1
2

# Refining the resolution of the yeast genotype-phenotype map using single-cell RNA-sequencing

3

4   **Authors**: Arnaud N'Guessan[1], Wen Yuan Tong[2,¥], Hamed Heydari[3,4], Alex N Nguyen Ba[1,2]

5   1. Department of Cell and Systems Biology, University of Toronto, Ramsay Wright Laboratories,
6   25 Harbord St, M5S3G5, Toronto, Ontario, Canada.

7   2. Department of Biology, University of Toronto at Mississauga, 3359 Mississauga Rd, L5L 1C5,
8   Mississauga, Ontario, Canada

9   3. Department of Molecular Genetics, University of Toronto, Medical Science Building, Room
10  4386, 1 King's College Cir, M5S1A8, Toronto, Ontario, Canada.

11  4. Donnelly Centre for Cellular and Biomolecular Research, University of Toronto, 160 College
12  Street, Room 230, M5S3E1, Toronto, Ontario, Canada.

13  ¥. Present address: Center of Molecular and Cellular Oncology, Yale University, 300 George
14  Street Suite 6400, New Haven, CT

15  **ABSTRACT**

16  Genotype-phenotype mapping (GPM) or the association of trait variation to genetic variation has
17  been a long-lasting problem in biology. The existing approaches to this problem allowed
18  researchers to partially understand within- and between-species variation as well as the emergence
19  or evolution of phenotypes. However, traditional GPM methods typically ignore the transcriptome
20  or have low statistical power due to challenges related to dataset scale. Thus, it is not clear to what
21  extent selection modulates transcriptomes and whether cis- or trans-regulatory elements are more
22  important. To overcome these challenges, we leveraged the cost efficiency and scalability of
23  single-cell RNA sequencing (scRNA-seq) by collecting data from 18,233 yeast cells from 4,489
24  segregants of a cross between the laboratory strain BY4741 and the vineyard strain RM11-1a.
25  More precisely, we performed eQTL mapping with the scRNA-seq data to identify single-cell
26  eQTL (sc-eQTL) and transcriptome variation patterns associated to fitness variation inferred from
27  the segregants' bulk fitness assay. Due to the larger scale of our dataset, we were able to
28  recapitulate results from decades of work in GPM from yeast bulk assays while revealing new
29  associations between phenotypic and transcriptomic variations. The multidimensionality of this
30  dataset also allowed us to measure phenotype and expression heritability and partition the variance
31  of cell fitness into genotype and expression components to highlight selective pressure at both
32  levels. Altogether these results suggest that integrating large-scale scRNA-seq data into GPM
33  improves our understanding of trait variation in the context of transcriptomic regulation.

34

**INTRODUCTION**

The process by which DNA encodes proteins via transcription and translation has been studied for decades to make sense of organisms' phenotypes. However, being able to explain organisms' phenotypes from their genetic material, i.e. genotype-phenotype mapping (GPM), has been a long-lasting problem with important applications (1,2). Indeed, making sense of genetic variation at the phenotypic level enables the understanding of trait variation between and within species as well as the emergence and evolution of phenotypes (3). For instance, reverse genetics approaches, e.g. gene knockout or transgenic technologies, and forward genetics approaches like GWAS and QTL mapping helped in determining the function of multiple genes and the effects of mutations on growth in different environments (4). However, reverse genetics approaches typically fail to account for natural variation and forward genetics approaches like QTL mapping typically focus on genetic and phenotypic variation so they cannot highlight selection on the transcriptome.

An essential characteristic of this problem is the multi-layered organization of the GPM. Indeed, GPM is not strictly restricted to the direct association between genotypes and phenotypes. This association is better resolved and complemented by understanding the intermediary transcriptome layer, e.g. cell mechanisms at the transcriptomic level are involved in diseases and pathogenicity (2,5–8). However, it is not clear to what extent transcriptomic changes relate to phenotypic changes or selection. Pioneering work from Mary-Claire King and Allan Charles Wilson set the tone for investigating this question by proposing that variations in morphological and behavioral traits arise more often through gene expression regulation than evolution at the protein-coding level (9). François Jacob then postulated an essay that stemmed from this theory in which he highlights how evolution acts as a tinkerer that works from already available material, i.e. through regulation of gene expression, to create new adaptations (10). This constituted the core of the evolutionary developmental biology which matured into the still-debated claim that new adaptations mainly emerge through cis-regulation of gene expression, i.e. through noncoding DNA regulating a neighbor gene contrarily to trans-regulators acting on distant genes (11–14). This debate has been reinforced by the technical difficulties and complexity of assessing the evolution and outcome of mutations in non-coding regions (11,12). Advances in sequencing technologies have clarified some of these hypotheses, particularly in the context of transcriptome analyses of the model organism *Saccharomyces cerevisiae*. For instance, Brem et al (2002) used microarray technology to relate the gene expression profiles of 40 yeast segregants from a lab (BY) and natural vineyard strain (RM) to their genetic markers (15). They found that cis-acting modulation is the main mechanism for regulating gene expression. Nearly two decades later, by greatly increasing statistical power, Albert and collaborators (2018) found that most of the expression variation arise through trans-regulation using non-multiplexed RNA-seq to analyze 5,720 genes in 1,012 yeast segregants generated by a crossing between RM and BY (16). The analysis method they used, i.e. expression quantitative trait loci (eQTL) mapping, consists in correlating allele frequencies to gene expression levels to find the loci modulating expression.

Although eQTL mapping is a traditional GPM analysis that accounts for the transcriptomic layer, it is typically realized through non-multiplexed RNA-seq which tends to have low statistical power due to challenges with experimental scale and confounding factors (17,18). Thus, eQTL

76  mapping traditionally cannot identify significant low-effect regulatory mutations that are
77  important for understanding the genetic bases of complex traits and diseases (19,20). Furthermore,
78  most eQTL studies only assess the average transcriptomic profile of bulk populations without
79  being able to capture the profile of rare cell lineages within a population. This is a critical limitation
80  in heterogenous populations such as cancer or microbial populations where rare lineages can drive
81  relapse or drug resistance (21).

82  Here, we sought to circumvent the challenges of non-multiplexed bulk RNA-seq imposed
83  by the scale and population heterogeneity by performing eQTL mapping through single-cell RNA
84  sequencing (scRNA-seq) of a pool of ~4500 well-characterized F1 segregants of a yeast cross
85  (16,22,23). In the same way that combinatorial indexing/barcoding and multiplexing enable the
86  collection of large-scale fitness and genotype data (24), we hypothesized that scRNA-seq can help
87  us collect both genotype and expression data on a large pool of segregants. We employ several
88  strategies to overcome previous obstacles of eQTL mapping studies: i) we pool cells from
89  thousands of segregants during the growth step and perform a single scRNA-seq run on the culture
90  to account for environmental effects, and ii) from the exome sequencing data of single-cells we
91  take advantage of the reference panel to validate that we accurately infer the genotype of each cell
92  from extremely low number of reads mapping to polymorphic sites per cell (effectively ~0.2x
93  coverage).

94  Using this approach, we integrated the resulting transcriptomic data from growth in rich
95  media with a pre-existing yeast GPM. We estimated the heritability of the transcriptome and the
96  extent at which transcriptome is associated with fitness. We show that this increased scale from
97  scRNA-seq enables eQTL mapping directly without the use of a reference genotype panel, and
98  relate identified single-cell eQTL (sc-eQTL) to previously identified QTL. We also exploit the
99  identified sc-eQTL to analyze the patterns of cis- and trans- regulation in the GPM.

100

**Our single-cell RNA-seq approach is consistent with yeast GPM results from non-multiplexed assays**

103  We initially aimed to show that performing scRNA-seq at a large scale can generate data that are
104  consistent with non-multiplexed DNA and RNA sequencing. To do so, we analyzed a dataset of
105  thousands of yeast lineages generated by Nguyen Ba and collaborators (2022) (24). To understand
106  the yeast GPM, they collected fitness and genotype data from ~100,000 segregants of an F1 cross
107  between a laboratory strain of yeast (BY) and a natural vineyard strain (RM) (**Figure 1A**).
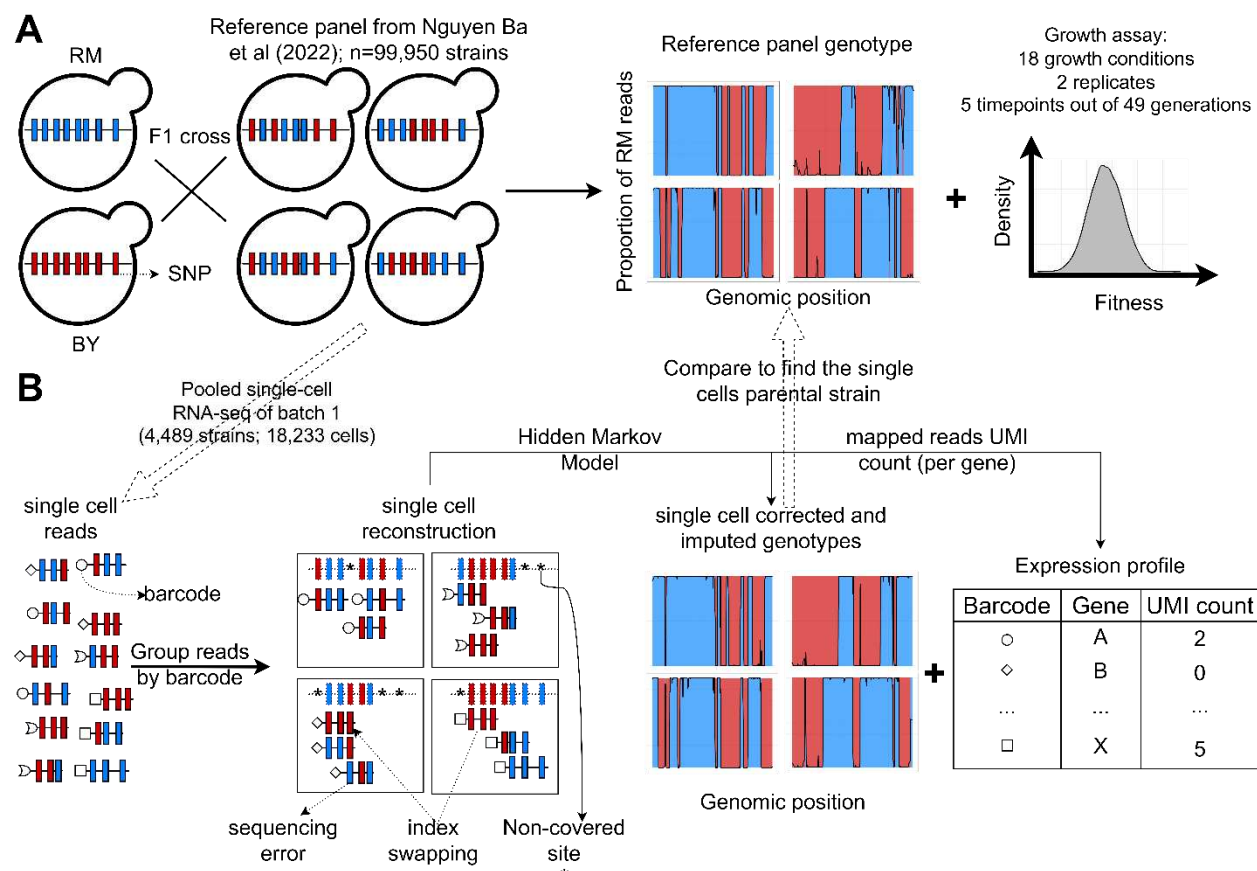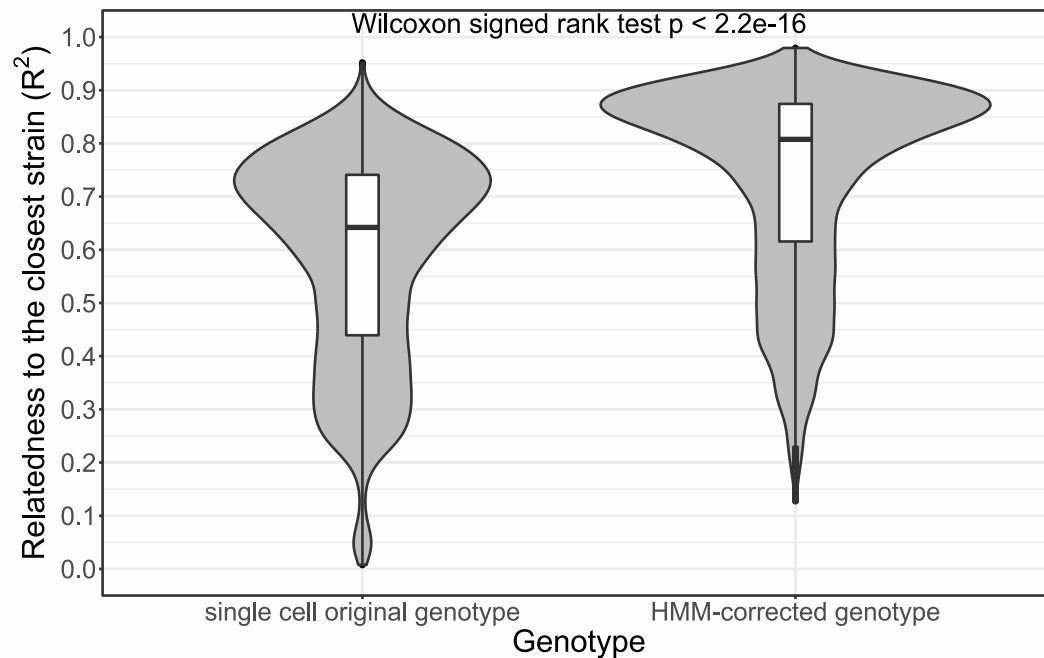
108

109 **Figure 1 Yeast segregants datasets**. A) Reference panel from the barcoded bulk sequencing. The
110 99,995 yeast segregants in the reference panel come from a F1 cross between a laboratory strain
111 of yeast (BY) and a natural vineyard strain (RM) (24). Thus, they only have 2 possible alleles at
112 each of the 41,594 polymorphic sites. The lineages barcodes enabled fitness estimation from
113 competition assays in 18 environments recapitulating the adaptation to temperature gradients, the
114 ability to process different sources of carbon and the resistance to antifungal compounds. B)
115 Pooled scRNA-seq dataset from a single batch. We performed scRNA-seq of the first batch of
116 segregants (n=4,489) to obtain genotypes that are similar to the reference panel and single cell's
117 expression profiles. Non-covered sites, sequencing errors and the presence of reads in the wrong
118 library (index swapping) are corrected for using the HMM described in Figure S1.

119 Using this approach named barcoded bulk QTL mapping or BB-QTL mapping, they revealed the
120 complex polygenic and pleiotropic nature of phenotypes as well as an unprecedented number of
121 pairwise epistatic interactions. To integrate transcriptomic data to that GPM, we performed
122 scRNA-seq using the 10X Genomics Chromium microfluidics platform and obtained both
123 genotype and expression profiles from 18,233 cells of the first batch of segregants (**Figure 1B**).
124 This short-read scRNA-seq method comes with challenges like low-coverage sites due to technical
125 sequencing biases and low sequencing depth in some cells (25,26). To overcome these challenges,
126 the unique molecular identifiers (UMIs) of the 10X Genomics platform provide a control for
127 technical biases by quantifying gene expression from unique transcribed molecule counts instead
128 of reads counts (25). In addition, Hidden Markov Models (HMMs) can infer accurate genotype

129     data even at sequencing depths as low as 0.1x (24). Nguyen Ba and collaborators (2022) designed
130     an HMM to infer the segregants genotypes from the observed reads at low depth of DNA
131     sequencing by accounting for sequencing error rate, recombination rate and index swapping rate
132     (24). As there are only two ancestral lineages, there are only two possible alleles for the strains at
133     each of the 41,594 polymorphic sites. Thus, the genotype of the segregants can be represented by
134     the frequency of only one of the parental alleles, which is RM in the dataset. Applying this model
135     to low-coverage segregants yielded genotypes that are significantly similar to high-coverage
136     replicates (24). We sought to use a similar model to infer genotypes from scRNA-seq data, but we
137     anticipated that some of these parameters may differ due to increased error rate of the reverse-
138     transcriptase, increased index swapping due to pooled-reaction, etc (**Figure S1**). In Nguyen Ba et
139     al, those rates were heuristically determined, but here we estimated these from the read mapping
140     data and found that re-estimated parameters from data increase the proportion of recovered strains
141     in the single cell data from 58.6% to 72.0%.

142         After adapting the HMM to the scRNA-seq data, we sought to validate that the resulting
143     cell genotypes relate well to their corresponding strain in the reference panel obtained by non-
144     multiplexed DNA sequencing strategies. Ideally, each single-cell barcode (from 10x Genomics
145     Chromium) should be associated with a single cell and a cell should have a clear match with a
146     unique strain in the reference panel. However, several factors can obscure these associations, e.g.
147     a single-cell droplet containing cells from 2 different strains, a low-coverage cell, uncertainty in
148     the allele of the reference genotype, etc. Thus, we designed an approach to clearly assign cells to
149     the correct reference panel strain (see **Methods**). This approach relies on two metrics of similarity
150     between the cells and the strains' genotypes, i.e. the expected distance between them, which should
151     be minimized for the best match, and the relatedness ($R^2$). The statistical significance of the
152     relatedness between single cells and reference lineages was determined by a permutation test
153     (**Figure S2**). From the read mapping alone, we obtained a mean $R^2$ of 0.59 ($\sigma = 0.19$ and median
154     = 0.64), which was significantly improved after applying our HMM to correct for mis-identified
155     alleles and imputing data in low-coverage sites using recombination probability. Indeed, the
156     single-cell HMM genotypes yield a mean $R^2$ of 0.73 ($\sigma = 0.18$ and median = 0.81; **Figure2A**). We
157     found that the distribution of relatedness after HMM was still left-skewed, with many cells
158     statistically significantly assigned to a reference genotype despite having what appeared to be low
159     relatedness. Upon investigation, it was found that these could be explained by genotyping
160     uncertainty either in the single-cell and/or in the reference panel genotype (**s**) (**Table S1**).
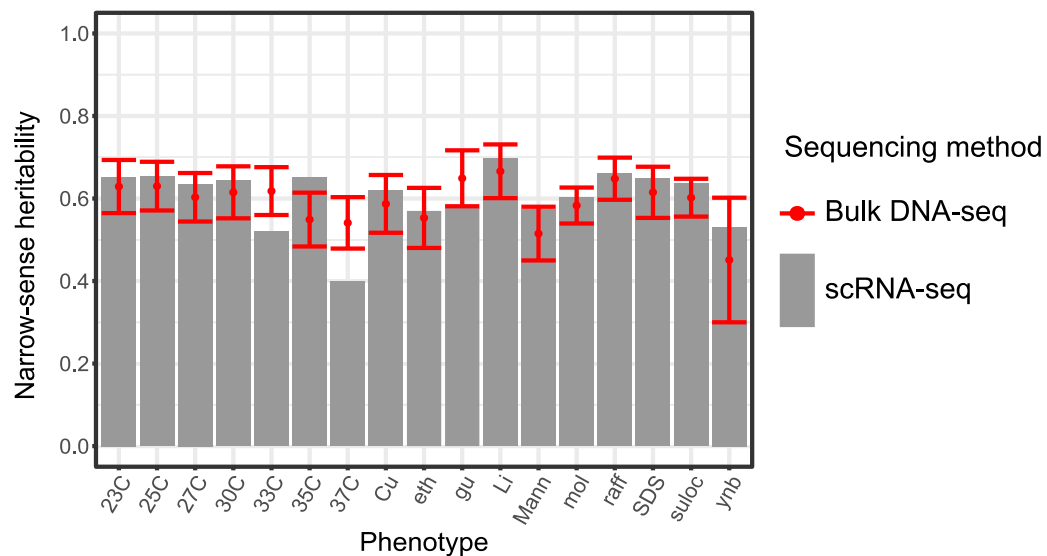
**A**



**B**



**Figure 2 Single-cell RNA-seq data recapitulate bulk DNA and RNA assays results.** A) Effect of the HMM on the relatedness between single cell genotypes and their closest reference lineage. The single-cell original genotype represents the genotype of the cells before the correction with the HMM. The relatedness to the closest lineage in batch1 has been measured with the adjusted $R^2$. To control for genotype uncertainty, only the 13,069 barcodes with a significant lineage assignment (lineage-barcode genotype correlation FDR<0.05) and a reference lineage with a lower uncertainty than the single cell HMM are selected, which represents 72.2% of the barcodes. We

169  then rounded the genotypes to remove the uncertainty during the comparison. Wilcoxon signed
170  test p-value is indicated above the violin plots. B) Narrow-sense heritability measured with non-
171  multiplexed DNA sequencing and scRNA-seq. The grey bars represent the scRNA-seq estimates
172  of narrow-sense heritability while the red dots represent the estimates from bulk DNA sequencing.
173  The interval of confidence of the bulk DNA sequencing is indicated by the red line around the red
174  dot and was obtained from genotype and phenotype measurement error in the BB-QTL paper (24).
175  The 23C-37C represents the temperature for the competition assay in YPD media while the other
176  phenotypes represent growth on YNB, molasses (mol), mannose (Mann) or raffinose (raff) and
177  chemical resistance to copper sulfate (Cu), ethanol (eth), guanidinium chloride (gu), lithium
178  acetate (Li), Sodium dodecyl sulfate (SDS) and suloctidil (suloc) (24).

179

180      To further establish that the genotyping obtained from scRNA-seq data was comparable to
181  previous non-multiplexed genotyping of the reference genotype panel, we estimated the
182  contribution of genetic variation to the phenotypic variation, i.e. fitness heritability. Nguyen Ba
183  and collaborators (2022) estimated the narrow- and broad-sense heritabilities of complex
184  phenotypes associated with temperature gradient, carbon source and chemical resistance for which
185  RM and BY segregants exhibit a significant level of diversity (24). We used our lineage assignment
186  to that panel to obtain fitness but used our single-cell genotyping to perform this association.
187  Encouragingly, most GCTA-REML estimates of narrow-sense heritability are within the
188  confidence intervals of Nguyen Ba and collaborators (2022) estimates (**Figure 2B**).

189      Although the variance partitioning is consistent with previous studies, it only provides a
190  broad view of the genotype-phenotype map as it does not allow to identify the loci that significantly
191  explain phenotype variation. If the genotypes obtained by scRNA-seq were of high-quality, then
192  we would expect that a QTL mapping model from scRNA-seq would yield a similar model than
193  non-multiplexed DNA sequencing data. To do so, we used a cross-validated stepwise forward
194  linear regression on the strain fitness and consensus genotypes data from single-cells that shared
195  the same lineage assignment (**Methods**). Performing the QTL mapping on the batch 1 scRNA-seq
196  dataset enabled the identification of 29 QTL compared to 31 QTL identified with the bulk barcoded
197  approach (**Tables S2 and S3**) (24). These QTL were largely similar as shown by the non-
198  significant difference between the effect sizes (Wilcoxon signed rank test $p = 0.29$) and by a model
199  similarity metric (24) that considers the recombination distance between matched QTL, the
200  similarity of the effect sizes and the allele frequencies (**Methods**). Using this approach, we
201  estimated that the similarity score between the batch 1 single cells QTL and the batch 1 BB-QTL
202  is 86.2% while each model respectively had a similarity score of 78.7% and 78.2% with the full
203  BB-QTL mapping performed on 99,950 segregants (24) (**Figure S4**). The QTL identified from the
204  scRNA-seq dataset also recapitulated several important biological features of the reference panel
205  such as an enrichment of non-synonymous and disordered region QTL (24) (**Figure S5**).

206      Finally, the variance partitioning model can also be modified to include gene expression as
207  the response variable and cell genotypes as the only random effect (**Methods**). This enables the
208  quantification of expression heritability, i.e. the variance of expression explained by genotype.
209  Using this approach, we estimated that genotype explains 72.3% of expression variance, which is

210 consistent with results from previous non-multiplexed eQTL mapping studies. Indeed, Albert and
211 collaborators (2018) estimated that genotype explains 70% of expression variance using a dataset
212 of 5720 genes in 1012 yeast segregants generated by the same parental strains (RM and BY).

**Integrating scRNA-seq data to an existing GPM highlights selection on the transcriptome**

214 Having shown that scRNA-seq is consistent with non-multiplexed assays while being more
215 scalable, we next sought to highlight new associations within the BY/RM GPM. Selection is often
216 highlighted at the genotype level through convergent evolution, increase in allele frequency within
217 a population or population genetics metric (26–28). However, the central dogma of molecular
218 biology and evolution tinkering entail that phenotype variation should be linked to transcriptomic
219 variation. As our dataset included all these variables, we sought to provide a variance partitioning
220 framework to evaluate the association between the transcriptome and trait variation (**Methods**)
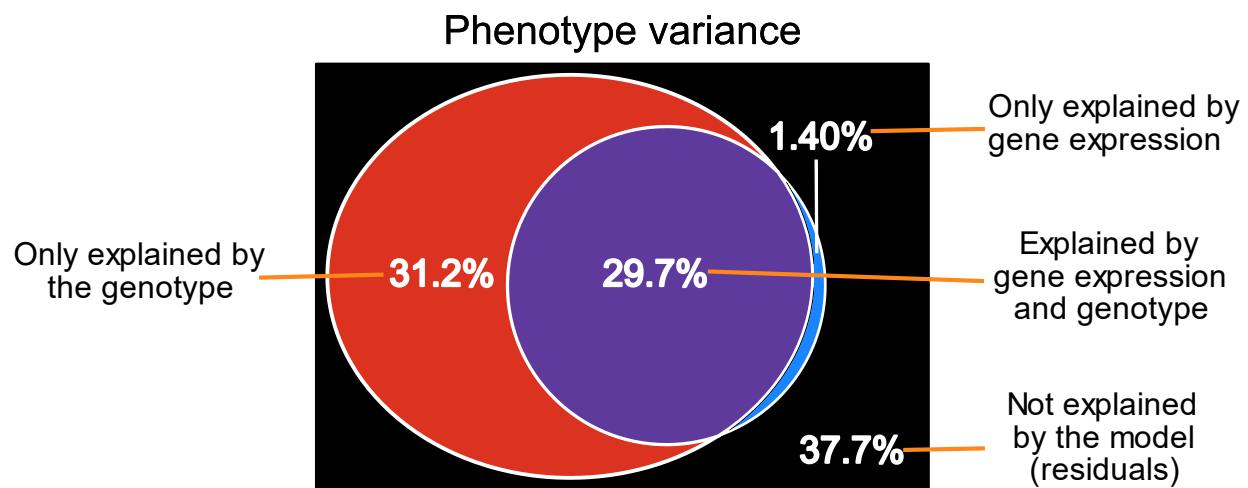221 with the 30C phenotype as an example (**Figure 3**).



**Figure 3 Variance partitioning of the 30C phenotype from scRNA-seq data.** The percentages
represent the proportion of fitness variance (whole rectangle area) explained by the components.
The ellipse area represents the phenotype variance explained by genotype variation and the circle
area represents the phenotype variance explained by expression variation. The black area of the
rectangle represents the residual of the model while the other colored areas represent the shared
and exclusive components explaining fitness variation.

230 The components of this variance partitioning all relate to at least one biological phenomenon.
231 Indeed, the portion of trait variation explained exclusively by the genotype variation (red in **Figure**
232 **3**) represents the effect of mutations on fitness via several biological phenomena such as protein
233 stability, enzymatic function etc, independent of expression level. For the 30C phenotype, this
234 component explains 31.2% of the fitness variation in the BY/RM background which is similar to
235 the 29.7% explained by the shared component between phenotype, genotype and expression
236 variations (purple in **Figure 3**). The latter represents the association between selection (fitness)
237 and the transcriptome either through loci influencing fitness via expression directly or through loci

238 affecting expression via an effect on cell fitness (indirectly) (29,30). Its considerable association
239 to fitness variation thus supports the evolution tinkering model. As for the phenotype variation
240 explained exclusively by gene expression (blue in **Figure 3**), it could represent epigenetics and
241 stochastic gene expression, which weakly explain variations in the 30C phenotype.

242 Although this model accurately estimates the narrow-sense heritability of 30C, the
243 residuals still represent 37.7% of fitness variation. This could be explained by unmeasured factors
244 like high-order epistasis, mitochondrial mutations or protein properties but the broad-sense
245 heritability of this phenotype is similar to the narrow-sense heritability, suggesting that the
246 residuals are mostly not explained by genotype and expression (24). Nguyen Ba et al. (2022) also
247 estimated that epistasis only explained around 5% of fitness (24). These results suggest that a
248 single run of scRNA-seq on a single batch of yeast segregants converge with bulk DNA sequencing
249 results while revealing previously hidden components of the GPM.

**250 Revealing hidden components of the yeast GPM with scRNA-seq**

251 Our integrative scRNA-seq approach is not limited to enabling the quantification of the association
252 between transcriptomic changes and trait variation. Indeed, the same approach we used to identify
253 QTL can be used to detect loci regulating gene expression which can reveal the cell mechanisms
254 underlying trait variation through transcriptomic changes. We thus modified the QTL mapping
255 framework such that the response variable is the level of expression of a single gene in the single
256 cells (**Methods**). This approach is a cost-efficient way to perform eQTL mapping from the
257 expression profile and genotype of cells from thousands of lineages in a multiplexed way (sc-
258 eQTL mapping).

259 Consistent with yeast non-multiplexed eQTL results, the genes with the highest expression
260 heritability are enriched in functions related to carbohydrate catabolic process (GO:0016052) and
261 cellular biosynthetic process (translation GO:0006412, organelle assembly GO:0070925,
262 ribosome biogenesis GO:0042254 and gene expression GO:0010467) (Fisher's exact test
263 FDR<0.05; **Methods**). In both datasets, these genes are also highly expressed, which reflects the
264 positive correlation between expression heritability and expression levels ($R^2 = 0.66$ and $p < 2.2e-$
265 16). Conversely, genes with the lowest expression heritability observed in the RM/BY background,
266 which we defined as the bottom 10% expression heritability, are enriched in functions related to
267 the cell cycle biological process (GO:0007049, Fisher's exact test FDR<0.05) (16,31).

268 Because of the increased scale of our collection, our approach is more powered to estimate
269 the gene heritability. We were thus able to detect new overrepresented biological processes, i.e.
270 DNA metabolic process (GO:0006259) and the response to nutrient levels (GO:0031667), for
271 which the variation of expression levels is weakly associated to the genetic variation observed
272 across the RM/BY segregants.

273

274 The functional enrichment analysis using scRNA-seq data revealed new associations
275 between expression heritability and biological processes in the RM/BY genetic background.
276 However, while it suggests that many eQTL are also QTL, it cannot accurately point to the specific

277    loci involved in trait variation and cannot address whether mutations on regulatory hubs have
278    stronger effects on traits. To investigate this, we mapped the QTL to hotspots of gene regulation
279    (or regulatory hubs), which we defined as 25 kb genomic windows that were repeatedly identified
280    in the eQTL mapping procedure (for different genes). This was done to acknowledge the
281    uncertainty in the exact position of the eQTL due to linkage disequilibrium and power. We then
282    ranked the 30C QTL identified by Nguyen Ba and collaborators (2022) based on their absolute
283    effect size and correlated it to the rank of the eQTL hotspots based on the number of regulated
284    genes. This resulted in a positive correlation (Spearman $\rho = 0.33$ and $p = 5.21e\text{-}5$), suggesting that
285    larger effects on the regulatory network translate into larger trait variation. Indeed, we observed
286    that some previously reported high-effect-size QTL genes are located in eQTL hotspots, eg MKT1,
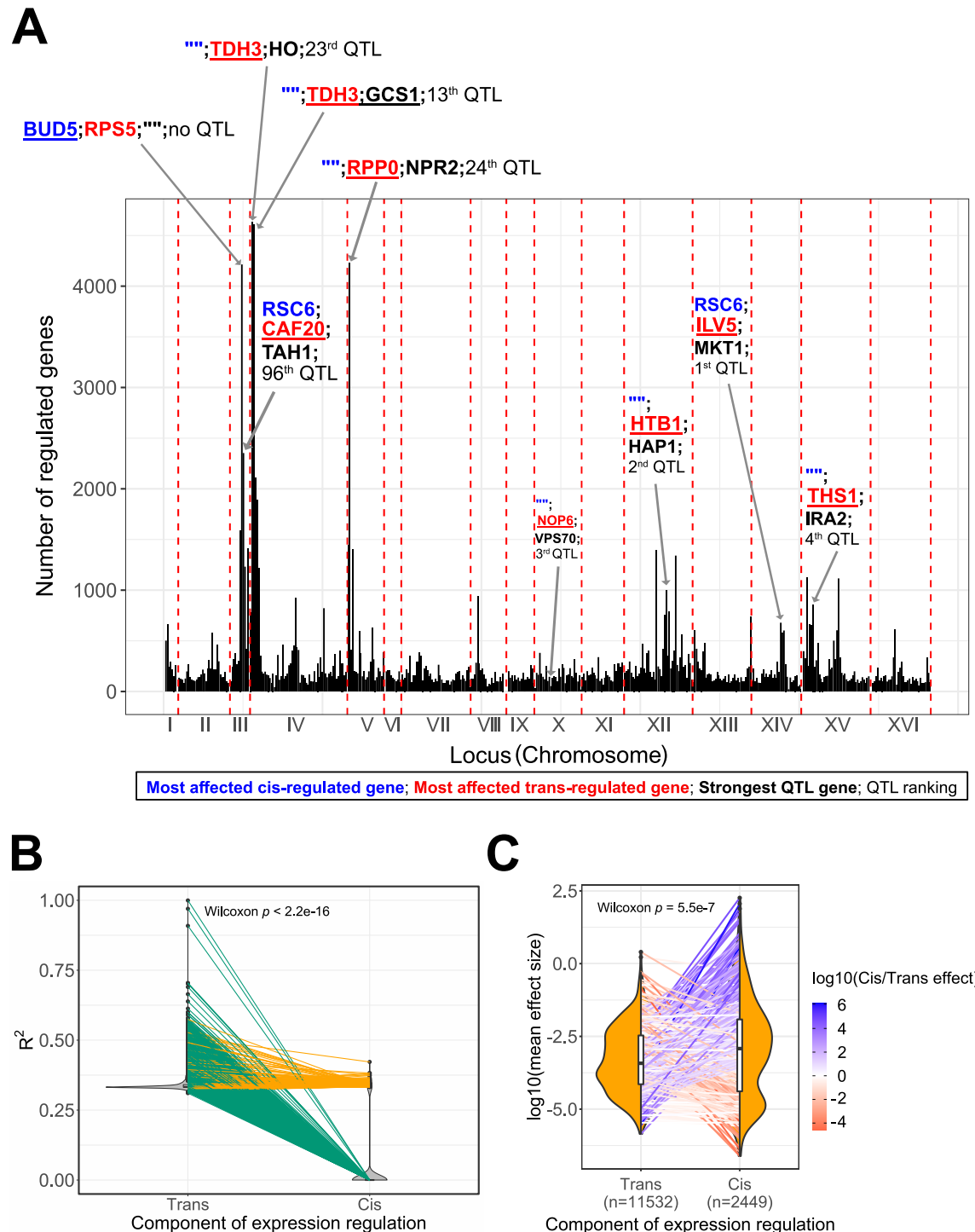287    HAP1, and IRA2 (**Figure 4A**).

**Figure 4: eQTL features underlying trait variation across the BY/RM segregants.** A) Mapping of the 30C QTL in the eQTL hotspots. We represent the hotspots of expression regulation as genomic windows (25 kb) to acknowledge the uncertainty around the real position of the eQTL due to linkage disequilibrium. We annotated the 5 top eQTL hotspots and the eQTL hotspots in which the top additive QTL identified by the BB-QTL mapping of the 30C phenotype are located. In these regions, we represented the most affected trans-regulated genes in red, the most affect cis-

294      regulated gene in blue and the genes of the top QTL in black. The double quotation characters
295      represent the absence of such genes in the associated region. We also represented the rank of the
296      QTL in the set of 159 QTL of the 30C phenotype. B) Partitioning of the expression heritability or
297      explained variance ($R^2$) among cis- and trans-eQTL. Each pair of points connected by a line
298      represents a gene. Green lines represent the genes that are only have trans-eQTL and orange lines
299      represent the genes that have both trans- and cis-eQTL. C) Comparison of the mean effect size
300      between cis- and trans-eQTL. Each pair of points connected by a line represents a gene. The ratio
301      of the average effect size between cis- and trans-eQTL is represented by the line color. The sample
302      size of each eQTL category is represented in the x axis. This is the number of trans-eQTL and cis-
303      eQTL used for calculating the average effect sizes per gene not the number of points per
304      distribution.

305      Performing this rank-test on individual genes also yielded the result that eQTL effect is correlated
306      with fitness effect for 35.1% of the genes (permutation test $p < 0.05$, see **Methods**). Although this
307      correlation does not apply to most genes, it reveals potential regulatory mechanisms explaining
308      the importance of the strongest growth loci or QTL. For instance, MKT1, i.e. the strongest growth
309      loci, is part of a regulation hotspot affecting genes that are important for yeast growth like ENP1
310      which is involved in RNA processing and HXT6 which is involved in glucose uptake (32–34).
311      Among the strongest growth loci, VPS70 is part of a hotspot of regulation that strongly affects the
312      expression of RSF2, a zinc-finger protein regulating glycerol-based growth and respiration (35).
313      Furthermore, the highest peak for expression regulation contains important growth loci in
314      chromosome IV around the mating type loci. This suggests the presence of cells with different
315      mating types in the dataset which we confirmed from the read mapping to Mat-a and Mat-α genes.
316      This is consistent with previous budding yeast eQTL mapping and is also expected because the
317      mating types in yeast express sets of genes that are "turned off" in other mating types (15,16,36).
318      This peak of expression regulation is also responsible for regulating TDH3 which is involved in
319      glycolysis and glucogenesis and can have important effect on fitness (37).

320         These hotspots suggest that expression differences in BY/RM would predominantly be due
321      to mutations in trans-regulatory elements. To test this, we partitioned the variation in gene
322      expression between cis- and trans- regulatory loci for each gene (see Methods). This analysis
323      revealed that all the genes are affected by at least one polymorphic trans-regulatory locus and that
324      these polymorphic trans-regulatory loci explain most of that gene's expression (**Figure 4B**). It is
325      well known that mutations in promoters and nearby enhancers can influence gene expression
326      (38,39). Indeed, we identified many genes that contained an allele in a cis-regulatory element that
327      strongly explain that gene's expression variation (n=750 genes out of 6088, **Figure 4B**). As
328      expected, mutations in cis-regulatory elements were of stronger effect size than trans-eQTL
329      individually, but the cumulative aggregate effect of all trans-eQTL acting on that gene was
330      comparable to the few cis-eQTL they had (**Figure 4C**). This can be explained by the fact that there
331      are more opportunities for mutations to arise in trans-regulatory elements. Finally, we found that
332      trans-eQTL have two times higher odds of affecting cell fitness than cis-eQTL ($\chi^2 p = 0.01$).

333  Taken together, the link between the genetic basis of transcription variation across RM/BY
334  segregants and fitness could only be revealed by integrating large-scale transcriptomic data to an
335  existing GPM, which scRNA-seq facilitates.

336

337  **CONCLUSION**

338  By leveraging the scalability of scRNA-seq, we obtained thousands of transcriptomes from a
339  reference pool of strains in a single experiment. This enabled the analysis of association between
340  genotype, transcriptome, and phenotype at an unprecedented scale. Questions surrounding
341  transcriptomic variation and phenotypic variation have been at the center of many previous
342  quantitative genetics studies (15,16,22,36,40). These ideas and discoveries all support the fact that
343  researchers can gain valuable insight about the evolution of traits by integrating the transcriptome
344  in GPM analyses, which can translate into fundamental knowledge or other important applications
345  where phenotypes evolve.

346  In this study, we took advantage of a previously characterized BY/RM cross where the genetic
347  basis of growth in various environments was examined in detail (24). By integrating transcriptomic
348  data in this genotype-phenotype map, we revealed how transcriptomic components are involved
349  in trait variation. Similar to a previous study, which obtained transcriptomes by individual strain
350  sequencing, we found that gene expression is highly heritable. Further, our study design also
351  allowed us to conclude that gene expression contributes to a significant portion of the phenotypic
352  variation in this strain collection.

353  This finding is corroborated by our findings that most eQTL detected in our study were previously
354  shown to be QTL. This is perhaps not surprising given that QTL in this cross were previously
355  inferred to be in regulatory genes, but this provides a more mechanistic view of the effect of an
356  allele on phenotype. Indeed, we find a bias for trans-regulation for generating transcription
357  innovation where the cumulative effect of trans-eQTL on gene expression are significant. That is
358  not to say that cis-regulatory alleles are dispensable as cis-regulatory alleles often have large effect
359  on gene expression. This genome-wide view of the genetic basis of transcriptional variation has
360  consequences for the evolution of phenotypes, as the target size afforded by trans-eQTL is far
361  larger than cis-eQTL. Thus, adaptation to small and fluctuation environmental changes may
362  proceed preferentially through allelic changes or recombination of many small-effect trans-eQTL,
363  but large expression changes are likely to require some cis-eQTL.

364  In this study, we leveraged the fact that our pool of strain was previously genotyped and
365  phenotyped. This was obtained by liquid handling robotics and pooled competitive growth assay
366  with barcode sequencing. While this was performed on a very large scale, it was essentially
367  obtained by brute-force and through approaches that are not necessarily applicable to other
368  systems. While it is clear from our results that genotyping single-cells can achieve the same
369  genotype quality as single-reaction genotyping, it is much harder to obtain phenotyping data from
370  scRNA-seq. Thus, our framework might not be readily translatable to other systems where similar
371  studies on the GPM are desirable. However, two observations from this cross can be used to
372  suggest an experimental approach. First, while epistasis is important, it contributes to a relatively

373 small portion of the phenotypic variance. Further, transcriptomic variation contributes little to the
374 missing heritability. Thus, it may be possible to use predicted fitness instead of observed fitness
375 and recapitulate essentially similar results as this study. Predicted fitness could be obtained from
376 bulk-segregant analysis where the additive effect of loci can be inferred from whole-genome
377 sequencing (23,41). While it is not clear if these observations are generalizable, it may be possible
378 to verify this for a study system of interest with some modest time-course single-cell based
379 sequencing where low-coverage genotyping is possible.

380 However, despite the study's limitation on generalizability, our scRNA-seq framework helps
381 bridge understanding of how genetic variation influences transcriptomic variation. Our framework
382 relies on identifying the genome of single-cells from the transcriptome, which is going to be
383 possible from low-coverage sequencing when genetic variation within the pool is high (such as
384 this cross, microbiome sequencing, or cancer cells with extensive copy number variation), and
385 from low cell diversity with sufficient transcriptomic variation such that aggregation of single-
386 cells with similar transcriptomes can afford pseudo-high coverage sequencing. Thus, integrating
387 genotype, transcriptome, and phenotype using scRNA-seq data can be particularly efficient for
388 developing a more fundamental understanding of other important traits or diseases.

389

390 **MATERIAL AND METHODS**

391 **Yeast strains and segregants**

392 We analyzed cells from a single batch (batch 1) of 4,489 segregants obtained from a F1 cross
393 between the yeast laboratory strain BY4741 and the vineyard strain RM11-1a generated in a
394 previous study (24). These strains have been selected to generate this collection of segregants
395 because they exhibit differences in multiple phenotypes including the adaptation to temperature,
396 the ability to process different sources of carbon and the ability to resist antifungal compounds.
397 Therefore, the genetic variation observed across the segregants can be correlated to the differences
398 in growth rate observed in the 18 environments recapitulating these phenotypes in the Nguyen et
399 al (2022) study (24). The selection of the batch is random and the fact that we performed the
400 analyses on a single batch eliminates batch effects that could obscure variable associations.
401 Genotypes and fitness data used were the same ones obtained in the previous study.

402 **Yeast growth and single-cell RNA-sequencing protocol**

403 To prepare strains for scRNA sequencing, we unfroze the batch of segregants and inoculated
404 approximately 5*10^6 cells in YPD (1% Yeast Extract, 2% Peptone, 2% Dextrose) to saturation.
405 The next day, about 10^7 cells were passaged to 5 mL of fresh YPD and grew for 4 hours to bring
406 cells to log-phase. We then pelleted 100 ul of cells and resuspended them in spheroplasting solution
407 (5 mg/mL zymolyase 20T, 10 mM DTT, 1 M Sorbitol, 100 mM Sodium Phosphate pH 7.4) at a
408 concentration of 10^7 cells/mL. The cells were incubated at 37 degrees Celcius for approximately
409 10 minutes at which point spheroplasting was verified by mixing a small aliquot of cells with
410 detergent to observe lysis. The cells at this point were quantified using a hemocytometer and
411 prepared using the standard 10x Genomics Gel Beads-in-emulsion (GEM) protocol. We used the
412 Chromium Next GEM Single-cell 3' Reagent Kit to prepare the sequencing libraries and
413 sequenced on a NextSeq 500 high-output flow cell.

414 We note that the cells analyzed here were grown in bulk and assayed for their transcriptome in
415 log-phase. Our fitness data was obtained from competitive bulk fitness assays which includes
416 several whole growth cycle over multiple days and thus captures lag phase, exponential growth,
417 and saturation. Nevertheless, previous experiments had shown that fitness was mostly determined
418 by exponential growth which suggests that our analysis is adequate even if the cells were prepared
419 for sequencing at a single time point.

420 **Single-cell RNA-sequencing data parsing**

421 From the scRNA-seq reads, we obtained gene expression levels and allele counts using the pipeline
422 count from CellRanger version 3.1.0 (42). For each of the ancestral strain, i.e RM11-1a and
423 BY4741, the pipeline mapped the scRNA-seq reads to the reference genome, filtered the barcodes
424 by comparing the UMI count per barcode distribution to a background model of empty gel-bead
425 in-emulsion, and counted the number of UMI per gene per barcode. The barcode filtering retained
426 18,233 barcodes. For each barcode, we then counted the number of RM and BY alleles at each
427 polymorphic site by parsing the RM and BY bam files using a python script
428 (https://github.com/arnaud00013/sc-eQTL/tree/main/II_scRNA-seq_genotyping). This script only

429 keeps reads that mapped at the same loci on both reference genomes to increase the level of
430 confidence of the mapping.

**Correction and imputation of single-cell genotypes with a Hidden Markov Model**

432 Because there are only two possible alleles at each polymorphic sites of the RM/BY segregants,
433 their genotype can be recapitulated by a quantitative variable measuring the proportion of reads
434 from one of the parental strains, which is RM in our dataset. The raw allele count data provides a
435 first estimate of this RM allele frequency at each polymorphic site. However, due to the low mean
436 depth of coverage of scRNA-seq data (0.2x), the absence of reads in some polymorphic sites and
437 the biases introduced during sequencing like index hopping/swapping, we expect that the raw data
438 can be imputed and corrected for errors and uncertainty in the observed alleles. Therefore, we
439 applied a Hidden Markov Model (HMM) on the observed allele count. Such model can infer
440 accurate genotype data at sequencing depths as low as 0.1x (24,25,43). Nguyen Ba and
441 collaborators (2022) designed an HMM to infer the segregants genotypes from bulk DNA
442 sequencing by accounting for sequencing error rate, recombination rate and index swapping rate
443 (24). Because scRNA-seq uses the reverse transcriptase, which has a higher error rate, and because
444 it is a pooled assay with higher chances of index swapping, we expected the HMM parameter to
445 differ for the single cell data. Therefore, we adapted the HMM to scRNA-seq data by measuring
446 its parameters in our dataset (**Figure S1**). The scripts are available on GitHub
447 (https://github.com/arnaud00013/sc-eQTL/tree/main/II_scRNA-seq_genotyping).

448

**Assigning single cells to the reference panel strains**

450 To evaluate the level of relatedness between the reference panel strains and the imputed single cell
451 genotypes, we used the expected distance to identify the strain that best relate to each single cell:

$$Expected\ distance(g_c, g_s) = \sum_{i=1}^{41594} g_c + g_s - 2g_c g_s \qquad (Eq.1)$$

453 where $g_c$ is the cell genotype and $g_s$ is the strain genotype. Next, we assigned the single cell to its
454 best match in the studied batch of 4,489 trains only if this match is better than the best match in
455 randomly generated batches of the same size (**Figure S2**). This procedure is implemented and
456 available at https://github.com/arnaud00013/sc-eQTL/tree/main/III_Genotype_analysis).

**Partitioning the phenotypic variance into genetic and transcriptomic components**

458 To analyze the yeast GPM at a broad scale and to evaluate the association between selection and
459 the transcriptome, we estimated the contribution of genetic and transcriptomic variations to
460 phenotypic variation from scRNA-seq data. More precisely, we performed a Genome-wide
461 Complex Trait Analysis (GCTA) by fitting a linear mixed model to the data using the restricted
462 maximum-likelihood (REML) method (44):

$$y = X\beta + W_g u_g + \varepsilon_g \qquad (Eq.2)$$

$$y = X\beta + W_e u_e + \varepsilon_e \qquad (Eq.3)$$

465
$$y = X\beta + W_g u_g + W_e u_e + \varepsilon \tag{Eq.4}$$

466 where $y$ is the fitness vector for the $n$ cells, $X$ is the $nxk$ matrix of $k$ fixed effects, $\beta$ is the vector of
467 $k$ coefficients of the fixed effects, $W_g$ is the $nxp$ genotype matrix, $u_g$ is the vector of $p$ SNP effects,
468 $W_e$ is the $nxm$ expression matrix, $u_e$ is the vector of $m$ gene expression effects and $\varepsilon$ is the error
469 term. Because the dataset does not include fixed effects, we set the fixed effect to a vector of ones
470 such that its coefficients represent the mean fitness while the genotype and expression data are the
471 random effects that explain the fitness variance along with the error terms. The REML solution
472 assumes that the data follow a Gaussian distribution, so the data are standardized before fitting the
473 model. We also divided the standardized expression counts by the cell sum of expression counts
474 to control for molecule count biases across cells. The cell fitness is based on the fitness of the
475 closest segregant in batch 1 as measured by the expected distance. Because this model is linear
476 and additive, it can be compared to the estimates of narrow-sense heritability obtained by Nguyen
477 Ba and collaborators (2022) (24). The difference between the variance explained in equation 4 and
478 equations 2 or 3 allow to infer the variance explained only by the genotype or the expression
479 component of the model. The code for the variance partitioning is available on GitHub
480 (https://github.com/arnaud00013/sc-eQTL/tree/main/IV_variance_partitioning).

481

482 **Estimating the expression heritability from scRNA-seq**

483 To obtain this estimate from scRNA-seq data, we needed to consider the fact that GCTA-REML
484 only takes a vector as a response variable while the gene expression matrix is multi-dimensional.
485 To solve this, we orthogonalized the gene expression matrix using principal component analysis
486 (PCA), and used each of the PC one at a time as a response variable of the model. Indeed, if the
487 expression PCs recapitulate the total expression variance and are orthogonal or independent to
488 each other, then the sum of the PCs variance explained by genotype should be the expression
489 heritability. To save time, we only used the 898 expression PCs that explain 99% of expression
490 variance:

491 $$Expression\ heritability = \sum_{i=1}^{898} PC_i\ eigen\ value * "PC_i \sim genotype"\ model\ R^2 \tag{Eq.5}$$

492

493 **QTL mapping**

494 To identify the loci that influence cell fitness, we performed a linear regression on the consensus
495 genotypes of the strains from the single cell data and the strain fitness. We decided to use the
496 consensus genotypes of the strains as they relate better to the bulk segregant genomes. To build
497 the consensus genotypes, we defined cells from the same lineage as the ones that shared the same
498 closest segregant in batch 1. Next, we used the median to obtain cells' consensus genotypes as it is
499 less sensitive to outliers and because it yields the best relatedness to the batch 1 reference
500 genotypes (median $R^2$ = 87.0%; $\mu$=79.5%; $\sigma$=18.2; **Figure S3**). We selected the QTL in the linear
501 models using cross-validation on the scRNA-seq data. This analysis consists in dividing the dataset
502 into 10 random partitions of similar sample sizes and running a cross-validated stepwise forward

503    linear regression on each partition. For each partition, the model starts with no QTL and a linear
504    model "Fitness ~ Genotype" is fitted using the genotype data at each polymorphic site, where the
505    correlation coefficient represents the effect size of the SNP. Then, the forward search starts and at
506    each iteration, a new locus with the minimum linear model residual sum of squares (RSS) is added
507    to the QTL model, which is updated with new effect sizes after the addition of a new SNP. Because
508    the order of addition of QTL matters in the forward search and because some QTL are linked or
509    collinear, the model can be refined by exploring different QTL around the local optima. These
510    steps are repeated until the model RSS cannot be improved anymore or until the number of QTL
511    reaches an arbitrary maximum far from the cross-validated number of QTL. After the forward
512    search is completed in each partition, the algorithm calculates the optimal $\lambda$ values that minimizes
513    the objective function $F_o$:

514 $$F_o(\beta) = RSS(\beta) + Lasso\ penalty(\beta)$$

515 $$\|Y - X\beta\|_2^2 + \lambda\|\beta\|_0 \text{ (Eq.6)}$$

516    where $\beta$ is the vector of SNP effect sizes in the QTL model, $\|Y - X\beta\|_2^2$ is the RSS of the linear
517    QTL model, $\lambda$ defines the penalty for adding a new SNP to the model and $\|\beta\|_0$ is the number of
518    SNPs in the QTL model. This objective function has the property to add sparsity in the QTL model
519    and thus avoid overestimating the number of QTL while being consistent (24). The optimal $\lambda$ has
520    a minimum of log(n) which corresponds to the Bayesian Information Criterion (BIC), which is
521    known to yield correct models asymptotically (45). This allows to consider the possibility that a
522    sparser model than the one found using the BIC could yield better predictive power on a test set
523    while avoiding overfitting. The optimal $\lambda$ values found in all the partitions are then averaged and
524    the resulting mean $\lambda$ is used to solve the objective function in the full dataset, which yields the
525    optimal QTL model. The cross-validation assumes that the partitions are independent, such that
526    the variance explained by the model and the number of relevant QTL are unbiased estimates.

527    **Highlighting hotspots of gene regulation through eQTL mapping**

528    To identify the loci regulating gene expression regulation, we adapted the QTL mapping
529    framework using expression as the predicted phenotype. Because this approach had to be repeated
530    for each of the 6,240 genes, we needed to modify it so that the execution time is convenient. To
531    do so, the parameter $\lambda$ was not estimated using cross validation but rather from the Bayesian
532    Inference Criterion (BIC), i.e. $\lambda = log(n)$ where n is the number of cells. We found that the BIC
533    was often selected by the cross-validation procedure when tested on a few genes and thus we do
534    not believe that this approach will significantly change our results.

535    To acknowledge the uncertainty around the exact position of eQTL due to linkage disequilibrium,
536    we define eQTL hotspots as 25 kb genomic windows that were repeatedly identified in the eQTL
537    mapping procedure. The code for the single cell eQTL mapping is available on GitHub
538    (https://github.com/arnaud00013/sc-eQTL/tree/main/V_sc_eQTL_mapping).

539    **Functional enrichment analysis by gene ontology annotation**

540    To highlight gene functions enriched at different levels of expression or expression heritability,
541    we performed the panther database binomial test for statistical overrepresentation of gene ontology

542 biological processes (31,46). A low level was defined as within the 25% bottom part of the
543 distribution (<Q1) while a high level was defined as within the top 25% part of the distribution
544 (>Q3). The *p*-values were corrected for multiple testing using the false discovery rate correction
545 (FDR).

**Matching QTL to eQTL**

547 To evaluate the contribution of gene expression regulation to fitness variation, we created a model
548 to match QTL and eQTL based on the similarity of loci and the similarity of predicted effect on
549 gene expression. More precisely, for each of the 6,088 genes for which we could detect eQTL, we
550 performed a new eQTL model by correlating the expression level of the gene to the genetic
551 variation at QTL positions. This allowed us to measure the predicted effect of the QTL on gene
552 expression. We then calculated the distance between the QTL and the real eQTL of the gene based
553 on recombination distance within each chromosome, which decreases exponentially with genetic
554 distance, and the difference in the predicted effect on the gene expression using the formulation
555 developed by Nguyen Ba et al (2022) (24). Next, we used the same Needleman-Wunsch algorithm
556 to find the most likely set of pairing between QTL and eQTL, where an unmatched QTL is also
557 possible but penalized. Finally, we determined the proportion of genes for which gene expression
558 regulation is associated with higher fitness. To do so, for each gene, we performed a permutation
559 test by comparing the average rank of the matched QTL of the gene to the average rank of 999
560 random subsets of unmatched QTL of the same size. The p-value is the proportion of random
561 subsets of unmatched QTL with a higher average QTL rank than the set of matched QTL.

**Comparing cis- and trans-eQTL contribution to expression variation**

563 We used the definition of local eQTL in Albert et al. (2018) to define cis-eQTL, i.e. any eQTL
564 between 1,000 bp upstream of the gene and 200 bp downstream of the gene. Thus, we defined
565 trans-eQTL as the eQTL that do not follow this criterion. For each gene, we then performed
566 variance partitioning using the GCTA:

$$y = X\beta + W_{g\_cis}u_{g\_cis} + \varepsilon_{cis} \tag{Eq.7}$$

$$y = X\beta + W_{g\_trans}u_{g\_trans} + \varepsilon_{trans} \tag{Eq.8}$$

$$y = X\beta + W_{g\_cis}u_{g\_cis} + W_{g\_trans}u_{g\_trans} + \varepsilon \tag{Eq.9}$$

570 where *y* is the vector of expression level of the gene across the *n* cells, *X* is the *nxk* matrix of *k*
571 fixed effects, $\beta$ is the vector of *k* coefficients of the fixed effects, $W_{g\_cis}$ is the *nxp* cis-eQTL
572 genotype matrix, $u_{g\_cis}$ is the vector of *p* cis-eQTL effects on expression, $W_{g\_trans}$ is the *nxm*
573 trans-eQTL expression matrix, $u_e$ is the vector of *m* trans-eQTL effects on expression and $\varepsilon$
574 represent the error terms. Because the dataset does not include fixed effects, we set the fixed effect
575 to a vector of ones such that its coefficients represent the mean expression level while the cis-
576 eQTL and trans-eQTL genotypes are the random effects that explain the expression variance along
577 with the error terms. We can infer the variance explained by the cis-eQTL by the difference in
578 variance explained between the models in equations 9 and 8. Likewise, the difference of variance
579 explained by the models in equations 9 and 7 can help us estimate the variance explained by the

580    trans-eQTL. Finally, we estimate the effect sizes using the absolute value of the correlation

581    coefficients of each loci and compare the mean between the cis- and trans-eQTL from the same

582    gene (paired data) with a Wilcoxon signed rank test.

583

**DATA AVAILABILITY**

The code used for this study is available and explained at https://github.com/arnaud00013/sc-eQTL and the original single-cell reads from the pooled segregants scRNA-seq assay have been uploaded in the NCBI BioProject database with the accession number PRJNA1022775. The single-cell barcodes expression data are also available at https://github.com/arnaud00013/sc-eQTL as an archive file named Matrix_gene_expression_barcodes_1_to_9000.csv.tar.gz or Matrix_gene_expression_barcodes_9001_to_18233.csv.tar.gz.

**ACKNOWLEDGEMENTS**

**BIBLIOGRAPHY**

1. Bartoli C, Roux F. Genome-Wide Association Studies In Plant Pathosystems: Toward an Ecological Genomics Approach. Front Plant Sci [Internet]. 2017 May 23;8. Available from: ://WOS:000402030200001

2. Ferreira MA, Gamazon ER, Al-Ejeh F, Aittomaki K, Andrulis IL, Anton-Culver H, et al. Genome-wide association and transcriptome studies identify target genes and risk loci for breast cancer. Nat Commun [Internet]. 2019 Apr 15;10. Available from: ://WOS:000464494100010

3. Aguet F, Alasoo K, Li YI, Battle A, Im HK, Montgomery SB, et al. Molecular quantitative trait loci. Nat Rev Methods Primer. 2023 Jan 25;3(1):1–22.

4. Tarantino LM, Eisener-Dorman AF. Forward Genetic Approaches to Understanding Complex Behaviors. Curr Top Behav Neurosci. 2012;12:25–58.

5. Casamassimi A, Federico A, Rienzo M, Esposito S, Ciccodicola A. Transcriptome Profiling in Human Diseases: New Advances and Perspectives. Int J Mol Sci. 2017 Jul 29;18(8):1652.

6. Wainberg M, Sinnott-Armstrong N, Mancuso N, Barbeira AN, Knowles DA, Golan D, et al. Opportunities and challenges for transcriptome-wide association studies. Nat Genet. 2019 Apr;51(4):592–9.

7. Transcriptome: Connecting the Genome to Gene Function | Learn Science at Scitable [Internet]. [cited 2023 Aug 31]. Available from: https://www.nature.com/scitable/topicpage/transcriptome-connecting-the-genome-to-gene-function-605/

8. Williams CG, Lee HJ, Asatsuma T, Vento-Tormo R, Haque A. An introduction to spatial transcriptomics for biomedical research. Genome Med. 2022 Jun 27;14(1):68.

9. King MC, Wilson AC. Evolution at Two Levels in Humans and Chimpanzees. Science. 1975 Apr 11;188(4184):107–16.

10. Jacob F. Evolution and Tinkering. Science. 1977 Jun 10;196(4295):1161–6.

11. Hoekstra HE, Coyne JA. The Locus of Evolution: Evo Devo and the Genetics of Adaptation. Evolution. 2007;61(5):995–1016.

12. Kratochwil CF, Meyer A. Evolution: Tinkering within Gene Regulatory Landscapes. Curr Biol. 2015 Mar 30;25(7):R285–8.

13. Primig M, Williams RM, Winzeler EA, Tevzadze GG, Conway AR, Hwang SY, et al. The core meiotic transcriptome in budding yeasts. Nat Genet. 2000 Dec;26(4):415–23.

14. Cavalieri D, Townsend JP, Hartl DL. Manifold anomalies in gene expression in a vineyard isolate of Saccharomyces cerevisiae revealed by DNA microarray analysis. Proc Natl Acad Sci U S A. 2000 Oct 24;97(22):12369–74.

15. Brem RB, Yvert G, Clinton R, Kruglyak L. Genetic Dissection of Transcriptional Regulation in Budding Yeast. Science [Internet]. 2002 Apr 26 [cited 2022 Apr 29]; Available from: https://www.science.org/doi/pdf/10.1126/science.1069516

16. Albert FW, Bloom JS, Siegel J, Day L, Kruglyak L. Genetics of trans-regulatory variation in gene expression. Wittkopp PJ, editor. eLife. 2018 Jul 17;7:e35471.

17. Schwarz T, Boltz T, Hou K, Bot M, Duan C, Loohuis LO, et al. Powerful eQTL mapping through low-coverage RNA sequencing. Hum Genet Genomics Adv. 2022 Jul 14;3(3):100103.

18. Fan Y, Zhu H, Song Y, Peng Q, Zhou X. Efficient and effective control of confounding in eQTL mapping studies through joint differential expression and Mendelian randomization analyses. Bioinformatics. 2020 Aug 13;37(3):296–302.

19. Bush WS, Moore JH. Chapter 11: Genome-Wide Association Studies. Plos Comput Biol [Internet]. 2012 Dec;8(12). Available from: ://WOS:000312901500028

20. Lorenz K, Cohen BA. Small- and Large-Effect Quantitative Trait Locus Interactions Underlie Variation in Yeast Sporulation Efficiency. Genetics. 2012 Nov;192(3):1123-+.

21. Shaffer SM, Dunagin MC, Torborg SR, Torre EA, Emert B, Krepler C, et al. Rare cell variability and drug- induced reprogramming as a mode of cancer drug resistance. Nature. 2017 Jun 15;546(7658):431-+.

22. Bloom JS, Ehrenreich IM, Loo WT, Lite TLV, Kruglyak L. Finding the sources of missing heritability in a yeast cross. Nature. 2013 Feb;494(7436):234–7.

23. Ehrenreich IM, Torabi N, Jia Y, Kent J, Martis S, Shapiro JA, et al. Dissection of genetically complex traits with extremely large pools of yeast segregants. Nature. 2010 Apr;464(7291):1039–42.

24. Nguyen Ba AN, Lawrence KR, Rego-Costa A, Gopalakrishnan S, Temko D, Michor F, et al. Barcoded Bulk QTL mapping reveals highly polygenic and epistatic architecture of complex traits in yeast. Verstrepen KJ, editor. eLife. 2022 février;11:e73983.

25. Hwang B, Lee JH, Bang D. Single-cell RNA sequencing technologies and bioinformatics pipelines. Exp Mol Med [Internet]. 2018 Aug 7;50. Available from: ://WOS:000441266700002

26. Bergström A, Simpson JT, Salinas F, Barré B, Parts L, Zia A, et al. A High-Definition View of Functional Genetic Variation from Natural Yeast Genomes. Mol Biol Evol. 2014 Apr 1;31(4):872–88.

27. Good BH, McDonald MJ, Barrick JE, Lenski RE, Desai MM. The dynamics of molecular evolution over 60,000 generations. Nature. 2017/10/19 ed. 2017 Nov 2;551(7678):45–50.

28. Johnson MS, Gopalakrishnan S, Goyal J, Dillingham ME, Bakerlee CW, Humphrey PT, et al. Phenotypic and molecular evolution across 10,000 generations in laboratory budding yeast populations. Verstrepen KJ, Wittkopp PJ, Verstrepen KJ, Hodgins-Davis A, editors. eLife. 2021 Jan 19;10:e63910.

29. Sun XM, Bowman A, Priestman M, Bertaux F, Martinez-Segura A, Tang W, et al. Size-Dependent Increase in RNA Polymerase II Initiation Rates Mediates Gene Expression Scaling with Cell Size. Curr Biol. 2020 Apr 6;30(7):1217-1230.e7.

30. Marguerat S, Bähler J. Coordinating genome expression with cell size. Trends Genet. 2012 Nov 1;28(11):560–5.

31. Mi H, Muruganujan A, Huang X, Ebert D, Mills C, Guo X, et al. Protocol Update for large-scale genome and gene function analysis with the PANTHER classification system (v.14.0). Nat Protoc. 2019 Mar;14(3):703–21.

32. Roos J, Luz JM, Centoducati S, Sternglanz R, Lennarz WJ. ENP1, an essential gene encoding a nuclear protein that is highly conserved from yeast to humans. Gene. 1997 Jan 31;185(1):137–46.

33. Chen W, Bucaria J, Band DA, Sutton A, Sternglanz R. Enp1, a yeast protein associated with U3 and U14 snoRNAs, is required for pre-rRNA processing and 40S subunit synthesis. Nucleic Acids Res. 2003 Jan 15;31(2):690–9.

34. Roy A, Dement AD, Cho KH, Kim JH. Assessing Glucose Uptake through the Yeast Hexose Transporter 1 (Hxt1). PLOS ONE. 2015 Mar 27;10(3):e0121985.

35. Lu L, Roberts GG, Oszust C, Hudson AP. The YJR127C/ZMS1 gene product is involved in glycerol-based respiratory growth of the yeast Saccharomyces cerevisiae. Curr Genet. 2005 Oct 1;48(4):235–46.

36. Brem RB, Storey JD, Whittle J, Kruglyak L. Genetic interactions between polymorphisms that affect gene expression in yeast. Nature. 2005 Aug;436(7051):701–3.

37. Vande Zande P, Hill MS, Wittkopp PJ. Pleiotropic effects of trans-regulatory mutations on fitness and gene expression. Science. 2022 Jul;377(6601):105–9.

38. Mattioli K, Oliveros W, Gerhardinger C, Andergassen D, Maass PG, Rinn JL, et al. Cis and trans effects differentially contribute to the evolution of promoters and enhancers. Genome Biol. 2020 Aug 20;21(1):210.

39. Romero IG, Ruvinsky I, Gilad Y. Comparative studies of gene expression and the evolution of gene regulation. Nat Rev Genet. 2012 Jun 18;13(7):505–16.

40. Bloom JS, Boocock J, Treusch S, Sadhu MJ, Day L, Oates-Barker H, et al. Rare variants contribute disproportionately to quantitative trait variation in yeast. Landry CR, Barkai N, editors. eLife. 2019 Oct 24;8:e49212.

41. Brauer MJ, Christianson CM, Pai DA, Dunham MJ. Mapping novel traits by array-assisted bulk segregant analysis in Saccharomyces cerevisiae. Genetics. 2006 Jul;173(3):1813–6.

42. Zheng GX, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, et al. Massively parallel digital transcriptional profiling of single cells. Nat Commun. 2017/01/17 ed. 2017 Jan 16;8:14049.

43. Stark R, Grzelak M, Hadfield J. RNA sequencing: the teenage years. Nat Rev Genet. 2019 Nov;20(11):631–56.

44. Yang J, Lee SH, Goddard ME, Visscher PM. GCTA: a tool for genome-wide complex trait analysis. Am J Hum Genet. 2010/12/21 ed. 2011 Jan 7;88(1):76–82.

45. Harrell FE. Regression Modeling Strategies: With Applications to Linear Models, Logistic Regression, and Survival Analysis [Internet]. New York, NY: Springer; 2001 [cited 2023 Mar 29]. (Springer Series in Statistics). Available from: http://link.springer.com/10.1007/978-1-4757-3462-1

46. PANTHER: Making genome-scale phylogenetics accessible to all [Internet]. [cited 2023 Sep 29]. Available from: https://onlinelibrary.wiley.com/doi/epdf/10.1002/pro.4218

47. Kanehisa M, Goto S. KEGG: Kyoto Encyclopedia of Genes and Genomes. Nucleic Acids Res. 2000 Jan 1;28(1):27–30.