

Randomized gates eliminate bias in sort-seq assays

Brian L. Trippe^{1,2,3,*}, Buwei Huang^{3,4}, Erika A. DeBenedictis^{3,4}, Brian Coventry^{3,4}, Nicholas Bhattacharya^{2,5}, Kevin K. Yang², David Baker^{3,4,6}, and Lorin Crawford^{2,*}

¹Massachusetts Institute of Technology, Cambridge, MA, USA

²Microsoft Research New England, Cambridge, MA, USA

³Institute for Protein Design, University of Washington, Seattle, WA, USA

⁴Department of Biochemistry, University of Washington, Seattle, WA, USA

⁵Department of Mathematics, University of California Berkeley, Berkeley, CA, USA

⁶Howard Hughes Medical Institute, University of Washington, Seattle, WA, USA

February 17, 2022

Sort-seq assays are a staple of the biological engineering toolkit, allowing researchers to profile many groups of cells based on any characteristic that can be tied to fluorescence. However, current approaches, which segregate cells into bins deterministically based on their measured fluorescence, introduce systematic bias. We describe a surprising result: one can obtain unbiased estimates by incorporating randomness into sorting. We validate this approach in simulation and experimentally, and describe extensions for both estimating group level variances and for using multi-bin sorters.

Quantitative, multiplexed assays relying on fluorescence activated cell sorting (FACS) followed by high-throughput sequencing are critical to modern biology and molecular engineering because they enable construction of large scale datasets connecting sequence to function. For example, these “sort-seq” assays are widely used to profile the strength of protein-protein binding interactions via yeast display^{4;5;11;17}. In particular one (i) synthesizes a *library* of 10^4 to 10^5 DNA sequences encoding proteins that may bind to a target of interest; (ii) transforms the library into yeast such that each putative binder is expressed on the surface of a population of cells; (iii) incubates cells with fluorescently labeled target protein; (iv) physically separates 10^6 to 10^8 cells based on binding affinity by FACS; and finally, (v) quantifies the prevalence, and thereby binding affinity, of each library member by high throughput sequencing. Due to biological and technical variability, there is a distribution over (log) fluorescence for each library sequence, and the challenge is to estimate the means of each of these distributions (Figure 1A-B). For example, for binding interactions, this mean fluorescence relates directly to biophysical quantities of interest including dissociation constants and binding energies^{1;14;17}.

In previous work, cells are deterministically segregated into one or more collection tubes (referred to as “bins”) based on their measured fluorescence, and the mean fluorescence of each population is estimated from the histogram of observed sequence counts in each bin (Figure 1C). Peterman and Levine¹³ compare the error associated with different strategies for collecting and analyzing such data, and they show that average squared error is the sum of contributions from *bias* and *variance* (e.g., Hastie et al.¹⁰, Chapter 7.3). The variance arises from experimental noise and variability across cells, and it can be reduced by increasing the number of cells screened. The bias arises from the discretization of the space of log fluorescence into bins (Figure 1B-C); for example, narrow distributions can be sorted all into the same bin but have means as different as the bin width. Because even the most sophisticated FACS machines can sort cells into at most six bins, resolution is

limited. This low resolution limits the value of sort-seq data in quantitative analyses, for instance, by prohibiting computation of precise binding energies. This challenge has spurred much work on how to effectively reduce histogram bias^{1;8;13;14}. One common approach seeks to overcome the resolution limits of histograms by assuming fluorescence is log-normally distributed for each population and using maximum likelihood estimation to estimate moments^{5;8;9;15}. However, on real data, this assumption is violated and the resulting estimates can have greater bias than the naive approach (Figure S1).

In this work, we show that the bias generated using histograms can be eliminated altogether by incorporating randomness into FACS collection strategies with as few as two bins (Figure 1D), thereby obtaining arbitrarily accurate estimates with many cells (Figure 1E). To do this, we take a statistical approach. We consider a population of cells that pass through a 2-bin sorter, each with log fluorescence F independently and identically distributed according to a density function p_F . Our target of interest is the mean log fluorescence, $\mu_F = \int f p_F(f) df$. Let B denote the bin (either 1 or 2) into which a cell is collected, and let Y_1 and Y_2 be the counts of cells in bins 1 and 2 after sorting, respectively. In multiplexed sort-seq assays, we obtain Y_1 and Y_2 for thousands of populations, and our goal is to accurately estimate the mean of each population simultaneously.

For standard binning, a *gate* is chosen for each bin that defines the range of values F for which cells are collected into that bin; so, the bin B is deterministic once F is measured (e.g., as in Figure 1C). We instead consider *randomized gates* which define for each bin the probability of collecting a cell at each fluorescence (as in Figure 1D) and rely on pseudo-random numbers to determine the bin. For estimating population means, when the fluorescence measurements fall between lower and upper bounds L and U , one first sorts using randomized gates such that for any f on the interval $[L, U]$,

$$\mathbb{P}(B=1 | F=f) = 1 - \frac{f-L}{U-L} \quad \text{and} \quad \mathbb{P}(B=2 | F=f) = \frac{f-L}{U-L}. \quad (1)$$

The counts are then combined into an empirical estimate of μ_F as $\hat{\mu} = (U-L) \cdot Y_2 / (Y_1 + Y_2) + L$.

While one might expect introducing randomness to decrease precision by introducing additional noise, $\hat{\mu}$ is directly informative to the mean fluorescence. In particular, $\hat{\mu}$ is an *unbiased* estimate of the true population mean in the sense that the average value we would expect for $\hat{\mu}$ if we repeated the sort-seq experiment many times is equal to μ_F (see Theorem 1 in Methods).

This unbiasedness theorem guarantees that, in contrast to the histogram approach, we can get arbitrarily accurate estimates by screening a larger numbers of cells (Figure 1E and Figure S2). More precisely, recalling that the mean squared error (MSE) is the sum of the bias squared and the variance¹⁰, unbiasedness implies

*To whom correspondence should be addressed: btrippe@mit.edu, lcrawford@microsoft.com

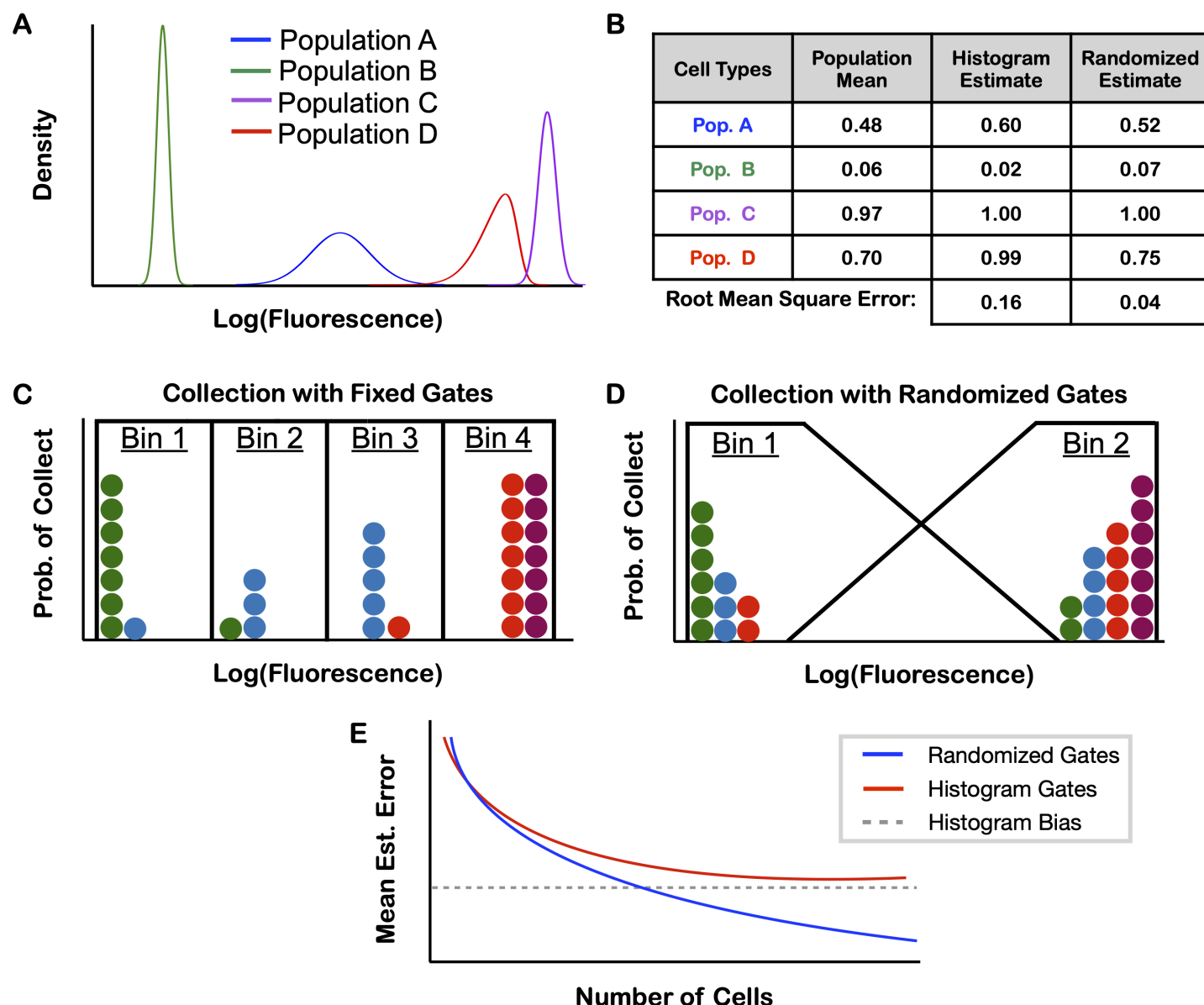


Figure 1. Schematic overview of randomized gates. (A) distributions of log fluorescence for different cell populations and (B) their hypothetical true and estimated means. (C) An example of histogram approach with deterministic collection into four bins and (D) an example of randomized collection approach with two bins. (E) Estimated means of the randomized gating scheme are more accurate than the histogram approach as the number of collected cells increases.

that the error of $\hat{\mu}$ is dictated solely by its variance. Moreover, $\hat{\mu}$ allows a transparent trade-off between the number of cells sorted per population and the precision of the estimates; notably, with as few as 400 cells, a 95% confidence interval for μ_F will cover at most 10% of the range from L to U (Methods).

We used a simulation study to explore the implications of unbiasedness on estimation accuracy with the randomized gate approach relative to the standard histogram approach. In this study, we simulated fluorescence of 250 cells from log-normal distributions with different means and variances (Figure 2A). We then simulated sorting these cells based on their fluorescence either with four deterministic gates of equal width or with two randomized gates as dictated by Equation (1). For the deterministic gates, we constructed histograms and computed estimates of the mean fluorescence as the average of the bin centers weighted by the fraction of cells they contained; and for the randomized gates, we estimated the mean as $\hat{\mu}$. Figure 2B and C report the performance of these estimates in terms of MSE, along with their bias and variance

components. As expected, the randomized gates approach has negligible bias except for broad distributions violating the conditions of our theorem (Methods).

With even as few as 250 cells per population, the MSE of the histogram approach is dominated by bias. Accordingly, the unbiased randomized approach typically provides more accurate estimates. Notably, 250 cells is fewer than is the typical in sort-seq assays; with larger samples, more pronounced improvements are obtained (Figure S2). Because the histogram estimates are systematically biased toward bin centers, they can however be more accurate for narrow distributions with means near bin centers (Figure 2B).

We next tested our approach experimentally. Current FACS software does not support randomized gate programming, so we devised an experimental approximation in which we manually changed the gating threshold 20 times during sorting at regular intervals (Methods). We tested this procedure in the context of a binding assay using yeast display³. We synthesized DNA encoding four mini-protein binders to the SARS-COV-2 receptor binding domain

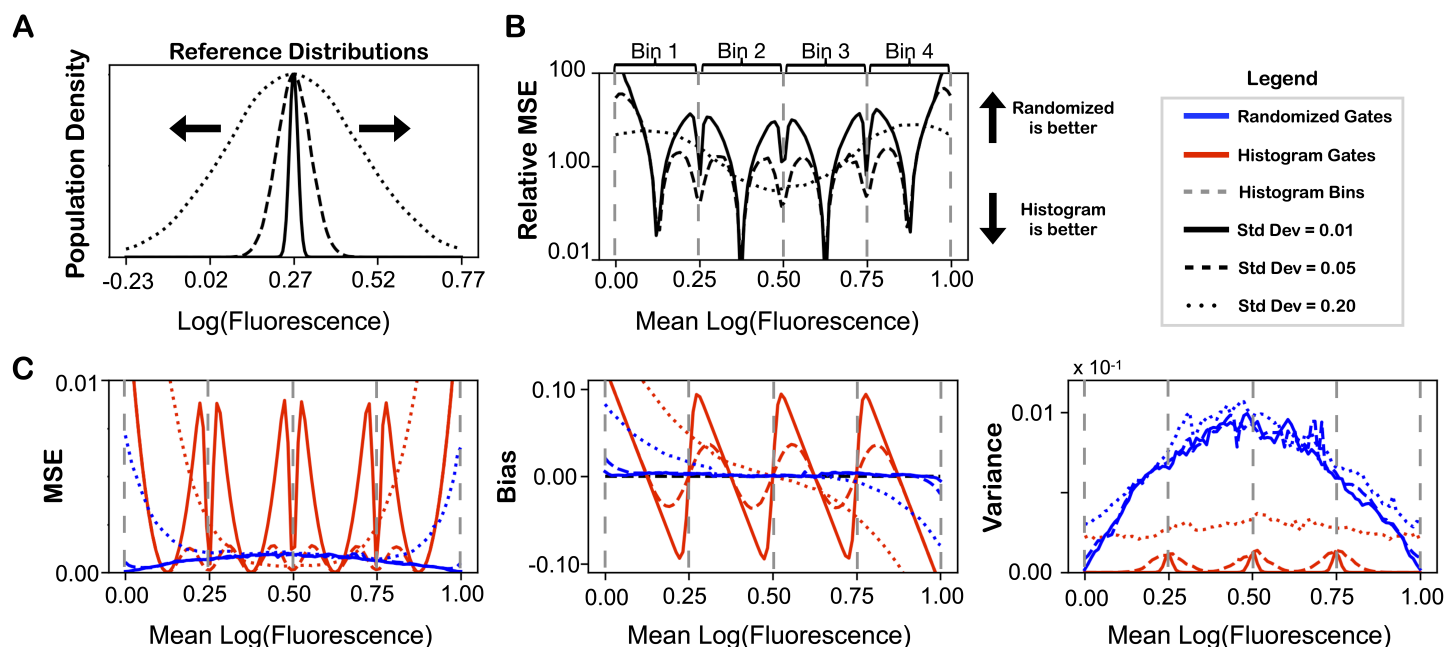


Figure 2. Simulation study reveals improved estimation properties obtained with randomized gates as compared to histograms. (A) Fluorescence values of cells are drawn independently from log-normal distributions with different scales. (B) The relative performances of estimates from histograms and randomized gates across a range of mean log fluorescences in terms of mean squared error (ratios greater than 1 reflect lower error with randomized gates and ratios below 1 reflect lower error with histograms). (C) The mean squared error (left) decomposed into bias (center) and variance (right) for both estimates. All points are the average across 200 replicates, each with $N = 250$ cells.

(RBD) with a range of binding affinities⁴. While the value of this approach is greatest for highly multiplexed assays with many thousands of sequences, we chose this small number so that we could also test each binder easily in serial. We separately transformed and expressed each design in yeast and then incubated the populations with RBD. Both the target and binders were fluorescently labeled, and we considered the log ratio of target to binder fluorescence as an expression normalized proxy for binder strength¹⁴. We measured each sample on a Sony SH800 cell sorter separately, recording the binding signal for each binder (Figure 3A). We then pooled the samples together and sorted 1,000,000 cells, collecting 50,000 cells at each of the 20 thresholds (Methods).

The multiplexed measurements largely recapitulate the ground-truth clonal measurements (Figure 3B), with the exception of design candidate 2018, for which the multiplexed estimate is below the clonal one. We suspect this is due to dissociation of some of the target protein in the time between the clonal and multiplexed measurements; kinetics experiments suggest dissociation occurs rapidly for this design⁴.

In the supplementary note, we additionally describe two extensions of this idea. First, because the differences in the variability of fluorescence across each population is often of interest (in addition to mean fluorescence), we show how to extend the approach to estimate the variance for each population. Second, we describe how to effectively take advantage of sorters that sort into more than two bins simultaneously to obtain more accurate estimates. We view these contributions as a starting point for future work of using randomness to obtain precise, multiplexed estimates.

We have shown how to obtain precise, multiplexed estimates in sort-seq experiments with a simple strategy that incorporates randomness. This mathematical technique allows better data to be collected using the same or less sophisticated hardware. While we have emphasized studies of binding affinity, we believe our strategy is applicable to a wider range of applications of sort-seq assays including studying transcriptional regulation^{9,16} and protein sta-

bility¹⁵, and building datasets for protein design². Widespread implementation of randomized gates in FACS and community adoption of this strategy, will greatly simplify and improve sort-seq assays by eliminating a common bias in this ubiquitous assay. We believe this will allow FACS to play a more central role in screening settings, for construction of reliable datasets for machine learning models in bio-design applications, and for building datasets for quantitative models in biology more generally.

Acknowledgements

We would like to thank Sarah Kate Nyquist for helpful conversations and suggestions. BLT would like to acknowledge support from the National Science Foundation Graduate Research Program. LC is supported by a David & Lucile Packard Fellowship for Science and Engineering. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of any of the funders.

References

- [1] Rhys M Adams, Thierry Mora, Aleksandra M Walczak, and Justin B Kinney. Measuring the sequence-affinity landscape of antibodies with massively parallel titration curves. *Elife*, 5: e23156, 2016.
- [2] Surojit Biswas, Gleb Kuznetsov, Pierce J Ogden, Nicholas J Conway, Ryan P Adams, and George M Church. Toward machine-guided design of proteins. *bioRxiv*, page 337154, 2018.
- [3] Eric T Boder and K Dane Wittrup. Yeast surface display for screening combinatorial polypeptide libraries. *Nature Biotechnology*, 15(6):553–557, 1997.
- [4] Longxing Cao, Inna Goreshnik, Brian Coventry, James Brett Case, Lauren Miller, Lisa Kozodoy, Rita E Chen, Lauren

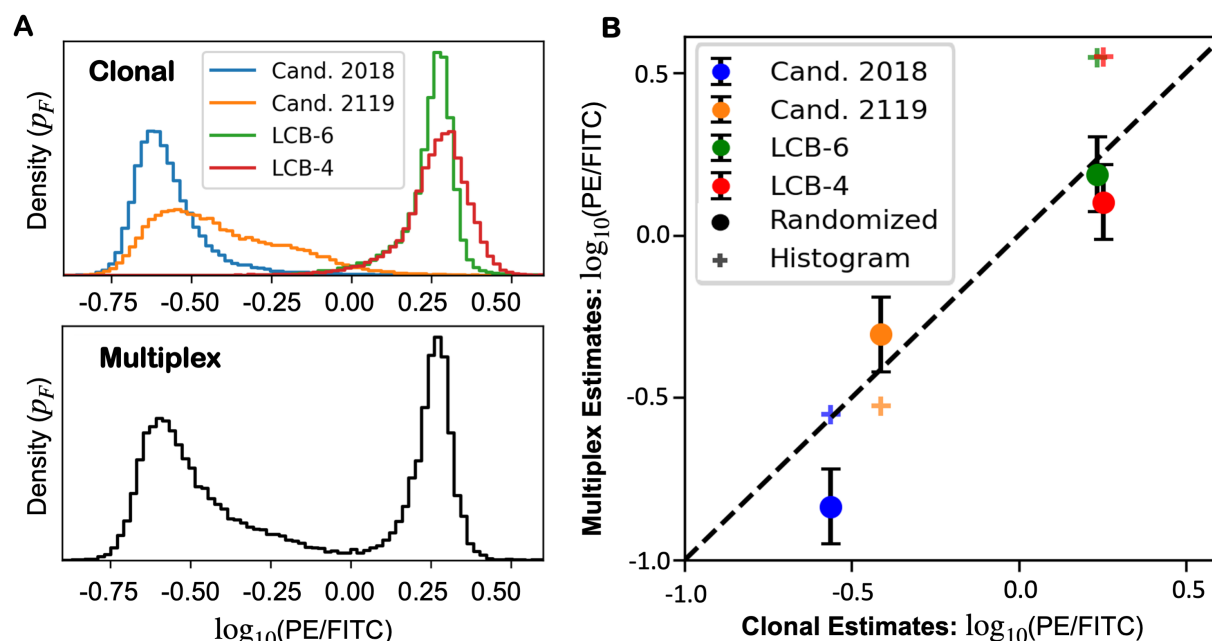


Figure 3. Agreement of binding signal of *de novo* designed binding proteins measured via yeast display in multiplex with ground truth values obtained in clonal yeast. (A) Distributions of samples measured clonally by flow cytometry, and distributions of pooled samples during sorting with a shifting gate boundary. (B) Agreement of clonal and multiplexed binding signal. The x-axis is measured by flow cytometry while the y-axis is a multiplexed measurement by next-generation sequencing. Error bars represent size of the steps used when shifting the threshold.

- Carter, Alexandra C Walls, Young-Jun Park, et al. De novo design of picomolar SARS-CoV-2 miniprotein inhibitors. *Science*, 370(6515):426–431, 2020.
- [5] Longxing Cao, Brian Coventry, Inna Goresnik, Buwei Huang, Joon Sung Park, Kevin M Jude, Iva Marković, Rameshwar U Kadam, Koen HG Verschueren, Kenneth Verstraete, et al. Robust de novo design of protein binding proteins from target structural information alone. *bioRxiv*, 2021.
- [6] Sebastian M Castillo-Hair, John T Sexton, Brian P Landry, Evan J Olson, Oleg A Igoshin, and Jeffrey J Tabor. Flowcal: a user-friendly, open source software tool for automatically converting flow cytometry data from arbitrary to calibrated units. *ACS synthetic biology*, 5(7):774–780, 2016.
- [7] Aaron Chevalier, Daniel-Adriano Silva, Gabriel J Rocklin, Derrick R Hicks, Renan Vergara, Patience Murapa, Steffen M Bernard, Lu Zhang, Kwok-Ho Lam, Guorui Yao, et al. Massively parallel de novo protein design for targeted therapeutics. *Nature*, 550(7674):74–79, 2017.
- [8] Carl G de Boer, John P Ray, Nir Hacohen, and Aviv Regev. MAUDE: inferring expression changes in sorting-based CRISPR screens. *Genome Biology*, 21:1–16, 2020.
- [9] Charles P Fulco, Joseph Nasser, Thouis R Jones, Glen Munson, Drew T Bergman, Vidya Subramanian, Sharon R Grossman, Rockwell Anyoha, Benjamin R Doughty, Tejal A Patwardhan, et al. Activity-by-contact model of enhancer–promoter regulation from thousands of crispr perturbations. *Nature Genetics*, 51(12):1664–1669, 2019.
- [10] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. The elements of statistical learning – data mining, inference, and prediction, 2001.
- [11] Justin B Kinney, Anand Murugan, Curtis G Callan, and Edward C Cox. Using deep sequencing to characterize the biophysical mechanism of a transcriptional regulatory sequence. *Proceedings of the National Academy of Sciences*, 107(20):9158–9163, 2010.
- [12] David G Luenberger. *Optimization by vector space methods*. John Wiley & Sons, 1969.
- [13] Neil Peterman and Erel Levine. Sort-seq under the hood: implications of design choices on large-scale characterization of sequence-function relations. *BMC Genomics*, 17(1):1–17, 2016.
- [14] Lothar Reich, Sanjib Dutta, and Amy E Keating. SORTCERY—a high-throughput method to affinity rank peptide ligands. *Journal of Molecular Biology*, 427(11):2135–2150, 2015.
- [15] Gabriel J Rocklin, Tamuka M Chidyausiku, Inna Goresnik, Alex Ford, Scott Houliston, Alexander Lemak, Lauren Carter, Rashmi Ravichandran, Vikram K Mulligan, Aaron Chevalier, et al. Global analysis of protein folding using massively parallel design, synthesis, and testing. *Science*, 357(6347):168–175, 2017.
- [16] Eilon Sharon, Yael Kalma, Ayala Sharp, Tali Raveh-Sadka, Michal Levo, Danny Zeevi, Leeat Keren, Zohar Yakhini, Adina Weinberger, and Eran Segal. Inferring gene regulatory logic from high-throughput measurements of thousands of systematically designed promoters. *Nature Biotechnology*, 30(6):521–530, 2012.
- [17] Tyler N Starr, Allison J Greaney, Sarah K Hilton, Daniel Ellis, Katharine HD Crawford, Adam S Dingens, Mary Jane Navarro, John E Bowen, M Alejandra Tortorici, Alexandra C Walls, et al.

Deep mutational scanning of SARS-CoV-2 receptor binding domain reveals constraints on folding and ACE2 binding. *Cell*, 182(5):1295–1310, 2020.

Methods

Unbiasedness of estimates from randomized gates. The advantage of the randomized gates presented in Equation (1) is that the resulting counts in each bin (Y_1 and Y_2) may be combined as $\hat{\mu} = (U - L) \cdot Y_2 / (Y_1 + Y_2) + L$ to estimate μ_F without bias. We make this statement precise and present a theorem that guarantees when this is the case.

For an estimator $\hat{\theta}$ of a fixed estimand θ , the estimator’s *bias* is the expected value of its error $\mathbb{E}[\hat{\theta} - \theta | \theta]$ conditioned on that particular value of θ . An estimator is called *unbiased* if, regardless of the value of the estimand, the bias is equal to zero—that is, if $\mathbb{E}[\hat{\theta} | \theta] = \theta$ for every θ . Theorem 1 states that this property holds for $\hat{\mu}$.

Theorem 1 (Unbiasedness with randomized gates). *If the support of p_F is bounded between L and U , then $\hat{\mu}$ is an unbiased estimator of mean fluorescence. That is, $\mathbb{E}[\hat{\mu}] = \mu_F$.*

Proof. We begin by rewriting the probability that a cell is collected into bin 2 to expose the connection between this quantity and μ_F :

$$\begin{aligned} \mathbb{P}(B=2) &= \int_L^U p_F(f) \mathbb{P}(B=2 | F=f) df \\ // \text{ via law of total probability \& support assumption} \\ &= \int_L^U p_F(f) (f-L)/(U-L) df \\ // \text{ by Equation (1)} \\ &= (\mu_F - L)/(U-L), \end{aligned}$$

If $N = Y_1 + Y_2$ total cells are collected, then the count in the second bin is distributed as $Y_2 | N \sim \text{Binomial}((\mu_F - L)/(U-L), N)$ and has mean $\mathbb{E}[Y_2] = N \cdot (\mu_F - L)/(U-L)$. Accordingly, for any N total number of cells, $\mathbb{E}[\hat{\mu} | Y_1 + Y_2 = N] = N \cdot (\mu_F - L)/(Y_1 + Y_2) + L = \mu_F$. When N is random as well, then by the law of iterated expectation, $\mathbb{E}[\hat{\mu}] = \mathbb{E}[\mathbb{E}[\hat{\mu} | Y_1 + Y_2 = N]] = \mu_F$ as desired. \square

Notably, this theorem holds for any distribution p_F satisfying the support condition and does not require any parametric assumptions such as log-normality.

Trade-off between number of cells sorted and precision of estimates. The relative simplicity of the estimate $\hat{\mu}$ leads to a transparent trade-off between the precision and scale of the experiment. Recalling that $Y_2 | N \sim \text{Binomial}((\mu_F - L)/(U-L), N)$, the variance of $\hat{\mu}$ is

$$\text{Var}[\hat{\mu} | N] = \frac{(U-L)^2}{N^2} \text{Var}[Y_2] = \frac{(U-L)^2}{N} \mathbb{P}(B=1) \mathbb{P}(B=2).$$

To construct a confidence interval for μ_F , we can therefore first approximate the standard error of $\hat{\mu}$ by $\frac{U-L}{\sqrt{N}} \frac{\sqrt{Y_1 Y_2}}{N}$, and appeal to approximate normality of the Binomial distribution for moderate to large N to report $\mu_F = \hat{\mu} \pm 2 \frac{U-L}{\sqrt{N}} \frac{\sqrt{Y_1 Y_2}}{N}$ with 95% confidence. Because $\sqrt{Y_1 Y_2}/N$ can be at most $1/2$ (if $Y_1 = Y_2$), the size of this interval is at most $2(U-L)/\sqrt{N}$. Therefore, to estimate μ_F to within one tenth of the range with high confidence, at most $N = 400$ cells are needed, since in this case $2(U-L)/\sqrt{N} = (U-L)/10$.

For scale, commercial machines sort on the order of ten thousand cells per second, and typical assays sort tens of millions of cells

divided amongst many populations. Thus, a library of one hundred thousand populations could be screened to high precision with on the order of 1 hour of sorting time.

Simulation details. In the simulations depicted in Figure 2, we compare against the standard approach of using a histogram to estimate μ_F . Consider a K bin histogram. For each bin k , if the range of fluorescences collected is from lower bound l_k to upper bound u_k , then $\mathbb{P}(B=k | F=f) = \mathbf{1}[l_k \leq f < u_k]$. The histogram estimate then corresponds to combining the resulting counts as

$$\hat{\mu}_{\text{Hist}} = \sum_{k=1}^K \frac{Y_k}{N} \left(\frac{u_k + l_k}{2} \right).$$

In order to use the unbiased estimator, both in simulation and in practice, we must slightly extend the randomized gate definition proposed in Equation (1). In particular, Theorem 1 assumes that the support of the fluorescence density p_F is bounded between L and U (i.e., that for $F \sim p_F$, $\mathbb{P}[L \leq F \leq U] = 1$). In practice, this may not be the case. But, as previously stated, Equation (1) returns negative “probabilities” outside of this range. Therefore, we propose to “clip” the collection probabilities at the boundaries, and instead define

$$\begin{aligned} \mathbb{P}(B=1 | F=f) &= \left(1 - \frac{f-L}{U-L} \right)_{\dagger} \quad \text{and} \\ \mathbb{P}(B=2 | F=f) &= \left(\frac{f-L}{U-L} \right)_{\dagger} \end{aligned}$$

where \dagger denotes clipping between zero and one such that, for a scalar x , $(x)_{\dagger} = \max(\min(x, 1), 0)$. This ensures that $\hat{\mu}$ is well-defined, but gives up unbiasedness in situations where the support assumption of Theorem 1 is violated. This bias is apparent, for example, at the right and left sides of the left panel of Figure 2C.

Experimental approximation of randomized gates with shifting thresholds. Because current FACS software does not support randomized gate programming, we devised an experimental approximation in which we manually changed the gating threshold 20 times during sorting at regular intervals. Specifically, we use a gate that collects all cells with fluorescence above a threshold into bin 2 and those below the threshold into bin 1, and we shift that threshold over the course of the collection from the lower limit L to the upper limit U . In theory, this approach exactly recovers Equation (1) in the limit that the threshold is shifted continuously from L to U at a constant rate. This is because for a cell with fluorescence f between L and U , the probability that it is collected into bin 2 is the fraction of the experimental time during which the threshold is below f , which is $(f-L)/(U-L)$. This approximation does not, however, account for possible changes in the distribution, p_F over time. Such changes occur in binding assays, for example, when nontrivial labeled target protein dissociates over time. This challenge is a disadvantage of the approximation relative to randomized gates that could in theory be implemented into sorters.

Yeast display and deep sequencing. EBY100 yeast cells expressing each of the four mini-protein binders were grown in C-Trp-Ura media. Binder protein expression was induced by replacing the growing buffer with SGCAA and incubating at 30° C for 24h⁷. The induced cells were labelled with 250 nM biotinylated receptor binding domain target protein, washed twice with PBSF (PBS+1% BSA), then labelled again with anti-c-Myc fluorescein isothiocyanate (FITC) and streptavidin-phycoerythrin (SAPE). The experiments were performed on a Sony SH800 cell sorter.

60,000 cells were recorded for each binder to reflect the individual distribution of baseline PE signal intensity. In the shifting gate experiment, a square area (**AreaTotal**) with side length (L) was pre-determined at the SH800 collection panel. The area was divided into 2 separate collection gates, **Gate1** and **Gate2** (corresponding to bin 1 and bin 2 in Equation (1)). **Gate2** was in an isosceles right triangle and started with a small area in the right-bottom corner of **AreaTotal** and **Gate1** took up the remaining. The yeast cells were run through the SH800 and each cell went into either the **Gate2** or **Gate1** collection tube if its log PE/FITC signal was in the range of **AreaTotal**. All other cells were discarded. After collecting 50,000 cells, the cell flow was paused, **Gate2** was shifted both leftwards and upwards for $L/10$ and cell flow continued. Because the proprietary software for operating the sorter allowed setting gate positions only through a point and click graphical user interface (rather than numerically), we measured out gate increments by pixel distance on the display using a ruler. The above shifting process repeated 19 times for a total of 20 collections. The cells collected in **Gate1** and **Gate2** were then grown, and 1×10^7 cells from each gate were barcoded and the sequences for each cell were determined by Illumina next-generation sequencing¹⁵. The number of cells collected by each gate for each population was estimated from the proportion of sequencing reads attributed to each population and the number of cells collected into the gates.

Because the number of cells collected by each gate was not made directly available through the proprietary software, we estimated this from the raw exported data. In particular, we imported the data using the **FlowCal** python package⁶ and computationally implemented the gates and filters (including for forward and backward scatter).

Sensitivity of maximum likelihood inference to non-normality of real data. Likelihood-based inference is a common strategy used with the intent to circumvent the resolution limitation of the histogram approach^{5;8;9;15}. However, this approach can fail on real data. In particular, existing likelihood methods rely on the assumption that for each of the cell populations the fluorescence values are log normal distributed, $\log F \sim \mathcal{N}(\mu, \sigma^2)$ where the mean log fluorescence $\mu = \mu_F$ is the target of inference and σ^2 is the typically unknown variance of the population.

We evaluate performance of maximum likelihood inference in this situation with simulations using data sub-sampled from a flow cytometry dataset of binding signal of a computationally designed mini-protein binder to ActRII. Data were collected using yeast display as previously described except with the addition of a supplemental binding protein, protein A, the binding signal log(FITC/PE) was recorded for approximately one million cells. The distribution of this signal is highly non-Gaussian (Figure S1A).

We first compared the performance of the maximum likelihood approach (described in greater detail below) to the randomized approach on downsampled datasets with $N = 250$ cells with the same set-up described in Figure 2. As in the earlier simulations, the randomized approach provides improved MSE across most simulation conditions (Figure S1B). This improvement is again explained by estimation bias, which is mitigated by the randomized approach (Figure S1C). Though one might expect the benefit of maximum likelihood would appear for larger sample sizes (e.g., due to the asymptotic efficiency of maximum likelihood estimation in theory), this is not the case. In fact, due to the bias of maximum likelihood, the relative improvement of the randomized approach is larger at $N = 1000$ cells (Figure S1D). Moreover, Figure S1E demonstrates that the maximum likelihood approach does not empirically provide more accurate estimates even under correct specification (with fluorescences sampled as in Figure 2A).

Maximum likelihood estimation. To estimate μ_F , likelihood-based approaches consider the counts in each of K bins (Y_1, Y_2, \dots, Y_K), since the measured fluorescence values cannot be disambiguated when multiple populations are sorted in multiplex. These counts follow a multinomial distribution as

$$Y_1, Y_2, \dots, Y_K \sim \text{Mult}(\pi(\mu, \sigma^2), N),$$

where $N = \sum_{k=1}^K Y_k$ is the total number of cells sorted into any bin and $\pi(\mu, \sigma^2) = (\pi_1, \pi_2, \dots, \pi_K)$ are the normalized bin probabilities. In particular, if for each bin k the range of fluorescences collected is from lower bound l_k to upper bound u_k , then

$$\pi_k = \frac{\Phi(\frac{u_k - \mu}{\sigma}) - \Phi(\frac{l_k - \mu}{\sigma})}{\sum_{k'=1}^K \Phi(\frac{u_{k'} - \mu}{\sigma}) - \Phi(\frac{l_{k'} - \mu}{\sigma})},$$

where $\Phi(\cdot)$ is the cumulative density function of the standard normal. The log likelihood function is then

$$\log p(Y_1, \dots, Y_K; \mu, \sigma^2) = \log N! - \sum_{k=1}^K \log Y_k! + \sum_{k=1}^K Y_k \log \pi_k,$$

where the dependence of each π_k on μ and σ^2 is left implicit. The maximum likelihood approach is to return μ that maximizes this expression,

$$\hat{\mu}_{\text{MLE}} = \arg \max_{\mu} \left[\max_{\sigma^2 > 0} \log p(Y_1, \dots, Y_K; \mu, \sigma^2) \right].$$

This optimization problem is not analytically tractable, and its constraints and non-convexity pose challenges for local, gradient-based optimizers. So we instead solve the optimization approximately with a grid search.