

Highly multiplexed design of an allosteric transcription factor to sense new ligands

Received: 18 May 2024

Accepted: 5 November 2024

Published online: 19 November 2024



Kyle K. Nishikawa^{1,7}, Jackie Chen^{1,7}, Justin F. Acheson¹, Svetlana V. Harbaugh², Phil Huss¹, Max Frenkel¹, Nathan Novy¹, Hailey R. Sieren^{1,3}, Ella C. Lodewyk^{1,3}, Daniel H. Lee^{1,3}, Jorge L. Chávez², Brian G. Fox^{1,4} & Srivatsan Raman^{1,4,5,6} ✉

Allosteric transcription factors (aTF) regulate gene expression through conformational changes induced by small molecule binding. Although widely used as biosensors, aTFs have proven challenging to design for detecting new molecules because mutation of ligand-binding residues often disrupts allostery. Here, we develop Sensor-seq, a high-throughput platform to design and identify aTF biosensors that bind to non-native ligands. We screen a library of 17,737 variants of the aTF TtgR, a regulator of a multidrug exporter, against six non-native ligands of diverse chemical structures – four derivatives of the cancer therapeutic tamoxifen, the antimalarial drug quinine, and the opiate analog naltrexone – as well as two native flavonoid ligands, naringenin and phloretin. Sensor-seq identifies biosensors for each of these ligands with high dynamic range and diverse specificity profiles. The structure of a naltrexone-bound design shows shape-complementary methionine-aromatic interactions driving ligand specificity. To demonstrate practical utility, we develop cell-free detection systems for naltrexone and quinine. Sensor-seq enables rapid and scalable design of new biosensors, overcoming constraints of natural biosensors.

Allosteric proteins pervade biology. They drive virtually all cellular processes as sensors, regulators, enzymes, and signaling proteins. Allostery describes the regulation of protein activity from a distance, where a perturbation at the allosteric site affects activity at a distant site within the same protein. Designing allosteric proteins is key to our efforts to engineer biology. One such class of allosteric proteins is allosteric transcription factors (aTFs). aTFs are universal molecular switches, governing gene expression in response to diverse cellular and environmental cues^{1–3}. These regulators modulate transcription by changing their affinity to specific operator sequences, often upon binding to small molecules. This simple yet powerful genetic control mechanism cements aTFs as foundational elements for small-molecule biosensing in synthetic biology^{4–12}.

Despite the versatility of aTFs as molecular switches, a critical limitation in the field is the ability to create aTFs specifically tailored to sense particular target molecules. This shortcoming hinders the precise customization of aTFs for applications such as metabolic engineering, circuit design, and biosensing, where the ability to selectively respond to defined ligands is paramount. A recently created database for potential ligand-inducible transcription factors called GroovDB consists of 199 regulators and 294 unique ligands at the time of writing¹³. GroovDB was created by parsing peer-reviewed literature and curating information related to aTFs such as ligand and operator identities. However, only a small subset of these have been experimentally validated as biosensors¹⁴. Although bioprospecting has uncovered new aTFs for specific molecules, this strategy is slow and

¹Department of Biochemistry, University of Wisconsin-Madison, Madison, WI, USA. ²711th Human Performance Wing, Air Force Research Laboratory, Wright Patterson Air Force Base, OH, USA. ³Dane County Youth Apprenticeship Program, State of Wisconsin Department of Workforce Development, Madison, WI, USA. ⁴Great Lakes Bioenergy Research Center, University of Wisconsin-Madison, Madison, WI, USA. ⁵Department of Bacteriology, University of Wisconsin-Madison, Madison, WI, USA. ⁶Department of Chemical and Biological Engineering, University of Wisconsin-Madison, Madison, WI, USA. ⁷These authors contributed equally: Kyle K. Nishikawa, Jackie Chen. ✉ e-mail: sraman4@wisc.edu

impractical for the discovery of aTFs designed to sense any arbitrary molecule of interest (Fig. 1a)^{15,16}. A promising alternative is to engineer the specificity of known aTFs toward new ligands. However, evolved aTFs with successfully altered specificities typically only bind ligands structurally similar to the native ligand, limiting the overall utility of directed evolution to design biosensors for entirely new chemical classes^{17–20}.

Altering the ligand specificity of allosteric proteins presents a formidable challenge in library design and screening. The challenge arises from the tight interconnection of residues involved in ligand binding with those crucial for allosteric actuation^{21–24}. Therefore, libraries designed solely to optimize protein–ligand interactions, as one might for designing binders, without considering the role of residues in allostery, are likely to abrogate the switch-like properties of aTFs. Furthermore, successful designs within a library are typically rare and may display weak activity toward non-native ligands. We need a high-throughput screening method with high sensitivity to effectively identify these rare, low-activity variants amid a vast pool of non-functional designs. These low-activity variants are useful scaffolds

from which high-activity variants can ultimately be derived by iterative design or directed evolution. Additionally, data from low-activity variants, which are commonly overlooked, can be leveraged to strengthen machine-learning models. Unfortunately, conventional enrichment techniques (like flow cytometry, plate-based screening, or growth selections) cannot identify these important variants owing to their lack of sensitivity, scale, or both^{17,25–29}.

In this work, we report Sensor-seq, a platform for creating aTF biosensors with high sensitivity and scale (Fig. 1b). Sensor-seq combines phylogeny-guided sequence diversification for library design with an RNA barcoding system to screen aTF variants through deep sequencing. As a starting scaffold for design, we seek a promiscuous aTF with a large binding pocket that could accommodate diverse ligands. We reason that designed mutations could customize the binding pocket with shape complementary interactions for a target ligand. Our rationale is supported by prior research demonstrating that enzyme evolvability and acquisition of novel functions benefit from a promiscuous starting point^{30,31}. Therefore, we choose the aTF TtgR, a multi-drug efflux regulator, as our starting scaffold because it is

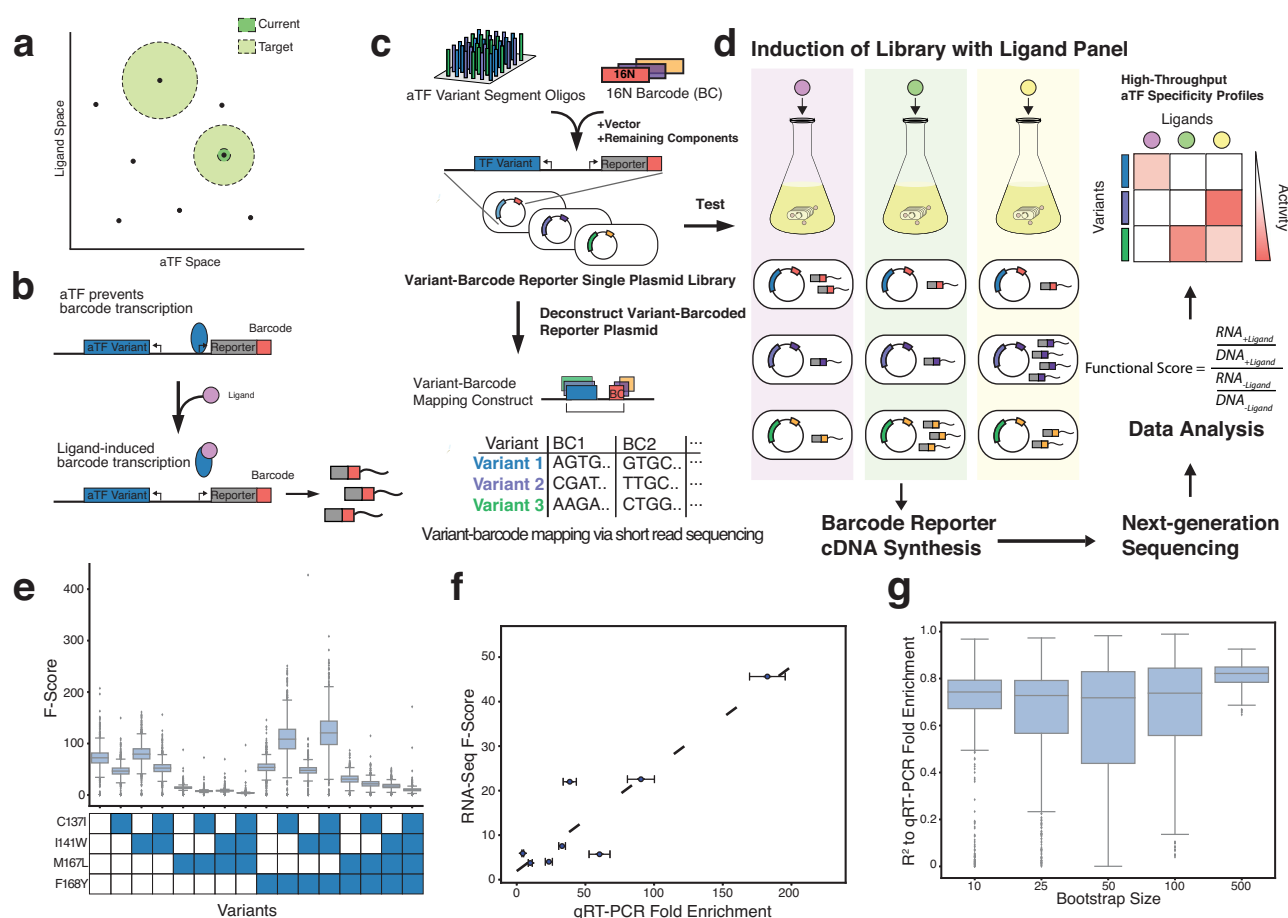


Fig. 1 | An RNA-Seq approach for high-throughput aTF characterization. **a** The space of allosteric transcription factors (aTFs) and small molecules that can be sensed with different approaches. Black points represent specific ligand:aTF pairs. The dark green circle represents the extent to which existing methods can currently expand aTF:ligand affinity. The light green circle represents the extent to which new methodologies must increase aTF:ligand pairs. **b** Construct design pairs aTF variant activity to transcription of randomized barcodes. **c** Construction and mapping of the aTF library to barcodes. **d** Workflow to acquire ligand specificity profiles for each aTF variant across different ligands. **e** Box plot showing distribution of F-scores for barcodes mapped to each TtgR variant via RNA sequencing (RNA-Seq) after induction with naringenin. The number of barcodes for each variant is described in Supplementary Fig. 1. The box represents the interquartile range

(IQR). Whiskers extend to 1.5 times the IQR. Fliers denote points that lie outside the whiskers. Center line represents median. Filled blue squares indicate the presence of the corresponding mutation in the variant. **f** Correlation of quantitative reverse transcription polymerase chain reaction (qRT-PCR) fold enrichment and RNA-Seq F-scores for 8 out of 16 clonal TtgR variants after induction with naringenin. qRT-PCR fold enrichment data is presented as the mean \pm SD of three biological replicates. The RNA-Seq fold enrichment value was calculated by summing the counts of all barcodes associated with a particular variant (see “Methods”). The R^2 for this dataset is 0.83. **g** Bootstrap correlation of qRT-PCR fold enrichment to RNA-Seq data for 8 of the 16 variants. Groups of 10, 25, 50, 100, or 500 barcodes were sampled for each variant across 500 cycles. Box plot elements same as (e). Source data are provided as a Source Data file.

known to bind to several antibacterial molecules and has a large binding pocket (1500 Å³)³². We evaluate 17,737 TtgR variants against a panel of seven non-native ligands and two native ligands. Sensor-seq identifies biosensors for all but one of these ligands, with high dynamic range and diverse specificity profiles. Moreover, our comprehensive dataset allows us to identify distinct linchpin positions and mutations driving specificity toward each ligand and unravel the molecular principles governing protein–ligand specificities for this scaffold protein. We obtain a crystal structure of a TtgR design bound to naltrexone, a non-native ligand, to elucidate the structural basis of the customizable specificity of TtgR to many ligands, a generally uncommon property among proteins. To illustrate the practical application of these engineered aTFs, we construct cell-free biosensing systems for detecting naltrexone and quinine, which have potential applications in detecting opioid overdose and wastewater contamination, respectively. In summary, Sensor-seq advances the on-demand design of biosensors for any target molecule. By broadening specificities to include non-native ligands, Sensor-seq liberates us from the restrictions imposed by solely relying on natural biosensors. Finally, the Sensor-seq workflow is adaptable to other proteins whose function can be linked to transcription.

Results

Validation of Sensor-seq on a pilot library

Sensor-seq assesses a pooled library of aTF variants by linking their ligand-induced responses to deep sequencing. Each cell expresses a unique aTF variant under a constitutive promoter regulating a reporter locus controlled by the aTF's native promoter (Fig. 1b). The activity of each aTF variant is quantified by measuring reporter transcript RNA levels through RNA sequencing (RNA-seq). The F-score, a normalized ratio of reporter transcript levels in the presence and absence of the ligand, is a quantitative measure of each variant's activity. The advantage of this approach is that aTF variants that are allosterically inactive either in the constitutively OFF (always repressing transcription) or constitutively ON state (unable to repress transcription) have an F-score of -1. These unproductive variants can be disregarded, allowing us to focus on the rare but potentially productive members of the library (F-score > 1).

A key challenge for high-throughput screens of proteins that regulate gene expression in *trans* (like aTFs) is the association of each variant's genotype with its transcriptional output. Linking the aTF variant to its function (e.g., reporter abundance) presents a technical challenge for mapping genotype to phenotype. We addressed this technical hurdle with a barcoding method. Each aTF variant is placed in *cis* with randomized barcodes with an intervening constant region comprising the unmutated part of the aTF and the promoter regulating the reporter (Fig. 1b). The randomized barcode is analogous to that in reporters used in most massively parallel reporter assays. However, sequencing the reporter alone would be insufficient for identifying which aTF variant was responsible for transcribing each reporter. To map the aTF to the transcript barcode, the intervening constant region is first removed by PCR amplification, and the aTF variant and reporter barcode are brought in close proximity with a Golden Gate Assembly. The variant and the associated barcode are mapped by traditional short-read deep sequencing (Fig. 1c). The transcript levels of each barcode of this screening construct reveal whether the associated aTF variant responds to the target ligand. Since the barcode region is short, we use high-volume, short-read sequencing to accurately quantify the activity levels of each variant (Fig. 1d). This platform, Sensor-seq, is generalizable and easily repeated to profile the variant library among any number of environmental perturbations (e.g., ligands). Our mapping scheme scales as a constant with the number of ligands tested. That is, only a single round of mapping is needed for users to test any number of ligands against thousands of aTF variants. *E. coli* containing the plasmid library are dosed with either the target ligand or a vehicle

control and harvested in log phase to obtain both total RNA (sequenced as cDNA) and the library plasmid DNA (Fig. 1d). The cDNA count provides a measure of function while the plasmid DNA count is used for normalization. These counts are used to calculate functional scores (F-scores) for each variant induced with different ligands. The library can be incubated with several ligands independently and evaluated in a single pooled sequencing to achieve scale.

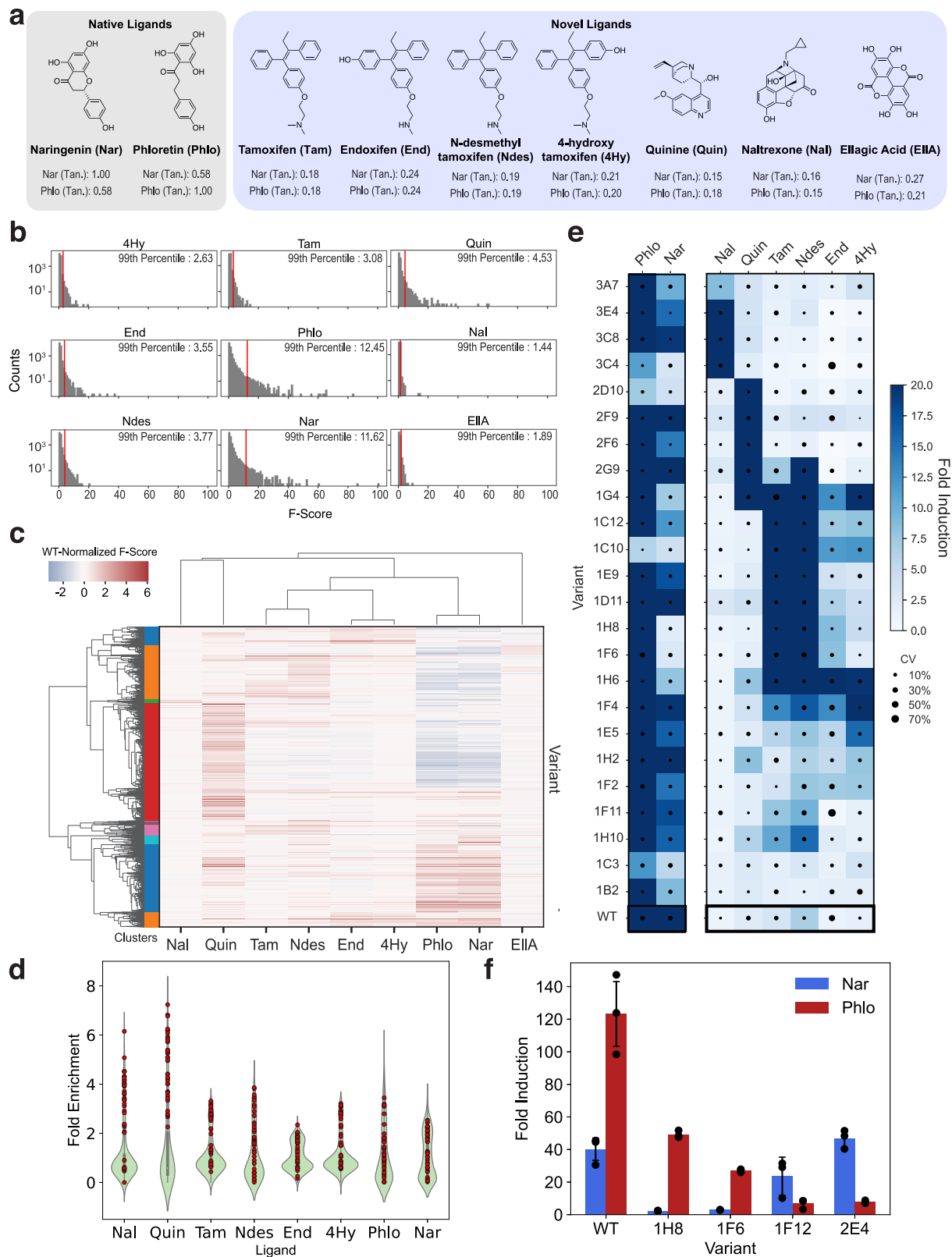
To verify the effectiveness of our proposed methodology, we applied Sensor-seq to a small library of 16 TtgR variants previously characterized for their responses to a native ligand, naringenin, using a GFP reporter protein³³. Gene fragments encoding these variants were incorporated with random 16N reporter barcodes into our screening construct and mapped through short-read sequencing, with each variant mapping to ~8000 barcodes (Supplementary Fig. 1). The F-score of each variant in response to 1 mM naringenin was determined using pooled RNA-seq, calculating cumulative barcode counts of the reporter relative to the vehicle control after normalizing to plasmid DNA (see “Methods”) (Fig. 1e). From the 16 variants, 8 spanning the full range of naringenin responses were selected for clonal quantification through qRT-PCR. Comparison of the qRT-PCR fold enrichment and the Sensor-Seq (RNA-seq based) F-score showed a high correlation ($R^2 = 0.83$) (Fig. 1f). In summary, Sensor-seq accurately reproduced activity differences observed via qRT-PCR for a small library.

When dealing with larger protein libraries, the requirement of 8000 barcodes per variant for sequencing becomes impractical. To evaluate the sensitivity of our F-score measurements to changes in the number of barcodes per variant, we down-sampled the number of barcodes per variant. Using a Monte Carlo sampling approach, we randomly selected 10, 25, 50, 100, and 500 barcodes per variant for 500 trials each and scored each sample based on its correlation to the qRT-PCR dataset (Fig. 1g). Each bootstrap group exhibited, on average, similar correlation to the qRT-PCR assay compared to the 8000 barcodes per variant. Our data suggest that sequencing volumes can be reduced ~800-fold and still preserve accuracy of measurements. With these reduced barcode-to-variant and sequencing requirements, Sensor-seq can increase scalability by accommodating larger protein libraries and have reduced cost.

Identifying sensors for non-native ligands

Having validated Sensor-seq, we next sought to apply this approach to redesign the specificity of TtgR toward non-native ligands. We used FuncLib to generate a library of TtgR variants in a ligand-agnostic manner (i.e., not targeted toward any particular ligand) with mutations around the binding pocket³⁴. FuncLib was developed to increase the thermodynamic stability of proteins by combining evolutionary and energy-guided design to introduce mutations^{34–37}. FuncLib delineates a sequence space of allowable point mutations within the ligand-binding site of TtgR by filtering out mutations that are unlikely to occur naturally as determined by a multiple sequence alignment of TtgR homologs. This important filtering step eliminates mutations likely to disrupt allostery but preserves mutations that provide diversity^{34–37}. Mutations passing the evolutionary curation test are computationally modeled with Rosetta into the TtgR structure to assess their stability. Only those that maintain protein stability are considered viable point mutations. Based on a list of allowable point mutations using FuncLib, we created a combinatorial library of 17,737 variants containing between 1 and 4 mutations per sequence.

We chose nine ligands (two native and seven non-native) with diverse chemical structures to develop new biosensors for (Fig. 2a). The native ligands of TtgR generally have the following characteristics: at least one aromatic ring, uncharged, and around 250–450 kDa³². Three factors guided our selection of non-native ligands: (1) dissimilarity to native ligands, (2) discriminating highly similar ligands, and (3) practical applicability. To balance the trade-off between dissimilarity to native ligands and the potential for discovering new



biosensors, we chose ligands with ring-like structures similar to native ligands but were chemically different, based on Tanimoto scores (Fig. 2a). We selected four derivatives of tamoxifen (Tam), a breast cancer therapeutic, to create specific and multi-specific sensors. Tamoxifen is a selective estrogen receptor-modulating prodrug that is converted into active metabolites, 4-hydroxy-tamoxifen (4Hy) and *N*-desmethyltamoxifen (Ndes). These two metabolites are then

catabolized to endoxifen (End)^{38,39}. These four metabolites share a core structure with three aromatic rings, one linked to a methylamine sidechain, but they differ in the derivatization of aromatic rings and/or the sidechain. The tamoxifen analogs served as a stringent test to assess our ability to design specificity switches between highly similar non-native ligands. Additionally, we selected quinine (Quin), naltrexone (Nal), and ellagic acid (EIIA) as non-native ligand targets. Quinine, a

Fig. 2 | Identifying sensors in a ligand agnostic library. **a** Ligands used in this study. Tanimoto scores (Tan.) are provided comparing each ligand to native ligands, naringenin (Nar) and phloretin (Phlo). **b** Histograms showing distribution of F-scores for each variant in the library with each tested ligand. Red line denotes the 99th percentile of F-scores. **c** RNA-Seq fold enrichment data for 16,191 variants which passed filters across nine ligands. Ligands and variants have been clustered via the UPGMA algorithm with a correlation distance metric and a target of 12 clusters (see “Methods”). The different clusters are denoted by the colored bars on the left of the heatmap. aTF function is shown as the $\log_2(\text{F-score})$ normalized to wildtype (WT). **d** Violin plot showing distribution of fold enrichment scores of a 251-member library consisting of the top 40 variants for each ligand (determined by F-

scores) after induction with indicated ligand and fluorescence-based sorting. The box represents the interquartile range (IQR). Whiskers extend to 1.5 times the IQR. Center dot represents the median. Red circles indicate scores of variants selected as top 40 for corresponding ligands. **e** Heatmap showing clonal measurements of fold induction of WT TtgR and variants that were selected as top 3 for each ligand based on data from **(d)**. Right panel consists of non-native ligands, and the left panel consists of the two native ligands. A score ceiling of 20 was imposed for visualization. Black circles within heatmap cells reflect the coefficient of variation (CV). **f** Bar plot of clonal measurements of variants with a specificity switch for the native ligands. Bars represent the mean fold induction \pm SD ($n = 3$ independent replicates). Source data are provided as a Source Data file.

malaria therapeutic, features a quinoline ring with two aromatic benzene rings, resulting in a distinctive non-planar structure⁴⁰. Naltrexone, an opioid analog used in addiction treatment, has a tetracyclic arrangement with three fused six-membered carbon rings⁴¹. Ellagic acid is a plant polyphenol with a planar, symmetric arrangement. (Fig. 2a)⁴². We included two native ligands of TtgR belonging to the flavonoid family, naringenin and phloretin, with the goal of identifying biosensors that could discriminate between the two molecules³². All seven non-native molecules differ significantly from the native flavonoid ligands (Fig. 2a). Six out of the seven non-native molecules are synthetic chemicals for which traditional genome-mining for natural biosensors would likely be ineffective. Since our primary goal was to identify gain-of-function variants for non-native ligands, we did not design against the native ligands, i.e., we aimed to expand TtgR's ligand-binding repertoire without consideration for native specificity. We reasoned that the potential applications of new biosensors would not likely involve the native ligands, which are uncommon plant metabolites. Biosensors for the tamoxifen family, naltrexone, and quinine could be used in practical applications such as measurement of drug metabolism, detection of drug overdose, and monitoring environmental levels in wastewater, respectively.

From a synthesized and cloned variant library of 17,737 members, we observed an average of 20 barcodes per variant for 17,533 variants (98.8%) (Supplementary Fig. 2, Supplementary Data 1, 2). We removed variants with a high degree of noise based on the coefficient of variation across biological replicates (see “Methods”). Variants that did not meet this threshold were excluded from further analysis. Reporter RNA was quantified in cells containing the library, which were incubated with each of the nine ligands or their respective vehicle controls (Supplementary Table 1). We calculated the distribution of F-scores for each ligand to quantitatively assess the library's performance (Fig. 2b, Supplementary Fig. 3). The F-score at the 99th percentile provided an estimate of the activity and prevalence of variants responding to a particular ligand. Naringenin and phloretin exhibited the highest F-scores at the 99th percentile (11.6 and 12.5, respectively), followed by quinine (4.5), *N*-desmethyltamoxifen (3.8), endoxifen (3.6), tamoxifen (3.1), and 4-hydroxytamoxifen (2.6). In contrast, both naltrexone and ellagic acid had the lowest F-scores (1.4 and 1.9) (Fig. 2b), suggesting that hits for these ligands, if present in the library, were exceedingly rare. As a reference, the wildtype TtgR's F-scores for its native ligands phloretin (2.0) and naringenin (1.9) were considerably lower than those observed for the library at the 99th percentile on the non-native target ligands (Supplementary Fig. 4). This suggests that hits for non-native ligands were likely real.

We performed unsupervised agglomerative clustering of 16,191 filtered variants against the nine ligands based on their functional profiles across all ligands (vertical) and the chemical relatedness of ligands (horizontal) with the activity scaled relative to wildtype TtgR (Fig. 2c, Supplementary Fig. 5). Activity on naringenin and phloretin (native ligands of TtgR) showed that nearly half the FuncLib-derived library retained ligand-induced allosteric response on par with or exceeding that of wildtype TtgR. Moreover, ~85% of TtgR variants

demonstrated the ability to repress transcription (Supplementary Fig. 6). Thus, some variants retained DNA binding activity but lost the ability to respond to native ligands. As a point of reference, the percentage of repression-competent variants was as low as 15% in our previous LacI variant library which was designed without FuncLib¹⁷. These results collectively suggest that FuncLib adeptly navigated the sequence space to preserve a high fraction of allosterically active and repression-competent variants in the library, thereby increasing the likelihood of successful designs.

Most TtgR variants showed concordant activities on the native ligands, naringenin, and phloretin, similar to wildtype TtgR (Supplementary Fig. 7). However, a small fraction exhibited differential activities toward one molecule or the other—either higher activity on naringenin than phloretin or vice versa. These included variants that lost activity on one ligand but maintained wildtype-like activity on the other (Fig. 2c, blue and white shades), and those that retained wildtype-like activity on one and outperformed wildtype on the other (Fig. 2c, white and red shades). We also observed differential activities among the four tamoxifen derivatives (Fig. 2c). There was minimal overlap between variants responding to the pair of hydroxylated derivatives, 4-hydroxytamoxifen and endoxifen, and the pair of non-hydroxylated derivatives, tamoxifen, and *N*-desmethyl-tamoxifen. Further, a small group of variants could discriminate between tamoxifen and *N*-desmethyl-tamoxifen. A substantial fraction of sequences responded to quinine, and these quinine-responsive sequences constituted an independent group that did not appear to overlap with variants responding to other ligands (Fig. 2c). A few variants appeared to respond to naltrexone and ellagic acid, albeit weakly in the case of ellagic acid (Fig. 2c). These results show that the Sensor-seq can identify allosterically responsive variants to non-native inducers and successfully distinguish structurally similar ligands to create specific aTF biosensors.

We carried out secondary and tertiary screens to progressively winnow candidates and to validate those that are most likely hits from the primary Sensor-seq screen. The secondary screen involved pooling the top 40 variants for each ligand in a cell-based reporter fluorescence assay, and the tertiary screen included clonally assaying individual variants from the secondary screen. For the secondary screen, we created a mini-library by resynthesizing 251 variants corresponding to the top-performing variants for each ligand, including overlapping hits between different ligands (e.g., naringenin vs. phloretin or tamoxifen derivatives) (Supplementary Fig. 8). From this mini-library, we sorted variants capable of repressing GFP expression in the absence of any small molecule. (Supplementary Figs. 9 and 10). These repression-competent variants were then exposed to each ligand and the high-fluorescence cells were isolated and sequenced (Supplementary Figs. 9 and 11). Each variant was scored based on the percentage fold change in abundance between the high fluorescence and repressed sorted populations. The secondary screen was consistent with the results of the primary Sensor-seq results (Fig. 2d). The top-performing variants for each ligand, as determined by F-scores obtained through Sensor-seq (indicated by red dots), consistently exhibited high fold

enrichment when evaluated using the GFP reporter assay (Fig. 2d). However, certain variants associated with tamoxifens, naringenin, phloretin, and naltrexone gave considerably lower fold enrichment changes in the GFP assay compared to their Sensor-seq scores. This observation raises the possibility that the primary Sensor-seq screen is susceptible to false positives. Quinine showed the most consistent response between the two assays in the top 40 variants. We found no hits for ellagic acid in the GFP assay, and thus this ligand was excluded from further analysis. As a tertiary screen, we clonally tested the top three unique variants for each ligand using the GFP reporter by measuring the fold induction: the ratio of mean cell fluorescence in the presence of the ligand and in the absence of the ligand (see “Methods”). As expected, WT TtgR had a strong response to native ligands, phloretin, and naringenin, with only weak activity on *N*-desmethyltamoxifen and little to no activity on the other non-native ligands (Fig. 2e). We observed variants with high dynamic range (i.e., fold induction) for each ligand. The greatest fold induction was 29 for 4-hydroxytamoxifen, 21 for endoxifen, 43 for naltrexone, 105 for naringenin, 108 for *N*-desmethyl tamoxifen, 234 for phloretin, 205 for quinine, and 108 for tamoxifen (Fig. 2e).

We observed remarkable diversity in the specificity profiles of different variants (Fig. 2e). Variants that acquired specificity for naltrexone (3A7, 3E4, 3C8, and 3C4) showed varied responses to the ligand, ranging from 9- to 43-fold induction. These monospecific variants did not respond to the other non-native ligands. Mono-specificity was also observed with 3 (2D10, 2F9, and 2F6) of the 5 strong quinine responders (Fig. 2e). In particular, one variant, 2D10, exhibited monospecificity toward quinine with minimal activation from even the native ligands. The two other quinine-responsive variants showed varying degrees of broader specificity. Specifically, 2G9 was additionally activated by tamoxifen and *N*-desmethyltamoxifen, while 1G4 was activated by all non-native ligands except naltrexone, with fold induction ranging from 13- (endoxifen) to 98-fold induction (tamoxifen). Among variants with broad specificity, the activities toward each ligand varied considerably. For instance, 1H6 shared a similar specificity profile as 1G4, except this variant had high activity toward endoxifen (21-fold induction) at the cost of reduced quinine response (8-fold induction) (Fig. 2d). 1F4 is another variant that shared a similar profile to 1G4 and 1H6 but showed high activity on 4-hydroxytamoxifen (29-fold induction) and a more muted response to endoxifen, *N*-desmethyltamoxifen, tamoxifen, and quinine (Fig. 2e). Eight variants (1C12, 1C10, 1E9, 1D11, 1H8, 1F6, 1F11, and 1H10) appeared to be largely bi-specific to *N*-desmethyltamoxifen and tamoxifen (both methylated), with activities ranging from 8- to 90-fold induction (Fig. 2e). 1C12 and 1C10 displayed weak to moderate activity on the other two tamoxifen derivatives, 4-hydroxytamoxifen and endoxifen (both hydroxylated) (Fig. 2e). However, we did find variants specific to only 4-hydroxytamoxifen and endoxifen in the secondary screen.

Most variants that had high activity on other ligands retained activity on naringenin and/or phloretin, which is consistent with our design criterion of not selecting against native ligand responders. Our data also suggest true specificity switches away from native function may be difficult because binding pockets that could accommodate the larger non-native ligands can also fit the native ligands. (Fig. 2e). 2F9 had an even greater response to naringenin and phloretin (101 and 234-fold, respectively) compared to wildtype TtgR (40 and 123-fold, respectively). Two variants (1H8 and 1F6) showed specificity with dramatically higher activity on phloretin over naringenin (49- vs. 2-, and 27- vs. 3-fold induction, respectively). To determine if specificity for naringenin can be favored over phloretin, we clonally evaluated two variants (1F12 and 2E4), which showed a preference for naringenin from our secondary screen. 1F12 and 2E4 retained near-WT responses for naringenin (24- and 47-fold induction, respectively) with a 15- and 18-fold loss in activation with phloretin compared to WT (7- and 8-fold induction, respectively) (Fig. 2f).

We next evaluated the predictive performance of Sensor-seq using the results from the tertiary screen as the ground truth (Supplementary Data 3). We observed a strong relationship between the F-scores from Sensor-seq and the fold enrichment scores from the clonal flow cytometry data (Spearman's $\rho = 0.78$) (Supplementary Fig. 12a). In comparison, the results from the secondary screen had only a moderate relationship to the clonal data (Spearman's $\rho = 0.57$) (Supplementary Fig. 12b). The poorer correlation likely reflects the inherent noise in screens based on cell sorting. We also examined the false positive and false negative rates at various F-score and fold enrichment thresholds (Supplementary Fig. 12c, d). The F-score and fold enrichment values reflect the activity level of individual variants in the RNAseq and sorting experiments, respectively. At a relatively stringent F-score threshold of 3 and a fold enrichment threshold of 1.5, we observed no false positives but ~41% false negatives. In contrast, a looser F-score threshold of 1.5 and fold enrichment threshold each at 1.5 resulted in ~7% false positives and ~19% false negatives. These results suggest a trade-off between false positive and false negative rates at each threshold. Using the same F-score threshold and testing multiple fold enrichment thresholds, we find an average precision of 0.94 (Supplementary Fig. 12d). Taken together, the pooled data from Sensor-seq largely summarizes the individual measurements acquired from clonal tests despite these errors.

In summary, Sensor-seq enabled the identification of TtgR variants with unique specificity profiles, including rare hits, from the FuncLib ligand-agnostic library. We identified variants exhibiting strong activity on six of the seven non-native ligands, showcasing varying degrees of specificity toward different ligand groups, and enhanced activity on native ligands. These results also highlight the inherent malleability of some proteins such as TtgR to accommodate a diverse family of ligands. We posit that the malleability is a natural consequence of TtgR's role as a regulator of a multidrug efflux, which necessitates the ability to detect and respond to different compounds, reflecting an evolutionary pressure to maintain a broad ligand-binding capability.

Unsupervised learning reveals key residues for ligand specificity

We next sought to determine amino acid sequence preferences for each ligand. To uncover and visualize these sequence determinants in the high dimensional sequence-ligand landscape, we generated a two-dimensional Uniform Manifold Approximation and Projection (UMAP) projection of 17,430 variants using physicochemical embeddings⁴³. The amino acid at each of the 11 mutated positions in TtgR is incorporated as a feature using a 19-length dimensionally reduced AAindex which captures diverse amino acid properties such as polarity, hydrophobicity, and alpha/beta propensities^{44,45}. Thus, each variant is physicochemically encoded with (11 residues \times 19 features) 209 total features^{45,46}. Because variants differ by only a few mutations, we chose AAindex as a simpler and lower-dimensional encoding than pre-trained transformer embeddings. We then applied Hierarchical Density-Based Spatial Clustering of Applications with Noise (HDBSCAN) to group variants that are similar⁴⁷. With our chosen hyperparameters, we were able to organize 16,882 variants (~97% of total variants) into 23 clusters containing between 93 and 3916 variants, where each cluster represents sequences with similar physicochemical properties (Fig. 3a, Supplementary Figs. 13 and 14).

To identify functional regions in sequence space for each ligand, we overlaid the F-scores obtained from Sensor-Seq for each ligand on the UMAP projection (Supplementary Fig. 15). An examination of the UMAP with overlaid F-scores revealed that functional sequences for each ligand are distributed across multiple clusters. This implies the potential existence of degenerate solutions, wherein several groups of sequences can bind to the same ligand. A notable feature is the high local ruggedness, where adjacent points within a cluster exhibit gains or losses of function, suggesting that small physicochemical changes

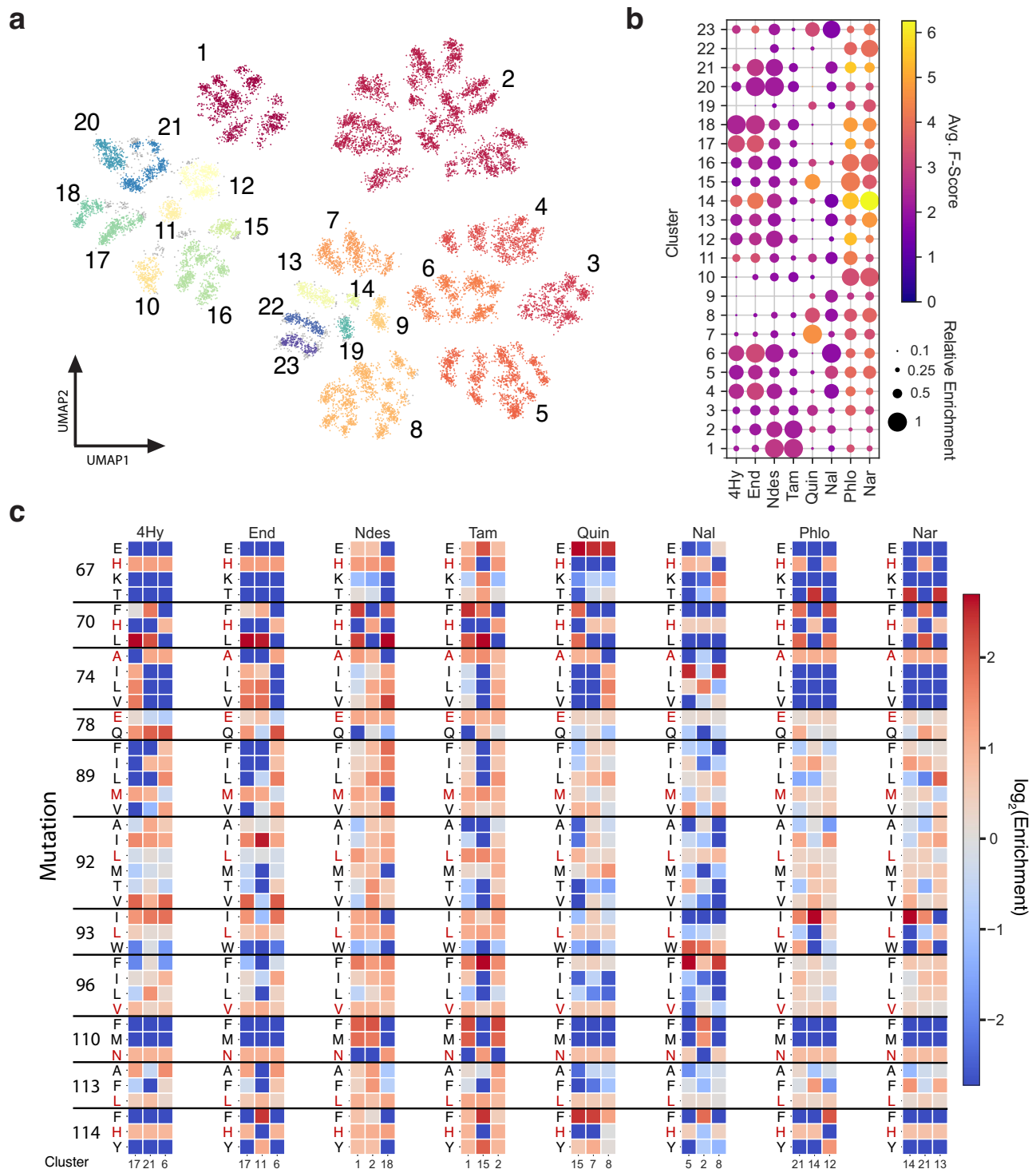


Fig. 3 | Unsupervised learning reveals amino acid preferences for ligand specificity. **a** Uniform Manifold Approximation and Projection (UMAP) 2D embedding of 17,430 variants with physicochemical properties of amino acids at each variable position in the TtgR library. Multicolored plot shows 23 clusters identified with Hierarchical Density-Based Spatial Clustering of Applications with Noise (HDBSCAN). Gray points correspond to points the algorithm identified as noise. **b** Dot plot showing the performance of each of the 23 clusters identified after UMAP-HDBSCAN. The color of each dot represents the average F-score of all variants within the cluster-ligand pair that has a minimum of 1.5 F-score. The size of the

dot represents the percentage of variants with >1.5 F-score within each cluster normalized to the highest percentage for that ligand. **c** Heatmaps for the top 3 performing clusters of each ligand showing \log_2 Enrichment of each possible amino acid at the variable positions of TtgR. Clusters are shown from rank 1 to 3 going from left to right. Enrichment was calculated by obtaining the F-score-weighted frequency of amino acids in the cluster using variants with a minimum of 1.5 F-score and normalizing to the DNA count-weighted frequencies of the initial library. Red letters denote the wild-type residues. See "Methods" for an in-depth description of analysis. Source data are provided as a Source Data file.

in closely related protein sequences are sufficient to alter function dramatically (Supplementary Fig. 15).

To assess the distribution of ligand-responsive sequences across clusters quantitatively, we analyzed the percentage of hits and the average F-score of hits within a cluster (Fig. 3b, Supplementary Fig. 16, Supplementary Data 4, 5). Here, a hit is defined as any variant with an F-score of 1.5 or higher (Fig. 3b, Supplementary Fig. 16). The choice of 1.5 as the threshold for activity represents the approximate value where incremental increases in the threshold do not significantly change the number of hits (Supplementary Fig. 17a, b). Using these metrics, we first compared the performance of each cluster for each ligand to gain a global view of the sequence-function landscape. Then we delved deeper to elucidate sequence determinants of TtgR driving a gain-of-function response to each ligand. This involved examining the enrichment or depletion residues among highly active variants within the top three clusters for each ligand. To ensure that the sequences within the top clusters are a good representation of sequences that would be highly active on a specific ligand, we ranked the clusters based on the average F-score of hits (variants with F-score ≥ 1.5) and imposed a minimum of 15 hits per cluster. Because naltrexone only had 4 clusters that passed these thresholds and to simplify analysis, we conservatively picked the top three for each ligand.

The hydroxylated tamoxifen derivatives, 4-hydroxytamoxifen, and endoxifen, share similar top-performing clusters (Fig. 3b). As expected, *N*-desmethyl tamoxifen and tamoxifen share high-performing clusters mostly distinct from the top-performing clusters for 4-hydroxytamoxifen and endoxifen, suggesting, in general, distinct groups of sequences respond to both ligand pairs (Fig. 3b). Wildtype residue H67 is strongly enriched and L93W is strongly depleted among variants active on all tamoxifen derivatives (Fig. 3c). Variants active on 4-hydroxytamoxifen and endoxifen showed preference for histidine at position 67, glutamine at position 78, and asparagine at position 110 (Fig. 3c). In contrast, the variants active on *N*-desmethyltamoxifen and tamoxifen were more permissive of different substitutions at position 67, favored glutamate at position 78, and had either a phenylalanine or methionine in top clusters at position 110 (Fig. 3c). These three residues, 67, 78 and 110, may impart specificity for 4-hydroxytamoxifen and endoxifen through hydrogen bonding with the 4-hydroxyl group of the ligand. Specificity could be engineered for 4-hydroxytamoxifen and endoxifen if positions 78 and 110 contain glutamine and asparagine, respectively (Fig. 3c). Clusters 20 and 21 contain a relatively large portion (12–18%) of sequences that respond to *N*-desmethyltamoxifen and endoxifen, albeit weakly (lower average F-score than highly active clusters) (Fig. 3b), indicating the presence of variants within these clusters that are capable of accommodating both methylated and hydroxylated forms of tamoxifen but not the tertiary amine group found in tamoxifen and 4-hydroxytamoxifen.

Quinine and naltrexone have unique cluster profiles due to their distinct structures (Fig. 3b). Thus, specific protein sequence characteristics govern specificity for each of these ligands. For quinine, we observed enrichment of H67E, wildtype residue N110, and either phenylalanine or tyrosine at position 114, which may stabilize quinine with a π -stacking interaction. Positions 78, 89, and 92 are tolerant to most substitutions except threonine at position 92 (Fig. 3c). At position 96, isoleucine and leucine are disfavored, but valine and phenylalanine are permissible suggesting that size of the amino acid may not be the determining factor (Fig. 3c). For naltrexone, the key specificity determining mutations are H70, W93, and F96 (Fig. 3c). Interestingly, tryptophan at position 93 is strongly favored for naltrexone, but strongly disfavored for quinine (Fig. 3c).

Naringenin and phloretin-responsive sequences are distributed throughout the 23 clusters suggesting that the response to these native ligands is often robust to mutation (Fig. 3b). Both shared similar mutation profiles with positions 78, 89, 92, 96, and 113 tolerant to several substitutions (Fig. 3c). We observed co-occurring residues at

positions 67 and 70 for both ligands, where H67 is accompanied by either F70 or L70 and T67 is paired with H70 (Fig. 3c). N110 and H114 are strongly preferred for both native ligands (Fig. 3b). Increased activity on the native ligands may be mediated by interplay between these 4 key positions, 67, 70, 110, and 114. Cluster 12, a top cluster for phloretin but not for naringenin, contains sequences with either an H114F or H114Y. We also observed a markedly lower hit percentage for cluster 12 with naringenin (20%) than with phloretin (32%) (Supplementary Data 4). Thus, variants with a substitution at position 114 to an aromatic residue may have shifted activity towards phloretin than naringenin. Indeed, variants 1H8 and 1F6 from our clonal flow cytometry assessment (Fig. 2d) carried these aromatic substitutions at position 114 and had specificity for phloretin but not naringenin.

In summary, these results highlight the power of Sensor-seq to provide a holistic view of the mutational adaptability of TtgR to accommodate diverse ligands, and the key determinants of specificity for each ligand. By revealing the distribution of ligand-responsive sequences across clusters and elucidating sequence determinants driving gain-of-function responses, Sensor-seq offers valuable insights into the underlying landscape of protein–ligand interactions. We anticipate that this information-rich dataset will serve as a valuable resource for benchmarking machine-learning algorithms aimed at designing and understanding protein–ligand interactions.

Crystal structure reveals protein-stabilizing motifs for a non-native ligand interaction

To gain structural insights into TtgR's ligand specificity, we obtained a high-resolution crystal structure of a quadruple variant (A74L, L93W, N110M, and H114F) bound to non-native ligand naltrexone at 2.05 Å resolution (Supplementary Table 2, Supplementary Fig. 18). This structure revealed that TtgR's ability to interact with various ligands is likely due to an unusually large ligand-binding pocket (1500 Å³) for a ~200 residue protein³². The substantial pocket volume enables TtgR to bind ligands in various orientations and with different shapes and features, broadening the potential ligands TtgR can bind. For example, in a previous study, we showed that TtgR's binding pocket is malleable enough to bind to resveratrol, in either vertical or horizontal orientations, with distinct sequences³³. Ligand specificity is achieved by selecting residue rotamer states within the binding pocket that enhance the tightness of ligand fit primarily through van der Waals interactions. This capability is prominently evident in the naltrexone-bound structure. As a charge-neutral molecule, naltrexone does not require complementary charged residues. All four substitutions, A74L, L93W, N110M, and H114F, are small-to-large changes that facilitate better packing of naltrexone within the binding pocket (Fig. 4a). This jigsaw puzzle fit customizes the pocket for naltrexone and excludes other ligands, as seen by the high specificity of this mutant (Fig. 2d, 3A7). For instance, the mutations observed in the crystal structure appear to be participating in a total of three Methionine-Aromatic (Met-Aro) interactions, a ubiquitous, protein-stabilizing motif, involving the sulfur atom of methionine and the aromatic ring of a partner residue at a molecular distance of ~4–6 Å, that can play crucial roles in high-affinity ligand binding^{48,49}. L93W forms two 5.7–6.2 Å Met-Aro interactions with M167 and M89 (Fig. 4b) in the upper section of the binding pocket. N110M and H114F form a 5.4 Å Met-Aro interaction at the lower section of the binding pocket (Fig. 4b).

We also observed a larger portal for ligand entry into the binding pocket (Fig. 4c). The entry appears to be widened through the rotation of two key residues, C137 and R75 (Fig. 4c). In the quadruple mutant, the sulfur atom of C137 is rotated maximally ~101° around the β -carbon from the wildtype rotamer position and towards the entry port, thereby increasing access to the binding pocket³². A similar effect can be observed with the R75 guanidino group rotated ~135° around the β -carbon from the wildtype rotamer position and away from the entry port (Fig. 4c). This rotation can be explained by the additional steric

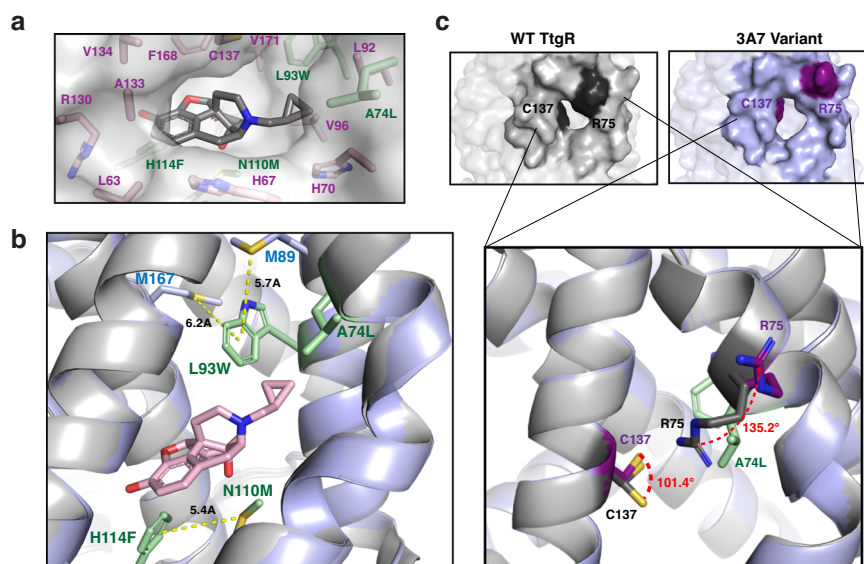


Fig. 4 | Crystal structure for naltrexone-bound Variant 3A7. a Binding pocket of Variant 3A7 with naltrexone. Gray molecule is naltrexone. Green side chains indicate mutant positions. Pink side chains indicate positions within 4 Å of the naltrexone molecule. **b** Three instances of Met-Aro interactions within the binding pocket, mediated by the mutant residues. WT TtgR (PDB: 7K1C) is shown in gray and 3A7 is shown in light purple. Green side chains indicate mutant positions.

Naltrexone is shown in pink. Molecular distances for each interaction are indicated. **c** Ligand entry port of WT TtgR (gray) and Variant 3A7 (light purple). Darkened surface residues denote positions of C137 and R75. Dotted red lines indicate rotation angles of C137 and R75 in WT (dark gray) TtgR and Variant 3A7 (dark purple). Green side chains indicate mutant positions.

bulk caused by the neighboring A74L substitution (Fig. 4c). The open conformation in 3A7 should allow access of differently-sized molecules to enter the ligand-binding pocket. In summary, the stabilizing energy of the jigsaw-fit ligand-aTF interaction and the larger entry portal likely facilitates the accommodation of naltrexone into the designed TtgR.

Design of cell-free biosensors for naltrexone and quinine

To demonstrate the practical utility of the designed biosensors, we sought to develop cell-free expression systems (CFE) for simple, economical analyte detection^{50,51}. We selected variants 3A7 (naltrexone) and 2F9 (quinine) for testing as these were able to repress transcription in the cell-free system (Fig. 5b, d, Supplementary Fig. 19). We first tested a biosensor for naltrexone given the compelling need for low-cost point-of-care approaches for detecting opioid use in rural communities without easy access to healthcare⁵². Naltrexone works by binding to and blocking the effects of mu-opioid receptors and is used to treat substance use disorders⁵³. As an opioid analog, naltrexone can serve as a stand-in substitute for opioids such as heroin or morphine. Two plasmids—one encoding the naltrexone-responsive 3A7 aTF (sensor plasmid) and one encoding GFP regulated by a TtgR promoter (reporter plasmid)—were added to a cell-free reaction containing processed *E. coli* extract together with the additional biological cofactors required for transcription and translation *in vitro*⁵⁴ (Fig. 5a). We first validated the ability of the sensor to repress expression of the fluorescent reporter by testing a range of sensor plasmid concentrations (0–10 nM) in the absence of analyte and quantifying fluorescence over 18 h (Fig. 5b). We observed dose-dependent repression of the GFP signal, with fluorescence being distinguishable across sensor plasmid concentrations within 2 h and reaching steady-state levels at ~10–12 h (Fig. 5b). The minimum fluorescence level was reached at 5 nM sensor plasmid (Fig. 5b). Having validated the 3A7 sensor as a functional repressor in a CFE, we next tested whether varying concentrations of naltrexone (0–100 μ M) can be detected by the sensor and abolish repression. Generally, by increasing the amount of naltrexone, we observed higher steady-state levels of fluorescence, with the largest stepwise increase obtained transitioning from 10 nM to 100 nM

(Fig. 5c). At 100 μ M naltrexone, we observed higher variance of the measurements, suggesting possible inhibition of the cell-free reaction at high concentrations of analyte (Fig. 5c). The fluorescence signal was readily distinguishable between 10 nM naltrexone and higher concentrations within 2 h, but took longer at lower concentrations (Fig. 5c).

As a second proof of concept, we aimed to create a cell-free sensor for quinine, another non-native ligand. One variant for quinine (2F9) had a 205-fold induction in our tertiary screen. We wanted to test if this strong activation can be translated into a cell-free system. Quinine has been used since the 17th century as an antimalarial drug and, despite its well-documented adverse side effects, is still in use today, primarily in underdeveloped countries with resource-limited healthcare systems or where safer modern alternatives are not available⁴⁰. In such communities, significant levels of quinine are likely to be detected in the wastewater, which is often recycled for irrigation or freshwater sustainability and can carry unforeseen environmental and health ramifications⁵⁵. Thus, an affordable biosensor for monitoring quinine metabolism and detection in water supplies may be of interest. We tested the 2F9 variant, one of the strongest quinine responders from our tertiary validation. We examined the time-dependent changes in fluorescence after adding a range of sensor plasmid (0–10 nM) in a similar experimental set-up with the naltrexone sensor. Like with the 3A7 sensor, we observed a dose-dependent repression that was evident within 2 h, and we found that 5 nM of sensor plasmid was sufficient for minimal fluorescence (Fig. 5d). Having validated that the 2F9 sensor can repress reporter expression, we next tested if expression can be restored with the addition of quinine. We tested 3 concentrations of quinine (0 μ M, 100 nM, and 200 nM) for 18 h. We found that at 100 nM, the initial fluorescence difference was distinguishable in 2 h and approximately doubled at steady-state, which took ~6–8 h to reach (Fig. 5e). We also observed a decrease in overall fluorescence when the quinine concentration was increased to 200 nM (Fig. 5e). The decrease in fluorescence can be explained by the weak DNA intercalating properties of quinine, which may have inhibitory effects on expression of the GFP reporter, or toxicity to cell-free reaction components⁵⁶.

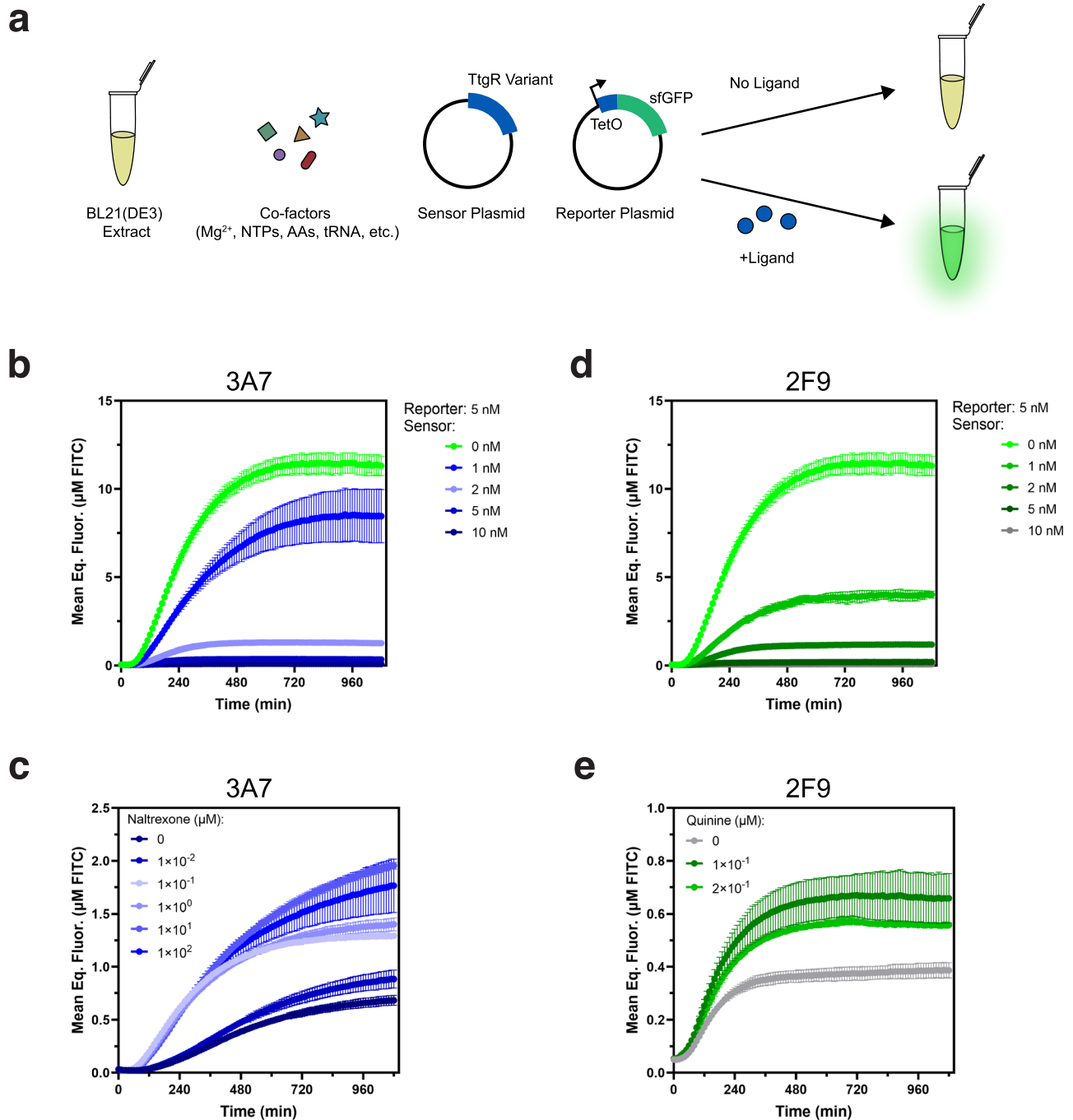


Fig. 5 | Cell-free biosensor for detection of naltrexone and quinine. **a** Cell-free expression system consisting of a sensor plasmid and an sfGFP reporter plasmid whose expression is mediated by ligand binding by the sensor. **b** The plasmid encoding 3A7 TtgR (sensor plasmid) was titrated against 5 nM of the plasmid encoding sfGFP regulated by 3A7 TtgR (reporter plasmid) in cell-free gene expression reactions. **c** 6 nM of the 3A7 sensor plasmid and 2 nM reporter plasmid were added to cell-free gene expression reactions and incubated with naltrexone at

log-fold increments. **d** 2F9 TtgR sensor plasmid was titrated against 5 nM reporter plasmid in cell-free gene expression reactions. **e** 4 nM 2F9 TtgR sensor plasmid and 2 nM reporter plasmid were added to cell-free gene expression reactions and incubated with indicated concentrations of quinine. Reported are the sfGFP concentrations normalized to a FITC standard over 18-h experiment at 30 °C, continuously monitored every 10 min. Data for (**b–e**) are presented as mean values \pm SD of three independent replicates. Source data are provided as a Source Data file.

We next compared the *in vivo* dose response of the 2F9 and 3A7 sensors to varying levels of quinine and naltrexone in a GFP reporter assay, respectively. We find that variant 2F9 had appreciable fluorescence with quinine even after a 25-fold dilution of the working ligand concentration used in previous screens (Supplementary Fig. 20). 3A7 had the same level of activation with naltrexone after a 5-fold dilution. The same concentrations used on WT did not produce a signal. Collectively, these data are a promising step toward

developing cell-free assays for naltrexone and quinine detection. Future steps could optimize the cell-free biosensors in matrices such as wastewater and biological fluids,

Discussion

Sensor-seq is a platform technology that can engineer aTFs with new ligand specificities. By integrating phylogeny-guided sequence diversification to preserve allosteric signaling with an RNA barcoding

system, Sensor-seq screens aTF variants through deep sequencing to achieve both sensitivity and scale. This approach enables assaying more variants without increasing experimental complexity in comparison to cell sorting and plate-based screening methods. When applied to redesign the specificity of a bacterial aTF, TtgR, Sensor-seq yielded variants exhibiting strong activity on six of seven non-native ligands, with distinct specificity profiles toward different ligand groups and enhanced activity and specificity toward native ligands. Statistical analysis of this dataset comprising nearly 160,000 sequence–function data points elucidated sequence determinants driving gain-of-function and provided detailed insights into the underlying landscape of protein–ligand interactions.

Sensor-seq is agnostic to the design methodology used to generate variants. To improve dynamic range or specificities, we recommend that a targeted design approach is employed (Rosetta/ProteinMPNN as examples)^{35–37,57}. At the screening step, multiple sequencing pools at different ligand concentrations can help identify variants with different dynamic ranges. Variants with higher dynamic ranges should have higher F-scores at higher ligand concentrations while those with lower will drop below the positive F-score threshold. Our current library and ligand selection highlights the capacity of Sensor-seq to pinpoint variants with altered specificity. These variants will have high F-scores for the target ligand (or ligand set) while having lower F-scores elsewhere. In these cases, a user should assay all the desired on-target and off-target ligands to cross-examine the F-score profile of the variants.

TtgR's role in regulating a multi-drug exporter necessitates its ability to detect and respond to different compounds, reflecting an evolutionary pressure to maintain a broad ligand-binding capability. This natural adaptive potential can be advantageous for the evolution of new functions. From a structural standpoint, TtgR's multifunctionality arises from an unusually large binding pocket volume that allows TtgR to accommodate ligands of various sizes and orientations^{32,33}. Ligand specificity is achieved by the selection of residue rotamer states within the binding pocket that enhances the tightness of ligand fit primarily through van der Waals interactions. Other bacterial regulators of multi-drug efflux pumps may emerge as potential scaffolds for future biosensor design studies⁵⁸. Sensor-seq itself is a generalizable approach that can be applied to different aTF scaffolds. aTFs with smaller binding pockets can be used as a scaffold for smaller ligands.

Redesigning the ligand specificity of allosteric proteins is challenging due to the interconnectedness of residues involved in ligand binding and those crucial for allosteric actuation. Sequence diversification must carefully balance both generating the necessary diversity for new ligand recognition as well as accounting for allosteric hotspots required for function⁵⁹. A key factor in this balance is the FuncLib method for library generation. FuncLib excludes mutations that are unlikely to occur naturally. This important step eliminates mutations that are likely to disrupt allostery and structure but preserves mutations that provide diversity. The Sensor-seq library of aTFs demonstrates ~85% repression competence compared to ~15% in a previous LacI library¹⁷. It is possible that mutations generated using FuncLib preferentially select for combinations that lead to allosterically functional proteins. These combinatorial mutations may allow TtgR to sample new conformational states that may confer new ligand specificities without abrogating native function and may explain the extraordinary diversity of ligands binding to TetR-family proteins^{1,60}. A previous deep mutational scan of TtgR done by our laboratory identified mutations that generate an allosterically dead variant, which stabilizes the inactive state and hinders ligand-induced state switching²³. We identified seven variants from Sensor-seq that contained dead mutants that were functional when combined with other mutations (Supplementary Fig. 21). This suggests that compensatory mutations elsewhere in the protein can rescue a dead variant and that

FuncLib can find these restorative mutations to generate allosterically functional variants.

Resources such as PDB or sequence databases have facilitated the development of machine learning tools for the prediction of structure and mutational effects on proteins^{61–63}. However, large-scale experimental datasets of designed protein–ligand interactions do not exist. The dataset of ~160,000 quantitative, sequence–function relationships from this work represents a useful large-scale experimental study of protein–ligand specificities of an allosteric protein. Despite focusing on a single protein scaffold, this dataset could have significant utility in training machine learning models to predict aTF variant function directly from sequence. These models could be leveraged to guide the design of aTFs capable of binding to ligands absent from the original training set. Additional guidance from models trained on alternative datasets such as ProteinMPNN⁵⁷ or MSA-based variational autoencoders⁶⁴ could further guide optimization. Technological advances in molecular docking and structure prediction tools could further reveal the molecular rules governing ligand specificity. As additional scaffolds are incorporated in future studies, we envision these datasets evolving into essential benchmarks for refining and evaluating machine learning algorithms dedicated to deciphering and designing protein–ligand interactions.

Sensor-seq is a general platform for discovery and design. Beyond its role in engineering functions, it serves as a useful tool for advancing basic science studies by integrating diverse protein libraries, generated through methodologies like deep mutational scanning, ancestral sequence reconstruction, and metagenomic analysis of protein domains to study sequence–function relationships and evolutionary biochemistry. Moreover, Sensor-seq's scope extends beyond one-component transcription factors (aTFs) to encompass multi-protein relays involved in transcription, such as two-component systems, chemoreceptors, and quorum sensing. This adaptability positions Sensor-seq as a comprehensive and adaptable resource for both applied engineering and the elucidation of fundamental biological principles. Future iterations of Sensor-seq will incorporate experimental data into machine-learning models to improve protein design.

Methods

Plasmid creation

PCR amplicons are generated using KAPA HiFi PCR kits (Roche, KK2102) following the manufacturer's protocol. Primers can be found in Supplementary Data 6. Amplicons are treated with 15 U of DpnI (New England Biolabs, R0176L) for 2.5 h at 37 °C followed by 20 min at 80 °C. Amplicons are then purified using EZNA Cycle Pure kits (Omega Bio-Tek, D6492-02). Isothermal assembly followed Gibson Assembly protocols (New England Biolabs), but contained 100 mM Tris-HCl pH 7.5, 20 mM MgCl₂, 0.2 mM dATP, 0.2 mM dCTP, 0.2 mM dGTP, 10 mM DTT, 5% PEG-8000, 1 mM NAD⁺, 4 U/mL T5 exonuclease (New England Biolabs, M0663L), 4 U/μL Taq DNA ligase (New England Biolabs, M0208L), and 25 U/mL Phusion polymerase (New England Biolabs, M0530L). Isothermal assembly reactions are diluted 10× in dH₂O prior to transformation. DH10B electrocompetent cells (New England Biolabs, C3020K) are transformed with 2 μL of diluted isothermal assembly reaction. Transformants are recovered in 700 μL SOC for 1 h at 37 °C. Dilutions are plated on LB-kanamycin (50 μg/mL) plates and incubated at 37 °C overnight. Colony PCR is performed using KAPA Robust (Roche, KK5005) using a single colony diluted in 100 μL of dH₂O. Plasmid purifications are performed using the ZR Plasmid Miniprep Classic kit (Zymo, D4016). All plasmid assemblies follow this methodology unless stated otherwise.

We first generated a sensor-reporter plasmid containing a TtgR gene under an apFAB61-BBaj61132 constitutive operator sequence, an sfGFP gene under control of the TtgR operator sequence, a kanamycin resistance cassette, and a ColE1 origin of replication⁶⁵. The sfGFP fragment with a modified TtgR operator sequence was amplified from

the TtgR_pBBR1_V2 plasmid. The TtgR gene was amplified from the TtgR_SCI01BBA plasmid. The sfColE1 backbone amplicon contains a kanamycin marker and a ColE1 origin. These fragments were assembled in a Gibson Assembly reaction as described above to form TtgR_ColE1_SPS. The sfGFP promoter was then modified to have the wildtype TtgR operator sequence. sfGFP with the wildtype operator sequence was amplified from a separate plasmid. The backbone amplicon was amplified from TtgR_ColE1_SPS and consisted of the TtgR gene, the kanamycin resistance marker, and the ColE1 origin. The assembled plasmid was labeled as TtgR_ColE1_SPS_V2 and is used as the vector for cell-based fluorescent assays. A third Gibson assembly reaction was required to insert stop codons and BsaI cut sites into the middle of the GFP gene to create the barcode insertion site. Primers KN_197 and KN_198 were used to linearize TtgR_ColE1_SPS_V2 and split the GFP gene after the L18 position. The primers have overhangs that contain a 32 bp complementary region, BsaI recognition sites for barcode insertion, and 2× ochre stop codons in one of the primers to truncate the GFP gene. The isothermal reaction was performed using 60 ng of the amplicon in a 50 °C incubation for 30 min. This construct was labeled TtgR_ColE1_SPS_V5 and used as the template for the RNA-Seq assays.

TtgR library generation

Plasmid libraries are generated using Golden Gate Assembly Kits (New England Biolabs, BsaI-HFv2, E1601L). The reactions undergo a cycling protocol of 30 alternating 5-min 37 °C and 16 °C cycles followed by a final 60 °C 5-min hold. The reactions are dialyzed against dH₂O on 0.025 µm semi-permeable membrane filters (Millipore, VSWP02500) for 1 h at room temperature. DH10B competent cells (New England Biolabs, C3020K) were transformed with 3 µL of dialyzed reaction via electroporation. Transformants were recovered in 1 mL of SOC and then diluted 2×, 5×, and 10× with fresh SOC. Each dilution recovered for 1 h shaking at 37 °C. 4 mL of LB-kanamycin (50 µg/mL) was added to each dilution after recovery and 50× and 500× dilutions were plated of each recovered dilution to calculate transformation efficiency. The remaining transformants were grown for 6 h shaking at 250 rpm. A frozen stock was made in 25% glycerol and stored at −80 °C for each dilution. Fresh cultures were created by diluting each 6-h growth 50× into fresh LB-kanamycin. These were grown overnight, and plasmids were harvested via ZR Plasmid Miniprep Classic kit (Zymo, D4016). All plasmid library assemblies follow this methodology unless stated otherwise.

Pre-defined or random barcodes were synthesized as a short primer (IDT). These barcode primers were combined separately with another constant primer to create short double-stranded fragments containing the barcode flanked by BsaI cut sites in a single cycle of PCR using KAPA HiFi (Roche, KK2102). 1 µL of this reaction was added into a second KAPA HiFi reaction with additional primers to increase the length of the amplicon over 18 cycles. The resulting amplicon was purified using the DNA Clean and Concentrator-5 kit (Zymo, D4004).

To generate the pilot 16-member mini-library, the TtgR gene variants were isolated from a set of 16 pre-existing plasmids each containing a single TtgR variant¹⁷. 100 ng of each amplicon was combined in a pool. The TtgR_ColE1_SPS_V5 backbone was amplified using primers that encompassed the sfGFP gene, the ColE1 origin, and the kanamycin resistance marker. The barcodes, TtgR gene variants, and backbone were assembled in a single Golden Gate reaction (New England Biolabs, BsaI-HFv2, E1601L) as described above.

FuncLib mutations were encoded into short oligos (Agilent) consisting of the TtgR gene region flanked by BsaI cut sites for Golden Gate assembly. Four pools of ~4400 variants were created by randomly combining between 1 and 4 tolerated mutations. Each pool had unique priming sequences to isolate from a pooled sample. The pooled library was diluted to 0.005 µM in Tris-HCl (pH 7.5). Each pool was amplified using Kapa HiFi and 1 µL of the diluted library in 15 cycles in triplicate.

The amplified reactions were pooled together and purified using the DNA Clean and Concentrator 5 kit (Zymo, D4004). The pooled oligos, barcodes, and TtgR_ColE1_SPS_V5 backbone were assembled using Golden Gate assembly (New England Biolabs, BsaI-HFv2, E1601L). The libraries with -15 barcodes per variant, calculated by CFU/mL, were selected for RNA-Seq.

Creating GFP control

To create a GFP positive control that will be used for quality control in Sensor-Seq, the TtgR gene was removed from the TtgR_ColE1_SPS_V2 plasmid. The backbone was amplified with primers that had complementary overlap with the sfGFP gene. The sfGFP gene was amplified with primers complementary to the backbone. The BsaI cut sites and early stop codons were inserted into sfGFP in the same fashion as the creation of TtgR_ColE1_SPS_V5. The plasmid was labeled TtgR_ColE1_SPS_V3_GFPControl. Three pre-defined 20nt barcodes (AAACCCTGTGCCAGAGGGTG, GAGTGACCTTAAGTCAGGGA, and GCTTCTGTCCAAGCAGGTTA) were generated according to standard protocols. The barcodes were inserted into the TtgR_ColE1_SPS_V3_GFPControl using Golden Gate assembly.

RNA preparation for Sensor-seq

For the pilot study, cells containing 8 different TtgR variants were streaked out on an LB-Kan plate and grown overnight at 37 °C. Three colonies were inoculated into LB-kanamycin for overnight growth. The overnight cultures were diluted 50× into fresh LB-kanamycin containing either 1 mM naringenin or DMSO as a control. DH10B containing each barcoded GFP Control plasmid was struck out on LB-kanamycin plates. One colony was selected from each barcoded DH10B and grown in 3 mL LB-kanamycin overnight. These barcoded control cultures were combined in equal ratio and added to the test library culture to a final composition of 0.25% control. The induced cultures were grown and prepared following standard protocols. The cultures were grown at 37 °C shaking at 250 rpm in an Innova 4230 (New Brunswick Scientific). At the targeted OD₆₀₀, cultures were placed on ice for 10 min. 5×10^8 cells were harvested by centrifugation at $5500 \times g$ based on the OD₆₀₀ and the assumption that 1.0 OD₆₀₀ cultures have 8×10^8 cells/mL. The pelleted cells were decanted and stored at −80 °C. This process was repeated in biological triplicate for each target OD₆₀₀ with new colonies. RNA was purified from cell pellets via Trizol reagent (Invitrogen, 15596026). 1 mL of Trizol reagent was added to each cell pellet and vortexed briefly. The samples were incubated at room temperature for 5 min. 200 µL of chloroform was added to each sample. The samples were incubated at room temperature for 2 min and were centrifuged at $12,000 \times g$ for 15 min at 4 °C. 300 µL of the aqueous phase was transferred to a clean 2 mL centrifuge tube and placed on ice. RNA was purified from the aqueous phase using the RNA Clean and Concentrator 5 kit (Zymo, R1014) and eluted in 15 µL Ultrapure RNase-free dH₂O (Invitrogen, 10977015). The purified RNA was digested using 4U DNase I (New England Biolabs, M0303L) in a 50 µL reaction incubated at 37 °C for 30 min. The digestion reactions were purified using the RNA Clean and Concentrator 5 kit and eluted in 15 µL Ultrapure RNase-free dH₂O. Concentrations were measured using a Nanodrop instrument (Thermo Fisher).

To test the agnostic library, the four pools were grown individually in 5 mL LB-kanamycin overnight in triplicate. The four pools were combined prior to inoculation in 25 mL LB-kanamycin for the RNA harvest. GFP Control barcoded cells were spiked into the combined agnostic replicates at a final concentration of 0.25%. The same pooled replicates were used for all ligand inductions. The agnostic libraries were induced with the ligand (Supplementary Table 1). DMSO, dH₂O, and EtOH were included as solvent controls. No more than 2% v/v (DMSO) or 1% v/v (EtOH and H₂O) of solvent were tolerated. These cultures were processed as previously done with the pilot library.

cDNA synthesis and sequencing

cDNA synthesis uses ~3 µg total RNA, a primer encoding a 16nt unique molecular identifier (UMI), and the Maxima H Minus Double-Stranded cDNA Synthesis Kit (Thermo Scientific, K2561). The cDNA is purified using the DNA Clean and Concentrator 5 kit (Zymo, D4004). The Illumina sequencing regions are added in 2 PCR reactions in the same manner as the MiSeq barcode-variant mapping reactions. Three sets of primers containing the Illumina sequencing primer and a predefined barcode (ATCG, CGAT, and GTCA) were used in the first PCR reaction to add the Illumina sequencing regions (11 cycles). One set of primers was used for each biological replicate. The first reaction is purified using the DNA Clean and Concentrator 5 kit. The second reaction uses 4 µL of the first reaction and primers that add i5 and i7 indices in 8 cycles. The final amplicons are purified again. All replicates were combined in an equal molar ratio after purification. Plasmids are harvested from the remaining culture of the RNA preparation step using the ZR Plasmid Miniprep—Classic kit (Zymo, D4016). The UMI is added to the plasmid-derived samples in a 2-cycle PCR reaction using 100 ng of template. The amplification of all DNA libraries followed an identical protocol to the RNA preparation. The cDNA and DNA samples are sequenced using either an Illumina NovaSeq SP chip (test library) or a NovaSeq S4 chip (agnostic libraries) by the UWBC. Read volumes were calculated by targeting 500 reads per barcode with the assumption that 50% of the reads will be lost due to filtering criteria.

Computational data analysis

All bioinformatic and computational analyses were done through the UW Madison's Center for High Throughput Computing cluster or locally on a MacBook Pro with an M1 processor and 32 Gb of memory.

RNA-Seq data analysis

Fastq files were merged using NGmerge v0.3 and filtered using Fastp v0.23.1 based on average Q-score > Q30 for reads^{66,67}. Reads containing the 5' and 3' constant regions were isolated using UMI-Tools v0.2.3 and counted using Tally v15-065^{68,69}. Reads containing the central constant region were isolated and UMI sequences were removed with UMI-Tools. The barcodes were then counted with Tally. RNA-Seq barcodes were matched to mapped barcode-variant pairs with a Hamming distance tolerance of 1 using the seal script from BBMap v38.94 (<https://sourceforge.net/projects/bbmap/>). If a barcode mapped to more than one TtgR variant, then the TtgR variant that had the most reads was selected if each other variant was less than 10% of the reads of the most abundant variant. RNA-Seq barcodes that were successfully mapped to known barcode-variant pairs were analyzed across the induced RNA, induced DNA, control RNA, and control DNA samples. A barcode both had to be found in all four groups to be included in downstream analysis. The read counts for a variant were then a sum of the barcode counts for all barcodes mapped and found in all four datasets. No read count threshold was imposed during analysis. The fold enrichment calculation uses Eq. 1.

$$\text{Fold enrichment} = \frac{\frac{RNA_{+Lig}}{DNA_{+Lig}}}{\frac{RNA_{-Lig}}{DNA_{-Lig}}} \quad (1)$$

If biological replicates were available for each condition, the fold enrichment per variant was curated based on the coefficient of variation (CV). Percent deviation is calculated with Eq. 2.

$$CV = \frac{\sigma}{\bar{x}} \quad (2)$$

In this equation, σ is the standard deviation of the fold enrichment and \bar{x} is the mean fold enrichment across replicates. A 30% CV cutoff was imposed (Supplementary Fig. 3). All variants were normalized to

wildtype fold enrichment for each replicate. Heatmaps were constructed using the average performance of each variant after normalization. Variants with data passing CV thresholds for more than 5 ligands and performing at least 1.5 times better than wildtype were selected for clustering. Missing data was imputed using KNN methods in SciKit Learn v1.0.1⁷⁰. The UPGMA algorithm with a correlation distance metric and a target of 12 clusters was used to cluster in SciPy⁷¹. The number of clusters was selected by plotting the silhouette score against the number of clusters (Supplementary Fig. 5).

Mapping barcode-variant pairs

A 60 nt spacer was created to bring the random barcode and TtgR variants physically adjacent on the same plasmid to enable short-read next-generation sequencing mapping of barcode-variant pairs. The library plasmids were amplified with primers encoding BsaI cut sites that would remove the intervening region and allow for integration of the spacer between the barcode and the TtgR variant region. The spacer was inserted into the backbone using Golden Gate (New England Biolabs, BsaI-HFv2, E1601L) following manufacturer protocols.

Two primer groups were used to add Illumina sequencing regions to the barcode-spacer-variant region of the mapping plasmid libraries. Each primer group consisted of three primers with different numbers of Ns (0N, 3N, or 6N) to increase positional base diversity during runs. The adapter primers had complementarity to the plasmid and contained Illumina sequencing primer binding regions. Stem primers had the i7 and i5 indices and the adapter sequence to anneal to the sequencing flow cell. The adapter regions were added using 1 ng of template, 0.6 µL of 10 µM primers, and KAPA HiFi mix (Roche, KK2102) for 14 cycles. These reactions were purified using the DNA Clean and Concentrator 5 kit (Zymo, D4004). The stem primers were used in a second PCR reaction using 4 µL of the first reaction for 10 cycles.

The pilot library was sequenced on a 15M 2 × 250 MiSeq chip (Illumina). The agnostic library was sequenced using a 2 × 250 NovaSeq SP chip (Illumina). For MiSeq-based sequencing, the proper band was isolated using gel extraction on a 0.5% agarose gel followed by purification with the EZNA gel extraction kit (Omega BioTek, D2500-02). The concentration of the DNA was measured using AccuClear (Biotium, 31028) following manufacturer protocols. The flow cell was loaded with 15 pM DNA with 5% PhiX. For NovaSeq-based sequencing, samples were purified using PippinHT (Sage Science), and the concentration was measured via 4200 TapeStation (Agilent). The size selection, concentration measurement, and NovaSeq runs were performed by the University of Wisconsin Madison Biotechnology Center (UWBC).

The fastq output was merged using PEAR v0.9.11⁷². A C++ script was used to filter poor-scoring reads based on Q-scores. The C++ script performs quality filtering by summing the error probabilities from the PHRED score at each position across the sequence. If the summed error is below 10, then the merged read is accepted for further analysis. Reads that passed the quality filter were then filtered on constant regions surrounding the barcode and TtgR variants. Barcodes that had read counts greater than 10 and were unique for a single TtgR variant were mapped to that variant. If a single barcode maps to multiple sequences, the barcode is assigned to one variant if the total reads for that variant is greater than 90% of total reads for that barcode.

qRT-PCR validation of Sensor-seq

The abundance of the sfGFP and rrsA transcripts was measured via qRT-PCR. Each biological triplicate RNA was run in a technical triplicate in a MicroAmp Fast Optical 96-well plate (Applied Biosystems, 4346907). 1 ng of RNA was added to Luna Universal One-Step qRT-PCR mix (New England Biolabs, E3005L) containing 4 µmol of each primer on ice. The standard cycling protocol was used according to the manufacturer's suggestion. Each sample consisted of a set of reactions containing sfGFP-specific primers and another set containing rrsA-

specific primers. The reactions were run on a CFX Connect Real-Time PCR Detection System (BioRad). Fold enrichment was calculated using Eqs. 3 and 4. The error was propagated from the technical replicates and biological replicates using Eq. 5.

$$\text{Fold enrichment} = 2^{-\Delta\Delta C_t} \quad (3)$$

$$\Delta\Delta C_t = (C_{t\text{GFP}} - C_{t\text{rrsA}})_{\text{Ligand}} - (C_{t\text{GFP}} - C_{t\text{rrsA}})_{\text{-Ligand}} \quad (4)$$

$$\text{error} = \sqrt{\sum (\sigma_i)} \quad (5)$$

Cell sorting

An overnight culture is diluted 50× in phosphate-buffered saline (137 mM NaCl, 2.7 mM KCl, 10 mM Na₂HPO₄, 1.8 mM KH₂PO₄) and placed on ice for 10 min prior to sorting. Sorting was performed on an SH800 (Sony) using the 488 nm laser and a 525 ± 25 nm filter. Sorted cells were grown for 1 h shaking at 37 °C in 5 mL LB. Kanamycin was added to a final concentration of 50 µg/mL and the culture was grown overnight. An aliquot of the sorted culture was stored at −80 °C in 25% glycerol. Plasmids were isolated from the remaining culture using the ZR Plasmid Miniprep—Classic kit (Zymo, D4016).

Secondary validation of top hits

Top performing variants were selected based on the mean rank of each variant across the three biological replicates. These variants were encoded in gene fragments (Twist) and synthesized in a 96-well plate format. The fragments were resuspended to a final concentration of 10 ng/µL, pooled together, and cloned into the TtgR_ColE1_SPS_V2 backbone using Golden Gate Assembly. The resulting library was sorted based on fluorescence. LB-kanamycin is inoculated with 50 µL of the frozen stock of the library and grown overnight shaking at 37 °C. Sorting was performed according to the Cell Sorting protocol. 500,000 cells were isolated from the lower 70% of the population based on fluorescence. Plasmids were isolated from the remaining culture using the ZR Plasmid Miniprep—Classic kit (Zymo, D4016). DH10B electrocompetent cells (New England Biolabs, C3020K) were transformed with the purified plasmid library according to the Library Creation protocol.

LB-kanamycin is inoculated with 50 µL of the frozen stock of the repressed library and grown overnight, shaking at 37 °C. The culture was diluted 50× into fresh LB-kanamycin and grown overnight at 37 °C shaking with the ligands (Supplementary Table 1). Sorting was performed according to the Cell Sorting protocol. 400,000 cells were isolated using a gate that encompassed the top 0.5% of the population based on the fluorescence distribution in the absence of any ligand. Plasmids were isolated from the remaining culture using the ZR Plasmid Miniprep—Classic kit (Zymo, D4016).

The abundance of variants was determined using next-generation sequencing. Sequencing amplicons were generated using primers that had complementarity to the TtgR gene around the barcode insertion site. The amplification process followed the Barcode-variant mapping via next-generation sequencing protocol. The concentration of the DNA was measured using Qubit Fluorometric Quantification (Thermo Fisher) following manufacturer protocols. The flow cell was loaded with 15 pM DNA with 5% PhiX. Sequencing was performed on a MiSeq instrument (Illumina). Fastq files were merged using NGmerge and filtered using Fastp based on average Q-score > Q30 for reads^{42,43}. The abundances of each variant in each sample, including the repressed population, were normalized to the total reads for that sample. The fold enrichment is the ratio of the induced and repressed normalized abundances.

Tertiary validation of top hits

Top three performing variants for each ligand were selected based on the secondary validation. If a top variant for a ligand exists in the current selection, the next best variant is selected until each ligand has three top performers. These variants are individually cloned into the TtgR_ColE1_SPS_V2 backbone using Golden Gate Assembly. 3 µL of overnight culture for each clone is inoculated into 144 µL LB-kanamycin spiked with 3 µL of ligand or vehicle control in triplicate such that the final working concentration is as described in Supplementary Table 1. These cultures were grown at 37 °C shaking overnight. Induced overnight cultures were diluted 50× in 1× PBS. 10²–10³ cells were collected using an Attune NxT Flow Cytometer with Auto-sampler (Thermo Fisher). The median fluorescence of the collected cells for each sample was determined. The fold induction is calculated as the median fluorescence in the presence of ligand divided by the median fluorescence in the absence of ligand.

Unsupervised learning of ligand-agnostic dataset

The 11 possible mutated positions of the TtgR agnostic library were each physicochemically encoded in a 19-length dimensionally-reduced AAindex for a total of 209 features for each variant^{44,45}. The variants were embedded in a 2-dimensional projection and subsequently clustered using a combination of UMAP and HDBSCAN^{43,47}. Hyperparameter space (n_neighbors: [15, 25, 50, 75, 100, 300, 500, 1000, 2000], n_components: [2], min_dist: [0, 0.1], metric: [manhattan], min_cluster_size: range(5, 110, step=5)) was randomly searched with 100 iterations to find the optimal hyperparameters that minimized a cost function. Here, the cost function is the percentage of points with <5% certainty of assignment. The best combination was found to be n_neighbors=500, min_dist=0.1, and min_cluster_size=80, which minimized the cost function to 0.03144 (i.e., ~3% of points had <5% certainty of assignment). These hyperparameters were used to cluster the points in the associated UMAP embedding.

Top three clusters for each ligand were chosen for further analysis. Cluster rankings were determined by the average F-score of hits (variants with F-score ≥ 1.5). In order to be considered as a top three cluster, we also require the cluster to have a minimum of 15 hits. To determine the enrichment of substitutions at each mutated position from these top clusters, we calculated the ratio of the F-score-weighted substitution frequency within hits with the DNA-count-normalized substitution frequency. The F-score-weighted substitution frequency is determined by finding the frequency of a substitution in a set of variants where each variant's abundance is weighted by its average F-score from three replicates. The DNA-count-normalized substitution frequency measures the average frequency of each substitution in the induced plasmid library where the abundance of each variant is dictated by the number of sequenced DNA barcodes associated with that variant and is used to control for bias that may exist in the plasmid library. Thus, the ratio of these two values describes the change in frequency of each substitution such that large values signify enrichment of a substitution in the presence of a ligand.

Cell-free extract preparation

Cell-free extract was prepared as previously described with some modifications⁵⁴. Briefly, 1 L of 2X YT + P media (16 g/L tryptone, 10 g/L yeast extract, 5 g/L NaCl, 7 g/L potassium phosphate dibasic, 3 g/L potassium phosphate monobasic) was inoculated with 20 mL of saturated overnight culture of BL21 Star (DE3) (Invitrogen, C601003) in LB and grown to optical density 2.6. The cells were harvested by centrifugation at 5000 × g for 10 min, resuspended and washed in Buffer B (14 mM Mg-glutamate, 60 mM K-glutamate, 5 mM Tris, pH 8.2) three times, then resuspended to a final concentration of 1 g/mL cell in Buffer B. The suspension was lysed using a QSonica Q125 sonicator with a 3.175 mm diameter probe at a frequency of 20 kHz and 50% amplitude by 10 s ON/OFF pulses until the lysed suspensions turned

brown and became less viscous (around 60 s and delivering ~350 J). The lysate was clarified by a 10-min centrifugation at $12,000 \times g$ and 4°C . The supernatant was removed and incubated, shaking at 220 rpm for 80 min at 37°C for the ribosomal runoff reaction. After a second $12,000 \times g$ spin at 4°C for 10 min, the supernatant from the runoff was dialyzed against Buffer B overnight at 4°C in a 3.5 K MWCO membrane. The dialysate was removed, centrifuged once more, aliquoted, and flash-frozen on liquid nitrogen for long-term storage at -80°C .

Cell-free gene expression reaction

CFE reactions were prepared as previously described⁵⁴. Briefly, the overall reaction composition was 6 mM magnesium glutamate; 10 mM ammonium glutamate; 130 mM potassium glutamate; 1.2 mM ATP; 0.850 mM each of GTP, UTP, and CTP; 0.034 mg/mL folinic acid; 0.171 mg/mL yeast tRNA; 2 mM amino acids; 30 mM PEP; 0.33 mM NAD; 0.27 mM CoA; 4 mM oxalic acid; 1 mM putrescine; 1.5 mM spermidine; 57 mM HEPES; 30% extract by volume; plasmid DNA to the desired concentration and water. 10 μL reactions were mixed on ice in replicates and then pipetted onto a black Corning clear-bottom 384-well plate for measurement of sfGFP (excitation/emission 470/510 nm) on a Biotek Synergy HI plate reader every 10 min at 30°C .

Dose-response assay

Overnight cultures of cells containing a GFP reporter and WT TtgR, 3A7, 2F9, or 3C8 were back-diluted by adding 3 μL of overnight culture and 3 μL of serially diluted ligand (naltrexone, naringenin, or quinine) to 144 μL of LB + kanamycin (50 $\mu\text{g}/\text{mL}$) in a 96-well plate. GFP (excitation/emission 470/510 nm) fluorescence and optical density measurements were obtained every 10 min at 37°C on a Biotek Synergy HTX plate reader.

X-ray crystallography

TtgR-pET31b variants were expressed in BL21 cells (NEB) as previously described with a few minor differences³³. Briefly, 1 L of TB-M-80155 using the following recipe: TB (24 g l^{-1} yeast extract, 12 g l^{-1} tryptone), 50 mL of $20 \times \text{M}$ (1 M NH_4Cl , 0.5 M $\text{Na}_2\text{H}(\text{PO}_4)$, 0.5 M $\text{KH}_2(\text{PO}_4)$, 0.1 M $\text{Na}_2(\text{SO}_4)$) per liter, autoclaved and supplemented with 25 mL of sterile filtered 40×80155 autoinduction solution (32% (w/v) glycerol, 0.6% (w/v) glucose, 2% (w/v) lactose) per liter of media along with 50 $\mu\text{g}/\text{mL}$ kanamycin⁷³ was inoculated from a starter culture and grown at 25°C for 24 h. The cells were harvested by centrifuging 20 min at $6250 \times g$. To purify, cells were lysed using an M110 Microfluidizer (Microfluidics International Corporation), as described with the following modifications³³. The lysate was centrifuged at $75,000 \times g$ for 45 min at 4°C and passed through a 0.8 μm filter. The lysate was then applied to a 5 mL HisTrap HP column (Cytiva) using an Akta purifier. The column was washed with 20 column volumes (CV) IMAC-A (500 mM NaCl, 25 mM Imidazole, 25 mM HEPES pH 7.5, 0.3 mM TCEP). MBP-6His-TtgR was eluted with a gradient of 100% IMAC-A to 100% IMAC-B (IMAC-A with 350 mM imidazole) over 10 CV and collected in 5 mL fractions. Peak fractions with the highest absorbance at 280 nm were analyzed using SDS-PAGE and combined for simultaneous dialysis and TEV cleavage. TEV cleavage was carried out using a 10 kDa MWCO snake-skin dialysis bag (Thermo) at 1:50 mg/mg TEV:TtgR in 5 L of IMAC-A. Dialysis occurred over a 16 h interval at 4°C while stirring at low speed. Following cleavage MBP and TEV were removed by passing over Ni-NTA resin (Pierce) equilibrated with IMAC-A. The flow through was concentrated and applied to a 16/60 S200 size exclusion column (Cytiva) equilibrated with 20 mM HEPES pH 7.5, 100 mM NaCl, and 0.3 mM TCEP. Peak fractions were analyzed using SDS-PAGE. Pure fractions were combined, concentrated immediately used for crystallization trials, and/or frozen in liquid nitrogen and stored at -80°C for later usage.

Crystals were screened based on conditions for TtgR with resveratrol (PDB 7K1A, 7K1C, and 7KD8)³³, however, only variant 3A7

yielded diffraction quality crystals in 100 mM Bis-Tris pH 6.5 18–20% MEPEG 2000, 200 mM MgSO_4 after extensive effort. To obtain naltrexone in the active site, crystals were both co-crystallized with 5 mM naltrexone and prior to cryoprotection, were incubated for 1 h in cryoprotectant (100 mM Bis-Tris pH 6.5 18–20% MEPEG 2000, 200 mM MgSO_4 20% glycerol) with saturating naltrexone and frozen in liquid N_2 .

X-ray diffraction data were collected at Advanced Photon Source (APS) beamlines LS-CAT ID-D. Diffraction data were reduced and scaled using XDS (Build = 20220820) and autoPROC (Build = 20230217)^{74,75}. The structure was solved by molecular replacement with Phenix.phaser⁷⁶ using PDB ID 7K1C³³. The model was manually inspected and built through iterative rounds of Phenix.refine (Phenix_1.20.1_4487) and manual inspection in COOT (0.9.8)⁷⁷.

Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

Data availability

The processed Sensor-seq datasets generated in this study have been deposited in Zenodo [<https://doi.org/10.5281/zenodo.13381000>]. The raw RNA-seq and amplicon sequencing files generated in this study have been deposited in the Sequence Read Archive (SRA) under accession code [SRP540614](https://www.ncbi.nlm.nih.gov/sra/SRP540614). The crystallography data generated in this study for the naltrexone-bound 3A7 variant are available in the RCSB Protein Data Bank (PDB) under accession code [8V90](https://www.rcsb.org/structure/8V90). The resveratrol-bound TtgR structure used for comparisons with the 3A7 mutant is available in the RCSB PDB under accession code [7K1C](https://www.rcsb.org/structure/7K1C). Source data are provided with this paper.

Code availability

Scripts used to analyze the Sensor-seq data are available at Github [https://github.com/raman-lab/sensor_seq].

References

- Cuthbertson, L. & Nodwell, J. R. The TetR family of regulators. *Microbiol. Mol. Biol. Rev.* **77**, 440–475 (2013).
- Gong, Z., Li, H., Cai, Y., Stojkoska, A. & Xie, J. Biology of MarR family transcription factors and implications for targets of antibiotics against tuberculosis. *J. Cell. Physiol.* **234**, 19237–19248 (2019).
- Meinhardt, S. et al. Novel insights from hybrid LacI/GalR proteins: family-wide functional attributes and biologically significant variation in transcription repression. *Nucleic Acids Res.* **40**, 11139–11154 (2012).
- Mitchler, M. M., Garcia, J. M., Montero, N. E. & Williams, G. J. Transcription factor-based biosensors: a molecular-guided approach for natural product engineering. *Curr. Opin. Biotechnol.* **69**, 172–181 (2021).
- Tellechea-Luzardo, J., Stiebritz, M. T., & Carbonell, P. Transcription factor-based biosensors for screening and dynamic regulation. *Front. Bioeng. Biotechnol.* **11**, 1118702 (2023).
- Rogers, J. K., Taylor, N. D. & Church, G. M. Biosensor-based engineering of biosynthetic pathways. *Curr. Opin. Biotechnol.* **42**, 84–91 (2016).
- Zhang, F., Carothers, J. M. & Keasling, J. D. Design of a dynamic sensor-regulator system for production of chemicals and fuels derived from fatty acids. *Nat. Biotechnol.* **30**, 354–359 (2012).
- Dietrich, J. A., McKee, A. E. & Keasling, J. D. High-throughput metabolic engineering: advances in small-molecule screening and selection. *Annu. Rev. Biochem.* **79**, 563–590 (2010).
- Mustafi, N., Grünberger, A., Kohlheyer, D., Bott, M. & Frunzke, J. The development and application of a single-cell biosensor for the detection of L-methionine and branched-chain amino acids. *Metab. Eng.* **14**, 449–457 (2012).

10. Xu, P., Li, L., Zhang, F., Stephanopoulos, G. & Koffas, M. Improving fatty acids production by engineering dynamic pathway regulation and metabolic control. *Proc. Natl. Acad. Sci. USA* **111**, 11299–11304 (2014).
11. Dietrich, J. A., Shis, D. L., Alikhani, A. & Keasling, J. D. Transcription factor-based screens and synthetic selections for microbial small-molecule biosynthesis. *ACS Synth. Biol.* **2**, 47–58 (2013).
12. Raman, S., Rogers, J. K., Taylor, N. D. & Church, G. M. Evolution-guided optimization of biosynthetic pathways. *Proc. Natl. Acad. Sci. USA* **111**, 17803–17808 (2014).
13. d'Oelsnitz, S., Love, J. D., Diaz, D. J. & Ellington, A. D. GroovDB: a database of ligand-inducible transcription factors. *ACS Synth. Biol.* **11**, 3534–3537 (2022).
14. Mahr, R. & Frunzke, J. Transcription factor-based biosensors in biotechnology: current state and future prospects. *Appl. Microbiol. Biotechnol.* **100**, 79–90 (2016).
15. Münch, R. et al. PRODORIC: prokaryotic database of gene regulation. *Nucleic Acids Res.* **31**, 266–269 (2003).
16. Novichkov, P. S. et al. RegPrecise 3.0—a resource for genome-scale exploration of transcriptional regulation in bacteria. *BMC Genom.* **14**, 745 (2013).
17. Taylor, N. D. et al. Engineering an allosteric transcription factor to respond to new ligands. *Nat. Methods* **13**, 177–183 (2016).
18. Wise, A. A. & Kuske, C. R. Generation of novel bacterial regulatory proteins that detect priority pollutant phenols. *Appl. Environ. Microbiol.* **66**, 163–169 (2000).
19. Galvão, T. C., Mencia, M. & de Lorenzo, V. Emergence of novel functions in transcriptional regulators by regression to stem protein types. *Mol. Microbiol.* **65**, 907–919 (2007).
20. Tang, S.-Y. & Cirino, P. C. Design and application of a mevalonate-responsive regulatory protein. *Angew. Chem. Int. Ed. Engl.* **50**, 1084–1086 (2011).
21. Süel, G. M., Lockless, S. W., Wall, M. A. & Ranganathan, R. Evolutionarily conserved networks of residues mediate allosteric communication in proteins. *Nat. Struct. Biol.* **10**, 59–69 (2003).
22. Leander, M., Yuan, Y., Meger, A., Cui, Q. & Raman, S. Functional plasticity and evolutionary adaptation of allosteric regulation. *Proc. Natl. Acad. Sci. USA* **117**, 25445–25454 (2020).
23. Leander, M., Liu, Z., Cui, Q. & Raman, S. Deep mutational scanning and machine learning reveal structural and molecular rules governing allosteric hotspots in homologous proteins. *eLife* **11**, e79932 (2022).
24. Chen, J., Vishweshwaraiah, Y. L. & Dokholyan, N. V. Design and engineering of allosteric communications in proteins. *Curr. Opin. Struct. Biol.* **73**, 102334 (2022).
25. Rohlhill, J., Sandoval, N. R. & Papoutsakis, E. T. Sort-seq approach to engineering a formaldehyde-inducible promoter for dynamically regulated *Escherichia coli* growth on methanol. *ACS Synth. Biol.* **6**, 1584–1595 (2017).
26. Peterman, N. & Levine, E. Sort-seq under the hood: implications of design choices on large-scale characterization of sequence-function relations. *BMC Genom.* **17**, 206 (2016).
27. Ding, N. et al. Programmable cross-ribosome-binding sites to fine-tune the dynamic range of transcription factor-based biosensor. *Nucleic Acids Res.* **48**, 10602–10613 (2020).
28. Li, J.-W., Zhang, X.-Y., Wu, H. & Bai, Y.-P. Transcription factor engineering for high-throughput strain evolution and organic acid bio-production: a review. *Front. Bioeng. Biotechnol.* **8**, 98 (2020).
29. Stainbrook, S. C. & Tyo, K. E. J. Model-guided mechanism discovery and parameter selection for directed evolution. *Appl. Microbiol. Biotechnol.* **103**, 9697–9709 (2019).
30. Aharoni, A. et al. The “evolvability” of promiscuous protein functions. *Nat. Genet.* **37**, 73–76 (2005).
31. Buda, K., Miton, C. M., Fan, X. C. & Tokuriki, N. Molecular determinants of protein evolvability. *Trends Biochem. Sci.* **48**, 751–760 (2023).
32. Alguet, Y. et al. Crystal structures of multidrug binding protein TtgR in complex with antibiotics and plant antimicrobials. *J. Mol. Biol.* **369**, 829–840 (2007).
33. Nishikawa, K. K., Hoppe, N., Smith, R., Bingman, C. & Raman, S. Epistasis shapes the fitness landscape of an allosteric specificity switch. *Nat. Commun.* **12**, 5562 (2021).
34. Khersonsky, O. et al. Automated design of efficient and functionally diverse enzyme repertoires. *Mol. Cell* **72**, 178–186.e5 (2018).
35. Alford, R. F. et al. The Rosetta all-atom energy function for macromolecular modeling and design. *J. Chem. Theory Comput.* **13**, 3031–3048 (2017).
36. Jha, R. K., Chakraborti, S., Kern, T. L., Fox, D. T. & Strauss, C. E. M. Rosetta comparative modeling for library design: engineering alternative inducer specificity in a transcription factor. *Proteins* **83**, 1327–1340 (2015).
37. Leaver-Fay, A. et al. ROSETTA3: an object-oriented software suite for the simulation and design of macromolecules. *Methods Enzymol.* **487**, 545–574 (2011).
38. Johnson, M. D. et al. Pharmacological characterization of 4-hydroxy-*N*-desmethyl tamoxifen, a novel active metabolite of tamoxifen. *Breast Cancer Res. Treat.* **85**, 151–159 (2004).
39. Lien, E. A., Solheim, E. & Ueland, P. M. Distribution of tamoxifen and its metabolites in rat and human tissues during steady-state treatment. *Cancer Res.* **51**, 4837–4844 (1991).
40. Achan, J. et al. Quinine, an old anti-malarial drug in a modern world: role in the treatment of malaria. *Malar. J.* **10**, 144 (2011).
41. Niciu, M. J. & Arias, A. J. Targeted opioid receptor antagonists in the treatment of alcohol use disorders. *CNS Drugs* **27**, 777–787 (2013).
42. Sharifi-Rad, J. et al. Ellagic acid: a review on its natural sources, chemical stability, and therapeutic potential. *Oxid. Med. Cell. Longev.* **2022**, e3848084 (2022).
43. McInnes, L., Healy, J., Saul, N. & Großberger, L. UMAP: uniform manifold approximation and projection. *J. Open Source Softw.* **3**, 861 (2018).
44. Kawashima, S., Ogata, H. & Kanehisa, M. AAindex: amino acid index database. *Nucleic Acids Res.* **27**, 368–369 (1999).
45. Gelman, S., Fahlberg, S. A., Heinzelman, P., Romero, P. A. & Gitter, A. Neural networks to learn protein sequence–function relationships from deep mutational scanning data. *Proc. Natl. Acad. Sci. USA* **118**, e2104878118 (2021).
46. Jones, E. M. et al. Structural and functional characterization of G protein-coupled receptors with deep mutational scanning. *eLife* **9**, e54895 (2020).
47. Campello, R. J. G. B., Moulavi, D., & Sander, J. Density-based clustering based on hierarchical density estimates. in *Advances in Knowledge Discovery and Data Mining* (eds Pei, J., Tseng, V. S., Cao, L., Motoda, H., & Xu, G.) *Lecture Notes in Computer Science* 160–172. https://doi.org/10.1007/978-3-642-37456-2_14 (Springer, 2013).
48. Valley, C. C. et al. The methionine-aromatic motif plays a unique role in stabilizing protein structure. *J. Biol. Chem.* **287**, 34979–34991 (2012).
49. Weber, D. S. & Warren, J. J. The interaction between methionine and two aromatic amino acids is an abundant and multifunctional motif in proteins. *Arch. Biochem. Biophys.* **672**, 108053 (2019).
50. Liu, X. et al. Design of a transcriptional biosensor for the portable, on-demand detection of cyanuric acid. *ACS Synth. Biol.* **9**, 84–94 (2020).
51. Jung, J. K. et al. Cell-free biosensors for rapid detection of water contaminants. *Nat. Biotechnol.* **38**, 1451–1459 (2020).
52. Li, Z. & Wang, P. Point-of-care drug of abuse testing in the opioid epidemic. *Arch. Pathol. Lab. Med.* **144**, 1325–1334 (2020).
53. Weerts, E. M. et al. Differences in delta- and mu-opioid receptor blockade measured by positron emission tomography in naltrexone-treated recently abstinent alcohol-dependent subjects. *Neuropsychopharmacology* **33**, 653–665 (2008).

54. Silverman, A. D., Kelley-Loughnane, N., Lucks, J. B. & Jewett, M. C. Deconstructing cell-free extract preparation for in vitro activation of transcriptional genetic circuitry. *ACS Synth. Biol.* **8**, 403–414 (2019).
55. Tappin, A. D., Loughnane, J. P., McCarthy, A. J. & Fitzsimons, M. F. Unexpected removal of the most neutral cationic pharmaceutical in river waters. *Environ. Chem. Lett.* **14**, 455–465 (2016).
56. Punihaole, D. et al. New insights into quinine–DNA binding using Raman spectroscopy and molecular dynamics simulations. *J. Phys. Chem. B* **122**, 9840–9851 (2018).
57. Dauparas, J. et al. Robust deep learning-based protein sequence design using ProteinMPNN. *Science* **378**, 49–56 (2022).
58. Huang, L. et al. Bacterial multidrug efflux pumps at the frontline of antimicrobial resistance: an overview. *Antibiotics* **11**, 520 (2022).
59. Anderson, D. W., McKeown, A. N. & Thornton, J. W. Intermolecular epistasis shaped the function and evolution of an ancient transcription factor and its DNA binding sites. *eLife* **4**, e07864 (2015).
60. Fernandez-Escamilla, A. M., Fernandez-Ballester, G., Morel, B., Casares-Atienza, S. & Ramos, J. L. Molecular binding mechanism of TtgR repressor to antibiotics and antimicrobials. *PLoS ONE* **10**, e0138469 (2015).
61. Pandurangan, A. P. & Blundell, T. L. Prediction of Impacts of mutations on protein structure and interactions: SDM, a Statistical approach, and mCSM, using machine learning. *Protein Sci.* **29**, 247–257 (2020).
62. Hopf, T. A. et al. Mutation effects predicted from sequence co-variation. *Nat. Biotechnol.* **35**, 128–135 (2017).
63. Mansoor, S., Baek, M., Juergens, D., Watson, J. L. & Baker, D. Zero-shot mutation effect prediction on protein stability and function using RoseTTAFold. *Protein Sci. Publ. Protein Soc.* **32**, e4780 (2023).
64. Riesselman, A. J., Ingraham, J. B. & Marks, D. S. Deep generative models of genetic variation capture the effects of mutations. *Nat. Methods* **15**, 816–822 (2018).
65. Kosuri, S. et al. Composability of regulatory sequences controlling transcription and translation in *Escherichia Coli*. *Proc. Natl. Acad. Sci. USA* **110**, 14024–14029 (2013).
66. Gaspar, J. M. NGmerge: merging paired-end reads via novel empirically-derived models of sequencing errors. *BMC Bioinform.* **19**, 536 (2018).
67. Chen, S., Zhou, Y., Chen, Y. & Gu, J. Fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics* **34**, i884–i890 (2018).
68. Smith, T., Heger, A. & Sudbery, I. UMI-tools: modeling sequencing errors in unique molecular identifiers to improve quantification accuracy. *Genome Res.* **27**, 491–499 (2017).
69. Davis, M. P. A., van Dongen, S., Abreu-Goodger, C., Bartonicek, N. & Enright, A. J. Kraken: a set of tools for quality control and analysis of high-throughput sequence data. *Methods* **63**, 41–49 (2013).
70. Pedregosa, F. et al. Scikit-learn: machine learning in Python. *J. Mach. Learn. Res.* **12**, 2825–2830 (2011).
71. Virtanen, P. et al. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat. Methods* **17**, 261–272 (2020).
72. Zhang, J., Kobert, K., Flouri, T. & Stamatakis, A. PEAR: a fast and accurate Illumina Paired-End reAd mergeR. *Bioinformatics* **30**, 614–620 (2014).
73. Acheson, J. F., Ho, R., Goularte, N. F., Cegelski, L. & Zimmer, J. Molecular organization of the *E. Coli* cellulose synthase macro-complex. *Nat. Struct. Mol. Biol.* **28**, 310–318 (2021).
74. Kabsch, W. XDS. *Acta Crystallogr. D Biol. Crystallogr.* **66**, 125–132 (2010).
75. Vonrhein, C. et al. Data processing and analysis with the autoPROC toolbox. *Acta Crystallogr. D Biol. Crystallogr.* **67**, 293–302 (2011).
76. Liebschner, D. et al. Macromolecular structure determination using X-rays, neutrons and electrons: recent developments in Phenix. *Acta Crystallogr. D Struct. Biol.* **75**, 861–877 (2019).
77. Emsley, P. & Cowtan, K. Coot: model-building tools for molecular graphics. *Acta Crystallogr. D Biol. Crystallogr.* **60**, 2126–2132 (2004).

Acknowledgements

This was supported by United States Army Research Office Grants W911NF20C0005 and W911NF1710043 (S.R.). K.K.N. was supported by NIH National Research Service Award T32 GM07215 (Molecular Biophysics Training Program) and the Robert and Katherine Burris Biochemistry Fund. J.C. and N.N. were supported by the National Institute of General Medical Sciences of the National Institute of Health under Award Number T32GM135066 (Biotechnology Training Program). S.H. and J.L.C. were supported by the AFRL 711th Human Performance Wing. The use of the Advanced Photon Source was supported by the U.S. Department of Energy, Basic Energy Sciences, Office of Science, under contract No. DE-AC02-06CH11357. LS-CAT Sector 21 was supported by the Michigan Economic Development Corporation and the Michigan Technology Tri-Corridor (Grant 085P1000817). The use of Life Science Collaborative Access Team at the Advanced Photon Source was supported by the College of Agricultural and Life Sciences, Department of Biochemistry, and the Office of the Vice Chancellor for Research and Graduate Education of the University of Wisconsin–Madison. Crystallization, data collection, and structure solution and refinement were supported by gift funds to B.G.F. Some sequencing was performed at the University of Wisconsin–Madison Biotechnology Center’s DNA Sequencing Facility (Research Resource Identifier—RRID:SCR_017759).

Author contributions

K.K.N. and S.R. designed the study, analyzed the data, and wrote the manuscript. K.K.N. assembled the libraries and performed the pilot screen, the RNA sequencing, and FACS. J.C. analyzed the data, performed the clonal flow cytometry validations, and wrote the manuscript. J.F.A., H.R.S., E.C.L., and D.H.L. performed the protein purification and crystallization. S.V.H. performed the cell-free experiments. P.H., M.F., and N.N. provided helpful advice in experimental design and data analysis. J.L.C., B.G.F., and S.R. supervised the study and acquired funding.

Competing interests

K.K.N., N.N., and S.R. have filed a provisional patent application with the title “Method of identifying allosteric biosensor proteins with new specificities” based on the results generated from this paper through the Wisconsin Alumni Research Foundation (P230343US02). The remaining authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41467-024-54260-8>.

Correspondence and requests for materials should be addressed to Srivatsan Raman.

Peer review information *Nature Communications* thanks Pablo Carbo-nell and the other, anonymous, reviewer(s) for their contribution to the peer review of this work. A peer review file is available.

Reprints and permissions information is available at <http://www.nature.com/reprints>

Publisher’s note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Open Access This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

© The Author(s) 2024