

Measuring the sequence-affinity landscape of antibodies with massively parallel titration curves

Rhys M. Adams,^{1,2,*} Thierry Mora,^{3,†} Aleksandra M. Walczak,^{1,†} and Justin B. Kinney^{2,†}

¹*Laboratoire de Physique Théorique, UMR8549, CNRS and École Normale Supérieure, 24, rue Lhomond, 75005 Paris, France*

²*Simons Center for Quantitative Biology, Cold Spring Harbor Laboratory,
1 Bungown Rd., Cold Spring Harbor, NY, 11724, USA*

³*Laboratoire de Physique Statistique, UMR8550, CNRS and École Normale Supérieure, 24, rue Lhomond, 75005 Paris, France*

(Dated: December 28, 2016)

Despite the central role that antibodies play in the adaptive immune system and in biotechnology, much remains unknown about the quantitative relationship between an antibody’s amino acid sequence and its antigen binding affinity. Here we describe a new experimental approach, called Tite-Seq, that is capable of measuring binding titration curves and corresponding affinities for thousands of variant antibodies in parallel. The measurement of titration curves eliminates the confounding effects of antibody expression and stability that arise in standard deep mutational scanning assays. We demonstrate Tite-Seq on the CDR1H and CDR3H regions of a well-studied scFv antibody. Our data shed light on the structural basis for antigen binding affinity and suggests a role for secondary CDR loops in establishing antibody stability. Tite-Seq fills a large gap in the ability to measure critical aspects of the adaptive immune system, and can be readily used for studying sequence-affinity landscapes in other protein systems.

I. INTRODUCTION

During an infection, the immune system must recognize and neutralize invading pathogens. B-cells contribute to immune defense by producing antibodies, proteins that bind specifically to foreign antigens. The astonishing capability of antibodies to recognize virtually any foreign molecule has been repurposed by scientists in wide variety of experimental techniques (immunofluorescence, western blots, ELISA, ChIP-Seq, etc.). Antibody-based therapeutic drugs have also been developed for treating many different diseases, including cancer [1].

Much is known about the qualitative mechanisms of antibody generation and function [2]. The antigenic specificity of antibodies in humans, mice, and most jawed vertebrates is primarily governed by six complementarity determining regions (CDRs), each roughly 10 amino acids (aa) long. Three CDRs (denoted CDR1H, CDR2H, and CDR3H) are located on the antibody heavy chain, and three are on the light chain. During B-cell differentiation, these six sequences are randomized through V(D)J recombination, then selected for functionality as well as against the ability to recognize host antigens. Upon participation in an immune response, CDR regions can further undergo somatic hypermutation and selection, yielding higher-affinity antibodies for specific antigens. Among the CDRs, CDR3H is the most highly variable and typically contributes the most to antigen specificity; less clear are the functional roles of the other CDRs, which often do not interact with the target antigen directly.

Many high-throughput techniques, including phage display [3–5], ribosome display [6], yeast display [7, 8], and mammalian cell display [9], have been developed for optimizing antibodies *ex vivo*. Advances in DNA sequencing technology have also made it possible to effectively monitor both antibody and T-cell receptor diversity within immune repertoires, e.g. in healthy individuals [10–21], in specific tissues [22], in individuals with diseases [23] or following vaccination [24–28]. Yet many questions remain about basic aspects of the quantitative relationship between antibody sequence and antigen binding affinity. How many different antibodies will bind a given antigen with specified affinity? How large of a role do epistatic interactions between amino acid positions within the CDRs have on antigen binding affinity? How is this sequence-affinity landscape navigated by the V(D)J recombination process, or by somatic hypermutation? Answering these and related questions is likely to prove critical for developing a systems-level understanding of the adaptive immune system, as well as for using antibody repertoire sequencing to diagnose and monitor disease.

Recently developed “deep mutational scanning” (DMS) assays [29] provide one potential method for measuring binding affinities with high enough throughput to effectively explore antibody sequence-affinity landscapes. In DMS experiments, one begins with a library of variants of a specific protein. Proteins that have high levels of a particular activity of interest are then enriched via one or more rounds of selection, which can be carried out in a variety of ways. The set of enriched sequences is then compared to the initial library, and protein sequences (or mutations within these sequences) are scored according to how much this enrichment procedure increases their prevalence.

Multiple DMS assays have been described for investigating protein-ligand binding affinity. But no DMS assay

* Current address: Francis Crick Institute, 1 Midland Rd, London NW1 1AT, United Kingdom.

† These authors contributed equally.

has yet been shown to provide absolute quantitative binding affinity measurements, i.e., dissociation constants in molar units. For example, one of the first DMS experiments [30] used phage display technology to measure how mutations in a WW domain affect the affinity of this domain for its peptide ligand. These data were sufficient to compute enrichment ratios and corresponding sequence logos, but they did not yield quantitative affinities. Analogous experiments have since been performed on antibodies using yeast display [31, 32] and mammalian cell display [9]. Yeast-display-based DMS assays have also proven particularly useful for mapping protein epitopes that are targeted by specific antibodies of interest [32–34]. Still, none of these approaches provides quantitative affinity values. SORTCERY [31, 35], a DMS assay that combines yeast display and quantitative modeling, has been shown to provide approximate rank-order values for the affinity of a specific protein for short unstructured peptides of varying sequence. Determining quantitative affinities from SORTCERY data, however, requires separate low-throughput calibration measurements [31]. Moreover, it is unclear how well SORTCERY, if applied to a library of folded proteins rather than unstructured peptides, can distinguish sequence-dependence effects on affinity from sequence-dependent effects on protein expression and stability. Other recent work has described a DMS assay, again based on yeast display, for measuring fold-changes in affinity relative to a reference protein [36]. This method, however, does not provide absolute values for dissociation constants, is vulnerable to the confounding effects of sequence-dependent expression and protein stability, and was observed to have only a 10-fold dynamic range.

To enable massively parallel measurements of absolute binding affinities for antibodies and other structured proteins, we have developed an assay called “Tite-Seq.” Tite-Seq, like SORTCERY, builds on the capabilities of Sort-Seq, an experimental strategy that was first developed for studying transcriptional regulatory sequences in bacteria [37]. Sort-Seq combines fluorescence-activated cell sorting (FACS) with high-throughput sequencing to provide massively parallel measurements of cellular fluorescence. In the Tite-Seq assay, Sort-Seq is applied to antibodies displayed on the surface of yeast cells and incubated with antigen at a wide range of concentrations. From the resulting sequence data, thousands of antibody-antigen binding titration curves and their corresponding absolute dissociation constants (here denoted K_D) can be inferred. By assaying full binding curves, Tite-Seq is able to measure affinities over many orders of magnitude [38]. Moreover, the resulting affinity values provided by Tite-Seq are not confounded by the (rather substantial) effect that sequence variation can have on either (a) the amount of protein expressed on the surface of cells or (b) the specific activity of displayed proteins (i.e., the fraction of protein molecules that are functional).

We demonstrated Tite-Seq on a protein library derived from a well-studied single-chain variable fragment (scFv)

antibody specific to the small molecule fluorescein [7, 39]. Mutations were restricted to CDR1H and CDR3H regions, which are known to play an important role in the antigen recognition of this scFv [39, 40]. The resulting affinity measurements were validated with binding curves for a handful of clones measured using standard low-throughput flow cytometry. Our Tite-Seq measurements reveal both expected and unexpected differences between the effects of mutations in CDR1H and CDR3H. These data also shed light on structural aspects of antigen recognition that are independent of effects on antibody stability.

II. RESULTS

A. Overview of Tite-Seq

Our general strategy is illustrated in Fig. 1. First, a library of variant antibodies is displayed on the surface of yeast cells (Fig. 1A). The composition of this library is such that each cell displays a single antibody variant, and each variant is expressed on the surface of multiple cells. Cells are then incubated with the antigen of interest, bound antigen is fluorescently labeled, and fluorescence-activated cell sorting (FACS) is used to sort cells one-by-one into multiple “bins” based on this fluorescent readout (Fig. 1B). Deep sequencing is then used to survey the antibody variants present in each bin. Because each variant antibody is sorted multiple times, it will be associated with a histogram of counts spread across one or more bins (Fig. 1C). The spread in each histogram is due to cell-to-cell variability in antibody expression, and to the inherent noisiness of flow cytometry measurements. Finally, the histogram corresponding to each antibody variant is used to compute an “average bin number” (Fig. 1C, dots), which serves as a proxy measurement for the average amount of bound antigen per cell.

It has previously been shown that K_D values can be accurately measured using yeast-displayed antibodies by taking binding titration curves, i.e., by measuring the average amount of bound antigen as a function of antigen concentration [8, 41]. The median fluorescence f of labeled cells is expected to be related to antigen concentration via

$$f = A \frac{c}{c + K_D} + B \quad (1)$$

where A is proportional to the number of functional antibodies displayed on the cell surface, B accounts for background fluorescence, and c is the concentration of free antigen in solution. Fig. 1D illustrates the shape of curves having this form. By using flow cytometry to measure f on clonal populations of yeast at different antigen concentrations c , one can infer curves having the sigmoidal form shown in Eq. 1 and thereby learn K_D . Such measurements, however, can only be performed in a low-throughput manner.

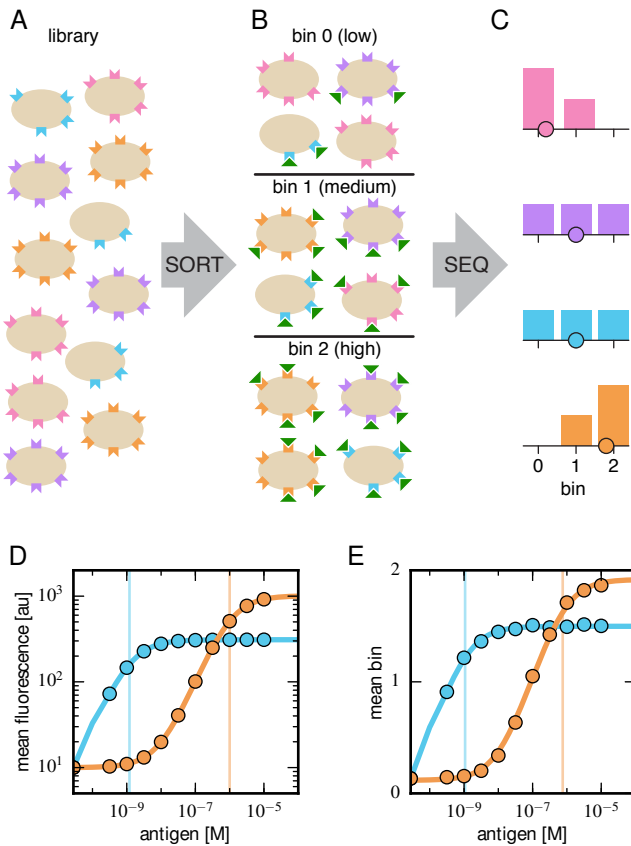


FIG. 1. Schematic illustration of Tite-Seq. (A) A library of variant antibodies (various colors) are displayed on the surface of yeast cells (tan). (B) The library is exposed to antigen (green triangles) at a defined concentration, cell-bound antigen is fluorescently labeled, and FACS is used to sort cells into bins according to measured fluorescence. (C) The antibody variants in each bin are sequenced and the distribution of each variant across bins is computed (histograms; colors correspond to specific variants). The mean bin number (dot) is then used to quantify the typical amount of bound antigen per cell. (D) Binding titration curves (solid lines) and corresponding K_D values (vertical lines) can be inferred for individual antibody sequences by using the mean fluorescence values (dots) obtained from flow cytometry experiments performed on clonal populations of antibody-displaying yeast. (E) Tite-Seq consists of performing the Sort-Seq experiment in panels A-C at multiple antigen concentrations, then inferring binding curves using mean bin number as a proxy for mean cellular fluorescence. This enables K_D measurements for thousands of variant antibodies in parallel. We note that the Tite-Seq results illustrated in panel E were simulated using three bins under idealized experimental conditions, as described in Appendix A. The inference of binding curves from real Tite-Seq data is more involved than this panel might suggest, due to the multiple sources of experimental noise that must be accounted for.

Tite-Seq allows thousands of binding titration curves to be measured in parallel. The Sort-Seq procedure illustrated in Fig. 1A-C is performed at multiple antigen concentrations, and the resulting average bin number for each variant antibody is plotted against concentration. Sigmoidal curves are then fit to these proxy measurements, enabling K_D values to be inferred for each variant.

We emphasize that K_D values cannot, in general, be accurately inferred from Sort-Seq experiments performed at a single antigen concentration. Because the relationship between binding and K_D is sigmoidal, the amount of bound antigen provides a quantitative readout of K_D only when the concentration of antigen used in the labeling procedure is comparable in magnitude to K_D . However, single mutations within a protein binding domain often change K_D by multiple orders of magnitude. Sort-Seq experiments used to measure sequence-affinity landscapes must therefore be carried out over a range of concentrations large enough to encompass this variation.

Furthermore, as illustrated in Figs. 1C and 1D, different antibody variants often lead to different levels of functional antibody expression on the yeast cell surface. If one performs Sort-Seq at a single antigen concentration, high affinity (low K_D) variants with low expression (blue variant) may bind less antigen than low affinity (high K_D) variants with high expression (orange variant). Only by measuring full titration curves can the effect that sequence has on affinity be deconvolved from sequence-dependent effects on functional protein expression.

B. Proof-of-principle Tite-Seq experiments

To test the feasibility of Tite-Seq, we used a well-characterized antibody-antigen system: the 4-4-20 single chain variable fragment (scFv) antibody [7], which binds the small molecule fluorescein with $K_D = 1.2$ nM [8]. This system was used in early work to establish the capabilities of yeast display [7], and a high resolution co-crystal structure of the 4-4-20 antibody bound to fluorescein, shown in Fig. 2A, has been determined [42]. An ultra-high-affinity ($K_D = 270$ fM) variant of this scFv, called 4m5.3, has also been found [39]. In what follows, we refer to the 4-4-20 scFv from [7] as WT, and the 4m5.3 variant from [39] as OPT.

The scFv was expressed on the surface of yeast as part of the multi-domain construct illustrated in Fig. 2B and previously described in [7]. Following [39], we used fluorescein-biotin as the antigen and labeled scFv-bound antigen with streptavidin-RPE (PE). The amount of surface-expressed protein was separately quantified by labeling the C-terminal c-Myc tag using anti-c-Myc primary antibodies, followed by secondary antibodies conjugated to Brilliant Violet 421 (BV). See Appendix B for details on this labeling procedure.

Two different scFv libraries were assayed simultane-

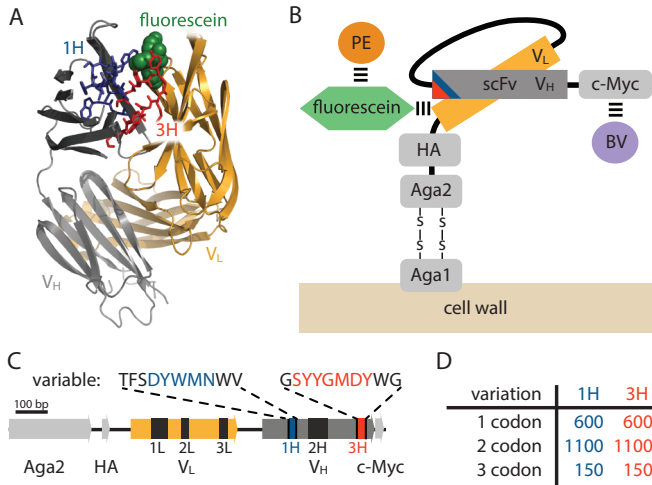


FIG. 2. Yeast display construct and antibody libraries. (A) Co-crystal structure of the 4-4-20 (WT) antibody from [42] (PDB code 1FLR). The CDR1H and CDR3H regions are colored blue and red, respectively. (B) The yeast display scFv construct from [7] that was used in this study. Antibody-bound antigen (fluorescein) was visualized using PE dye. The amount of surface-expressed protein was separately visualized using BV dye. Approximate location of the CDR1H (blue) and CDR3H (red) regions within the scFv are illustrated. (C) The gene coding for this scFv construct, with the six CDR regions indicated. The WT sequence of the two 10 aa variable regions are also shown. (D) The number of 1-, 2-, and 3-codon variants present in the 1H and 3H scFv libraries. Fig. 2 – figure supplement 1 shows the cloning vector used to construct the CDR1H and CDR3H libraries, as well as the form of the resulting expression plasmids.

ously. In the “1H” library, a 10 aa region encompassing the CDR1H region of the WT scFv (see Fig. 2C) was mutagenized using a microarray-synthesized oligos (see Appendix C for details on library generation). The resulting 1H library consisted of all 600 single-codon variants of this 10 aa region, 1100 randomly chosen 2-codon variants, and 150 random 3-codon variants (Fig. 2D). An analogous “3H” library was generated for a 10 aa region containing the CDR3H region of this scFv. In all of the Tite-Seq experiments described below, these two libraries were pooled together and supplemented with WT and OPT scFvs, as well with a nonfunctional scFv referred to as Δ .

Tite-Seq was carried out as follows. Yeast cells expressing scFv from the mixed library were incubated with fluorescein-biotin at one of eleven concentrations: 0 M, $10^{-9.5}$ M, 10^{-9} M, $10^{-8.5}$ M, 10^{-8} M, $10^{-7.5}$ M, 10^{-7} M, $10^{-6.5}$ M, 10^{-6} M, $10^{-5.5}$ M, and 10^{-5} M. After subsequent PE labeling of bound antigen, cells were sorted into four bins using FACS (Fig. 3A). Separately, BV-labeled cells were sorted according to measured scFv expression levels (Fig. 3B). The number of cells sorted into each bin is shown in Fig. 3C. Each bin of cells was re-grown and bulk DNA was extracted. The 1H and 3H variable regions were then PCR amplified and sequenced

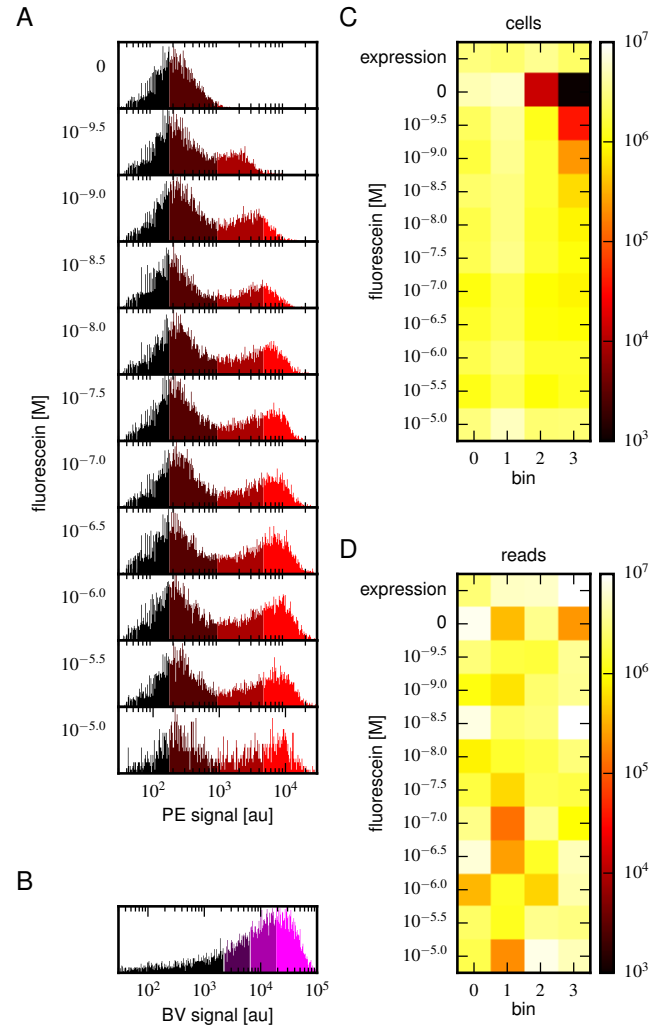


FIG. 3. Details of our Tite-Seq experiments. (A) Gates used to sort cells based on PE fluorescence, which provides a readout of bound antigen. Cells were labeled at the eleven different antigen concentrations. Shades of red indicate the four fluorescence gates used to sort cells; these correspond to bins 0,1,2, and 3 (from left to right). (B) Gates, indicated in shades of purple, used to sort cells based on BV fluorescence, which provides a readout of antibody expression. (C) The number of cells sorted into each bin. (D) The number of Illumina reads obtained from each bin of sorted cells after quality control measures were applied. The data shown in this figure corresponds to a single Tite-Seq experiment. Fig. 3 – figure supplement 1 and Fig. 3 – figure supplement 2 show data for two independent replicates of this experiment.

using paired-end Illumina sequencing, as described in Appendix D. The final data set consisted of an average of 2.6×10^6 sequences per bin across all 48 bins (Fig. 3D). Three independent replicates of this experiment were performed on three different days.

For each variant scFv gene, a K_D value was inferred by fitting a binding curve to the resulting Tite-Seq data, with separate curves independently fit to data from each

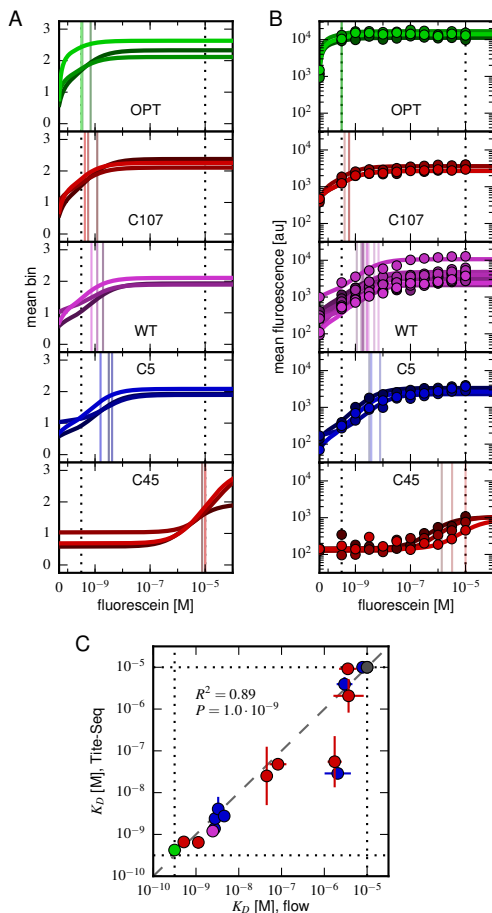


FIG. 4. Accuracy and precision of Tite-Seq. (A) Binding curves and K_D measurements inferred from Tite-Seq data. (B) Mean fluorescence values (dots) and corresponding inferred binding curves (lines) obtained by flow cytometry measurements for five selected scFvs (WT, OPT, C5, C45, and C107). In (A,B), values corresponding to 0 M fluorescein are plotted on the left-most edge of the plot, dotted lines show the upper (10^{-5} M) and lower ($10^{-9.5}$ M) limits on K_D sensitivity, vertical lines show inferred K_D values, and different shades correspond to different replicate experiments. (C) Comparison of the Tite-Seq-measured and flow-cytometry-measured K_D values for all clones tested. Colors indicate different scFv protein sequences as follows: WT (purple), OPT (green), Δ (black), 1H clones (blue), and 3H clones (red). Each K_D value indicates the mean $\log_{10} K_D$ value obtained across all replicates, with error bars indicating standard error. Clones with K_D outside of the affinity range are drawn on the boundaries of this range, which are indicated with dotted lines. The coefficient of determination (R^2) between \log Tite-seq and \log flow mean K_D includes clones outside of the affinity range; in such cases, the corresponding boundary value ($10^{-9.5}$ M or $10^{-5.0}$ M) has been used. The amino acid sequences and measured K_D values for all clones tested are provided in Table I. Fig. 4 – figure supplement 1 provides plots, analogous to panels A and B, for all of the assayed clones. Fig. 4 – figure supplement 2 compares K_D and E values obtained across all three Tite-Seq replicates. Fig. 4 – figure supplement 3 quantifies measurement error using synonymous mutants. Fig. 4 – figure supplement 4 provides information about the library compositions. Fig. 4 – figure supplement 5 illustrates the poor correlation between scFv enrichment and Tite-seq measured K_D values. Fig. 4 – figure supplement 6 shows a 2-fold difference in the specific activities of OPT and WT scFvs. Fig. 4 – figure supplement 7 illustrates the simulations we used in Fig. 4 – figure supplement 8 to illustrate the ability of our pipeline to correctly infer K_D values.

Tite-Seq experiment (Fig. 4A). As illustrated in Fig. 1E, this fitting procedure uses the sigmoidal function in Eq. 1 to model mean bin number as a function of antigen concentration. However, the need to account for multiple sources of noise in the Tite-Seq experiment necessitates a more complex procedure than Fig. 1E might suggest; the details of this inference procedure are described in Appendix E.

Separately, the Sort-Seq data obtained by sorting the BV-labeled libraries were used to determine the expression level of each scFv. Specifically, we use E to denote (for each scFv in the library) the mean bin number that results from this expression-based sorting; this E value provides a measurement of the surface expression level of that scFv. All E values have been scaled so that the mean of such measurements for all synonymous WT scFv gene variants is 1.0.

C. Low-throughput validation experiments

To judge the accuracy of Tite-Seq, we separately measured binding curves for individual scFv clones as described for Fig. 1D. In addition to the WT, OPT, and Δ scFvs, we assayed eight clones from the 1H library (named C3, C5, C7, C18, C22, C132, C133 and C144) and eight clones from the 3H library (C39, C45, C93, C94, C102, C103, C107, C112). Each clone underwent the same labeling procedure as in the Tite-Seq experiment, after which median fluorescence values were measured using standard flow cytometry. K_D values were then inferred by fitting binding curves of the form in Eq. 1 using the procedure described in Appendix F. These curves, which can be directly compared to the Tite-Seq measurement (Fig. 4A), are plotted in Fig. 4B; at least three replicate binding curves were measured for each clone. See Fig. 4 – figure supplement 1 for the titration curves of all the tested clones.

D. Tite-Seq can measure dissociation constants

Fig. 4C reveals a strong correspondence between the K_D values measured by Tite-Seq and those measured using low-throughput flow cytometry. The robustness of Tite-Seq is further illustrated by the consistency of K_D values measured for the WT scFv. Using Tite-Seq, and averaging the results from the 33 synonymous variants and over all three replicates, we determined $K_D = 10^{-8.87 \pm 0.02}$ M for the WT scFv. These measurements are consistent with the measurement of $K_D = 10^{-8.61 \pm 0.07}$ M obtained by averaging low-throughput flow cytometry measurements across 10 replicates, and coincides with the previously measured value of 1.2 nM $= 10^{-8.9}$ M reported in [8]. The three independent replicate Tite-Seq experiments give reproducible results as measured by direct comparison (Fig. 4 – figure supplement 2), from synonymous mutant variation (Fig. 4 –

figure supplement 3) and library composition 4 – figure supplement 4) with Pearson coefficients ranging from $r = 0.82$ to $r = 0.89$ for all the measured K_D values between replicates; note that K_D values outside of the sensitivity range are included in the calculation of these Pearson coefficients as described in the Fig. 4 caption.

The error bars for K_D values in Fig. 4C calculated from the variability of the fits to different replicates therefore support the reproducibility of the experiment. The main discrepancy in these error bar calculations occurred for clones c22 and c102 (see also Fig. 4 – figure supplement 1). The reason for this discrepancy is currently unclear. We note that Tite-Seq-measured K_D values for these two clones are close to 10^{-7} M, and that the analysis of synonymous variants (Fig. 4 – figure supplement 3) found that Tite-Seq-measured K_D s in this region exhibited the largest variations.

The necessity of performing K_D measurements over a wide range of antigen concentrations is illustrated in Fig. 4 – figure supplement 5. At each antigen concentration used in our Tite-Seq experiments, the enrichment of scFvs in the high-PE bins correlated poorly with the K_D values inferred from full titration curves. Moreover, at each antigen concentration used, a detectable correlation between K_D and enrichment was found only for scFvs with K_D values close to that concentration.

Fig. 4 – figure supplement 6 suggests a possible reason for the weak correlation between K_D values and enrichment in high-PE bins. We found that, at saturating concentrations of fluorescein (2 μ M), cells expressing the OPT scFv bound twice as much fluorescein as cells expressing the WT scFv. This difference was not due to variation in the total amount of displayed scFv, which one might control for by labeling the c-Myc epitope as in [35]. Rather, this difference in binding reflects a difference in the specific activity of displayed scFvs. Yeast display experiments performed at a single antigen concentration cannot distinguish such differences in specific activity from differences in scFv affinity.

To further test the capability of Tite-Seq to infer dissociation constants from sequencing data over a wide range of values, as well as to validate our analysis procedures, we simulated Tite-Seq data *in silico* and analyzed the results using the same analysis pipeline that we used for our experiments. Details about the simulations are given in Appendix G. The simulated data is illustrated in Fig. 4 – figure supplement 7. K_D values inferred from these simulated data agreed to high accuracy with the K_D used in the simulation (Fig. 4 – figure supplement 8), thus validating our analysis pipeline.

E. Properties of the affinity and expression landscapes

Fig. 5 shows the effect that every single-amino-acid substitution mutation within the 1H and 3H variable regions has on affinity and on expression; histograms of

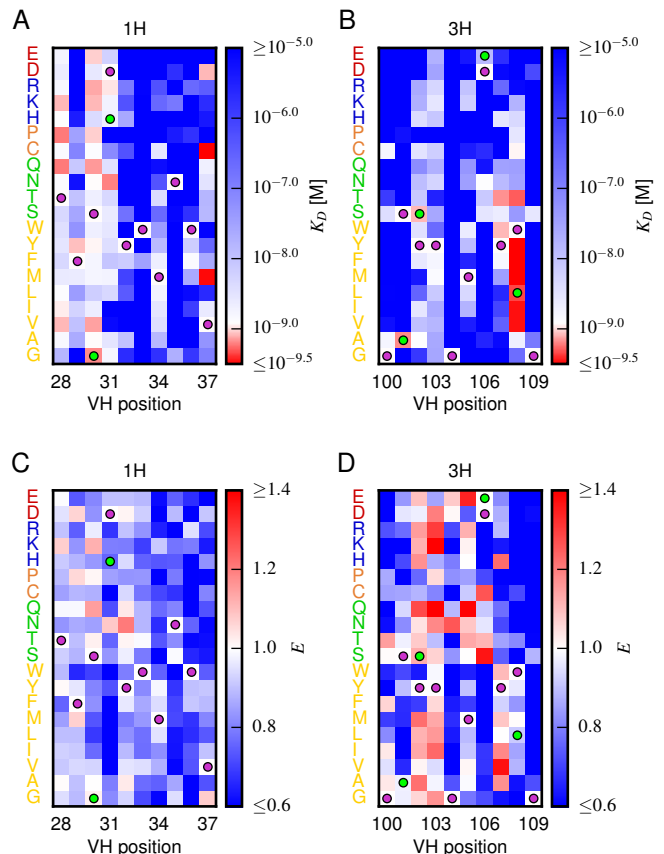


FIG. 5. Effects of substitution mutations on affinity and expression. Heatmaps show the measured effects on affinity (A,B) and expression (C,D) of all single amino acid substitutions within the variable regions of the 1H (A,C) and 3H (B,D) libraries. Purple dots indicate residues of the WT scFv. Green dots indicate non-WT residues in the OPT scFv. Fig. 5 – figure supplement 1 provides histograms of the non-WT values displayed in panels A-D. Fig. 5 – figure supplement 2 compares the effects on K_D of both single-point and multi-point mutations.

these effects are provided in Fig. 5 – figure supplement 1. In both regions, the large majority of mutations weaken antigen binding (1H: 88%; 3H: 93%), with many mutations increasing K_D above our detection threshold of 10^{-5} M (1H: 36%; 3H: 52%). Far fewer mutations reduced K_D (1H: 12%; 3H: 7%), and very few dropped K_D below our detection limit of $10^{-9.5}$ M (1H: 0%; 3H: 3%). Histograms of the effect of 2 or 3 amino acid changes relative to WT, shown in Fig. 5 – figure supplement 2A, show that multiple random mutations tend to further deteriorate affinity. We also observed that mutations within the 3H variable region have a larger effect on affinity than do mutations in the 1H variable region. Specifically, single amino acid mutations in 3H were seen to increase K_D more than mutations in 1H (1H median $K_D = 10^{-6.84}$; 3H median $K_D \gtrsim 10^{-5.0}$; $P = 4.7 \times 10^{-4}$, one-sided Mann-Whitney U test). This result suggests that binding affinity is more sensitive to variation in CDR3H than

to variation in CDR1H, a finding that is consistent with the conventional understanding of these antibody CDR regions [43, 44].

Our observations are thus fully consistent with the hypothesis that the amino acid sequences of the CDR1H and CDR3H regions of the WT scFv have been selected for high affinity binding to fluorescein. We know this to be true, of course; still, this result provides an important validation of our Tite-Seq measurements.

To further validate our Tite-Seq affinity measurements, we examined positions in the high affinity OPT scFv (from [39]) that differ from WT and that lie within the 1H and 3H variable regions. As illustrated in Figs. 5A and 5B, five of the six OPT-specific mutations reduce K_D or are nearly neutral. Previous structural analysis [40] has suggested that D106E, the only OPT mutation that we find significantly increases K_D , may indeed disrupt antigen binding on its own while still increasing affinity in the presence of the S101A mutation.

Next, we used our measurements to build a “matrix model” [45] (also known as a “position-specific affinity matrix,” or PSAM [46]) describing the sequence-affinity landscape of these two regions. Our model assumed that the $\log_{10} K_D$ value for an arbitrary amino acid sequence could be computed from the $\log_{10} K_D$ value of the WT scFv, plus the measured change in $\log_{10} K_D$ produced by each amino acid substitution away from WT. We evaluated our matrix models on the 1H and 3H variable regions of OPT, finding an affinity of $10^{-9.16}$ M. Our simple model for the sequence affinity landscape of this scFv therefore correctly predicts that OPT has higher affinity than WT. The quantitative affinity predicted by our model does not match the known affinity of the OPT scFv ($K_D = 10^{-12.6}$ M), but this is unsurprising for three reasons. First the OPT scFv differs from WT in 14 residues, only 6 of which are inside the 1H and 3H variable regions assayed here. Second, one of the OPT mutations (W108L) reduces K_D below our detection threshold of $10^{-9.5}$ M; in building our matrix model, we set this value equal to $10^{-9.5}$, knowing it would likely underestimate the affinity-increasing effect of the mutation. Third, our additive model ignores potential epistatic interactions. Still, we thought it worth asking how likely it would be for 6 random mutations within the 1H and 3H variable regions to reduce affinity as much as our model predicts for OPT. We therefore simulated a large number (10^7) of variants having a total of 6 substitution mutations randomly scattered across the 1H and 3H variable regions. The fraction of these random sequences that had an affinity at or below our predicted affinity for OPT was 4.7×10^{-5} . This finding is fully consistent with the fact that the mutations in OPT relative to WT were selected for increased affinity, an additional confirmation of the validity of our Tite-Seq measurements.

The sequence-expression landscape measured in our separate Sort-Seq experiment yielded qualitatively different results (Figs. 5C and 5D). We observed no significant difference in the median effect that mutations in

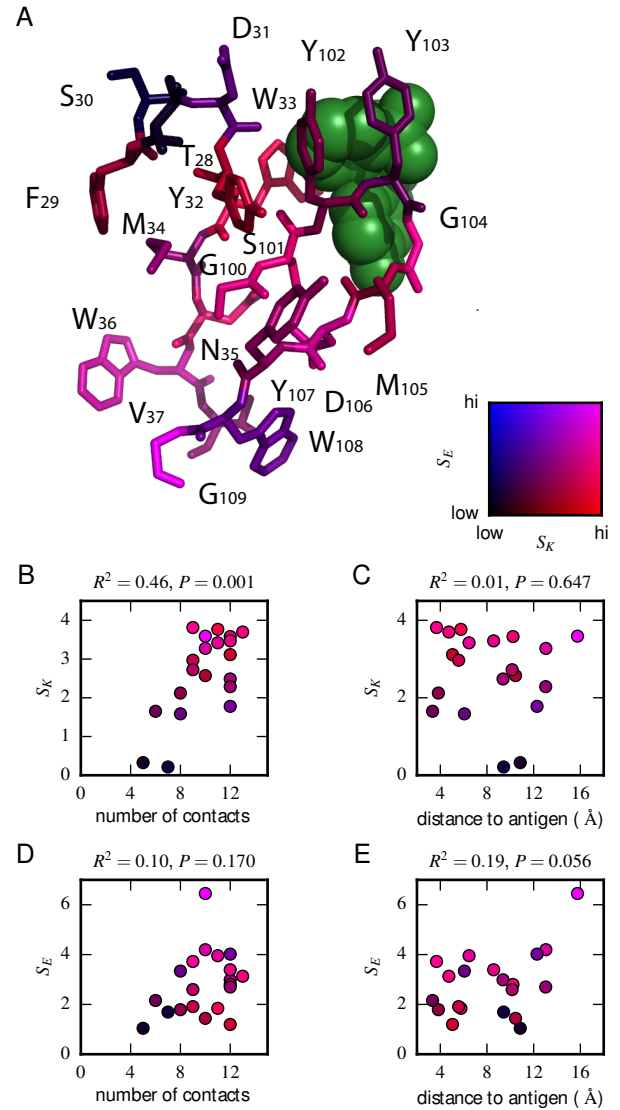


FIG. 6. Structural context of mutational effects. (A) Crystal structure [42] of the CDR1H and CDR3H variable regions of the WT scFv in complex with fluorescein (green). Each residue (CDR1H: positions 28-37; CDR3H: positions 100-109) is colored according to the S_K and S_E values computed for that position. These variables, S_K and S_E , respectively quantify the sensitivity of K_D and E to amino acid substitutions at each position, with larger values greater sensitivity; see Eqs. 2 and 3 for definitions of these quantities. (B,C) For each position in the CDR1H and CDR3H variable regions, S_K is plotted against either (B) the number of contacts the WT residue makes within the protein structure, or (C) the distance of the WT residue to the fluorescein molecule. (D,E) Similarly, S_E is plotted against either (D) the number of contacts or (E) the distance to the antigen. R^2 is the coefficient of determination.

the variable regions of 1H (median $E = 0.826$) versus 3H (median $E = 0.822$) have on expression ($P = 0.96$, two-sided Mann-Whitney U test); see also Fig. 5 – figure

supplement 1. The variance in these effects, however, was larger in 3H than in 1H ($P = 9.9 \times 10^{-16}$, Levene’s test). These results suggest two things. First, the 3H variable region appears to have a larger effect on scFv expression than the 1H variable region has. At the same time, since we observe fewer beneficial mutations in 1H (Fig. 5 C) than in 3H (Fig. 5 D), the WT sequence appears to be more highly optimized for expression in CDR1H than in CDR3H. The effect of double or triple mutations further reduced expression in both CDRs (Fig. 5 – figure supplement 2B), similar to what was observed for affinity.

F. Structural correlates of the sequence-affinity landscape

We asked if the sensitivity of the antibody to mutations could be understood from a structural perspective. To quantify sensitivity of affinity and expression at each position i , we computed two quantities:

$$S_K^i = \sqrt{\left\langle (\log_{10} K_D^{ia} - \log_{10} K_D^{\text{WT}})^2 \right\rangle_{a|i}}, \quad (2)$$

$$S_E^i = \sqrt{\left\langle (E^{ia} - E^{\text{WT}})^2 \right\rangle_{a|i}}. \quad (3)$$

Here, K_D^{WT} and E^{WT} respectively denote the dissociation constant and expression level measured for the WT scFv, K_D^{ia} and E^{ia} denote analogous quantities for the scFv with a single substitution mutation of amino acid a at position i , and $\langle \cdot \rangle_{a|i}$ denotes an average computed over the 19 non-WT amino acids at that position.

Fig. 6A shows the known structure [42] of the 1H and 3H variable regions of the WT scFv in complex with fluorescein. Each residue is colored according to the S_K and S_E values computed for its position. To get a better understanding of what aspects of the structure might govern affinity, we plotted S_K values against two other quantities: the number of amino acid contacts made by the WT residue within the antibody structure (Fig. 6B), and the distance between the WT residue and the antigen (Fig. 6C). We found a strong correlation between S_K and the number of contacts, but no significant correlation between S_K and distance to antigen. By contrast, S_E did not correlate significantly with either of these structural quantities (Figs. 6D and 6E).

III. DISCUSSION

We have described a massively parallel assay, called Tite-Seq, for measuring the sequence-affinity landscape of antibodies. The range of affinities measured in our Tite-Seq experiments ($10^{-9.5}$ M to $10^{-5.0}$ M) includes a large fraction of the physiological range relevant to affinity maturation (10^{-10} M to $\sim 10^{-6}$ M) [47–49]. Expanding the measured range of affinities below $10^{-9.5}$ M might require larger volume labeling reactions, but would

be straight-forward. Tite-Seq therefore provides a potentially powerful method for mapping the sequence-affinity trajectories of antibodies during the affinity maturation process, as well as for studying other aspects of the adaptive immune response.

The details of our Tite-Seq experiments (e.g., 11 antigen concentrations, 4 sorting bins per concentration, etc.) were chosen largely for experimental convenience. The effects of varying these parameters have not been systematically explored, and a future investigation of these effects might be valuable. Fig. 4 – figure supplement 8 does illustrate, via simulation, the effect of read depth on the precision of measured K_D values. These simulations, along with an analysis of synonymous variants (Fig. 4 – figure supplement 3), suggest that the primary source of noise in our experiments came not from a lack of sorted cells or Illumina reads, but rather from the inefficient post-sort recovery of antibody sequences. We therefore suggest that improvements to our post-sort DNA recovery protocol might substantially improve the resolution of Tite-Seq.

Tite-Seq fundamentally differs from prior DMS experiments in that full binding titration curves, not two-bin enrichment statistics, are used to determine binding affinities. The measurement of binding curves provides three major advantages. First, binding curves provide absolute K_D values in molar units, not just rank-order affinities, like those provided by SORTCERY [31], or relative affinity ratios, like those provided by the method of [36]. Second, because ligand binding is a sigmoidal function of affinity, DMS experiments performed at a single ligand concentration (e.g., [36]) are insensitive to receptor K_D s that differ substantially from this ligand concentration. Yet mutations within a protein’s binding domain often change K_D by multiple orders of magnitude. Binding curves, by contrast, integrate measurements over a wide range of concentrations and are therefore sensitive to a wide range of K_D s.

The third advantage of measuring binding curves pertains to the fact that protein sequence determines not just ligand-binding affinity, but also the quantity and specific activity of surface-displayed proteins. Our data (Fig. 4 – figure supplement 5 and Fig. 4 – figure supplement 6) suggest that these confounding effects can be large and that they can distort yeast display affinity measurements computed from enrichment statistics gathered at a single antigen concentration. Strong sequence-dependent effects on both the expression and specific activity of yeast-displayed proteins has been reported by other groups as well (e.g., [50]), although the absence of such effects has also been reported (e.g., [36]). Ultimately, the magnitude of these effects is likely to vary substantially from protein to protein. It should also be noted that many DMS studies using yeast display (e.g., epitope mapping studies [32–34]) might not suffer from these potentially confounding effects, and in such cases it probably makes sense to employ a simpler experimental design than is required for Tite-Seq. Nevertheless, Tite-Seq or experi-

mental methods that assay full binding curves are probably essential if one wants to quantitatively and reliably measure K_D values in a massively parallel fashion.

We wish to emphasize, more generally, that changing a protein’s amino acid sequence can be expected to change multiple biochemical properties of that protein. Our work illustrates the importance of designing massively parallel assays that can disentangle these multiple effects so that measurements of a specific activity of interest can be obtained. Tite-Seq provides a general solution to this problem for massively parallel studies of protein-ligand binding. Indeed, the Tite-Seq procedure described here can be readily applied to any protein binding assay that is compatible with yeast display and FACS. Many such assays have been developed [51]. We expect that Tite-Seq can also be readily adapted for use with other expression platforms, such as mammalian cell display [9].

Our Tite-Seq measurements reveal interesting distinctions between the effects of mutations in the CDR1H and CDR3H regions of the anti-fluorescein scFv antibody studied here. As expected, we found that variation in and around CDR3H had a larger effect on affinity than variation in and around CDR1H. We also found that CDR1H is more optimized for protein expression than is CDR3H, an unexpected finding that appears to be novel. Yeast display expression levels are known to correlate with thermostability [52]. Our data is limited in scope, and we remain cautious about generalizing our observations to arbitrary antibody-antigen interactions. Still, this finding suggests the possibility that secondary CDR regions (such as CDR1H) might be evolutionarily optimized to help ensure antibody stability, thereby freeing up CDR3H to encode antigen specificity. If this hypothesis holds, it could provide a biochemical rationale for why CDR3H is more likely than CDR1H to be mutated in functioning receptors [44] and why variation in CDR3H is often sufficient to establish antigen specificity [43].

Tite-Seq can also potentially shed light on the structural basis for antibody-antigen recognition. By comparing the effects of mutations with the known antibody-fluorescein co-crystal structure [42], we identified a strong correlation between the effect that a position has on affinity and the number of molecular contacts that the residue at that position makes within the antibody. By contrast, no such correlation of expression with this number of contacts is observed. Again, we are cautious about generalizing from observations made on a single antibody. If our observation were to hold for other antibodies, however, it would suggest that the functional geometry of paratopes might be governed by networks of residues whose positions and orientations are strongly interdependent.

IV. METHODS

Tite-Seq was performed as follows. Variant 3H and 1H regions were generated using microarray-synthesized

oligos (LC Biosciences, Houston TX, USA). These were inserted into the 4-4-20 scFv of [7] using cassette-replacement restriction cloning as in [37]; see Appendix C. Yeast display experiments were performed as previously described [39] with modifications; see Appendix B. Sorted cells were regrown and bulk DNA was extracted using standard techniques, and amplicons containing the 1H and 3H variable regions were amplified using PCR and sequenced using the Illumina NextSeq platform; see Appendix D. Three replicate experiments were performed on different days. Raw sequencing data has been posted on the Sequence Read Archive under BioProject ID PRJNA344711. Low-throughput flow cytometry measurements were performed on clones randomly picked from the Tite-Seq library. Sequence data and flow cytometry data were analyzed using custom Python scripts, as described in Appendices E and F. Processed data and analysis scripts are available at github.com/jbkinney/16_titeseq.

Acknowledgements. We would like to thank Jacklyn Jansen, Amy Keating, Lothar Reich, and Bruce Stillman for helpful discussions. We would also like to thank Dane Wittrup for sharing plasmids and yeast strains. RMA, TM and AMW were supported by European Research Council Starting Grant n. 306312. JBK was supported by the Simons Center for Quantitative Biology at Cold Spring Harbor Laboratory.

The authors declare that they have no conflict of interests.

Appendix A: Schematic simulations

For panels D and E of Fig. 1, data was simulated (using Eq. 1) for two hypothetical scFvs: one similar to WT, with $K_D = 1.2 \times 10^{-9}$ M, $A = 300$, and $B = 10$, and one similar to a typical mutant, with $K_D = 10^{-6}$ M, $A = 1000$, $B = 10$. Simulated sorts were performed at the eleven antigen concentrations used in our experiments ($c = 0$ M, $10^{-9.5}$ M, $10^{-9.0}$ M, \dots , $10^{-5.0}$ M). For each clone at each antigen concentration, fluorescence signals were simulated for 1000 cells by multiplying the f quantity in Eq. 1 by a factor of $\exp(\eta)$ where η is a normally distributed random number. Fig. 1D shows the mean values of these simulated fluorescence signals. Curves of the form in Eq. 1 were fit to these data by minimizing the square deviation between predicted $\log_{10} f$ values and the \log_{10} mean of the simulated fluorescence values. The Tite-Seq measurements illustrated in Fig. 1E were simulated by sorting 1000 cells, using fluorescence values generated in the same manner as above, into three bins defined by the following fluorescence boundaries: (0, 30) for bin 0, (30, 300) for bin 1, and (300, ∞) for bin 2. The mean bin number for each clone at each antigen concentration was then computed. Curves having the form in Eq. 1 were then fit to these data by minimizing the square deviation of predicted $\log_{10} f$ values from these mean bin values.

Appendix B: Yeast display

To help ensure consistency across samples, the yeast display cultures used in our low-throughput flow cytometry measurements and in our Tite-Seq experiments were inoculated with carefully prepared frozen liquid culture inocula. Specifically, inoculation cultures were grown at 30°C in SC-trp + 2% glucose to an OD600 value between 0.9 and 1.1, then stored at -80° in aliquots containing 10% glycerol and either 0.4 ml-OD of cells (for clones) or 1 ml-OD of cells (for libraries).

The expression of yeast-displayed scFvs was induced as follows. Liquid cultures of SC-trp + 2% glucose were inoculated using single frozen inocula, yielding an approximate starting OD of 0.05. These cultures were grown at 30°C for 8 hours; the final OD of these cultures was approximately 0.7. Cells were then spun down at 1932 g for 8 minutes at 4°C, resuspended in SC-trp + 2% galactose + 0.1% glucose at 0.2 OD, and incubated for 16 hr at 20°C. We note that adding 0.1% glucose to these galactose induction cultures was essential for reliably achieving scFv expression in a large fraction of yeast cells.

Induced yeast were fluorescently labeled as follows. Galactose induction cultures were spun down and washed with ice cold TBS-BSA (0.2 mg/ml BSA, 50 mM Tris, 25 mM NaCl, pH 8). This yielded approximately 5.3 ml-OD of cells for tite-seq FACS. For antigen binding reactions, cells were then resuspended in a primary labeling reaction containing 40 ml TBS-BSA and biotinylated fluorescein (ThermoFisher B1370) at a concentration between 0 M and 10^{-5} M, then incubated with shaking for 1 hour at room temperature. Reaction volumes were large enough to ensure that $\gtrsim 10$ antigen molecules per scFv were present, assuming $\sim 10^5$ scFvs per cell [7]. Cells were then washed twice with 40 ml ice cold TBS-BSA, suspended in a secondary labeling reaction containing 1 ml ice-cold TBS-BSA and 7 μ g/ml streptavidin R-PE (ThermoFisher S866), and incubated for 30 min at 4°C while shaking. Cells were then spun down and resuspended in ice cold TBS-BSA and saved for FACS later that day. Expression labeling reactions proceeded in the same manner, except that the primary labeling reaction contained 1.4 μ g/ml rabbit anti-c-Myc antibody (Sigma-Aldrich C3956) in place of the antigen, and the secondary labeling reaction contained 0.8 μ g/ml BV421-conjugated donkey anti-rabbit antibody (BioLegend 406410) in place of streptavidin R-PE. The labeling reactions used to filter out improperly cloned scFvs (as described in Appendix C) proceeded in the same manner as the expression labeling reaction, except that 0.8 μ g/ml mouse anti-HA antibody (Roche 11583816001) was added to the primary labeling reaction, while 0.4 μ g/ml APC-conjugated anti-mouse antibody (BD Biosciences 550826) was added to the secondary labeling reaction. For clonal flow cytometry measurements, excluding secondary labeling we kept reagent and cell concentrations the same as described above, but reduced reaction volumes 27-fold. Secondary labeling reactions with streptavidin R-PE were done at

4 μ g/ml 112.5 μ l to facilitate mixing. Secondary labeling reactions with 0.8 μ g/ml BV421-conjugated donkey anti-rabbit antibody were performed in 60 μ l.

Appendix C: Cloning strategy

Amplicons containing variable CDR1H or CDR3H regions were generated as follows. An oligonucleotide library containing mutagenized 1H and 3H variable regions (see Table II) was generated by LC Sciences using microarray-based synthesis. The specific oligos used are provided at github.com/jbkinney/16_titeseq. 1H and 3H library oligos were separately amplified via PCR using primers oRAL10 and oRAR10 (for 1H) or oRAL11 and oRAR11 (for 3H). Oligos containing the WT sequence were amplified from plasmid pCT302 [7] using primers 1H2F and 1H1R (for the 1H region) or 3H1F and 3H2R (for the 3H region). Overlap-extension PCR using primers oRA10 and oRA11, one oligo library (1H or 3H) and the complementary WT oligo (3H or 1H, respectively), and plasmid pCT302, were then used to create the iRA11 amplicon library (Fig. 2 – figure supplement 1A). Note that each amplicon in this library has mutations only in the 1H variable region or in the 3H variable region, but not in both of these regions.

The pRA10 cloning vector (Fig. 2 – figure supplement 1B) was assembled using Gibson cloning [53] with template plasmids pCT302 [7] and pJK14 [37]. pCT302 is the yeast display expression plasmid containing the WT scFv. pJK14 contains a ccdB cloning cassette flanked by outward-facing BsmBI restriction sites. pRA10 closely resembles pCT302, except that it contains the ccdB cassette from pJK14 in place of the region of the scFv gene that we aimed to mutagenize. Multiple spurious BsmBI restriction sites present pCT302 were also removed in pRA10. pRA10 was propagated in *Escherichia coli* strain DB3.1, which is resistant to the CcdB toxin.

The pRA11 plasmid library (Fig. 2 – figure supplement 1C) was generated by digesting pRA10 with BsmBI, digesting the iRA11 amplicon library with BsaI, and subsequent ligation with T4 DNA ligase. Ligation reactions were desalted and transformed into DH10B *E. coli* via electroporation, yielding $\gtrsim 10^8$ transformants. The 1H and 3H libraries were cloned separately.

The pRA11 libraries were introduced into the EBY100 strain of *Saccharomyces cerevisiae* using high-efficiency LiAc transformation [54]. This yielded $\gtrsim 10^5$ transformants. To filter out yeast containing improperly cloned scFvs, we induced scFv expression, immuno-affinity labeled the HA and c-Myc epitopes on the scFv, and used FACS to recover $8 \times 10^5 - 2 \times 10^6$ cells that registered positive for both epitopes. The scFv induction and labeling procedures used to do this are described in Appendix B. 144 yeast clones were picked at random from this library and submitted for low-throughput Sanger sequencing of the 1H and 3H variable regions of the scFv. Based on preliminary Tite-Seq experiments, 19 of these clones were

then chosen for low-throughput K_D measurements.

Appendix D: Tite-Seq procedure

The inocula used for our Tite-Seq experiments comprised yeast harboring the 1H and 3H pRA11 plasmid libraries, mixed in equal proportions, and spiked at 0.625% with OPT-containing yeast (as a positive control) and at 0.625% with Δ -containing yeast (negative control). Cells were then grown, induced, and labeled with antigen at eleven different concentrations (0 M, $10^{-9.5}$ M, $10^{-9.0}$ M, \dots , $10^{-5.0}$ M) as described in Appendix B.

Each batch of labeled cells was then sorted, using FACS, over a period of approximately 20 min. During FACS, cells were first filtered based on forward scatter and side scatter to help ensure exactly one live cell per droplet. Cells passing this criterion were sorted into 4 bins based on R-PE fluorescence. The fluorescence gates used in these sorts were kept the same across all antigen concentrations (see Figs. 3, 3 – figure supplement 1, and 3 – figure supplement 2). Cells were sorted into a rounded 5 ml polypropylene tube containing 1 ml 2X YPAD media. In our separate Sort-Seq experiments assaying scFv expression levels, cells were prepared and sorted in the same way, save for the changes to the labeling reaction described in Appendix B and the use of gates on BV421 fluorescence instead of R-PE fluorescence.

Each of the 48 bins of sorted cells, as well as a sample of unsorted cells, were then deposited in 5 ml of SC-trp + 2% glucose and regrown overnight at 30°C. Approximately 25 ml-OD of cells were spun down, resuspended in a lysis reaction containing 200 μ l 0.5 mm glass beads, 200 μ l of Phenol/chloroform/isoamyl alcohol and 200 μ l of yeast lysis buffer (10 mM NaCl, 1 mM Tris, 0.1 mM EDTA, 0.2% Triton X-100, 0.1% SDS), and vortexed for 30 min. 200 μ l of water was added, cells were spun down, and the aqueous layer was extracted. Four subsequent extractions were performed, the first two using 200 μ l of Phenol/chloroform/isoamyl alcohol, the second two using 200 μ l for chloroform/isoamyl alcohol. Bulk nucleic acid was then ethanol precipitated and resuspended in 100 μ l of IDTE (Integrated DNA Technologies).

Two rounds of PCR were then performed on each of the 49 samples of bulk nucleic acid. In the first round of PCR, primers L1AF_XX and L2AF_XX were used to amplify the 1H-to-3H region and to add a bin-specific barcode (numbered XX = 01, 02, \dots , 64) on either end of the 1H-to-3H region; see Fig. 2 – figure supplement 1. To keep PCR crossover to a minimum, only 15 PCR cycles were used. These 49 PCR reactions were then pooled, purified using a QIAquick PCR purification kit (Qiagen), and used as template for second round of PCR with primers PE1v3ext and PE2v3. Again, to keep crossover to a minimum, only 25 PCR cycles were used. This PCR reaction was again purified, mixed with PhiX DNA (at \sim 25% molarity) and submitted for sequencing using the Illumina NextSeq platform.

Analysis of the resulting sequence data across the three replicate Tite-Seq experiments revealed that some of the 147 FACS bins were highly under-sampled. This under-sampling likely resulted from the use of a non-saturating number of PCR cycles. The different barcodes incorporated into the PCR primers also appear to have affected amplification efficiency to different extents. To even out the distribution of reads across bins, we selected the 27 most poorly sampled bins, re-amplified the 1H-to-3H regions in these bins using primers with different barcodes than before, and submitted the resulting amplicons for a fourth round of Illumina NextSeq sequencing.

Appendix E: Inference of K_D from Tite-Seq data

We modeled the binding titration curve of each sequence – i.e., the curve describing how mean cellular fluorescence depends on antigen concentration – using a non-cooperative Hill function. Making the dependence on scFv sequence s and antigen concentration c explicit, Eq. 1 of the main text becomes

$$f_{sc} = A_s \frac{c}{c + K_{D,s}} + B, \quad (\text{E1})$$

where f_{sc} denotes the mean fluorescence of cells carrying sequence s and labeled with antigen at concentration c . B represents the autofluorescence of cells, and was set equal to the mean fluorescence of cells labeled at 0 M antigen. A_s is the increase in fluorescence due to saturation of all surface-displayed scFvs, and $K_{D,s}$ is the dissociation constant for sequence s . We inferred A_s and $K_{D,s}$, for all sequences s , as follows.

Tite-Seq does not provide direct measurements of the fluorescence f_{sc} . Instead, we approximated this quantity using a weighted averaged over sorting bins. Specifically, we assumed that

$$\log f_{sc} \approx \sum_b p_{b|sc} F_{bc}. \quad (\text{E2})$$

Here, F_{bc} is the mean log fluorescence of the cells that were sorted at concentration c into bin b , and $p_{b|sc}$ is the probability that a cell having sequence s and labeled at concentration c , if sorted, would be found in bin b . Values for F_{bc} were computed directly from the FACS data log. The probabilities $p_{b|sc}$, by contrast, were inferred from Tite-Seq read counts. These probabilities are closely related to R_{bsc} , the number of sequence reads for each sequence s from bin b at antigen concentration c . This relationship is complicated by additional factors that arise from variability in sequencing depth from bin to bin. Moreover, because there were often a small number of reads for any particular sequence in a given bin, it was necessary in our inference procedure to treat the relationship between $p_{b|sc}$ and R_{bsc} probabilistically.

We therefore inferred values for the probabilities $p_{b|sc}$ through the following maximum likelihood procedure.

First, we assumed that the number of reads $R_{b|sc}$ is related to an “expected” number of reads $r_{b|sc}$ via a Poisson distribution. The log likelihood of observing a specific set of read counts $R_{b|sc}$ over all bins b and concentrations c for a given sequence s is therefore given by

$$L_s = \log \left[\prod_{b,c} \frac{1}{R_{b|sc}!} r_{b|sc}^{R_{b|sc}} e^{-r_{b|sc}} \right]. \quad (\text{E3})$$

The expected number of reads $r_{b|sc}$ is, in turn, related to the probability $p_{b|sc}$ via

$$r_{b|sc} = \frac{R_{bc}}{C_{bc}} C_c P_s p_{b|sc}. \quad (\text{E4})$$

Here, $R_{bc} = \sum_s R_{b|sc}$ is the total number of reads from bin b at antigen concentration c , C_{bc} is the number of cells sorted into bin b at antigen concentration c (obtained from the FACS data log), $C_c = \sum_b C_{bc}$ is the total number of cells sorted at concentration c , and P_s is the fraction of cells in the library with sequence s . The factor $R_{bc}C_c/C_{bc}$ in Eq. E4 accounts for differences in the depth with which each bin was sequenced. Note: Eq. E3 assumes that each final read arose from a different sorted cell. This assumption is clearly violated if $R_{bc} > C_{bc}$. In cases where this inequality was found to hold, we rescaled all $R_{b|sc} \rightarrow h R_{b|sc}$ where $h = C_{bc}/R_{bc}$ before undertaking further analysis.

For each sequence s , we inferred $K_{D,s}$, A_s , P_s , and all probabilities $p_{b|sc}$ by maximizing the likelihood L_s subject to the constraint that

$$\sum_b p_{b|sc} F_{bc} = \ln \left(A_s \frac{c}{c + K_{D,s}} + B \right) \quad (\text{E5})$$

at every concentration c . Note that, in this procedure, the sorting probabilities $p_{b|sc}$ are not modeled explicitly as a function of the putative mean fluorescence. Instead, they are inferred from the data along with the parameters of the non-cooperative Hill function. Doing this dispenses with the need for a detailed characterization of the noise in the Tite-Seq sorting procedure. The validity of this procedure is evinced by the analysis of simulated data, shown in Fig. 4 – figure supplement 8.

The maximum likelihood optimization problem described above was solved as follows. For each concentration c , we created a grid of 100 equally spaced points for $\ln f_{sc} \in [F_{0a}, F_{3a}]$. For each possible value of f_{sc} , we then used Nelder-Mead optimization of $p_{b|sc}$ to minimize L_s under the constraint in Eq. E5. Akima interpolation was then used to create a function of the optimized probability $\hat{p}_{b|sc}$ as a function of f_{sc} . We then scanned a 321×201 grid of values for the pair $(K_{D,s}, A_s)$ and selected the pair that minimized $L_s(K_{D,s}, A_s, \{\hat{p}_{b|sc}(K_{D,s}, A_s), P_s\}_{b,c})$. We repeated this scan by varying P_s over 95 different values. The final inferred values for $K_{D,s}$, A_s , and P_s were those so found to maximize L_s . Python code for this inference procedure is provided at github.com/jbkinney/16_titeseq.

Appendix F: Inference of K_D from flow cytometry experiments on individual clones

Low-throughput flow cytometry measurements performed on clonal cell populations were used to measure $f_{sc,\text{flow}}$, the mean fluorescence of cells carrying sequence s and labeled at antigen concentration c . As for the Tite-Seq inference procedure described in Appendix E, it was assumed that $f_{sc,\text{flow}}$ could be modeled using the non-cooperative Hill function $A_s c / (c + K_{D,s}) + B$, where A_s is the increase in mean fluorescence due to fully labeled scFvs of sequence s , $K_{D,s}$ is the corresponding dissociation constant, and B is background fluorescence. B was computed from the average fluorescein of clone s at 0 M fluorescein. A_s and $K_{D,s}$ were then inferred by minimizing the square deviation between measured $\ln f_{sc,\text{flow}}$ values and log Hill function predictions, i.e.,

$$\sum_c \left[\ln f_{sc,\text{flow}} - \ln \left(A_s \frac{c}{c + K_{D,s}} + B \right) \right]^2. \quad (\text{F1})$$

This optimization procedure was performed using a grid search algorithm in which $K_{D,s}$ was restricted to the interval $[10^{-10}M, 10^{-3}M]$ and A_s was restricted to the interval $[\bar{A}, 100\bar{A}]$ where \bar{A} denotes the average range of fluorescence values over the 4-to-8 clones assayed per flow cytometry session.

Appendix G: Realistic Tite-Seq simulations

In order to test our analysis pipeline, we simulated realistic Tite-Seq data and analyzed it with the same scripts that we used on real data. For each sequence s , a $K_{D,s}$ value was randomly drawn from the interval $[10^{-10}M, 10^{-4}M]$ using a uniform distribution in log space, and an A_s value was drawn uniformly from a uniform distribution spanning the bulk of experimentally observed A values. At each of the eleven antigen concentrations c , we then modeled the distribution of cellular fluorescence using a Gaussian Mixture Model (GMM) in log space. Specifically, letting x denote \log_{10} cellular fluorescence values, we assumed that the probability density describing x to be

$$P_{sc}(x) = \frac{\alpha}{\sigma_0 \sqrt{2\pi}} e^{-\frac{(x-\mu_0)^2}{2\sigma_0^2}} + \frac{1-\alpha}{\sigma_1 \sqrt{2\pi}} e^{-\frac{(x-\mu_{sc})^2}{2\sigma_1^2}}, \quad (\text{G1})$$

where $\alpha = 0.2$ is the fraction of non-expressing cells, $\mu_0 = 4.77$ and $\sigma_0 = 1$ are the mean and standard deviation of x values for dark cells, and $\sigma_1 = 0.5$ is the standard deviation of x values for scFv-expressing cells. The mean x value of scFv-expressing cells, here denoted μ_{cs} , was chosen so that the population average of x is given by the Hill function in Eq. E1, i.e., so that

$$\langle x \rangle_{P_{sc}} = A_s \frac{c}{c + K_{D,s}} + B. \quad (\text{G2})$$

The left hand side of Eq. G2 can be computed analytically using Eq. G1. Doing this and solving for μ_{sc} gives

$$\mu_{sc} = \ln \left(\frac{A_s c}{c + K_{D,s}} + B - \alpha e^{\mu_0 + \sigma_0^2/2} \right) - \ln(1 - \alpha) - \frac{\sigma_1^2}{2}; \quad (\text{G3})$$

this is the specific formula we used to compute μ_{sc} as a function of c , A_s , and $K_{D,s}$. Next we computed exact $p_{b|sc}$ values using

$$p_{b|sc} = \int_{b^-}^{b^+} dx P_{sc}(x), \quad (\text{G4})$$

where b^+ and b^- are the upper and lower fluorescence bounds used for bin b in our Tite-Seq experiment (replicate number 1). These values were then used to draw read counts $R_{b|sc}$ for each sequence s values via

$$R_{b|sc} \sim \text{Bionomial}(n = k_s R_{bc}, p = p_{b|sc}), \quad (\text{G5})$$

where k_s is a random variable, uniformly distributed on a log scale between 0.01 and 100, used to represent noise due to PCR jackpotting. Data thus simulated for WT values of K_D and A are shown in Fig. 4 – figure supplement 7.

-
- [1] A. C. Chan and P. J. Carter, Nat. Rev. Immunol. **10**, 301 (2010).
 - [2] K. P. Murphy, P. Travers, and M. Walport, *Janeway's Immunobiology*, 7th ed. (Garland Science, 2008).
 - [3] G. P. Smith, Science **228**, 1315 (1985).
 - [4] T. J. Vaughan, A. J. Williams, K. Pritchard, J. K. Osbourn, A. R. Pope, J. C. Earnshaw, J. McCafferty, R. A. Hodits, J. Wilton, and K. S. Johnson, Nat. Biotechnol. **14**, 309 (1996).
 - [5] T. Schirrmann, T. Meyer, M. Schütte, A. Frenzel, and M. Hust, Molecules **16**, 412 (2011).
 - [6] Y. Fujino, R. Fujita, K. Wada, K. Fujishige, T. Kanamori, L. Hunt, Y. Shimizu, and T. Ueda, Biochemical and Biophysical Research Communications **428**, 395 (2012).
 - [7] E. T. Boder and K. D. Wittrup, Nat. Biotechnol. **15**, 553 (1997).
 - [8] S. A. Gai and K. D. Wittrup, Curr. Opin. Struc. Biol. **17**, 467 (2007).
 - [9] C. M. Forsyth, V. Juan, Y. Akamatsu, R. B. DuBridge, M. Doan, A. V. Ivanov, Z. Ma, D. Polakoff, J. Razo, K. Wilson, and D. B. Powers, mAbs **5**, 523 (2013).
 - [10] S. D. Boyd, E. L. Marshall, J. D. Merker, J. M. Maniar, L. N. Zhang, B. Sahaf, C. D. Jones, B. B. Simen, B. Hanczaruk, K. D. Nguyen, K. C. Nadeau, M. Egholm, D. B. Miklos, J. L. Zehnder, and A. Z. Fire, Sci Transl Med **1**, 12ra23 (2009).
 - [11] J. A. Weinstein, N. Jiang, R. A. White, D. S. Fisher, and S. R. Quake, Science **324**, 807 (2009).
 - [12] H. S. Robins, P. V. Campregher, S. K. Srivastava, A. Wachter, C. J. Turtle, O. Kahsai, S. R. Riddell, E. H. Warren, and C. S. Carlson, Blood **114**, 4099 (2009).
 - [13] H. S. Robins, S. K. Srivastava, P. V. Campregher, C. J. Turtle, J. Andriesen, S. R. Riddell, C. S. Carlson, and E. H. Warren, Sci. Transl. Med. **2**, 47ra64 (2010).
 - [14] T. Mora, A. M. Walczak, W. Bialek, and C. G. Callan, Proc. Natl. Acad. Sci. USA **107**, 5405 (2010).
 - [15] V. Venturi, M. F. Quigley, H. Y. Greenaway, P. C. Ng, Z. S. Ende, T. McIntosh, T. E. Asher, J. R. Almeida, S. Levy, D. A. Price, M. P. Davenport, and D. C. Douek, J. Immunol. **186**, 4285 (2011).
 - [16] A. Murugan, T. Mora, A. M. Walczak, and C. G. Callan, Proc. Natl. Acad. Sci. USA **109**, 16161 (2012).
 - [17] I. V. Zvyagin, M. V. Pogorelyy, M. E. Ivanova, E. A. Komech, M. Shugay, D. A. Bolotin, A. A. Shelenkov, A. A. Kurnosov, D. B. Staroverov, D. M. Chudakov, Y. B. Lebedev, and I. Z. Mamedov, Proc. Natl. Acad. Sci. USA **111**, 5980 (2014).
 - [18] Y. Elhanati, A. Murugan, C. G. Callan, T. Mora, and A. M. Walczak, Proc. Natl. Acad. Sci. USA **111**, 9875 (2014).
 - [19] Q. Qi, Y. Liu, Y. Cheng, J. Glanville, D. Zhang, J.-Y. Lee, R. A. Olshen, C. M. Weyand, S. D. Boyd, and J. J. Goronzy, Proc Natl Acad Sci USA **111**, 13139 (2014).
 - [20] N. Thomas, K. Best, M. Cinelli, S. Reich-Zeliger, H. Gal, E. Shifrut, A. Madi, N. Friedman, J. Shawe-Taylor, and B. Chain, Bioinformatics **30**, 3181 (2014).
 - [21] Y. Elhanati, Z. Sethna, Q. Marcou, C. G. Callan Jr, T. Mora, and A. M. Walczak, arXiv:1212.3647 [q-bio.QM] (2015), 1502.03136v2.
 - [22] A. Madi, E. Shifrut, S. Reich-Zeliger, H. Gal, K. Best, W. Ndifon, B. Chain, I. R. Cohen, and N. Friedman, Genome Res. **24**, 1603 (2014).
 - [23] P. Parameswaran, Y. Liu, K. M. Roskin, K. K. Jackson, V. P. Dixit, J.-Y. Lee, K. L. Artilles, S. Zompi, M. J. Vargas, B. B. Simen, *et al.*, Cell Host Microbe **13**, 691 (2013).
 - [24] N. Jiang, J. He, J. a. Weinstein, L. Penland, S. Sasaki, X.-S. He, C. L. Dekker, N.-Y. Zheng, M. Huang, M. Sullivan, P. C. Wilson, H. B. Greenberg, M. M. Davis, D. S. Fisher, and S. R. Quake, Sci. Transl. Med. **5**, 171ra19 (2013).
 - [25] C. Vollmers, R. V. Sit, J. a. Weinstein, C. L. Dekker, and S. R. Quake, Proc. Natl. Acad. Sci. USA **110**, 13463 (2013).
 - [26] U. Laserson, F. Vigneault, D. Gadala-Maria, G. Yaari, M. Uduman, J. a. Vander Heiden, W. Kelton, S. Taek Jung, Y. Liu, J. Laserson, R. Chari, J.-H. Lee, I. Bachelet, B. Hickey, E. Lieberman-Aiden, B. Hanczaruk, B. B. Simen, M. Egholm, D. Koller, G. Georgiou, S. H. Kleinstein, and G. M. Church, Proc. Natl. Acad. Sci. USA **111**, 4928 (2014).
 - [27] J. D. Galson, A. J. Pollard, J. Trüch, and D. F. Kelly, Trends Immunol. **35**, 319 (2014).
 - [28] C. Wang, Y. Liu, M. M. Cavanagh, S. Le Saux, Q. Qi, K. M. Roskin, T. J. Looney, J.-Y. Lee, V. Dixit, C. L. Dekker, G. E. Swan, J. J. Goronzy, and S. D. Boyd, Proc Natl Acad Sci USA **112**, 500 (2015).
 - [29] D. M. Fowler and S. Fields, Nat. Methods **11**, 801 (2014).
 - [30] D. M. Fowler, C. L. Araya, S. J. Fleishman, E. H. Kellogg, J. J. Stephany, D. Baker, and S. Fields, Nat. Methods **7**, 741 (2010).

- [31] L. . Reich, S. Dutta, and A. E. Keating, J. Mol. Biol. **427**, 2135 (2014).
- [32] C. A. Kowalsky, M. S. Faber, A. Nath, H. E. Dann, V. W. Kelly, L. Liu, P. Shanker, E. K. Wagner, J. A. Maynard, C. Chan, and T. A. Whitehead, J Biol Chem **290**, 26457 (2015).
- [33] K. M. Doolan and D. W. Colby, J Mol Biol **427**, 328 (2015).
- [34] T. Van Blarcom, A. Rossi, D. Foletti, P. Sundar, S. Pitts, C. Bee, J. Melton Witt, Z. Melton, A. Hasa-Moreno, L. Shaughnessy, D. Telman, L. Zhao, W. L. Cheung, J. Berka, W. Zhai, P. Strop, J. Chaparro-Riggers, D. L. Shelton, J. Pons, and A. Rajpal, J Mol Biol **427**, 1513 (2015).
- [35] L. L. Reich, S. Dutta, and A. E. Keating, J. Mol. Biol. **427**, 2135 (2015).
- [36] C. A. Kowalsky and T. A. Whitehead, Proteins: Structure, Function, and Bioinformatics **84**, 1914 (2016).
- [37] J. B. Kinney, A. Murugan, C. G. Callan, and E. C. Cox, Proc. Natl. Acad. Sci. USA **107**, 9158 (2010).
- [38] We note that Kowalsky et al. [32] have described yeast display DMS experiments performed at multiple concentrations. These data, however, were not used to reconstruct titration curves or infer quantitative K_D values.
- [39] E. T. Boder, K. S. Midelfort, and K. D. Wittrup, Proc. Natl. Acad. Sci. USA **97**, 10701 (2000).
- [40] K. S. Midelfort, H. H. Hernandez, S. M. Lippow, B. Tidor, C. L. Drennan, and K. D. Wittrup, J. Mol. Biol. **343**, 685 (2004).
- [41] J. J. VanAntwerp and K. D. Wittrup, Biotechnol. Prog. **16**, 31 (2000).
- [42] M. Whitlow, A. J. Howard, J. F. Wood, E. W. Voss, and K. D. Hardman, Protein Eng. **8**, 749 (1995).
- [43] J. L. Xu and M. M. Davis, Immunity **13**, 37 (2000).
- [44] G. Liberman, J. Benichou, L. Tsaban, J. Glanville, and Y. Louzoun, Front. Immunol. **4**, 274 (2013).
- [45] W. T. Ireland and J. B. Kinney, bioRxiv doi: <http://dx.doi.org/10.1101/054676> (2016).
- [46] B. Foat, A. Morozov, and H. Bussemaker, Bioinformatics **22**, e141 (2006).
- [47] F. D. Batista and M. S. Neuberger, *Immunity*, Immunity **8**, 751 (1998).
- [48] J. Foote and H. N. Eisen, Proc. Natl. Acad. Sci. USA **92**, 1254 (1995).
- [49] H.-P. Roost, M. F. Bachmann, A. Haag, U. Kalinke, V. Pliska, H. Hengartner, and R. M. Zinkernagel, Proc. Natl. Acad. Sci. USA **92**, 1257 (1995).
- [50] M. L. Burns, T. M. Malott, K. J. Metcalf, B. J. Hackel, J. R. Chan, and E. V. Shusta, Appl Environ Microbiol **80**, 5732 (2014).
- [51] B. Liu, ed., *Yeast Surface Display*, Methods, Protocols, and Applications (Humana Press, 2015).
- [52] E. V. Shusta, M. C. Kieke, E. Parke, D. M. Kranz, and K. D. Wittrup, J. Mol. Biol. **292**, 949 (1999).
- [53] D. G. Gibson, L. Young, R.-Y. Chuang, J. C. Venter, C. A. Hutchison, and H. O. Smith, Nat Methods **6**, 343 (2009).
- [54] R. D. Gietz and R. H. Schiestl, Nat Protoc **2**, 31 (2007).

Name	1H variable region	3H variable region	no. replicates (flow)	K_D [M] (flow)	K_D [M] (Tite-Seq)
OPT	TFghYWMNWV	GasYGMeYlG	3	$\lesssim 10^{-9.5}$	$\lesssim 10^{-9.5}$
C107	TFSDYWMNWV	GaYYGMDYWG	3	$10^{-9.28 \pm 0.04}$	$10^{-9.18 \pm 0.11}$
C112	TFSDYWMNWV	GSYYGMDYcG	3	$10^{-8.95 \pm 0.07}$	$10^{-9.19 \pm 0.14}$
WT	TFSDYWMNWV	GSYYGMDYWG	10	$10^{-8.61 \pm 0.07}$	$10^{-8.92 \pm 0.10}$
C144	vFSDYWMNWV	GSYYGMDYWG	3	$10^{-8.57 \pm 0.03}$	$10^{-8.86 \pm 0.04}$
C133	aFSDYWMNWV	GSYYGMDYWG	3	$10^{-8.55 \pm 0.06}$	$10^{-8.62 \pm 0.09}$
C132	TFmDYWlNWV	GSYYGMDYWG	3	$10^{-8.48 \pm 0.08}$	$10^{-8.38 \pm 0.29}$
C94	TFSDYWMNWV	GSYYGMDsWG	3	$10^{-8.46 \pm 0.06}$	$10^{-8.50 \pm 0.04}$
C5	TFSDYWiNWV	GSYYGMDYWG	3	$10^{-8.34 \pm 0.10}$	$10^{-8.55 \pm 0.09}$
C93	TFSDYWMNWV	GSYrGMDYWG	3	$10^{-7.35 \pm 0.08}$	$10^{-7.60 \pm 0.70}$
C39	TFSDYWMNWV	GSYYGMDYWa	3	$10^{-7.08 \pm 0.20}$	$10^{-7.28 \pm 0.17}$
C102	TFSDYWMNWV	sSkYGMDYWG	3	$10^{-5.76 \pm 0.16}$	$10^{-7.25 \pm 0.60}$
C22	ssSDYWMNWV	GSYYGMDYWG	3	$10^{-5.69 \pm 0.31}$	$10^{-7.53 \pm 0.07}$
C7	hFSDYWMNWl	GSYYGMDYWG	3	$10^{-5.53 \pm 0.18}$	$10^{-5.39 \pm 0.18}$
C45	TFSDYWMNWV	GSYdGnDYWG	3	$10^{-5.40 \pm 0.24}$	$\gtrsim 10^{-5.0}$
C103	TFSDYWMNWV	GSYYGMDlWG	3	$10^{-5.15 \pm 0.47}$	$10^{-5.44 \pm 0.55}$
C3	TFSDYWMsWV	GSYYGMDYWG	3	$\gtrsim 10^{-5.0}$	$\gtrsim 10^{-5.0}$
C18	TFSDYsMNWV	GSYYGMDYWG	3	$\gtrsim 10^{-5.0}$	$\gtrsim 10^{-5.0}$
Δ	—	—	12	$\gtrsim 10^{-5.0}$	$\gtrsim 10^{-5.0}$

TABLE I. **Clones measured using flow cytometry and Tite-Seq.** List of scFv clones, ordered by their flow-cytometry-measured K_D values. With the exception of OPT and Δ , these clones differed from WT only in their 1H and 3H variable regions. WT amino acids within these regions are capitalized; variant amino acids are shown in lower case. No sequence is shown for Δ because this clone contained a large deletion, making identification of the 1H and 3H variable regions meaningless. K_D values saturating our lower detection limit of $10^{-9.5}$ M or upper detection limit of $10^{-5.0}$ are written with a \lesssim or \gtrsim sign to emphasize the uncertainty in these measurements. Tite-Seq K_D values indicate mean and standard errors computed across the three replicate Tite-Seq experiments; they are not averaged across synonymous variants.

Name	Sequence
1H library	GTGTTGCCTCTGGATTCA ACTTTTAGTGACTACTGGATGAAC TGGGTCCGCCAGTCTCCAGA
3H library	GTGACTGAGGTTCCCTTG ACCCAGTAGTCCATACC ATAGTAAGA ACCC GTACAGTAATAGATACCCAT
oRAL10	TTCTGAGGAGACGGTGACTGAGGTTCCCTTG
oRAR10	TGAAGACATGGGTATCTATTACTGTACG
oRAL11	CAGTCCTTTCTCTGGAGACTGGCG
oRAR11	ATGAAACTCTCCTGTGTTGCCTCTGGATTG
3H1F	TTCTGAGGAGACGGTGACT
3H2R	TGAAGACATGGGTATCTATTACTGTAC
1H2F	CAGTCCTTTCTCTGGAGACTG
1H1R	ATGAAACTCTCCTGTGTTGCCT
oRA10	GCATATCTAAGGTCTCGTTCTGAGGAGACGGTGAC
oRA11	GCCGATTGTTGGTCTCCATGAAACTCTCCTGTGTTGC
PE1v3ext	AATGATACGGCGACCACCGAGATCTACACTCTTCCCTACACGACG
PE2v3	AAGCAGAAGACGGCATACGAGATCGGTCTCGGCATTCTGCT
L1AF_XX	ACACTCTTTCCCTACACGACGCTCTTCCGATCT [XX] AGTCTTCTTCAGAAATAAGC
L1AR_XX	CTCGGCATTCTGCTGAACCGCTCTTCCGATCT [XX] GCTTGGTGCAACCTG

TABLE II. **Primers.** Oligonucleotide sequences are written 5' to 3'. Bold sequences indicate variable regions. The “1H library” and “3H library” primers respectively contained the 1H and 3H variable regions (bold) analyzed in this paper. These primer libraries were synthesized by LC Biosciences using microarray-based DNA synthesis. All other primers were ordered from Integrated DNA Technologies. The “[XX]” portion of L1AF_XX and L1AR_XX indicates the location of each of 64 different barcodes (i.e., $XX = 01, 02, \dots, 64$), which ranged in length from 7 bp to 10 bp and which differed from each other by at least two substitution mutations.

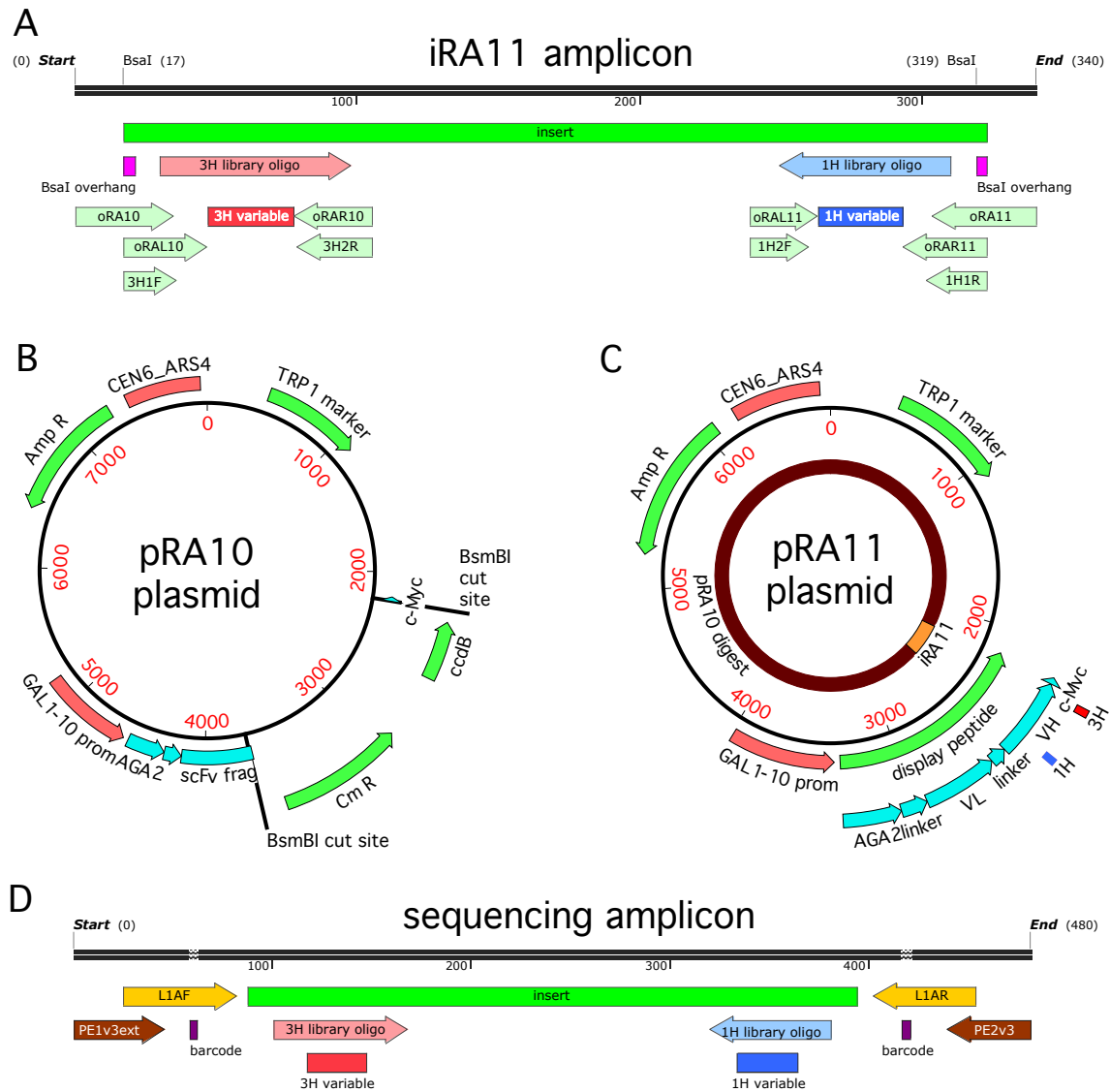


FIG. 2 – figure supplement 1. **Cloning strategy.** (A) The iRA11 amplicon library, which was prepared from microarray-synthesized oligos containing variant CDR1H or variant CDR3H regions. This amplicon is flanked by inward-facing BsaI restriction sites. (B) The pRA10 cloning vector, which contains the ccdB selection gene within a cassette flanked by outward-facing BsmBI restriction sites. (C) The pRA11 plasmid library, which was cloned by ligating BsaI-digested iRA11 amplicons and BsmBI-digest pRA10 vector. (D) The sequencing amplicon that was amplified from sorted cells after Tite-Seq and Sort-Seq experiments and submitted for ultra-high-throughput DNA sequencing. Appendix C provides more details about iRA11 amplicons, the pRA10 vector, and the pRA11 plasmid library. Appendix D provides more information about the creation of sequencing amplicons.

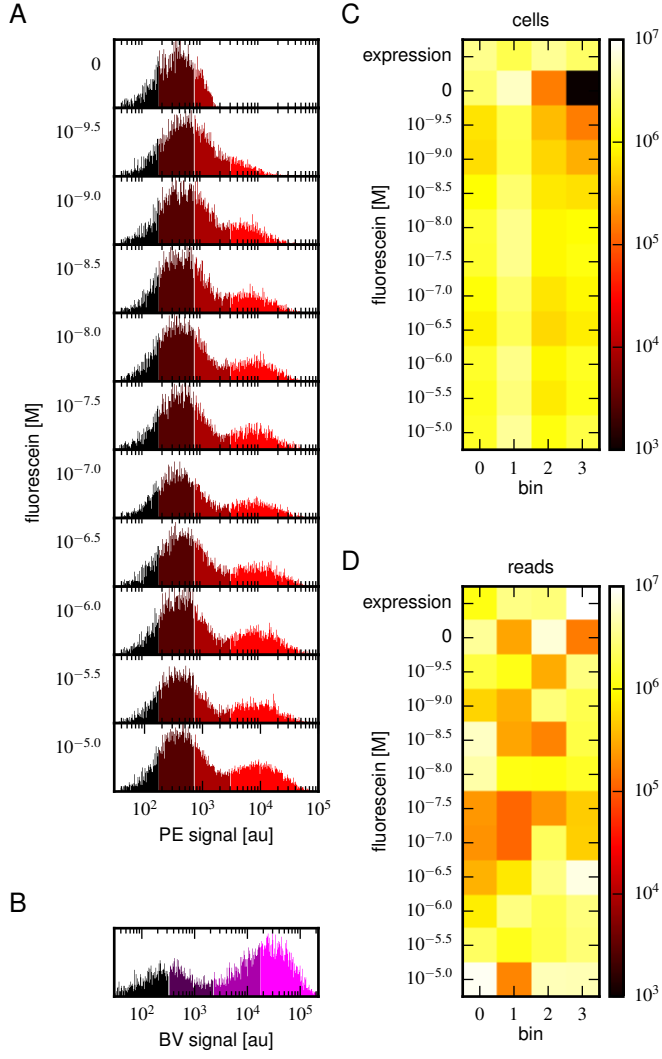


FIG. 3 – figure supplement 1. **Tite-Seq experiment, replicate 2.** Analog of Fig. 3 in the main text, but for the replicate 2 Tite-Seq experiment.

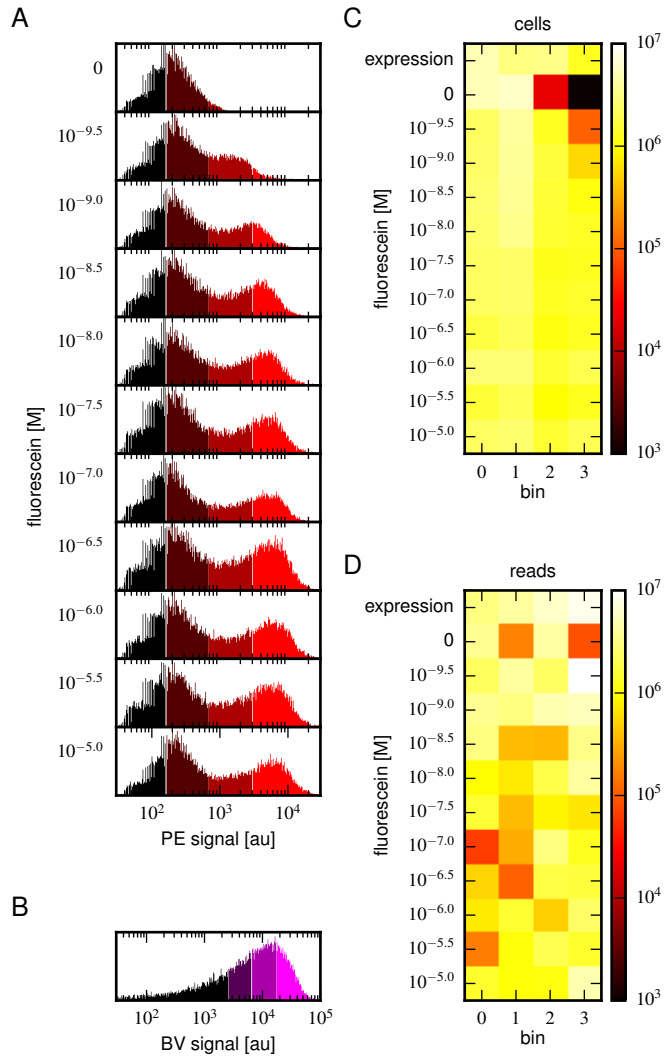


FIG. 3 – figure supplement 2. **Tite-Seq experiment, replicate 3.** Analog of Fig. 3 in the main text, but for the replicate 3 Tite-Seq experiment.

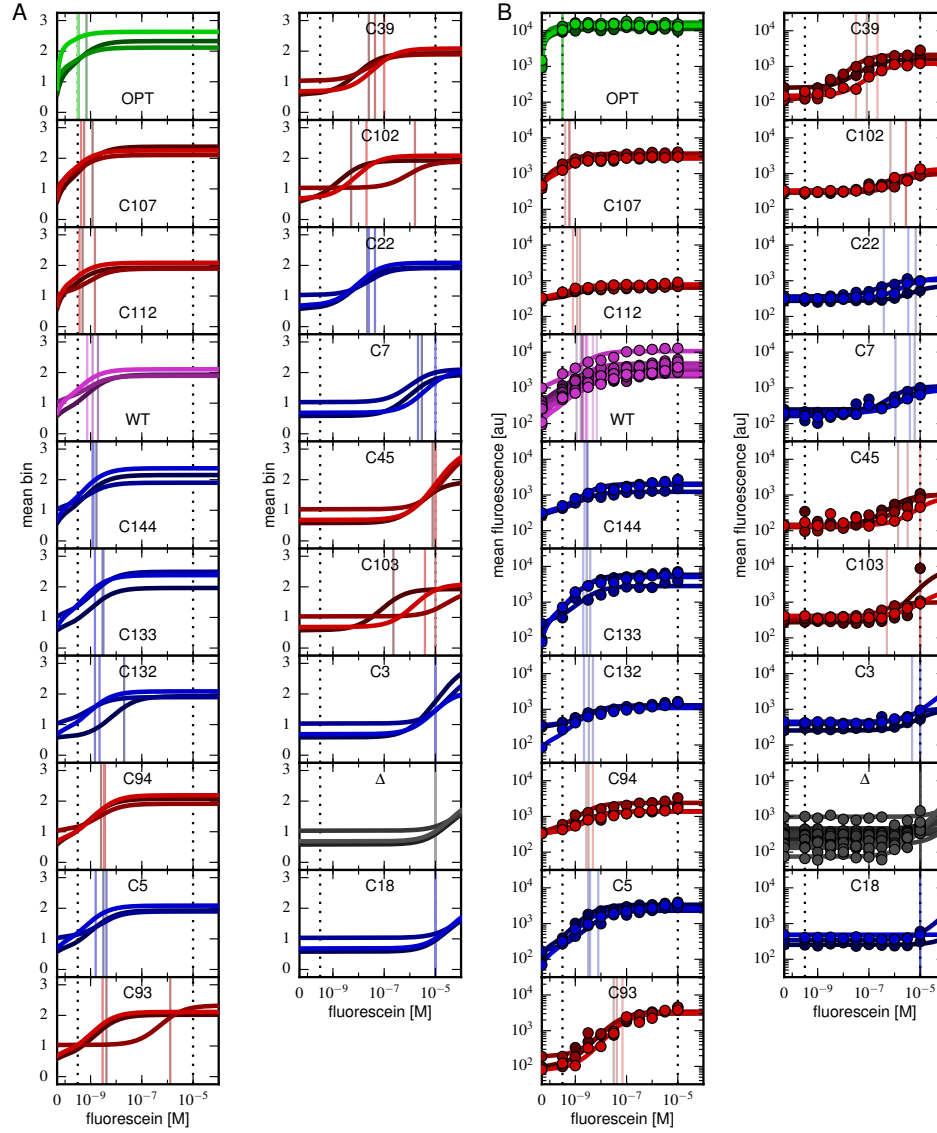


FIG. 4 – figure supplement 1. **Binding curves for all clones.** Binding curves, measured using (A) Tite-Seq or (B) flow cytometry, for all clones analyzed in this paper and described in Table I. Plots are drawn as in Fig. 4, panels A and B.

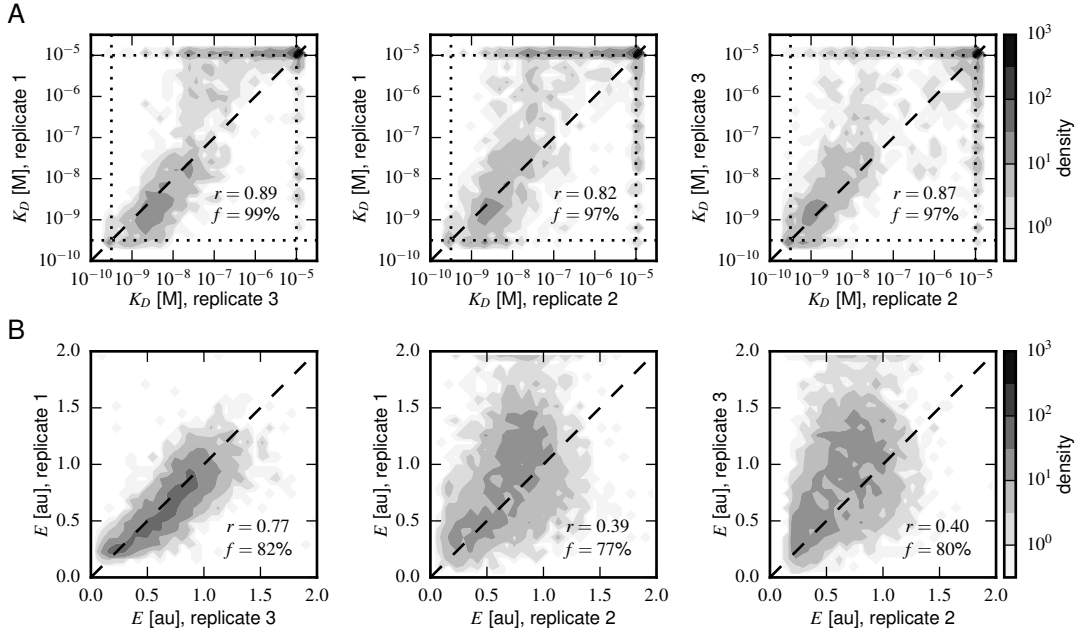


FIG. 4 – figure supplement 2. **Concordance between replicate experiments.** Density plots of (A) Tite-Seq-measured K_D values and (B) Sort-Seq-measured E values between all pairs of replicate experiments. Measurements for these quantities that were judged to be of low precision due to low sequence counts are not plotted. f indicates the percentage of total assayed sequences plotted; r is the Pearson correlation and includes clonal measurements outside the boundaries of our measurable ranges ($10^{-9.5} - 10^{-5}$ M for K_D , 0-2 for expression). Clones outside of these ranges were given values at the closest boundary.

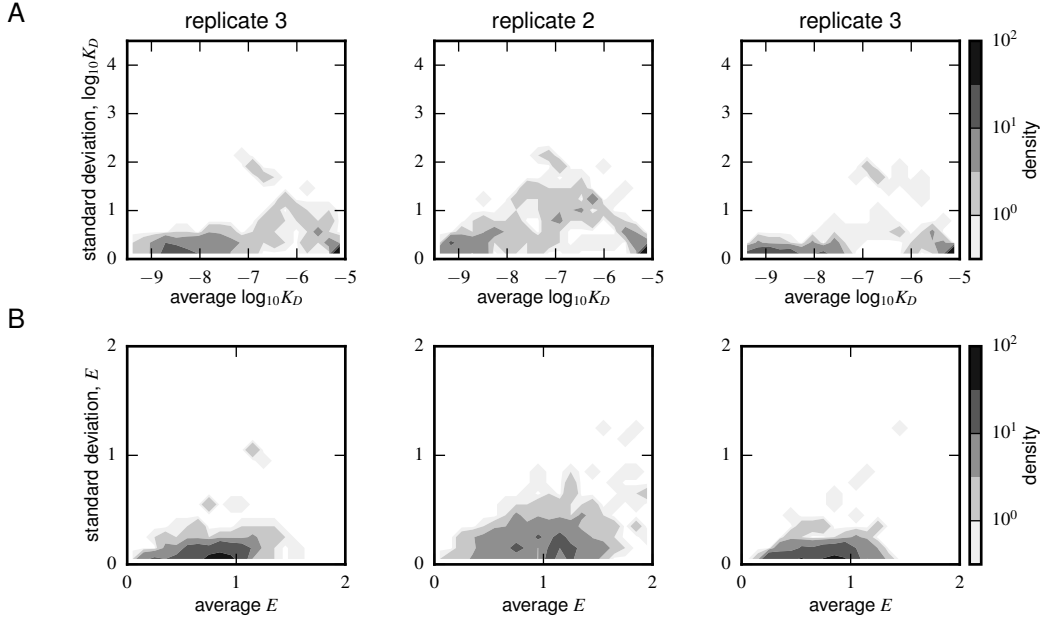


FIG. 4 – figure supplement 3. **Error estimates from synonymous mutants.** Density plots for (A) Tite-Seq-measured $\log_{10} K_D$ standard deviation and average $\log_{10} K_D$ and (B) Sort-Seq-measured E standard deviation and average E are shown for each scFv sequence with more than 1 synonymous mutant for each of the replicate experiments. The K_D error peaked between $10^{-7} - 10^{-6}$ M. The expression error peaked at or above WT expression (i.e. 1) levels.

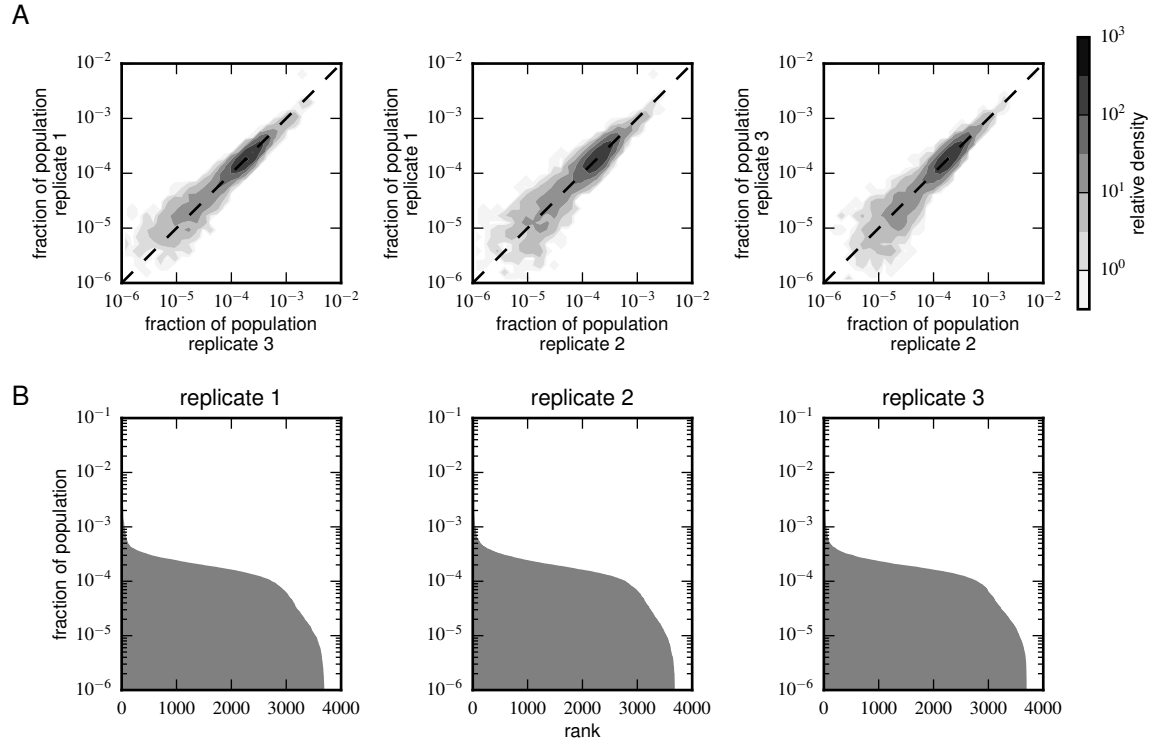


FIG. 4 – figure supplement 4. **Composition of scFv libraries.** (A) Comparison of library composition between all pairs of replicate experiments. (B) Zipf plots showing the library composition in each replicate experiment. In both panels, the prevalence of each scFv sequence in each replicate experiment was determined as part of the Tite-Seq curve fitting procedure, as described in Appendix E.

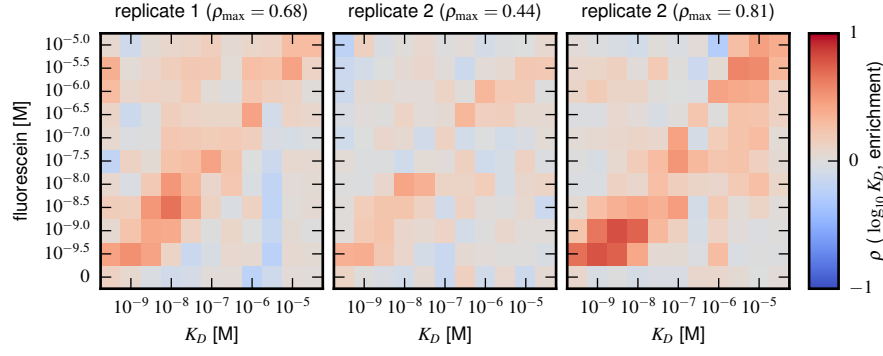


FIG. 4 – figure supplement 5. **Sort-Seq enrichment correlates poorly with Tite-Seq-measured affinity.** To assess how well simple enrichment calculations might reproduce the K_D values measured by Tite-Seq, we did the following calculation. For each of the two libraries (1H and 3H), we partitioned scFvs into seven groups based on their measured K_D s (columns). For each group at each antigen concentration (rows), we then computed the enrichment of each scFv in the high PE bins (bins 2,3) relative to the low PE bins (bins 0,1). In these enrichment calculations, the number of counts in each bin was re-weighted to accurately reflect the fraction of library cells falling within the fluorescence range of that bin. This figure shows the resulting Spearman rank correlation (ρ) between enrichment and $\log K_D$ values computed for each scFv group at each antigen concentration. In both libraries, we see that correlation values above background (which can be assessed from the values in the 0 M fluorescein row) only occur close to the diagonal, i.e., when K_D is close to the fluorescein concentration used.

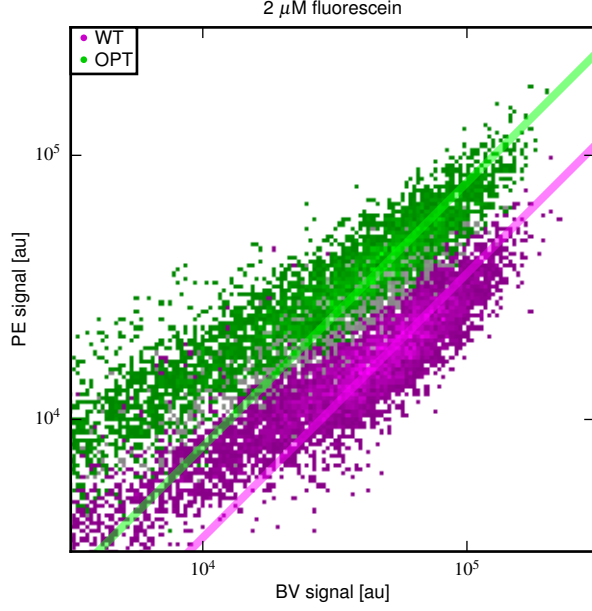


FIG. 4 – figure supplement 6. **Differing specific activities of OPT and WT.** 2D flow cytometry histograms showing both OPT- and WT-expressing cells labeled with PE and BV after incubation at 2 μ M fluorescein. At this fluorescein concentration, nearly all functional WT and OPT scFvs are bound. Regression lines (fixed to have slope 1) were fit to data points with BV signal between $10^{4.5}$ and 10^5 . The vertical shift of the OPT data relative to the WT data indicates a factor of 2.03 ± 0.07 difference (computed from four replicate experiments) in the amount labeled antigen. This difference is not due to a difference in the number of surface-displayed scFvs, as this would cause the OPT and WT clouds to lie along the same diagonal. Rather, this difference between WT and OPT is due to variation in specific activity.

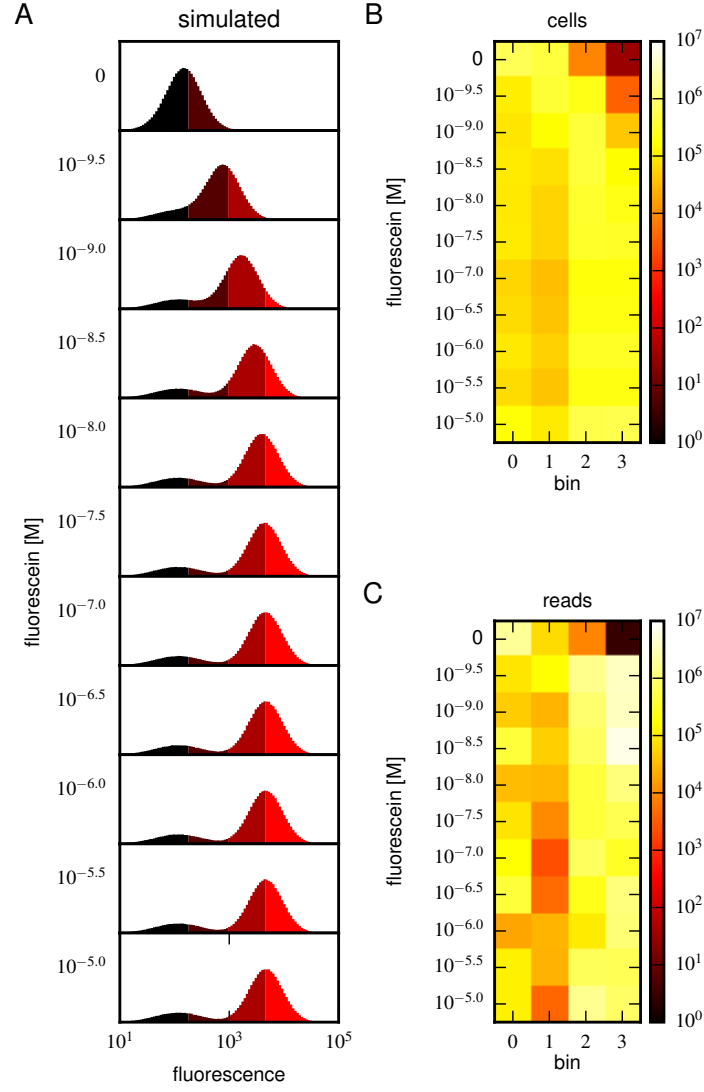


FIG. 4 – figure supplement 7. **Realistic Tite-Seq simulations.** Realistic Tite-Seq data were simulated separately for each distinct pair of affinity (K_D) and amplitude (A) values, as described in Appendix G. This figure shows simulated data, akin to the data displayed in Fig. 4 – figure supplement 6, for WT values of K_D and A .

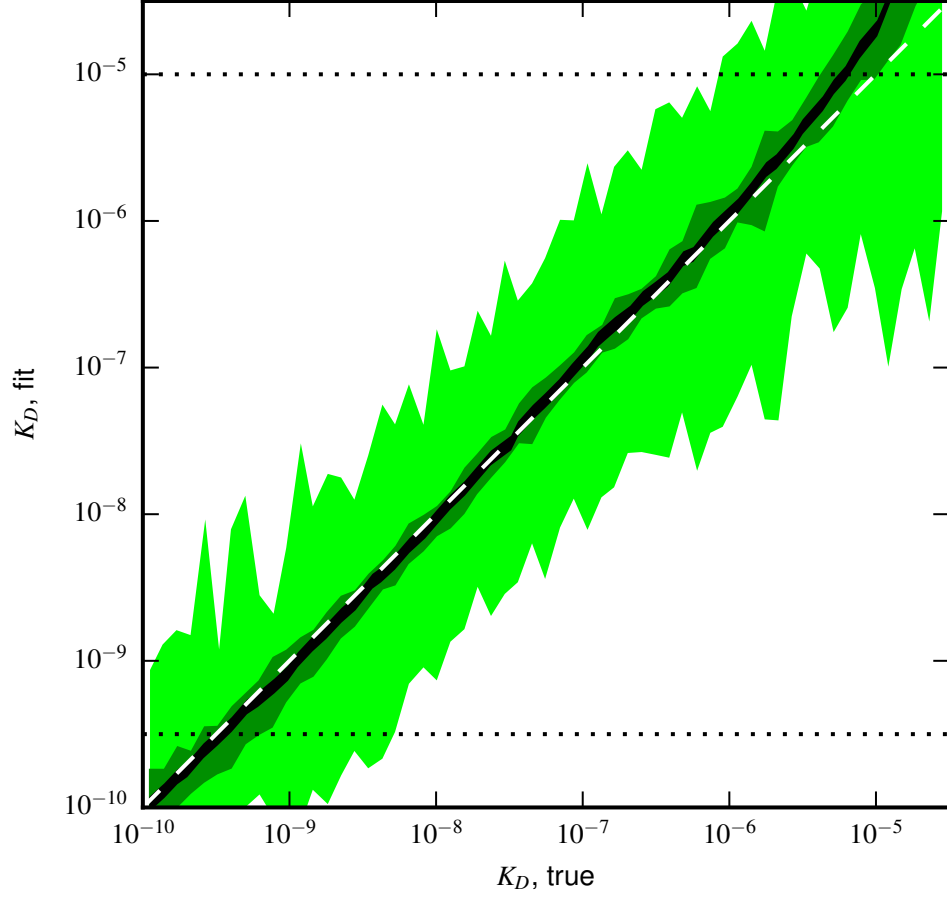


FIG. 4 – figure supplement 8. **Validation of analysis pipeline.** K_D values were inferred for Tite-Seq data simulated using (green) the same number of cells, (light green) 10^{-3} times as many cells, or (black) 10^4 times as many sorted cells as in our experiments. Areas indicate approximately plus or minus one standard deviation in the fitted K_D values obtained for each true K_D value.

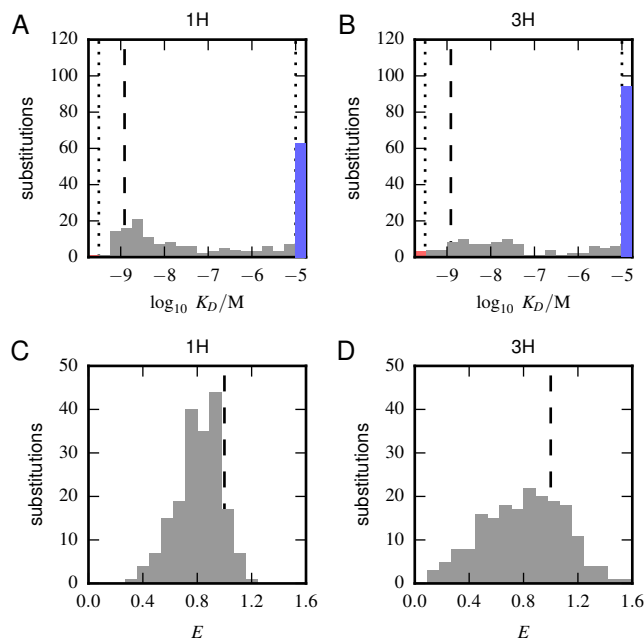


FIG. 5 – figure supplement 1. **Histograms of substitution effects on affinity and expression.** (A,B) Histogram showing the K_D values measured for all substitution mutations in the 1H (A) and 3H (B) libraries. Note that these are the values plotted in panels A and B of Fig. 5, except that the WT K_D value is not included. Dashed lines indicate the K_D of the WT scFv; dotted lines indicate thresholds just within our detection boundaries, $10^{-9.49}$ M and $10^{-5.01}$ M, while the colored bars outside this interval indicate the number of substitution mutations with K_D above (blue) and below (red) this range. (C,D) Histogram of E values for all single-substitution variants in the 1H (C) or 3H (D) libraries. These values, save those of the WT scFv, are plotted in panels C and D of Fig. 5. Dashed lines indicate the WT expression level of $E = 1.0$.

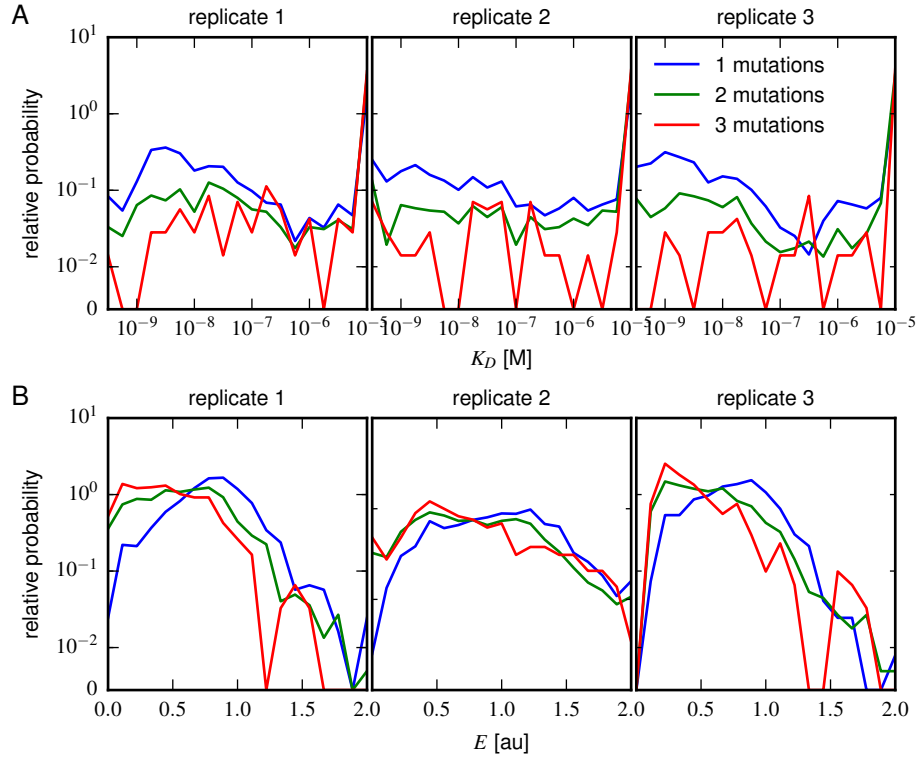


FIG. 5 – figure supplement 2. **Effects of multi-point mutations on affinity and expression.** The effect of 1, 2, or 3 mutations on (A) Tite-Seq-measured K_D values or (B) Sort-Seq-measured E values. Plots show the relative probability density (over 30 bins along the K_D or E axes) observed for variants in each class.