

1 **Computational prediction of protein interactions on single cells by**
2 **proximity sequencing**

3
4 Junjie Xia^{1†}, Hoang Van Phan^{1,2†}, Luke Vistain^{1,3†}, Mengjie Chen^{4,5}, Aly A. Khan⁶, Savaş
5 Tay^{1*}

6
7 ¹Pritzker School of Molecular Engineering, The University of Chicago, Chicago, IL, 60637,
8 USA

9 ²Present address: Division of Infectious Disease, University of California, San Francisco, CA,
10 94143, USA

11 ³Present address: Lymphocyte Biology Section, Laboratory of Immune Systems Biology,
12 NIAID, NIH, Bethesda, MD, 20892, USA

13 ⁴Section of Genetic Medicine, Department of Medicine, The University of Chicago, Chicago, IL,
14 60637, USA

15 ⁵Department Human Genetics, The University of Chicago, Chicago, IL, 60637, USA

16 ⁶Department of Pathology, The University of Chicago, Chicago, IL, 60637, USA

17 †These authors contributed equally

18 *Correspondence: tays@uchicago.edu

19

20 **Abstract**

21 Proximity sequencing (Prox-seq) measures gene expression, protein expression, and protein
22 complexes at the single cell level, using information from dual-antibody binding events and a
23 single cell sequencing readout. Prox-seq provides multi-dimensional phenotyping of single cells
24 and was recently used to track the formation of receptor complexes during inflammatory
25 signaling in macrophages and to discover a new interaction between CD9/CD8 proteins on naïve
26 T cells. The distribution of protein abundance affects identification of protein complexes in a
27 complicated manner in dual-binding assays like Prox-seq. These effects are difficult to explore
28 with experiments, yet important for accurate quantification of protein complexes. Here, we
29 introduce a physical model for protein dimer formation on single cells and computationally
30 evaluate several different methods for reducing background noise when quantifying protein
31 complexes. Furthermore, we developed an improved method for analysis of Prox-seq single-cell
32 data, which resulted in more accurate and robust quantification of protein complexes. Finally,
33 our model offers a simple way to investigate the behavior of Prox-seq under various biological
34 conditions and guide users toward selecting the best analysis method for their data.

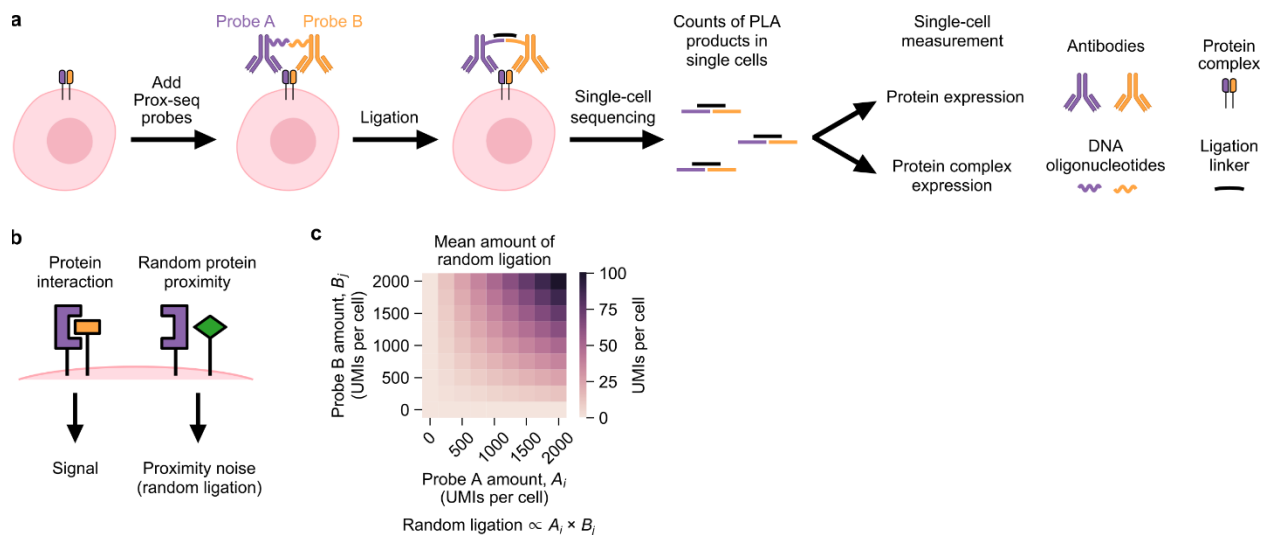
35

36 **Introduction**

37 Advances in single cell sequencing have enabled unprecedented analyses of cellular
38 heterogeneity in complex biological systems^{1,2}. Single-cell RNA sequencing³ (scRNA-seq) is
39 among the most widely used methods. However, because proteins are the effector molecules for
40 the majority of biological functions, RNA data alone is not sufficient to investigate these protein
41 functions thoroughly. Signaling events, for example, typically begin with receptor clustering,

42 protein phosphorylation, and other protein-protein interactions, all of which occur prior to
43 transcription.

44 To investigate the roles of protein interactions in greater depth, we recently developed a
45 method called proximity sequencing (Prox-seq) for simultaneous quantification of mRNA,
46 surface proteins and protein complexes at the single-cell level⁴. Prox-seq captures protein
47 complex information in barcoded DNA oligonucleotides (oligos) using a proximity ligation
48 assay^{5,6} (PLA). Each protein in Prox-seq is targeted by two DNA-conjugated antibodies, called
49 Prox-seq probes A and B (Figure 1a). The DNA oligos on probes A and B are ligated only if two
50 protein molecules are sufficiently close to each other. The result of this ligation is referred to as a
51 “PLA product.” The ligation distance is expected to be 50-70nm⁷. In order to generate a PLA
52 product, the oligo belonging to a Prox-seq probe A must ligate to the oligo belonging to a Prox-
53 seq probe B. Importantly, unligated probes do not contribute to the signal because both library
54 preparation and sequence alignment require barcodes from both the A and B probe. Upon
55 sequencing, the number of PLA products can be determined by counting the number of unique
56 molecular identifiers (UMIs). Because of this design, the number of PLA products measured for
57 a protein is a reflection of both the abundance of that protein and the availability of nearby Prox-
58 seq probes. By combining Prox-seq with scRNA-seq, these PLA products can be sequenced
59 alongside complementary DNA (cDNA) libraries, providing information on gene expression,
60 protein abundances, and protein complex formation from single cells⁴.



61

62 **Figure 1. Working principle of Prox-seq and identification of proximity noise.** (a) Schematic
 63 showing the main steps of Prox-seq. (b) Schematic showing the background in Prox-seq that is
 64 caused by proximity noise (random ligation of non-interacting protein molecules). (c) Heatmap
 65 showing the expected amount of proximity noise created from simulations of two protein
 66 molecules at varying expression levels. By modeling the mean amount of proximity noise with a
 67 binomial distribution (see Methods), we found that it was proportional to the product of the
 68 abundances of the two protein molecules.

69

70 Prox-seq protein data contains a unique source of background noise, namely the ligation
 71 of two protein molecules that do not functionally interact but are nevertheless sufficiently close
 72 to each other by random chance. We call this effect “proximity noise” (Figure 1b). Proximity
 73 noise exists because the average distance between probes on the cell surface decreases with
 74 increasing protein abundance (see Methods). A previous study showed that proximity noise led
 75 to false positive detection of protein interactions for in situ PLA⁸. A theoretical model showed
 76 that the mean amount of proximity noise is proportional to the product of the expression levels of
 77 the two proteins that made up the PLA product (Figure 1c). In short, the presence of PLA
 78 products for a specific pair of proteins does not guarantee that the two proteins functionally
 79 interact and form stable complexes.

80 To account proximity noise, we previously proposed and used a statistical method,
81 termed the iterative method, to differentiate protein complexes from random ligation in PLA
82 product counts⁴. Initially, this method establishes an "expected value" for each PLA product,
83 representing the number of PLA products that would exist if Prox-seq probes were randomly
84 distributed across the cell surface. Subsequently, the method subtracts the expected background
85 from the PLA product counts. If a PLA product's count exceeds its expected value, the difference
86 between observed and expected PLA products is attributed to non-random protein complexes.
87 This procedure is iteratively executed for every type of PLA product in each individual cell.
88 Although this method successfully recovered positive controls of known protein complexes,
89 assessing its performance on experimental data is challenging, as the generation of new Prox-seq
90 datasets often lacks comprehensive knowledge of the entire set of protein complexes and their
91 expression levels.

92 In this study, we present a simulation model for single-cell proteomic data in proximity
93 sequencing experiments and use it to computationally benchmark the performance of several
94 new and existing protein complex prediction methods. After calibrating the model with
95 experimental data, the simulation model allowed us to quantitatively analyze proximity noise and
96 its effects on the measured PLA product counts. We compare the performance of three methods:
97 the iterative method, a new linear regression-based method, and a new ensemble method that
98 combines the two. We find that, while the iterative and linear regression-based methods perform
99 well in several different scenarios, combining them into a single method yielded the most
100 accurate and robust quantification of protein complexes. These results shed insight onto how the
101 spatial organization of surface proteins translate into Prox-seq data and provides guidelines for
102 use of Prox-seq and related dual-binding technologies for multi-omic analysis of single cells.

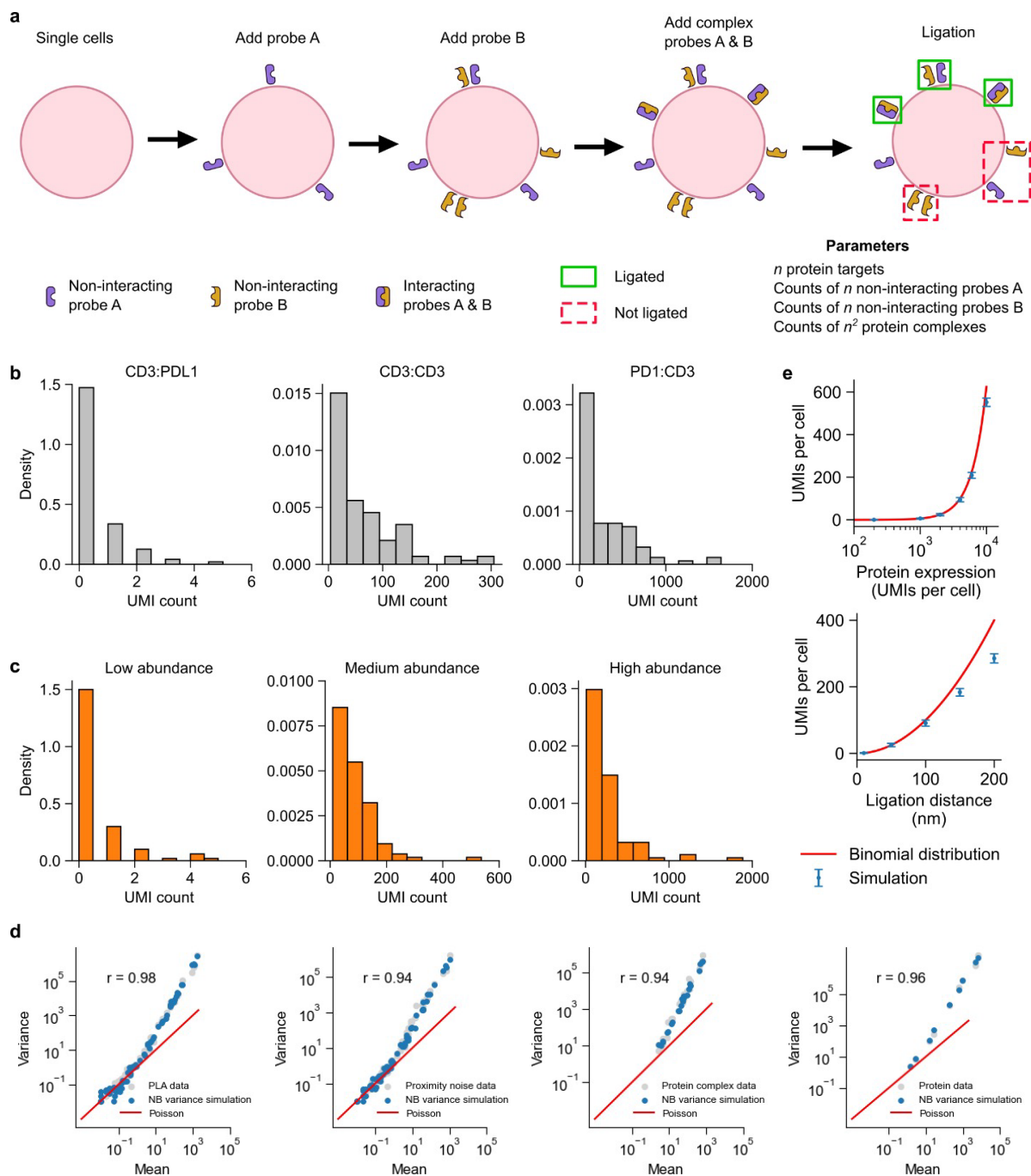
103 **Results**

104 **Overview of the simulation model**

105 Based on a physical model of how PLA products are formed in each single cell, we
106 created a simulation model of PLA product count data. We reasoned that proximity alone would
107 determine if a Prox-seq probe A and a Prox-seq probe B ligate and produce a PLA product. We
108 constructed the simulation model in a way that allowed us to simulate probes that bind to non-
109 interacting protein molecules (proteins that are not part of a complex) separately from probes that
110 bind to interacting molecules (proteins that are part of a complex). This procedure enabled us to
111 independently tune the abundance of proteins and protein complexes in the simulation, and to
112 observe how these properties affected Prox-seq data.

113 First, we generated the non-interacting Prox-seq probes A as random points on a sphere
114 (Figure 2a). These points indicate that the protein molecules exist as monomers; thus, any PLA
115 products they form would be caused by proximity noise and as a result of being in a protein
116 complex. Further, we assumed such protein monomers were distributed randomly on the cell
117 surface. Then, we repeated the process to generate the non-interacting Prox-seq probe B signal.
118 Second, we generated the interacting Prox-seq probes A and B by generating a sphere of random
119 points. These points corresponded to detectable protein complexes. Because these two probes A
120 and B both bound to the same protein complexes, the Prox-seq probe A points would necessarily
121 be in proximity with their corresponding Prox-seq probe B points. Finally, any pairs of probe A
122 and B with Euclidean distances less than the ligation distance were considered ligated and
123 produced PLA products (see Methods). If a probe A was within the ligation distance with more
124 than one probes B, one such probe B was chosen at random to ligate with said probe A.

125



126

127 **Figure 2. Overview and calibration of simulated Prox-seq data.** (a) Schematic for the
 128 simulation model of PLA products. The simulation was separately performed on a cell-by-cell
 129 basis. First, a number of non-interacting probes A and non-interacting probes B were added as
 130 random points on a sphere. Next, a number of protein complexes were added as random points on
 131 a sphere. These points corresponded to probes A and B that bound to interacting protein molecules.

132 Finally, probes A and B that had a Euclidean distance lower than the ligation distance were ligated,
133 thus creating PLA products. (b) Histograms showing the UMI counts of three example PLA
134 products in single Jurkat cells. (c) Histograms showing the UMI counts of three example simulated
135 PLA products with NB variance. (d) Scatter plots of mean-variance relationship show how
136 negative binomial variance captures overdispersion in PLA data, proximity noise data, protein
137 complex data, and protein data. (e) The relationship between proximity noise (measured as UMI's)
138 and protein abundance (top) or ligation distance (bottom).

139

140 We next compared the simulated Prox-seq data to experimental data. We analyzed T cells
141 (Jurkat cell line) and B cells (Raji cell line) with a panel of Prox-seq probes that targeted both T
142 cell and B cell markers from a previously reported study⁴. Simulated counts of PLA product and
143 protein expression followed the Poisson distribution, whereas the experimental data exhibited
144 overdispersion (Figure S1a, S2a). We found that adding variance in the form of a negative
145 binomial distribution (NB) for non-interacting probes and protein complexes was sufficient to
146 capture the overdispersion of the real data (NB variance, see Methods). With the added NB
147 variance, the simulated data, like the experimental data, had a right-skewed distribution across
148 different PLA product abundances (Figure 2b, c). Notably, the simulation model with added
149 variance captured the positive correlation between observed PLA product count and non-
150 proximal probe count in real data (Figure S1b-g). The simulation model with no variance,
151 however, showed a negative correlation between PLA product count and non-proximal probe
152 count (Figure S1d and S1e). The NB variance model also produced non-proximal probe counts
153 with similar distributions to those observed in experimental data (Figure S2).

154 We generated replicated datasets by sampling from the fitted model for posterior
155 predictive checks (PPCs)⁹. We then assessed how well these data samplings maintained the
156 properties of the observed data with two metrics. First, we measured the similarity between the
157 coefficient of variation per PLA product, proximity noise, protein complex and protein. This

158 comparison enables evaluation of how well the mean-variance relationship of real data is
159 preserved (Figure 2d & Figure S3a). Second, we perform Mann-Whitney U-test statistic to
160 measure the extent to which the replicated data and raw data come from the same distribution
161 (Figure S3b). Finally, we characterized the amount of proximity noise in the most basic scenario
162 when there were no protein complexes detectable by the Prox-seq probe panel. The simulation
163 demonstrated that the amount of PLA product produced by random ligation scales quadratically
164 with both protein abundance and ligation distance (Figure 2e). These results show that our model
165 and simulations faithfully capture key aspects of real Prox-seq data in single cells and reiterates
166 the importance of identifying and removing proximity noise, which can especially be large for
167 highly expressed proteins.

168

169 **Iterative prediction of protein complex abundance**

170 An iterative method was used to previously identify the existence of stable protein complexes in
171 Prox-seq measurements. This method proposed that when there were no protein complexes, the
172 observed count of a PLA product $i:j$ could be calculated from the abundance of the probe A
173 targeting protein i , and the probe B targeting protein j (see Methods). This calculation resulted in
174 an expected random count for PLA products that represents the PLA count caused by proximity
175 noise. We reasoned that if the observed count of PLA product $i:j$ was higher than the calculated
176 expected random count, then $i:j$ indicated a non-random protein interaction. To quantify the
177 protein complexes on each single cell, we calculated the difference between the observed and
178 expected random PLA product count (Figure 3a). This method was called the iterative method,
179 because it involved solving a system of quadratic equations (describing all possible protein
180 dimers) iteratively (see Methods)⁴. This method relied on the fact that Prox-seq can measure

181 protein abundance, similar to flow cytometry and CITE-seq¹⁰. The abundance of a protein was
182 the amount of protein molecules that were present on the cell surface, and therefore included
183 both molecules in monomeric and complex forms. In our previous study⁴, we proposed that the
184 protein abundance could be estimated from Prox-seq data by summing the appearances of each
185 protein across its associated PLA products (see Methods). Here, we find by using our simulated
186 data that such an estimate is a good approximation of the true protein abundance, as they are
187 strongly correlated (Figure S4).

188 To further examine the assumptions underlying the iterative method, we now constructed
189 the following simulation scenario: The simulation had three protein targets, called protein 1,
190 protein 2 and protein 3. These proteins did not interact with themselves, nor with any other
191 proteins. Furthermore, protein 3 had a lower non-interacting probe count (mean of 100 UMIs/cell
192 compared to 1000 UMIs/cell for proteins 1 and 2, Table S1). Simulated data showed that our
193 assumptions behind the iterative method were correct. When there were no interactions between
194 the proteins, the observed PLA product counts were similar to the expected random count
195 (Figure S5a). When we introduced the protein complex 1:1 to the simulation while keeping the
196 other parameters the same, the observed counts of the PLA product 1:1 was higher than its
197 expected random count (Figure S5b).

198 One weakness of the iterative method is complexity of hyper-parameter tuning, which
199 can result in sub-optimal convergence. The key parameter is the initialization setting, which are
200 the initial estimates of protein complex abundances. By default, the algorithm assigns an initial
201 value of 0 to all protein complexes. However, different initialization settings will influence
202 iterative behaviors to convergence, as well as tolerance (Figure S6a). Unsensible initialization

203 tends to generate nonsensical predictive outputs. (Figure S6b). This led us to consider more
204 robust methods for protein complex quantification.

205

206 **Prediction of protein complex abundance using linear regression**

207 To address the weakness of the iterative method we developed a new approach (the linear
208 regression - LR method). This method uses an experimentally modified Prox-seq procedure that
209 enables direct measurement of Prox-seq probes that were not ligated because they were not
210 proximal to another Prox-seq probe (we refer to these as non-proximal probes)⁴. The proximity
211 noise for a PLA product $i:j$ should be proportional to the product of the non-proximal probe A
212 targeting protein i , and the non-proximal probe B targeting protein j . We reasoned that if linear
213 regression is used to model the observed PLA product count onto the estimated random ligation
214 amount, true protein complexes would have positive intercepts (see Methods). The slope was
215 then used to estimate the amount of random ligation, and the count of a protein complex was
216 calculated by subtracting the estimated random ligation from the observed PLA product count
217 (Figure 3b). Experimentally, we observed strong heteroscedasticity in the PLA product count
218 when regressed on to the random ligation amount (Figure S7). Therefore, we performed linear
219 regression using weighted least squares instead of ordinary least squares (see Methods).

220 We created a new simulation to directly compare the iterative and LR methods. The simulation's
221 parameters were set to approximate the experimental data. More specifically, the simulation had
222 three protein targets: protein 1, protein 2 and protein 3. Proteins 1 and 2 interacted both with
223 themselves and each other (Figure 3c, Table S1). Protein 3 did not interact with itself, nor with
224 protein 1 or protein 2. Furthermore, protein 3 had very low non-interacting protein count (mean

225 of 2 UMIs/cell compared to 20 and 15 for proteins 1 and 2, respectively). We found that the
226 iterative method correctly identified protein complexes 1:1, 1:2, 2:1 and 2:2 (Figure 3c).

227 To determine if we can statistically infer the enrichment of PLA products, we performed
228 a one-sided Fisher's exact test on the counts of PLA products (Figure 3d, see Methods). This
229 analysis correctly identified the four protein complexes present in the sample, independently
230 confirming that the generated protein complexes occur at a higher frequency than random and
231 can be statistically inferred (Figure 3d, see Methods). With regards to quantification of protein
232 complexes on single cells, we observed that the iterative method consistently underestimated the
233 true protein complex count (Figure 3c, e). Conversely, the LR method not only correctly
234 identified the four true protein complexes (complexes 1:1, 1:2, 2:1 and 2:2), but also produced
235 much more accurate counts for them (Figure 3c, d, f). Overall, the results of the two methods
236 were correlated on the single-cell level (Figure 3g).

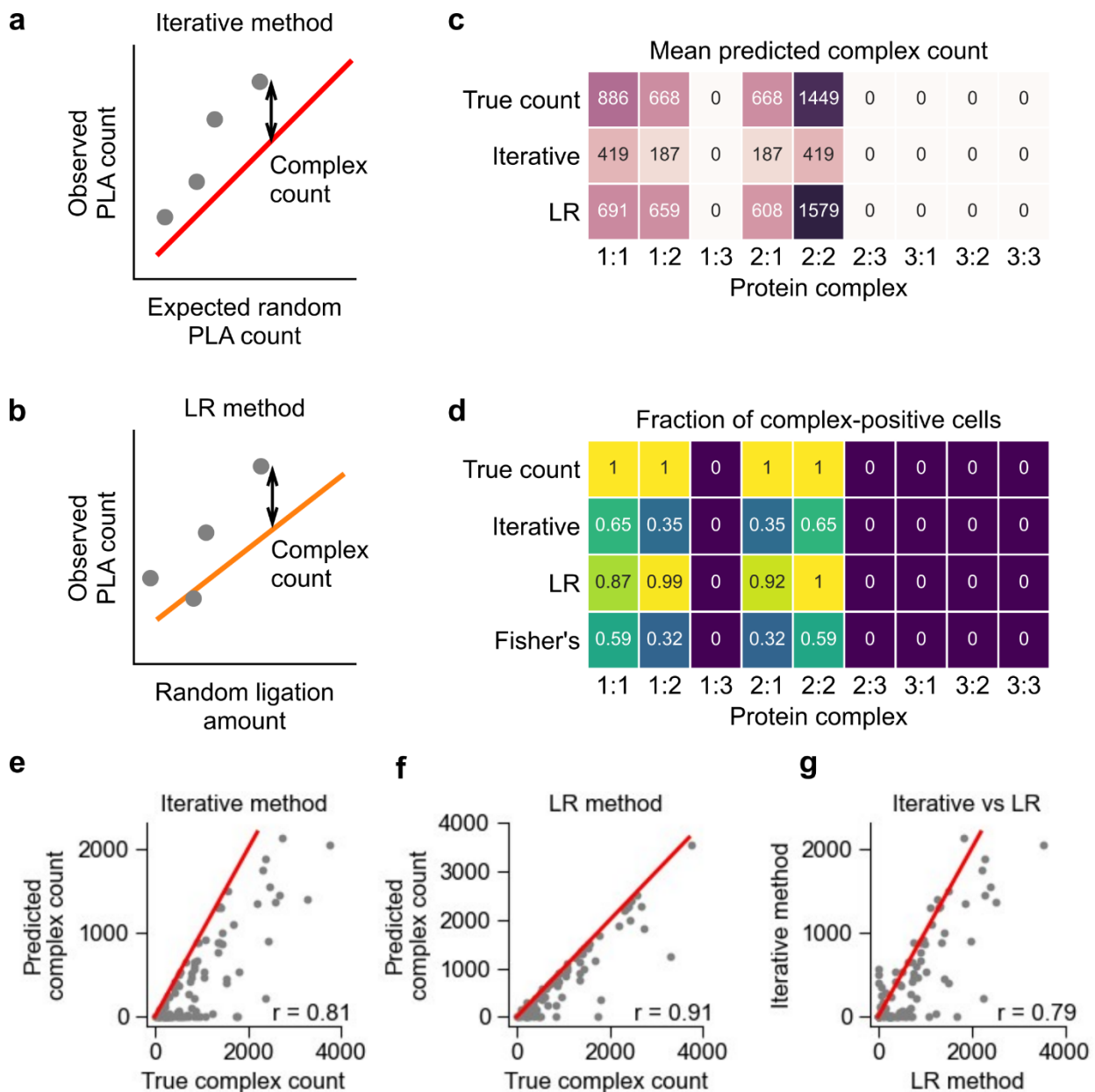
237

238 **Ensemble method that combines both the LR and iterative methods for analysis of Prox-** 239 **seq data**

240 We next chose to explore a method that had the potential to outperform both the LR and iterative
241 methods. As shown previously, the major weakness of the iterative method is its sensitivity to
242 initialization conditions. We reasoned that the output from the LR method could be used as a
243 sensible initialization for the iterative method (Figure 4a). Starting the iteration close to the
244 correct result would make it less likely that the method would fall into a spurious local
245 optimization. The performance of all three methods was compared in two simulations: one in

246 which a high percentage of proteins were in complex with other proteins (high signal) and one in
 247 which a low percentage of proteins were in complex (low signal) (Table S1).

248



249

250 **Figure 3. Comparison between the iterative and linear regression (LR) methods for protein**
 251 **complex prediction in simulated data.** (a, b) Schematics showing the working principle of (a)
 252 the iterative method and (b) the LR method. In the iterative method, the protein complex is
 253 the difference between the observed and expected PLA product count. In the LR method, the

254 protein complex count is the difference between the observed PLA product count and its expected
255 amount of random ligation, which is calculated from the non-proximal probe count. In (a), the red
256 line indicates $y = x$. In (b), the orange line indicates the linear regression fit. (c) Heatmap showing
257 the mean complex count of simulated data, and of the iterative and LR methods' prediction results.
258 (d) Heatmap showing the fraction of cells expressing a protein complex, as predicted by the
259 iterative method, the LR method, and Fisher's exact test. In (c, d), the true count represents the
260 ground truth of protein complex count in the simulation. (e, f) Scatter plots showing the simulated
261 and predicted count of protein complex 1:1 using (e) the iterative and (f) the LR method. (g) Scatter
262 plot comparing the predicted count of protein complex 1:1 from the iterative and the LR methods.
263 In (e-g), the red lines indicate $y = x$, and each dot represents a single cell.

264

265 The iterative method performed well when signal was high, but generated false positives
266 when signal was low (Figure 4b-c). The LR method performed better in the high noise
267 simulation but suffered from false positives when noise was low (Figure 4c). This is not
268 surprising because LR method depends on performing regression with product of non-proximal
269 probes as the explanatory variable and if there are few non-proximal probes across all single
270 cells (or low noise in our simulation), LR method will become unstable. Both methods
271 consistently underestimated the abundance of protein complexes. For the iterative method, this is
272 partly because expected PLA count we assumed is the maximal proximity noise it might have.
273 The slope we use to quantify protein complexes from LR method tends to be larger than the true
274 value because counts of non-proximal probes we can measure are inevitably lower than real
275 counts both in experiment and simulation, which would give us a smaller positive intercept and
276 protein complex count. In contrast, the ensemble method was able to maintain strong
277 performance in both scenarios. It was less likely to produce a false positive, assigned fewer reads
278 to false positives than other methods, and was closer to the true count for most of the protein
279 complexes (Figure 4b-c). Finally, for a given PLA product, the ensemble method was more
280 accurate in quantifying the abundance of true-positive complexes in single cells (Figure 4d).

281 **A quantitative scoring strategy to comprehensively evaluate prediction methods**

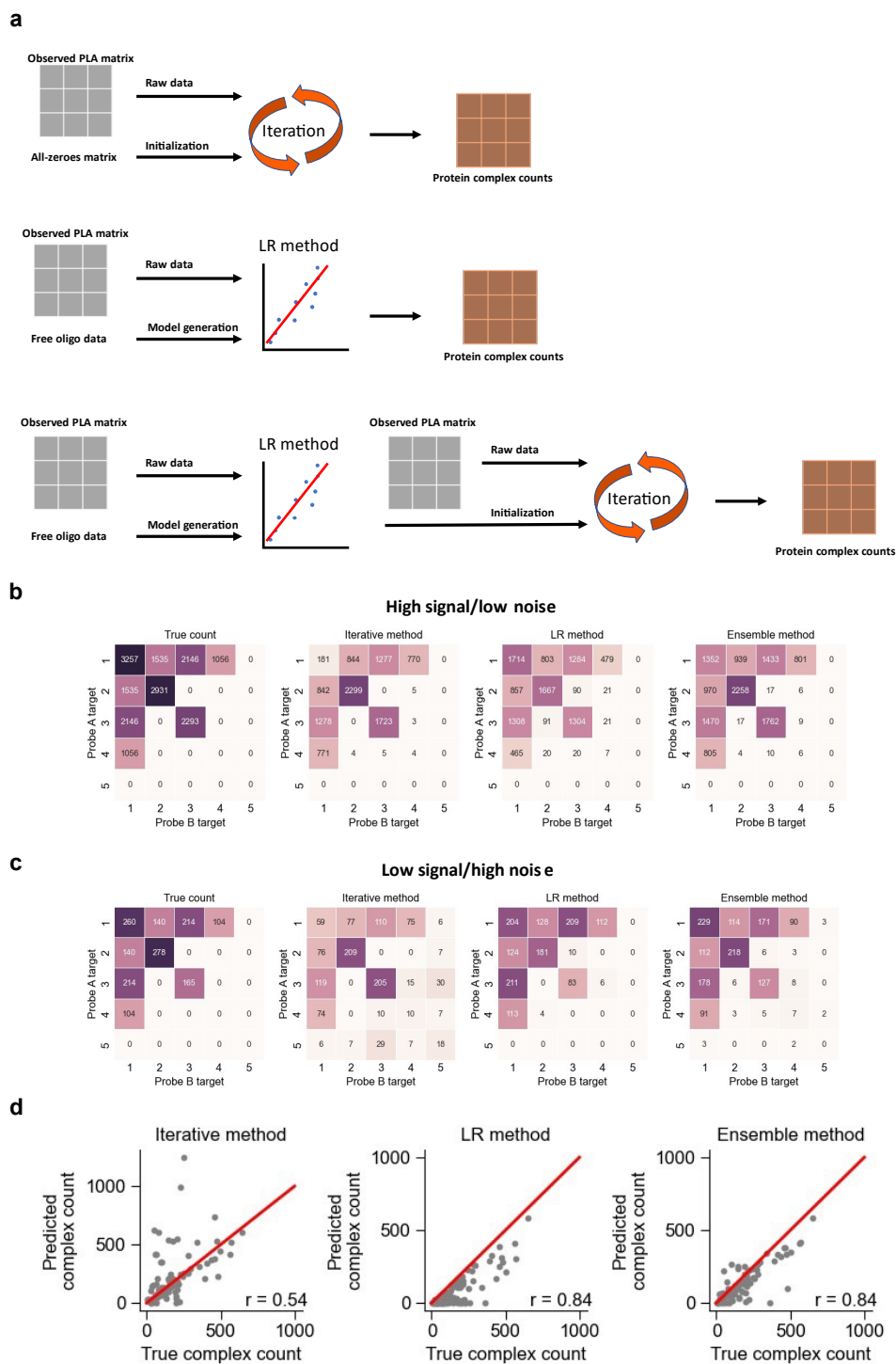
282 To evaluate the predictive performance of these methods more comprehensively, we further
283 propose a quantitative scoring strategy to assign a prediction score for every prediction (Figure
284 S8a). We simulate different biological scenarios with our model and score the overall prediction
285 performance of each method by considering sum of absolute deviation between mean true counts
286 and predicted counts ($\sum Mean_{deviation}$), sum of Pearson correlation coefficient ($\sum Pearson$)
287 across single cells (Figure S8b), and sum of ratios of false positive prediction ($\sum FPrate$)
288 across single cells (Figure S8c) (see Methods). Comparing the methods across all scenarios
289 showed that the ensemble method had the highest average prediction score and the lowest
290 variance (Figure S8d & Table S2). The ensemble approach effectively improves iterative method
291 and LR method's generalization to different biological scenarios.

292 **Comparison of all three analytical methods to real data and performance evaluations**

293 Next, we evaluated the concordance between all three methods on experimental data from single
294 Jurkat and Raji cells. Overall, we found that each method largely agreed on which PLA products
295 were predicted to be protein complexes (Figures 5a-f). While the bulk measurements of protein
296 complexes showed good agreement between methods, the three methods had varying levels of
297 correlation for single cells (Figure 5g, h). In addition, we observed all three method, along with
298 the Fisher's Exact test, largely identified the same protein complexes (Figure S9).

299 All methods predicted protein complexes CD3:CD3 and CD28:CD28 in Jurkat cells, both
300 of which are known protein complexes^{11,12}. All three methods also predicted protein complex
301 ICAM1:ICAM1 in Raji cells, which was shown to dimerize on the cell surface¹³. We also
302 evaluated our methods against a simulation designed to more closely represent the experimental

303 data. Protein expression levels were estimated from the experimental data and used to create
304 simulation models for Jurkat and Raji cells (Table S1). Then, protein complexes corresponding
305 to CD3:CD3, CD28:CD28, and CD3:CD28 were added to Jurkats, whereas HLADA:HLADR
306 and PDL1:PDL1 were added to Rajis. Once again, we observe largely similar performance for all
307 methods (Figure S10).



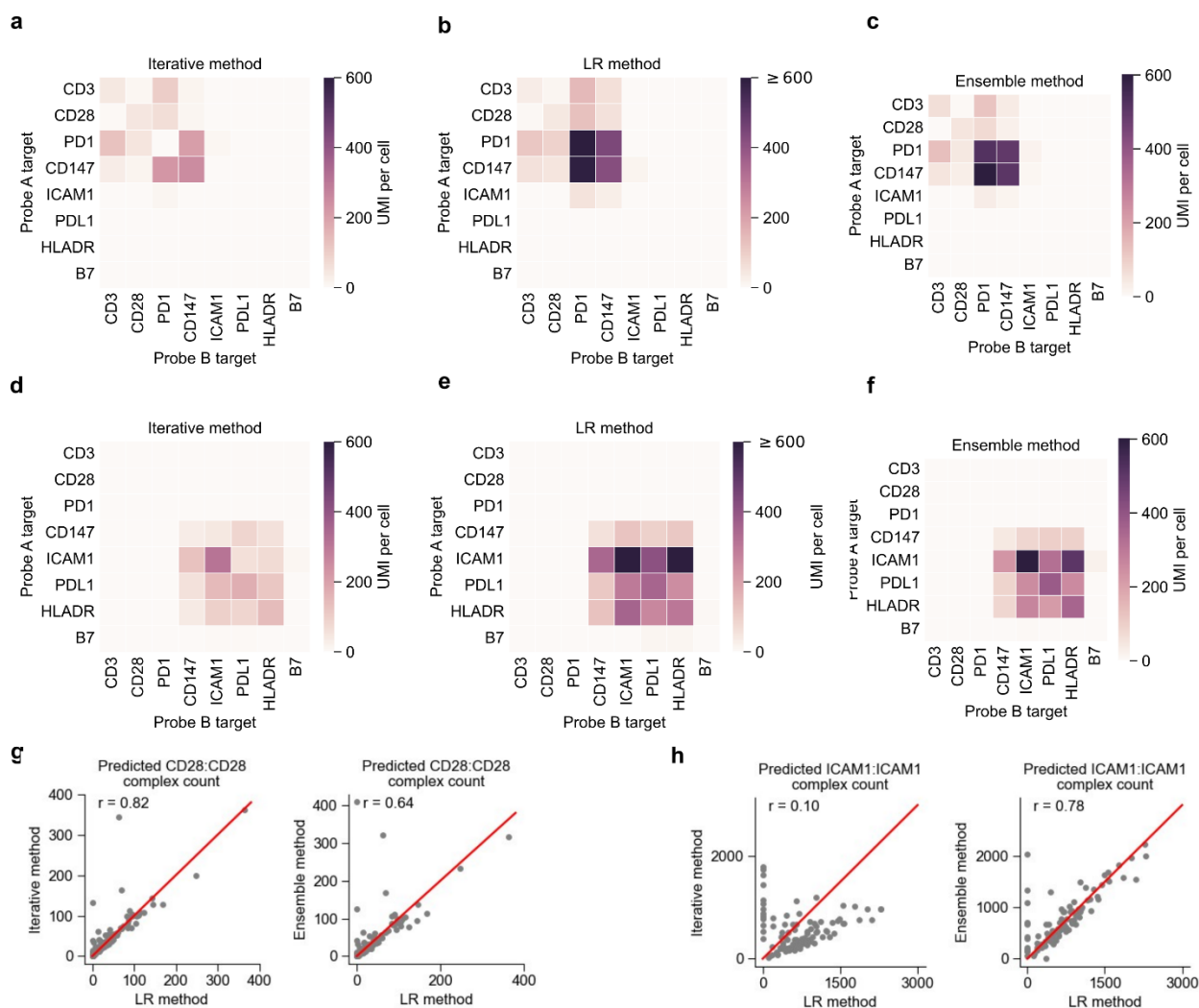
308

309 **Figure 4. The ensemble method for improved analysis of Prox-seq data.** We combine the
 310 iterative and LR methods for better prediction of protein complexes. (a) Schematic showing how
 311 all three methods arrive at protein complex estimation. The iterative method combines raw data
 312 and an initialization with an all-zeroes matrix to quantify protein complexes. The LR method uses
 313 raw data and free-oligo data to construct a linear regression model that quantifies protein

314 complexes. The ensemble method begins with applying the LR workflow and uses the output of it
 315 to initialize the iterative method. (b) Comparison of all three methods in a regime of high signal
 316 and low noise, compared to the true counts. (c) Comparison of all three methods in a regime of
 317 low signal and high noise, compared to the true counts. (d) The Pearson's correlation between true
 318 counts and the outputs for each method across single cells. Each example shows complex 3:3 from
 319 the low signal/high noise regime.

320

321



322

323 **Figure 5. Comparison between the iterative and LR methods on experimental data.** (a-c)
 324 Heatmaps showing the average of protein complex count, predicted by (a) the iterative method,
 325 (b) the LR method, and (c) the ensemble method in Jurkat cells. (d-f) Heatmaps showing the
 326 average of protein complex count, predicted by (a) the iterative method, (b) the LR method, and
 327 (c) the ensemble method in Raji cells. (g) Comparison of methods for predicting counts of protein
 328 complexes of CD28:CD28 and in Jurkat cells. (h) Comparison of methods for predicting counts of

329 protein complexes of ICAM1:ICAM1 and in Raji cells. In (g, h), the red lines indicate $y = x$, and
330 r indicates the Pearson's correlation coefficient.

331

332 **Discussion**

333 Here, we presented a comprehensive computational framework for simulating Prox-seq data, and
334 for predicting protein complex count from Prox-seq data. We studied how the quantification of
335 protein complexes was affected by proximity noise, which is caused by proteins that are not
336 functionally interacting but are sufficiently close to each other by random chance to produce
337 valid ligation products. Our simulation model showed that the amount of proximity noise is
338 strongly depended on the protein abundance. Similar results have been observed in commercial
339 *in situ* PLA⁸.

340 We showed that with respect to protein complex prediction, the iterative method, LR
341 method, and ensemble method largely agree on real experimental data. Therefore, we propose
342 that each of these methods could be used for protein complex detection and quantification, and
343 any protein complexes that were predicted by these methods were highly likely to be true protein
344 complexes. However, in head-to-head comparisons using simulated data, the ensemble method
345 performed well over a larger range of data types than the other methods.

346 Our simulation model had some limitations. First, it did not consider interactions higher
347 than dimers, diffusion of the protein molecules, their physical sizes, and the technical variability
348 of the Prox-seq assay. Second, the simulation model requires the user to independently select the
349 abundance of a protein complex and its constituents' non-interacting counterpart. In real cells,
350 these abundances are likely highly correlated. Finally, it assumed that the protein complexes and
351 the non-interacting proteins were uniformly distributed on the cell surface. Despite these

352 limitations, we showed that the overall structure of simulated Prox-seq data is very similar to real
353 Prox-seq data.

354 Currently, application of each method requires a relatively homogeneous population of
355 single cells. In practice, this requires that simultaneously acquired mRNA data is first used to
356 cluster cell types, and then either method can be applied to individual clusters. This requirement
357 is because each method relied on a statistic of the whole population (the difference between
358 observed and expected random PLA product count for the iterative method, and the linear
359 regression's intercept and slope coefficient for the LR method) and having different complex
360 expression levels would lower the power the methods. Further study is required to extend these
361 methods to a population of heterogeneous cell types without the use of mRNA data.

362 We envision that the Ensemble method will be particularly useful when Prox-seq is
363 extended to intracellular proteins. Indeed, since non-specific antibody binding is much more
364 severe in intracellular staining than extracellular staining, random ligation is an even more
365 important source of noise given common macromolecular crowding effect within cells. The
366 simulation model can also be further extended to model Prox-seq data of intracellular proteins. In
367 short, we have validated the protein complex prediction algorithm that was proposed previously⁴,
368 proposed two additional independent methods for protein complex prediction, and introduced a
369 model for simulating Prox-seq data.

370

371 **Methods**

372 **Theoretical calculation of proximity noise**

373 Suppose there are A_i probes A and B_j probes B on the cell surface. Assume that the probes are
374 random points on a spherical surface, and proteins i and j do not interact. Because the ligation
375 distance is significantly shorter than the cell's radius, we assume that a probe A and a probe B
376 can be ligated if and only if the Euclidean distance between them, L , is less than or equal to the
377 ligation distance, $d_{ligation}$. The Euclidean distance L between any pair of random points has the
378 following probability distribution¹⁴:

$$379 \quad P(L) = \frac{L}{2R^2}$$

380 where R is the cell radius.

381 Then, the probability of ligation between two random points on the cell surface is:

$$382 \quad P(L \leq d_{ligation}) = \frac{d_{ligation}^2}{4R^2}$$

383 Assume that each probe could be ligated as many times as possible, the mean counts of
384 ligated PLA product $i:j$, $X_{i,j}$, follow a binomial distribution:

$$385 \quad X_{i,j} \sim \text{Binomial} \left(n = A_i \times B_j, p = P(L \leq d_{ligation}) \right)$$

386 The expected count of PLA product that is created from random ligation of non-
387 interacting probes is:

$$388 \quad E(X_{i,j}) = \frac{d_{ligation}^2}{4R^2} A_i B_j$$

389 Note that this approximation assumes that each probe can be ligated many times, while
390 the simulation model assumes that each probe can only be ligated at most once.

391 **Simulation model**

392 Assume that each protein molecule and the Prox-seq probe that binds to it are point particles. Let
393 there be n protein targets. Let A_1, A_2, \dots, A_n be the simulation parameters that represent the count
394 of probe A that targets proteins 1, 2, ..., n . Let B_1, B_2, \dots, B_n be the simulation parameters that
395 represent the count of probe B that targets proteins 1, 2, ..., n . Let $c_{1,1}, c_{1,2}, \dots, c_{1,n}, c_{2,1}, c_{2,2}, \dots,$
396 $c_{n,n}$ be the simulation parameters that represent the counts of protein complexes 1:1, 1:2, ..., 1: n ,
397 2:1, 2:2, ..., n : n .

398 The simulation is performed separately on each single cell. For the single cell t , we first
399 generate $A_i^{(t)}$ number of random points on a sphere surface, which correspond to the number of
400 detected probe A that targets protein i on cell t . The coordinates of each point are¹⁵:

$$401 \quad x = R\sqrt{1 - u^2} \cos \theta$$

$$402 \quad y = R\sqrt{1 - u^2} \sin \theta$$

$$403 \quad z = Ru$$

404 where R is the radius of the sphere (taken to be 5 μm , or 5000 units, in our study), u is uniformly
405 distributed over $[-1, 1)$, and θ is uniformly distributed over $[0, 2\pi)$.

406 Without added variance, $A_i^{(t)} = A_i$. With added negative binomial variance:

$$407 \quad A_i^{(t)} \sim \text{NegativeBinomial}(n_{NB}, p_{NB})$$

408 where $n_{NB} = 1.5$ in our study, and $p_{NB} = \left(1 + \frac{A_i}{n_{NB}}\right)^{-1}$. The negative binomial distribution

409 formulated this way provides the probability of getting $A_i^{(t)}$ failures, given n_{NB} successes and p_{NB}
410 is the probability of success. n_{NB} is used to control the variance of the probe count, and p_{NB} is

411 calculated such that the mean of $A_i^{(t)}$ is equal to A_i .

412 Second, we randomly generate $B_i^{(t)}$ number of points on a surface of a sphere, which
413 correspond to the number of detected probe B that targets protein i on cell t. The coordinates of
414 each point are generated identically to above.

415 Without added variance, $B_i^{(t)} = B_i$. With added variance:

$$416 \quad B_i^{(t)} = \frac{B_i}{A_i} \times A_i^{(t)}$$

417 This is to ensure that the counts of detected probe A and probe B that target the same protein are
418 proportional to each other.

419 Third, we randomly generate $c_{i,j}^{(t)}$ number of points on a surface of a sphere, which
420 correspond to the count of protein complex i:j on cell t. Then, these $c_{i,j}^{(t)}$ points are added to the
421 previously generated probe A points targeting protein i $A_i^{(t)}$, and also to the previously generated
422 probe B targeting protein j $B_j^{(t)}$.

423 Without added variance, $c_{i,j}^{(t)} = c_{i,j}$. With added variance:

$$424 \quad c_{i,j}^{(t)} \sim \text{NegativeBinomial}(n_{NB}, p_{NB})$$

425 where $n_{NB} = 1.5$ in our study, and $p_{NB} = \left(1 + \frac{c_{i,j}}{n_{NB}}\right)^{-1}$.

426 Fourth, we calculated the pairwise Euclidean distances between all generated probe A
427 points and all generated probe B points. Finally, we randomly go through the pairs of points that
428 are within a ligation distance threshold (chosen to be 50 nm, or 50 units, in our study), and add
429 the corresponding PLA product to the simulated count matrix. Any probe A and probe B points

430 that are chosen are excluded from future PLA products. In other words, each probe A and each
431 probe B can only be ligated at most once.

432 The number of probe A and probe B points that are not ligated are returned as the
433 simulated non-proximal probe count that is measured by the free oligo modification.

434 The simulation is repeated 100 times to simulate PLA product counts of 100 single cells.
435 The parameters for all simulations are listed in Table S1. All simulations include negative
436 binomial variance, unless stated otherwise.

437 **Calculation of protein count and expected PLA product count**

438 The count of a protein i in a single cell is equal to the total number of detected PLA products that
439 are related to the protein i :

$$440 \quad \text{Protein } i = \sum_{l=1}^n X_{i,l} + \sum_{k=1}^n X_{k,i}$$

441 where $X_{i,l}$ and $X_{k,i}$ indicate the observed (i.e., measured) counts of PLA products $i:l$ and $k:i$,
442 respectively. The PLA product $i:i$ is counted twice towards the protein count to account for the
443 fact that two molecules are present in a homodimer.

444 The expected count of a PLA product $i:j$, $E_{i,j}$, is:

$$445 \quad E_{i,j} = \frac{\sum_{l=1}^n X_{i,l} \times \sum_{k=1}^n X_{k,j}}{\sum_{k=1}^n \sum_{l=1}^n X_{k,l}}$$

446 **Protein complex prediction: iterative method**

447 The count of protein complex $i:j$ is calculated iteratively using the following equation:

$$448 \quad Y_{i,j}^{(m+1)} = X_{i,j} - \frac{\left(\sum_{l=1}^n X_{i,l} - \sum_{l=1}^n Y_{i,l}^{(m)}\right) \times \left(\sum_{k=1}^n X_{k,j} - \sum_{k=1}^n Y_{k,j}^{(m)}\right)}{\sum_{k=1}^n \sum_{l=1}^n X_{k,l} - \sum_{k=1}^n \sum_{l=1}^n Y_{k,l}^{(m)}}$$

449 where $Y_{i,j}^{(m)}$ is the predicted count of protein complex i:j at the mth iteration. The initial values for
 450 all protein complexes are 0.

451 The second term of the right hand side represents the count of PLA product i:j that is
 452 caused by random ligation.

453 After each iteration, a one-sided t-test is performed on the values of $Y_{i,j}^{(m+1)}$ across all
 454 single cells. The alternative hypothesis is that the mean of $Y_{i,j}^{(m+1)}$ is greater than 1. Next, any
 455 $Y_{i,j}^{(m+1)}$ with Benjamini-Hochberg-corrected P-values above 0.05 are set to 0. In other words, any
 456 such PLA products were determined to not represent true protein interactions.

457 There is also a symmetry condition, such that if i:j is a protein complex, then j:i should
 458 also be a protein complex, even if $Y_{j,i}^{(m+1)}$ fails the t-test. This is done by setting $Y_{j,i}^{(m+1)}$ as a
 459 fraction of $Y_{i,j}^{(m+1)}$:

$$460 \quad Y_{j,i}^{(m+1)} = \text{sym_weight} \times Y_{i,j}^{(m+1)}$$

461 where sym_weight is arbitrarily chosen to be 1 in our study.

462 **Protein complex prediction: linear regression (LR) method**

463 For each PLA product i:j, its observed count is regressed onto the product of its corresponding
 464 non-proximal probe A count and non-proximal probe B count, using weighted least squares:

$$465 \quad X_{i,j} \sim \beta_0 + \beta_1 A_i' B_j'$$

466 where A'_i and B'_j are the count of non-proximal probe A targeting protein i, and non-proximal
467 probe B targeting protein j, respectively. The weight for a sample (ie, a single cell) p is:

$$468 \quad w_p = \frac{1}{A'_i B'_j}$$

469 For simulated data, we also scale the interaction term by 10^6 whenever necessary, such that it is
470 close to the orders of magnitude of $X_{i,j}$. A'_i and B'_j are obtained from PLA products that contain
471 the added free oligos. For example, the count of non-proximal CD3 probe A is equal to the count
472 of PLA product CD3:free_oligo_B, and the count of non-proximal CD28 probe B is equal to the
473 count of PLA product free_oligo_A:CD28.

474 Next, we performed a one-sided t-test on the intercept coefficient, and the alternative
475 hypothesis is that $\beta_0 > \beta_{cutoff}$. For simulated data, $\beta_{cutoff} = 1$. For experimental data
476 $\beta_{cutoff} = 10$. All PLA products with Benjamini-Hochberg-corrected P-values below 0.05 are
477 considered to be true protein complexes. The protein complex count, $Y_{i,j}$, is calculated as the
478 difference between the observed PLA product count and the interaction term:

$$479 \quad Y_{i,j} = X_{i,j} - \beta_1 A'_i B'_j$$

480 The LR method is related to the binomial approximation of the random ligation signal
481 above. If the counts of non-proximal probes are perfect proxies for the count of non-interacting
482 probes, then we have the following relationship:

$$483 \quad \beta_1 = \frac{d_{ligation}^2}{4R^2}$$

484 **Protein complex prediction: Ensemble method**

485 The ensemble method relies on solving same quadratic equations as iterative method to
 486 approximate counts of protein complex. The only difference is that it takes protein complex
 487 matrix calculated from LR method as initial values. There is an argument called `df_guess`
 488 embedded in predictive function which is set to be all zeros by default. Note that LR method
 489 should be applied in advance in order to perform ensemble method.

490 **Protein complex prediction: Fisher's exact test**

491 For each PLA product $i:j$, we construct a 2x2 table:

	Probe B = j	Probe B \neq j
Probe A = i	$X_{i,j}$	$\sum_{\substack{l=1 \\ l \neq j}}^n X_{i,l}$
Probe A \neq i	$\sum_{\substack{k=1 \\ k \neq i}}^n X_{k,j}$	$\sum_{\substack{k=1 \\ k \neq i}}^n \sum_{\substack{l=1 \\ l \neq j}}^n X_{k,l}$

492 We perform a one-sided Fisher's exact test using the table. The alternative hypothesis is
 493 that $X_{i,j}$ is higher than the expected value. We perform a Benjamini-Hochberg correction for each
 494 cell on the P-values of all PLA products calculated for that cell. The cell is considered to be
 495 displaying the protein complex if the corrected P-value is below 0.05.

496 **Prediction score mechanism**

$$497 \quad \text{Score} = w_1 * \sum \text{Pearson} - w_2 * \sum \text{Mean}_{\text{deviation}} - w_3 * \sum \text{FPrate}$$

498 where w_1, w_2, w_3 are chosen to be 0.5, 0.4, 0.1 in our study. $\sum \text{Pearson}$ equals to the sum of
 499 Pearson correlation coefficients for every real protein complex between true complex counts and

500 predicted complex counts across all single cells. $\sum Mean_{deviation}$ equals to the sum of absolute
501 difference between mean true counts and predicted counts for every PLA product:

$$502 \quad Mean_{deviation} = \frac{|Mean_{pred} - Mean_{true}|}{Mean_{true}}$$

503 $\sum FPrate$ equals to the sum of ratios of false positive prediction across single cells for every
504 non-existing PLA product.

505 For quantification accuracy evaluation where there are true protein complex counts, we consider
506 parameters $\sum Pearson$ and $\sum Mean_{deviation}$. Pearson correlation coefficient takes single cells
507 into consideration while mean counts can give us information about bulk abundance of different
508 PLA products. We found that poor prediction of PLA counts in single cells might still contribute
509 to seemingly good mean counts estimation, which shed lights on us that Pearson correlation
510 should be a more important and robust parameter than mean counts. For $\sum FPrate$ evaluation
511 where there is no true complex, we use fraction of complex-positive cells to represent how many
512 ratios of single cells are wrongly assigned at least a complex count. According to our multiple
513 tests, each method tends to assign only few false positive reads, mostly only one in some single
514 cells to PLA products. So that we assume false positive rate a minor metric to be considered in
515 our scoring strategy. In conclusion, we arbitrarily choose effector weight for each parameter
516 given relative importance discussed above.

517 **Software implementation**

518 All code is implemented in Python3/Anaconda3 (v4.10.3). The code is deposited at
519 https://github.com/tay-lab/Prox-seq_computation.

520 **Data availability**

521 The raw sequencing data and processed PLA product count data are deposited in NCBI's Gene
522 Expression Omnibus (accession number GSE196130).

523 **Acknowledgements**

524 S.T. was awarded an NIH R01 grant GM127527, NIH MIRA/R35 grant R35GM148231, and a
525 Paul G. Allen Distinguished Investigator Award, which supported this work. M.C. was awarded
526 NIH R01 grants GM126553 and HG011883, and an NSF grant 2016307 which supported this
527 work. This work was supported in part by the Intramural Research program of NIAID, NIH.

528

529

530 **Author Contributions Statement**

531 L.V., H.V.P., J. X. and S.T. conceived of and designed the project. H.V.P., J. X., M.C., and A. K.
532 performed statistical and computational analysis. L.V., H.V.P., J. X. and S.T. wrote the
533 manuscript. S.T. supervised the project. All authors reviewed the manuscript.

534

535 **Competing Interests Statement**

536 The authors declare no competing financial interest.

537

538 **References**

539 1. Tay, S. *et al.* Single-cell NF- κ B dynamics reveal digital activation and analogue information
540 processing. *Nature* **466**, 267–271 (2010).

- 541 2. Patel, A. P. *et al.* Single-cell RNA-seq highlights intratumoral heterogeneity in primary
542 glioblastoma. *Science* **344**, 1396–1401 (2014).
- 543 3. Ziegenhain, C. *et al.* Comparative Analysis of Single-Cell RNA Sequencing Methods. *Mol.*
544 *Cell* **65**, 631-643.e4 (2017).
- 545 4. Vistain, L. *et al.* Quantification of extracellular proteins, protein complexes and mRNAs in
546 single cells by proximity sequencing. *Nat. Methods* **19**, 1578–1589 (2022).
- 547 5. Fredriksson, S. *et al.* Protein detection using proximity-dependent DNA ligation assays. *Nat.*
548 *Biotechnol.* **20**, 473–477 (2002).
- 549 6. Söderberg, O. *et al.* Characterizing proteins and their interactions in cells and tissues using the
550 in situ proximity ligation assay. *Methods* **45**, 227–232 (2008).
- 551 7. Chi, Q., Wang, G. & Jiang, J. The persistence length and length per base of single-stranded
552 DNA obtained from fluorescence correlation spectroscopy measurements using mean field
553 theory. *Phys. Stat. Mech. Its Appl.* **392**, 1072–1079 (2013).
- 554 8. Alsemarz, A., Lasko, P. & Fagotto, F. *Limited significance of the in situ proximity ligation*
555 *assay*. 411355 <https://www.biorxiv.org/content/10.1101/411355v2> (2018)
556 doi:10.1101/411355.
- 557 9. Gayoso, A. *et al.* Joint probabilistic modeling of single-cell multi-omic data with totalVI. *Nat.*
558 *Methods* **18**, 272–282 (2021).
- 559 10. Stoeckius, M. *et al.* Simultaneous epitope and transcriptome measurement in single cells.
560 *Nat. Methods* **14**, 865–868 (2017).
- 561 11. van der Merwe, P. A. & Dushek, O. Mechanisms for T cell receptor triggering. *Nat. Rev.*
562 *Immunol.* **11**, 47–55 (2011).

- 563 12. Esensten, J. H., Helou, Y. A., Chopra, G., Weiss, A. & Bluestone, J. A. CD28
564 Costimulation: From Mechanism to Therapy. *Immunity* **44**, 973–988 (2016).
- 565 13. Miller, J. *et al.* Intercellular adhesion molecule-1 dimerization and its consequences for
566 adhesion mediated by lymphocyte function associated-1. *J. Exp. Med.* **182**, 1231–1241
567 (1995).
- 568 14. Weisstein, E. W. Sphere Line Picking. *Wolfram MathWorld Sphere Line Picking*
569 <https://mathworld.wolfram.com/SphereLinePicking.html>.
- 570 15. Weisstein, E. W. Sphere Point Picking. *Wolfram MathWorld Sphere Point Picking*
571 <https://mathworld.wolfram.com/>.
- 572