

Deep-learning-assisted Sort-Seq enables high-throughput profiling of gene expression characteristics with high precision

Huibao Feng^{1,*#}, Fan Li^{1#}, Tianmin Wang³, Xin-hui Xing¹, An-ping Zeng^{4,5}, Chong Zhang^{1,2,*}

¹ MOE Key Laboratory for Industrial Biocatalysis, Institute of Biochemical Engineering, Department of Chemical Engineering, Tsinghua University, Beijing 100084, China

² Center for Synthetic and Systems Biology, Tsinghua University, Beijing 100084, China

³ Tsinghua-Peking Center for Life Sciences, School of Medicine, Tsinghua University, Beijing 100084, China

⁴ Institute of Bioprocess and Biosystems Engineering, Hamburg University of Technology, Hamburg 21073, Germany

⁵ Center of Synthetic Biology and Integrated Bioengineering, School of Engineering, Westlake University, Hangzhou 310024, China

These authors contributed equally to this work.

* To whom correspondence should be addressed.

Huibao Feng, Email: fhb_14@163.com

Chong Zhang, Email: chongzhang@tsinghua.edu.cn

ABSTRACT

As an essential physiological process, gene expression determines the function of each cell. However, owing to the complex nondeterministic and nonlinear nature of gene expression, the steady-state intracellular protein abundance of a clonal population forms a distribution. The characteristics of this distribution, including expression strength and noise, are closely related to cellular behavior. Therefore, quantitative description of these characteristics is an important goal in biology. This task, however, has so far relied on arrayed methods, which are time-consuming and labor-intensive. To address this issue, we propose a deep-learning-assisted Sort-Seq approach (dSort-Seq) in this work, enabling high-throughput profiling of expression properties with high precision. We demonstrated the validity of dSort-Seq for large-scale assaying of the dose–response relationships of biosensors. In addition, we comprehensively investigated the contribution of transcription and translation to noise production in *E. coli*, from which we discovered that the expression noise is strongly coupled with the mean expression level instead of translation strength, even in the case of weak

1 transcription. We also discovered that the transcriptional interference caused by overlapping
2 RpoD-binding sites contributes to noise production, which suggested the existence of a
3 simple and feasible noise control strategy in *E. coli*. Overall, dSort-Seq is able to efficiently
4 determine the strength-noise landscape, which has promising applications in studies related to
5 gene expression.

6

7 INTRODUCTION

8 Cells are sophisticated instruments driven by the central dogma that build varieties of lives.
9 For each cell, gene expression is a vital process by which information from genes flows to
10 RNA and then to proteins, determining the traits of the cell. However, gene expression is
11 often stochastic, as it involves many random events requiring the participation of various
12 low-copy-number chemical components¹⁻⁶. In addition, this process can be chaotic due to the
13 high complexity of the regulatory network^{7,8}. As a result, phenotypic heterogeneity exists
14 among genetically identical cells even under the same environmental conditions¹. Therefore,
15 steady-state protein production in a clonal population exhibits a distribution, wherein the
16 mean of the distribution (Mean) indicates the expression strength, and the squared coefficient
17 of variation (CV^2) exhibits the expression noise. These two characteristics are both important
18 indicators that are closely related to the phenotypes of a population, such as the bioproduction
19 efficiency^{9,10}, drug resistance^{11,12} and antibiotic persistence¹³. To date, the quantitative
20 description of expression strength and noise has been an important goal in biology to
21 illustrate cellular behavior¹⁴. However, this task has relied on fluorescence microscopy^{1,15} and
22 flow cytometry^{14,16} (FCM) assays of individual clonal populations, which are time-
23 consuming and labor-intensive when testing large amounts of genetic variants. Therefore, a
24 general, precise and high-throughput method for the profiling of expression properties is
25 urgently needed.

26 To address the above issue, we focused on Sort-Seq¹⁷⁻¹⁹ (also named FlowSeq, FACS-seq),
27 by which a library of cells with different expression intensities can be sorted into different
28 bins and then quantified through next-generation sequencing (NGS) to derive the expression
29 pattern of each genotype. This approach has been broadly used in profiling sequence-function
30 relationships associated with transcriptional regulation^{17,20-22}, translational regulation^{17,18,23},
31 regulatory RNAs^{24,25}, protein-sequence interactions²⁶, etc. In addition, the validity of Sort-Seq
32 has been demonstrated in a wide range of organisms, including bacteria, yeast and

1 mammalian cells¹⁹. However, it remains difficult to derive precise expression characteristics
2 from Sort-Seq data. Existing methods have focused on fitting the binned distribution to a log-
3 normal^{18,19,24,27} or gamma distribution^{20,21,28}, which are limited by the inexact representation
4 capability of these probability densities^{6,29,30}. On the other hand, the parameter learning
5 process of these methods still needs to be improved. For instance, apart from the binned
6 distribution, other data, such as the overall fluorescence intensity density, should be
7 considered. Hence, to obtain expression properties with high throughput and high precision, a
8 common, rigorous data processing method for Sort-Seq is needed.

9 Therefore, we have developed dSort-Seq, a deep-learning-assisted Sort-Seq approach (**Fig.**
10 **1**). In this method, instead of using log-normal or gamma distribution, we applied a two-
11 component log-Gaussian mixture model (LGMM) to match the steady-state gene expression
12 density, which is more precise and robust in fitting the real data. To decode Sort-Seq data, for
13 the first time, we adopted a Bayesian neural network to perform parameter learning. These
14 innovations significantly improve the accuracy of Sort-Seq to derive expression
15 characteristics for thousands of variants. We demonstrated the validity of this pipeline from
16 two aspects. First, dSort-Seq enables large-scale assays of dose–response relationships of
17 biosensors with high precision, with which the optimal design can be efficiently identified.
18 Second, it also supports the high-throughput exploration of noise production mechanisms.
19 For instance, we applied dSort-Seq to determine the effects of transcription and translation on
20 expression noise in *E. coli* and found them to have comparable contributions, contradicting
21 the commonly accepted translational bursting mechanism³. In addition, we also revealed that
22 overlapping RpoD-binding sites would lead to high expression noise, which suggested an
23 effective noise regulation strategy. Overall, our method, which provides significant
24 mathematical and biological insights, can serve as a promising high-throughput tool for use
25 in various studies associated with gene expression.

26

27 **RESULTS**

28 **Framework and superior performance of dSort-Seq**

29 Recent research on the stochastic nature of gene expression has shown that steady-state
30 protein production in a clonal population follows a gamma (negative binomial)^{3,4} or log-
31 normal^{5,6} distribution. However, neither of them can precisely match the real expression data

1 (Fig. 2a-2c). To address this issue, dSort-Seq applied a two-component log-Gaussian mixture
2 model (LGMM) to represent the steady-state protein production density. This distribution
3 was selected for several reasons, the first and foremost of which is that the mixture of
4 Gaussians can theoretically approximate any continuous density given enough components³¹,
5 ensuring its ability to fit more complex densities compared with conventionally used models.
6 In addition, the outliers (in more extreme cases, one peak of the bimodal expression
7 densities³⁰), which have a great impact on matching^{19,29}, can be viewed as being generated by
8 a Gaussian component³²⁻³⁴. To verify the model's ability to match expression distributions,
9 we compared it with gamma and log-normal distributions in fitting quantitative datasets from
10 independent resources^{27,35} (Fig. 2a-2c). Our method exhibited higher precision in
11 representation, and hence, we used the LGMM for subsequent analyses.

12 Next, to derive gene expression characteristics from Sort-Seq data, we revisited the
13 experimental procedure and considered incorporating more data into the parameter learning
14 step (Fig. 1). To intuitively represent the dSort-Seq method, we defined the following terms:
15 (1) the mixing coefficients, denoted by $\boldsymbol{\pi} = (\pi_1, \pi_2, \dots, \pi_n)$, where π_i is the proportion of the
16 i th variant in the library; (2) the log-scaled sorting boundaries, denoted by $\mathbf{b} =$
17 $(b_0 = -\infty, b_1, \dots, b_K = +\infty)$; (3) the parameters involved in LGMM, denoted by $\boldsymbol{\lambda} =$
18 $(\lambda_1, \lambda_2, \dots, \lambda_n)$, $\boldsymbol{\mu} = (\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_n)$ and $\boldsymbol{\sigma} = (\boldsymbol{\sigma}_1, \boldsymbol{\sigma}_2, \dots, \boldsymbol{\sigma}_n)$, where $\boldsymbol{\mu}_i = (\mu_{1i}, \mu_{2i})^T$, $\boldsymbol{\sigma}_i =$
19 $(\sigma_{1i}, \sigma_{2i})^T$; and (4) the probability of sorting the i th variant into the k th bin, denoted by P_{ik} .
20 Based on these definitions, we built a Bayesian network to show the data generative process
21 and dependencies among variables (see **Methods, Fig. 2d**). For parameter learning, instead
22 of only matching the binned distribution via maximum likelihood estimation as in previous
23 methods^{19,21,27}, we constructed a Bayesian neural network to fit both the binned distribution
24 and the overall fluorescence intensity density. Specifically, two objective functions were
25 designed, where the first was defined as the cross-entropy of the observed binned distribution
26 relative to the theoretical binned distribution derived from LGMM (see **Methods, Fig. 2e**),
27 by minimizing which the parameters of each LGMM can be optimized for approximation to
28 the observed sorting data. The second objective was aimed at matching the overall
29 fluorescence intensity distribution of the whole library. For this purpose, we applied a
30 generative adversarial network³⁶ (Fig. 2f). To elaborate, a generator was designed based on
31 the data generative process (see **Methods**). A fully connected neural network was applied as
32 the discriminator. During training, data generated from the generator are sent to the
33 discriminator along with the real fluorescence intensity values, and then the discriminator

1 determines whether each piece of data is real or not. Hence, a two-player game is played
2 between the generator and the discriminator, and the overall fluorescence intensity
3 distribution can be matched by the generator when they are in equilibrium. We included these
4 two objectives in the Bayesian neural network, with which the parameters can be learned
5 through backpropagation (**Supplementary Fig. 1**).

6 Subsequently, we tested the validity of dSort-Seq with data from our previous Sort-Seq
7 profiling of a comprehensive codon-level mutagenesis library of *tnaC*²⁷. This experiment was
8 performed under 3 different ligand concentrations (0, 100 and 500 μ M Ala-Trp), each with
9 two biological replicates (**Supplementary Fig. 2a**). However, by fitting each binned
10 distribution to the nonrobust log-normal density, their results were obviously affected by
11 outliers and could not precisely match the experimental observations (see **Methods**) in terms
12 of both the binned distribution (**Fig. 2g**) and the overall fluorescence intensity distribution
13 (**Fig. 2h**). In addition, the expression characteristics derived from the log-normal distribution
14 were also subject to error (see **Methods, Supplementary Fig. 5**). Therefore, we applied
15 dSort-Seq in this case to derive the expression properties (see **Methods, Supplementary**
16 **Data 1**). As a result, the strong correlations of the mean (**Supplementary Figs. 2b-2d**,
17 Pearson's $r = 0.989, 0.979, \text{ and } 0.978$ for 0, 100, and 500 μ M Ala-Trp, respectively) and
18 standard deviation (SD; **Supplementary Figs. 2e-2g**, Pearson's $r = 0.914, 0.903, \text{ and } 0.891$
19 for 0, 100, and 500 μ M Ala-Trp, respectively) of expression between biological replicates
20 indicated the reliability of dSort-Seq profiling. In addition, the individual validation data
21 (**Supplementary Figs. 3 and 4**), even if measured via another flow cytometer, were highly
22 consistent with the calculation results (**Fig. 2i**, mean for 0 μ M Ala-Trp, Pearson's $r = 0.991$;
23 **Fig. 2j**, SD for 0 μ M Ala-Trp, Pearson's $r = 0.942$; **Fig. 2k**, mean for 100 μ M Ala-Trp,
24 Pearson's $r = 0.994$; **Fig. 2l**, SD for 100 μ M Ala-Trp, Pearson's $r = 0.931$), which proved the
25 model's ability to precisely capture the expression characteristics.

26

27 **DSort-Seq enables the screening of biosensors with desired response features**

28 Given the superior performance of dSort-Seq in characterizing expression properties, we
29 applied it to practical scenarios to highlight its applicability. First, as an example of
30 expression strength mining, we focused on the metabolite biosensor, through which the
31 intracellular concentration could be converted to a change in gene expression. The key

1 performance indicators of a biosensor include sensitivity, specificity, dynamic range and
2 operational range³⁷, most of which can be determined from dose–response relationships.
3 However, to the best of the authors’ knowledge, a method for large-scale profiling of the
4 dose–response curves with high precision is lacking. Therefore, we applied dSort-Seq to
5 address this problem. For instance, we tested it in our previously reported dataset by Zhou et
6 al.³⁸, which contained Sort-Seq results for 5,184 FapR-*fapO*-based biosensors, consisting of
7 combinations of 6 transcription factor dosages (*pGPD*, *pENO2*, *pHSP12*, *pEXG1*, *pCYC1*,
8 *pULI1*), 4 operator insertion schemes (TATA_OP, OP_TATA, OP_TATA_OP, N30_OP),
9 and 216 arrangements of upstream enhancer sequences (UASs; 3 tandem UASs selected from
10 UAS_A, UAS_B, UAS_C, UAS_D, UAS_E and UAS_F). Each combination was encoded by a specific
11 DNA barcode to ensure its identification via NGS. The library was transformed into
12 *Saccharomyces cerevisiae* BY4700 and assayed through Sort-Seq under 6 different cerulenin
13 concentrations (0, 1, 2, 3, 5, 8 mg/L), each with two biological replicates (**Fig. 3a**). However,
14 owing to the limited precision and robustness of log-normal-based analysis, the dose–
15 response curves derived from Sort-Seq were imprecise (**Supplementary Fig. 6**) and were
16 inconsistent with the individual characterization data³⁸. Therefore, we applied dSort-Seq to
17 this case to determine whether it could accurately evaluate the response performance. The
18 resulting responses showed strong correlations between biological replicates at different
19 concentrations (**Supplementary Fig. 7**, Pearson’s $r > 0.950$ for all experimental conditions),
20 indicating the reliability of the calculation.

21 As the library was nonuniform, we could obtain only 12,779 expression strengths of 2,616
22 combinations through dSort-Seq. To obtain the rest of the data, we applied a machine
23 learning approach to generate predictions. Specifically, each combination was encoded as a
24 27-dimensional vector (see **Methods**), along with the cerulenin concentration as the input
25 feature. Gradient boosting regression was applied to fit the log-scaled expression strength.
26 Note that as combinations containing the promoter *pGPD* suffered from a heavy metabolic
27 burden, leading to them being underrepresented in the library, we excluded them from the
28 machine learning analysis. Therefore, the data used to train the model contained 11,375
29 responses, covering 43.9% of the whole combinatorial space (**Supplementary Data 2**). To
30 avoid overfitting, we randomly split the dataset into two subgroups, with 80% of the data
31 used as the training dataset to optimize the hyperparameters through 5-fold validation as well
32 as train the model parameters. The remaining 20% were used as the test dataset to check the
33 generalization capacity of the model (**Fig. 3b**). The performances in the test dataset ($r^2 =$

1 0.989, **Fig. 3c**) indicated that the model had reasonable generalization capacity and captured
2 biological signals. Subsequently, we trained the model on the whole dataset and predicted the
3 uncharacterized responses. The dSort-Seq data accompanied by the predicted results were
4 then applied to generate the dose–response curves for all combinations (**Supplementary**
5 **Data 2**). We validated these dose–response relationships using 92 individual characterization
6 results³⁸, and linear regression was applied to fit the data values within the same scale. The
7 resulting high consistencies (Pearson’s $r > 0.970$ for all cases, **Supplementary Figs. 8 and 9**)
8 demonstrated that with dSort-Seq and machine learning, the expression properties of the
9 enormous combinatorial space could be effectively explored.

10 We then analyzed the features that contributed most to model predictions via Gini importance
11 (**Fig. 3d**). Overall, consistent with previous discoveries³⁸, the responses of the biosensor were
12 mostly affected by the operator insertion schemes. In addition, a strong determinant of the
13 result was observed if the third UAS was UAS_C. Next, we focused on the dynamic range,
14 which measures the signal-to-noise ratio of a biosensor, by increasing which the true signal is
15 more likely to be discerned from noise. Hence, we fitted each dose–response relationship to
16 the Hill equation to derive the corresponding dynamic range (see **Methods**). The top 10
17 combinations with the highest dynamic ranges were individually constructed and assayed
18 through FCM (**Fig. 3e**). Their response performances were consistent with dSort-Seq and
19 machine learning calculations, of which pHSP12-TATA_OP-UAS_FAC achieved the highest
20 dynamic range of 3.5. Notably, the third UAS of most of the 10 combinations was UAS_C,
21 indicating that UAS_C is important for the interactions of yeast synthetic promoters with FapR
22 when it is located at the third position. In addition to dynamic range, other indicators,
23 including operational range and sensitivity, can also be evaluated and optimized in a similar
24 manner. Therefore, with dSort-Seq, the optimal design with desired response features can be
25 effectively identified.

26

27 **DSort-Seq profiling of the mean noise landscape of *E. coli* endogenous promoters**

28 In addition to the expression strength, expression noise is also an important factor affecting
29 gene expression that leads to phenotypic diversity among genetically identical individuals.
30 Previous association studies have found that expression noise is a heritable trait³⁹ and is
31 determined by expression modules^{14–16,30}. Hence, for a given organism, how different
32 expression modules shape the patterns of noise is a fundamental question. On the other hand,

1 in terms of noise production mechanisms, the commonly accepted translational bursting
2 model suggests that the protein within a cell is produced in bursts, where the burst size (noise
3 strength, $\eta = SD^2 / \text{Mean}$) is related to only translation and is independent of transcription^{3,4}
4 (**Fig. 4b**). However, relevant experimental evidence is still scarce. The problem with the
5 translational bursting mechanism is not only that the gamma distribution cannot accurately
6 represent the gene expression distribution⁶ but also that it cannot interpret the dependence
7 between transcription and noise strength¹⁵. However, limited by the low throughput of classic
8 quantitative methods, research on transcriptional and translational contributions to expression
9 noise is always based on the analysis of a small amount of data^{15,40}, which is susceptible to
10 experimental error as well as outlier samples. To address these issues, our approach may
11 serve as a promising method due to its ability to produce high-quality data in a massively
12 parallel manner. Therefore, as a proof of concept, we performed systematic profiling of
13 transcriptional effects on expression noise in *E. coli* based on dSort-Seq.

14 For library construction, 3,804 endogenous promoters of *E. coli* K12 MG1655 were collected
15 from the EcoCyc database⁴¹ (<https://www.ecocyc.org/>), for which the 60-nt sequence
16 upstream of each transcription start site was regarded as the promoter region. The
17 oligonucleotide library composed of collected promoters was high-throughput synthesized
18 and assembled into a low-copy-number, dual-fluorescence plasmid,
19 pMPTPV_dual_fluorescence, in which a superfolder green fluorescent protein (sfGFP) was
20 used as the response reporter that was under the control of a particular promoter with a fixed
21 ribosome-binding site (RBS; BBa_J61106). In addition, a constitutively expressed reporter,
22 mCherry, served as an internal reference to eliminate cell-to-cell variations such as cell
23 volume and plasmid copy number. After electroporation into MG1655 cells, a cell library
24 with broad levels of sfGFP expression was obtained. The cultivated cell library was
25 characterized by FCM with three biological replicates and then sorted into 12 bins based on
26 the fluorescence intensity of sfGFP relative to mCherry, followed by NGS to quantify the
27 proportion of each variant in each bin. The acquired datasets were then processed and
28 analyzed by our method (**Fig. 4a**). The results showed that 2,920 (76.8% of the total library)
29 promoters were highly consistent among all three replicates (**Supplementary Fig. 15**,
30 **Supplementary Data 3**). Validation of this result was carried out through individual
31 cytometry assays of 60 randomly picked single colonies (**Supplementary Fig. 16**). The
32 strong consistency with the dSort-Seq results (**Supplementary Fig. 17**, Pearson's $r = 0.981$
33 and 0.921 for the mean and SD, respectively) indicated the reliability of the profiling.

1 Autofluorescence was quantified by assaying the pMPTPV strain with only mCherry
2 expression and no sfGFP expression, and the result showed that autofluorescence could be
3 neglected relative to the fluorescence intensity of each candidate of the library
4 (**Supplementary Fig. 12**).

5 The bulk data generated a comprehensive landscape of promoter strength and expression
6 noise along the *E. coli* genome (**Supplementary Data 3**, we also visualized it through
7 D3GB⁴² at http://www.thu-big.net/Escherichia_coli_K12_MG1655_promoters), which was
8 beneficial for understanding the transcriptional strategies for different genes. For instance, we
9 investigated whether essential and nonessential genes of *E. coli* possess different expression
10 patterns (see **Methods**). As a result, the essential genes showed greater transcriptional
11 intensities than the nonessential genes (**Supplementary Fig. 18**, $P = 4.22e-16$, one-tailed t
12 test). Given that the high transcriptional strength is usually related to low expression noise¹⁵,
13 these functionally important genes are more likely to confer lower levels of noise, which is
14 consistent with the results of a previous genome-wide association study¹⁶. Subsequently, we
15 investigated the relationship between noise strength and mean expression level. The results
16 showed that the noise strength was linearly correlated with the expression strength when the
17 transcription module varied (**Fig. 4c**, Pearson's $r = 0.745$); hence, transcription contributed to
18 expression noise. This discovery, however, is inconsistent with the inference of the
19 translational bursting model (**Fig. 4b**), suggesting the limitation of the model in interpreting
20 noise production mechanisms in *E. coli*.

21

22 **Transcription and translation make comparable contributions to noise production**

23 Although the translational bursting mechanism is unable to account for the contribution of
24 transcription to noise production, the hierarchical Bayesian model, developed by introducing
25 transcriptional and translational fluctuations into the translational bursting model, can
26 successfully explain this phenomenon³⁰ (**Fig. 4c**). However, the model showed that different
27 translation modules would lead to varying intercepts^{15,30} in the relationship between noise
28 strength and the mean expression level (**Figure 4c**, $\eta = C_1 \cdot \text{Mean} + C_2 \cdot b$, where b
29 represents the translation strength). Hence, at low levels of transcription, the expression noise
30 is still dominated by translation. This conclusion, although confirmed by a fluorescence
31 microscopy experiment that analyzed 40 *B. subtilis* strains expressing GFPmut3 with
32 different combinations of transcription and translation modules¹⁵, still needs to be verified, as

1 it was based on regression analysis of a small sample. To date, the contribution of
2 transcription and translation to noise production has not been comprehensively and directly
3 observed. As our method has greatly expanded the test throughput of expression noise, we
4 applied dSort-Seq here to examine these features.

5 Therefore, we designed a combination library comprising different combinations of 300
6 promoters and 13 RBSs. The promoters were randomly selected from the EcoCyc database,
7 whereas the RBSs, which were chosen for their varying translational strengths, were from the
8 apFAB# series⁴³. The combination library was prepared in the same manner as the promoter
9 library in the pMPTPV_dual_fluorescence plasmid. After electroporation, we performed a
10 dSort-Seq assay of the cell library, with three independent biological replicates to ensure the
11 reliability of the results. After data processing, 2,733 combinations (70.1% of the whole
12 library) were highly consistent among the replicates and were retained for subsequent
13 analysis (**Supplementary Fig. 20, Supplementary Data 4**). Subsequently, 60 single colonies
14 with different genotypes were randomly picked and assayed individually with FCM
15 (**Supplementary Fig. 21**). Their means and SDs of expression were strongly correlated with
16 the dSort-seq results (**Supplementary Fig. 22**, Pearson's $r = 0.976$ and 0.937 for the mean
17 and SD, respectively), proving the validity of the profiling.

18 We then performed regression analysis between noise strength and the expression mean for
19 different translational modules (**Fig. 4e, Supplementary Fig. 24**) to test the hierarchical
20 Bayesian model. However, their correlation barely changed when the translation module
21 varied in terms of both slope and intercept (**Fig. 4f, Supplementary Table 6**). Instead, our
22 results showed that the noise strength was highly coupled with the mean expression level,
23 indicating the difficulty of adjusting expression noise independently of the mean protein
24 abundance by tuning the strength of the transcription and translation modules. Furthermore,
25 to determine whether translation bursting dominates noise production at low transcription
26 levels, we constructed several weak expression combinations and performed cytometry
27 assays. As a result, the combinations with comparable mean expression levels showed similar
28 fluorescence intensity distributions (**Supplementary Fig. 25**), suggesting that the
29 contributions of transcription and translation to noise are comparable, even in the case of
30 weak expression strength.

31

32 **Overlapping RpoD-binding sites can lead to high expression noise**

1 We then analyzed the relationship between expression noise ($CV^2 = SD^2/Mean^2$) and
2 strength (**Fig. 5a and 5b**). The expression noise exhibited a strong negative correlation with
3 the mean protein abundance at low levels of expression and then reached a plateau after a
4 critical point, which is consistent with previous observations^{15,30}. In addition to the general
5 correlation, some unique expression features also piqued our interest, especially for those
6 sequences exhibiting high expression noise at their corresponding mean expression levels. To
7 ensure that these outliers were not the results of experimental error, we reconstructed and
8 assayed 20 high-noise candidates from the promoter library and 25 from the combination
9 library and then performed an individual cytometry assay (see **Methods**), proving the
10 credibility of the discovery. Moreover, the expression noise of these variants was apparently
11 higher than that of randomly selected colonies (**Fig. 5c and 5d**). Notably, among the 25 high-
12 noise combinations, several promoters appeared frequently (e.g., *fliEp1*, *ileSp2*, *yihVp3*,
13 *folEp*, **Fig. 5d**), suggesting that the extra noise may be derived from transcription rather than
14 translation.

15 Next, to identify the factor that contributed to the additional expression noise, we divided the
16 *E. coli* promoters into 3 groups based on the noise values (high-noise, medium-noise and
17 low-noise groups; see **Methods**). Subsequently, we analyzed the sequence features of the
18 three groups and found that the thymidine proportion was significantly higher in the high-
19 noise group (average 37.5%, $n = 138$) than in the low-noise group (**Fig. 5e**, average 28.0%, n
20 $= 138$; $P = 1.97e-21$, one-tailed t test) as well as in the medium-noise group (average 29.9%,
21 $n = 2,481$; $P = 2.24e-29$). Furthermore, there was no position preference for this phenomenon
22 (**Supplementary Fig. 30**), and no common regulator was found to be associated with the
23 sequences within the same group (**Supplementary Data 5**). Therefore, we focused on
24 transcription initiation factors, especially RpoD (σ^{70} factor), which transcribes most genes in
25 *E. coli*. Promoters recognized by RpoD generally contain two consensus hexamers centered
26 at 10 and 35 nucleotides upstream of the transcription start site. These two regions are rich in
27 adenosine and thymine, especially thymine (average of 4.73 per promoter compared to 3.75
28 for adenosine, 1.80 for cytosine and 1.70 for guanine⁴⁴). Based on this, we hypothesized that
29 the high-noise group contains more RpoD-binding sites than other groups. To test this
30 hypothesis, we searched the DPinteract database⁴⁵ for potential RpoD-binding sites in three
31 groups (see **Methods**). The sequences with no RpoD-binding hit were excluded from
32 subsequent analysis. As a result, the high-noise group (average 2.10, $n = 127$) showed more

1 potential RpoD-binding sites than the low-noise group (**Fig. 5f**, average 1.65, $n = 106$; $P =$
2 $1.86e-5$, one-tailed t test) and medium-noise group (average 1.82, $n = 2,091$; $P = 5.34e-5$).

3 Subsequently, to determine whether overlapping RpoD-binding sites would result in high
4 expression noise, we constructed 25 tandem promoters based on combinations of 5 Anderson
5 promoters (J23101/103/107/109/115; **Fig. 5g**) to drive the expression of *sfgfp*. In addition, the
6 5 constitutive promoters were also individually constructed as controls. To exclude the effect
7 of promoter length and the -35 region-proximal sequence on the results, we also constructed
8 5 promoters of the same length as the tandem promoter, while preserving only the
9 downstream RpoD-binding site (**Fig. 5h**). Subsequently, we performed an individual
10 cytometry assay of the 35 promoters (**Supplementary Fig. 32**). As a result, the gene
11 expression driven by tandem promoters showed higher noise compared to a single promoter
12 (**Fig. 5i**, $P = 2.40e-4$, one-tailed t test), especially when the stronger promoter was located
13 upstream of the weaker promoter (**Fig. 5i**), suggesting that the transcriptional interference
14 caused by the occlusion of promoters contributed to noise production. Hence, our results
15 uncovered a feasible and simple noise modulation strategy in *E. coli* by tuning the number
16 and relative positions of sigma factors upstream of the transcription start site.

17

18 **DISCUSSION**

19 Gene expression dosage is directly associated with a variety of phenotypes of a population⁴⁶;
20 hence, there is no doubt that high-throughput profiling will deepen our understanding of
21 cellular behavior. In this paper, we focused on biosensors that can sense metabolite
22 concentrations and regulate gene expression. According to the different output signals,
23 biosensors have various applications, including in high-throughput screening⁴⁷, medical
24 diagnosis⁴⁸, and cell imaging⁴⁹. For each application, the dose–response relationship is a key
25 indicator that needs to be tuned to meet practical needs. Fortunately, dSort-Seq was shown to
26 be a powerful tool for characterizing and optimizing the biosensor response performance in a
27 high-throughput and high-precision manner, enabling the engineering of biosensors with
28 desired properties. Compared to positive and negative screening^{50,51}, by which only the
29 dynamic range could be optimized, dSort-Seq can yield more comprehensive information to
30 meet the needs of various situations. In addition, it is much more efficient than traditional
31 trial-and-error approaches. Therefore, dSort-Seq provides a solution for profiling the
32 expression landscape of the combinatorial sequence space. On the other hand, the high-

1 quality dSort-Seq dataset also has the potential to serve as a basis for deciphering
2 physiological mechanisms^{25,27}.

3 Noise in biological systems has been widely demonstrated to influence various intracellular
4 processes⁵² and the physiological properties of a population^{13,53}, although the effects of noise
5 strength vary across different situations. For instance, low noise can ensure stable
6 biosynthesis pathways and robust synthetic gene circuits⁵⁴. In contrast, high phenotypic
7 variability promotes evolvability⁵⁵⁻⁵⁷. Therefore, it is necessary to understand the origin of
8 the noise, as well as control noise rationally for various applications. Regarding noise
9 regulation, various strategies have been proposed to control gene expression noise
10 independently, including engineering transcription and translation in synthetic gene
11 circuits^{58,59}, introducing pulsatile input to control the promoter activation frequency and
12 transcription rate independently⁶⁰ and expressing two copies of the target gene from separate
13 circuits with different characteristics⁶¹, among others. Through dSort-Seq profiling of
14 different combinations of promoters and RBSs in *E. coli*, the transcriptional effect on gene
15 expression noise was revealed. Specifically, a higher thymidine proportion without position
16 preference in the promoter sequence would lead to a higher level of noise. One hypothesis to
17 explain this phenomenon is that RpoD-binding sites, which are rich in thymine, influence
18 noise production. We have proven that promoters with overlapping RpoD-binding sites
19 contribute to noise production due to occlusion of promoters. Moreover, in *E. coli*, 831 genes
20 have been found to be under the control of tandem promoters⁶², suggesting the broadness of
21 the regulatory scheme. Hence, in-depth research on modeling molecular events in the
22 transcription process is needed to elucidate the effect of promoter architecture on expression
23 noise.

24 From the methodology point of view, the design-build-test-learn (DBTL) cycle is emerging
25 as a key workflow in synthetic biology, where the test is the rate-limiting step due to its low
26 throughput⁶³. The development of Sort-Seq has undoubtedly greatly extended the test
27 throughput, enabling more efficient characterization and optimization of biological parts. The
28 dSort-Seq approach shows that the learn can be encapsulated into the test to improve its
29 capability by modeling the data generative process for the high-throughput experiment.
30 Moreover, it is worth mentioning that as the ability of the Gaussian mixture models in
31 distribution matching can be improved by increasing the number of mixture components,
32 dSort-Seq can be easily transferred to more complex situations in which multiple feedback
33 circuits are involved⁶⁴. Thus, this pipeline has great potential to determine the mean-noise

1 space for various gene expression modules, providing diverse synthetic parts that can be
2 applied to different fields, such as biosynthesis⁶⁵, laboratory-based adaptive evolution⁵⁶,
3 transcriptional regulation⁶⁶, and protein–protein interactions^{67,68}. Overall, owing to the
4 flexibility, high precision and high throughput of this method, we believe dSort-Seq can serve
5 as a powerful tool that provides a wide range of novel research opportunities.

6

7 METHODS

8

9 Parameter-learning algorithm of dSort-Seq

10 We represent the log-scaled expression density of each variant by a two-component Gaussian
11 mixture model (where x_i denotes the log-scaled intensity value of the i th variant):

$$12 \quad x_i \sim f(x_i|\lambda_i, \boldsymbol{\mu}_i, \boldsymbol{\sigma}_i) = \lambda_i N(\mu_{1i}, \sigma_{1i}^2) + (1 - \lambda_i) N(\mu_{2i}, \sigma_{2i}^2) \quad (i = 1, 2, \dots, n). \quad (1)$$

13 Hence, the overall logarithmic fluorescence intensity distribution can be modeled as a
14 mixture of Gaussian mixture models:

$$15 \quad x \sim f(x) = \sum_{i=1}^n \pi_i f(x_i). \quad (2)$$

16 The generative process for each fluorescence intensity value should be as follows: (1) choose
17 a variant $z_1 \sim \text{Categorical}(\boldsymbol{\pi})$, (2) choose a Gaussian component z_2 from $\text{Bernoulli}(z_2|z_1)$
18 and (3) choose a log-scaled intensity value x from $N(x|z_1, z_2)$ (**Fig. 3d**). The distributions of
19 variables involved in the model are listed as follows:

$$20 \quad P(z_1 = i) = \pi_i \quad (i = 1, 2, \dots, n); \quad (3)$$

$$21 \quad P(z_2 = 1|z_1 = i) = \lambda_i; \quad (4)$$

$$22 \quad P(z_2 = 0|z_1 = i) = 1 - \lambda_i; \quad (5)$$

$$23 \quad f(x_i|z_1 = i, z_2 = 1) = N(\mu_{1i}, \sigma_{1i}^2); \quad (6)$$

$$24 \quad f(x_i|z_1 = i, z_2 = 0) = N(\mu_{2i}, \sigma_{2i}^2); \quad (7)$$

$$1 \quad P_{ik} = \int_{b_{k-1}}^{b_k} \lambda_i N(x_i | \mu_{1i}, \sigma_{1i}^2) + (1 - \lambda_i) N(x_i | \mu_{2i}, \sigma_{2i}^2) dx_i \quad (i = 1, 2, \dots, n; k = 1, 2, \dots, K). \quad (8)$$

2 Among the parameters involved in the model, the sets λ , μ and σ cannot be identified
 3 experimentally. To estimate them, we designed a probabilistic artificial neural network in
 4 which a double-objective optimization is performed. The first objective function is defined as
 5 the cross-entropy (H) of the observed binned distribution relative to the integral of the
 6 probability density over adjacent boundaries (**Fig. 3e**), which is shown as follows:

$$7 \quad \min H = - \sum_{i=1}^n \sum_{k=1}^K \left(-P_{ik} \log \left(\int_{b_{k-1}}^{b_k} f(x_i | \lambda_i, \mu_i, \sigma_i) dx_i \right) \right). \quad (9)$$

8 By minimizing the above loss function, the binned distribution can be fitted. The other
 9 objective is to match the overall fluorescence intensity density. To this end, a generative
 10 adversarial network³⁶ is applied. Specifically, a generator is constructed based on the
 11 abovementioned generative process. For the discriminator, a fully connected neural network
 12 is used to determine whether the data are real or fake (**Fig. 3f**). During training, a two-player
 13 game is played between the generator G and the discriminator D with value function $V(G, D)$:

$$14 \quad \min_G \max_D V(G, D) = E_{x_{true}} (\log(D(x_{true}))) + E_{x_{fake} \sim f(x|\pi, \lambda, \mu, \sigma)} (\log(1 - D(x_{fake}))). \quad (10)$$

15 Combining the above two parts, we can obtain the whole algorithm for parameter learning, as
 16 shown in **Supplementary Fig. 1**.

17

18 **Obtaining expression characteristics from cytometry data**

19 For each individual cytometry assay, the \log_{10} -transformed fluorescence intensity distribution
 20 was fitted by a two-component Gaussian mixture model (**Supplementary Figs. 3, 4, 12, 16,**
 21 **21, 25, 26, 28 and 32**) via the expectation–maximization (EM) algorithm, which resulted in a
 22 representation of $f(x_i) = \lambda_i N(\mu_{1i}, \sigma_{1i}^2) + (1 - \lambda_i) N(\mu_{2i}, \sigma_{2i}^2)$. The mean expression strength
 23 was calculated with Eq. 11.

$$24 \quad Mean = Mean_1 \times Mean_2 = \exp(m_1 + V_1/2) \times \exp(m_2 + V_2/2) \quad (11)$$

25 where $m_1 = \lambda_i \mu_{1i} \log(10)$, $V_1 = (\lambda_i \sigma_{1i} \log(10))^2$, $m_2 = (1 - \lambda_i) \mu_{2i} \log(10)$ and $V_2 =$
 26 $((1 - \lambda_i) \sigma_{2i} \log(10))^2$. The SD of each expression density was calculated with Eq. 12.

1
$$SD = \sqrt{Var_1 \times Var_2 + Var_1 \times Mean_2^2 + Var_2 \times Mean_1^2} \quad (12)$$

2 where $Var_1 = \exp(2m_1 + V_1) (\exp(V_1 - 1))$, $Var_2 = \exp(2m_2 + V_2) (\exp(V_2 - 1))$.

3

4 **Comparison of dSort-Seq and the log-normal-based method**

5 The dSort-Seq results were calculated as mentioned above, and the log-normal results were
6 obtained from previously reported data²⁷ (note that since the actual slope of the sorting
7 boundary lines on the log-log plot of eGFP-mCherry in the experiment was 0.8810 instead of
8 1, each boundary value was shifted to the right by 0.3801 compared to the previous analysis,
9 **Supplementary Table 4**). Next, as an example, we compared the performances of the two
10 methods in matching the binned distribution of variant V8A_GCC under 0 μ M Ala-Trp and
11 calculated the Kullback–Leibler divergences (Eq. 13) of the observation from the theoretical
12 binned distributions derived from the log-normal-based method and dSort-Seq. The results
13 showed that dSort-Seq is more precise and robust than log-normal (**Fig. 2g**).

14
$$D_{KL}(P||Q) = - \sum_x P(x) \log \left(\frac{Q(x)}{P(x)} \right) \quad (13)$$

15 Subsequently, we compared these two methods in fitting the overall fluorescence intensity
16 distribution. For instance, we calculated the theoretical log-scaled overall distribution (100
17 μ M Ala-Trp, replicate 1) derived from the log-normal-based method (Eq. 14) and dSort-Seq
18 (Eq. 2). As a result, dSort-Seq also showed better performance (**Fig. 2h**).

19
$$f_{log-normal}(x) = \sum_{i=1}^n \pi_i N(x_i | \mu_i, \sigma_i^2) \quad (14)$$

20 Moreover, as mentioned above, the log-normal-based method is unable to fit the individual
21 cytometry data, which usually serve as criteria for validating the Sort-Seq results
22 (**Supplementary Figs. 3 and 4**). To measure the error, we calculated the expression strength
23 and SD of individual validation data with both log-normal (Eq. 15 and Eq. 16) and the
24 LGMM (Eq. 11 and Eq. 12), where the LGMM results were applied as ground truth to
25 evaluate the precision of log-normal results. As a result, the response and SD inferred from
26 log-normal showed significant deviations (**Supplementary Fig. 5**). Therefore, with log-
27 normal, it is difficult to infer accurate expression properties from Sort-Seq experiments.

28
$$Mean_{log-normal(\mu_i, \sigma_i)} = \exp(\log(10) \cdot \mu_i + (\log(10) \cdot \sigma_i)^2 / 2) \quad (15)$$

$$SD_{log-normal}(\mu_i, \sigma_i) = \sqrt{[\exp(\sigma_i \log(10))^2 - 1] \exp(2\mu_i \log(10) + (\sigma_i \log(10))^2)} \quad (16)$$

2

3 **DNA manipulation and reagents**

4 Plasmid extraction and DNA fragment purification were performed using kits from Omega
5 Bio-Tek. PCRs were carried out using a KAPA HiFi PCR Kit from KAPA Biosystems. The
6 restriction enzyme FastDigest *Esp3I* (namely, *BsmBI*) and T4 DNA ligase were purchased
7 from Thermo Scientific. Cerulenin was ordered from Yuanye Bio-Technology. All strains
8 and plasmids used in this work are summarized in **Supplementary Table 1**. All
9 oligonucleotides (**Supplementary Table 2**) were ordered from Azenta. Molecular cloning
10 was performed with *E. coli* DH5 α (BioMed) as the host. The concentrations of the antibiotics
11 kanamycin and ampicillin were 50 mg/L and 100 mg/L, respectively. In all experiments,
12 bacteria and yeast were grown at 37 and 30°C, respectively.

13

14 **Featurization and gradient boosting regression**

15 We applied one-hot encoding to transform each biosensor combination into a 27-dimensional
16 vector. Among these dimensions, 5 of them represent the promoters of the transcription factor
17 (*pULII*, *pHSP12*, *pEXG1*, *pENO2*, *pCYCI*), 4 of them represent the operator insertion
18 schemes (OP_TATA_OP, TATA_OP, N30_OP, OP_TATA), and 18 of them represent the
19 sequences of the tandem UAS (UAS_1A/B/C/D/E/F, UAS_2A/B/C/D/E/F,
20 UAS_3A/B/C/D/E/F). These vectors then served as input features along with the cerulenin
21 concentration (0/1/2/3/5/8). Gradient boosting regression was applied to predict the log-
22 scaled expression strength. During training, the hyperparameters were optimized following
23 the given order (min_samples_split, max_depth, min_samples_leaf, max_features,
24 subsample, learning_rate and n_estimators) through the grid search method.

25

26 **Fitting the dose–response relationship to the Hill equation**

27 Each dose–response relationship was fitted by Eq. 17 via nonlinear least squares.

1
$$S = S_0 + \frac{S_m - S_0}{1 + (C/C_{1/2})^h} \quad (17)$$

2 where S_0 and S_m are the values of the sensor response at zero and saturating ligand
3 concentrations, $C_{1/2}$ is the concentration at half saturation, and h is the Hill coefficient. The
4 lower bounds and upper bounds of $(S_0, S_m, C_{1/2}, h)$ were set to $(0, 0, 0, 1)$ and $(1, 2, 8, 3)$,
5 respectively. The dynamic range was calculated with Eq. 18.

6
$$d = \frac{S_m - S_0}{S_0} \quad (18)$$

7

8 **Individual characterization of the dose–response relationships for malonyl-CoA** 9 **biosensors**

10 One variant (pHSP12-TATA_OP-UAS_FAC) was obtained from library stock, and the other
11 9 biosensor variants (pCYC1-OP_TATA-UAS_FAC, pCYC1-OP_TATA-UAS_DDC,
12 pCYC1-OP_TATA-UAS_EBC, pHSP12-TATA_OP-UAS_BDC, pCYC1-OP_TATA-
13 UAS_BEC, pEXG1-N30_OP-UAS_FDA, pCYC1-TATA_OP-UAS_EAC, pCYC1-
14 OP_TATA-UAS_FDC, pCYC1-OP_TATA-UAS_BDC) were constructed via Golden Gate
15 Assembly (**Supplementary Table 3**). After transformation of these plasmids into BY4700,
16 the strains were inoculated into 48-well deep-well plates with 1 mL of SC-Ura medium
17 (synthetic complete medium lacking uracil) in each well. After culturing for 12 h at 30°C and
18 250 rpm, 2 µL of cerulenin solutions of six distinct concentrations (0.5, 1.0, 1.5, 2.5, 4
19 mg/mL) was added to the corresponding well. The strains were then cultured for another 12
20 h. For sample preparation, cells were collected by centrifugation (4°C; 8,000 × g for 10 min)
21 and resuspended in prechilled phosphate-buffered saline (PBS) to an OD₆₀₀ of 2. BY4700
22 was used as a negative control. BY4700/POT1-pTEF2-mCherry-tADH1 and BY4700/POT1-
23 pCYC1-YPet-tPGK1 were used as positive controls for mCherry and YPet, respectively.
24 These control samples were prepared the same way as above. The fluorescence intensities of
25 the cells were characterized on an LSRFortessa (BD Biosciences). The double-positive area,
26 named Q2, was determined by the control samples, as described in a previous work³⁸. For
27 each sample, 100,000 events in the Q2 area were analyzed.

28

1 Construction of the two-reporter plasmid

2 The two-reporter plasmid pMPTPV_dual_fluorescence was derived from the common vector
3 pACYCDuet-1 by replacing the chloramphenicol resistance gene with the kanamycin
4 resistance gene, replacing the *lacI* expression cassette with *mcherry*, and inserting the *sfgfp*
5 cassette into the opposite strand of *mcherry*. The *mcherry* gene is controlled by a constitutive
6 promoter, pL_M1-37⁶⁹. To facilitate library construction, *sfgfp* is controlled by a variable
7 region containing two *BsmBI* restriction sites.

8

9 Feasibility verification of the plasmid

10 Ten promoters were randomly selected from the EcoCyc database. The sequence of each
11 promoter was defined as the 60 nt preceding the transcriptional start site. In addition, a
12 medium-strength RBS, BBa_J61106 (TCTAGAGAAAGATAGGAGACACTAGT), was
13 chosen for all strains to ensure the survival of strains (note that the combination of a strong
14 promoter and a strong RBS is lethal to *E. coli*¹⁷). After transformation of the plasmids
15 containing different promoters into *E. coli* K12 MG1655, the resulting strains were
16 individually cultured in LB medium containing kanamycin (initial OD₆₀₀ = 0.02), with three
17 biological replicates for each promoter. During cultivation, the growth rate and the
18 expression of fluorescent protein were monitored by sampling and testing every hour. The
19 OD₆₀₀ was measured by a microplate reader (Tecan Infinite 200Pro), and the fluorescence
20 intensity was assayed via flow cytometry (BD LSRFortessa). The 10 strains showed no
21 apparent differences in growth (**Supplementary Fig. 10a**), and the median value of
22 sfGFP/mCherry remained stable after culturing for 16 h (**Supplementary Fig. 10b**). Hence,
23 we chose 16 h for cultivation in subsequent experiments.

24

25 Preparation of library cells

26 The two plasmid libraries (the promoter library and combination library) were both ordered
27 from Genewiz. We transformed each library into *E. coli* K12 MG1655 via a BTX Harvard
28 ECM 630 High Throughput Electroporation System using optimized parameter settings (2.1
29 kV, 1 kΩ, 25 μF, 100 ng plasmids/100 μL competent cells). The transformed cells were
30 incubated in LB medium (four times the volume of the competent cells) for 1 h at 37°C for

1 recovery and then plated onto 37 Φ 150 LB agar plates containing kanamycin with an
2 EasySpiral Pro (Interscience). Generally, $\sim 10^4$ single colonies per plate can be harvested with
3 this protocol (data not shown), enabling ~ 100 times coverage of the designated library. All
4 colonies on the plates were rinsed off using sterile LB medium supplemented with
5 kanamycin, collected by centrifugation (4°C; 8,000 \times g for 10 min) and then resuspended and
6 thoroughly mixed to an OD₆₀₀ of 10 using fresh sterile LB medium containing kanamycin.
7 The cell suspension was stored at -80°C in glycerol (final OD₆₀₀ = 5).

8

9 **Characterization of transcriptional impact on growth**

10 We further tested whether different promoters have varying influences on growth. To this
11 end, the stored promoter library cells were cultured in LB medium containing kanamycin
12 (initial OD₆₀₀ = 0.02) at 37°C for 16 h. Cell samples were collected before and after
13 cultivation, followed by plasmid extraction. The promoter regions were amplified by PCR
14 (KAPA HiFi PCR Kit; 95°C for 3 min, 25 cycles [98°C for 20 s, 63°C for 15 s, 72°C for 5 s],
15 72°C for 30 s) using the Lib_F and Lib-R primers (**Supplementary Table 2**).

16 In a 50- μ L reaction, 5 ng of the plasmid library was added as the PCR template. The
17 sequencing library was prepared according to the NEBNext Ultra II DNA Library Prep Kit
18 for Illumina (NEB). Specifically, 30 ng of each purified PCR product was used to prepare the
19 sequencing library. The DNA fragments were treated with NEBNext End Prep for end repair,
20 5' phosphorylation and dA-tailing. Then, the fragments were ligated to NEBNext Adaptors,
21 followed by USER Enzyme excision. Subsequently, the products were purified using
22 NEBNext Sample Purification Beads and amplified by PCR for six cycles using the P5 and
23 P7 primers. The products were again purified using NEBNext Sample Purification Beads,
24 validated with an Agilent 2100 Bioanalyzer (Agilent Technologies) and quantified with a
25 Qubit 4 Fluorometer (Invitrogen). Subsequently, the libraries were delivered to Novogene for
26 sequencing.

27 Two biological replicates were analyzed in parallel in this experiment, which generated three
28 NGS raw datasets. After the production of clean data by demultiplexing and removing
29 adaptor regions, pairs of paired-end data were merged by FLASH script⁷⁰, and those reads
30 without detected pairs were removed. Python scripts generated in house were then used to
31 search for the 'GGATN86ATGC' 94-mer in the sequencing reads (and the reverse

1 complementary sequence), and those carrying mutations within the upstream (GGAT) or
2 downstream (ATGC) flanking regions (4 nt each) were removed. The read counts were then
3 adjusted using Eq. 19, where n is the number of sequencing libraries, to normalize the
4 different sequencing depths of each library.

$$5 \quad \text{Normalization factor}_i = \frac{\sum_{i=1}^n \text{Read count}_i}{n \times \text{Read count}_i} \quad (19)$$

6 The library showed negligible variation in growth (**Supplementary Fig. 11**), which ensured
7 the feasibility of using the library in subsequent Sort-Seq experiments.

8

9 **Sort-Seq experiments**

10 For both the promoter library and the combination library, a frozen glycerol stock of library
11 cells (*E. coli* MG1655) was inoculated into 100-mL flasks containing 20 mL of LB medium
12 with kanamycin to an initial OD₆₀₀ of 0.02. Library cells were grown for 16 h at 37°C and
13 220 rpm. The grown cells were transferred to fresh LB medium to an initial OD₆₀₀ of 0.02
14 and grown again under the same conditions as above. A third round of dilution and growth
15 was carried out to improve the expression stability of the fluorescent proteins. After growth,
16 500 μL of culture medium was chilled on ice immediately, and the cells were collected by
17 centrifugation (4°C; 8,000 × g for 10 min). The cells were resuspended in 500 μL of
18 prechilled PBS. Each cell suspension was diluted 150-fold in PBS to prepare samples
19 appropriate for sorting. Three biological replicates were prepared for Sort-Seq experiments.

20 Sorting was performed on a FACSAria SORP (BD Biosciences). Gating based on FSC-Area
21 and SSC-Area was carried out to exclude noncell particles. The population in this gated area
22 is referred to as P1. The fluorescence background noise for the two relevant wavelengths was
23 calibrated using the blank untransfected MG1655 strain. Note that the blank strain was
24 completely negative for both sfGFP and mCherry expression. The resulting double-positive
25 area in the region corresponding to the FITC-Area (sfGFP) and the PE-Texas Red-Area
26 (mCherry) is referred to as Q2 after P1. The prepared library cells were analyzed by
27 cytometry to determine the density distribution contour of the fluorescence in Q2.

28 Subsequently, in the histogram of sfGFP/mCherry, 12 bins were set to evenly split the overall
29 distribution of the population in Q2 (**Supplementary Figs. 13 and 19, Supplementary**
30 **Tables 4 and 5**), referred to as P2 to P13 after Q2, to ensure that the number of cells in each

1 bin was equal and improve the sorting efficiency. For calibration, $\sim 2 \times 10^6$ unsorted cells in
2 gate Q2 were first collected for each sample. In the main sorting process, the three replicates
3 were individually sorted into the 12 bins as described above. Each sample was successively
4 sorted three times using four-way sorting. In each of these sorting runs, cells falling in
5 nonadjacent bins were collected to eliminate the conflicting events between them. Thus, P2,
6 P5, P8 and P11 were simultaneously sorted in one run, as were P3, P6, P9 and P12. During
7 sorting, the cell flow rate was kept at ~ 8000 events/s, and $\sim 5 \times 10^5$ cells were collected in
8 each bin.

9 The sorted cells were collected in 36 (3 samples \times 12 bins) 5-mL polystyrene round-bottom
10 centrifuge tubes (BD Falcon), each of which contained 500 μ L of PBS. The entire contents of
11 each tube were then each transferred to 100-mL flasks containing 20 mL of LB medium with
12 kanamycin and cultured at 37°C for 7 h. These cells were then subjected to plasmid library
13 extraction. Together with the cells from gate Q2 of each of the three samples mentioned
14 above, we obtained 39 plasmid libraries in total.

15 The promoter and RBS regions of *sfGFP* in each library were amplified through PCR (KAPA
16 HiFi PCR Kit; 95°C for 3 min, 25 cycles [98°C for 20 s, 63°C for 15 s, 72°C for 5 s], 72°C
17 for 30 s), using 12 8-nt barcoded primers to identify different sorting bins (primers
18 sorting_P2 to _P13, **Supplementary Table 2**). The barcodes were designed according to the
19 following principles. (1) The Levenshtein distance between every two barcodes was ≥ 4 ; (2)
20 the GC content was 20% to 80%; and (3) there were no more than four consecutive identical
21 bases. In a 25- μ L PCR, 5 ng of plasmid library was added as a template. PCR products from
22 the 12 bins for each sample were mixed, thus obtaining three sorted PCR products (from
23 sorted library cells in different sorting bins for three samples) and three unsorted PCR
24 products (from unsorted library cells in the Q2 gate for three samples). The resulting PCR
25 products were analyzed and purified by electrophoresis. The sequencing libraries were
26 prepared as described above and were then delivered to Novogene for sequencing.

27

28 **Sort-Seq data processing**

29 According to NGS data, the read count $R_{i,k}$ for *Variant_i* in *bin_k* can be observed.
30 Additionally, by analyzing the cytometry data, we can obtain the ratio of cells sorted into
31 *bin_k* against all cells, which is denoted by C_k . Hence, assuming an unbiased NGS

1 quantification process, the proportion of $Variant_i$ in bin_k is $Q_{ik} = R_{i,k}/R_k$. Here, R_k is the
2 total read count for the NGS library derived from bin_k . Therefore, the probability of sorting
3 $Variant_i$ into bin_k should be $P_{ik} = Q_{ik}C_k / \sum_{k=1}^K Q_{ik}C_k$.

4 For the library of *tnaC* variants and the malonyl-CoA biosensors, data processing was
5 performed as described above. However, for the promoter library, a sorting error did exist
6 (**Supplementary Figs. 14a, 14b, 14d and 14f**). We ascribed this error to random screening
7 with an error rate, ϵ , and hence, we modified P_{ik} as shown in Eq. 20.

$$8 \quad \mathbf{P}'_i = [P'_{i1}, P'_{i2}, \dots, P'_{iK}]^T = ReLU(\mathbf{E}^{-1}\mathbf{P}_i). \quad (20)$$

9 where

$$10 \quad \mathbf{E} = (1 - K\epsilon)\mathbf{I}_K + \epsilon\mathbf{J}_K. \quad (21)$$

11 Here, \mathbf{I}_K is an identity matrix of size K , and \mathbf{J}_K is a matrix of ones of size K . Finally, the
12 binned distribution was obtained by $P''_{ik} = P'_{ik} / \sum_{k=1}^K P'_{ik}$. For both the promoter and
13 combination libraries, ϵ was set to 0.05, which made the binned distribution more precise
14 (**Supplementary Figs. 14c, 14e and 14g**). The strains with $P''_{i1} > 0.5$ or $P''_{iK} > 0.5$ were ruled
15 out, as they were not effectively sorted. Moreover, to ensure the quality of the results, we
16 eliminated the data with low consistency among replicates. Specifically, if the CV of the
17 calculated mean or SD among biological replicates was greater than 0.5, the related data were
18 removed from the dataset.

19

20 **Evaluation of the expression patterns of essential and nonessential genes**

21 The essential genes were identified based on a comprehensive pooled CRISPRi screening
22 dataset⁷¹ (threshold: fitness ≤ -6), whereas other genes were regarded as nonessential. We
23 calculated the transcriptional strength of each gene. Specifically, the genes belonging to the
24 operons closest to the downstream side of a promoter were considered to be driven by this
25 promoter, and the transcription strength of a gene was calculated as the summation of its
26 promoter strengths.

27

28 **Identification of RBS strengths**

1 To identify the translational strengths of the 13 RBSs, we defined each promoter strength as
2 the expression strength in the promoter library, and then we divided each expression strength
3 from the combination library by the corresponding promoter strength. The median of the
4 calculated results for combinations with the same RBS is defined as the translational strength
5 of that RBS. The resulting order of the RBS strengths was as follows: apFAB872(0.15) <
6 apFAB914(0.23) < apFAB864(0.36) < apFAB865(0.67) < apFAB927(0.75) <
7 apFAB827(0.95) < apFAB894(1.02) < apFAB909(1.04) < apFAB839(1.36) <
8 apFAB833(1.38) < apFAB834(1.41) < apFAB820(1.57) < apFAB916(2.13) (**Supplementary**
9 **Fig. 23**).

10

11 **Grouping promoters according to expression noise**

12 Each noise-mean relationship was fitted by the empirical formula $CV^2 = C_1 + C_2/Mean$
13 (**Fig. 5a and 5b**), and the residual of each expression pattern was calculated and sorted. In
14 addition, to ensure the reliability of the analysis, we only considered the sequences with
15 appropriate mean expression levels (> 0.15 for the promoter library; > 0.1 for the
16 combination library). The 20 candidates from the promoter library and 25 from the
17 combination library with the largest residues were reconstructed. For the promoter library,
18 the 10% with the largest residuals was classified as the high-noise group, the 10% with the
19 smallest residuals was classified as the low-noise group, and the remaining sequences were
20 grouped as medium-noise.

21

22 **Identification of potential RpoD-binding sites**

23 The DPinteract database contains computational predictions of possible RpoD-binding sites
24 with 15-19 nucleotide spacing in the *E. coli* genome⁴⁵. We searched these sequences in the
25 promoter library and counted the number of potential RpoD-binding sites of each promoter.
26 To avoid redundant results, we only accounted for independent hexamer pairs. Specifically, if
27 the -35 or -10 region of two RpoD-binding sites overlapped, we only considered the RpoD-
28 binding site with a higher z score; otherwise, both of them were retained (**Supplementary**
29 **Fig. 31**).

30

1 DATA AVAILABILITY

2 Raw NGS data of Sort-Seq have been deposited into the NCBI Short Read Archive with
3 BioProject accession number PRJNA800535. The plasmid maps related to this work can be
4 accessed via Github (<https://github.com/fenghuibao/dSort-Seq>). The dSort-Seq calculation
5 tool can be accessed through our laboratory website (<http://www.thu-big.net/dsort-seq/>).

6

7 ACKNOWLEDGMENTS

8 We would like to thank Drs. X. Zheng and Y. Zhou for kindly sharing their experimental
9 skills and data. We thank Dr. Y. Wu for construction of the pMPTPV plasmid. We thank B.
10 Yu for her help with the FACS experiments and F. Liu for her help with NGS library
11 construction. We thank Prof. F. Zhang (Washington University in St. Louis) for critical
12 discussions regarding this work. This work was supported by the National Key Research and
13 Development Program of China (2019YFA0904800), the National Natural Science
14 Foundation of China (U2032210), and the Foshan-Tsinghua Innovation Special Fund
15 (THFS01).

16

17 Author Contributions

18 H.F. conceived the general framework of this work, including the model construction and
19 experimental design. H.F. and F.L. prepared the promoter library and performed the raw Sort-
20 Seq experiment and the subsequent validation experiment. F.L. did the experiment for the
21 combination library. H.F. analyzed all experimental results. H.F., F.L. and C.Z wrote the
22 manuscript based on discussions and contributions of all authors. T.W. provided valuable
23 opinions for the work and helped to revise the manuscript. C.Z., A.Z. and X.X. supervised the
24 project.

25

26 CONFLICT OF INTEREST

27 The authors declare that there is no conflict of interest regarding the publication of this
28 article.

1

2 REFERENCES

- 3 1. Elowitz MB, Levine AJ, Siggia ED, Swain PS. Stochastic gene expression in a single
4 cell. *Science* (80-). 2002;297(5584):1183-1186. doi:10.1126/science.1070919
- 5 2. Li GW, Xie XS. Central dogma at the single-molecule level in living cells. *Nature*.
6 2011;475(7356):308-315. doi:10.1038/nature10315
- 7 3. Friedman N, Cai L, Xie XS. Linking stochastic dynamics to population distribution:
8 An analytical framework of gene expression. *Phys Rev Lett*. 2006;97(16):1-4.
9 doi:10.1103/PhysRevLett.97.168302
- 10 4. Shahrezaei V, Swain PS. Analytical distributions for stochastic gene expression. *Proc*
11 *Natl Acad Sci U S A*. 2008;105(45):17256-17261. doi:10.1073/pnas.0803850105
- 12 5. Salman H, Brenner N, Tung CK, et al. Universal protein fluctuations in populations of
13 microorganisms. *Phys Rev Lett*. 2012;108(23):1-5.
14 doi:10.1103/PhysRevLett.108.238105
- 15 6. Beal J. Biochemical complexity drives log-normal variation in genetic expression. *Eng*
16 *Biol*. 2017;1(1):55-60. doi:10.1049/enb.2017.0004
- 17 7. Heltberg ML, Krishna S, Jensen MH. On chaotic dynamics in transcription factors and
18 the associated effects in differential gene regulation. *Nat Commun*. 2019;10(1):1-10.
19 doi:10.1038/s41467-018-07932-1
- 20 8. Raj A, van Oudenaarden A. Nature, Nurture, or Chance: Stochastic Gene Expression
21 and Its Consequences. *Cell*. 2008;135(2):216-226. doi:10.1016/j.cell.2008.09.050
- 22 9. Toya Y, Shimizu H. Flux controlling technology for central carbon metabolism for
23 efficient microbial bio-production. *Curr Opin Biotechnol*. 2020;64:169-174.
24 doi:10.1016/j.copbio.2020.04.003
- 25 10. Xiao Y, Bowen CH, Liu D, Zhang F. Exploiting nongenetic cell-to-cell variation for
26 enhanced biosynthesis. *Nat Chem Biol*. 2016;12(5):339-344.
27 doi:10.1038/nchembio.2046
- 28 11. Fojo T, Bates S. Strategies for reversing drug resistance. *Oncogene*. 2003;22(47 REV.
29 ISS. 6):7512-7523. doi:10.1038/sj.onc.1206951
- 30 12. Saunders NA, Simpson F, Thompson EW, et al. Role of intratumoural heterogeneity in
31 cancer drug resistance: Molecular and clinical perspectives. *EMBO Mol Med*.
32 2012;4(8):675-684. doi:10.1002/emmm.201101131

- 1 13. Balaban NQ, Merrin J, Chait R, Kowalik L, Leibler S. Bacterial persistence as a
2 phenotypic switch. *Science* (80-). 2004;305(5690):1622-1625.
3 doi:10.1126/science.1099390
- 4 14. Newman JRS, Ghaemmaghami S, Ihmels J, et al. Single-cell proteomic analysis of *S.*
5 *cerevisiae* reveals the architecture of biological noise. *Nature*. 2006;441(7095):840-
6 846. doi:10.1038/nature04785
- 7 15. Deloupy A, Sauveplane V, Robert J, Aymerich S, Jules M, Robert L. Extrinsic noise
8 prevents the independent tuning of gene expression noise and protein mean abundance
9 in bacteria. *Sci Adv*. 2020;6(41). doi:10.1126/sciadv.abc3478
- 10 16. Silander OK, Nikolic N, Zaslaver A, et al. A genome-wide analysis of promoter-
11 mediated phenotypic noise in *Escherichia coli*. *PLoS Genet*. 2012;8(1).
12 doi:10.1371/journal.pgen.1002443
- 13 17. Kosuri S, Goodman DB, Cambray G, et al. Composability of regulatory sequences
14 controlling transcription and translation in *Escherichia coli*. *Proc Natl Acad Sci U S A*.
15 2013;110(34):14024-14029. doi:10.1073/pnas.1301301110
- 16 18. Townshend B, Kennedy AB, Xiang JS, Smolke CD. High-throughput cellular RNA
17 device engineering. *Nat Methods*. 2015;12(10):989-994. doi:10.1038/nmeth.3486
- 18 19. Peterman N, Levine E. Sort-seq under the hood: Implications of design choices on
19 large-scale characterization of sequence-function relations. *BMC Genomics*.
20 2016;17(1):1-17. doi:10.1186/s12864-016-2533-5
- 21 20. Sharon E, Kalma Y, Sharp A, et al. Inferring gene regulatory logic from high-
22 throughput measurements of thousands of systematically designed promoters. *Nat*
23 *Biotechnol*. 2012;30(6):521-530. doi:10.1038/nbt.2205
- 24 21. Sharon E, Van Dijk D, Kalma Y, et al. Probing the effect of promoters on noise in
25 gene expression using thousands of designed sequences. *Genome Res*.
26 2014;24(10):1698-1706. doi:10.1101/gr.168773.113
- 27 22. de Boer CG, Vaishnav ED, Sadeh R, Abeyta EL, Friedman N, Regev A. Deciphering
28 eukaryotic gene-regulatory logic with 100 million random promoters. *Nat Biotechnol*.
29 2020;38(1):56-65. doi:10.1038/s41587-019-0315-8
- 30 23. Noderer WL, Flockhart RJ, Bhaduri A, et al. Quantitative analysis of mammalian
31 translation initiation sites by FACS -seq . *Mol Syst Biol*. 2014;10(8):748.
32 doi:10.15252/msb.20145136

- 1 24. Peterman N, Lavi-Itzkovitz A, Levine E. Large-scale mapping of sequence-function
2 relations in small regulatory RNAs reveals plasticity and modularity. *Nucleic Acids*
3 *Res.* 2014;42(19):12177-12188. doi:10.1093/nar/gku863
- 4 25. Rutherford ST, Valastyan JS, Taillefumier T, Wingreen NS, Bassler BL.
5 Comprehensive analysis reveals how single nucleotides contribute to noncoding RNA
6 function in bacterial quorum sensing. *Proc Natl Acad Sci U S A.* 2015;112(44).
7 doi:10.1073/pnas.1518958112
- 8 26. Hawkins JS, Silvis MR, Koo BM, et al. Mismatch-CRISPRi Reveals the Co-varying
9 Expression-Fitness Relationships of Essential Genes in Escherichia coli and Bacillus
10 subtilis. *Cell Syst.* 2020;11(5):523-535.e9. doi:10.1016/j.cels.2020.09.009
- 11 27. Wang T, Zheng X, Ji H, Wang TL, Xing XH, Zhang C. Dynamics of transcription-
12 translation coordination tune bacterial indole signaling. *Nat Chem Biol.*
13 2020;16(4):440-449. doi:10.1038/s41589-019-0430-3
- 14 28. Schmitz A, Zhang F. Massively parallel gene expression variation measurement of a
15 synonymous codon library. *BMC Genomics.* 2021;22(1):1-12. doi:10.1186/s12864-
16 021-07462-z
- 17 29. Bishop CM. *Pattern Recognition and Machine Learning.*; 2006.
- 18 30. Taniguchi Y, Choi PJ, Li G-W, et al. Quantifying E. coli Proteome and Transcriptome
19 with Single-Molecule Sensitivity in Single Cells. *Science (80-).* 2010;329(5991):533-
20 538. doi:10.1126/science.1188308
- 21 31. Ferguson TS. *BAYESIAN DENSITY ESTIMATION BY MIXTURES OF NORMAL*
22 *DISTRIBUTIONS.* ACADEMIC PRESS, INC.; 1983. doi:10.1016/b978-0-12-589320-
23 6.50018-6
- 24 32. Marks RG, Rao P V. An estimation procedure for data containing outliers with a one-
25 directional shift in the mean. *J Am Stat Assoc.* 1979;74(367):614-620.
26 doi:10.1080/01621459.1979.10481657
- 27 33. Aitkin M, Wilson GT. Mixture models, outliers, and the em algorithm. *Technometrics.*
28 1980;22(3):325-331. doi:10.1080/00401706.1980.10486163
- 29 34. Beckman RJ, Cook RD. Outlier s. *Technometrics.* 1983;25(2):119-149.
30 doi:10.1080/00401706.1983.10487840
- 31 35. Nielsen AAK, Der BS, Shin J, et al. Genetic circuit design automation. *Science (80-).*
32 2016;352(6281). doi:10.1126/science.aac7341

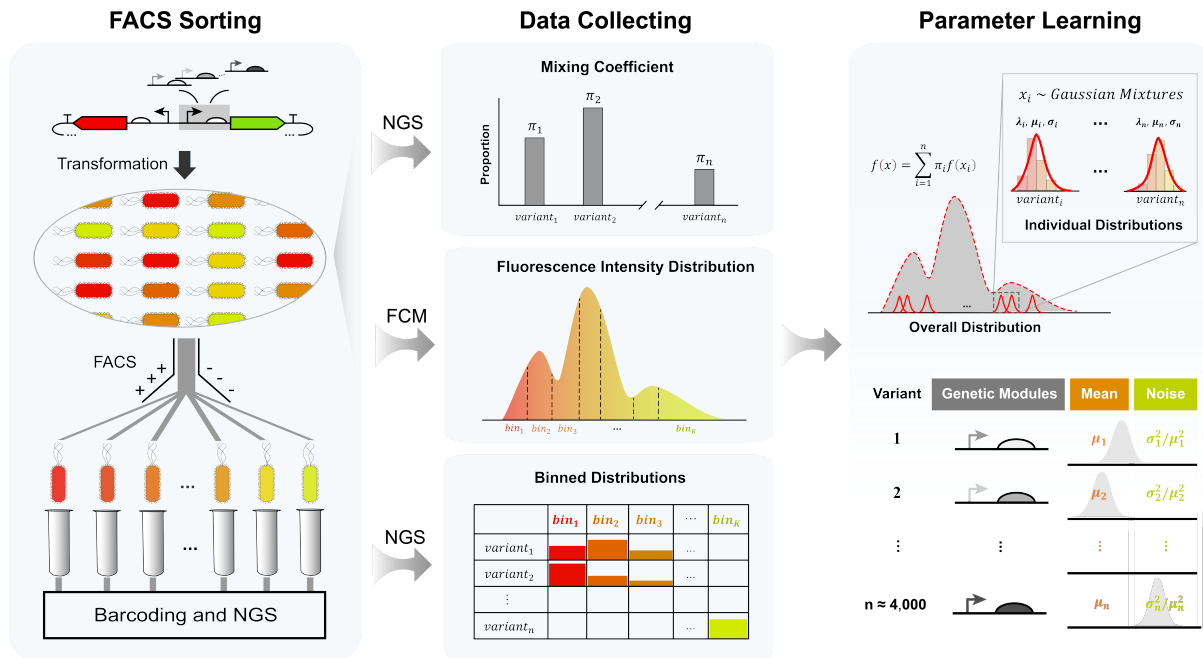
- 1 36. Goodfellow IJ, Pouget-Abadie J, Mirza M, et al. Generative Adversarial Networks.
2 *Adv Neural Inf Process Syst.* 2014;2017-Decem:4089-4099.
3 <http://arxiv.org/abs/1406.2661>
- 4 37. Rogers JK, Taylor ND, Church GM. Biosensor-based engineering of biosynthetic
5 pathways. *Curr Opin Biotechnol.* 2016;42:84-91. doi:10.1016/j.copbio.2016.03.005
- 6 38. Zhou Y, Yuan Y, Wu Y, et al. Encoding Genetic Circuits with DNA Barcodes Paves
7 the Way for Machine Learning-Assisted Metabolite Biosensor Response Curve
8 Profiling in Yeast. *ACS Synth Biol.* 2022;11(2):977-989.
9 doi:10.1021/acssynbio.1c00595
- 10 39. Ansel J, Bottin H, Rodriguez-Beltran C, et al. Cell-to-cell stochastic variation in gene
11 expression is a complex genetic trait. *PLoS Genet.* 2008;4(4).
12 doi:10.1371/journal.pgen.1000049
- 13 40. Ozbudak EM, Thattai M, Kurtser I, Grossman AD, Van Oudenaarden A. Regulation of
14 noise in the expression of a single gene. *Nat Genet.* 2002;31(1):69-73.
15 doi:10.1038/ng869
- 16 41. Keseler IM, Gama-Castro S, Mackie A, et al. The EcoCyc Database in 2021. *Front*
17 *Microbiol.* 2021;12(July):1-10. doi:10.3389/fmicb.2021.711077
- 18 42. Barrios D, Prieto C. D3GB: An Interactive Genome Browser for R, Python, and
19 WordPress. *J Comput Biol.* 2017;24(5):447-449. doi:10.1089/cmb.2016.0213
- 20 43. Mutalik VK, Guimaraes JC, Cambray G, et al. Precise and reliable gene expression via
21 standard transcription and translation initiation elements. *Nat Methods.*
22 2013;10(4):354-360. doi:10.1038/nmeth.2404
- 23 44. Lisser S, Margalit H. Compilation of E.coli mRNA promoter sequences. *Nucleic Acids*
24 *Res.* 1993;21(7):1507-1516. doi:10.1093/nar/21.7.1507
- 25 45. Robison K, McGuire AM, Church GM. A comprehensive library of DNA-binding site
26 matrices for 55 proteins applied to the complete Escherichia coli K-12 genome. *J Mol*
27 *Biol.* 1998;284(2):241-254. doi:10.1006/jmbi.1998.2160
- 28 46. Feng H, Yuan Y, Yang Z, Xing X hui, Zhang C. Genome-wide genotype-phenotype
29 associations in microbes. *J Biosci Bioeng.* 2021;132(1):1-8.
30 doi:10.1016/j.jbiosc.2021.03.011
- 31 47. Yang D, Park SY, Park YS, Eun H, Lee SY. Metabolic Engineering of Escherichia coli
32 for Natural Product Biosynthesis. *Trends Biotechnol.* 2020;38(7):745-765.
33 doi:10.1016/j.tibtech.2019.11.007

- 1 48. Chang HJ, Voyvodic PL, Zúñiga A, Bonnet J. Microbially derived biosensors for
2 diagnosis, monitoring and epidemiology. *Microb Biotechnol.* 2017;10(5):1031-1035.
3 doi:10.1111/1751-7915.12791
- 4 49. Sabatini BL, Tian L. Imaging Neurotransmitter and Neuromodulator Dynamics In
5 Vivo with Genetically Encoded Indicators. *Neuron.* 2020;108(1):17-32.
6 doi:10.1016/j.neuron.2020.09.036
- 7 50. Muranaka N, Sharma V, Nomura Y, Yokobayashi Y. An efficient platform for genetic
8 selection and screening of gene switches in *Escherichia coli*. *Nucleic Acids Res.*
9 2009;37(5):1-9. doi:10.1093/nar/gkp039
- 10 51. Liang JC, Chang AL, Kennedy AB, Smolke CD. A high-throughput, quantitative cell-
11 based screen for efficient tailoring of RNA device activity. *Nucleic Acids Res.*
12 2012;40(20):1-14. doi:10.1093/nar/gks636
- 13 52. Eldar A, Elowitz MB. Functional roles for noise in genetic circuits. *Nature.*
14 2010;467(7312):167-173. doi:10.1038/nature09326
- 15 53. Farquhar KS, Charlebois DA, Szenk M, Cohen J, Nevozhay D, Balázsi G. Role of
16 network-mediated stochasticity in mammalian drug resistance. *Nat Commun.*
17 2019;10(1):1-14. doi:10.1038/s41467-019-10330-w
- 18 54. Hooshangi S, Thiberge S, Weiss R. Ultrasensitivity and noise propagation in a
19 synthetic transcriptional cascade. *Proc Natl Acad Sci U S A.* 2005;102(10):3581-3586.
20 doi:10.1073/pnas.0408507102
- 21 55. Payne JL, Wagner A. The causes of evolvability and their evolution. *Nat Rev Genet.*
22 2019;20(1):24-38. doi:10.1038/s41576-018-0069-z
- 23 56. Bódi Z, Farkas Z, Nevozhay D, et al. Phenotypic heterogeneity promotes adaptive
24 evolution. *PLoS Biol.* 2017;15(5):1-26. doi:10.1371/journal.pbio.2000644
- 25 57. Schmutzer M, Wagner A. Gene expression noise can promote the fixation of beneficial
26 mutations in fluctuating environments. *PLoS Comput Biol.* 2020;16(10):1-24.
27 doi:10.1371/journal.pcbi.1007727
- 28 58. Aranda-Díaz A, Mace K, Zuleta I, Harrigan P, El-Samad H. Robust Synthetic Circuits
29 for Two-Dimensional Control of Gene Expression in Yeast. *ACS Synth Biol.*
30 2017;6(3):545-554. doi:10.1021/acssynbio.6b00251
- 31 59. Mundt M, Anders A, Murray SM, Sourjik V. A System for Gene Expression Noise
32 Control in Yeast. *ACS Synth Biol.* 2018;7(11):2618-2626.
33 doi:10.1021/acssynbio.8b00279

- 1 60. Benzinger D, Khammash M. Pulsatile inputs achieve tunable attenuation of gene
2 expression variability and graded multi-gene regulation. *Nat Commun.* 2018;9(1).
3 doi:10.1038/s41467-018-05882-2
- 4 61. Gerhardt KP, Rao SD, Olson EJ, Igoshin OA, Tabor JJ. Independent control of mean
5 and noise by convolution of gene expression distributions. *Nat Commun.*
6 2021;12(1):1-10. doi:10.1038/s41467-021-27070-5
- 7 62. Chauhan V, Bahrudeen MNM, Palma CSD, et al. Analytical kinetic model of native
8 tandem promoters in *E. coli*. *PLoS Comput Biol.* 2022;18(1):1-23.
9 doi:10.1371/journal.pcbi.1009824
- 10 63. Liu R, Bassalo MC, Zeitoun RI, Gill RT. Genome scale engineering techniques for
11 metabolic engineering. *Metab Eng.* 2015;32:143-154.
12 doi:10.1016/j.ymben.2015.09.013
- 13 64. Chalancon G, Ravarani CNJ, Balaji S, et al. Interplay between gene expression noise
14 and regulatory network architecture. *Trends Genet.* 2012;28(5):221-232.
15 doi:10.1016/j.tig.2012.01.006
- 16 65. Hartline CJ, Schmitz AC, Han Y, Zhang F. Dynamic control in metabolic engineering:
17 Theories, tools, and applications. *Metab Eng.* 2021;63(August 2020):126-140.
18 doi:10.1016/j.ymben.2020.08.015
- 19 66. Belliveau NM, Barnes SL, Ireland WT, et al. Systematic approach for dissecting the
20 molecular mechanisms of transcriptional regulation in bacteria. *Proc Natl Acad Sci U*
21 *S A.* 2018;115(21):E4796-E4805. doi:10.1073/pnas.1722055115
- 22 67. Whitehead TA, Chevalier A, Song Y, et al. Optimization of affinity, specificity and
23 function of designed influenza inhibitors using deep sequencing. *Nat Biotechnol.*
24 2012;30(6):543-548. doi:10.1038/nbt.2214
- 25 68. Podgornaia AI, Laub MT. Pervasive degeneracy and epistasis in a protein-protein
26 interface. *Science (80-).* 2015;347(6222):673-677. doi:10.1126/science.1257360
- 27 69. Lu J, Tang J, Liu Y, Zhu X, Zhang T, Zhang X. Combinatorial modulation of galP and
28 glk gene expression for improved alternative glucose utilization. *Appl Microbiol*
29 *Biotechnol.* 2012;93(6):2455-2462. doi:10.1007/s00253-011-3752-y
- 30 70. Magoč T, Salzberg SL. FLASH: Fast length adjustment of short reads to improve
31 genome assemblies. *Bioinformatics.* 2011;27(21):2957-2963.
32 doi:10.1093/bioinformatics/btr507

- 1 71. Wang T, Guan C, Guo J, et al. Pooled CRISPR interference screening enables
- 2 genome-scale functional genomics study in bacteria with superior performance-net.
- 3 *Nat Commun.* 2018;9(1). doi:10.1038/s41467-018-04899-x
- 4
- 5

1



2

3

4 **Figure 1** Schematic overview of the dSort-Seq data workflow. (a) During Sort-Seq, a library

5 with different expression patterns is sorted into customized bins based on the fluorescence

6 intensity value. (b) The mixing coefficients are quantified via next-generation sequencing

7 (NGS). (c) The overall fluorescence density is measured by flow cytometry (FCM), and the

8 sorting boundaries are specified based on the overall fluorescence intensity density. (d) The

9 read count number across all bins as quantified by NGS reveals the binned distribution of

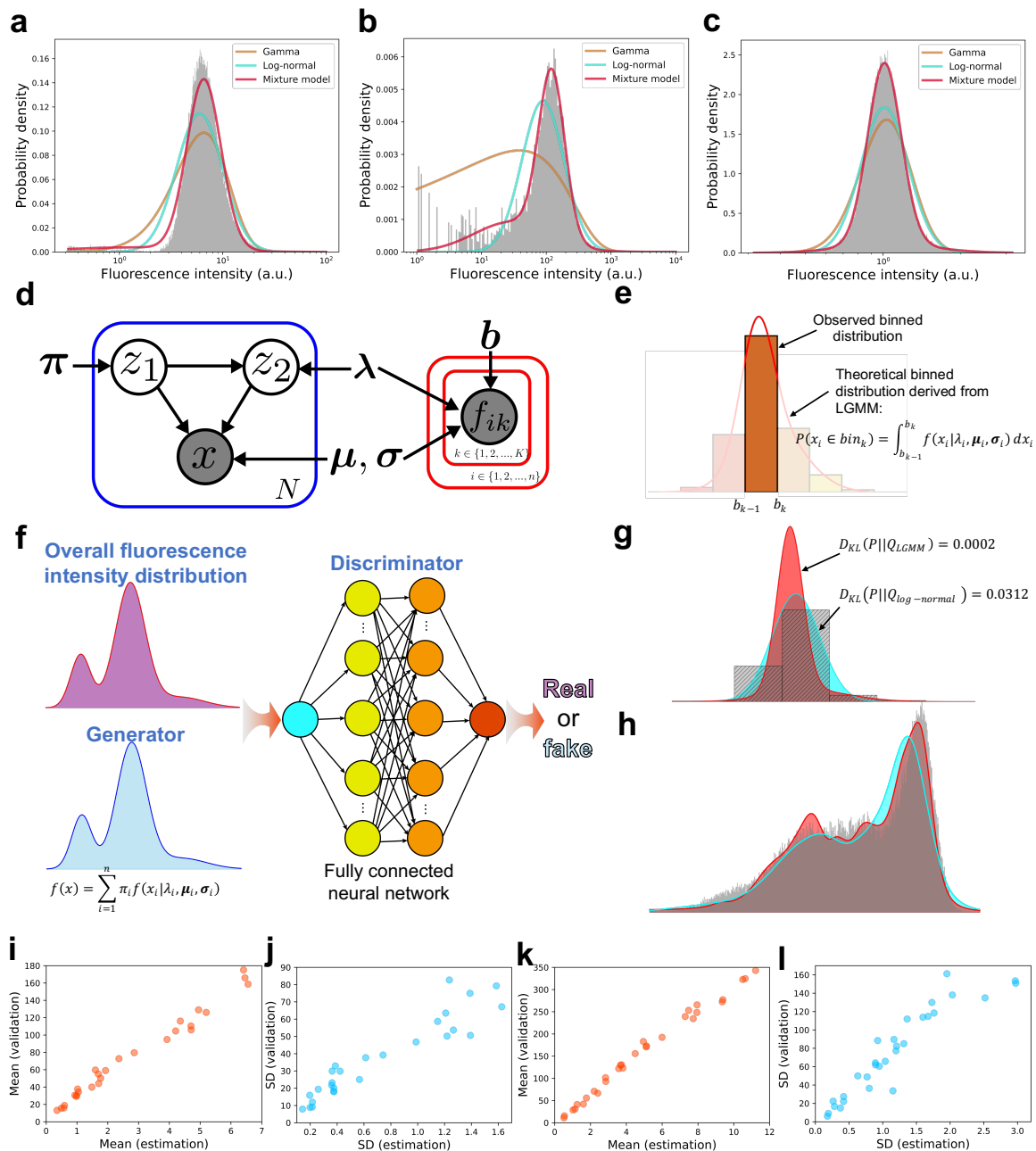
10 each variant in the library. (e) Through parameter learning, the mean, expression noise and

11 their relationships can be precisely identified.

12

13

1



2

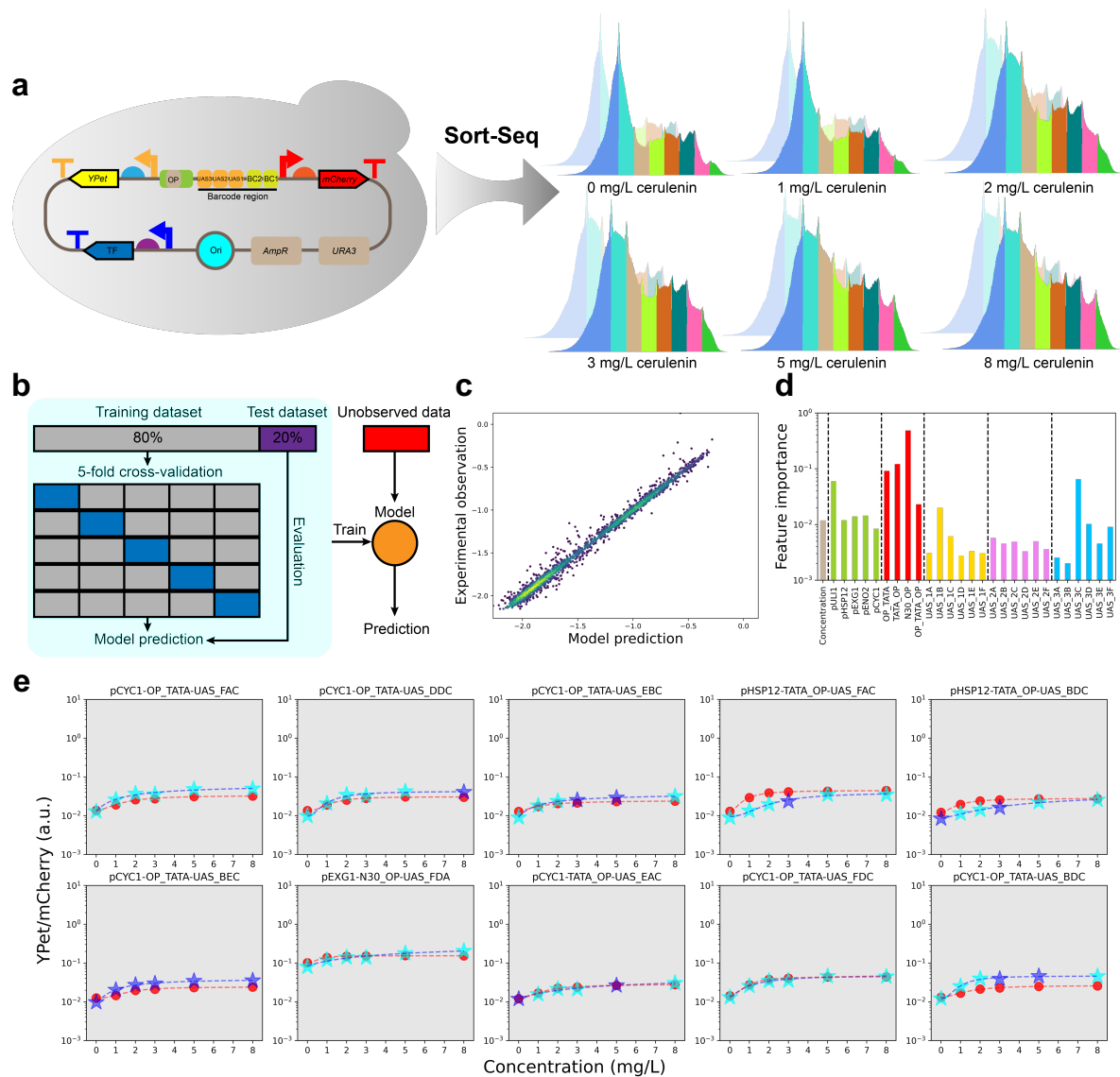
3

4 **Figure 2** Framework and performance of dSort-Seq. (a-c) Two-component log-mixture of
 5 Gaussians can better represent the gene expression distribution compared with conventionally
 6 used model-driven methods. (a) Gene expression controlled by the LmrA repressor³⁵; the
 7 histogram denotes the cytometry data of the unrepresed state. (b) Density of gene expression
 8 under the control of the *tnaC* variant K11R_CGC²⁷. The data were measured under 100 μ M
 9 Ala-Trp. (c) Gene expression driven by the promoter *yebVp2* (this study). In a-c, the red,
 10 cyan and brown lines represent the fitting result of the two-component log-mixture of
 11 Gaussian, log-normal and gamma distributions, respectively. (d) Graphical representation of

1 the model. **(e)** Theoretical fraction of the probability density within the corresponding
2 boundaries. **(f)** Matching the mixture of two-component Gaussian mixture models to the
3 overall fluorescence intensity distribution. The real data are sampled from experimental
4 cytometry data; the fake data are generated from the LGMM. A fully connected neural
5 network is used as a discriminator to determine whether the data are real or fake. **(g)** An
6 example (V8A_GCC, 0 μ M Ala-Trp, replicate 1) to illustrate the superior performance of
7 dSort-Seq in matching the binned distribution compared to the log-normal-based method.
8 The Kullback–Leibler divergence shows the performance of each fit. **(h)** An example (100
9 μ M Ala-Trp, replicate 1) to illustrate the superior performance of dSort-Seq in matching the
10 overall fluorescence distribution compared to the log-normal-based method. In **(g)** and **(h)**,
11 the red and cyan distributions refer to the results derived from dSort-Seq and the log-normal-
12 based method, respectively. The gray distribution refers to the real data. **(i-l)** Individually
13 analyzed expression characteristics of reconstructed *tnaC* variants by cytometry were highly
14 correlated with those estimated via dSort-Seq in terms of both their means (**i**, 0 μ M Ala-Trp,
15 $n = 26$; **k**, 100 μ M Ala-Trp, $n = 30$) and SDs (**j**, 0 μ M Ala-Trp; **l**, 100 μ M Ala-Trp).

16

17

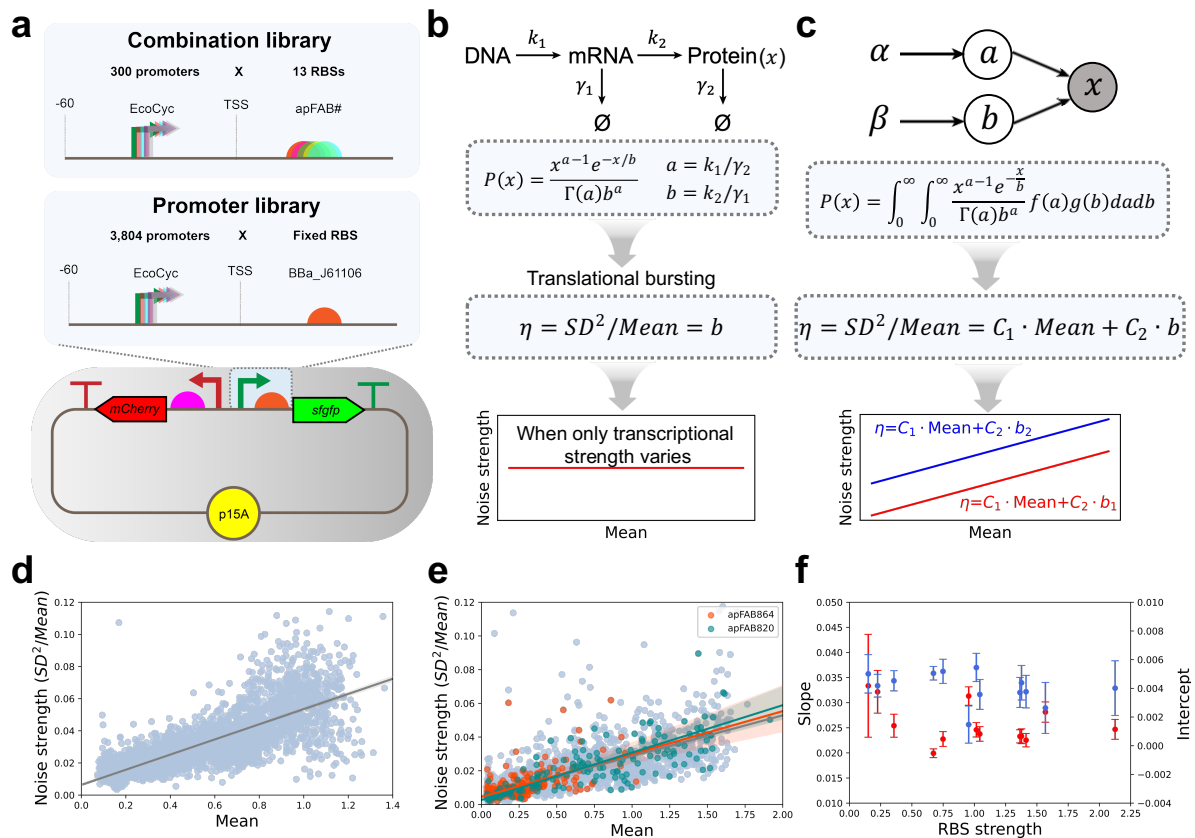


1
2

3 **Figure 3** The dSort-Seq profiling of FapR-*fapO*-based malonyl-CoA-dependent gene
4 expression. **(a)** Sort-Seq characterization of the malonyl-CoA biosensor library under 6
5 different cerulenin concentrations (0, 1, 2, 3, 5, 8 mg/L). Cells were sorted into 8 bins
6 according to their responses to ligand. Two biological replicates were examined for each
7 Sort-Seq experiment. **(b)** Schematic diagram of the machine learning process. Gradient
8 boosting regression was used here to interpret the relationship between features and
9 expression strengths. The hyperparameters were optimized through 5-fold cross-validation;
10 then, the whole training dataset was used to train the model parameters, and the test dataset
11 was used to evaluate the generalization capacity of the model. Finally, the model was trained
12 on the entire observed dataset to obtain predictions for unobserved data. **(c)** The model
13 performance in the test dataset showed a good generalization capacity ($n = 2,275$). **(d)** Gini
14 importance that contributes to the gradient boosting regression tree. **(e)** Dose–response curves

1 of the top 10 combinations with the highest dynamic ranges. Data points represent the mean
2 values of YPet/mCherry under different cerulenin concentrations, where red dots represent
3 individual characterization data, cyan stars represent data from dSort-Seq characterizations,
4 and blue stars denote data from machine learning predictions. The dashed lines represent
5 response curves fitted by the Hill equation (see **Methods**).
6

1

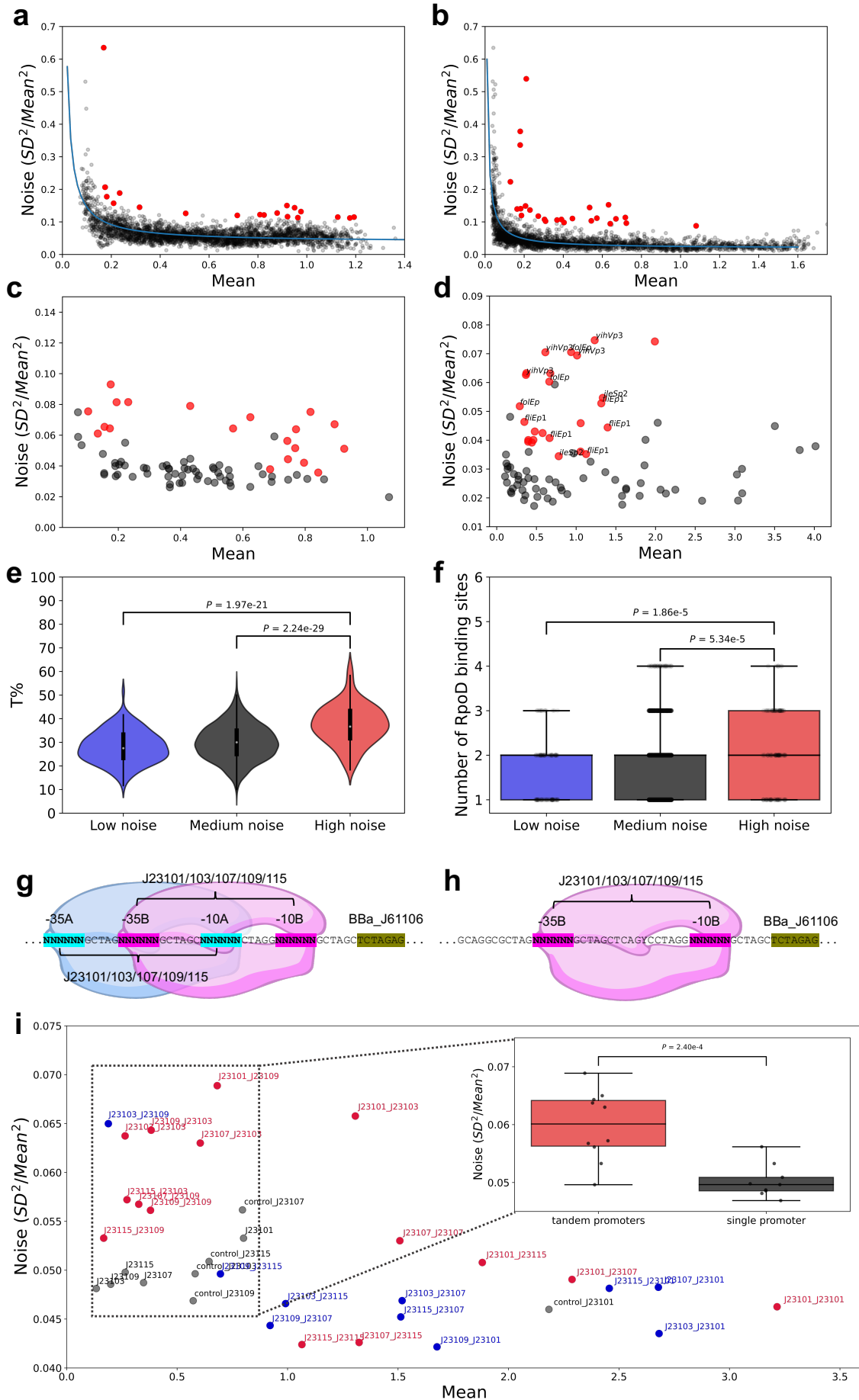


2

3

4 **Figure 4** The dSort-Seq profiling of transcriptional and translational effects on noise
5 production in *E. coli* K12 MG1655. **(a)** The design schemes of the promoter and the
6 combination library. **(b)** According to the translational bursting mechanism, steady-state
7 protein production follows a gamma distribution³; as a corollary, the burst size, denoted by
8 the Fano factor, is linearly correlated with the translation rate and independent of the
9 transcription rate. **(c)** According to the hierarchical Bayesian model, the intercept in the
10 relationship between noise strength and the mean expression level is proportional to
11 translational strength, indicating that translational bursting still dominates noise production at
12 low expression levels³⁰. **(d)** The noise strength is linearly correlated with the mean expression
13 level when only transcriptional strength varies. The gray line exhibits the linear regression
14 result, which is shaded to show the 95% confidence interval. **(e)** The relationships between
15 noise strength and mean expression level are similar when the translation module varies. The
16 gray, orange, and green lines represent the regressions of all combinations and combinations
17 with RBS apFAB864 and apFAB820, respectively. (RBS strength: apFAB820 (1.57) >
18 apFAB864 (0.36), see **Methods**). **(f)** The linear regression slopes and intercepts of noise

- 1 strength and mean expression level are not significantly correlated with translational strength.
- 2 The error bars indicate the 95% confidence intervals.
- 3
- 4



1 **Figure 5** Overlapping RpoD-binding sites result in high expression noise. **(a and b)**
2 Correlation of the expression noise with expression strength in the **(a)** promoter library and
3 **(b)** combination library. At low mean expression levels, the noise decreases as the expression
4 strength increases; at high mean expression levels, the noise converges to a constant value.
5 The blue lines show the regression results (see **Methods**). Twenty promoters and 25
6 combinations exhibiting high expression noise are marked as red dots. **(c)** Twenty sequences
7 from the promoter library and **(d)** 25 sequences from the combination library showing high
8 expression noise were constructed and assayed through FCM. As a result, their expression
9 noise (indicated by red dots) is higher than that of randomly selected variants (indicated by
10 black dots) at their corresponding mean expression levels. **(e)** The high-noise group had a
11 significantly higher thymine content than the low-noise group ($P = 1.97e-21$, one-tailed t test)
12 and the medium-noise group ($P = 2.24e-29$). **(f)** The number of potential RpoD-binding sites
13 in the high-noise group was significantly higher than that in the low-noise group ($P = 1.86e-$
14 5 , one-tailed t test) as well as in the medium-noise group ($P = 5.34e-5$). **(g)** Design scheme of
15 25 tandem promoters, each containing two overlapping RpoD-binding sites. **(h)** Design
16 scheme of 5 constitutive promoters with the same length as the tandem promoter, each with
17 only one RpoD-binding site. **(i)** Compared to promoters with a single RpoD-binding site, the
18 tandem promoters exhibited significantly higher expression noise ($P = 2.40e-4$, one-tailed t
19 test), especially when the stronger promoter was located upstream of the weaker promoter
20 (indicated by red dots).