genetic loci changing the fraction of problematic cells are likely modulators of the prognosis. If the functional properties (expression level, phosphorylation status, subcellular localization) of a key molecular player, such as a critical tumor-suppressor gene, can be monitored in numerous individual cells, then scPTL mapping, as presented here, may help identify genetic factors that modulate the activity of this gene in a probabilistic manner. Once identified, the association of these loci with the macroscopic phenotype can then be tested directly, avoiding at least partly the statistical challenges of whole-genome scans.

To illustrate this, we considered an idealized case where three scales are bridged: at the molecular level, a scPTL affects the expression of a protein X (same regulation as in Fig 4); at the cellular level, cells have higher probability to divide if their level of X is low (Fig 8A); and at the whole-organism level, disease appears if too many cells are present. Using a stochastic model of this scenario, we simulated a cohort of individuals and recorded the state and number of cells in every individual over time (Fig 8B, see methods for details). Disease appeared in all individuals, between age 22 and 29. Using the data at age 23, we compared the power of GWAS and scPTL mapping. For GWAS, the trait of individuals was whether they had declared the



**Fig 8. The added value of scPTL mapping.** A simple model is considered where a population of cells evolves, with a death rate that is constant and a division rate that depends on the intra-cellular concentration of a tumor-suppressor protein X (A). The population of cells is considered to be pathogenic if it exceeds 120,000 cells (over-proliferation). The expression of X follows the model described in Fig 4, with extrinsic noise strength  $\eta = 0.16$ . B) Time-evolution of the number of cells in 124 simulations corresponding to 124 distinct individuals (see methods). Each line represents one individual. Color: same genotypic groups as in Fig 4, which correspond to distinct sets of parameters for the regulation of X, which are linked to an scPTL located in the middle of a chromosome. Disease onset is earlier in group B (red curves). C) Using the macroscopic data (disease vs. healthy, 124 individuals) from panel B at age 23, the simulated disease-modifying locus was not detected by classical association mapping. Linkage was searched with an exact Fisher test. Dashed line: genome-wide 5% significance threshold determined by permutations. D) Using single-cell data (level of X in 10,000 cells, 124 individuals) from panel B at age 23, the locus was detected by scPTL mapping. Dashed line: genome-wide 5% significance threshold determined by permutations.

doi:10.1371/journal.pgen.1006213.g008

disease or not. For scPTL mapping, the trait was the expression level of X in 10,000 of their cells. As expected by the moderate effect on disease frequency, GWAS failed to detect the locus (Fig 8C). In contrast, scPTL detection was highly significant from the same cohort of individuals (Fig 8D). Importantly, although not significant genome-wide, the GWAS score at the locus had a nominal  $p$ -value lower than 0.01 (Fig 8C). The locus would therefore be considered significant if it had been the only one tested. This illustrates the added value of scPTL mapping: while keeping cohort size constant, it can highlight candidate loci of the genome that can then be tested individually for association to the disease. This power clearly results from i) additional traits (cellular ones) that are included before scanning the genome and ii) relaxation of multiple-testing correction when testing association to disease. Note that other system genetics methods, such as expression QTL (eQTL) mapping, improve power in a similar way: they highlight relevant candidates via the addition of intra-individual traits (molecular ones) [44]. Note also that recruiting large cohorts remains important: Methods detecting scPTL and eQTL can improve genetic mapping but their detection power remain strictly dependent on the number of individuals available in the study.

### What cellular trait should be quantified?

In real studies, external knowledge is needed on the link between the cellular trait and the disease: what single-cell trait should be measured? Can it be measured in a sufficiently large number of cells? If a reporter system of *de novo* mutations, for example based on the intracellular distribution of a fluorescently tagged repair protein [45,46], can be introduced in a relevant and large population of cells, then the high number of cell measurements may allow to detect loci that modify even slightly the mutation rate. For non-genetic features of problematic cells, choice of the trait can be driven by investigations at the molecular level, such as stochastic profiling [47], and at the cellular level, such as recording the response of cell populations to treatment or differentiation signals [9,10]. For example, the distribution of the biomarker JARID1B (a histone demethylase) in populations of melanoma cells is indicative of an intra-clonal heterogeneity that is important for tumour progression [7], biomarkers CD24 and CD133 can

distinguish rare cells that persist anti-cancer drug treatments [10] and multiplexed markers of signalling response can reveal patterns of population heterogeneity that predict drug sensitivity [48]. When relevant markers are not known, a possibility is to screen for them using stochastic profiling [47]. This method interrogates the transcriptomic variability between pools of few cells in order to identify transcripts displaying elevated cell-to-cell variability in specific biological contexts. It allowed the discovery of two molecular states of extracellular matrix-attached cells that can be distinguished by the *jun D proto-oncogene* and markers of TGF- $\beta$  signalling [8]. Such markers of isogenic cellular subtypes may allow the development of scPTL mapping in humans.

An important statistical requirement to identify scPTL is the abundance of cells on which the probabilistic trait is quantified. For human studies, peripheral blood offers access to many cells but, unfortunately, many internal organs do not. This requirement also implies using technologies where the throughput of quantitative acquisitions is high. This is the case for flow-cytometry and, although at higher costs, for high-content image analysis [14,15] and digital microfluidics [17,18]. For these practical reasons, it is possible that mouse immunological studies will help making progress in mammalian scPTL mapping. For example, the work initiated by Prince et al. [49] describing pre- and post-infection flow-cytometry profiles of F2 offsprings from different mouse strains may provide an interesting pilot framework.

### scPTL mapping and developmental instability

The interest of scPTL mapping is not restricted to cancer biology. Developmental processes and cellular differentiation are also vulnerable to mis-regulations happening in few cells or during short time intervals. Their macroscopic outcome can therefore be affected by probabilistic events at the cellular scale. For example, stochastic variation in the expression of the stem cell marker Sca-1 is associated with different cellular fates in mouse hematopoietic lineages [50], suggesting that genetic factors changing this stochastic variation may impact blood composition. Similarly, embryonic stem cells co-exist in at least two distinct molecular states that are sensitive to epigenetic and reprogramming factors [51]. Genetic variants modulating these factors may change the statistical partitioning of these states. Two observations made on flies remarkably support the existence of natural genetic factors altering developmental processes in a probabilistic manner. The first one is the fact that high levels of fluctuating asymmetry can be fixed in a wild population of *D. melanogaster* under artificial selection [52]. The second one comes from a comparative study of *Drosophila* species [53]. Embryos of *D. santomea* and *D. yakuba* display high inter-individual variability of expression of the signal transducer pMad at the onset of gastrulation, as compared to *D. melanogaster* embryos. This increased variability was attributed to a reduced activity of the homeobox gene *zerknüllt* thirty minutes before this stage. Very interestingly, it is accompanied by phenotypic variability (inter-individual variance of the number of amnioserosa cells) in *D. santomea* but not in *D. yakuba*. These and other examples [54] illustrate how developmental variability and phenotypic noise can evolve in natural populations. Applying scPTL mapping may allow to dissect the genetic factors responsible for this evolution.

### Methods for scPTL mapping

Our new method based on the Kantorovich distance is not the only one by which scPTL can be identified. Applying classical QTL mapping to summary statistics of the cellular traits can also be efficient. We emphasize that the two approaches are complementary. For example, our method missed to detect linkage for 9 yeast morphological traits for which cvQTL scans were successful, but it detected several significant scPTL that were missed by the QTL-based



approach (Figs 5 and 6). Second, we observed that scPTL detection was often efficient when the mean value of cellular trait differed among genotypic categories. As shown on Fig 5B, traits successfully mapped tended to display high heritability of the mean. Thus, after a scPTL is detected, it is necessary to examine the effect on the trait distributions and to determine if it is a QTL or not. Third, alternative ways of mapping scPTL are open and may prove more appropriate in some contexts. For example, if a cellular trait becomes preoccupying when it exceeds a certain threshold value, then the fraction of cells above this threshold can be used as a macro-trait to be mapped by QTL analysis. This way, the focus is made on the relevant aspect of the cellular trait, avoiding variation in other parts of the distribution. We therefore recommend conducting Kantorovich-based scPTL mapping in addition to classical methods and not as a replacement strategy.

### Future statistical improvements in scPTL mapping

While the principle of genotype-phenotype genetic linkage dates back to several decades ago, the statistical methods that test for linkage are still being improved, especially regarding multi-loci interactions or population structure corrections [55,56]. The present study provides a priming of a generic scPTL mapping approach (exploiting thousands of single-cell trait values) and demonstrates its feasibility and potential (new loci were detected). Since it is new, we anticipate that it will also evolve in the future. It is currently based on three steps: (i) computing pairwise distances between individuals by using the Kantorovich metric, (ii) using the resulting distance matrix to construct a relevant phenotypic space and (iii) testing for genetic linkage by LDA. A number of methodological considerations can be made in anticipation to future developments and applications.

Estimating the proportion of variance explained by scPTL is not straightforward: the '*captured variance*' as quantified by the eigenvalues of the LDA is not the same as the '*explained variance*' which must be re-computed by regression; and if linearity of the data is questionable, the method remains a useful tool if it detects scPTL but interpreting variance proportions need justifications.

A phenotypic space can be constructed by alternative ways that do not require the Kantorovich metric. For example, we considered representing individuals in a "space of moments", where the coordinate of every individual on the  $i$ -th axis was the  $i$ -th moment of the cellular trait distribution associated to this individual. We applied this to the yeast morphological data and we searched for genetic linkages by linear discriminant analysis as described above. This approach detected many significant scPTL but we encountered a difficulty that was avoided by our Kantorovich-metric based method. When searching for significant linear discrimination, the dimensionality of the phenotypic space is important. At high dimensionality, discriminant axes are more likely to be found. This improves detection in the actual data but at the expense of increasing the degrees of freedom and therefore the false positive hits estimated from the permuted data. In a "space of moments", the properties of the single-cell trait distributions are very important because they define which axis (moments) are relevant to separate individuals. Keeping the 4-th axis may be crucial even if all individuals have very similar first, second or third moments. Choosing the appropriate dimension for LDA is then arbitrary and it becomes difficult to keep a good detection power while still controlling the FDR. In fact, applying QTL mapping on the 3rd and 4th moments of all traits was fruitless because the FDR could not be controlled at the genome-by-phenome scale. This issue is avoided in the case of Kantorovich distances because multi-dimensional scaling can be applied without normalization and the axes of the phenotypic space are ranked by descending order of their contribution to the inter-individual differences. The 4-th axis, for example, contributes less than the first

three axes to the separation of individuals in the space. If keeping the 4-th dimension prior to LDA is beneficial for linkage, then keeping the first three axes is also highly relevant, and this is true regardless of the properties of the single-cell trait distributions. We found this very useful: our algorithm adds dimensions one by one and evaluates the benefit of each increase (see [methods](#)).

There are at least three lines along which our method may be further improved. First, LDA is only appropriate if genotypic categories can be distinguished along linear axis. If individuals in the phenotypic space are separated in non-linear patterns, other methods such as those based on kernel functions [57] may be more appropriate. Second, we propose to compute confidence intervals of scPTL position by bootstrap, following a method sometimes applied to QTL positions [58]. As expected, resampling not only affected scPTL position but also the optimal dimensionality retained (S2A Fig). A deeper investigation of the simultaneous variation of these two outputs could help improve the precision of mapping. And third, single-cell data acquisitions often generate multiple trait values for each individual cells. This is the case for morphological profiling as in the dataset we used here, but also for gene expression [59] or parameters describing the micro-environment of the cells [60]. It would therefore be interesting to search for scPTL affecting multiple cellular traits simultaneously instead of treating cellular traits one by one. A multidimensional analysis could be performed in order to extract a set of informative meta-traits, such as principal components or representative medoids and scPTL of these meta-traits could be searched using our method. This dimension-reduction approach would benefit from the redundant information available from correlated traits (e.g. the perimeter of a cell and its area are two measurements of its size), but the biological interpretation of a probabilistic effect on a meta-trait may not be straightforward. Alternatively, one might want to identify scPTL affecting the joint probability distribution of multiple cellular traits. In this case, a natural extension of our method would be to compute Kantorovich distances between multivariate distributions. However, the Kantorovich metric cannot be easily computed for more than two marginals (i.e. cellular traits in our case). In fact, its existence as a unique solution to the multi-dimensional transportation problem was itself a subject of research [61]. A possible alternative could be to compute a Euclidean distance in the "space of moments" mentioned above and then apply multi-dimensional scaling.

Furthermore, although our study was focused on probability density functions, steps (ii), constructing the phenotypic space, and (iii), testing for genetic linkage, could in principle be applied to other types of functions, provided that a relevant metric estimating the dissimilarity between such functions exists. This could be interesting in the case of function-valued traits, such as speech sound or other time-series functions. The evolution of these functions is being studied using phylogenetic methods that present challenging statistical issues [62,63]. Extending our approach to such functions may open the possibility to study them from a (complementary) quantitative genetics angle.

Finally, we can anticipate that gene-gene and gene-environment interactions also shape the probability density function of cellular traits. Our results on the activation of the yeast galactose network remarkably illustrate this: the effect of the scPTL on chromosome V is apparent only transiently, and in response to a change of environmental conditions. It is tempting to extrapolate that signaling pathways in plants and animals may be affected by scPTL that act at various times and steps along molecular cascades.

In conclusion, our study provides a novel method that can detect genetic loci with probabilistic effects on single-cell phenotypes, with no prior assumption on their mode of action. By exploiting the power of single-cell technologies, this approach has the potential to detect small-effect genetic variants that may underlie incomplete trait penetrance at the multicellular scale.

## Methods

### Stochastic modeling of a positively auto-regulated gene

Single-cell gene activity was modeled by a stochastic variable  $X$  that represented the number of proteins in one cell at a given time. Under the model, the dynamics of  $X$  is controlled by two processes: (1) protein production with rate  $\alpha$  and (2) protein degradation with rate  $\beta$ . We assume that the gene is positively auto-regulated by a 4-mer complex, meaning that  $\alpha$  is an increasing function of  $X$  with a typical Hill-like shape

$$\alpha = \frac{\alpha_0 + \alpha_1 \left(\frac{X}{K}\right)^4}{1 + \left(\frac{X}{K}\right)^4}$$

with  $\alpha_0$  the leaky production rate in absence of  $X$  4-mers at the promoter,  $\alpha_1$  the production rate in presence of 4-mers, and  $K$  the dissociation constant of the 4-mer. We set  $\beta = 1$ , which corresponds to scaling time units. The dynamics of the mean value of  $X$  in a population of isogenic cells follows the equation shown in Fig 4A. To obtain the probability distribution of  $X$ , we performed exact stochastic simulations of the chemical system defined by the two reactions rate  $\alpha$  and  $\beta$ , using the Stochastic Simulation Algorithm [64].

To generate two groups of individuals, we assumed that the set of parameters ( $\alpha_0, \alpha_1, K$ ) was controlled by one locus that could exist in two alleles (A and B) with mean values ( $\mu_0^{A/B}, \mu_1^{A/B}, \mu_K^{A/B}$ ) and, for simplicity, that the individuals were haploids. To account for sources of inter-individual variability within genotypic groups, the values of the parameters for one individual were drawn from normal distributions of mean values  $\mu_0^{A/B}, \mu_1^{A/B}$  and  $\mu_K^{A/B}$  and of standard deviations  $\eta\mu_0^{A/B}, \eta\mu_1^{A/B}$  and  $\eta\mu_K^{A/B}$  where  $\eta$  represented the strength of inter-individual variability.  $\eta$  was assumed to be the same for A and B alleles. Values were:  $\mu_0^A = 6.3, \mu_0^B = 0.1, \mu_1^A = 12, \mu_1^B = 10, \mu_K^A = 10$  and  $\mu_K^B = 1.6$ .

### Genetic heritabilities of cellular traits' mean and CV

All statistical analysis were done using R (version 3.1.2) [65]. The data from Nogami et al. [22] consisted of 220 traits, acquired on >200 cells per sample. Note that most traits are related to one of three division stages. Each trait was therefore measured on a subset of cells of the sample (less than 200). There were nine samples of the BY strain, nine of the RM strain, and three of each of 59 segregants of the BY x RM cross. For each trait, we computed the genetic heritabilities of the mean and CV as follows. The mean and CV of the cellular trait in each sample were computed, leading to two scalar values per sample that we call macro-traits hereafter (to distinguish them from the single-cell values). The broad-sense genetic heritability of each macro-trait was  $H^2 = (var_T - var_E) / var_T$ , with  $var_T$  and  $var_E$  being the total and environmental variance, respectively. For Fig 2B, we estimated  $var_T$  by randomly choosing one of the three replicate sample of each segregant and computing the variance across these 59 values. This was repeated 100 times and the estimates were averaged. Our estimate of  $var_E$  was the pooled variance of  $var_{BY}, var_{RM}, var_{Seg1}, var_{Seg2}, \dots, var_{Seg59}$  which were the between-replicates variance of each strain. Confidence intervals on  $H^2$  values were computed by bootstrapping the strains. For the filtering step prior to linkage,  $H^2$  was computed slightly differently in order to be consistent across mapping methods (see below).

### scPTL mapping of yeast morphological traits

We first normalized the distributions as densities (division of all bin counts by half the total number of cells). Following [27], we then computed the Kantorovich distance between two distributions  $f_1$  and  $f_2$  as the area under the absolute value of the cumulative sum of the difference

between the two distributions:

$$KD(f_1, f_2) = \int_{-\infty}^{+\infty} \left| \int_{-\infty}^x f_1(t) - f_2(t) dt \right| dx$$

Multi-dimensional scaling of the resulting distance matrix was then performed using the R function `cmdscale()` from the `stats` package. The number of dimensions retained (*ndim*) was the number of eigenvalues exceeding the expected value under the hypothesis of no structure in the data (i.e. mean of all eigenvalues, Kaiser criterion). We computed the heritability of each yeast morphological trait in this multidimensional space. This was done as above for one dimension, by computing the total variance of the data, and estimating the environmental variance from the replicated experiments made on the parental strains. For 147 traits, heritability was greater than 0.5 and `scPTL` were searched. Details on how these steps were implemented in R are described in [S1 Methods](#), and the code is available in the open source `ptlmapper` R package (<https://github.com/fchuffar/ptlmapper>).

The yeast genotypes we used were from Smith and Kruglyak [66]. For the morphological traits, we pooled triplicates together in order to increase the number of cells per sample. The data then corresponded to 220 traits, measured on >600 cells per sample, with 3 BY samples, 3 RM samples, and 1 sample per segregant. We scanned the genetic map with two methods. First, we considered the coordinates of each segregant on the first axis of the multi-dimensional scaling, and we considered this coordinate as a quantitative trait that we used for interval mapping using `R/qtl` [67]. Secondly, we applied a linear discrimination analysis (LDA) on the phenotypes data, using the genotype at every marker as the discriminating factor. An important issue in this step is the multidimensionality of the data: axis 2, 3 and more may contain useful information to discriminate genotypic groups, but if too many dimensions are retained, a highly-discriminant axis may be found by chance only. To deal with this issue, we evaluated the output of LDA at all dimensions *d* ranging from 2 to *ndim*. For each value of *d*, we applied LDA at every marker position and we recorded the Wilks' lambda statistics:

$$\Lambda = \prod_{j=1}^d \frac{1}{1 + \lambda_j}$$

where  $\lambda_j$  was the *j*-th eigenvalue of the discriminant analysis. Low values of this statistics allow to reject the null hypothesis of no discrimination by the factor of interest [68] which, in our case, is the genotype. We defined a linkage score (W score) as:

$$W = -\text{Log}_{10}(P)$$

where *P* is the *p*-value of the Wilk's test (deviation of  $\Lambda$  from the *F*-statistics with relevant degrees of freedom). Note that *P* is not interpreted directly as a significance value for linkage (see the permutation test below).

We then quantified how much the best marker position was distinguished from the rest of the genome by computing a Z-score:

$$Z = \frac{W_{best} - \langle W \rangle}{\sigma_w}$$

where  $W_{best}$ ,  $\langle W \rangle$  and  $\sigma_w$  were the highest, the mean and the standard deviation of all W scores found on the genome, respectively. Finally, we chose the dimension that maximized this Z-score (i.e. dimension where the linkage peak had highest contrast). Very importantly, the same degrees of freedom (exploration of the results at various dimensionalities) were allowed

when applying the permutation test of significance (see below). The distribution of the dimensionalities retained for the morphological traits is shown in [S2B Fig](#). Additional details are provided in [S1 Methods](#) and the code is available in the open source *ptlmapper* R package (<https://github.com/fchuffar/ptlmapper>).

## QTL, cvQTL and varQTL mapping

QTL-based mapping was performed as follows. A quantitative trait was considered at the cell-population level. This macro-trait was either the mean (for QTL), the coefficient of variation (standard deviation divided by mean, for cvQTL) or the variance (for varQTL) of the cellular trait in the population of cells. For the yeast morphology data, we selected the traits with  $H^2 > 0.5$  prior to linkage. To do so, we re-computed  $H^2$  values in a way that was consistent with the heritability calculation of the phenotypic space prior to scPTL mapping, where replicates of segregants were pooled together before analysis of inter-strain variation (see above). In this case, only 3 replicates of BY and 3 replicates of RM are then available to estimate the environmental variance. Therefore, we estimated  $var_T$  as the variance of the 59 macro-trait values of the segregants and  $var_E$  as  $(var_{BY} + var_{RM}) / 2$ , with  $var_{BY}$  (resp.  $var_{RM}$ ) being the variance of the three macro-trait values of the BY (resp. RM) strain. We then scanned the genome using the *scanone* function from *r/ctl* [67] with a single QTL model and the multiple imputation method [69]. Our code implementing the calls to *r/ctl* is available in the open source *ptlmapper* R package (<https://github.com/fchuffar/ptlmapper>).

## Permutation tests for statistical significance

We first explain the case where a single trait is studied. When the trait was mapped using R/ctl, significance was assessed by the permutation test implemented in function *scanone()* of the package [67]. For scPTL, we implemented our own permutation test as follows. The significance of an scPTL is the type one error when rejecting the following null hypothesis: "there is no marker at which the genotype of individuals discriminates their location in the phenotypic space", where one 'individual' refers to one population of isogenic cells, and where the 'phenotypic space' is the multi-dimensional space built above by computing Kantorovich distances and applying multi-dimensional scaling. The relevant permutation is therefore to randomly re-assign the phenotypic positions to the individuals before scanning genetic markers for discrimination. We did this 1,000 times. Each time, LDA was applied at dimensions 2 to  $ndim$ , the dimension showing the best contrast (high Z score) was retained, and the highest W score obtained at this dimension was recorded. The empirical threshold corresponding to genome-wide error rates of 0.1%, 1% and 5% were the 99.9<sup>th</sup>, 99<sup>th</sup> and 95<sup>th</sup> percentiles of the 1,000 values produced by the permutations, respectively. These thresholds are typically those employed in whole-genome scans for a single trait.

We now explain the case of the morphological study, where multiple traits (220) were considered. This case is similar to system genetics studies, where the FDR must be controlled. Keeping it below 10% ensures that 9 out of 10 results are true positives, which is often considered as acceptable. Four different methods were used. For three of them, single-cell trait values were resumed to a scalar macro-trait and QTL was searched. The three methods differed by the choice of this macro-trait, which was either the mean or the coefficient of variation of single-cell traits, or the coordinate of individuals on the first axis of the phenotypic space. For each of the three methods, morphological traits with less than 50% genetic heritability (see above) were not considered further, and QTL was searched for the remaining  $N_{traits}$  traits only. For each of these traits, LOD scores were computed on the genome by interval mapping using the macro-trait value as the quantitative phenotype of interest. Significance was assessed by random re-

assignment of the macro-trait values to the individuals (yeast segregants). We did 1,000 such permutations. For each one, the genome was scanned as above and the highest LOD score on the genome was retained. This generated a 1,000 x  $N_{traits}$  matrix  $M_{perm}$  of the hits expected by chance. At a LOD threshold  $L$ , the FDR was computed as:

$$FDR = N_{FalseL} / N_{ActualL}$$

where  $N_{ActualL}$  was the number of linkages obtained from the actual dataset at  $LOD > L$ , and  $N_{FalseL}$  was the expected number of false positives at  $LOD > L$ , which was estimated by the fraction of elements of  $M_{perm}$  exceeding  $L$ .

The fourth method considered all coordinates of the individuals in the phenotypic space. At this step, for each morphological trait, a phenotypic space of  $ndim$  dimensions had been built as explained above by computing Kantorovich distances and applying multi-dimensional scaling. Let P1, P2 and PS be the phenotypic matrices of parent 1 (strain BY), parent 2 (strain RM) and segregants, respectively, with rows being the samples (replicates for P1 and P2, and segregants for PS) and columns being the  $ndim$  coordinates of each sample in the phenotypic space. These matrices had dimensions  $3 \times ndim$  for P1 and P2 and  $59 \times ndim$  for PS. Genetic heritability was computed as  $H^2 = (var_T - var_E) / var_T$ , where the total variance  $var_T$  was the variance of the samples in PS, and where the environmental variance  $var_E$  was estimated as  $(var_{P1} + var_{P2}) / 2$ , with  $var_{P1}$  (resp.  $var_{P2}$ ) being the variance of the samples in P1 (resp. P2). Morphological traits showing  $H^2 < 0.5$  were discarded, and scPTL mapping was applied to the remaining  $N_{traits}$  traits as described above (choice of dimensionality with highest contrast and recording of the best W score obtained on the genome at this dimensionality). Significance of W scores was assessed as described above for the LOD scores, by performing 1,000 permutations and determining the FDR associated to various thresholds of W scores.

## Stochastic model of cell proliferation

For each cell, the probability to divide depended on the concentration of gene product X according to the following Hill-like function (Fig 8A):

$$P(X) = \frac{\beta_0}{1 + (\frac{X}{\vartheta})^n} + \beta_\infty$$

with  $\beta_0 = 0.2$ ,  $\vartheta = 2.5$ ,  $\beta_\infty = 0.05$  and  $n = 2.5$ .

The regulation of X was governed by the same model as above, with  $\eta = 0.16$ . For each individual, parameters  $\alpha_0$ ,  $\alpha_1$  and K were drawn from the same normal distributions as above, where mean and variance depended on the genotype (A or B). At age 0, a population of 1,000 cells was initiated with  $X = 5$ . This population was then evolved by Stochastic Simulation Algorithm [64], with a constant rate of cell death of 0.0001 until the age of 30. The python code implementing this simulation is provided in [S2 Methods](#).

## Yeast strains and plasmids

The yeast strains and oligonucleotides used in this study are listed in [S2 Table](#).

To construct the Gal-GFP reporter, we first removed the MET17 promoter of plasmid pGY8 [19] by digestion with restriction enzymes BspEI and SpeI followed by Klenow fill-in and religation. This generated plasmid pGY10. The GAL1 promoter fragment was digested (BglII-BamHI) from pFA6a-His3MX6-PGAL1 [70] and cloned in the BamHI site of pGY10. A small artificial open reading frame upstream GFP was then removed by digestion with EcoRV and BamHI, Klenow fill-in end blunting and religation. This generated plasmid pGY37, carrying a  $P_{GALI}-yEGFP-NatMX$  cassette that could be integrated at the *HIS3* genomic locus.

Plasmid pGY37 was linearized at NheI and integrated at the *HIS3* locus of strain BY4716 (isogenic to S288c), YEF1946 (a non-clumpy derivative of RM11-1a) and in 61 F1 non-clumpy segregants from BY471xRM11-1a described in [22] to generate strains GY221, GY225, and the S288c x RM11-1a *HIS3:P<sub>GALI</sub>-yEGFP-NatMX:HIS3* set, respectively.

In parallel, we also constructed a *GAL1-GFP<sub>PEST</sub>* reporter coding for a destabilized fluorescent protein [71]. We derived it from pGY334, where *GFP<sub>PEST</sub>* was under the control of the *PGK* promoter. pGY334 was constructed in several steps. The *PGK* promoter was PCR-amplified from pJL49 (gift from Jean-Luc Parrou) using primers 1A23 and 1A24, digested by BamHI and cloned into the BamHI site of pGY10. The resulting plasmid was digested with EcoRV and XbaI, subjected to Klenow fill-in end blunting and religated, generating plasmid pGY13 carrying a *HIS3:P<sub>PGK</sub>-yEGFP-NatMX:HIS3* cassette. The *lox-CEN/ARS-lox* sequence from pALREP [20] was amplified by PCR using primers 1I27 and 1I28 and cloned by homologous recombination into pGY13, generating plasmid pGY252. The *GFP<sub>PEST</sub>* sequence was PCR-amplified from pSVA18 [71] using primers 1I92 and 1I93 and cloned *in vivo* into pGY252 (digested by MfeI and DraIII), leading to pGY334. The *GAL1* promoter fragment was amplified by PCR from pGY37 using primers 1J33 and 1I42 and cloned into plasmid pGY334 by recombination at homologous sequences flanking the BamHI site of the plasmid. The *CEN/ARS* cassette of the resulting plasmid was excised by transient expression of the Cre recombinase in bacteria [20], generating the final integrative plasmid pGY338 carrying the *HIS3:P<sub>GALI</sub>-GFP<sub>PEST</sub>-NatMX:HIS3* cassette.

pGY338 was linearized by NheI and integrated at the *HIS3* locus of BY4724 (isogenic to S288c) and GY1561 to create GY1566 and GY1567 strains, respectively. Strain GY1561 is a non-clumpy derivative of RM11-1a where the KanMX4 cassette was removed. It was obtained by first transforming RM11-1a with an amplicon from plasmid pUG73 [72] obtained with primers 1E75 and 1E76 and selecting a G418-sensitive and LEU<sup>+</sup> transformant (GY739) which was then transformed with pSH47 [73] for expression of the CRE recombinase. After an episode of galactose induction, a LEU<sup>-</sup> derivative was chosen and cultured in non-selective medium (URA<sup>+</sup>) for loss of pSH47, leading to strain GY744, which was then crossed with GY689 [74] to generate GY1561.

## Galactose response measurements

Liquid cultures in synthetic medium with 2% raffinose were inoculated with a single colony and incubated overnight, then diluted to OD<sub>600</sub> = 0.1 (synthetic medium, 2% raffinose) and grown for 3 to 6 hours. Cells were then resuspended in synthetic medium with 2% raffinose and 0.5% galactose and grown for the desired time (0, 10, 20, 30, 40, 60, 80, 100, 130, 160, 205 and 250 minutes). Cells were then washed with PBS1X, incubated for 8 min in 2% paraformaldehyde (PFA) at room temperature, followed by 12 min of incubation in PBS supplemented with Glycine 0.1M at room temperature and finally resuspended in PBS. They were then analyzed on a FACSCalibur (BD Biosciences) flow cytometer to record 10,000 cells per sample.

Flow cytometry data was analysed using the *flowCore* package from Bioconductor [75]. Cells of homogeneous size were dynamically gated and normalized as follows: (i) removal of events with saturated signals (FSC, SSC or FL1  $\geq 1023$  or  $\leq 0$ ), (ii) correction of FL1 values by subtracting the mean(FL1) observed on the same strain at  $t = 0$ , (iii) computation of a density kernel of FSC,SSC values to define a perimeter of peak density containing 60% of events, (iv) cell gating using this perimeter and (v) removal of samples containing less than 3,000 cells at the end of the procedure. The GFP expression values were the corrected FL1 signal of the retained cells.

## Dot6 allele replacement

The *DOT6*<sup>RM</sup> allele was amplified by PCR from genomic DNA of the RM strain using primers 1K87 and 1K88. It was then cloned into plasmid pALREP [20] by homologous recombination at sequences flanking the HpaI site of the plasmid. The CEN/ARS cassette of the resulting plasmid was excised by transient expression of the Cre recombinase in bacteria, as previously described [20], generating plasmid pGY389, which was linearized at EcoRI (a unique site within the DOT6 gene) and integrated in strain GY1566 (isogenic to BY, and carrying the *HIS3:P<sub>GALI</sub>-GFP<sub>PEST</sub>:HIS3* cassette). The pop-in pop-out strategy was applied as previously described [20] and four independent transformants were selected (GY1604, GY1605, GY1606 and GY1607) where PCR and sequencing validated the replacement of the DOT6 allele.

## Data and method availability

The yeast morphological data corresponds to the experiments described in [22]. For the present study, raw images were re-analyzed using CalMorph 1.0. The single-cell values and genotypes used are provided in [S1 Dataset](#) of this article.

The flow cytometry data corresponding to yeast galactose response is made available from <http://flowrepository.org> under accession number FR-FCM-ZZPA.

The simulated data of [Fig 4](#) is available as an R package (ptldata) from <https://github.com/fchuffar/ptldata>.

The scPTL mapping method is made available as an open source R package (ptlmapper) which can be downloaded from <https://github.com/fchuffar/ptlmapper>. A tutorial of this package explains how to run the analysis on the simulated dataset.

## Supporting Information

**S1 Fig. DOT6 allele-replacement experiment.** Strains GY1566 (BY), GY1567 (RM) and GY1604, GY1605, GY1606, GY1607 (BY-*DOT6*<sup>RM</sup>) were cultivated in raffinose 2% and were shifted to a medium containing Raffinose 2% and Galactose 0.5%. After the indicated time, cultures were fixed with paraformaldehyde and analysed by flow cytometry. Histograms correspond to the fluorescent values obtained on cells gated for cell-size (see [methods](#)). (PDF)

**S2 Fig. A)** Confidence interval of scPTL location. An example of confidence intervals obtained by bootstrap is shown. In this case, we used the data from 90 individuals shown in [Fig 4](#) (simulations with noise strength  $\eta = 0.2$ ). We generated 1,000 bootstrapped samples by randomly choosing individuals, with replacement, and we applied scPTL mapping to each sample. In each case, we recorded the peak position of the scPTL (x-axis) and the dimensionality of the phenotypic space that was retained (y-axis). Lower and upper segment boundaries correspond to the 2.5th and 97.5th percentiles of the observed positions for each dimensionality. **B)** Dimensionalities retained for the yeast morphological traits. Let  $n_1$  be the Kaiser's based number of dimensions retained after MDS, and  $n_2$  the Z-score based number of dimensions retained for linkage. The figure shows the distribution of these values for yeast morphological traits where multi-dimensional scPTL outperformed the one-dimension only. On this plot, the number printed at position  $x = n_2, y = n_1$  is the number of morphological traits for which the corresponding  $(n_1, n_2)$  values were chosen. (PDF)

**S1 Table. List of QTL, cvQTL and scPTL of yeast morphological traits.** (XLS)



**S2 Table. Yeast strains and DNA oligonucleotides used in this study.**  
(DOCX)

**S1 Methods. Description and R code of the major steps of scPTL mapping.**  
(DOCX)

**S2 Methods. Python code implementing the stochastic model of cell proliferation.**  
(TGZ)

**S1 Dataset. Single-cell values of yeast morphological traits.**  
(ZIP)

## Acknowledgments

We thank Tamiki Komatsuzaki for mentioning the potential usefulness of the Kantorovich metric, Fabien Crauste, Olivier François and Emmanuel Grenier for helpful discussions, Satoru Nogami for re-analysis of CalMorph images, Arezki Boudaoud, Fabien Duveau, Olivier Gandrillon, Marie Sémon and Gérard Triqueneaux for critical reading of the manuscript, Sandrine Mouradian and SFR Biosciences Gerland-Lyon Sud (UMS344/US8) for access to flow cytometers and technical assistance, the Pôle Scientifique de Modélisation Numérique (Lyon, France) and the Grid'5000 testbed ([www.grid5000.fr](http://www.grid5000.fr)) for computer resource, developers of R, bioconductor and Ubuntu for their software, and three anonymous reviewers for their comments.

## Author Contributions

Conceived and designed the experiments: GY. Performed the experiments: MR HDB. Analyzed the data: FC DJ MR. Contributed reagents/materials/analysis tools: FC DJ YO CB. Wrote the paper: MR GY. Developed the scPTL mapping method: FC GY. Developed the analysis code: FC. Generated simulated data of Figure 4: DJ. Generated simulated data of Figure 8: CB. Produced figures: FC MR CB GY. Interpreted results: FC MR DJ GY.

## References

1. Rockman MV. The QTN program and the alleles that matter for evolution: all that's gold does not glitter. *Evol Int J Org Evol*. 2012; 66: 1–17.
2. Ikram MK, Xueling S, Jensen RA, Cotch MF, Hewitt AW, Ikram MA, et al. Four Novel Loci (19q13, 6q24, 12q24, and 5q14) Influence the Microcirculation In Vivo. *PLoS Genet*. 2010; 6: e1001184. doi: [10.1371/journal.pgen.1001184](https://doi.org/10.1371/journal.pgen.1001184) PMID: [21060863](https://pubmed.ncbi.nlm.nih.gov/21060863/)
3. Wray NR, Middeldorp CM, Birley AJ, Gordon SD, Sullivan PF, Visscher PM, et al. Genome-wide linkage analysis of multiple measures of neuroticism of 2 large cohorts from Australia and the Netherlands. *Arch Gen Psychiatry*. 2008; 65: 649–658. doi: [10.1001/archpsyc.65.6.649](https://doi.org/10.1001/archpsyc.65.6.649) PMID: [18519823](https://pubmed.ncbi.nlm.nih.gov/18519823/)
4. Winkler TW, Justice AE, Graff M, Barata L, Feitosa MF, Chu S, et al. The Influence of Age and Sex on Genetic Associations with Adult Body Size and Shape: A Large-Scale Genome-Wide Interaction Study. *PLoS Genet*. 2015; 11: e1005378. doi: [10.1371/journal.pgen.1005378](https://doi.org/10.1371/journal.pgen.1005378) PMID: [26426971](https://pubmed.ncbi.nlm.nih.gov/26426971/)
5. Yvert G. "Particle genetics": treating every cell as unique. *Trends Genet*. 2014; 30: 49–56. doi: [10.1016/j.tig.2013.11.002](https://doi.org/10.1016/j.tig.2013.11.002) PMID: [24315431](https://pubmed.ncbi.nlm.nih.gov/24315431/)
6. Raj A, Rifkin SA, Andersen E, van Oudenaarden A. Variability in gene expression underlies incomplete penetrance. *Nature*. 2010; 463: 913–8. doi: [10.1038/nature08781](https://doi.org/10.1038/nature08781) PMID: [20164922](https://pubmed.ncbi.nlm.nih.gov/20164922/)
7. Roesch A, Fukunaga-Kalabis M, Schmidt EC, Zabierowski SE, Brafford PA, Vultur A, et al. A Temporarily Distinct Subpopulation of Slow-Cycling Melanoma Cells Is Required for Continuous Tumor Growth. *Cell*. 2010; 141: 583–594. doi: [10.1016/j.cell.2010.04.020](https://doi.org/10.1016/j.cell.2010.04.020) PMID: [20478252](https://pubmed.ncbi.nlm.nih.gov/20478252/)
8. Wang C-C, Bajikar SS, Jamal L, Atkins KA, Janes KA. A time- and matrix-dependent TGFB $\beta$ 3–JUND–KRT5 regulatory circuit in single breast epithelial cells and basal-like premalignancies. *Nat Cell Biol*. 2014; 16: 345–356. doi: [10.1038/ncb2930](https://doi.org/10.1038/ncb2930) PMID: [24658685](https://pubmed.ncbi.nlm.nih.gov/24658685/)

9. Spencer SL, Gaudet S, Albeck JG, Burke JM, Sorger PK. Non-genetic origins of cell-to-cell variability in TRAIL-induced apoptosis. *Nature*. 2009; 459: 428–432. doi: [10.1038/nature08012](https://doi.org/10.1038/nature08012) PMID: [19363473](https://pubmed.ncbi.nlm.nih.gov/19363473/)
10. Sharma SV, Lee DY, Li B, Quinlan MP, Takahashi F, Maheswaran S, et al. A Chromatin-Mediated Reversible Drug-Tolerant State in Cancer Cell Subpopulations. *Cell*. 2010; 141: 69–80. doi: [10.1016/j.cell.2010.02.027](https://doi.org/10.1016/j.cell.2010.02.027) PMID: [20371346](https://pubmed.ncbi.nlm.nih.gov/20371346/)
11. Pisco AO, Brock A, Zhou J, Moor A, Mojtahedi M, Jackson D, et al. Non-Darwinian dynamics in therapy-induced cancer drug resistance. *Nat Commun*. 2013; 4.
12. Tegos GP, Evangelisti AM, Strouse JJ, Ursu O, Bologna C, Sklar LA. A high throughput flow cytometric assay platform targeting transporter inhibition. *Drug Discov Today Technol*. 2014; 12: e95–103. doi: [10.1016/j.ddtec.2014.03.010](https://doi.org/10.1016/j.ddtec.2014.03.010) PMID: [25027381](https://pubmed.ncbi.nlm.nih.gov/25027381/)
13. Bendall SC, Simonds EF, Qiu P, Amir ED, Krutzik PO, Finck R, et al. Single-Cell Mass Cytometry of Differential Immune and Drug Responses Across a Human Hematopoietic Continuum. *Science*. 2011; 332: 687–696. doi: [10.1126/science.1198704](https://doi.org/10.1126/science.1198704) PMID: [21551058](https://pubmed.ncbi.nlm.nih.gov/21551058/)
14. Lin J-R, Fallahi-Sichani M, Sorger PK. Highly multiplexed imaging of single cells using a high-throughput cyclic immunofluorescence method. *Nat Commun*. 2015; 6: 8390. doi: [10.1038/ncomms9390](https://doi.org/10.1038/ncomms9390) PMID: [26399630](https://pubmed.ncbi.nlm.nih.gov/26399630/)
15. Samsel L, McCoy JP Jr.. Imaging flow cytometry for the study of erythroid cell biology and pathology. *J Immunol Methods*. 2015; 423: 52–59. doi: [10.1016/j.jim.2015.03.019](https://doi.org/10.1016/j.jim.2015.03.019) PMID: [25858229](https://pubmed.ncbi.nlm.nih.gov/25858229/)
16. Ohya Y, Kimori Y, Okada H, Ohnuki S. Single-cell phenomics in budding yeast. *Mol Biol Cell*. 2015; 26: 3920–3925. doi: [10.1091/mbc.E15-07-0466](https://doi.org/10.1091/mbc.E15-07-0466) PMID: [26543200](https://pubmed.ncbi.nlm.nih.gov/26543200/)
17. Klein AM, Mazutis L, Akartuna I, Tallapragada N, Veres A, Li V, et al. Droplet Barcoding for Single-Cell Transcriptomics Applied to Embryonic Stem Cells. *Cell*. 2015; 161: 1187–1201. doi: [10.1016/j.cell.2015.04.044](https://doi.org/10.1016/j.cell.2015.04.044) PMID: [26000487](https://pubmed.ncbi.nlm.nih.gov/26000487/)
18. Macosko EZ, Basu A, Satija R, Nemesh J, Shekhar K, Goldman M, et al. Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell*. 2015; 161: 1202–1214. doi: [10.1016/j.cell.2015.05.002](https://doi.org/10.1016/j.cell.2015.05.002) PMID: [26000488](https://pubmed.ncbi.nlm.nih.gov/26000488/)
19. Ansel J, Bottin H, Rodriguez-Beltran C, Damon C, Nagarajan M, Fehrmann S, et al. Cell-to-cell stochastic variation in gene expression is a complex genetic trait. *PLoS Genet*. 2008; 4: e1000049. doi: [10.1371/journal.pgen.1000049](https://doi.org/10.1371/journal.pgen.1000049) PMID: [18404214](https://pubmed.ncbi.nlm.nih.gov/18404214/)
20. Fehrmann S, Bottin-Duplus H, Leonidou A, Mollereau E, Barthelaix A, Wei W, et al. Natural sequence variants of yeast environmental sensors confer cell-to-cell expression variability. *Mol Syst Biol*. 2013; 9: 695. doi: [10.1038/msb.2013.53](https://doi.org/10.1038/msb.2013.53) PMID: [24104478](https://pubmed.ncbi.nlm.nih.gov/24104478/)
21. Metzger BPH, Yuan DC, Gruber JD, Duveau F, Wittkopp PJ. Selection on noise constrains variation in a eukaryotic promoter. *Nature*. 2015; 521: 344–347. doi: [10.1038/nature14244](https://doi.org/10.1038/nature14244) PMID: [25778704](https://pubmed.ncbi.nlm.nih.gov/25778704/)
22. Nogami S, Ohya Y, Yvert G. Genetic complexity and quantitative trait loci mapping of yeast morphological traits. *PLoS Genet*. 2007; 3: e31. PMID: [17319748](https://pubmed.ncbi.nlm.nih.gov/17319748/)
23. Levy SF, Siegal ML. Network hubs buffer environmental variation in *Saccharomyces cerevisiae*. *PLoS Biol*. 2008; 6: e264. doi: [10.1371/journal.pbio.0060264](https://doi.org/10.1371/journal.pbio.0060264) PMID: [18986213](https://pubmed.ncbi.nlm.nih.gov/18986213/)
24. Yvert G, Ohnuki S, Nogami S, Imanaga Y, Fehrmann S, Schacherer J, et al. Single-cell phenomics reveals intra-species variation of phenotypic noise in yeast. *BMC Syst Biol*. 2013; 7: 54. doi: [10.1186/1752-0509-7-54](https://doi.org/10.1186/1752-0509-7-54) PMID: [23822767](https://pubmed.ncbi.nlm.nih.gov/23822767/)
25. Rönnegård L, Valdar W. Detecting Major Genetic Loci Controlling Phenotypic Variability in Experimental Crosses. *Genetics*. 2011; 188: 435–447. doi: [10.1534/genetics.111.127068](https://doi.org/10.1534/genetics.111.127068) PMID: [21467569](https://pubmed.ncbi.nlm.nih.gov/21467569/)
26. Cao Y, Maxwell TJ, Wei P. A Family-Based Joint Test for Mean and Variance Heterogeneity for Quantitative Traits. *Ann Hum Genet*. 2015; 79: 46–56. doi: [10.1111/ahg.12089](https://doi.org/10.1111/ahg.12089) PMID: [25393880](https://pubmed.ncbi.nlm.nih.gov/25393880/)
27. Baba A, Komatsuzaki T. Construction of effective free energy landscape from single-molecule time series. *Proc Natl Acad Sci*. 2007; 104: 19297–19302. PMID: [18048341](https://pubmed.ncbi.nlm.nih.gov/18048341/)
28. Vershik AM. Kantorovich Metric: Initial History and Little-Known Applications. *J Math Sci*. 2006; 133: 1410–1417.
29. Kantorovich LV. Mathematics in Economics: Achievements, Difficulties, Perspectives. *Am Econ Rev*. 1989; 79: 18–22.
30. Kantorovich LV. On the Translocation of Masses. *J Math Sci*. 2006; 133: 1381–1382.
31. Anderson MJ, Willis TJ. Canonical analysis of principal coordinates: a useful method of constrained ordination for ecology. *Ecology*. 2003; 84: 511–525.
32. Anderson MJ, Robinson J. Generalized discriminant analysis based on distances. *Aust N Z J Stat*. 2003; 45: 301–318.
33. Cox TF, Cox MAA. *Multidimensional Scaling*, Second Edition. CRC Press; 2000.

34. Hirsch JP, Dietzel C, Kurjan J. The carboxyl terminus of Scg1, the G alpha subunit involved in yeast mating, is implicated in interactions with the pheromone receptors. *Genes Dev.* 1991; 5: 467–74. PMID: [1848203](#)
35. Brown AJ, Dyos SL, Whiteway MS, White JH, Watson MA, Marzioch M, et al. Functional coupling of mammalian receptors to the yeast mating pathway using novel yeast/mammalian G protein alpha-subunit chimeras. *Yeast.* 2000; 16: 11–22. PMID: [10620771](#)
36. Yvert G, Brem RB, Whittle J, Akey JM, Foss E, Smith EN, et al. Trans-acting regulatory variation in *Saccharomyces cerevisiae* and the role of transcription factors. *Nat Genet.* 2003; 35: 57–64.
37. Lang GI, Murray AW, Botstein D. The cost of gene expression underlies a fitness trade-off in yeast. *Proc Natl Acad Sci.* 2009; 106: 5755–5760. doi: [10.1073/pnas.0901620106](#) PMID: [19299502](#)
38. New AM, Cerulus B, Govers SK, Perez-Samper G, Zhu B, Boogmans S, et al. Different Levels of Catabolite Repression Optimize Growth in Stable and Variable Environments. *PLoS Biol.* 2014; 12: e1001764. doi: [10.1371/journal.pbio.1001764](#) PMID: [24453942](#)
39. Wang J, Atolia E, Hua B, Savir Y, Escalante-Chong R, Springer M. Natural Variation in Preparation for Nutrient Depletion Reveals a Cost–Benefit Tradeoff. *PLoS Biol.* 2015; 13: e1002041. doi: [10.1371/journal.pbio.1002041](#) PMID: [25626068](#)
40. Cherry JM, Hong EL, Amundsen C, Balakrishnan R, Binkley G, Chan ET, et al. *Saccharomyces Genome Database: the genomics resource of budding yeast.* *Nucleic Acids Res.* 2012; 40: D700–D705. doi: [10.1093/nar/gkr1029](#) PMID: [22110037](#)
41. Dalal CK, Cai L, Lin Y, Rahbar K, Elowitz MB. Pulsatile Dynamics in the Yeast Proteome. *Curr Biol.* 2014; 24: 2189–2194. doi: [10.1016/j.cub.2014.07.076](#) PMID: [25220054](#)
42. Garber JE, Offit K. Hereditary cancer predisposition syndromes. *J Clin Oncol.* 2005; 23: 276–92. PMID: [15637391](#)
43. Rahman N. Realizing the promise of cancer predisposition genes. *Nature.* 2014; 505: 302–308. doi: [10.1038/nature12981](#) PMID: [24429628](#)
44. Schadt EE, Lamb J, Yang X, Zhu J, Edwards S, Guhathakurta D, et al. An integrative genomics approach to infer causal associations between gene expression and disease. *Nat Genet.* 2005; 37: 710–7. PMID: [15965475](#)
45. Elez M, Murray AW, Bi L-J, Zhang X-E, Matic I, Radman M. Seeing Mutations in Living Cells. *Curr Biol.* 2010; 20: 1432–1437. doi: [10.1016/j.cub.2010.06.071](#) PMID: [20674359](#)
46. Uphoff S, Lord ND, Okumus B, Potvin-Trottier L, Sherratt DJ, Paulsson J. Stochastic activation of a DNA damage response causes cell-to-cell mutation rate variation. *Science.* 2016; 351: 1094–1097. doi: [10.1126/science.aac9786](#) PMID: [26941321](#)
47. Wang L, Janes KA. Stochastic profiling of transcriptional regulatory heterogeneities in tissues, tumors and cultured cells. *Nat Protoc.* 2013; 8: 282–301. doi: [10.1038/nprot.2012.158](#) PMID: [23306461](#)
48. Singh DK, Ku C-J, Wichaidit C, Steininger RJ, Wu LF, Altschuler SJ. Patterns of basal signaling heterogeneity can distinguish cellular populations with different drug sensitivities. *Mol Syst Biol.* 2010; 6.
49. Prince J, Lundgren A, Stadnisky MD, Nash WT, Beeber A, Turner SD, et al. Multiparametric Analysis of Host Response to Murine Cytomegalovirus in MHC Class I–Disparate Mice Reveals Primacy of Dk-Licensed Ly49G2+ NK Cells in Viral Control. *J Immunol.* 2013; 191: 4709–4719. doi: [10.4049/jimmunol.1301388](#) PMID: [24068668](#)
50. Chang HH, Hemberg M, Barahona M, Ingber DE, Huang S. Transcriptome-wide noise controls lineage choice in mammalian progenitor cells. *Nature.* 2008; 453: 544–7. doi: [10.1038/nature06965](#) PMID: [18497826](#)
51. Singer ZS, Yong J, Tischler J, Hackett JA, Altinok A, Surani MA, et al. Dynamic Heterogeneity and DNA Methylation in Embryonic Stem Cells. *Mol Cell.* 2014; 55: 319–331. doi: [10.1016/j.molcel.2014.06.029](#) PMID: [25038413](#)
52. Carter AJ, Houle D. Artificial selection reveals heritable variation for developmental instability. *Evolution.* 2011; 65: 3558–64. doi: [10.1111/j.1558-5646.2011.01393.x](#) PMID: [22133225](#)
53. Gavin-Smyth J, Wang Y-C, Butler I, Ferguson EL. A Genetic Network Conferring Canalization to a Bistable Patterning System in *Drosophila*. *Curr Biol.* 2013; 23: 2296–2302. doi: [10.1016/j.cub.2013.09.055](#) PMID: [24184102](#)
54. Richard M, Yvert G. How does evolution tune biological noise? *Syst Biol.* 2014; 5: 374.
55. Zhou X, Stephens M. Genome-wide efficient mixed-model analysis for association studies. *Nat Genet.* 2012; 44: 821–824. doi: [10.1038/ng.2310](#) PMID: [22706312](#)
56. Stephan J, Stegle O, Beyer A. A random forest approach to capture genetic effects in the presence of population structure. *Nat Commun.* 2015; 6: 7432. doi: [10.1038/ncomms8432](#) PMID: [26109276](#)

57. Kurita T. Discriminant Kernels derived from the optimum nonlinear discriminant analysis. The 2011 International Joint Conference on Neural Networks (IJCNN). 2011. pp. 299–306.
58. Visscher PM, Thompson R, Haley CS. Confidence Intervals in QTL Mapping by Bootstrapping. *Genetics*. 1996; 143: 1013–1020. PMID: [8725246](#)
59. Flatz L, Roychoudhuri R, Honda M, Filali-Mouhim A, Goulet J-P, Kettaf N, et al. Single-cell gene-expression profiling reveals qualitatively distinct CD8 T cells elicited by different gene-based vaccines. *Proc Natl Acad Sci*. 2011; 108: 5724–5729. doi: [10.1073/pnas.1013084108](#) PMID: [21422297](#)
60. Snijder B, Sacher R, Rämö P, Damm E-M, Liberali P, Pelkmans L. Population context determines cell-to-cell variability in endocytosis and virus infection. *Nature*. 2009; 461: 520–523. doi: [10.1038/nature08282](#) PMID: [19710653](#)
61. Gangbo W, Świąch A. Optimal maps for the multidimensional Monge-Kantorovich problem. *Commun Pure Appl Math*. 1998; 51: 23–45. doi: [10.1002/\(SICI\)1097-0312\(199801\)51:1<23::AID-CPA2>3.0.CO;2-H](#)
62. Group TFP. Phylogenetic inference for function-valued traits: speech sound evolution. *Trends Ecol Evol*. 2012; 27: 160–166. doi: [10.1016/j.tree.2011.10.001](#) PMID: [22078766](#)
63. Hadjipantelis PZ, Jones NS, Moriarty J, Springate DA, Knight CG. Function-valued traits in evolution. *J R Soc Interface*. 2013; 10: 20121032. doi: [10.1098/rsif.2012.1032](#) PMID: [23427095](#)
64. Gillespie DT. Exact stochastic simulation of coupled chemical reactions. *J Phys Chem*. 1977; 81: 2340–2361.
65. R Core Team. R: A Language and Environment for Statistical Computing [Internet]. Vienna, Austria: R Foundation for Statistical Computing; 2015. Available: <https://www.R-project.org>
66. Smith EN, Kruglyak L. Gene-Environment Interaction in Yeast Gene Expression. *PLoS Biol*. 2008; 6: e83. doi: [10.1371/journal.pbio.0060083](#) PMID: [18416601](#)
67. Broman KW, Wu H, Sen S, Churchill GA. R/qtl: QTL mapping in experimental crosses. *Bioinformatics*. 2003; 19: 889–890. PMID: [12724300](#)
68. Muller KE, Peterson BL. Practical methods for computing power in testing the multivariate general linear hypothesis. *Comput Stat Data Anal*. 1984; 2: 143–158.
69. Sen S, Churchill GA. A Statistical Framework for Quantitative Trait Mapping. *Genetics*. 2001; 159: 371–387. PMID: [11560912](#)
70. Longtine MS, McKenzie A, Demarini DJ, Shah NG, Wach A, Brachat A, et al. Additional modules for versatile and economical PCR-based gene deletion and modification in *Saccharomyces cerevisiae*. *Yeast*. 1998; 14: 953–61. PMID: [9717241](#)
71. Mateus C, Avery SV. Destabilized green fluorescent protein for monitoring dynamic changes in yeast gene expression with flow cytometry. *Yeast*. 2000; 16: 1313–23. PMID: [11015728](#)
72. Gueldener U, Heinisch J, Koehler GJ, Voss D, Hegemann JH. A second set of loxP marker cassettes for Cre-mediated multiple gene knockouts in budding yeast. *Nucleic Acids Res*. 2002; 30: e23. PMID: [11884642](#)
73. Gueldener U, Heck S, Fielder T, Beinhauer J, Hegemann JH. A new efficient gene disruption cassette for repeated use in budding yeast. *Nucleic Acids Res*. 1996; 24: 2519–24. PMID: [8692690](#)
74. Abraham AL, Nagarajan M, Veyrieras JB, Bottin H, Steinmetz LM, Yvert G. Genetic modifiers of chromatin acetylation antagonize the reprogramming of epi-polymorphisms. *PLoS Genet*. 2012; 8: e1002958. doi: [10.1371/journal.pgen.1002958](#) PMID: [23028365](#)
75. Hahne F, LeMeur N, Brinkman RR, Ellis B, Haaland P, Sarkar D, et al. flowCore: a Bioconductor package for high throughput flow cytometry. *BMC Bioinformatics*. 2009; 10: 106. doi: [10.1186/1471-2105-10-106](#) PMID: [19358741](#)