

1 **Title:** Encoding of odors by mammalian olfactory receptors

2 **Authors:** Aashutosh Vihani<sup>1</sup>, Maira H. Nagai<sup>2</sup>, Conan Juan<sup>2</sup>, Claire A. de March<sup>2</sup>, Xiaoyang S. Hu<sup>2</sup>, John  
3 Pearson<sup>1,3</sup>, Hiroaki Matsunami<sup>1,2,3\*</sup>

4 1. Department of Neurobiology, Neurobiology graduate program, Duke University Medical Center,  
5 Durham, NC 27710, USA

6 2. Department of Molecular Genetics and Microbiology, Duke University Medical Center, Durham,  
7 NC 27710, USA

8 3. Duke Institute for Brain Sciences, Duke University, Durham, NC 27710, USA

9 \*Correspondence to HM ([hiroaki.matsunami@duke.edu](mailto:hiroaki.matsunami@duke.edu))

10

11 **Contributions:**

12 AV, MHN, XSH, and HM did *in vivo* experiments. CJ and CAdM did *in vitro* experiments. AV analyzed  
13 data. CAdM did OR homology modeling. JP advised data modeling analysis. AV drafted the paper. All  
14 authors reviewed and edited the paper. HM supervised the work.

15

16 **Acknowledgements:** We thank Michael Schmuker for critical reading and comments, Mengjue Jessica  
17 Ni for expert technical assistance, Alex Koulakov for helpful discussions on data analysis, Priya Meesa  
18 for manuscript edits, and members of the Matsunami lab for helpful discussions. AV thanks Jack J. Zhan  
19 for advising how to analyze data in python during initial stages.

20

21 **Funding Sources:** This work was funded by NIH (DC014423 and DC016224) and NSF (1556207 and  
22 1555919) to HM, and NIH (DC018333) to CAdM.

- 23
- 24 **Competing Interests:** HM has received royalties from ChemCom, has received research grants from
- 25 Givaudan, and has received consultant fees from Kao.

26 **Abstract:**

- 27       **1.** Identified ligands for > 500 mouse ORs
- 28       **2.** ORs are specifically tuned towards individual odorants and their molecular properties
- 29       **3.** Odor molecular properties are informative of odor responses
- 30       **4.** Predictive modeling and convergent evolution analyses suggest specific residues within a
- 31               canonical location for odorant binding

32 Olfactory receptors (ORs) constitute the largest multi-gene family in the mammalian genome, with

33 hundreds to thousands of loci in humans and mice respectively<sup>1</sup>. The rapid expansion of this massive

34 family of genes has been generated by numerous duplication and diversification events throughout

35 evolutionary history. This size, similarity, and diversity has made it challenging to define the principles

36 by which ORs encode olfactory stimuli. Here, we performed a broad surveying of OR responses, using

37 an *in vivo* strategy, against a diverse panel of odorants. We then used the resulting interaction profiles

38 to uncover relationships between OR responses, odorants, odor molecular properties, and OR

39 sequences. Our data and analyses revealed that ORs generally exhibited sparse tuning towards

40 odorants and their molecular properties. Odor molecular property similarity between pairs of odorants

41 was informative of odor response similarity. Finally, ORs sharing response to an odorant possessed

42 amino acids at poorly conserved sites that exhibited both, predictive power towards odorant selectivity

43 and convergent evolution. The localization of these residues occurred primarily at the interface of the

44 upper halves of the transmembrane domains, implying that canonical positions govern odor selectivity

45 across ORs. Altogether, our results provide a basis for translating odorants into receptor neuron

46 responses for the unraveling of mammalian odor coding.

47 **Introduction:**

48 Stimulus encoding and feature extraction are fundamental tasks performed by all sensory systems.  
49 Therefore, a central problem in neurobiology is defining how aspects of a stimulus are represented by  
50 the activity of sensory receptors<sup>2-5</sup>. This problem is particularly intriguing in the case of olfactory  
51 stimuli, which do not vary along a single, continuous dimension, such as wavelength or amplitude.  
52 Odorants, rather, have discrete molecular structures that determine their physical-chemical properties.  
53 An inability to relate how these discrete molecular structures and their associated physical-chemical  
54 properties influence receptor responses represents a major gap in knowledge. Consequently, one  
55 cannot robustly predict the neural activity patterns nor the perceptual attributes<sup>6</sup> of an odorant  
56 starting from its physical-chemical properties.

57  
58 A major hindrance in deciphering the coding of olfactory information by olfactory receptors (ORs) has  
59 been the historic inability to comprehensively identify ORs that respond to an odorant. Various *in vivo*,  
60 *ex vivo*, and *in vitro* methods have generally suffered from either a lack of insight into receptor identity  
61 or have been too low throughput for a comprehensive surveying of OR selectivity<sup>5,7-9</sup>. With the mouse  
62 genome encoding over 1000 intact ORs, and odor reception following a combinatorial coding scheme,  
63 where one OR can be activated by a set of odorants and one odorant can activate a combination of  
64 ORs, defining a logic for peripheral odor coding is dependent on a comprehensive surveying while  
65 tracking receptor identity over a large odor panel<sup>1,10-12</sup>.

66  
67 Here, we performed a broad surveying of odorants *in vivo* to identify odorant-OR interactions in *Mus*  
68 *musculus*. By leveraging phosphorylated S6 ribosomal subunit capture (pS6-IP) coupled to RNA-Seq  
69 (pS6-IP-Seq), we were able to identify ORs expressed by recently active olfactory sensory neurons  
70 (OSNs; receptor deorphanization)<sup>11,13-15</sup>. Then, using a library of molecular property descriptors, we

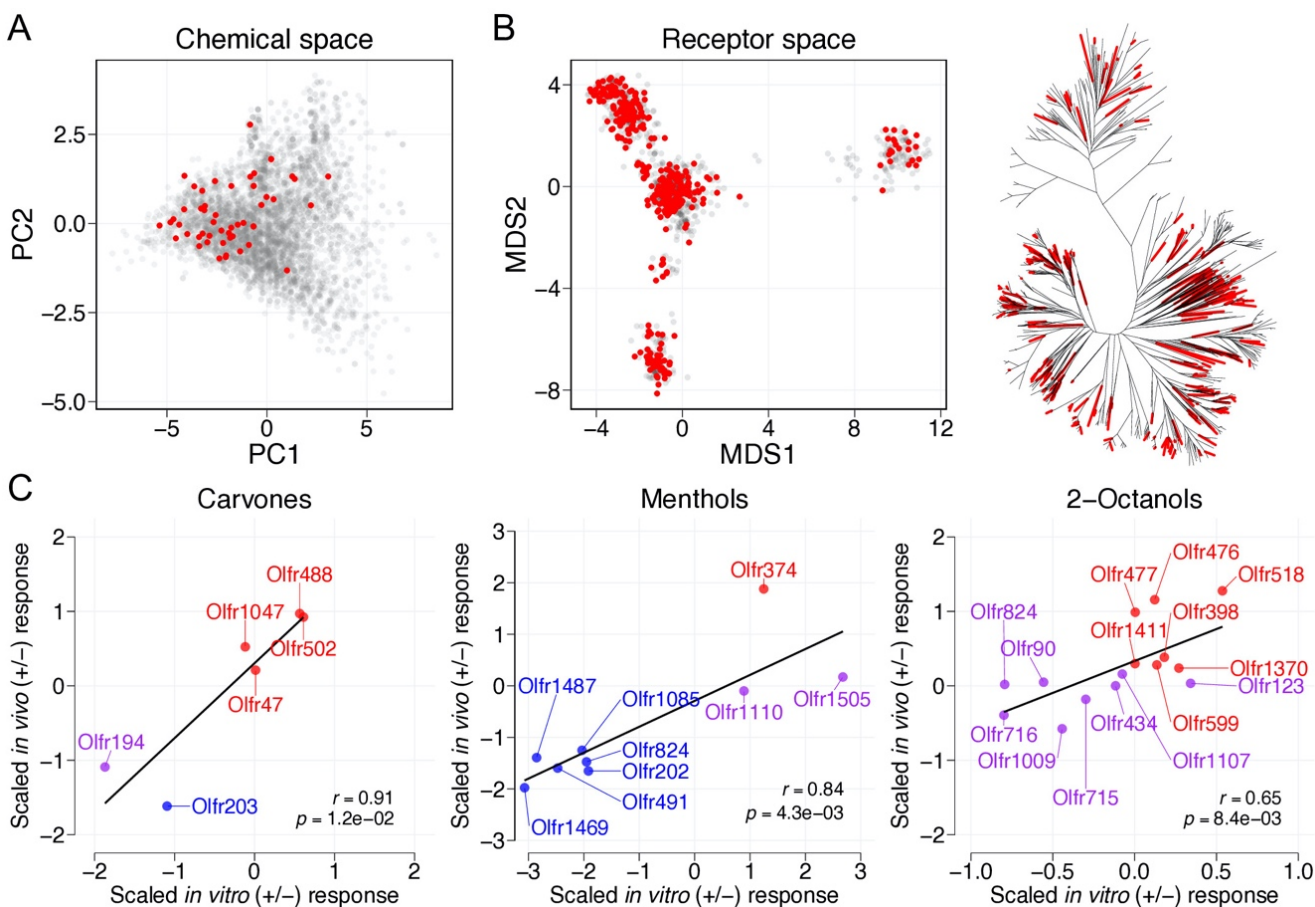
71 parameterized the physical-chemical properties of the tested odorants to uncover relationships to the  
72 responses they elicited from cognate receptors. Finally, using our data, we asked 1) how well does  
73 odor molecular property similarity predict receptor response similarity and 2) if there are specific  
74 amino acid positions that influence odorant selectivity amongst receptors. Our results and analyses  
75 provide a foundational framework for understanding the molecular logic by which the quality of an  
76 odor molecule is encoded across a mammalian receptor repertoire.

77 **Results:**

78 **Estimation of chemical and receptor space sampling**

79 First, we set out to identify ORs activated by a set of 61 odorants at various concentrations by  
80 leveraging pS6-IP-Seq. Immunoprecipitation of phosphorylated ribosomes from activated neurons  
81 followed by associated mRNA profiling by RNA-Seq, and differential expression analysis, enabled us to  
82 identify ORs expressed by OSNs activated by specific odorants (**Supplementary figure 1A**)<sup>11,14,15</sup>. ORs  
83 were considered odor-responsive if enrichment values ( $\log_2FC$ ) were positive with a false discovery  
84 rate (FDR) < 0.05. Considering all odorants at all tested concentrations, this approach deorphanized a  
85 total of 555 ORs across 72 conditions (**Supplementary table 1**). Considering unique odorants yielding at  
86 least one activated OR at the lowest tested concentration, this approach deorphanized a total of 375  
87 ORs across 52 odorants.

88  
89 To examine the bias in our odorant set, we built an 1811-dimensional (1811D) space in which each  
90 dimension represented a molecular property descriptor<sup>5</sup>, such as molecular weight, number of atoms,  
91 or aromatic ratio, parameterizing the physical-chemical properties of the odor molecule. We then  
92 plotted our 52 uniquely tested odorants together with 4680 other small molecules<sup>16</sup> other small  
93 molecules commonly found in foods and fragrances in this 1811D space to construct a chemical space  
94 consisting of a total of 4732 small molecules. Visualization of the first two principal components (PCs)  
95 did not reveal any obvious segregation of the test odorants, suggesting a broad sampling of chemical  
96 space by our test odor panel (**Figure 1A, Supplementary figure 1B-C**). To examine bias in our resulting  
97 deorphanized OR cohort, we computed pairwise OR Grantham distances<sup>17</sup>, an index of amino acid  
98 similarity, and visualized the results using multidimensional scaling (MDS). Examination of the first two  
99 MDS coordinates did not reveal any obvious segregation of the deorphanized 375 ORs, suggesting a  
100 broad sampling of receptor space (**Figure 1B**).



**Figure 1. Data bias and pS6-IP-Seq validation.** **A**, A total of 1811 molecular descriptors were calculated for a total of 4732 small molecules. The small molecules were projected onto a 2D chemical space made of the first and second principal components. The 52 unique odorants tested by pS6-IP-Seq at low concentrations, each yielding at least one activated OR, are colored in red. **B**, Left, Grantham distances were used to calculate a distance matrix for intact ORs. The matrix was visualized in two dimensions with multidimensional scaling to represent receptor space. ORs responding to at least one of the 52 tested odorants are colored red ( $n = 375$ ). Right, a phylogenetic tree of intact ORs. Tree edges with identified and unidentified agonists at low concentrations are colored red and black respectively. **C**, ORs responsive to tested enantiomers were evaluated by heterologous expression. ORs enriched by pS6-IP-Seq against (+)-odorant are colored red while ORs enriched by (-)-odorant are colored blue. ORs enriched by both enantiomers are colored purple. Linear regression reveals *in vivo* and *in vitro* responses to be highly correlated (carvones  $r = 0.91$ ,  $p = 0.012$ ; menthols  $r = 0.84$ ,  $p = 0.0043$ ; 2-octanols  $r = 0.65$ ,  $p = 0.0084$ ).

101

102 To validate the receptor specificity of the pS6-IP-Seq dataset, we selected enantiomers (carvones,

103 menthols, and 2-octanols) for *in vitro* testing. We transiently expressed ORs responsive to tested

104 enantiomers in Hana3A cells and challenged with individual odorants to generate dose response

105 curves. Comparison of *in vitro* responses to *in vivo* responses revealed the data to be highly correlated

106 (carvones  $r = 0.91$ ,  $p = 1.2E-2$ ; menthols  $r = 0.84$ ,  $p = 4.3E-3$ ; 2-octanols  $r = 0.65$ ,  $p = 8.3E-3$ ). Altogether,

107 these results substantiated the pS6-IP-Seq dataset and yielded confidence that the pS6-IP-Seq strategy  
108 would provide an index of receptor selectivity even amongst structurally similar odorants (Figure 1C,  
109 Supplementary figure 2A-F).

110

### 111 **Describing receptor tuning**

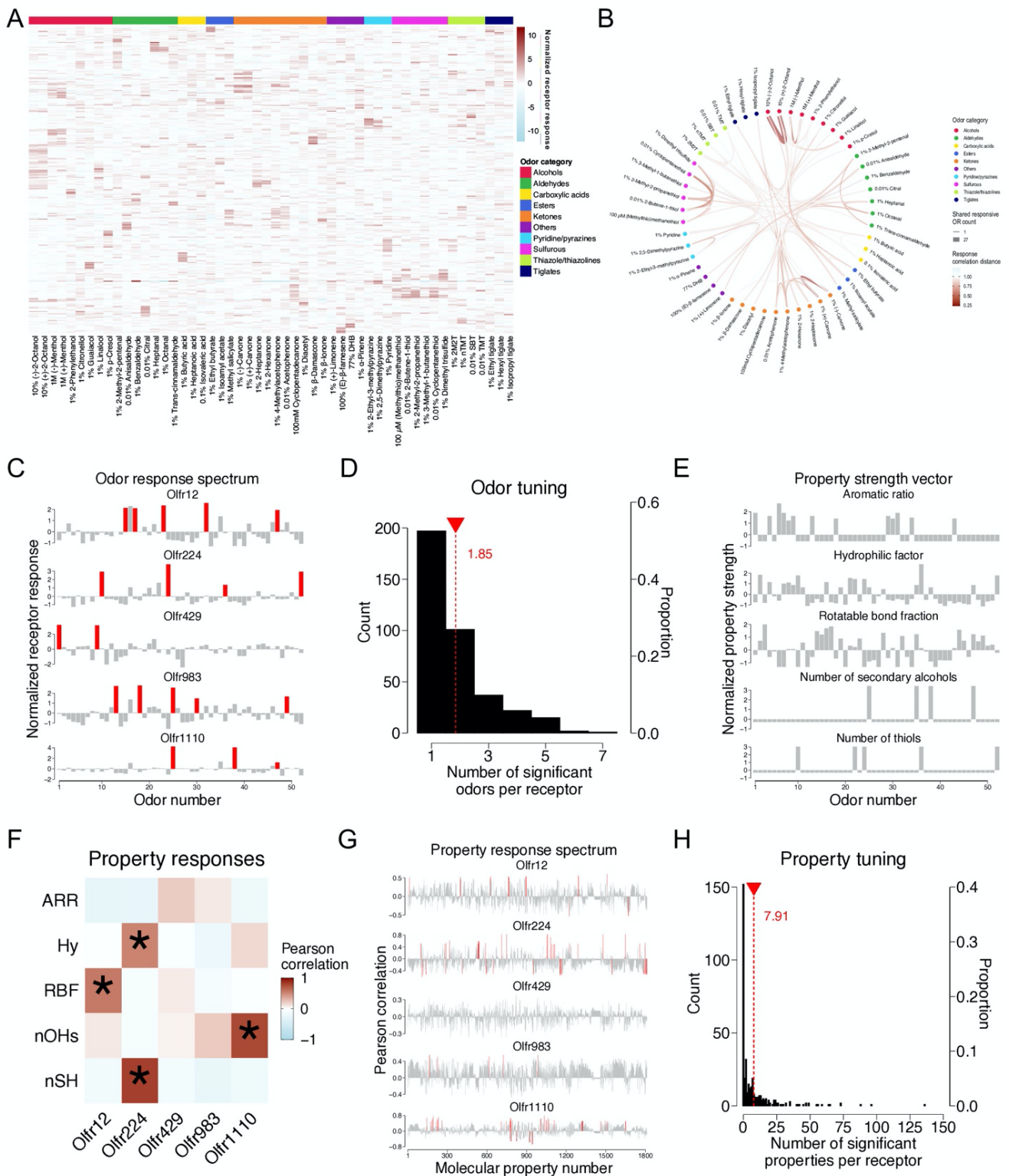
112 Having broadly sampled chemical and receptor spaces, we next sought to quantify the relative  
113 responses of individual receptors to the test odor panel. Individual receptors displayed unique  
114 response profiles across the odorants. Examining receptor tuning did not reveal a bimodal distribution  
115 of narrowly and broadly tuned receptors, but rather a continuum of tuning breadths with an average of  
116 1.85 cognate odorants per significantly responding receptor (Figure 2A-D, Supplementary figure 3A-B).

117

118 To describe the tuning of ORs towards specific molecular properties, we next generated property  
119 strength vectors (PSVs) for each of the molecular descriptors (Figure 2E)<sup>18</sup>. The responsiveness of each  
120 OR to each molecular property was then characterized as a Pearson's correlation between the odor  
121 response spectrum and the values taken by the PSV across the 52 odorants tested (Figure 2F). The  
122 array of such correlations (hereby termed property response spectrum) taken across the molecular  
123 property descriptor set defined the molecular receptive range and property tuning of each OR (Figure  
124 2G). For example, several ORs that displayed robust responses towards thiol odorants yielded tuning  
125 towards the "number of thiol groups" molecular property (Supplementary figure 4A-D). The number of  
126 properties that single receptors responded to significantly (FDR < 0.05) varied from receptor to  
127 receptor with a range of 0 to 136 (mean = 7.91, median = 2). Indeed, the majority of deorphanized ORs  
128 (223/375) displayed significant correlations to at least one of the molecular property descriptors  
129 (Figure 2H). Within the subset of significant OR response-property pairs (2967/ 679125), correlations



130 spanned both negative (-0.77 to -0.49) and positive (0.49 to 0.82) values with an absolute average of  
 131 0.55 (Supplementary figure 3C).



**Figure 2. Combinatorial coding and OR tuning.** **A**, Odor responses across the 52 unique odors, tested at low concentrations, visualized by heatmap. A total of 375 ORs were determined to be responsive to at least one of the tested odors ( $\log_2FC > 0$  and  $FDR < 0.05$ ). Responses are

normalized such that each odor has zero response mean and one standard deviation (z-scored). Responses are color coded from negative, to zero, to positive responses in light blue to white to red. The following odorants are abbreviated: 2-methyl-2-thiazoline (2M2T), 2,4,5-trimethyl-4,5-dihydrothiazole (TMT), 2-*sec*-butyl-4,5-dihydrothiazole (SBT), 2,4,5-trimethylthiazole (nTMT), and 3,4-dehydro-*exo*-brevicommin (DHB). ORs were sorted by correlation distance. Odorants were sorted by odor category. **B**, Odor responses across the same 52 odorants visualized by chord plot. Odorants were sorted by odor category and associations were visualized. Association band thickness corresponds to the number of ORs shared, while color corresponds to overall response similarity across the 375 deorphanized ORs using correlation distance ( $1-r$ ). **C**, Z-scored odor response spectra of five example ORs across the 52 tested odorants. Responses were normalized such that each OR is z-scored. **D**, Histogram of the number of significant odorants per receptor. On average, each responding OR was activated by 1.85 odorants ( $n = 692$ ). **E**, Z-scored PSVs of five example molecular properties across the same set of odorants as **B**. **F**, Property responses given by the Pearson's correlation coefficients between the odor responses (**B**) and PSVs (**D**) calculated over the 52 odorants in the panel. The following molecular properties are abbreviated: aromatic ratio (ARR), hydrophilic factor (Hy), rotatable bond fraction (RBF), number of secondary alcohols (nOHs), and number of thiols (nSH). **G**, Property response spectra characterized by Pearson's correlation coefficients between the five example ORs and 1811 molecular properties. **H**, Histogram of the number of significant molecular properties per receptor ( $n = 2967$ ). On average, each odor-responsive OR was significantly correlated with 7.91 molecular properties with a median of 2 (FDR < 0.05). Out of the 375 deorphanized ORs, 223 ORs displayed significant correlations to at least one molecular property descriptor.

132

### 133 **Odor molecular properties are informative of odor response patterns**

134 Having identified receptor responses to a large and diverse set of odorants, we next sought to  
135 determine the effectiveness of using odor molecular properties to predict receptor responses via  
136 similarities<sup>5,18,19</sup>. To describe the similarity between odorants in molecular property space, we  
137 calculated distances of normalized property strength values between odorant pairs. To represent odor  
138 similarity in OR response space, we similarly calculated pairwise distances between normalized  
139 receptor responses. Linear regression between odorant similarity distances and response similarity  
140 distances revealed a significant relationship ( $r = 0.29$ ,  $p = 3.4E-27$ ; **Figure 3A, Supplementary figure 5A-**  
141 **B**), implying odor molecular properties were informative of receptor response patterns.

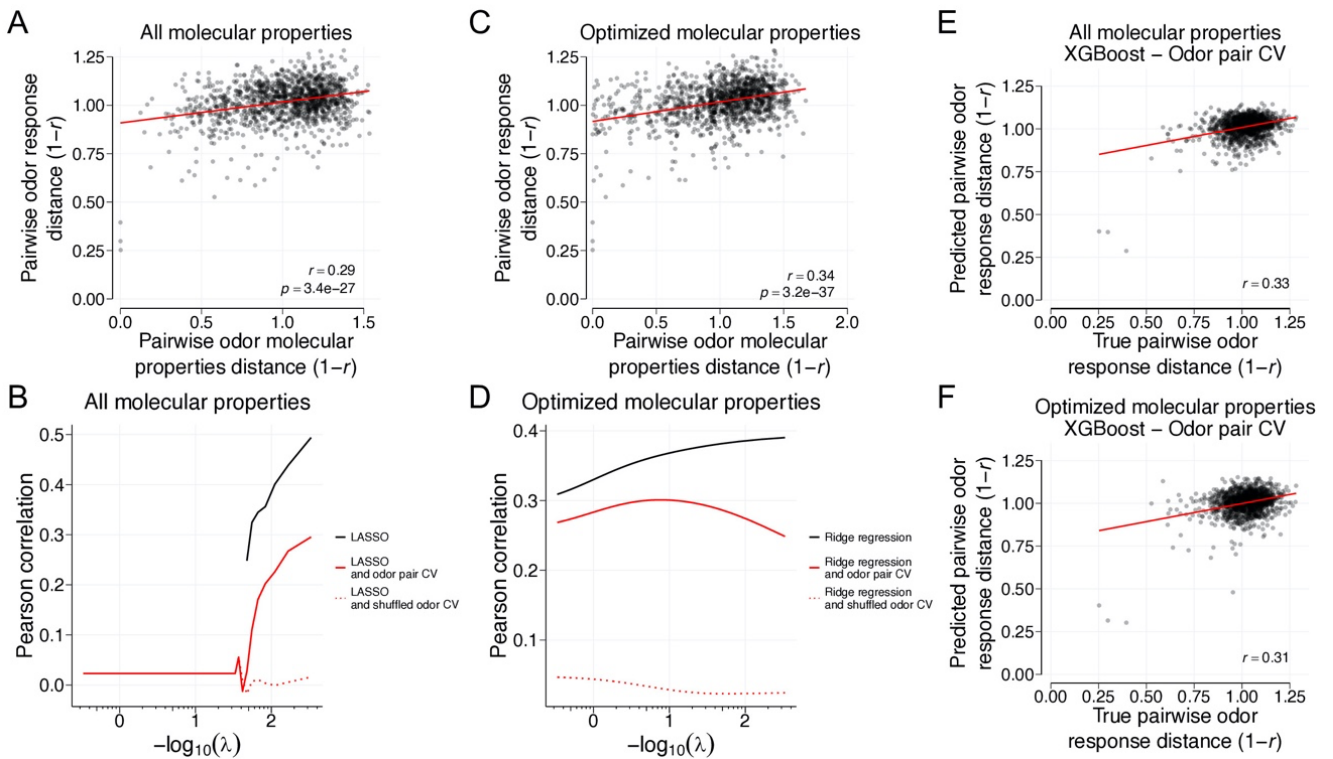
142

143 We next considered the possibility that a subset of the molecular property descriptors may be better  
144 able to relate odor molecular property similarities to receptor response similarities. To test this

145 possibility, we built a sparse regression and performed feature selection using the Least Absolute  
146 Shrinkage and Selection Operator (LASSO). By varying the LASSO loss function ( $\lambda$ ) to influence the  
147 number and relative contribution of the selected molecular properties, we observed improved  
148 correlations with increasing numbers of weighted molecular properties. Importantly, by performing  
149 odor pair cross-validation, we observed that parsimonious combinations of odor molecular properties  
150 selected by LASSO yielded positive predictive abilities (optimal correlation distance odor pair cross-  
151 validation  $r = 0.30$ , [Figure 3B](#), [Supplementary figure 6A-B](#)). To complement these findings, we also  
152 selected an “optimized” set of 65 molecular properties; which included descriptions of aromaticity,  
153 functional group, and molecular geometry; that could be individually linearly decoded and regressed  
154 from OR response patterns alone ([Supplementary figure 7A](#), [Supplementary table 2](#)). Using ridge  
155 regression with the “optimized” set of 65 molecular properties, we could again predict response  
156 similarities from molecular property similarities (optimal correlation distance odor pair cross-validation  
157  $r = 0.30$ , [Figure 3C-D](#), [Supplementary figure 7B-C](#)). Altogether, we interpreted these “optimized” set of  
158 65 molecular property descriptors as both, being capable of explaining OR response variance, and  
159 contributing to the natural statistics of odorants.

160

161 To further validate the predictive abilities of molecular properties and our molecular property  
162 optimization, we also trained and cross-validated a feed-forward non-linear model (XGBoost). In the  
163 first cross-validation scheme, we performed odor-pair cross-validation using all calculated molecular  
164 properties as predictors. In the second, we limited molecular properties to the “optimized” set. In both  
165 cross-validation schemes, predicting response similarities from molecular properties outperformed  
166 shuffled controls ([Figure 3E-F](#), [Supplementary figure 8A-B](#)). Altogether, these results show that odor  
167 responses can be explained in part by combinations of molecular property descriptors.



**Figure 3. Pairwise odor similarity comparisons in response and property spaces.** **A**, Pairwise correlation distance measurements between odorants in response space regressed against molecular property space ( $r = 0.29$ ,  $p = 3.4E-27$ ,  $n = 1326$ ). **B**, LASSO regression as a function of varying the loss function ( $\lambda$ ). Large  $\lambda$  values (leftmost edge) generally correspond to the selection of a few weighted molecular properties. Small  $\lambda$  values (rightmost edge) generally correspond to the selection of many weighted molecular properties. By increasing the number of weighted molecular properties selected by the LASSO algorithm, odor responses become increasingly well fit by molecular properties (black line). To evaluate the generalizability of the sparse property selection, pairs of odorants were iteratively held-out during training and added back intact (solid red line) or shuffled (dashed red line) for cross-validation. Euclidean distances were used to quantify differences between pairs of odorants in molecular property space. Correlation distances were used to quantify differences between pairs of odorants in response space. **C**, Pairwise correlation distance measurements between odorants in response space against optimized molecular property space ( $r = 0.34$ ,  $p = 3.2E-37$ ,  $n = 1326$ ). **D**, Ridge regression results as a function of varying the  $\lambda$  loss function using the optimized set of molecular properties. Large  $\lambda$  values (leftmost edge) generally correspond to dependence on a few weighted molecular properties. Small  $\lambda$  values (rightmost edge) generally correspond to dependence on many weighted molecular properties. **E**, Results of odor pair cross-validation using the XGBoost model framework with default hyperparameters with all molecular properties ( $r = 0.33$ ,  $n = 1326$ ). **F**, Results of odor pair cross-validation using the XGBoost model framework with default hyperparameters with the optimized set of molecular properties ( $r = 0.31$ ,  $n = 1326$ ).

168

### 169 Specific receptor residues predict ligand selectivity of ORs

170 Comprehensive identification of ORs responsive to many odorants prompted us to next search for

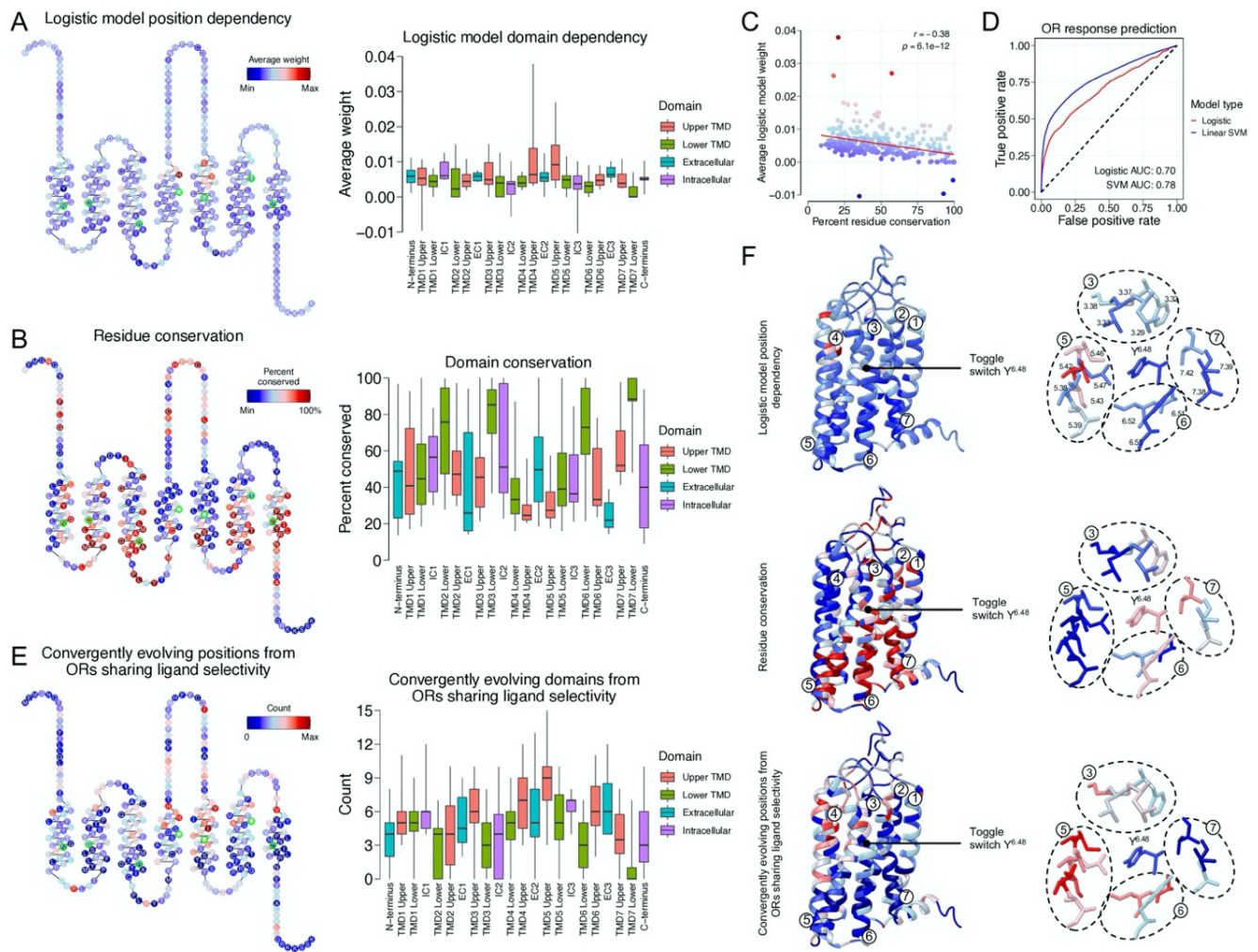
171 generalizable relationships between odorants, ORs, and receptor residues<sup>11,14</sup>. To do so, we built

172 logistic models using aligned ORs. For each odor fit with a regularized logistic model, receptors were  
173 randomly split into 90% training and 10% testing sets for 100 repetitions. Iterating this process over the  
174 set of tested odorants identified a series of weighted positions, harboring amino acids with predictive  
175 power, occurring primarily at the upper halves of the fourth and fifth transmembrane domains (TMDs)  
176 (Figure 4A). For example, visualizing amino acids occurring at these positions identified an enrichment  
177 of cysteine and methionine residues in TMD5 amongst ORs responding to sulfurous odorants  
178 (Supplementary figure 9A). Regressing the average weight assigned to each position, from odorants  
179 solvable by logistic regression (area under receiver operating characteristic curve, AUROC > 0.5), by  
180 percent conservation revealed an anti-correlation ( $r = -0.38$ ,  $p = 6.1E-12$ , Figure 4B-C)<sup>1</sup>. Using a Support  
181 Vector Machine (SVM) classifier, with a linear kernel, led to similar predictions regarding the response  
182 likelihoods of held-out ORs (Logistic regression AUROC = 0.70, Linear SVM AUROC = 0.78, Figure 4D,  
183 Supplementary figure 10A).

184  
185 The massive expansion and rapid evolution of the OR gene family posits opportunities for the  
186 convergent evolution of distantly related ORs to evolve odorant selectivity independently. To search  
187 for receptor sequence positions exhibiting convergent evolution, we asked if ORs sharing response to  
188 an odorant possessed positions harboring amino acids with physical-chemical properties, measured by  
189 Grantham's distance<sup>17</sup>, which deviated from comparable but odor-unresponsive ORs. Iterating over the  
190 set of tested odorants, this analysis identified a series of poorly conserved positions ( $r = -0.54$ ,  $p = 6.5E-$   
191  $25$ , Figure 4E, Supplementary figure 10B), especially localized to the upper half of TMD5. Regressing  
192 the average weight assigned to each position, via regularized logistic regression, by the number of  
193 times each position displayed convergent evolution, revealed similar findings between the two  
194 approaches ( $r = 0.42$ ,  $p = 1.5E-14$ , Supplementary figure 10C). Importantly, the localization of these  
195 positions, and those identified by logistic models, was consistent with a region implicated in ligand

196 binding in other class A GPCRs<sup>20-25</sup>. Altogether these results show the odorant selectivity of ORs are in  
197 part explained by convergently evolving residues occurring at a common site of poorly conserved  
198 residues within the TMDs.

199  
200 To visualize the results of our analyses in 3D, we next built an OR homology model. Focusing on the  
201 conserved “toggle switch” Y<sup>6.48</sup> residue previously reported to reside at the bottom of the ligand-  
202 binding cavity of other class A GPCRs<sup>20,26-28</sup>, we consistently observed nearby residues in the upper  
203 halves of TMD3, TMD5, and TMD6 as exhibiting heavy weights in our logistic models, poor  
204 conservation, and convergent evolution, implying a canonical cavity for odorant binding across our  
205 tested odorants. Altogether, these results are consistent with the idea that few mutations within the  
206 ligand binding site of ORs can broadly reconfigure chemical tuning, a feature that is likely to have  
207 facilitated the rapid evolution of receptors with distinct ligand specificities.



**Figure 4. Sequence-function relationships of ORs.** **A**, ORs were split into 90% training and 10% testing sets for 100 repetitions of regularized logistic regression using aligned sequences as predictors. Odorants which yielded AUROC values  $> 0.5$  (Supplementary table 3) were subset and non-zero weights were averaged across positions, repetitions, and odorants. Individual weights were visualized by snakeplot (left) with the most commonly occurring amino acid described at each position. Ballesteros-Weinstein  $\times 0.50$  numbers for each TMD are highlighted in green. Weight distributions were visualized across domains by box-and-whisker plot (right). **B**, Left, the most commonly occurring amino acid was quantified by its percent conservation across ORs and visualized by snake plot. Right, percent conservation distributions were visualized across receptor domains by box-and-whisker plot. **C**, Regressing the average weight assigned to a position, by regularized logistic regression, by its percent conservation reveals an anti-correlation ( $r = -0.38$ ,  $p = 6.1E-12$ ). **D**, Using a SVM classifier with a linear kernel to predict response likelihoods of held-out ORs leads to similar results as logistic regression (logistic regression AUROC = 0.70, linear SVM AUROC = 0.78). **E**, ORs sharing response to an odorant were subset and pairwise compared to a null set of ORs with comparable protein sequences. Pairwise Grantham distance distributions between the responsive and null sets were compared by Kolmogorov-Smirnov statistical test to determine if amino acid similarity distributions were different. Count measurements reflect the number of times a position harbored amino acids with differing distributions between the two groups at an FDR  $< 0.05$ . Left, these results were visualized by snakeplot. Right, these results were visualized by box-and-whisker plot describing domain distributions. **F**, Left, visualization of the data in 3D using homology models. Color schemes are identical those in panels A, B, and E. TMDs are indicated. Right, zoomed in view

focusing on residues found directly above the conserved Y<sup>6.48</sup> toggle switch residue. A consensus of poorly conserved residues found in TMD3, 5, and 6 can be seen to exhibit both higher weights by regularized logistic regression models, and convergent evolution. Ballesteros-Weinstein numbers associated with displayed residues are also described.



209 **Discussion:**

210 Given the inordinate complexity of the chemical world, large repertoires of ORs appear to be necessary  
211 for the detection and discrimination of diverse chemicals in the environment, as exemplified by the  
212 significant genomic space that is systematically subjugated to ORs across numerous species. These  
213 findings are compounded by the identification of OR-specific chaperone proteins which may allow  
214 functional expression of ORs with cryptic mutations, further underscoring the high degree of sequence  
215 diversification ORs are enabled to possess<sup>29,30</sup>. Using a diverse set of odorants, here we have  
216 performed a functional *in vivo* characterization of the OR repertoire of *Mus musculus*. Linking the  
217 activity of the receptor repertoire to an extensive set of molecular property descriptors parameterizing  
218 the physical-chemical properties of the odorants, we learned that ORs displayed a continuum of tuning  
219 breadths. Similarities between sparse sets of molecular properties could be used to predict receptor  
220 response patterns. Finally, analyses linking odorant selectivity and amino acid residues most  
221 consistently identified a series of poorly conserved residues located primarily in the upper half of the  
222 transmembrane domains.

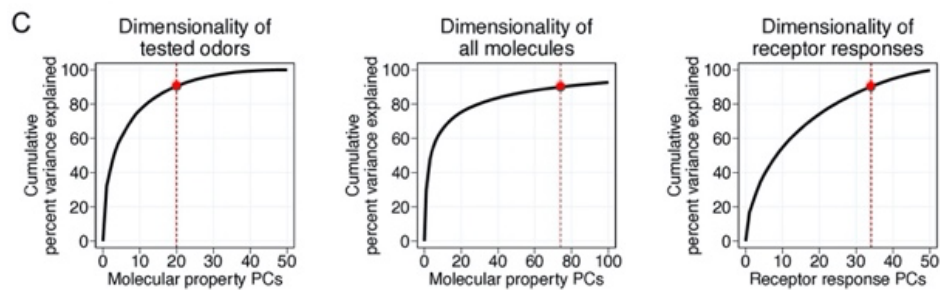
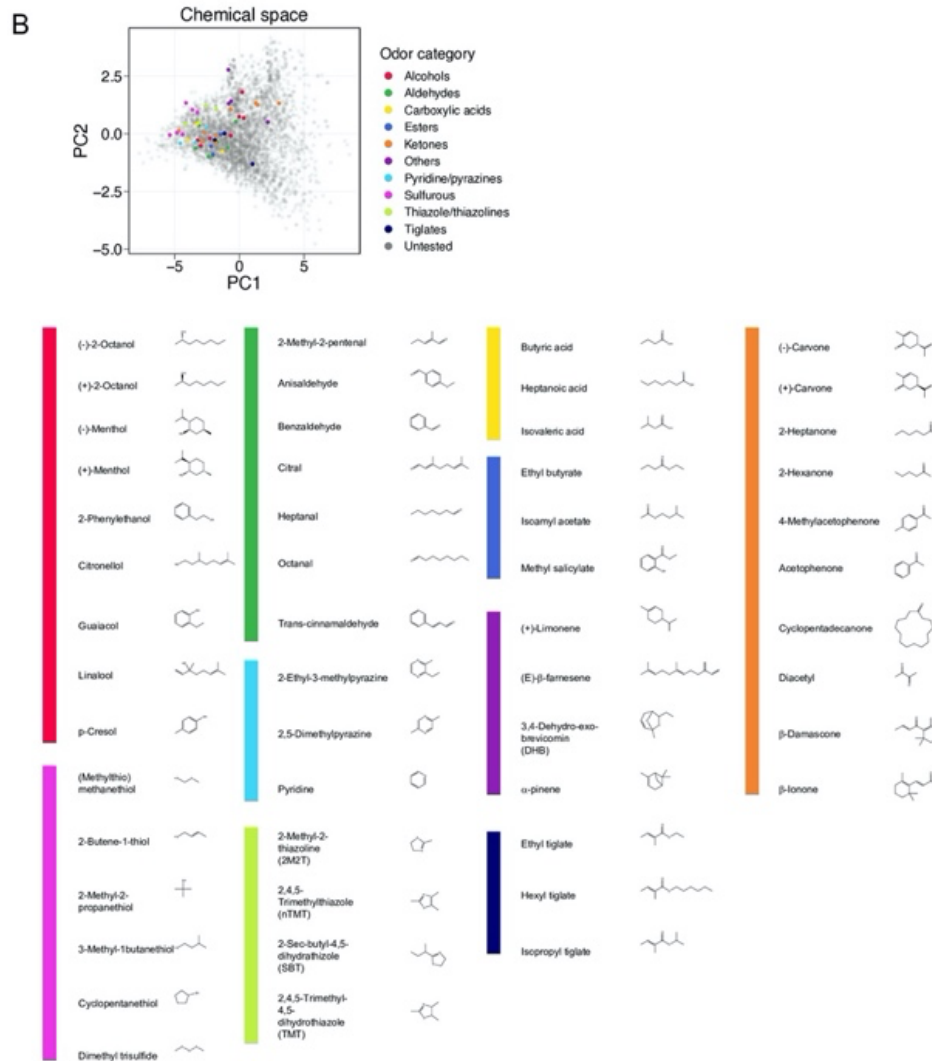
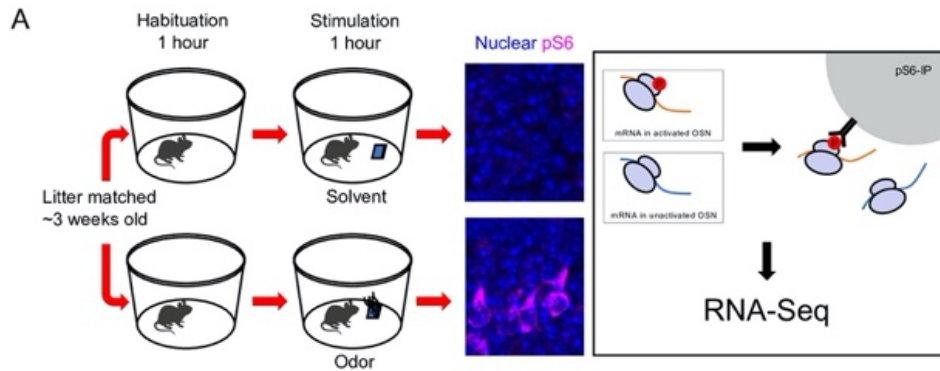
223  
224 While our test odor panel, and resulting deorphanized receptor set, was broad in the coverage of odor  
225 and receptor spaces, the data was by no means all-encompassing. 20 PCs were required to cover 90%  
226 of variance in the 52 set of tested odorants, whereas the full set of 4732 small molecules required 74  
227 PCs to achieve comparable coverage (Figure 1 – Supplementary figure 1B). Our analyses therefore  
228 likely reflect lower bound estimates of chemical and receptor response spaces. Nevertheless, these  
229 results also imply that a substantial amount of the information in the molecular property descriptors is  
230 highly redundant. Furthermore, we note that the dimensionality of the tested odorants in receptor  
231 response space is higher than the dimensionality of the odorants in chemical space, with 34 PCs  
232 needed to explain more than 90% of the response variance (Figure 1 – Supplementary figure 1B). This

233 increased dimensionality indicates there are facets of odor response by ORs that are poorly explained  
234 by a similar number of flat surfaces in chemical space described by molecular property descriptors.  
235 Although the molecular property descriptors used in this study can explain and predict response  
236 similarities, further searches for latent descriptors capable of better associating odor properties to  
237 their receptor responses may improve these associations<sup>31,32</sup>.

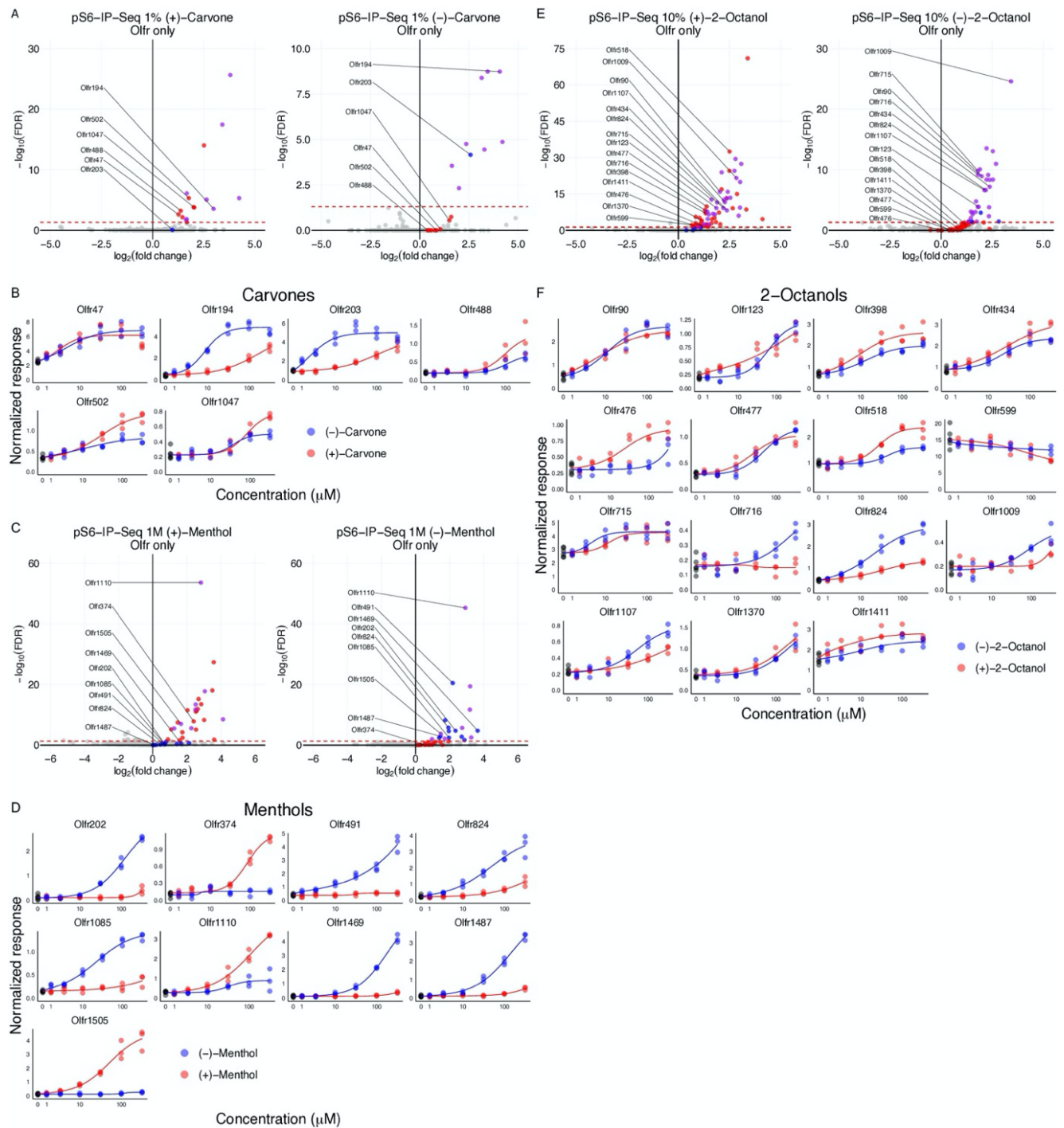
238  
239 Several implications arise from our observation that diverse odorants share commonalities that relate  
240 their receptor responses to amino acid residues. First, the poor conservation of positions harboring  
241 residues exhibiting predictive power and convergent evolution suggest a mechanism by which flexible  
242 chemical recognition can be achieved by a family of proteins while maintaining a degree of  
243 conservation necessary for functional protein integrity and activation of conserved downstream  
244 signaling cascades. Second, the association of the third, fifth, and sixth transmembrane domains with  
245 odor selectivity are also consistent with site-directed mutagenesis efforts on single ORs that have been  
246 shown to influence OR responses<sup>33-37</sup>. Finally, these results are consistent with recent evidence from  
247 structural elucidation of an ionotropic insect OR, which revealed a single binding pocket for a  
248 structurally diverse odorants<sup>38</sup>. Future structural elucidation of mammalian ORs will enable direct  
249 addressing of the modes of odorant-OR interactions.

250  
251 In summary, we have provided a systematic, quantitative analysis of the primary representation of an  
252 odor, as registered by the differential responses of individual ORs. Our results and analyses provide a  
253 foundational framework for investigating how these primary odorant representations are transformed  
254 into subsequent representations to ultimately guide behavioral outputs.

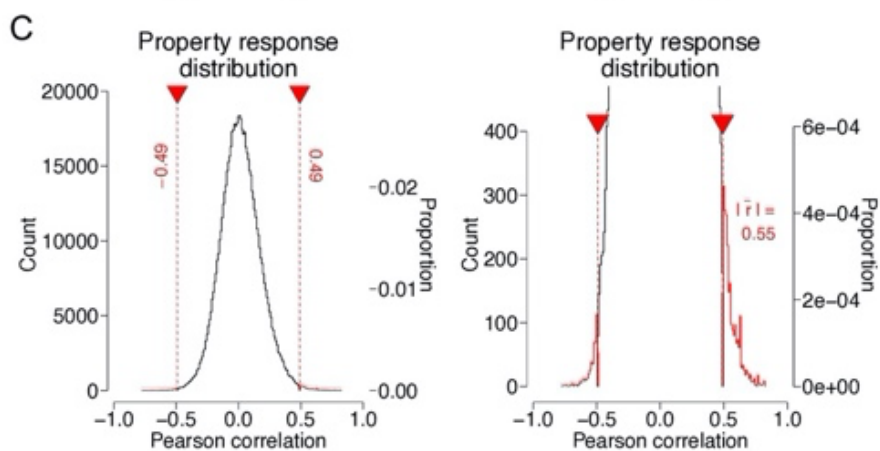
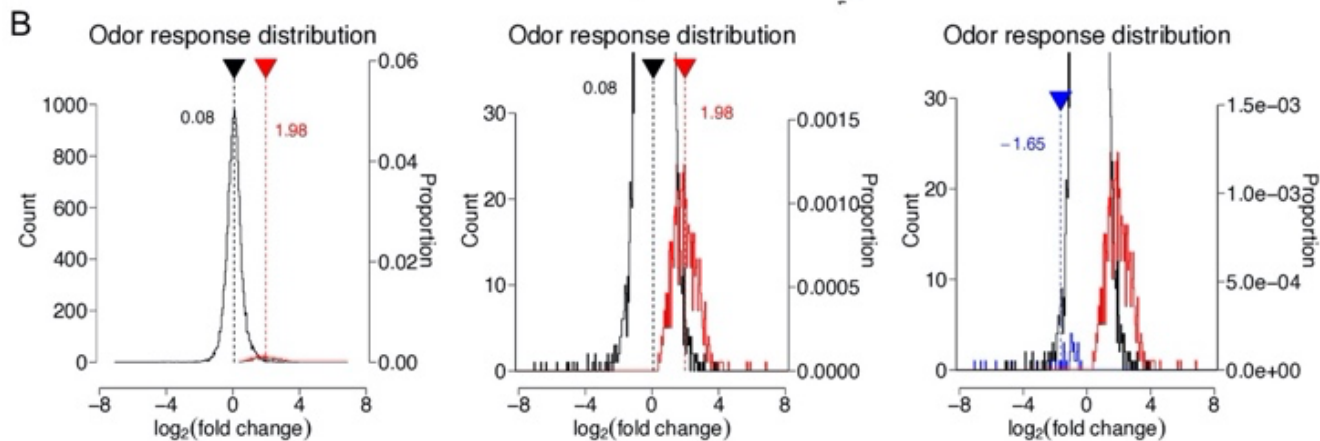
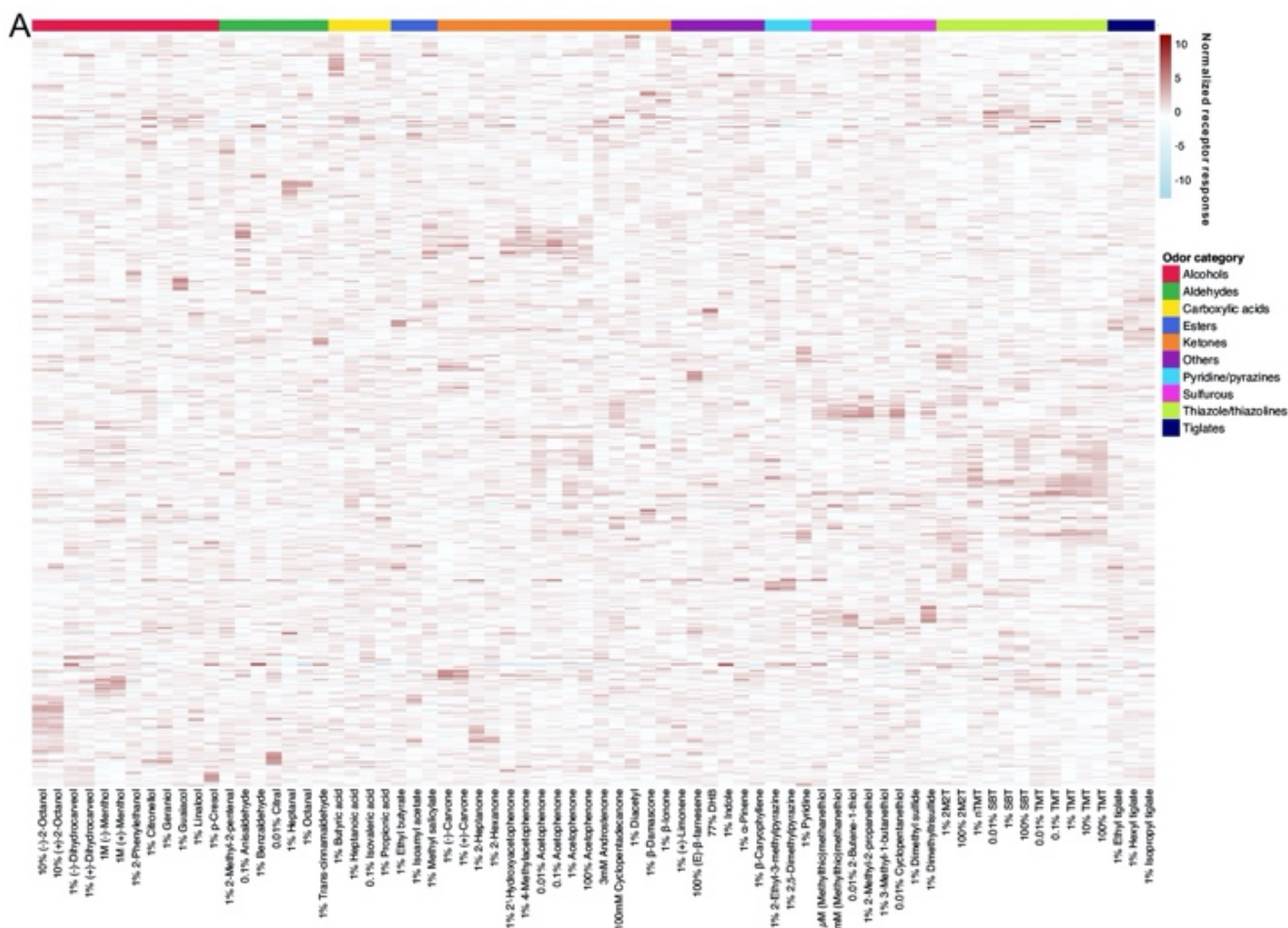
255 **Supplementary figures:**



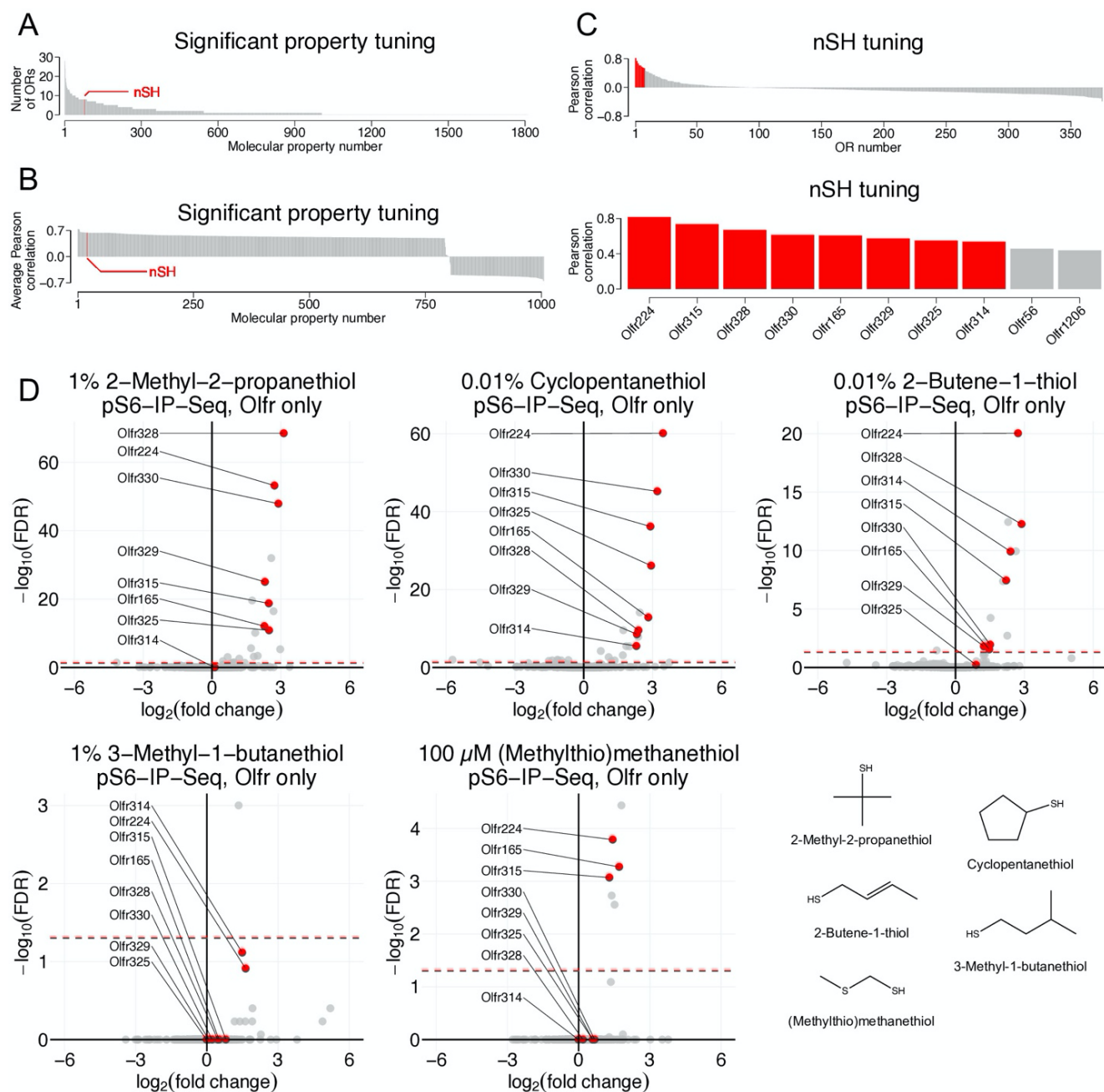
**Supplementary figure 1. Details of chemical and receptor space sampling.** **A**, Schematic of the pS6-IP-Seq experiment. Litter matched, ~3 weeks old mice were used. Mice were first habituated to an odor-free environment for 1 hour. One mouse then received exposure to an odor stimulus, while another received exposure to the solvent, again for 1 hour. Whole olfactory mucosa was then harvested and immunoprecipitated using an antibody against pS6 and subjected to RNA-Seq. **B**, Top, tested odorants in chemical space colored by odor category. Bottom, molecular structures of tested odorants sorted by odor category. **C**, Left, the cumulative percent variance of the tested odorants in chemical space explained as a function of included PCs of molecular properties. A minimum of twenty PCs was required to capture at least 90% of the variance for the test odorants in chemical space. Middle, the cumulative percent variance of all molecules in chemical space explained as a function of included PCs of molecular properties. A minimum of seventy-four PCs was required to capture at least 90% of the variance for all molecules in chemical space. Right, the cumulative percent variance of the tested odorants in response space explained as a function of included PCs of receptor responses. A minimum of thirty-four PCs was required to capture at least 90% of the variance for the test odorants in response space.



**Supplementary figure 2. Details of *in vitro* validation.** **A**, Volcano plots for pS6-IP-Seq results by exposure of mice to 1% (+)-carvone and 1% (-)-carvone. A red line at FDR = 0.05 is drawn. ORs enriched by just (+)-odorant are colored red while ORs enriched by just (-)-odorant are colored blue. ORs enriched by both enantiomers are colored purple. Labeled ORs were validated *in vitro*. **B**, Dose response curves of ORs displaying *in vitro* response to at least one of the tested carvone enantiomers. **C**, Volcano plots for pS6-IP-Seq results using 1M (+)-menthol and 1M (-)-menthol. **D**, Dose response curves of ORs displaying *in vitro* responses towards menthol enantiomers. **E**, Volcano plots for pS6-IP-Seq results using 10% (+)-2-octanol and 10% (-)-2-octanol. **F**, Dose response curves of ORs displaying *in vitro* responses towards 2-octanol enantiomers.

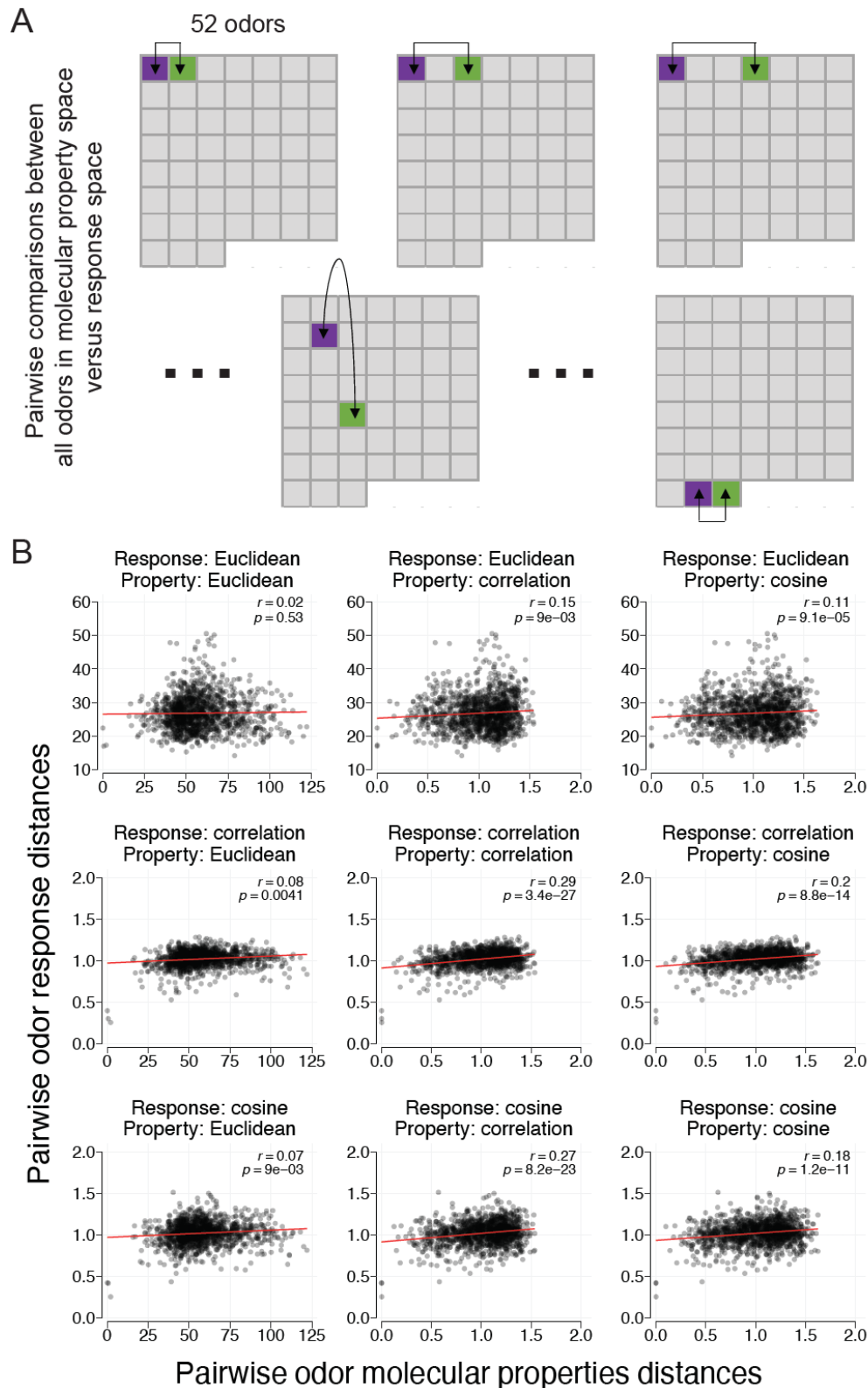


**Supplementary figure 3. Details of OR response properties. A**, Odorant responses determined by pS6-IP-Seq across all tested 72 odorants/concentrations visualized by heatmap. A total of 555 ORs are determined to respond to at least one of these odorants/concentrations. Responses are z-scored by odor. The following odorants are abbreviated: 2-methyl-2-thiazoline (2M2T), 2,4,5-trimethyl-4,5-dihydrothiazole (TMT), 2-sec-butyl-4,5-dihydrothiazole (SBT), 2,4,5-trimethylthiazole (nTMT), and 3,4-dehydro-exo-brevicommin (DHB). Odorants are sorted by functional group while ORs are sorted using correlation distance. **B**, Left, histogram (bin size = 0.05) of the distributions of OR responses across the panel of 52 odorants tested at low concentrations. Significant OR-odor pairs ( $\log_2FC > 0$  and  $FDR < 0.05$ ) are colored red, while non-significant pairs are colored black. OR-odor pairs that were classified as nonsignificant had an average  $\log_2FC$  enrichment of 0.08 by pS6-IP-Seq. OR-odor pairs classified as significant had an average  $\log_2FC$  enrichment of 1.98 by pS6-IP-Seq (nonsignificant responses  $n = 18808$ , significant responses  $n = 692$ ). Middle, zoomed in. Right, A small number of inhibitory responses ( $\log_2FC < 0$  and  $FDR < 0.05$ ) were observed in the data. These responses were otherwise classified as nonsignificant ( $n = 44$ ). **C**, Left, distribution of the OR-molecular property pairwise Pearson correlation coefficients using the 52 odorants tested at low concentrations (bin size = 0.01). At an  $FDR < 0.05$ , the Pearson correlation coefficient cutoffs were -0.49 and 0.49 for negative and positive correlations respectively, with an absolute average of 0.55 for significant correlations (nonsignificant correlations  $n = 676158$ , significant correlations  $n = 2967$ ). Right, zoomed in.

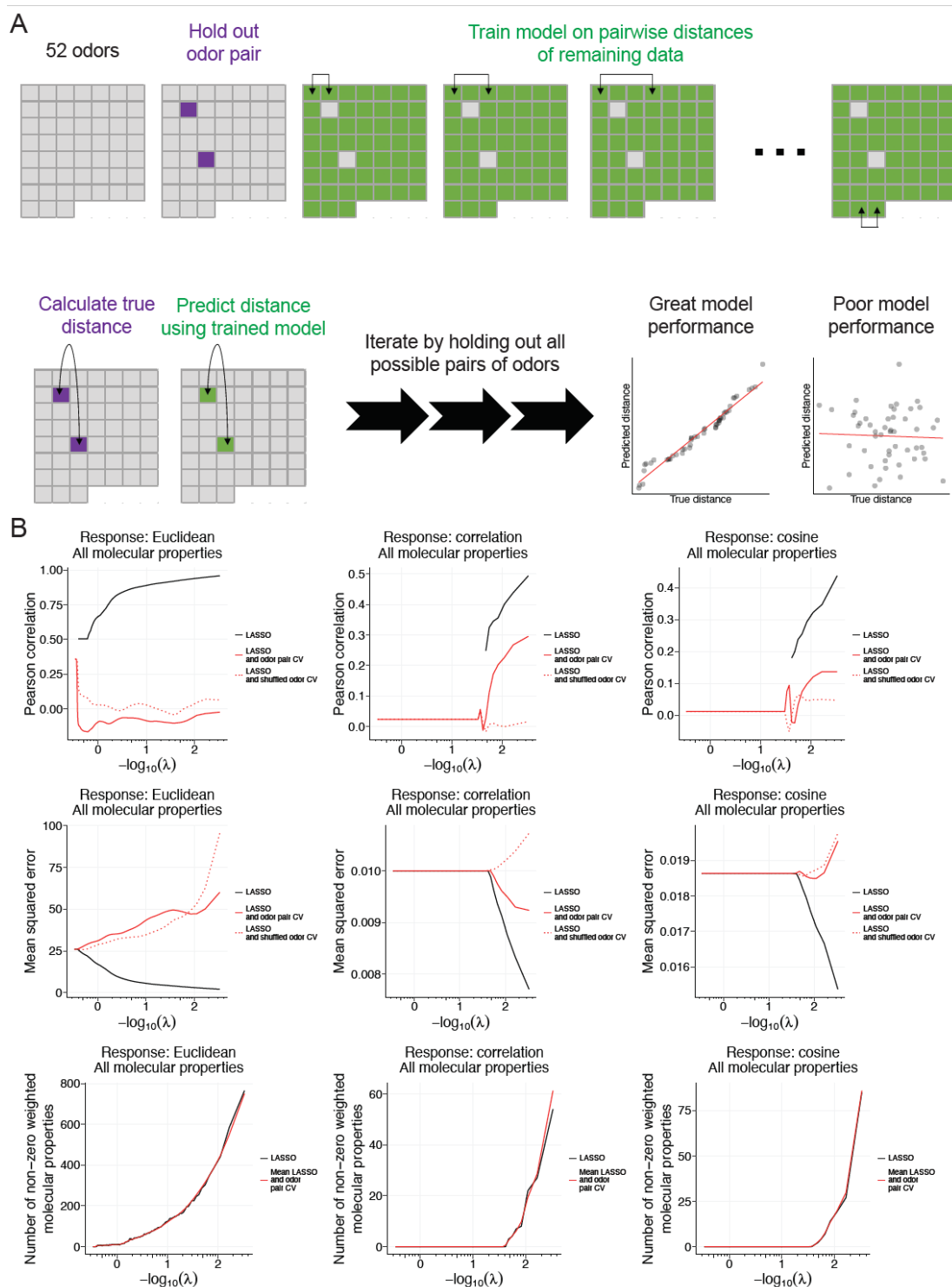


**Supplementary figure 4. Details of thiol property tuning.** **A**, Tuning towards molecular properties described by the number of significant OR-molecular property associations. Out of the 1811 molecular properties, 1005 were tuned towards by at least one OR. Molecular property nSH (number of thiol groups) is highlighted in red. **B**, Visualization of the average Pearson correlation tuning value for ORs that were significantly tuned toward each of the 1005 molecular properties. **C**, Top, tuning of individual ORs, measured by Pearson correlation, towards nSH. Bottom, the top ten ORs displaying nSH tuning. Eight of these ten ORs were significant (FDR < 0.05). **D**, Raw pS6-IP-Seq differential expression data, visualized by volcano plot, with chemical structures of the tested thiol odorants. ORs that were tuned towards thiol are highlighted in red. A consensus of thiol odorant response can be observed for the ORs tuned towards nSH.

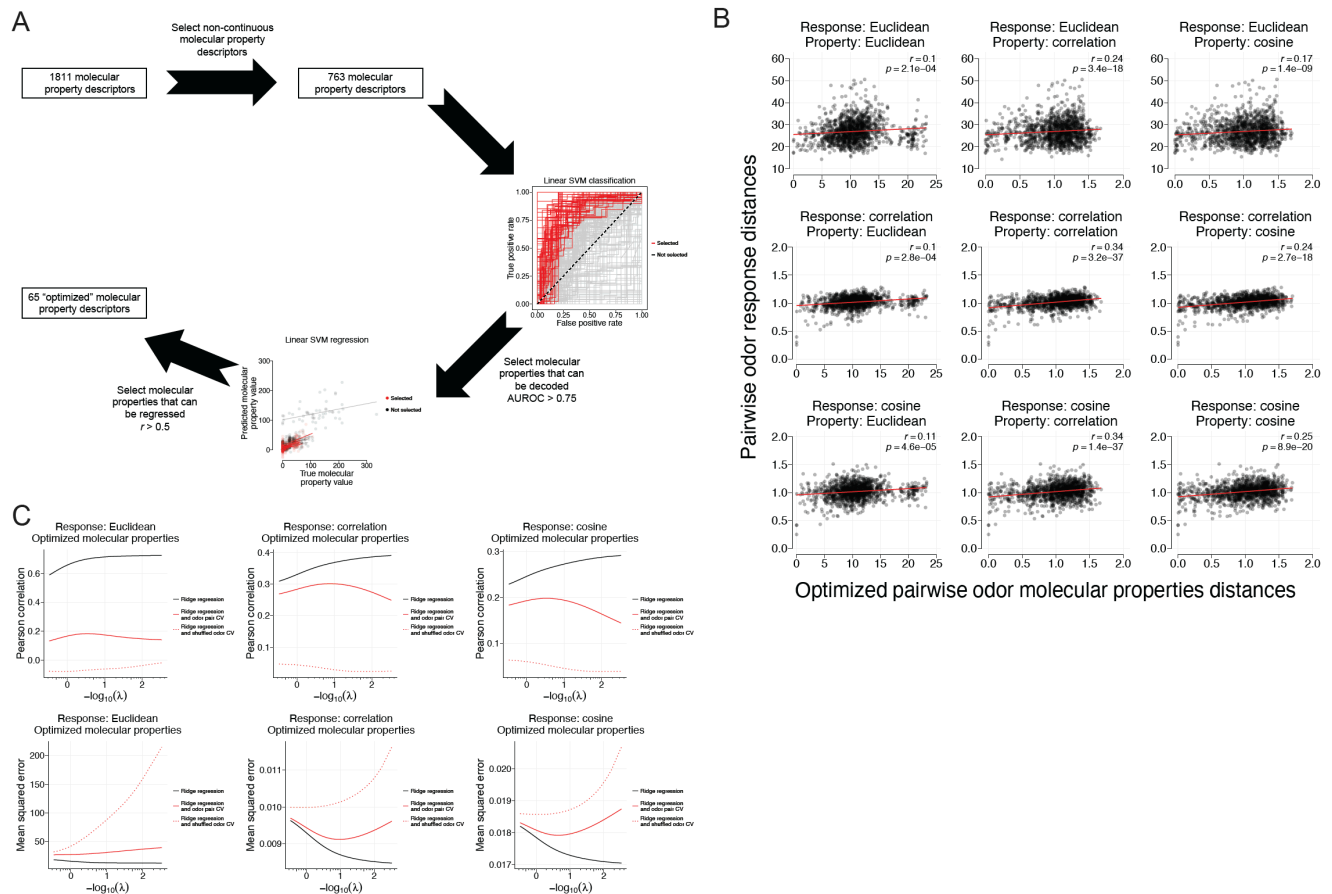


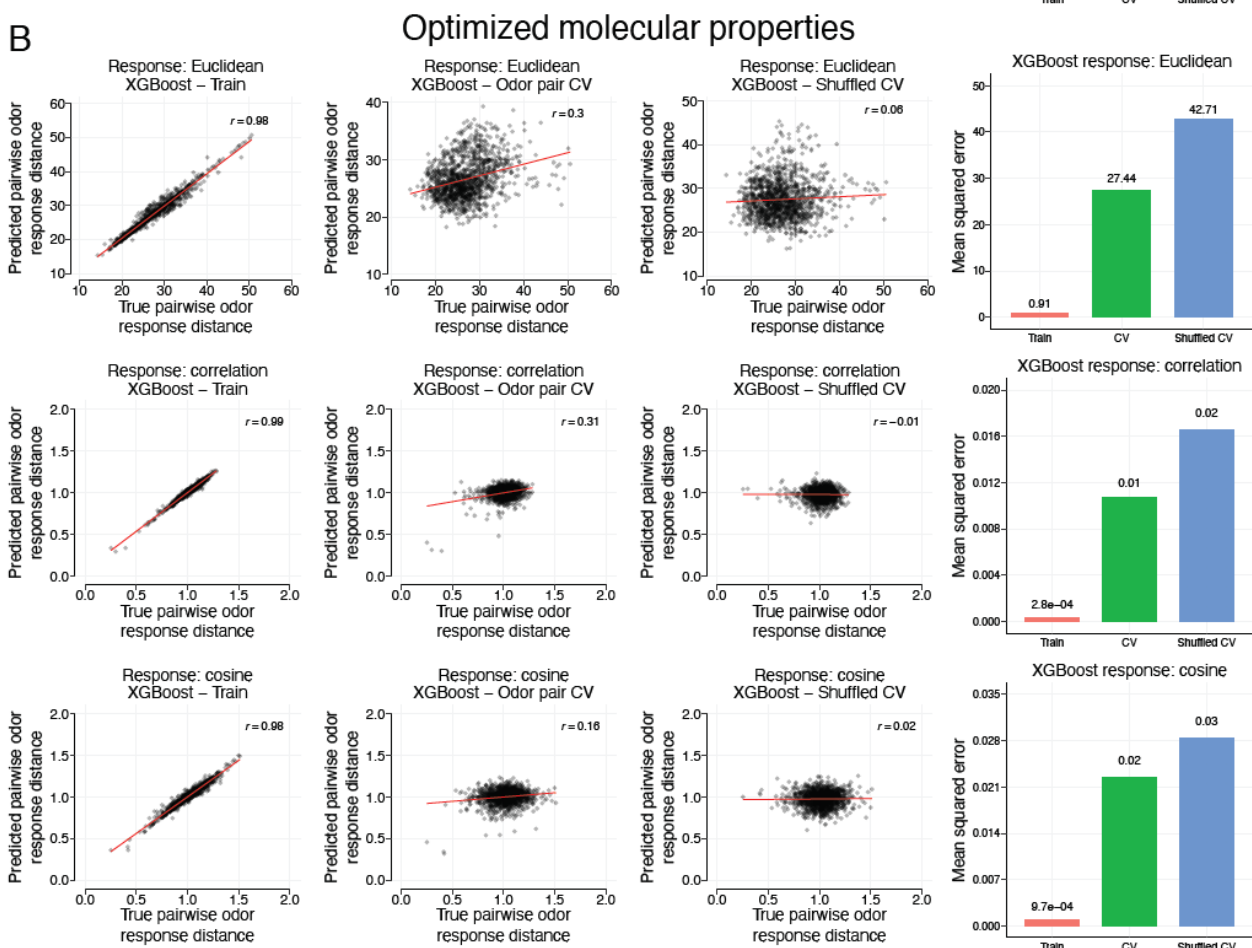
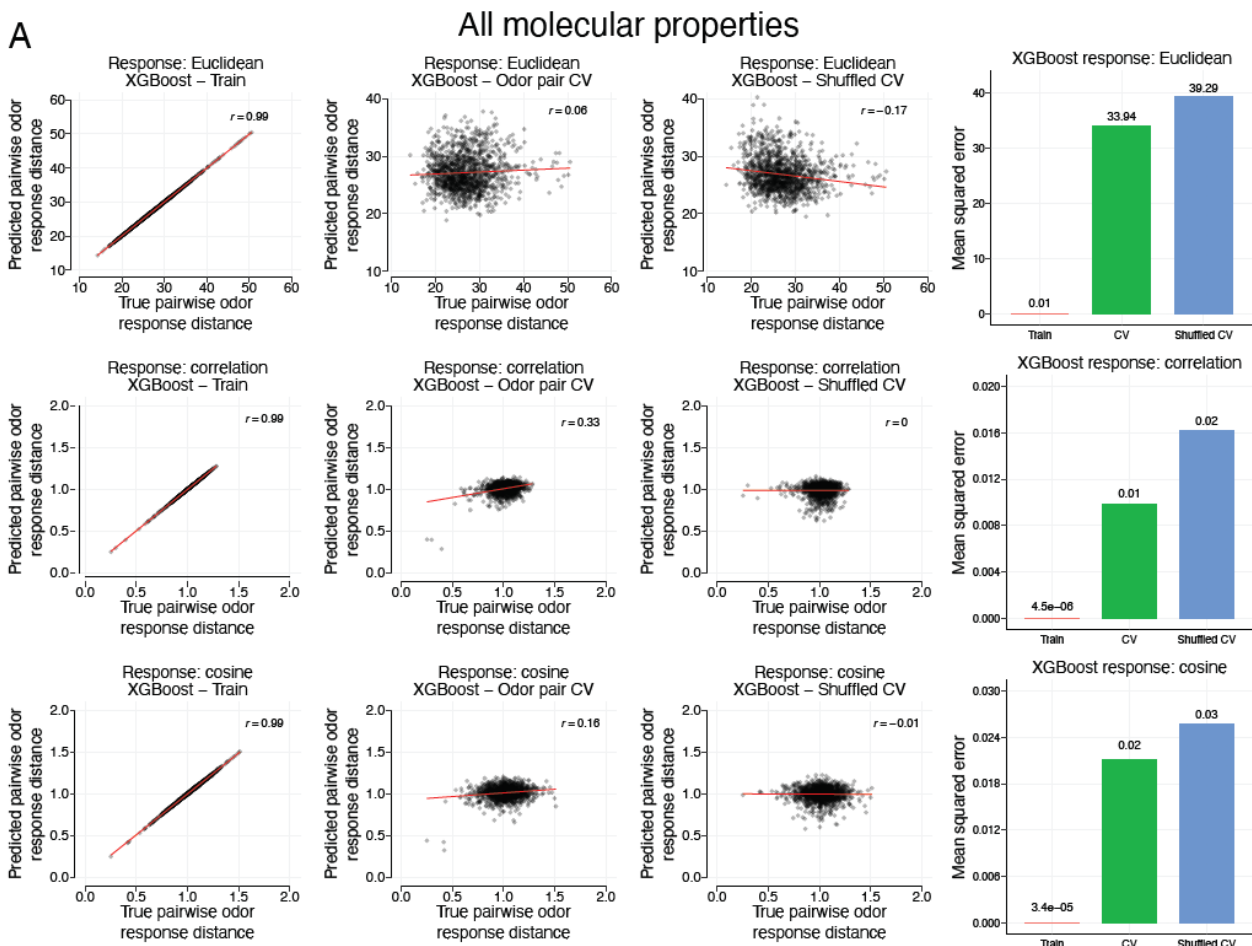


**Supplementary figure 5. Details of the pairwise comparisons between odorants in response and molecular property space. A**, Schematic of how pairwise distance comparisons between odorants were calculated for figures 3A and 3C. **B**, Correspondence between odorants in response and molecular property spaces. Three (Euclidean, correlation, and cosine) different distance metrics were used. Pearson correlation coefficients and  $p$ -values are reported for each combination of distance metric.

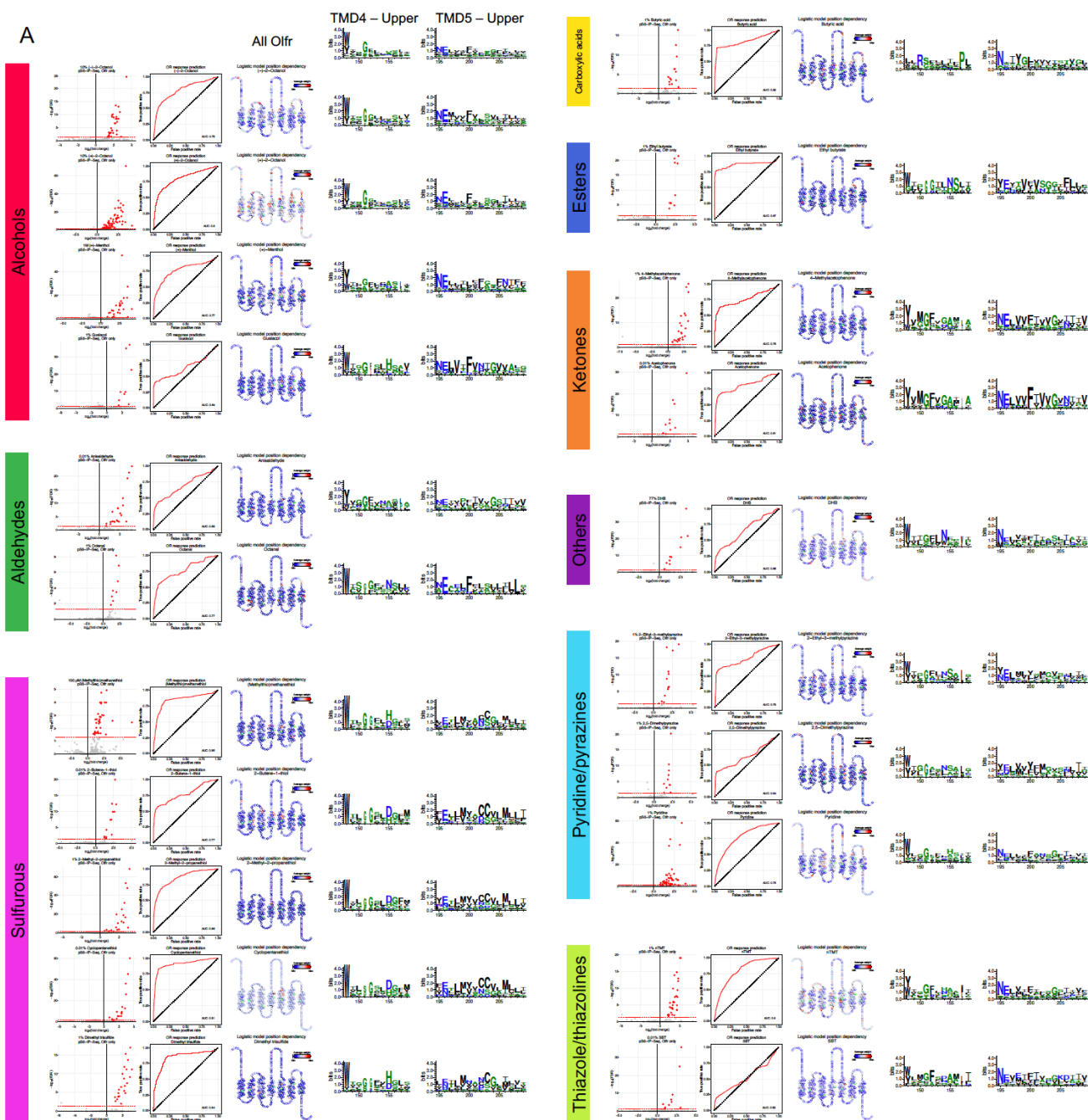


**Supplementary figure 6. Details of the odor pair cross-validation for regression models. A,** Schematic of how odor pair cross-validation was performed. Pairs of odorants were iteratively held-out from training. Distances between held-out odorants were iteratively predicted and regressed against true distances to report a Pearson correlation and mean squared error. **B,** Results of LASSO regression using various metrics to quantify distances between odorants in response space. Molecular property distances were quantified by Euclidean distances. Reported are Pearson correlation coefficients, mean squared error, and the number of non-zero weighted molecular properties as a function of varying the loss function ( $\lambda$ ).

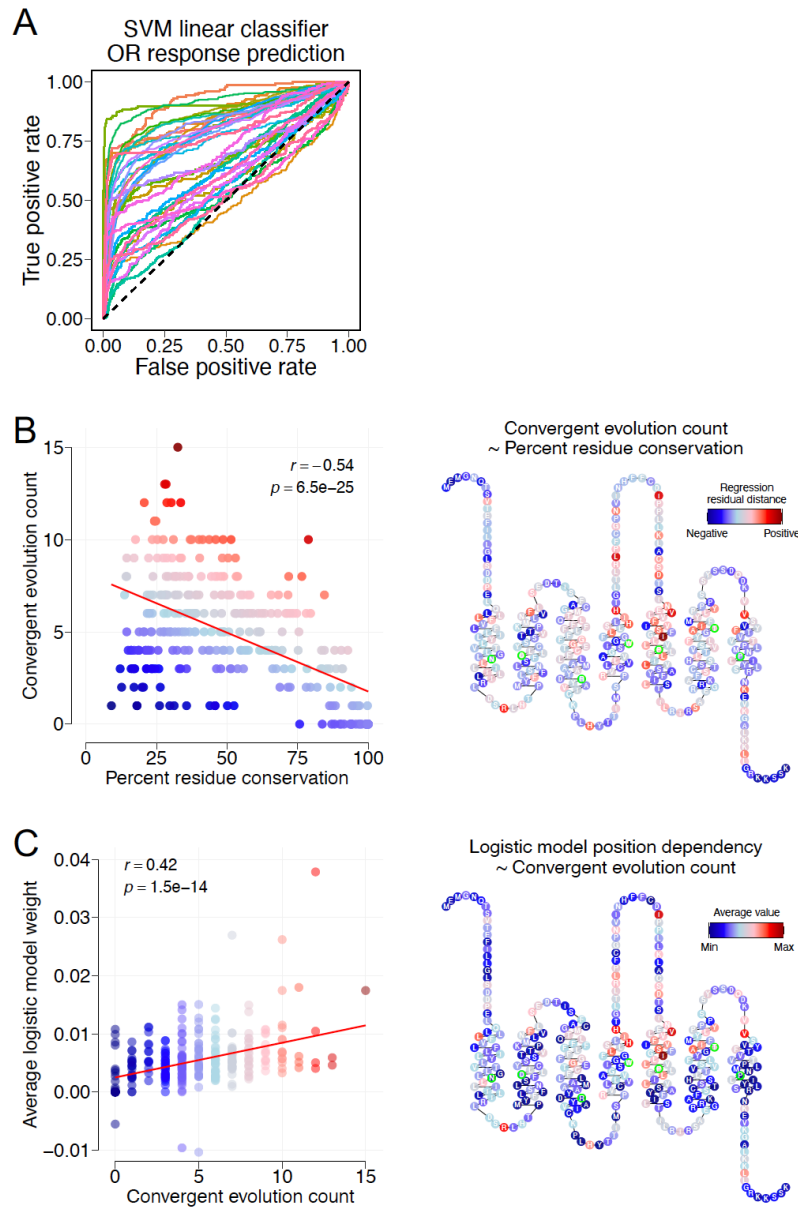




**Supplementary figure 8. Results of using the XGBoost model framework. A,** Using default XGBoost hyperparameters with 1811 molecular property descriptors, we asked how well does odor molecular property similarity predict receptor response similarity. Response similarities were calculated using Euclidean, correlation, and cosine distances. Odor pair cross-validation was performed to evaluate the generalizability of the models. **B,** Results of the XGBoost models when only the 65 set of “optimized” odor molecular properties were used as predictors and odor pair cross-validation is performed. The positive predictive abilities of these sparse 65 odor molecular properties, independent of the distance metric used) demonstrates odor similarity in response space can be approximated with parsimonious combinations of odor molecular properties.



**Supplementary figure 9. Details of logistic regression to uncover sequence-function relationships of ORs. A,** Data for odorants solvable by regularized logistic regression (AUROC > 0.5) is shown. For each odor, shown is a volcano plot highlighting responsive ORs in red ( $\log_2FC > 0$  and  $FDR < 0.05$ ) and non-responsive ORs in gray. Aligned amino acid sequences were used as inputs to predict response likelihoods of a held-out 10% of ORs for 100 repetitions. The receiver operating characteristic (ROC) curve for held-out data is shown. Positions harboring amino acids assigned non-zero weights were averaged and visualized by snakeplot. Ballesteros-Weinstein x.50 numbers for each TMD are highlighted in green. The consistency of high weights assigned to residues localized to the upper halves of the fourth and fifth transmembrane domains motivated visualization of the amino acid distributions of responsive ORs by WebLogo. An enrichment of cysteine residues can be seen amongst ORs responsive to sulfurous odorants at positions 202 and 203 in TMD5. Similarly, an enrichment of methionine residues can be seen amongst ORs responsive to sulfurous odorants at positions 199 and 206.



**Supplementary figure 10. Details of SVM classifiers and comparing between approaches. A,** Results of using a linear SVM classifier to predict OR response likelihoods across 100 repetitions of splitting ORs into 90% training and 10% testing sets. Models were trained and optimized using ten-fold cross-validation. Prediction likelihoods across 100 repetitions were compounded to generate single ROC curves for single odorants. AUROC values for each odorant are reported in supplementary table 4. **B,** Regressing convergent evolution counts by percent residue conservation reveals an anti-correlation ( $r = -0.54$ ,  $p = 6.5E-25$ ). Snakeplot is colored by the distance of each data point from the line of regression. Ballesteros-Weinstein x.50 numbers for each TMD are highlighted in green. **C,** Regressing logistic model dependency, evaluated by average positional weight, by convergent evolution count reveals consistency between approaches ( $r = 0.42$ ,  $p = 1.5E-14$ ).

266 **Methods:**

267 **Phosphorylated S6 ribosomal capture (pS6-IP)**

268 Mice used for pS6-IP were ~3 weeks old, mixed sex, and littermates. Mice were killed by  
269 CO<sub>2</sub> asphyxiation and cervical dislocation. Olfactory tissue was rapidly dissected in Buffer B (2.5 mM  
270 HEPES KOH pH 7.4, 0.63% glucose, 100 µg/mL cycloheximide, 5 mM sodium fluoride, 1 mM sodium  
271 orthovanadate, 1 mM sodium pyrophosphate, 1 mM β-glycerophosphate, in Hank's balanced salt  
272 solution). Tissue pieces were then minced in 1.35 mL Buffer C (150 mM KCl, 5 mM MgCl<sub>2</sub>, 10 mM  
273 HEPES KOH pH 7.4, 0.100 µM Calyculin A, 2 mM DTT, 100 U/mL RNAsin, 100 µg/mL cycloheximide,  
274 protease inhibitor cocktail, 5 mM sodium fluoride, 1 mM sodium orthovanadate, 1 mM sodium  
275 pyrophosphate, 1 mM β-glycerophosphate) and subsequently transferred to homogenization tubes for  
276 steady homogenization at 250 rpm three times and at 750 rpm nine times at 4 °C. Samples were then  
277 transferred to a 1.5 mL LoBind tube (Eppendorf 022431021) and clarified at 2000xg for 10 min at 4 °C.  
278 The low-speed supernatant was transferred to a new tube on ice, and 90 µL of NP40 (Sigma  
279 11332473001) and 90 µL of 1,2-diheptanoyl-sn-glycero-3-phosphocholine (DHPC, Avanti Polar Lipids  
280 850306P, 100 mg/0.69 mL) were added to this solution. This solution was mixed and then clarified at a  
281 max speed (17,000xg) for 10 min at 4 °C. The resulting high-speed supernatant was transferred to a  
282 new tube where 20 µL was saved and transferred to a tube containing 350 µL buffer RLT. To the  
283 remainder of the sample, 1.3 µL of 100 µg/mL cycloheximide, 27 µL of phosphatase inhibitor cocktail  
284 (250 mM sodium fluoride, 50 mM sodium orthovanadate, 50 mM sodium pyrophosphate, 50 mM β-  
285 glycerophosphate) and 6 µL of anti-pS6 antibody (Cell Signaling D68F8) were added. The sample was  
286 gently rotated for 90 min at 4 °C. To prepare beads, 100 µL of beads (Invitrogen 10002D) was washed  
287 three times with 900 µL of buffer A (150 mM KCl, 5 mM MgCl<sub>2</sub>, 10 mM HEPES KOH pH 7.4, 10% NP40,  
288 10% BSA), and once with 500 µL of buffer C. Sample homogenate was added to the beads and  
289 incubated with gentle rotation for 60 min at 4 °C. Following incubation, beads were washed with four



290 times with 700  $\mu$ L of buffer D (350 mM KCl, 5 mM MgCl<sub>2</sub>, 10 mM HEPES KOH pH 7.4, 10% NP40, 2 mM  
291 DTT, 100 U/mL RNAsin, 100  $\mu$ g/mL cycloheximide, 5 mM sodium fluoride, 1 mM sodium  
292 orthovanadate, 1 mM sodium pyrophosphate, 1 mM  $\beta$ -glycerophosphate). During the final wash,  
293 beads were moved to room temperature, wash buffer was removed, and 350 mL of buffer RLT was  
294 added. Beads were incubated in buffer RLT for 5 min at room temperature. Buffer RLT containing  
295 immunoprecipitated RNA was then eluted and stored at  $-80^{\circ}\text{C}$  until clean up using a kit (Qiagen  
296 74004). cDNA was generated using 11 rounds of amplification with 10 ng RNA input. DNA libraries  
297 were prepared using a half-sized Nexterra XT DNA Library Preparation Kit (Illumina 15032354) protocol  
298 as per the manufacturer's guidelines. Libraries were sequenced on either HiSeq 2000/2500 (50 base  
299 pair single read mode) or NextSeq 500 (75 base pair single read mode) with 6–12 pooled indexed  
300 libraries per lane.

301

### 302 **RNA-Seq alignment, quantification, and differential expression analysis**

303 Reads were aligned against a modified GRCm38.p6 (M25) reference, in which we deleted  
304 ENSMUSG00000116179 (*Olf290*), using STAR<sup>39</sup> with `--outFilterMultimapNmax 10`. Reads mapping to  
305 *Olf290* were inferred from ENSMUSG00000070459, with the rationale that this gene model included  
306 ENSMUSG00000116179 plus untranslated regions. Gene-level read quantification was done using  
307 RSEM<sup>40</sup>. Differential expression analysis was performed against all genes using EdgeR<sup>41</sup>. Gene  
308 nomenclature was retrieved from BioMart<sup>42</sup>. Intact *Olf* genes with identifiable sequences were  
309 filtered, and p-values were then re-corrected by FDR. Only ORs exhibiting odor response to at least one  
310 of the tested odorants ( $\log_2\text{FC} > 0$  and  $\text{FDR} < 0.05$ ) were considered. A total of 555 ORs responded  
311 across the 72 different odorants at various concentrations. A total of 375 ORs were responsive to  
312 unique odorants at the lowest tested concentrations. Raw and processed RNA-Seq datasets generated  
313 as part of this study are available from NCBI GEO at accession GSE185415.

314

315 **Source of odorants**

316 The following odors/concentrations were used for comparing molecular properties to receptor

317 responses: 1% 2-methyl-2-pentenal (Sigma 294667), 1% *trans*-cinnamaldehyde (Sigma C80687), 1% 2-

318 heptanone (Sigma 537683<sup>15</sup>), 1% linalool (Sigma L2602), 1% ethyl butyrate (Sigma W242713), 1%

319 guaiacol (Sigma G10903), 1% diacetyl (Sigma W237027), 1% 2-ethyl-3-methylpyrazine (Sigma

320 W315508), 1% 2,5-dimethylpyrazine (Sigma 175420<sup>15</sup>), 1% benzaldehyde (Sigma W212717), 1% (+)-

321 limonene (Sigma 183164), 1%  $\beta$ -damascone (Sigma W324300), 1%  $\alpha$ -pinene (Sigma W290267), 1% 2-

322 methyl-2-thiazoline (Sigma M83406), 1% citronellol (Sigma W230915), 1% dimethyl trisulfide (Sigma

323 W327506), 1% *p*-Cresol (Sigma C85751), 0.01% citral (Sigma W230316), 1 M (+)-menthol (Sigma

324 224464), 1 M (-)-menthol (Sigma M2780), 0.01% anisaldehyde (Sigma A88107), 1% 4-

325 methylacetophenone (Sigma W267708), 1% methyl salicylate (Sigma W274502), 1% (+)-carvone (Sigma

326 22070), 1% (-)-carvone (Sigma 22060), 1%  $\beta$ -ionone (Sigma W259525), 1% isopropyl tiglate (Sigma

327 W322903), 1% hexyl tiglate (Sigma W500909), 1% pyridine (Sigma 270970), 1% butyric acid (Sigma

328 W222119), 0.01% cyclopentanethiol (Sigma W326208), 0.01% 2-butene-1-thiol (1717 CheMall Corp

329 OR116574), 100 mM cyclopentadecanone (Sigma C111201), 1% 2-methyl-2-propanethiol (Sigma

330 109207), 0.01% acetophenone (Sigma W200910), 0.1% isovaleric acid (Sigma 129542), 1% isoamyl

331 acetate (Sigma 306967), 1% ethyl tiglate (Sigma W246000), 1% heptanoic acid (Sigma W334812), 10%

332 (+)-2-octanol (Sigma O4504), 10% (-)-2-octanol (Sigma 147990), 1% 2-hexanone (Sigma 103004), 1% 2-

333 phenylethanol (Sigma 77861), 1% 3-methyl-1-butanethiol (Sigma W385808), 1% octanal (Sigma

334 O5608), 1% heptanal (Sigma W254002), 1% 2,4,5-trimethylthiazole (nTMT, Sigma 219185), 100% (E)- $\beta$ -

335 Farnesene (Bedoukian P3500-90<sup>14</sup>), 100  $\mu$ M (methylthio)methanethiol (MTMT, synthesized<sup>15</sup>), 0.01%

336 2-*sec*-butyl-4,5-dihydrothiazole (SBT, synthesized<sup>15</sup>), 77% 3,4-dehydro-*exo*-brevicommin (DHB,

337 synthesized<sup>15</sup>), and 0.01% 2,4,5-trimethyl-4,5-dihydrothiazole (TMT, synthesized<sup>14</sup>).

338

339 For logistic regression and identifying residues with predictive power towards ligand selectivity,  
340 odorants tested at the lowest concentration with at least 8 activated ORs ( $\log_2FC > 0$  and  $FDR < 0.05$ )  
341 were used to promote class stability. Thus, following odors were removed from consideration using  
342 logistic regression compared to above: 100 mM cyclopentadecanone, 1% 2-heptanone, 1% 2-  
343 hexanone, 1% 3-methyl-1-butanethiol, 1%  $\alpha$ -pinene, 1% benzaldehyde, 1%  $\beta$ -ionone, 1% ethyl tiglate,  
344 1% heptanoic acid, 1% hexyl tiglate, 1% isopropyl tiglate, 1% linalool, 1% methyl salicylate, 1% (+)-  
345 limonene, 1% *trans*-cinnamaldehyde, and 0.1% isovaleric acid. The following odors were considered at  
346 a modified concentration from above: 0.1% TMT, 0.1% acetophenone, and 10 mM MTMT.

347

348 Odorants were excluded from all analysis if no ORs were identified as responsive at the tested  
349 concentrations: 1%  $\beta$ -Caryophyllene (Sigma W225207<sup>15</sup>), 1% dimethyl sulfide (Sigma 274380), 1%  
350 geraniol (Sigma W250716), 1% indole (Sigma W259378), 1% (-)-dihydrocarveol (Sigma 37278), 1% (+)-  
351 dihydrocarveol (Sigma 37277), 1% propionic acid (Sigma 109797), 3mM androstenone (Sigma 284998),  
352 or if the number of ORs identified as responsive were more than five standard deviations away from  
353 the mean: 1% 2'-hydroxyacetophenone (Sigma H18607).

354

### 355 **Chemical space estimation**

356 To estimate chemical space, we first identified 4680 small molecules commonly found in foods and  
357 fragrances from <http://www.thegoodscentscompany.com/><sup>16</sup>. Three dimensional structures for these  
358 molecules and the 52 in the test odor set were then downloaded from PubChem, and 5666 molecular  
359 properties were calculated using AlvaDesc (v2.0.10). From the 5666 calculated molecular properties,  
360 3855 were discarded because they were either not calculated for all molecules or exhibited zero

361 variance across all molecules, leaving behind 1811 molecular descriptors. Chemical space was

362 estimated by PCA dimensionality reduction on all molecules in R.

363

#### 364 **Receptor alignment and space estimation**

365 Mouse ORs were aligned to one another using the MAFFT E-INS-I method with manual refinements<sup>43</sup>.

366 The resulting alignment file was subjected to ModelTest-NG to identify ideal amino acid substitution

367 models<sup>44</sup>. Phylogenetic trees were generated with RAxML-NG using the JTT+I+G4 amino acid

368 substitution model with 100 bootstraps<sup>45</sup>. Receptor pairwise similarity matrices for multidimensional

369 scaling were generated from an alignment in which positions with amino acids in at least 60% of the

370 receptors were considered. Receptor pairwise similarity was calculated by summing amino acid

371 differences at each position by Grantham's amino acid distances<sup>17</sup>. Multidimensional scaling was done

372 in R.

373

#### 374 **Generating response spectra**

375 To generate odor response spectra, we first began with the log<sub>2</sub>FC values of each odor-responsive OR.

376 Each OR,  $r$ , was then centered and scaled (z-scored) by mean subtraction and standard deviation

377 division across the odorants,  $o$ , in the test panel. The resulting matrix is denoted as  $\tilde{\Delta}_{ro}$ . To generate

378 property strength vectors, each molecular property,  $p$ , was z-scored across the odorants in the test

379 panel. The resulting matrix is denoted as  $\tilde{P}_{op}$ . To calculate property responses and thereby property

380 response spectra (Pearson correlation coefficients), we used the following formula:

$$381 \quad \Phi_{rp} = \sum_o \tilde{\Delta}_{ro} \tilde{P}_{op}$$

382 where  $\Phi_{rp}$  refers to Pearson correlation coefficients between individual receptors,  $r$ , and molecular

383 properties,  $p$ .

384

385 To evaluate the significance of the correlation between a property and the response pattern of an OR,  
386 we used an FDR cutoff of 0.05.  $P$ -values were obtained by first calculating the  $t$ -statistic using  $t =$   
387  $\frac{r\sqrt{n-2}}{\sqrt{1-r^2}}$ , where  $r$  is the correlation coefficient and  $n$  is the number of data points. The two-tailed  $P$ -  
388 value was then calculated as twice the probability a  $t$ -distributed variable exceeds  $t$  using the python  
389 `scipy.stats.t.sf` function.  $P$ -values were adjusted by FDR correction in R.

390

### 391 **Odor distance calculation in property space and response space**

392 We calculated the distances between odorants in property and response space by calculating Euclidean  
393 distances, Pearson correlation coefficients, and cosine similarities between all possible unique pairs of  
394 odorants. Molecular properties and receptor response data were normalized to their respective zero  
395 mean and unit standard deviation. Correlation distances were reported as  $1 - r$ , and cosine distances  
396 were reported as  $1 - \cos(\theta)$ .

397

### 398 **Regression models with odor pair cross-validation**

399 Linear models (LASSO and ridge regression) were implemented with the `glmnet` package (v4.1) in R<sup>46</sup>.  
400 XGBoost was implemented with the `xgboost` module (v1.4.2) in python. Distances (Euclidean, cosine,  
401 and correlation) between each unique pair of odorants were calculated in normalized receptor  
402 response space. Then, Euclidean distances between each unique pair of odorants were calculated for  
403 each feature in normalized feature space. Regularization was then applied as either the L1 (LASSO) or  
404 L2 (ridge regression) norm. The  $\lambda$  loss function, which controls the number and relative contribution of  
405 selected features, was sequentially varied from zero to three by length 1000. Pearson correlation  
406 values were reported for the varied  $\lambda$  hyperparameters of the models by comparing model predicted  
407 response distances to true response distances. Shuffled controls consisted of using 52 fictitious

408 odorants whose individual feature vectors were generated by resampling with replacement across the  
409 52 test odorants.

410  
411 Default XGBoost model hyperparameters were used as follows: base\_score=0.5, booster='gbtree',

412 colsample\_bylevel=1, colsample\_bynode=1, colsample\_bytree=1, gamma=0, importance\_type='gain',  
413 interaction\_constraints='', learning\_rate=0.300000012, max\_delta\_step=0, max\_depth=6,

414 min\_child\_weight=1, missing=nan, monotone\_constraints='()', n\_estimators=100,

415 num\_parallel\_tree=1, random\_state=42, reg\_alpha=0, reg\_lambda=1, scale\_pos\_weight=1,

416 subsample=1, tree\_method='exact', validate\_parameters=1, and verbosity=None.

417  
418 Odor pair cross-validation was performed by iteratively holding out each unique pair of odorants from

419 the normalized 52 odor dataset (test data). Features with zero variance from the remaining 50 odor set  
420 were dropped (train data). Distances were calculated between the test data for each remaining feature

421 (xtest) and response pattern (ytest). Train data were normalized by z-scoring independent of test data.

422 Distances were then calculated between pairwise combinations of the 50 train odorants for each

423 feature (xtrain) and response pattern (ytrain). Feature distances between the held-out odor pair (xtest)

424 were then used to predict response pattern distances (ypred). Pearson correlation and mean squared

425 error values were reported from comparing model predicted response distances (ypred) to true

426 distances (ytest).

427

#### 428 **Optimized molecular property selection**

429 To select a subset of molecular properties that were well represented in the data, we utilized Support

430 Vector Machine (SVM) classifiers and regressors with linear kernels in the python sklearn.svm module

431 (v0.24.2). Beginning with the 1811 molecular properties, we first considered those that were non-

432 continuous (at least one zero entry, ex. molecular weight is a continuous molecular property). Non-  
433 zero values were set to one and zero values were kept. Classifiers were trained and cross-validated  
434 across the 52 odorant molecules using the normalized 375 deorphanized receptor responses as  
435 predictors in a leave-one-out scheme. Data normalization was first performed including the test data.  
436 After removing test data, training data were normalized independently to prevent contamination.  
437 Classifier area under receiver operating characteristic (AUROC) thresholds of 0.75 were applied.  
438 Molecular properties passing this threshold were next subjected to regression with non-zero entries  
439 restored. A Pearson correlation cutoff of 0.5 was applied to finally select the 65 “optimized” molecular  
440 properties.

441

#### 442 **Protein sequence analysis of ORs by logistic regression and SVM classifiers**

443 Regularized logistic regression was used to build models linking OR-protein sequence properties to OR-  
444 odor responses with the glmnet package (v4.1) in R<sup>46</sup>. ORs were classified as responders if they  
445 exhibited  $\log_2FC > 0$  and  $FDR < 0.05$  following pS6-IP-Seq and differential expression analysis. For  
446 odorants tested at multiple concentrations, the lowest concentration that activated at least 8 ORs was  
447 used to promote class stability. Predictors were generated from converting the FASTA alignment file  
448 into categorical variables reflecting the presence/absence of specific amino acids at each position.

449

450 To evaluate model performance, fitted odorants were randomly split into 90% training and 10% testing  
451 receptor sets. Each test set contained at least one responding receptor. Predictors with zero variance  
452 in the training set were dropped. The grid-search optimized  $\alpha$  hyperparameter (setting the ratio of the  
453 L1 and L2 norms) was set by ten-fold cross-validation with ten-fold cross-validation to set the  $\lambda$  (loss  
454 function) value.  $\lambda$  values one standard error of mean greater than optimal were selected to encourage  
455 statistically identical but sparser solutions. Model weighted predictors were then used to determine

456 the response likelihood of the test receptors. This procedure was repeated 100 times. Non-zero  
457 weights were averaged across repetitions and odorants to report positions with residues contributing  
458 predictive power towards odor selectivity. WebLogo visualizations were prepared at  
459 <http://weblogo.threeplusone.com/><sup>47</sup>.

460  
461 SVM classifier response probabilities were calculated using the same inputs as logistic regression using  
462 100 repetitions of 90% training (with ten-fold cross-validation for hyperparameter tuning) and 10%  
463 testing. Each repetition's response likelihoods and true outcomes were aggregated to generate a single  
464 ROC curve for a single odor, which were then combined to generate an aggregate ROC curve.

465  
466 **Protein sequence analysis of ORs by comparison to convergently evolved ORs**  
467 As an alternative strategy, we also performed a statistical evaluation of amino acid properties of ORs  
468 sharing responsiveness to an odor against convergently evolved receptors. First, responsive ORs  
469 ( $\log_2FC > 0$  and  $FDR < 0.05$  from differential expression) were subset, and pairwise Grantham distances  
470 were calculated at each position to generate Grantham distance distributions within the responsive OR  
471 alignment. Pairwise comparisons between gaps were considered to have zero distance while pairwise  
472 comparisons between gaps and amino acids were considered to have the average Grantham distance  
473 across all pairwise comparisons between all ORs at that position. Null distributions were generated  
474 similarly from convergently evolved odor-unresponsive ORs. To identify convergently evolved odor-  
475 unresponsive sets of ORs, odor-specific receptors with  $\log_2FC < 0$  or  $FDR > 0.25$  were first subset. Then,  
476 for each unique pairwise comparison between the odor-responsive ORs, full protein sequence  
477 Grantham distances were calculated. For each receptor in each pairwise comparison, the closest  
478 receptor was selected from the odor-unresponsive subset with the most similar absolute full protein  
479 sequence Grantham distance to the pairwise comparison. This meant, for each odor with some



480 number of responsive receptors, there was twice as many receptors identified as convergently evolved  
481 and odor-unresponsive. Distributions were compared using the Kolmogorov-Smirnov statistical test.  
482 FDR correction was applied across all calculated  $P$ -values with a cutoff of 0.05. The number of times  
483 responding receptors displayed statistically significant deviations in the distribution of Grantham  
484 distances from the null set, at each position, was counted and summed across all odorants.

485

#### 486 **Residue conservation calculation**

487 Using the 313 length alignment file, in which each position was occupied by an amino acid in at least  
488 60% of the responsive ORs (387 ORs that were responsive to the lowest concentration of tested  
489 odorants yielding response to at least 8 ORs each), we first identified the most common amino acid at  
490 each position. We term this the reduced consensus OR sequence. The percent presence of the most  
491 commonly occurring amino acid at each position was then reported as conservation percentage for  
492 said position.

493

#### 494 **Homology models**

495 To build an OR homology model, we adapted previously published methods<sup>48,49</sup>. The reduced  
496 consensus OR sequence was manually re-aligned to pre-aligned sequences of the bovine rhodopsin  
497 (PDB ID 1U19), the human chemokine receptors CXCR4 (3ODU) and CXCR1 (2LNL), and human  
498 adenosine A2A receptor (2YDV) using Jalview. Experimental GPCR structures of these receptors were  
499 then used as templates to build the homology model of the reduced consensus sequence with  
500 Modeller. Visualization and analysis of the homology model was done using VMD and Chimera.

501

#### 502 **Heterologous luciferase assay**

503 Hana3A cells; which stably express  $G_{olf}$ , RTP1, RTP2, and REEP1; were grown in minimum essential  
504 medium eagle (MEM; Corning 10-010-CV) containing 10% Fetal Bovine Serum (FBS; vol/vol; Gibco  
505 16000-044), penicillin-streptomycin (Sigma-Aldrich P4333), and amphotericin B (Gibco 15290018). Cells  
506 were cultured and incubated at 37°C, 5% CO<sub>2</sub>, and saturated humidity for use with the Dual-Glo  
507 Luciferase Assay (Promega E2980)<sup>29,50</sup>. Cells were plated at 20-25% confluence on poly-D-lysine-coated  
508 96-well plates (Corning 3843) overnight. After overnight incubation, cells were transfected with 6 mL of  
509 MEM containing 10% FBS, 0.5 µg SV40-RL (Promega E2980), 1 µg CRE-Luc (Promega E2980), 0.5 µg  
510 mouse RTP1s, 0.25 µg M3 muscarinic receptor<sup>51</sup>, 0.5 µg of Rho-tagged receptor plasmid DNA, and 20  
511 µg Lipofectamine 2000 (Invitrogen 11668019) per plate. Transfection medium was divided equally  
512 among the wells so that each OR-odorant combination could be conducted in triplicates. The following  
513 day, cells were incubated with 25µL of odorant solution diluted in CD-293 (Gibco 11913-019)  
514 containing 30 µM CuCl<sub>2</sub> (Sigma-Aldrich C-6641) and 2 mM glutamine (Gibco 25030-081) for 3.5 hours.  
515 cAMP-driven firefly Luciferase luminescence (Luc) was used to assess OR activation, and SV40-driven  
516 *Renilla* Luciferase luminescence (Ren) was used to control for variation in cell viability within wells. Cell  
517 luminescence was read by a POLARstar OPTIMA (BMG Labtech) luminometer, and normalized response  
518 values were calculated using the formula (Luc-400)/(Ren-400). ORs were considered responsive *in vitro*  
519 if ANOVA *p*-value was < 0.05 and ANOVA with post-hoc Dunnet's test correction *p*-adjusted was < 0.05  
520 for at least 2 of the tested odor concentrations using the R package DescTools (v0.99.42). Log-logistic  
521 4-parameter dose response curves were fit to the data using the R package drc (v3.0-1). *In vitro*  
522 responses were compared to *in vivo* responses by subtracting mean ligand-independent activity  
523 (luciferase values of ORs with no odor stimulation) from each of the ligand stimulated data points and  
524 summing. Scaled summed (+)-enantiomer responses were divided by scaled summed (-)-enantiomer  
525 responses and log<sub>2</sub> transformed for comparison to log<sub>2</sub>FC (+)/(-) *in vivo* enrichments.

526

527 **Data and code availability:**

528 Data and code are available upon reasonable request.

529 **References:**

- 530 1 Buck, L. & Axel, R. A novel multigene family may encode odorant receptors: a molecular basis  
531 for odor recognition. *Cell* **65**, 175-187, doi:10.1016/0092-8674(91)90418-x (1991).
- 532 2 Hallem, E. A. & Carlson, J. R. Coding of odors by a receptor repertoire. *Cell* **125**, 143-160,  
533 doi:10.1016/j.cell.2006.01.050 (2006).
- 534 3 Wang, G., Carey, A. F., Carlson, J. R. & Zwiebel, L. J. Molecular basis of odor coding in the  
535 malaria vector mosquito *Anopheles gambiae*. *Proc Natl Acad Sci U S A* **107**, 4418-4423,  
536 doi:10.1073/pnas.0913392107 (2010).
- 537 4 Carey, A. F., Wang, G., Su, C. Y., Zwiebel, L. J. & Carlson, J. R. Odorant reception in the malaria  
538 mosquito *Anopheles gambiae*. *Nature* **464**, 66-71, doi:10.1038/nature08834 (2010).
- 539 5 Saito, H., Chi, Q., Zhuang, H., Matsunami, H. & Mainland, J. D. Odor coding by a Mammalian  
540 receptor repertoire. *Sci Signal* **2**, ra9, doi:10.1126/scisignal.2000016 (2009).
- 541 6 Keller, A. *et al.* Predicting human olfactory perception from chemical features of odor  
542 molecules. *Science* **355**, 820-826, doi:10.1126/science.aal2014 (2017).
- 543 7 Xu, L. *et al.* Widespread receptor-driven modulation in peripheral olfactory coding. *Science* **368**,  
544 doi:10.1126/science.aaz5390 (2020).
- 545 8 Zhao, H. *et al.* Functional expression of a mammalian odorant receptor. *Science* **279**, 237-242,  
546 doi:10.1126/science.279.5348.237 (1998).
- 547 9 Araneda, R. C., Kini, A. D. & Firestein, S. The molecular receptive range of an odorant receptor.  
548 *Nat Neurosci* **3**, 1248-1255, doi:10.1038/81774 (2000).
- 549 10 Malnic, B., Hirono, J., Sato, T. & Buck, L. B. Combinatorial receptor codes for odors. *Cell* **96**, 713-  
550 723, doi:10.1016/s0092-8674(00)80581-4 (1999).
- 551 11 Jiang, Y. *et al.* Molecular profiling of activated olfactory neurons identifies odorant receptors for  
552 odors in vivo. *Nat Neurosci* **18**, 1446-1454, doi:10.1038/nn.4104 (2015).
- 553 12 von der Weid, B. *et al.* Large-scale transcriptional profiling of chemosensory neurons identifies  
554 receptor-ligand pairs in vivo. *Nat Neurosci* **18**, 1455-1463, doi:10.1038/nn.4100 (2015).
- 555 13 Knight, Z. A. *et al.* Molecular profiling of activated neurons by phosphorylated ribosome  
556 capture. *Cell* **151**, 1126-1137, doi:10.1016/j.cell.2012.10.039 (2012).
- 557 14 Hu, X. S. *et al.* Concentration-Dependent Recruitment of Mammalian Odorant Receptors.  
558 *eNeuro* **7**, doi:10.1523/eneuro.0103-19.2019 (2020).
- 559 15 Vihani, A. *et al.* Semiochemical responsive olfactory sensory neurons are sexually dimorphic and  
560 plastic. *Elife* **9**, doi:10.7554/eLife.54501 (2020).
- 561 16 Ravia, A. *et al.* A measure of smell enables the creation of olfactory metamers. *Nature* **588**, 118-  
562 123, doi:10.1038/s41586-020-2891-7 (2020).
- 563 17 Grantham, R. Amino acid difference formula to help explain protein evolution. *Science* **185**,  
564 862-864, doi:10.1126/science.185.4154.862 (1974).
- 565 18 Chae, H. *et al.* Mosaic representations of odors in the input and output layers of the mouse  
566 olfactory bulb. *Nat Neurosci* **22**, 1306-1317, doi:10.1038/s41593-019-0442-z (2019).
- 567 19 Pashkovski, S. L. *et al.* Structure and flexibility in cortical representations of odour space. *Nature*  
568 **583**, 253-258, doi:10.1038/s41586-020-2451-1 (2020).
- 569 20 Zhou, Q. *et al.* Common activation mechanism of class A GPCRs. *Elife* **8**,  
570 doi:10.7554/eLife.50279 (2019).
- 571 21 Liapakis, G. *et al.* The forgotten serine. A critical role for Ser-2035.42 in ligand binding to and  
572 activation of the beta 2-adrenergic receptor. *J Biol Chem* **275**, 37779-37788,  
573 doi:10.1074/jbc.M002092200 (2000).

- 574 22 Strader, C. D., Candelore, M. R., Hill, W. S., Sigal, I. S. & Dixon, R. A. Identification of two serine  
575 residues involved in agonist activation of the beta-adrenergic receptor. *J Biol Chem* **264**, 13572-  
576 13578 (1989).
- 577 23 Katritch, V. & Abagyan, R. GPCR agonist binding revealed by modeling and crystallography.  
578 *Trends Pharmacol Sci* **32**, 637-643, doi:10.1016/j.tips.2011.08.001 (2011).
- 579 24 Munk, C., Harpsoe, K., Hauser, A. S., Isberg, V. & Gloriam, D. E. Integrating structural and  
580 mutagenesis data to elucidate GPCR ligand binding. *Curr Opin Pharmacol* **30**, 51-58,  
581 doi:10.1016/j.coph.2016.07.003 (2016).
- 582 25 Surgand, J. S., Rodrigo, J., Kellenberger, E. & Rognan, D. A chemogenomic analysis of the  
583 transmembrane binding cavity of human G-protein-coupled receptors. *Proteins* **62**, 509-538,  
584 doi:10.1002/prot.20768 (2006).
- 585 26 Shi, L. *et al.* Beta2 adrenergic receptor activation. Modulation of the proline kink in  
586 transmembrane 6 by a rotamer toggle switch. *J Biol Chem* **277**, 40989-40996,  
587 doi:10.1074/jbc.M206801200 (2002).
- 588 27 Eddy, M. T. *et al.* Allosteric Coupling of Drug Binding and Intracellular Signaling in the A2A  
589 Adenosine Receptor. *Cell* **172**, 68-80 e12, doi:10.1038/nbt.4096 (2018).
- 590 28 de March, C. A. *et al.* Conserved Residues Control Activation of Mammalian G Protein-Coupled  
591 Odorant Receptors. *J Am Chem Soc* **137**, 8611-8616, doi:10.1021/jacs.5b04659 (2015).
- 592 29 Saito, H., Kubota, M., Roberts, R. W., Chi, Q. & Matsunami, H. RTP family members induce  
593 functional expression of mammalian odorant receptors. *Cell* **119**, 679-691,  
594 doi:10.1016/j.cell.2004.11.021 (2004).
- 595 30 Larsson, M. C. *et al.* Or83b encodes a broadly expressed odorant receptor essential for  
596 *Drosophila* olfaction. *Neuron* **43**, 703-714, doi:10.1016/j.neuron.2004.08.019 (2004).
- 597 31 Sanchez-Lengeling, B. *et al.* Machine Learning for Scent: Learning Generalizable Perceptual  
598 Representations of Small Molecules. arXiv:1910.10685 (2019).  
599 <https://ui.adsabs.harvard.edu/abs/2019arXiv191010685S>.
- 600 32 Tran, N., Kepple, D., Shuvaev, S. & Koulakov, A. in *Proceedings of the 36th International*  
601 *Conference on Machine Learning* Vol. 97 (eds Chaudhuri Kamalika & Salakhutdinov Ruslan)  
602 6305--6314 (PMLR, Proceedings of Machine Learning Research, 2019).
- 603 33 Katada, S., Hirokawa, T., Oka, Y., Suwa, M. & Touhara, K. Structural basis for a broad but  
604 selective ligand spectrum of a mouse olfactory receptor: mapping the odorant-binding site. *J*  
605 *Neurosci* **25**, 1806-1815, doi:10.1523/jneurosci.4723-04.2005 (2005).
- 606 34 Yu, Y. *et al.* Responsiveness of G protein-coupled odorant receptors is partially attributed to the  
607 activation mechanism. *Proc Natl Acad Sci U S A* **112**, 14966-14971,  
608 doi:10.1073/pnas.1517510112 (2015).
- 609 35 Sekharan, S. *et al.* QM/MM model of the mouse olfactory receptor MOR244-3 validated by site-  
610 directed mutagenesis experiments. *Biophys J* **107**, L5-l8, doi:10.1016/j.bpj.2014.07.031 (2014).
- 611 36 Baud, O. *et al.* Exchanging ligand-binding specificity between a pair of mouse olfactory receptor  
612 paralogs reveals odorant recognition principles. *Sci Rep* **5**, 14948, doi:10.1038/srep14948  
613 (2015).
- 614 37 Gelis, L., Wolf, S., Hatt, H., Neuhaus, E. M. & Gerwert, K. Prediction of a ligand-binding niche  
615 within a human olfactory receptor by combining site-directed mutagenesis with dynamic  
616 homology modeling. *Angew Chem Int Ed Engl* **51**, 1274-1278, doi:10.1002/anie.201103980  
617 (2012).
- 618 38 Del Marmol, J., Yedlin, M. A. & Ruta, V. The structural basis of odorant recognition in insect  
619 olfactory receptors. *Nature*, doi:10.1038/s41586-021-03794-8 (2021).
- 620 39 Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21,  
621 doi:10.1093/bioinformatics/bts635 (2013).

- 622 40 Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-Seq data with or  
623 without a reference genome. *BMC Bioinformatics* **12**, 323, doi:10.1186/1471-2105-12-323  
624 (2011).
- 625 41 Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for differential  
626 expression analysis of digital gene expression data. *Bioinformatics* **26**, 139-140,  
627 doi:10.1093/bioinformatics/btp616 (2010).
- 628 42 Smedley, D. *et al.* BioMart--biological queries made easy. *BMC Genomics* **10**, 22,  
629 doi:10.1186/1471-2164-10-22 (2009).
- 630 43 Katoh, K., Kuma, K., Toh, H. & Miyata, T. MAFFT version 5: improvement in accuracy of multiple  
631 sequence alignment. *Nucleic Acids Res* **33**, 511-518, doi:10.1093/nar/gki198 (2005).
- 632 44 Darriba, D. *et al.* ModelTest-NG: A New and Scalable Tool for the Selection of DNA and Protein  
633 Evolutionary Models. *Mol Biol Evol* **37**, 291-294, doi:10.1093/molbev/msz189 (2020).
- 634 45 Kozlov, A. M., Darriba, D., Flouri, T., Morel, B. & Stamatakis, A. RAxML-NG: a fast, scalable and  
635 user-friendly tool for maximum likelihood phylogenetic inference. *Bioinformatics* **35**, 4453-  
636 4455, doi:10.1093/bioinformatics/btz305 (2019).
- 637 46 Friedman, J., Hastie, T. & Tibshirani, R. Regularization Paths for Generalized Linear Models via  
638 Coordinate Descent. *J Stat Softw* **33**, 1-22 (2010).
- 639 47 Crooks, G. E., Hon, G., Chandonia, J. M. & Brenner, S. E. WebLogo: a sequence logo generator.  
640 *Genome Res* **14**, 1188-1190, doi:10.1101/gr.849004 (2004).
- 641 48 de March, C. A., Kim, S. K., Antonczak, S., Goddard, W. A., 3rd & Golebiowski, J. G protein-  
642 coupled odorant receptors: From sequence to structure. *Protein Sci* **24**, 1543-1548,  
643 doi:10.1002/pro.2717 (2015).
- 644 49 Bushdid, C., de March, C. A., Matsunami, H. & Golebiowski, J. Numerical Models and In Vitro  
645 Assays to Study Odorant Receptors. *Methods Mol Biol* **1820**, 77-93, doi:10.1007/978-1-4939-  
646 8609-5\_7 (2018).
- 647 50 Zhuang, H. & Matsunami, H. Evaluating cell-surface expression and measuring activation of  
648 mammalian odorant receptors in heterologous cells. *Nat Protoc* **3**, 1402-1413,  
649 doi:10.1038/nprot.2008.120 (2008).
- 650 51 Li, Y. R. & Matsunami, H. Activation state of the M3 muscarinic acetylcholine receptor  
651 modulates mammalian odorant receptor signaling. *Sci Signal* **4**, ra1,  
652 doi:10.1126/scisignal.2001230 (2011).
- 653