

Enhanced Sequence-Activity Mapping and Evolution of Artificial Metalloenzymes by Active Learning

Tobias Vornholt^{1,2*}, Mojmír Mutný^{3*}, Gregor W. Schmidt¹, Christian Schellhaas¹, Ryo Tachibana⁴, Sven Panke^{1,2}, Thomas R. Ward^{2,4‡}, Andreas Krause^{3‡}, Markus Jeschek^{1,5‡}

¹Department of Biosystems Science and Engineering, ETH Zurich, Mattenstrasse 26, 4058 Basel, Switzerland

²National Centre of Competence in Research (NCCR) Molecular Systems Engineering, Switzerland; Web: www.nccr-mse.ch

³Department of Computer Science, ETH Zurich, Andreasstrasse 5, 8092 Zurich, Switzerland

⁴Department of Chemistry, University of Basel, Mattenstrasse 24a, 4058 Basel, Switzerland

⁵Institute of Microbiology, University of Regensburg, Universitätsstraße 31, 93053 Regensburg, Germany

* These authors contributed equally

‡ Corresponding authors

Abstract

Tailored enzymes hold great potential to accelerate the transition to a sustainable bioeconomy. Yet, enzyme engineering remains challenging as it relies largely on serendipity and is, therefore, highly laborious and prone to failure. The efficiency and success rates of engineering campaigns may be improved substantially by applying machine learning to construct a comprehensive representation of the sequence-activity landscape from small sets of experimental data. However, it often proves challenging to reliably model a large protein sequence space while keeping the experimental effort tractable. To address this challenge, we present an integrated pipeline combining large-scale screening with active machine learning and model-guided library design. We applied this strategy to efficiently engineer an artificial metalloenzyme (ArM) catalysing a new-to-nature hydroamination reaction. By combining lab automation and next-generation sequencing, we acquired sequence-activity data for several thousand ArM variants. We then used Gaussian process regression to model the activity landscape and guide further screening rounds according to user-defined objectives. Crucial characteristics of our enhanced enzyme engineering pipeline include i) the cost-effective generation of information-rich experimental data sets, ii) the integration of an explorative round to improve the performance of the model, as well as iii) the consideration of experimental noise during modelling. Our approach led to an order-of-magnitude boost in the hit rate of screening while making efficient use of experimental resources. Smart search strategies like this should find broad utility in enzyme engineering and accelerate the development of novel biocatalysts.

33 Introduction

34 Biocatalysis and metabolic engineering offer sustainable production routes for many compounds of
35 interest and thus hold the potential to transform various industries. However, extensive enzyme
36 engineering is typically required to obtain a suitable biocatalyst for a desired application. This is often
37 a time-consuming, empirical process whose outcome is subject to chance, as classical methods are
38 agnostic to the topology of the underlying sequence-activity landscape. Engineering strategies that
39 incorporate machine learning to model this landscape could render enzyme engineering more efficient
40 and increase the likelihood of identifying an optimal solution. Accordingly, machine learning-assisted
41 directed evolution (MLDE) has attracted significant attention in recent years¹⁻³.

42 In general, MLDE starts with an initial screening round in which both sequence and activity are
43 recorded for a number of enzyme variants. These sequence-activity data are then used to train a
44 machine learning model, with the objective of predicting the activity of untested variants directly from
45 their sequence. If successful, such models can suggest variants that are likely to be highly active and
46 thus support further screening rounds by *in silico* library design¹. Further, the model can be iteratively
47 updated with new data to improve its predictive performance, a strategy referred to as active learning.
48 While several studies have demonstrated the general feasibility of such approaches⁴⁻¹², there are still
49 various challenges that need to be addressed to maximize the success rate and efficiency of MLDE and
50 enable its widespread implementation. This pertains to various aspects such as library design,
51 experimental data acquisition, model development, and the strategy for sampling the sequence space.

52 With regard to library design, the crucial challenge is to create a library that is as information-dense as
53 possible to allow for the development of accurate models while keeping the screening effort
54 manageable. In the initial stages of model development, this calls for libraries that exhibit a high degree
55 of sequence diversity to provide adequate information on the underlying sequence space, while at the
56 same time containing a sufficient number of active mutants¹³. These requirements can be difficult to
57 reconcile, as simultaneous randomization of multiple residues commonly results in a large fraction of
58 inactive mutants, from which little to no meaningful information for model training can be extracted.

59 Once a library has been generated, it is often challenging to measure a sufficiently large set of
60 sequence-activity data. In some cases, high-throughput assays such as fluorescence-activated cell
61 sorting can be combined with deep sequencing to obtain very large data sets^{14,15}. However, most
62 enzymatic reactions of industrial relevance require more laborious analytical procedures to obtain a
63 readout for activity. Moreover, the need to also obtain sequence information on all tested variants can
64 lead to prohibitive costs if conventional Sanger sequencing is used. Consequently, most studies to date
65 have relied on small data sets (10^1 - 10^2 variants)⁴⁻¹⁰. While this has led to several successful
66 demonstrations of MLDE, larger data sets are likely to lead to more accurate machine learning models
67 and improve the chances of identifying variants with the desired properties¹¹, particularly as the search
68 space increases in size.

69 Beyond these experimental considerations, several critical decisions have to be made regarding the
70 machine learning strategy. Prominent examples in this regard include the encoding strategy for the
71 protein sequences and the choice of a suitable machine learning algorithm. Many encoding strategies
72 have been suggested for creating a meaningful representation of protein variants, ranging from simple
73 one-hot encoding and descriptors based on amino acid properties¹⁶⁻¹⁸ to structure-based
74 descriptors^{19,20} and learned embeddings^{21,22}. Similarly, various machine learning algorithms have been
75 employed or suggested for MLDE, including linear regression²³⁻²⁵, Gaussian processes^{4,7-9,25,26}, and
76 neural networks¹². While the best strategy depends on the data set and task at hand, Gaussian
77 processes have repeatedly revealed their utility for active learning^{8,9,25}.

78 Less attention has been devoted to other aspects of the machine learning process, such as the handling
79 of experimental noise or the sampling strategy during ML-guided screening rounds, both of which are
80 critical to the success and efficiency of MLDE. With regard to the sampling strategy, many studies have
81 relied on a single training phase followed by greedy sampling of the top predictions of the resulting
82 model. Due to inevitable biases in library generation and the limitations in generating sufficient
83 sequence-activity data, this is unlikely to result in a comprehensive and accurate representation of the
84 sequence-activity landscape. Consequently, such models may be “blind” for promising regions of the
85 sequence space, leading to suboptimal outcomes such as low hit rates. Active learning strategies that
86 improve the model in iterative cycles of experiments and machine learning may help to develop a
87 better representation of the sequence-activity landscape, as these can converge to the optimal
88 solution over time²⁷. However, the aforementioned bottleneck in experimental data generation makes
89 performing many iterations undesirable. Thus, resources invested into model improvement (i.e.,
90 exploration) must be carefully weighed against the focus on regions of the sequence space that are
91 likely to contain active variants but might only comprise local optima (exploitation). In addition, activity
92 may not be the only selection criterion during exploitation. Instead, it is often desirable to sample
93 various potential optima to obtain a diverse set of variants, which requires more elaborate approaches
94 than simple greedy selection of top predictions²⁸. Hence, smart sampling strategies for active learning
95 are required to maximize the chances of success at a given experimental budget.

96 In this study, we introduce an integrated experimental and computational pipeline that addresses
97 critical limitations in the MLDE of enzymes. Specifically, we combine informed library design with large-
98 scale screening and a novel active machine-learning strategy. As an impactful testbed, we selected an
99 artificial metalloenzyme (ArM) for gold-catalysed hydroamination, a new-to-nature reaction for atom-
100 economical C-N bond formation. We simultaneously engineered five crucial amino acid residues in this
101 ArM, corresponding to a search space of 3,200,000 possible variants. To sample this space, we
102 combined lab automation with a cost-efficient next-generation sequencing (NGS) strategy, which
103 allowed us to acquire sequence-activity data on more than 2,000 ArM variants. Furthermore, we
104 developed a machine learning model based on Gaussian process regression that incorporates
105 optimized descriptors and estimates of experimental noise to efficiently navigate the sequence space.
106 Guided by the model’s uncertainty estimates, we performed a second screening round focused on
107 exploration and model refinement. Importantly, our results demonstrate that this targeted exploration
108 substantially improved the model’s performance. The optimized model reliably proposed highly active
109 ArM variants in a final exploitation round, as illustrated by a 12-fold increased hit rate compared to
110 the initial library.

111 **Results**

112 **Design of an information-dense ArM library**

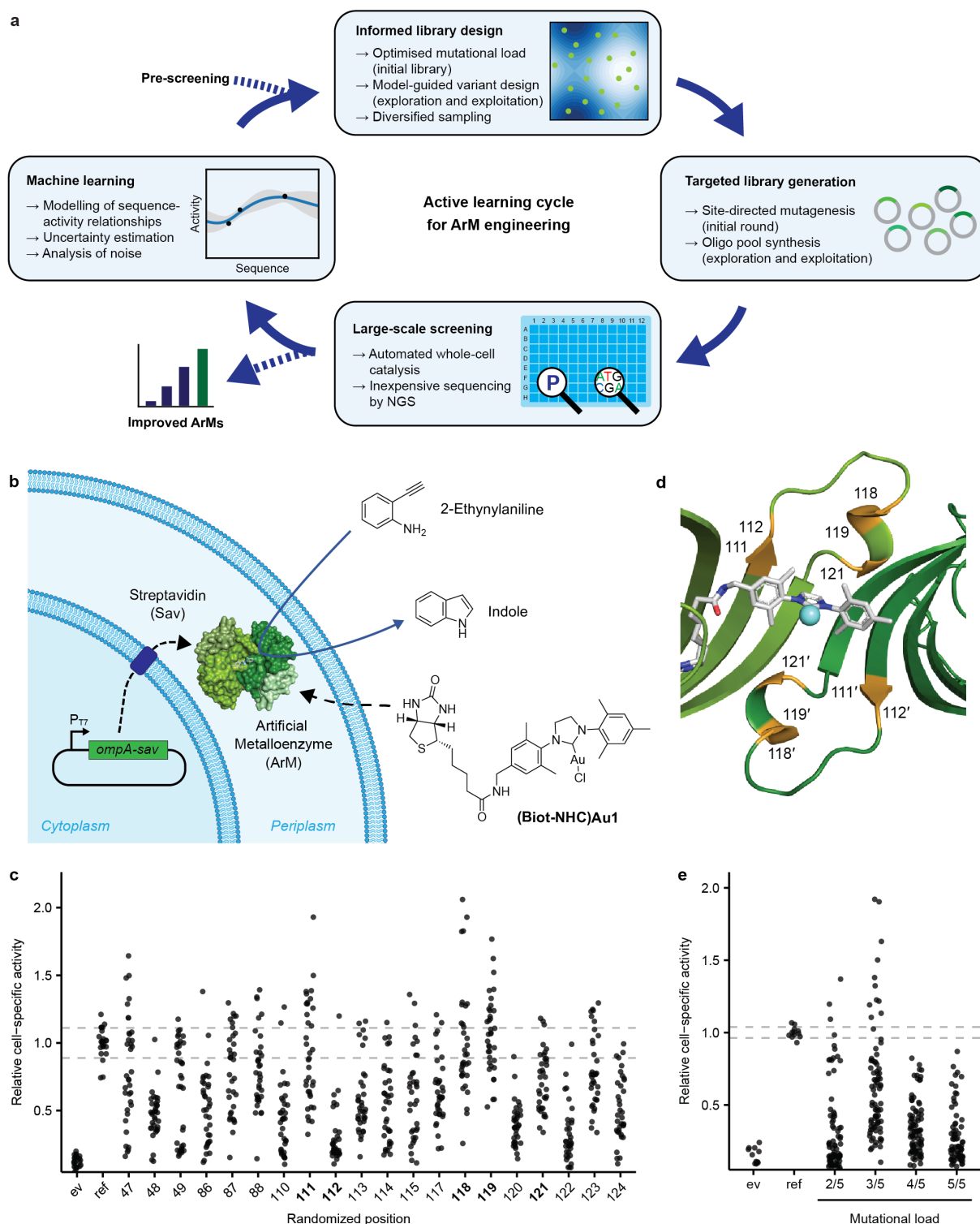
113 ArMs are hybrid catalysts that promise to significantly increase the number of reactions available in
114 biocatalysis by equipping enzymes with the catalytic versatility of abiotic transition metal cofactors²⁹.
115 ArMs have been created for a variety of natural and non-natural reactions^{30–35}, and some have
116 demonstrated catalytic prowess comparable to that of natural enzymes^{36–39}. However, most ArMs
117 initially display a low activity, and extensive protein engineering is required to identify catalytically
118 proficient variants. This engineering is typically a labour-intensive and slow process. Therefore, ArMs
119 represent an impactful yet challenging use case for MLDE.

120 A particularly versatile strategy for creating ArMs is to incorporate an organometallic cofactor into the
121 tetrameric protein streptavidin (Sav) using a biotin moiety as the anchor. Using this approach, we have
122 previously engineered an ArM for gold-catalysed hydroamination by exhaustively screening a library
123 of 400 Sav double mutants (Sav S112X K121X) using a whole-cell assay in 96-well plates⁴⁰. While this

124 represents an attractive starting point, extending the search space to more positions offers the
125 opportunity to achieve further improvements, which will be crucial for adapting ArMs for real-world
126 applications. However, exhaustive screening quickly becomes intractable in this case, and smart
127 heuristics for the efficient exploration of the underlying sequence-activity landscape are essential⁴¹.

128 To navigate the sequence-activity landscape of the ArM, we devised an iterative active learning cycle
129 involving library design, cloning, screening, and machine learning (Fig. 1a). With regard to library
130 design, the first step is to choose the target residues and a randomization scheme. To maximize the
131 potential impact of the screening campaign, we aimed to find important positions in Sav besides the
132 previously identified residues S112 and K121⁴⁰. Thus, we individually randomized the 20 residues
133 closest to the biotinylated gold cofactor in Sav S112F K121Q, which is the most active variant we had
134 observed before⁴⁰ (referred to as “reference variant” herein). Randomization was performed using
135 degenerate NDT (N = A, C, G or T; D = A, G or T) codons, which encode 12 amino acids covering all
136 chemical classes of amino acids, a strategy that has revealed high success rates at a reduced screening
137 effort⁴⁰. Subsequently, we measured hydroamination activity using our previously established protocol
138 relying on periplasmic catalysis in *Escherichia coli* (Fig. 1b)⁴⁰. We tested 36 clones per randomized
139 position to achieve a statistical library coverage of approximately 95 %⁴². As expected, most variants
140 displayed reduced activity compared to the reference variant (Fig. 1c). Notably, positions 111, 118, and
141 119 revealed the highest potential for improvement upon mutagenesis, with several variants
142 outperforming the reference variant. Consequently, we selected these positions for further
143 engineering. In addition, we chose to also randomize positions 112 and 121 again, as our observations
144 had indicated that epistatic effects play an important role in highly active ArM mutants⁴⁰.

145 Next, we sought to create a combinatorial library of the five selected positions (111, 112, 118, 119,
146 and 121, Fig. 1d), which, upon full randomization, corresponds to a search space of $20^5 = 3,200,000$
147 variants. This greatly exceeds the capacity of typical activity assays and well plate-based screenings.
148 Thus, navigating the underlying sequence-activity landscape represents a significant challenge. In
149 order to model this space for MLDE, it is crucial to design a library that offers a good coverage of the
150 targeted sequence space and at the same time maintains a sufficient proportion of active variants¹³.
151 While simultaneous randomization of all five residues would fulfil the first criterion, we anticipated
152 that the high mutational load would likely lead to a large fraction of inactive variants. This would not
153 only diminish the chances of identifying improved variants but, importantly, would be uninformative
154 for machine learning. Upon initial tests, we indeed observed a marked drop in the activity distribution
155 when randomizing more than three of the five positions simultaneously (Fig. 1e). Accordingly, we set
156 out to construct a library with three to four mutations distributed across the five target residues as a
157 good compromise between high sequence-diversity and sufficient residual activity. In other words, the
158 constructed library covers all five target positions, but individual variants contain at most four amino
159 acid substitutions relative to the reference variant Sav S112F K121Q, which served as the parent of the
160 library (Supplementary Fig. 1). This was achieved by site-directed mutagenesis PCR using various sets
161 of primers containing degenerate NNK (K = G or T) codons at different positions and subsequent mixing
162 of the resulting sub-libraries (see Methods).



163
 164 **Fig. 1 | Engineering strategy and library design for ArMs catalysing hydroamination.** **a**, Illustration of the active
 165 learning strategy for ArM engineering. An iterative process of library design, cloning, large-scale screening and
 166 machine learning was used to model the sequence-activity landscape and identify improved ArMs. Crucial steps
 167 and considerations are highlighted and are explained in the main text. **b**, Illustration of whole-cell biocatalysis
 168 using an ArM in the periplasm of *E. coli*. Sav is exported to the periplasm by means of an N-terminal OmpA signal
 169 peptide, where it binds the biotinylated cofactor **(Biot-NHC)Au1**. The resulting ArM converts 2-ethynylaniline to
 170 indole in a new-to-nature hydroamination reaction. Indole can subsequently be quantified using a colorimetric
 171 assay. **c**, Single site-saturation mutagenesis to identify influential amino acid residues with respect to ArM
 172 activity. Starting from the reference variant Sav S112F K121Q, 20 residues in Sav were individually mutated using
 173 degenerate NDT codons. The activity of the resulting variants is displayed relative to the mean activity of the
 174 reference variant (“ref”). Dashed lines indicate one standard deviation around the mean activity of the reference
 175 variant, which was measured in triplicate in each 96-well plate. A strain lacking Sav, i.e., containing an empty

176 vector (“ev”), was included as a control (n = 3 per 96-well plate). The five positions selected for combinatorial
177 randomization are highlighted in bold. Note that no improvement was expected at positions 112 and 121, as the
178 reference variant had already been optimized with regard to these positions⁴⁰. **d**, Residues selected for
179 randomization (highlighted in orange) in a ribbon model of Sav harbouring a metathesis catalyst (PDB 5IRA). For
180 clarity, only two biotin-binding sites of two opposing Sav monomers (a so-called functional dimer) are displayed.
181 **e**, Effect of different multi-site randomization strategies on the activity distribution of ArM libraries. Starting from
182 the reference variant, either two, three, four or five residues amongst positions 111, 112, 118, 119, and 121 were
183 randomized simultaneously. Hydroamination activity is displayed relative to the average activity of the reference
184 variant (“ref”, n = 3 per 96-well plate) for 90 variants from each library. A strain containing an empty vector (“ev”)
185 was included as a control (n = 3 per 96-well plate).

186 **Large-scale acquisition of sequence-activity data**

187 Our previously established whole-cell screening protocol for ArMs relied on periplasmic Sav
188 expression, ArM assembly and catalysis in 96-well plate format. By combining this protocol with
189 conventional Sanger sequencing, we were able to obtain sequence-activity data for a few hundred
190 variants⁴⁰. Although this platform was more flexible and simpler than comparable screening strategies
191 involving protein purification, it still required considerable manual labour, particularly for product
192 quantification. Additionally, when larger data sets are required, Sanger sequencing rapidly leads to
193 prohibitively high sequencing costs. To facilitate the generation of larger data sets for MLDE, we thus
194 sought to minimize manual intervention in the activity assay and develop more cost-efficient means
195 of obtaining the sequence information for each functionally characterised variant.

196 First, we automated all steps in the assay protocol that are labour-intensive (and thus limiting in terms
197 of throughput) or critical for reproducibility. Specifically, we made use of a Tecan EVO 200 platform for
198 all steps from colony picking to product quantification, with the exception of Sav expression in 96-deep
199 well plates, which only requires a small number of pipetting steps (Fig. 2a). The most important
200 addition to our previous semi-automated pipeline⁴⁰ is the photometric quantification of the product
201 indole. While this is a laborious procedure when carried out manually, the automated version
202 simplifies screenings and proved to be very reproducible (Supplementary Fig. 2). As the robotic
203 platform can handle up to eight 96-well plates at the same time, it greatly accelerates the acquisition
204 of large data sets.

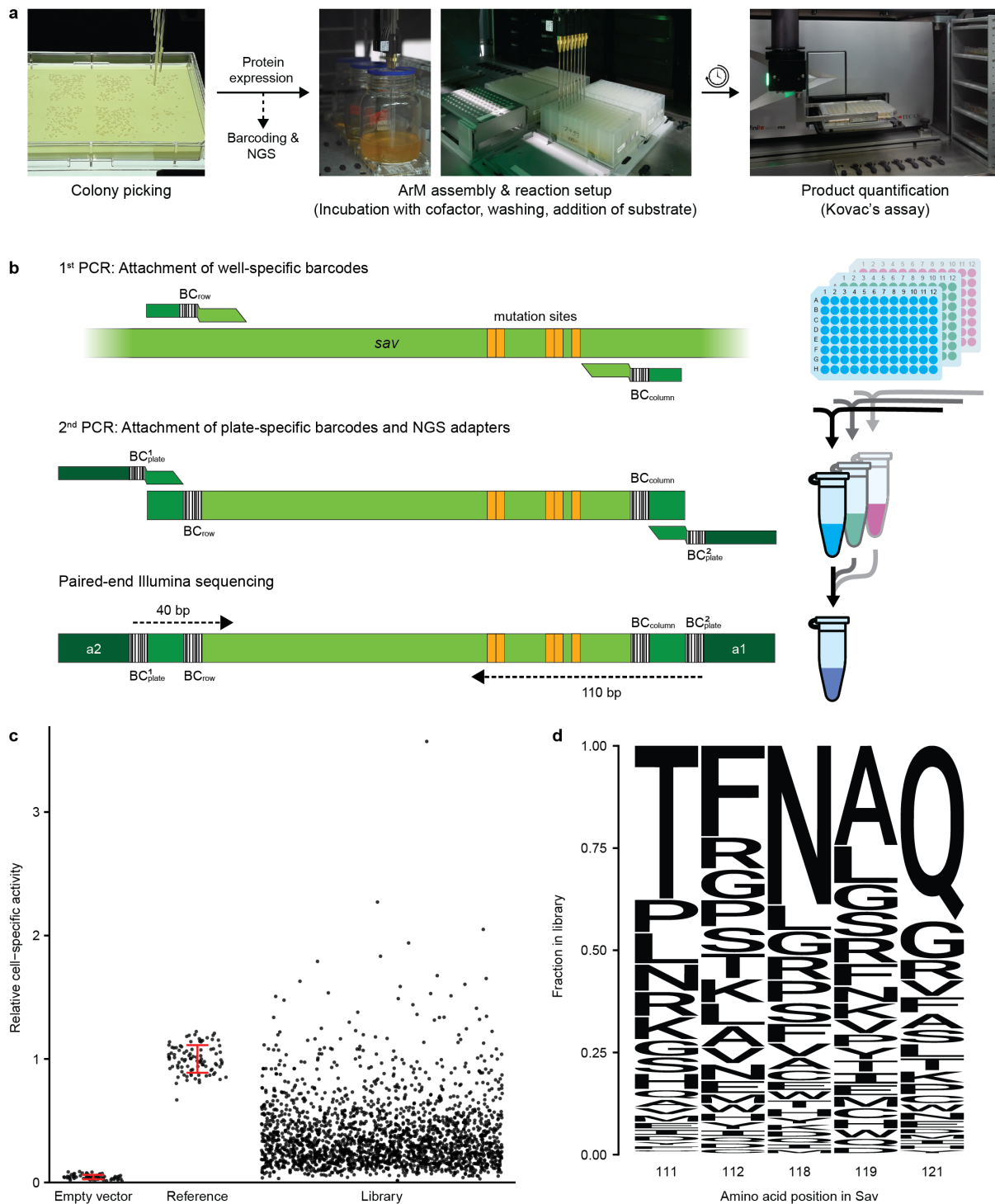
205 Besides the activity assay, another critical barrier to obtaining sufficiently large sets of sequence-
206 activity data can be the cost of sequencing. Obtaining the sequences of several thousand protein
207 variants by Sanger sequencing typically costs more than USD 10,000, which is prohibitive for most
208 academic labs. In principle, the cost per variant can be reduced significantly by relying on NGS, which
209 quickly becomes more cost-efficient than Sanger sequencing as the library size increases. However, in
210 NGS all variants are sequenced in bulk, which means a method to retroactively link each sequence to
211 the corresponding activity measurement is required. Previously, the use of DNA barcodes has been
212 suggested to enable NGS of protein variants distributed across 96-well plates^{43–45}. Building on these
213 strategies, we established a two-step PCR protocol for the barcoding of Sav variants that is compatible
214 with the Illumina NGS platforms (Fig. 2b). In the first step, which is carried out in 96-well plates, the
215 randomized region of the Sav gene is amplified using primers that append a well-specific barcode
216 combination as well as constant regions to the ends of the PCR products. This is achieved using eight
217 forward (representing the plate’s rows) and twelve reverse primers (representing the columns). For
218 simplicity, heat-treated samples of bacterial cultures serve as templates, avoiding the need for
219 laborious and costly plasmid purification.

220 Subsequently, PCR products are pooled by plate, and each pool is gel-purified and used as a template
221 for a second PCR. In this step, primers binding to the previously added terminal constant regions are
222 used for amplification. These primers contain overhangs to append plate-specific barcodes as well as
223 the adapters required for NGS. Through the combination of well- (1st step) and plate-specific (2nd step)

224 barcodes, it is possible to sequence thousands of variants from multiple plates in a single, low-cost
225 NGS run and to assign the obtained sequences to the corresponding activity value obtained in the
226 functional assay. In our specific case, paired-end sequencing of 40 bp from one end and 110 bp from
227 the other end of the final PCR product was sufficient to read all well- and plate-specific barcodes as
228 well as the five mutation sites in the Sav gene at a high read coverage (average of >100-fold per variant)
229 and low cost (see Discussion).

230 Relying on the combination of automated activity assay and NGS, we screened 32 96-well plates
231 containing variants from the aforementioned library of Sav. As each plate contained six controls (empty
232 vector and reference variant in triplicate), this amounts to a total of 2,880 variants. Excluding mutants
233 that failed to grow, we obtained activity data on 2,790 variants. Most of these displayed an
234 intermediate activity between the background level of cells lacking Sav (empty vector) and the
235 reference variant Sav S112F K121Q (Fig. 2c). Notably, approximately 3 % of all mutants were more
236 active than the reference. Using the NGS-based strategy, we retrieved the sequences for 2,663 out of
237 2,880 wells containing Sav mutants. After excluding variants with nonsense mutations and wells
238 containing more than one variant, sequence-activity data for 2,164 clones were obtained, of which
239 2,035 were distinct variants. Notably, for variants appearing in multiple wells, the deviation between
240 these replicate activity measurements was generally low, corroborating the high robustness of the
241 assay (Supplementary Fig. 3). Importantly, the library displayed a high degree of sequence diversity,
242 with every amino acid appearing in every position (Fig. 2d) and an average Hamming distance of 4.3
243 between the mutants. Note that the amino acids of the reference variant were the most abundant in
244 each position, as we did not randomize all five positions simultaneously. Thus, the library exhibited a
245 high degree of variability both in terms of activity distribution (including a low fraction of inactive
246 variants) as well as sequence diversity. This indicated that the aforementioned design goals for the
247 library were met, providing a promising data basis for modelling the sequence-activity landscape by
248 machine learning.

249 As we had previously recorded sequence-activity data for 400 Sav double mutants (S112X K121X) that
250 are part of the same sequence space⁴⁰, we added these older data to the measurements obtained
251 herein. As a result, a total of 2,992 data points covering 2,435 distinct ArM variants were available as
252 initial training data for machine learning.



253
254
255
256
257
258
259
260
261
262
263
264
265
266

Fig. 2 | Large-scale acquisition of sequence-activity data for ArMs. **a**, Depiction of the critical automated steps in the screening workflow. Colony picking, ArM assembly, reaction setup, and product quantification were performed on a lab automation platform. The less labour-intensive protein expression protocol was performed manually. In parallel to the activity assay, samples of the starter cultures were processed further for NGS. **b**, PCR-based barcoding strategy for cost-effective sequencing of Sav variants in 96-well plates by NGS. First, the mutated region of the Sav gene is amplified using primers with row- (BC_{row}) and column-specific (BC_{column}) DNA barcodes. This step is performed in PCR plates using heat-treated bacterial cultures as templates. After pooling all samples from one plate, a second PCR is performed to add two plate-specific barcodes (BC_{plate}) as well as adapters required for Illumina sequencing (a1 and a2). Subsequently, all samples are pooled and sequenced via paired-end reading to cover all barcodes and mutation sites. **c**, Cell-specific hydroamination activity of 2,164 ArM variants from the initial library obtained by automated screening of 32 96-well plates. Only variants that were included for model training are displayed. Controls (empty vector and reference variant) are displayed with their standard deviation in red. **d**, Fraction of amino acids at the five randomized positions in Sav. Note that the amino

267 acids of the reference variant (Sav 111T 112F 118N 119A 121Q, abbreviated Sav TFNAQ) are the most abundant,
268 as the library was derived from this variant and contained at most four amino acid substitutions per variant.

269 **Development of an initial machine learning model of ArM activity**

270 To construct a model that can reliably predict the activity of untested ArM variants and guide further
271 screening rounds, we relied on Gaussian process (GP) regression⁴⁶. This machine learning technique
272 can capture highly non-linear relationships and has the distinct advantage of being probabilistic, which
273 means that it predicts a probability distribution rather than a point estimate, and thus provides an
274 estimate for the confidence of each prediction. This feature can not only help users assess the
275 uncertainty of individual predictions, but is ideally suited for active learning strategies. In this scenario,
276 the model's uncertainty estimates can be used to guide subsequent screening rounds towards
277 uncertain regions of sequence space with the goal of improving the model (i.e., exploration), before
278 suggesting highly active variants in later rounds (i.e., exploitation).

279 GPs are characterized by a mean and a covariance function, which is commonly referred to as kernel.
280 In our case, as we operate on the space of protein sequences, the kernel measures the similarity
281 between different ArM variants. Since the selection of a suitable kernel is of paramount importance
282 for good performance and sample efficiency (i.e., predicting accurately with little data), we performed
283 a benchmarking process and found that the non-linear Matérn kernel⁴⁶ performed best in our case
284 (see Methods).

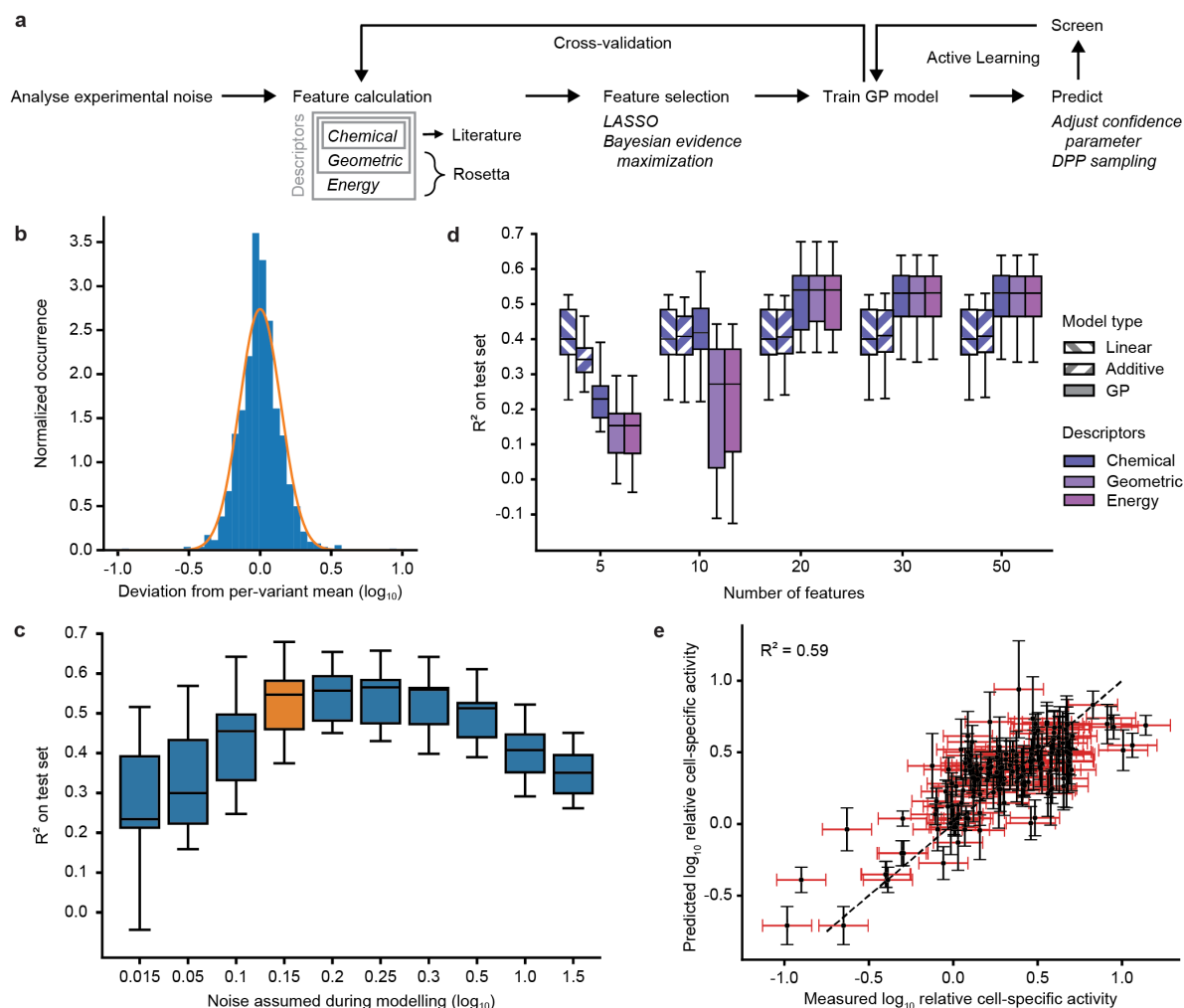
285 Moreover, our model development pipeline included steps to account for experimental noise and to
286 select suitable descriptors (Fig. 3a). Considering the inherent noise in biological experiments during
287 modelling is crucial to ensure that decisions are not influenced by random fluctuations. To distinguish
288 the genuine signal from these fluctuations, it is necessary to define a probabilistic model for data
289 generation, known as the likelihood. This step involves specifying the likelihood and its parameters,
290 which is essential for applying Bayes' theorem to calculate the posterior distribution (see Methods).
291 To elucidate the form of the likelihood, we relied on the variants appearing multiple times in the
292 screening. This revealed that the deviation of these replicates from the per-variant mean closely
293 follows a log-normal distribution, which can be viewed as a conservative estimate of the experimental
294 noise in the data (Fig. 3b). Considering the log-transformed values, this implies a Gaussian likelihood.
295 Next, we used the replicate measurements to determine a standard deviation, which is a key element
296 in defining the data likelihood. We made the simplifying assumption that the variance of the
297 measurement remains constant across the different ArM variants and repeated this analysis after each
298 round of screening. As illustrated in Fig. 3c, under- or overestimating the experimental noise leads to
299 a drastically reduced performance of the resulting model, likely due to overfitting to noise in the data.
300 In contrast, the procedure applied here results in a robust performance in the face of noisy data.

301 With regard to the descriptors that represent the ArM variants during training, we considered features
302 that reflect chemical properties of amino acids¹¹ as well as features that were extracted from Sav
303 mutant structures predicted with the Rosetta software⁴⁷. The latter included both geometric features
304 (e.g. solvent accessible surface area, number of hydrogen bonds, partial charge, dihedral angles, etc.)
305 and energy terms. Note that the geometric descriptors were compiled to be strict supersets of the
306 chemical descriptors (i.e., they also included the chemical descriptors), and similarly the energy-based
307 descriptors are strict supersets of the geometric descriptors. Given the large number of features (125
308 chemical, 682 geometric, and 161 energy features), we sought to select subsets that are parsimonious
309 while still highly predictive to ensure data efficiency and eliminate redundancy. To this end, we relied
310 on Bayesian evidence maximization (see Methods). Due to the non-linearity of the optimization
311 challenge, we first reduced the feature sets using LASSO, which performed best in a benchmarking test
312 (Supplementary Fig. 4). More precisely, we fitted a linear model and selected features with non-zero
313 coefficients for automatic relevance detection using Bayesian evidence maximization with a Gaussian

314 process. This allowed us to reduce the initial pool of features to 20-100 and speed up the evidence
315 maximization step, which required multiple optimisation restarts to ensure that an adequate
316 maximum was achieved.

317 Finally, we trained GP models using the different reduced feature sets on the available sequence-
318 activity data and evaluated model performance using 15-fold cross-validation. For comparison, we
319 included a linear and an additive, non-linear model based on chemical descriptors. The latter is
320 restricted to treating potentially non-linear effects on the activity additively and is therefore not
321 capable of modelling epistatic effects. Notably, the linear and additive models performed considerably
322 worse than the GP models (Fig. 3d), confirming that advanced methods such as GP models are required
323 to accurately capture the sequence-activity relationships in the data. Interestingly, the chemical,
324 geometric, and energy-based descriptors displayed a comparable performance, and a set of 20
325 features proved to be sufficient in all cases. The most influential features based on automatic relevance
326 detection are listed in Supplementary Table 1 (see Supplementary Fig. 5 for an analysis of their
327 influence).

328 As computationally expensive structural calculations are required to generate the geometric and
329 energy-based features and no clear benefit over models relying only on chemical descriptors was
330 observable, we chose to continue with the subset of 20 chemical features as our primary encoding
331 strategy for further modelling. The resulting model displayed a good predictive performance, with a
332 median R^2 of 0.54 based on 15-fold cross-validation (see Fig. 3e and Supplementary Fig. 6 for exemplary
333 validation splits). While leaving room for improvement, this degree of correlation has previously been
334 shown to be suitable for guiding directed evolution campaigns¹¹. Moreover, the median Spearman
335 correlation of 0.68 demonstrates that the relative ranking of variants was largely reproduced by the
336 model (Supplementary Fig. 7), which is important for confident selection of high-activity variants.



337
 338 **Fig. 3 | Development of the initial GP model.** **a**, Overview of the machine learning pipeline. Initially, the standard
 339 deviation of the activity measurements was estimated to account for experimental noise. Subsequently, three
 340 feature sets were calculated and reduced sets were obtained by applying LASSO and Bayesian evidence
 341 maximization. The resulting descriptors were then used to train GP models. Model selection and model fitting
 342 were benchmarked using cross-validation. Ultimately, the GP model can be used to navigate the sequence space
 343 in active learning cycles. **b**, Histogram of the deviation between replicates in the initial library. The distribution
 344 of residuals can be conservatively approximated by a normal distribution with a specific variance (orange). **c**,
 345 Influence of the noise estimate on the predictive performance of the resulting GP model. The value chosen based
 346 on Fig. 3b is highlighted in orange. The models used here were based on chemical descriptors with 20 features
 347 (see Fig. 3d) and were evaluated using 15-fold cross-validation. The box plots display the 25th, 50th and 75th
 348 percentile with whiskers denoting the 1.5-fold interquartile range. **d**, Influence of feature number (x-axis), model
 349 type (fill pattern), and descriptors (colour) on the performance of machine learning models analysed by 15-fold
 350 cross-validation. The box plots display the 25th, 50th and 75th percentile with whiskers denoting the 1.5-fold
 351 interquartile range. **e**, Performance of the GP model using chemical descriptors and 20 features on an exemplary
 352 cross-validation split. The measurement uncertainty (one standard deviation) is displayed in red, while the
 353 uncertainty of the model is in black. The R^2 value of this particular cross-validation split is displayed.

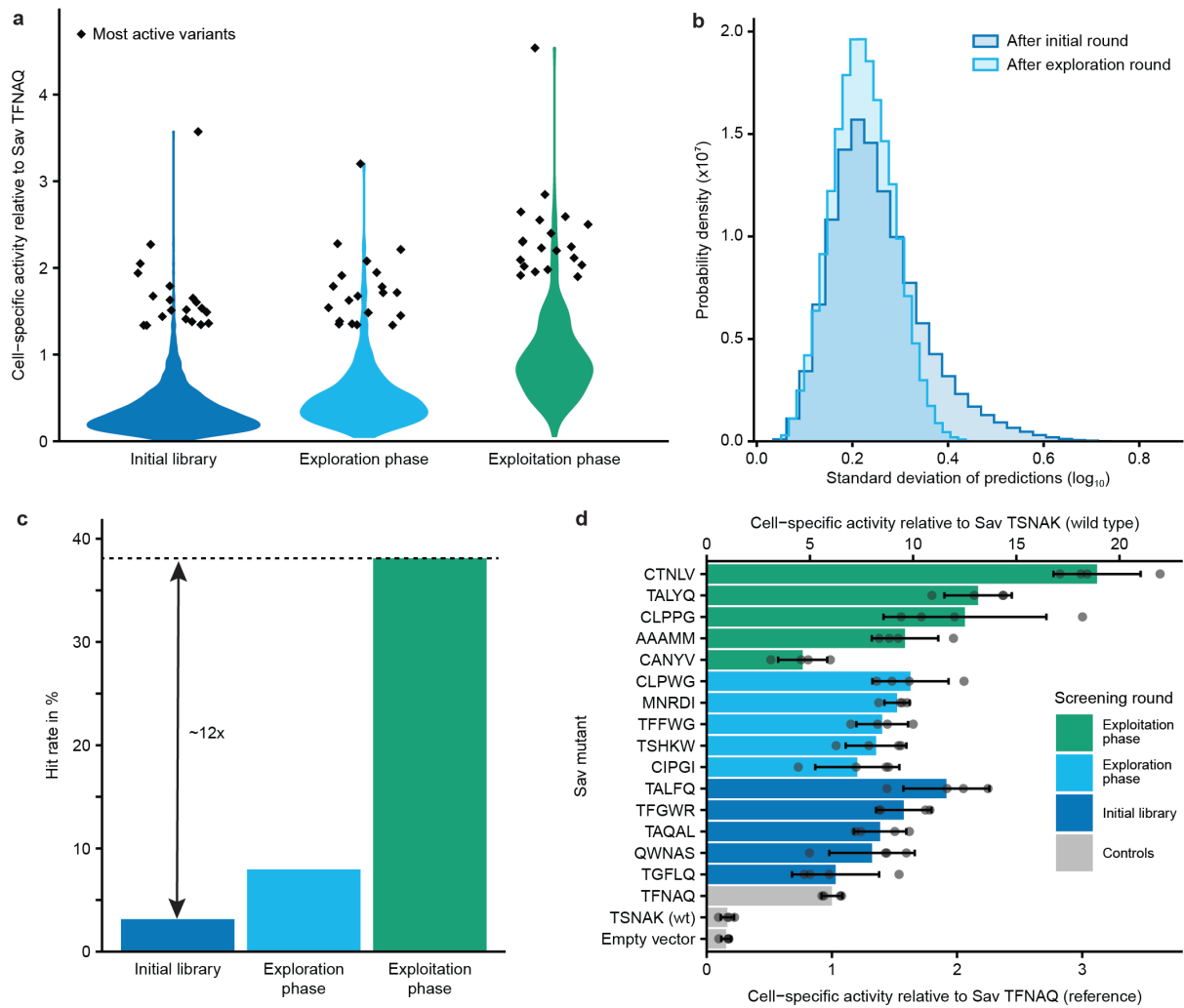
354 Model refinement by active learning

355 The aforementioned performance parameters indicate that the initial GP model can predict ArM
 356 activity with reasonable accuracy. However, due to the vast sequence space, the random sampling
 357 from this space during the generation of training data, as well as inevitable biases in experimental
 358 library construction, it is likely that this initial model will not generalize well across the entire sequence-
 359 activity landscape. Consequently, it may be “blind” for certain underexplored regions containing highly
 360 active ArMs. Therefore, we performed a second, exploratory screening round with the goal of
 361 improving the model’s accuracy and ability to generalize across the entire sequence space. To this end,

362 we designed a new library consisting of 720 variants that were primarily selected to be “informative”.
363 Specifically, we utilised the uncertainty estimates of the GP model and selected the variants with the
364 highest uncertainty in the predicted activity among all 3.2 million mutants^{48,49}.

365 We generated these variants based on a pool of oligonucleotides obtained through commercial
366 synthesis on arrays, a method that allows for the cost-efficient construction of large and targeted
367 libraries⁵⁰ and is therefore highly useful for active learning with large batch sizes. After cloning the
368 oligonucleotides into the Sav expression plasmid, we screened the resulting exploration library relying
369 on the automated pipeline in combination with NGS as described above. This exploratory round
370 yielded sequence-activity data on 465 additional variants. It should be noted that this library also
371 contained chimeric variants with amino acid combinations that were not planned in the computational
372 design, likely due to PCR-mediated recombination between variants^{51,52}. While unintended, these
373 additional variants can also be used to augment the machine learning model and were therefore
374 included for training. If desired, chimera formation can be minimized by optimizing the PCR
375 conditions^{51,52}.

376 The exploration library displayed a similar activity distribution as the initial training data (Fig. 4a), which
377 is in line with the focus on informative instead of active variants. Importantly, these new data led to a
378 decrease in the standard deviation of the predictions, most prominently for variants that had
379 previously exhibited a high uncertainty (Fig. 4b). While this observation alone is not a proof of
380 increased accuracy, it hints towards an improved representation of previously underexplored regions
381 of the sequence space, which we examined in more detail in subsequent analyses (see below).



382
 383 **Fig. 4 | ArM engineering by means of active learning.** **a**, Activity distributions in the three screening rounds
 384 displayed as violin plots. The 20 most active variants in each round are depicted as diamonds. Activity is displayed
 385 relative to the reference variant (Sav TFNAQ). **b**, Normalized histograms of the standard deviations of predictions
 386 across all 3.2 million variants after the first and second round of screening. **c**, Hit rate in the three screening
 387 rounds. Here, any variant with a higher cell-specific activity than the reference variant is considered a hit. The hit
 388 rate represents the fraction of hits amongst all variants screened in the respective round. Note that the hit rate
 389 in the initial library was calculated based on the triple and quadruple mutants, excluding the double mutants that
 390 had been tested previously⁴⁰. In the third round, chimeric variants that were not part of the computationally
 391 designed library were excluded to provide a better analysis of the models' performance. **d**, The five most active
 392 variants from each screening round were tested again in four replicates. The five-letter codes denote the amino
 393 acids in positions 111, 112, 118, 119, and 121 for the respective variants.

394 Active learning increases the efficiency of directed evolution

395 Following model refinement in the exploration round, we set out to test whether our model-guided
 396 approach can indeed aid in the discovery of active ArMs. With this goal in mind, we designed a third
 397 library of 720 variants predicted to be of high activity. Additionally, we employed an *in silico*
 398 diversification step to avoid choosing only variants with highly similar sequences. This provides a
 399 safeguard against inaccuracies in the top predictions and increases the likelihood of obtaining variants
 400 with diverse properties besides activity (e.g. thermostability, solubility, or activity under alternative
 401 conditions). To this end, we used a notion of diversity known as determinantal point processes
 402 (DPPs)^{48,49}, which use the GP kernel to determine which variants are similar to each other (see Methods
 403 and Supplementary Fig. 8a). In short, this approach treats the descriptors of the Sav variants as vectors
 404 in Euclidian space and attempts to select a set of vectors that are as orthogonal to each other as
 405 possible. We applied this process to a set of variants with the highest predicted activity to obtain a

406 subset of active and yet sequence-diverse variants. This led to a more diverse set of variants compared
407 to a simple greedy selection of the variants with the highest predicted activity as assessed by three
408 different metrics of diversity (Supplementary Fig. 8b).

409 As described for the exploration round, we obtained the designed library based on an oligonucleotide
410 pool and acquired experimental data for 349 distinct variants. Gratifyingly, this third library displayed
411 a clear shift towards higher activities compared to the first two rounds, both in terms of the average
412 as well as the top activities (Fig. 4a). We further analysed the hit rate in the screening rounds, which
413 we define here as the fraction of ArM variants with higher activity than the reference variant, which is
414 the most active variant identified in a previous study⁴⁰. While only 3 % of the initial library were hits,
415 this rate reached 38 % in the exploitation phase, amounting to an approximately 12-fold increase (Fig.
416 4c). This demonstrates that the model acquired a meaningful representation of the activity landscape
417 and can reliably predict active ArMs.

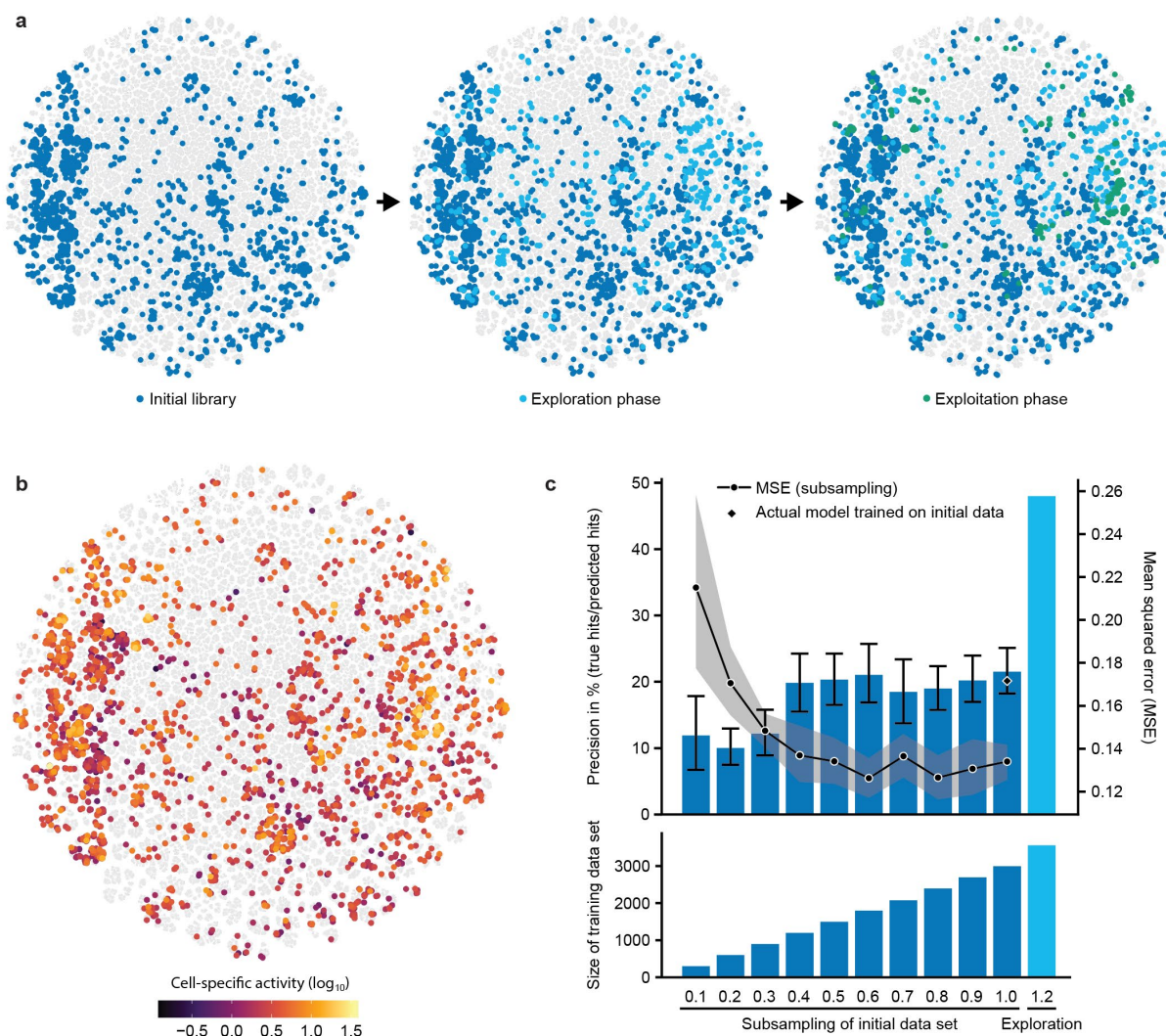
418 To confirm the results from the different screening rounds, which were performed in single
419 measurements, we tested the most promising variants from all three rounds again in four replicates
420 (Fig. 4d). This revealed that Sav 111C 112T 118N 119L 121V (abbreviated Sav CTNLV) was the most
421 active variant, reaching an 18-fold higher cell-specific hydroamination activity than the wild type (Sav
422 TSNAK) and a three-fold higher cell-specific activity than the reference variant (Sav TFNAQ). In addition,
423 we purified the most active variants from our whole-cell screening to test whether they also display
424 an increased total turnover number *in vitro*, which was the case for five of the seven variants tested
425 (Supplementary Fig. 9). As observed before⁴⁰, the ranking of the variants changed *in vitro*, which can
426 be expected due to the different reaction environments and varying expression levels in the
427 periplasmic screening.

428 Notably, the Sav CTNLV mutant does not retain the S112F K121Q mutations that were found to be
429 optimal in the previous double mutant screening⁴⁰. Likewise, all other variants evaluated in the
430 validation experiment (Fig. 4d) retain neither or only one of these two mutations. This highlights the
431 importance of epistatic effects, which can only be adequately considered through combinatorial library
432 designs and non-additive models. Strikingly, several highly active variants contain a cysteine at position
433 111, which seems counter-intuitive as cysteine has been repeatedly shown to have a pronounced
434 inhibitory effect on gold-catalysed hydroamination⁵³. However, residue 111 is pointed away from the
435 metal, presumably preventing the thiol from interfering with catalysis. Notably, the beneficial impact
436 of this mutation was not obvious from the initial data set, but became increasingly apparent in
437 subsequent rounds. This indicates that active learning can traverse the mutational space more broadly
438 than alternative methods and enable the identification of counter-intuitive effects on activity.

439 To further corroborate this hypothesis, we performed more detailed analyses to investigate whether
440 the active learning strategy with a model-guided exploration round indeed led to a better
441 representation of the available sequence space. We visualised the sequence space (using t-SNE⁵⁴ on
442 the kernel matrix, see Methods) to analyse how the tested variants are distributed across this space
443 (Fig. 5a, b). While care must be taken when interpreting such low-dimensional projections, this analysis
444 indicates that the initial library did indeed not cover the sequence space uniformly. The subsequent
445 exploration round filled in several of the “gaps” in accordance with the design goal of this phase. The
446 exploitation phase focused on a number of regions of high activity, indicating that the selection criteria
447 of high activity and sequence diversity were met. The emergence of multiple clusters of active variants
448 is compatible with the notion of a “rugged” activity landscape with many local optima. Such landscapes
449 can be challenging to navigate using classical methodologies, which frequently follow a single “uphill”
450 trajectory. In contrast, the GP model developed here acquires a holistic understanding of the entire
451 space of 3.2 million ArM variants and allows us to sample various potential optima, increasing the
452 chances of finding suitable variants.

453 Lastly, we sought to quantify the effect of the applied sampling strategy in relation to the size of the
454 training data set. A crucial question in this regard is whether the active learning strategy suggested
455 here provides a significant benefit over a comparable increase in the size of the training data set by
456 random sampling of variants. To investigate this, we trained models on different fractions of the initial
457 data set using the same model development pipeline as before. As a proxy for an experimentally
458 determined hit rate, we analysed the models' precision in identifying hits among the variants tested in
459 the exploitation phase (i.e., the percentage of true hits among variants predicted to be hits). As
460 illustrated in Fig. 5c, this analysis indicates that acquiring training data by random sampling is
461 accompanied by strong diminishing returns: Approximately 40 % of the initial data set size (equivalent
462 to ~1200 data points) is sufficient to achieve a similar performance (in terms of precision and mean
463 squared error (MSE)) as a model trained on the entire initial data set (~3000 data points). This suggests
464 that additional random screening rounds of similar size would not have led to noteworthy
465 improvements of the model. In contrast, the model-guided exploration round, which consisted of only
466 564 additional data points (an increase of less than 20 % in data volume), improved the precision in
467 identifying hits from ~20 % to 48 %. This increase is significantly beyond any improvement that can be
468 anticipated due to the mere increase in data volume, emphasizing the fact that this round was
469 substantially more informative than random sampling. This confirms the validity of the suggested
470 active learning and model-guided exploration strategies, pointing to a high potential for enhancing
471 MLDE campaigns while at the same time minimizing the experimental effort.

472



473
474 **Fig. 5 | Enhanced sequence-activity mapping through active learning.** **a**, t-SNE visualisation of the sequence
475 space. ArM variants that were tested in the three screening rounds are highlighted in different colours. To
476 generate this visualisation, all 3.2 million mutants were considered, and a uniform subsample of untested
477 variants was plotted in grey. The similarity metric used was derived from the GP model (see Methods for details).
478 **b**, t-SNE visualisation of the sequence space with colour encoding the activity of experimentally tested variants.
479 The clustering is identical to that in Fig. 5a. **c**, Precision in identifying hits and mean squared error (MSE) of
480 predictions as a function of the size of the training data set. The dark-blue bars in the upper graph indicate the
481 average precision of models that were trained on different fractions of the initial data set (screening round 1).
482 The diamond at 1.0 represents the precision of the model used to inform experiments. The light-blue bar on the
483 right represents the model refined by model-guided exploration (screening round 2). Note that the precision is
484 not identical to the experimentally determined hit rate (see Methods). The lower graph depicts the size of the
485 data sets used to train the respective models.

486 Discussion

487 MLDE is a highly promising strategy for engineering enzymes and other proteins. However, the success
488 and efficiency of such engineering campaigns hinges on the ability to generate sufficiently large and
489 informative data sets, the use of smart sampling strategies, and the choice of suitable machine learning
490 techniques that optimally leverage the resulting data.

491 Many studies on MLDE have relied on small data sets⁴⁻¹⁰ and a single training phase^{4,5,10,55,56}, which
492 may be attributed to experimental limitations. This bears the risk that the resulting models do not
493 accurately represent the sequence space, and thus are likely to leave significant potential hidden
494 within this space untapped. Here, we applied lab automation and NGS to acquire large data sets in a

495 simple and cost-efficient manner, and directed our sampling to the most informative data by means
496 of advanced active learning techniques.

497 Lab automation greatly increases the throughput of screenings and is, at the same time, highly
498 adaptable to various reactions and target proteins. In this study, we performed some experimental
499 steps manually, but a fully automated workflow could also be implemented. Similarly, the
500 computational pipeline is largely automated, and thus it is conceivable to conduct protein engineering
501 with minimal human intervention. Importantly, recent developments such as academic biofoundries
502 and cloud labs are making such approaches more widely accessible^{57,58}.

503 The NGS strategy employed here enables the sequencing of thousands of protein variants for the cost
504 of a small Illumina run and PCR reagents. The former is available for a few hundred dollars (e.g. MiSeq
505 Nano, yielding approx. 1 million reads) and will likely continue to get cheaper. If combined with other
506 samples and run on an instrument with a large capacity, the prorated costs may even be in the range
507 of a few dollars. Regarding the PCR reagents, primer synthesis costs are low as only 20 primers are
508 required to address all 96 positions in a well plate. Similarly, the use of two plate barcodes means that
509 12 primers for the second PCR are sufficient to distinguish 36 well plates. Overall, this means that
510 sequencing is possible at a cost of less than one cent per variant.

511 Combined, automation and NGS are ideally suited to generate large data sets for MLDE. At the same
512 time, it is also crucial to design information-dense libraries to maximize the efficiency of experimental
513 screening rounds. In the initial round, we achieved this by optimizing the mutational load in the library,
514 which is a straightforward and broadly applicable strategy. Alternatively, approaches such as zero-shot
515 methods relying on $\Delta\Delta G$ calculations¹³ can be applied as well. In subsequent rounds, library design can
516 be guided by the machine learning model. While it may seem attractive to apply an exploitation-
517 focused strategy to quickly identify active variants, we hypothesized that a model-guided exploration
518 round could substantially improve the predictive performance and thus increase the chances of
519 identifying suitable variants in large and rugged sequence spaces in a subsequent round. Indeed, we
520 observed that the exploration round improved the model's ability to identify active variants far beyond
521 what would be expected due to the increase in data volume alone. This demonstrates that active
522 learning is a highly effective and efficient strategy for developing accurate models of sequence-activity
523 landscapes. Moreover, the separation into exploration and exploitation phases provides a transparent
524 and practical solution to the exploration-exploitation dilemma, as it allows for a clear and plannable
525 resource allocation. In addition, our study introduces DPP sampling as a strategy for diversifying the
526 selection of active variants, which increases the robustness of MLDE to possible model inaccuracies
527 and may be beneficial with regard to secondary properties beyond activity.

528 In terms of the machine learning approach, this study corroborates that Gaussian process regression is
529 an attractive choice for MLDE, particularly when strong epistatic effects are present in the sequence-
530 activity landscape. Moreover, it is well-suited for active learning strategies, as the uncertainty
531 quantification is computationally simple, which constitutes an advantage over alternative methods
532 such as deep learning. Our results demonstrate that simple and computationally efficient descriptors
533 are sufficient for non-trivial improvements to engineering campaigns, which is in line with other
534 literature on the subject^{59,60}. Nonetheless, it might be possible to further boost the predictive
535 performance, for example by employing improved structure prediction algorithms or descriptors from
536 modern protein language models^{61,62}. Lastly, our results highlight that accurately accounting for
537 experimental noise is crucial during model development, an aspect that has frequently been
538 neglected⁶³.

539 The application of these strategies to the engineering of ArMs for gold-catalysed hydroamination led
540 to the identification of a variant with 18-fold higher cell-specific activity than the wild type. Compared

541 to our previous screening of double mutants⁴⁰, extending the search space to five positions led to a
542 three-fold improvement. Further rounds of active learning could potentially lead to the discovery of
543 even more active variants. Moreover, the methods developed here could be used to target additional
544 positions. However, it should be noted that this ArM is likely a challenging engineering target due to
545 the relatively exposed location of the cofactor in Sav. Therefore, applying this engineering strategy to
546 alternative scaffolds with a more shielded active site might enable larger improvements⁶⁴. Currently,
547 artificial (metallo)enzymes are typically limited by their rather modest activity. Thus, the field could
548 profit greatly from advanced machine learning-guided engineering strategies, as demonstrated here.
549 Similarly, the active learning approach described here could be applied to tailor natural enzymes for
550 industrial applications, or to engineer other proteins such as antibodies, biosensors, or transporters.

551 **Materials and Methods**

552 **Chemicals and reagents**

553 **(Biot-NHC)Au1** was synthesized as previously described⁴⁰. All other chemicals were obtained from
554 Sigma-Aldrich. Primers were synthesized by Sigma-Aldrich, and enzymes for molecular cloning were
555 obtained from New England Biolabs.

556 **Plasmids**

557 All plasmids were based on a previously described expression plasmid that contains a T7-tagged Sav
558 gene with an N-terminal OmpA signal peptide for export to the periplasm under control of the T7
559 promoter in a pET30b vector³³. This plasmid is available from Addgene (#138589). A version of this
560 plasmid encoding the Sav S112F K121Q mutant was used as the starting point for library generation.

561 **Cloning of Sav libraries**

562 **Site-saturation mutagenesis at 20 positions:** To individually randomize 20 positions in Sav, the plasmid
563 encoding Sav S112F K121Q was amplified in two parts in order to create two overlapping fragments
564 for each position, with mutations being introduced by an NDT codon in one of the primer overhangs.
565 The PCRs were conducted using the primer pairs SSM_X_NDT_fwd and kanR_rev, and kanR_fwd and
566 SSM_X_rev (X denotes the position to be randomized, see Supplementary Table 2). PCRs were carried
567 out using Q5 High-Fidelity DNA Polymerase (New England Biolabs). Following DpnI digest and PCR
568 purification, the corresponding fragments were assembled by Gibson assembly and transformed into
569 *E. coli* BL21-Gold(DE3). Three clones per position were sequenced by Sanger sequencing to verify
570 correct assembly and diversity at the desired position.

571 **Double, triple, quadruple, and quintuple mutant libraries:** To generate sets of double, triple,
572 quadruple, and quintuple mutants, the plasmid encoding Sav S112F K121Q was amplified in two parts.
573 One part included the Sav positions 111 and 112, and the other part included positions 118, 119, and
574 121. To generate fragments with variable but defined numbers of mutations, the primers from
575 Supplementary Table 3 were used in several PCR reactions according to Supplementary Table 4.
576 Following DpnI digest and PCR purification, the fragments were assembled in several Gibson assembly
577 reactions as summarized in Supplementary Table 5. The reactions were then transformed separately
578 into chemocompetent *E. coli* Top10. Plasmids were isolated from the transformants and transformed
579 into the expression strain BL21-Gold(DE3). When picking colonies for screening, the theoretical
580 diversity of the individual sub-libraries (Supplementary Table 5) was taken into account in order to
581 obtain balanced sets of double, triple, quadruple and quintuple mutants.

582 **Active learning libraries:** To create libraries of specific Sav variants that were suggested by the machine
583 learning models, oligo pools were ordered from Twist Bioscience. These oligos were used as primers
584 that bind immediately downstream of position 121 in Sav. The 5'-overhang contained the five mutation
585 sites with the desired changes as well as a constant region for Gibson assembly (see Supplementary
586 Table 6). For the first library of ML-designed variants, insert and backbone were generated according
587 to Supplementary Table 7. For the second library, the PCRs were run according to Supplementary Table
588 8. Following DpnI digest and PCR purification, the fragments were assembled by Gibson assembly and
589 transformed into chemocompetent *E. coli* Top10. Plasmids were isolated from the transformants and
590 transformed into the expression strain BL21-Gold(DE3).

591 **Sav expression in 96-well plates**

592 96-deep well plates were filled with 500 μ L of LB (+ 50 mg L⁻¹ kanamycin) per well. Cultures were
593 inoculated from glycerol stocks and grown overnight at 37 °C and 300 revolutions per minute (rpm) in
594 a Kuhner LT-X shaker (50-mm shaking diameter). 20 μ L per culture was used to inoculate expression
595 cultures in 1 mL of LB with kanamycin. These cultures were grown at 37 °C and 300 rpm for 1.5 h. At

596 this point, the plates were placed at room temperature for 20 min, and subsequently, Sav expression
597 was induced by addition of isopropyl- β -D-thiogalactopyranoside (IPTG, final concentration 50 μ M).
598 Expression was carried out at 20 °C and 300 rpm for an additional 16 h.

599 **Whole-cell screening**

600 Following the expression of Sav mutants in deep-well plates, the OD₆₀₀ of the cultures was determined
601 in a plate reader using 50 μ L of samples diluted with an equal volume of PBS. Afterwards, the plates
602 were centrifuged (3,220 rcf, 15 °C, 10 min), the supernatant was discarded and the pellets were
603 resuspended in 400 μ L of incubation buffer (10 μ M **(Biot-NHC)Au1** in 50 mM MES, 0.9 % NaCl, 10 mM
604 diamide, pH 6.1). Cells were incubated with the cofactor for 1 h at 15 °C and 300 rpm. Afterwards,
605 plates were centrifuged (2,000 rcf, 15 °C, 10 min), the supernatant was removed and the pellets were
606 resuspended in 500 μ L of washing buffer (50 mM MES, 0.9 % NaCl, 10 mM diamide, pH 6.1). Following
607 another centrifugation step, cell pellets were resuspended in 200 μ L of reaction buffer (5 mM 2-
608 ethynylaniline in 50 mM MES, 0.9 % NaCl, 10 mM diamide, pH 6.1). Reactions were performed at 37 °C
609 and 300 rpm for 20 h before determining the product concentration. To account for differences in cell
610 density and plate-to-plate variations, the product concentrations were divided by the OD₆₀₀ of the
611 culture and normalized to the mean of the cell-specific product concentrations measured for the Sav
612 S112F K121Q controls in the respective plate.

613 **Kovac's assay**

614 Indole was quantified using the photometric Kovac's assay (adapted from Piñero-Fernandez et al.⁶⁵).
615 For measurements in culture supernatant, plates were centrifuged (3,220 rcf, 20 °C, 10 min) and 110 μ L
616 supernatant was mixed with 165 μ L of Kovac's reagent (50 g L⁻¹ 4-(dimethylamino)benzaldehyde,
617 710 g L⁻¹ isoamyl alcohol, 240 g L⁻¹ hydrochloric acid) in a separate plate. After 5 min of incubation,
618 these plates were centrifuged (3,220 rcf, 20 °C, 10 min). Subsequently, 75 μ L of the upper phase was
619 transferred to a new transparent plate and the absorbance at 540 nm was measured in a plate reader
620 (Tecan Infinite M1000 PRO).

621 **Lab automation**

622 Colony picking, reaction setup and product quantification were implemented using an automation
623 platform featuring two Tecan EVO 200 (Tecan Group AG) robotic platforms coupled to each other. Both
624 platforms were controlled using the EVOware standard software (Tecan Group AG). Colony picking was
625 performed using the integrated Pickolo system (SciRobotics). For shaking, incubation, and
626 resuspension of cultures, the platform was equipped with a Kuhner ES-X shaking platform (Adolf
627 Kühner AG) running at 300 rpm at 50-mm shaking radius. The shaking platform was surrounded by a
628 custom-made box made of aluminum plastic composite panels (Tecan Group AG). The temperature
629 inside the box was maintained at 15 °C using an "Icecube" (Life imaging services) heater/cooler device.
630 Centrifugation of the samples was performed using the integrated Rotanta 46 RSC Robotic centrifuge
631 (Hettich AG). All buffer exchanges during sample preparation were performed using the integrated
632 liquid-displacement pipetting system equipped with eight 2500 μ L dilutors and fixed stainless steel
633 needles. Absorbance measurements were performed using a Tecan Infinite M200 PRO plate reader.
634 The automation method files are available upon request.

635 **Barcoding of mutants**

636 Following colony picking, cultures were grown overnight at 37 °C and 200 rpm in 96-deep well plates.
637 On the following day, 150 μ L per culture was transferred to a 96-well PCR plate. The plates were sealed
638 and placed in a thermal cycler for 5 min at 95 °C to lyse the bacteria. Subsequently, the plates were
639 centrifuged (3200 rcf, 5 min) and 0.5 μ L of the supernatant was used as template for the first PCR. This
640 PCR step was done in 96-well plates, with each well containing a distinct combination of barcoded
641 primers (see Supplementary Table 9). 30 cycles were performed with 30 s denaturation at 98 °C, 20 s

642 annealing at 71 °C and 30 s elongation by Pfu DNA polymerase at 72 °C. The products from each plate
643 were pooled, run on a 2.5 % (w/v) agarose gel at 100 V for 2 h and purified using a gel extraction kit
644 (Sigma-Aldrich). The products were then used as templates for a second PCR with distinct
645 combinations of barcoded primers (Supplementary Table 10) to generate a plate-specific labelling. The
646 primer overhangs also contained the adapters required for Illumina sequencing. 30 additional cycles
647 were performed, consisting of 30 s denaturation at 98 °C, 20 s annealing at 63 °C and 30 s elongation
648 by Q5 High-Fidelity DNA Polymerase (New England Biolabs) at 72 °C. Ultimately, all products were
649 pooled, run on a 2.5 % (w/v) agarose gel at 100 V for 2 h, and purified using a gel extraction kit.

650 **Illumina sequencing**

651 NGS was performed by the Genomics Facility Basel using an Illumina MiSeq platform and a Reagent
652 Kit v2 Nano (150 cycles, PE 110/40) using ~20 % genomic PhiX library as spike-in to increase sequence
653 diversity.

654 **NGS data analysis**

655 NGS data were analyzed using a custom R script. Forward and reverse reads retrieved from fastq files
656 were paired and target fragments were selected based on several constant regions
657 (GTCACACGTAGCATGTGG, GAGACCTTGTGTGCATGG, GGCCTCGGTGGTGCC, no mismatches). Mutation
658 sites as well as barcodes were extracted based on their distance to these regions. All reads with a Q-
659 score < 30 at the mutation sites were discarded, as well as those for which a barcode did not match
660 any of the expected sequences. The codons at the mutation sites were translated to amino acids in
661 order to identify the Sav variants and the barcodes were used to identify the plate and well for each
662 read. For each plate, the entries were then grouped by variant and only the combinations of variant
663 and well with the highest number of reads was kept. This eliminates combinations of variants and
664 barcodes that result from chimera formation during the second PCR step. Subsequently, variants that
665 accounted for less than 80 % of reads for a given barcode combination were discarded in order to
666 eliminate cases where more than one variant had been present in a well.

667 **Sav expression for purification**

668 A single colony of *E. coli* BL21-Gold(DE3) harbouring a plasmid for periplasmic expression of the desired
669 Sav variant was used to inoculate a starter culture (4 mL of LB with 50 mg L⁻¹ kanamycin), which was
670 grown overnight at 37 °C and 200 rpm. On the following day, 100 mL of LB with kanamycin in a 500 mL
671 flask was inoculated to an OD₆₀₀ of 0.01. The culture was grown at 37 °C and 200 rpm until it reached
672 an OD₆₀₀ of 0.5. At this point, the flask was placed at room temperature for 20 min and 50 µM IPTG
673 (final concentration) was added to induce Sav expression. Expression was performed at 20 °C and
674 200 rpm overnight, and cells were harvested by centrifugation (3,220 rcf, 4 °C, 15 min). Pellets were
675 stored at -20 °C until purification.

676 **Sav purification**

677 Cell pellets were resuspended in 10 mL of lysis buffer (50 mM tris, 150 mM NaCl, 1 g L⁻¹ lysozyme,
678 pH 7.4). After 30 min of incubation at room temperature, cell suspensions were subjected to three
679 freeze-thaw cycles. Subsequently, nucleic acids were digested by addition of 10 µL of DNaseI (2000
680 units/mL, New England Biolabs) and CaCl₂ to a final concentration of 10 mM, followed by incubation at
681 37 °C for 45 min. After centrifugation, the supernatant was transferred to a new tube and mixed with
682 40 mL of binding buffer (50 mM ammonium bicarbonate, 500 mM NaCl, pH 11). Pierce iminobiotin
683 agarose (Thermo Fisher Scientific) was equilibrated in falcon tubes and used to pack a PD-10 column
684 up to a bed height of approximately 1 cm. The lysate was loaded onto the column relying on gravity
685 flow. Subsequently, the column was washed twice with 10 mL binding buffer. Ultimately, Sav was
686 eluted using 10 mL of elution buffer (50 mM ammonium acetate, 500 mM NaCl, pH 4). Amicon Ultra

687 filters (10 kDa molecular weight cut-off) were then used to concentrate the samples and exchange the
688 buffer against the reaction buffer (50 mM MES, 0.9 % NaCl, pH 6.1).

689 **Quantification of biotin-binding sites**

690 The concentration of Sav biotin-binding sites was determined using a modified version of the assay
691 described by Kada et al.⁶⁶, which relies on the quenching of the fluorescence of a biotinylated
692 fluorophore upon binding to Sav. Specifically, 190 μL of the binding site buffer (1 μM biotin-4-
693 fluorescein, 0.1 g L^{-1} bovine serum albumin in PBS) was mixed with 10 μL of purified Sav. After
694 incubation at room temperature for 90 min, the fluorescence intensity was measured (excitation at
695 485 nm, emission at 525 nm), and a calibration curve produced with lyophilized Sav was used to
696 calculate the concentration of Sav biotin-binding sites.

697 ***In vitro* catalysis**

698 *In vitro* reactions were performed with 2.5 μM purified Sav (tetrameric; corresponding to 10 μM biotin-
699 binding sites), 5 μM (**Biot-NHC**)**Au1** and 5 mM 2-ethynylaniline in MES buffer (50 mM MES, 0.9 % NaCl,
700 pH 6.1). The reactions were performed in a volume of 200 μL in glass vials and were incubated at 37 °C
701 and 200 rpm for 20 h. Subsequently, the indole concentration was determined using the Kovac's assay.

702 **Machine learning**

703 All machine learning methods were implemented in Python using scikit-learn⁶⁷, Pytorch⁶⁸, Biotite⁶⁹,
704 pyRosetta⁷⁰ and SciPy⁷¹.

705 **Calculation of descriptors:** In this work, we encoded the Sav mutants by three different classes of
706 descriptors: chemical descriptors, geometric descriptors, and energy-based descriptors. To obtain the
707 chemical descriptors, we utilized amino-acid descriptors from four different sources: Z-scores¹⁶,
708 VSHE¹⁷, Barley score¹⁸, and PCscores⁵⁵. All of these are based on physical amino-acid properties (see
709 Supplementary Table 11) and principal component analysis (PCA) was used to construct a reduced
710 representation. Here, we concatenated these features, resulting in 25 values per amino-acid position.
711 As we considered quintuple mutants, each Sav variant is thus described by 125 features.

712 The geometric and energy-based features were created using the Rosetta software. First, we
713 calculated the approximate dimeric structure of each mutant with a fixed seed using the *mutate*
714 function with the default distance for post-mutational changes. The mutations were performed in the
715 order of the five sites in the primary protein sequence (111, 112, 118, 119, 121). We calculated all 3.2
716 million approximate Sav dimer structures. Next, we used the package Biotite to calculate charge,
717 distance to the centre of mass, and radii of each amino-acid residue. Additionally, we calculated the
718 solvent accessible surface area of each residue, the number of hydrogen bonds per residue, and the
719 dihedral angles. A summary of the features can be found in Supplementary Table 12. We discarded
720 variables that did not vary across the 3.2 million structures, leaving us with 682 features. The energy-
721 based features were calculated in the same manner as the geometric features using the approximate
722 structure of the variant and correspond to the ref2015 set of 31 features per mutant from the Rosetta
723 suite (see Supplementary Table 13). A common pre-processing step applied to all features involved
724 subtracting the mean of each descriptor across the 3.2 million variants and scaling by the absolute
725 value of the maximum value of that descriptor. This process ensured that the descriptors fell within
726 the range [-1,1] and that their average value was zero.

727 **Likelihood elucidation:** The first step of any data analysis is to understand its randomness and
728 generation process. In our case, the likelihood specified the experimental error introduced by
729 biological variability, the measurement procedure, etc. In other words, we assumed that our
730 measurements were corrupted by additive noise under log transformation. To justify this hypothesis,
731 we analysed the distribution of the differences between replicates from their mean value. As a normal

732 distribution appeared to be a good and conservative approximation for these data, we used a Gaussian
733 likelihood with a standard deviation determined from the aforementioned distribution. In the first
734 round, this value was determined to be 0.15, rounded to two decimal points in the log-transformed
735 cell-specific activity. We repeated the same procedure for the subsequent screening rounds to account
736 for variability between experiments. The standard deviations determined for the second and third
737 round were 0.20 and 0.12, respectively.

738 **Model section:** For further analysis and Gaussian process fitting, we did not use the full set of features
739 due to the complexity of the initial fitting procedure, which involves optimizing the marginal
740 likelihood⁷². To simplify this process, we preprocessed the initial set of descriptors using one of three
741 straightforward machine learning models: LASSO, elastic net, and random forests. We evaluated the
742 effectiveness of this procedure through cross-validation on the entire feature space. In all cases, we
743 utilized the scikit-learn implementation of these methods. Both the LASSO and elastic net methods
744 employed an adaptive selection of the regularization parameter, which involved an additional layer of
745 cross-validation within the training split. For random forests, we used a configuration of 500 trees with
746 a maximum depth of 15 and a minimum split size of 5. After training, we selected k descriptors with
747 either the largest coefficients or the highest feature importance for further analysis. We varied k across
748 20, 40, 60, 80, and 100. This range was chosen as the maximum set of descriptors that we believed
749 would allow the Gaussian process library to reliably optimize the marginal likelihood.

750 **Gaussian process:** The functional relationship between the Sav sequence and ArM activity was
751 modelled using Gaussian processes (GPs). This Bayesian method is versatile in capturing a wide range
752 of structures, and is defined by its mean and covariance function, also known as the kernel. In our case,
753 we found that kernels of the following form performed best among selected statistical models with
754 calibrated uncertainty:

755
$$k(p, \tilde{p}) = \kappa(\text{poly}(d(p, \tilde{p}))) \exp(-d(p, \tilde{p})^2) \text{ where } d_\gamma(p, \tilde{p}) \propto \left(\sqrt{\sum_{j=1}^m 1/\gamma_j^2 (\Phi_j(p) - \Phi_j(\tilde{p}))^2} \right).$$

756 This kernel is known as Matérn kernel with regularity parameter $\eta=5/2$ and is commonly used to model
757 twice differentiable smooth response surfaces⁴⁶. The letters p and \tilde{p} denote different protein variants
758 of which we want to calculate similarity. The function Φ corresponds to the feature representation of
759 the protein p . In this work, this is a function that maps the protein sequence or structure to a fixed
760 length vector. The parameters γ_i are usually referred to as length scales and are used for automatic
761 relevance detection⁷³. They guide the importance of a certain variable, i.e., if γ is very large, this part
762 of the descriptor vector Φ has less impact if changed than a coordinate Φ_j with larger γ_j . The length
763 scales can be selected based on Bayesian evidence maximization, which is a well-tested methodology
764 to select length scales that most likely explain the activity data⁷². The parameter κ was selected using
765 the expected maximal achievable improvement of the protein, in this case $\kappa = 3$, meaning that the
766 maximum achievable improvement is 1000-fold over the wild-type variant (due to modelling \log_{10}).

767 **Bayesian evidence maximization:** Hyperparameters, specifically the length scales of the Matérn
768 kernel, were optimized for each of the chosen features using the maximization of evidence, a common
769 Bayesian approach⁴⁶. As before, we denote length scales γ . By maximization of evidence, we mean

770
$$\gamma^* = \text{argmax}_\gamma P(D|\gamma) \text{ and } P(D|\gamma) = \int P(D|f, \gamma)P(f, \gamma)df,$$

771 where $P(f|\gamma)$ is the Gaussian process prior parametrized by length scales, and $P(D|f, \gamma)$ is the
772 Gaussian likelihood as specified in the prior section on likelihood elucidation. The integration in the
773 prior formula represents marginalization of the Gaussian process, and strictly speaking integration

774 requires certain mathematical regularity conditions, which we omit here. Upon finding the right length
775 scales from the initial data, these were fixed, and the posterior $P(f|\gamma, D)$ was calculated after each
776 experimentation round without changing them. To implement the Bayesian posterior calculation, we
777 used a custom implementation in Python.

778 **Active Learning**

779 To employ active learning, we used a technique similar to the upper confidence bound method as
780 described by Srinivas et al.²⁷, or greedy information maximization. In the exploration round, we
781 generated predictions using the GP model based on chemical descriptors with 20 features. To select
782 informative variants, the confidence parameter was set to infinity. In addition, we allocated a smaller
783 part of the experimental budget to variants predicted to be active to validate the model. The latter
784 budget was split equally into three categories: A conservative set representing the Sav mutants for
785 which the mean prediction minus two standard deviations was highest, as well as balanced and
786 optimistic predictions chosen based on the mean and the mean plus two standard deviations as
787 ranking mechanisms, respectively. See Supplementary Table 14 for an overview of the budget
788 allocation in the exploration round. We obtained additional data points through a small random
789 mutagenesis as well as chimeric variants, which were not part of the designed library.

790 In the exploitation round, we aimed to select active and diverse ArM variants. To this end, we trained
791 three GP models on the new data set (including the exploration round). The three models employed
792 different descriptors (chemical descriptors with 20 features, geometric and energy –based descriptors
793 with 50 features) to possibly obtain more diverse predictions. We split the experimental budget equally
794 among the three models. Further, we split the experimental budget per model into conservative and
795 balanced predictions (see above). The experimental budget allocation can be found in Supplementary
796 Table 15. The confidence parameter was set to 2 for the exploitation round. Additionally, a diversifying
797 principle based on determinantal point processes⁴⁸, a mathematical model of diversity, was employed
798 to choose a diverse subset of variants, following the principles described by Nava et al.⁴⁹ (see below).
799 Upon retrieval of the above budget, we performed a validation step. As part of it, we augmented the
800 chemical descriptor model with the new data and proposed 30 additional Sav variants to test for
801 potential improvements. These were selected to be conservative or balanced (10 variants each), and
802 10 variants were selected to be the best predicted according to the balanced prediction metric.

803 **DPP sampling**

804 When selecting Sav variants for experimental testing, it is advisable that these are diverse, especially
805 in the context of the exploitation round. For example, if we were to identify the best x candidates using
806 the machine learning pipeline, it is very likely that all these top x candidates are highly similar to each
807 other for small x . If the model happens to be incorrect with regard to the top predictions, this will lead
808 to failure to identify any active mutants. A more principled approach is to pick a diverse subset.
809 Namely, select a set of promising mutants, and then further select a subset of these which is diverse.
810 This ensures robustness to potential misspecification errors. The model of diversity we employed here
811 is the inverse of the similarity model we used to train the GP regressor, namely the kernel. We
812 measured the diversity of the selected subset by the determinant of the kernel matrix. This is a
813 common approach in the machine learning literature⁴⁸, as it has an intuitive interpretation where the
814 determinant between two vectors is proportional to the volume that the two vectors span (see
815 Supplementary Fig. 7a). The more orthogonal (dissimilar) these two vectors are, the larger the volume.
816 A natural extension to non-parametric models such as GP models is to use the kernel matrix instead of
817 the inner product between vectors. Finding a subset of maximum determinant is an NP-hard
818 problem⁷⁴. Hence, often a probabilistic method is employed to find the subsets^{49,75}.

819 Suppose that the probability of sampling a set is proportional to the value of the determinant for this
820 set. This probabilistic object is known as determinantal point process (DPP)⁴⁸ and can be sampled very
821 efficiently. In order to diversify our top-x batches, we select a top y number of candidates, where y is
822 bigger than x, from which we choose a diverse set of size x using DPP sampling. The value of y = 500
823 was chosen arbitrarily for our experiments. The value of x depends on the available experimental
824 budget in each round. The explorative round does not require diversification as the goal to select
825 informative Sav variants already leads to diversity. In fact, it is related to the greedy search for a set
826 with the largest determinant⁷⁵.

827 In order to compare the diversity of the measurements, we use the isometry score, which is a ratio
828 determinant and trace of a kernel matrix defined via the batch of sequences. The score equates to the
829 normalized ratio of trace and determinant.

830

$$I(K) = \left(\frac{\det(K)^{\frac{1}{n}}}{\frac{1}{n} \text{trace}(K)} \right).$$

831 The score is valued between 0 and 1, where 1 is achieved once K forms essentially a diagonal matrix.
832 If this is the case, this means the implicit features (defined via the kernel) are orthogonal to each other.
833 On the other hand, 0 indicates that the implicit features defined via the kernel are very closely aligned
834 to each other. Of course, this score depends on the kernel metric we use. The DPP method practically
835 maximizes this metric under the models' kernel in expectation.

836 **Clustering of ArM variants**

837 The clustering shown in Fig. 5 was created using the t-SNE (t-distributed stochastic neighbour
838 embedding)⁵⁴ clustering methodology. For this analysis, we used the kernel matrix of the chemical
839 descriptor model. This model is based on a Gaussian process with ARD (automatic relevance
840 determination) kernel length scales. The t-SNE algorithm clusters the data based on a similarity metric
841 that includes exponentiated negative Euclidean distances. This is very similar to our machine learning
842 model, with the exception that instead of a pure exponential, we use the Matérn kernel. However, this
843 should qualitatively lead to similar results. Hence, to generate the clustering, we took the chemical
844 descriptors, scaled them with appropriate length scales, and used the scikit-learn implementation of
845 the t-SNE algorithm to generate the clusters. We tested several values of complexity, and the plotted
846 clusterings correspond to a value of 150, as it appeared to generate the most structured results.

847 **Subsampling analysis**

848 To analyse the effect of data set size on the predictive ability of the model, we created 20 random
849 subsamples of the original data set for each subsampling fraction (0.1 - 1 in intervals of 0.1). We then
850 applied the previously described machine learning pipeline, starting with the feature selection. To
851 analyse the performance of the models, we used them to predict the activity of all ArM variants that
852 were tested in the exploitation round, and calculated the mean squared error of the predictions as
853 well as the precision in predicting hits (i.e., ArM variants with a higher activity than the reference
854 variant). Precision is defined as the percentage of true hits among predicted hits. To investigate the
855 effect of the exploration round, we calculated the precision of a model that was trained on all data
856 from the initial library and the exploration round. In the latter case, the precision is different from the
857 experimentally determined hit rate as not all experimentally tested variants were predicted to be hits
858 by the model used here.

859 **Data and code availability**

860 The data and code will be made available upon publication of the manuscript.

861 **Acknowledgments**

862 The authors thank Fadri Christoffel for synthesizing the gold cofactor. This work was created as part of
863 the NCCR Catalysis and the NCCR Molecular Systems Engineering, both National Centres of
864 Competence in Research funded by the Swiss National Science Foundation. R.T. acknowledges a grant
865 from the Naito Foundation.

866 **Author contributions**

867 T.V. and M.J. conceived the project. T.V. and G.S. developed the automated screening methods. T.V.
868 and C.S. performed experiments. T.V. analysed screening results and NGS data. M.M. developed,
869 applied, and analysed the machine learning pipeline. R.T. developed initial computational models. M.J.,
870 S.P. and T.R.W. supervised experimental work. A.K. supervised machine learning aspects. T.V., M.M.
871 and M.J. wrote the manuscript with input from all authors.

872 **Competing interests**

873 The authors declare no competing interests.

874 References

- 875 1. Yang, K. K., Wu, Z. & Arnold, F. H. Machine-learning-guided directed evolution for protein
876 engineering. *Nat. Methods* **16**, 687–694 (2019).
- 877 2. Mazurenko, S., Prokop, Z. & Damborsky, J. Machine learning in enzyme engineering. *ACS Catal.*
878 **10**, 1210–1223 (2020).
- 879 3. Freschlin, C. R., Fahlberg, S. A. & Romero, P. A. Machine learning to navigate fitness landscapes
880 for protein engineering. *Curr. Opin. Biotechnol.* **75**, 102713 (2022).
- 881 4. Saito, Y. *et al.* Machine-learning-guided mutagenesis for directed evolution of fluorescent
882 proteins. *ACS Synth Biol* **7**, 2014–2022 (2018).
- 883 5. Biswas, S., Khimulya, G., Alley, E. C., Esvelt, K. M. & Church, G. M. Low-N protein engineering
884 with data-efficient deep learning. *Nat. Methods* **18**, 389–396 (2021).
- 885 6. Bedbrook, C. N., Yang, K. K., Rice, A. J., Gradinaru, V. & Arnold, F. H. Machine learning to design
886 integral membrane channelrhodopsins for efficient eukaryotic expression and plasma
887 membrane localization. *PLoS Comput. Biol.* **13**, e1005786 (2017).
- 888 7. Bedbrook, C. N. *et al.* Machine learning-guided channelrhodopsin engineering enables
889 minimally invasive optogenetics. *Nat. Methods* **16**, 1176–1184 (2019).
- 890 8. Romero, P. A., Krause, A. & Arnold, F. H. Navigating the protein fitness landscape with Gaussian
891 processes. *Proc. Natl. Acad. Sci.* **110**, E193–E201 (2013).
- 892 9. Greenhalgh, J. C., Fahlberg, S. A., Pflieger, B. F. & Romero, P. A. Machine learning-guided acyl-
893 ACP reductase engineering for improved in vivo fatty alcohol production. *Nat. Commun.* **12**, 1–
894 10 (2021).
- 895 10. Li, G. *et al.* Machine learning enables selection of epistatic enzyme mutants for stability against
896 unfolding and detrimental aggregation. *ChemBioChem* **22**, 904–914 (2021).
- 897 11. Wu, Z., Kan, S. B. J., Lewis, R. D., Wittmann, B. J. & Arnold, F. H. Machine learning-assisted
898 directed protein evolution with combinatorial libraries. *Proc. Natl. Acad. Sci.* **116**, 8852–8858
899 (2019).
- 900 12. Gelman, S., Fahlberg, S. A., Heinzelman, P., Romero, P. A. & Gitter, A. Neural networks to learn
901 protein sequence–function relationships from deep mutational scanning data. *Proc. Natl. Acad.*
902 *Sci.* **118**, e2104878118 (2021).
- 903 13. Wittmann, B. J., Yue, Y. & Arnold, F. H. Informed training set design enables efficient machine
904 learning-assisted directed protein evolution. *Cell Syst.* **12**, 1026-1045.e7 (2021).
- 905 14. Atwal, G. S. & Kinney, J. B. Learning quantitative sequence–function relationships from
906 massively parallel experiments. *J. Stat. Phys.* **162**, 1203–1243 (2016).
- 907 15. Höllerer, S., Desczyk, C., Muro, R. F. & Jeschek, M. From sequence to function and back – High-
908 throughput sequence-function mapping in synthetic biology. *Curr. Opin. Syst. Biol.* **37**, 100499
909 (2024).
- 910 16. Sandberg, M., Eriksson, L., Jonsson, J., Sjöström, M. & Wold, S. New chemical descriptors
911 relevant for the design of biologically active peptides. A multivariate characterization of 87
912 amino acids. *J. Med. Chem.* **41**, 2481–2491 (1998).

- 913 17. Mei, H., Liao, Z. H., Zhou, Y. & Li, S. Z. A new set of amino acid descriptors and its application in
914 peptide QSARs. *Biopolymers* **80**, 775–786 (2005).
- 915 18. Barley, M. H., Turner, N. J. & Goodacre, R. Improved descriptors for the quantitative structure-
916 activity relationship modeling of peptides and proteins. *J. Chem. Inf. Model.* **58**, 234–243 (2018).
- 917 19. Somnath, V. R., Bunne, C. & Krause, A. Multi-scale representation learning on proteins. In
918 *Advances in Neural Information Processing Systems* (2021).
- 919 20. Gainza, P. *et al.* Deciphering interaction fingerprints from protein molecular surfaces using
920 geometric deep learning. *Nat. Methods* **17**, 184–192 (2020).
- 921 21. Alley, E. C., Khimulya, G., Biswas, S., AlQuraishi, M. & Church, G. M. Unified rational protein
922 engineering with sequence-based deep representation learning. *Nat. Methods* **16**, 1315–1322
923 (2019).
- 924 22. Yang, K. K., Wu, Z., Bedbrook, C. N. & Arnold, F. H. Learned protein embeddings for machine
925 learning. *Bioinformatics* **34**, 2642–2648 (2018).
- 926 23. Fox, R. J. *et al.* Improving catalytic function by ProSAR-driven enzyme evolution. *Nat.*
927 *Biotechnol.* **25**, 338–344 (2007).
- 928 24. Cadet, X. F., Gelly, J. C., van Noord, A., Cadet, F. & Acevedo-Rocha, C. G. Learning strategies in
929 protein directed evolution. In *Directed Evolution: Methods and Protocols* (Springer, 2022).
- 930 25. Saito, Y. *et al.* Machine-learning-guided library design cycle for directed evolution of enzymes:
931 The effects of training data composition on sequence space exploration. *ACS Catal.* **11**, 14615–
932 14624 (2021).
- 933 26. Büchler, J. *et al.* Algorithm-aided engineering of aliphatic halogenase WelO5* for the
934 asymmetric late-stage functionalization of soraphens. *Nat. Commun.* **13**, 371 (2022).
- 935 27. Srinivas, N., Krause, A., Kakade, S. M. & Seeger, M. W. Information-theoretic regret bounds for
936 Gaussian process optimization in the bandit setting. *IEEE Trans. Inf. Theory* **58**, 3250–3265
937 (2012).
- 938 28. Hie, B. L. & Yang, K. K. Adaptive machine learning for protein engineering. *Curr. Opin. Struct.*
939 *Biol.* **72**, 145–152 (2022).
- 940 29. Vornholt, T. & Jeschek, M. The quest for xenobiotic enzymes: From new enzymes for chemistry
941 to a novel chemistry of life. *ChemBioChem* **21**, 2241–2249 (2020).
- 942 30. Kan, S. B. J., Lewis, R. D., Chen, K. & Arnold, F. H. Directed evolution of cytochrome c for carbon-
943 silicon bond formation: Bringing silicon to life. *Science* **354**, 1048–1051 (2016).
- 944 31. Zhou, Z. & Roelfes, G. Synergistic catalysis in an artificial enzyme by simultaneous action of two
945 abiological catalytic sites. *Nat. Catal.* **3**, 289–294 (2020).
- 946 32. Yang, H. *et al.* Evolving artificial metalloenzymes via random mutagenesis. *Nat. Chem.* **10**, 318–
947 324 (2018).
- 948 33. Jeschek, M. *et al.* Directed evolution of artificial metalloenzymes for in vivo metathesis. *Nature*
949 **537**, 661–665 (2016).
- 950 34. Key, H. M., Dydio, P., Clark, D. S. & Hartwig, J. F. Abiological catalysis by artificial haem proteins
951 containing noble metals in place of iron. *Nature* **534**, 534–537 (2016).

- 952 35. Song, W. J. & Tezcan, F. A. A designed supramolecular protein assembly with in vivo enzymatic
953 activity. *Science* **346**, 1525–1528 (2014).
- 954 36. Bordeaux, M., Tyagi, V. & Fasan, R. Highly diastereoselective and enantioselective olefin
955 cyclopropanation using engineered myoglobin-based catalysts. *Angew. Chem. Int. Ed.* **54**, 1744–
956 1748 (2015).
- 957 37. Kan, S. B. J., Huang, X., Gumulya, Y., Chen, K. & Arnold, F. H. Genetically programmed chiral
958 organoborane synthesis. *Nature* **552**, 132–136 (2017).
- 959 38. Dydio, P. *et al.* An artificial metalloenzyme with the kinetics of native enzymes. *Science* **354**,
960 102–106 (2016).
- 961 39. Studer, S. *et al.* Evolution of a highly active and enantiospecific metalloenzyme from short
962 peptides. *Science* **362**, 1285–1288 (2018).
- 963 40. Vornholt, T. *et al.* Systematic engineering of artificial metalloenzymes for new-to-nature
964 reactions. *Sci. Adv.* **7**, eabe4208 (2021).
- 965 41. Currin, A., Swainston, N., Day, P. J. & Kell, D. B. Synthetic biology for the directed evolution of
966 protein biocatalysts: navigating sequence space intelligently. *Chem. Soc. Rev.* **44**, 1172–1239
967 (2015).
- 968 42. Reetz, M. T., Kahakeaw, D. & Lohmer, R. Addressing the numbers problem in directed evolution.
969 *ChemBioChem* **9**, 1797–1804 (2008).
- 970 43. Chen, Y. *et al.* Barcoded sequencing workflow for high throughput digitization of hybridoma
971 antibody variable domain sequences. *J. Immunol. Methods* **455**, 88–94 (2018).
- 972 44. Glenn, T. C. *et al.* Adapterama II: Universal amplicon sequencing on Illumina platforms
973 (TaggiMatrix). *PeerJ* **7**, e7786 (2019).
- 974 45. Wittmann, B. J., Johnston, K. E., Almhjell, P. J. & Arnold, F. H. evSeq: Cost-effective amplicon
975 sequencing of every variant in a protein library. *ACS Synth. Biol.* **11**, 1313–1324 (2022).
- 976 46. Rasmussen, C. E. & Williams, C. K. I. Gaussian Processes for Machine Learning. (The MIT Press,
977 2005).
- 978 47. Kellogg, E. H., Leaver-Fay, A. & Baker, D. Role of conformational sampling in computing
979 mutation-induced changes in protein structure and stability. *Proteins* **79**, 830–838 (2011).
- 980 48. Kulesza, A. & Taskar, B. Determinantal point processes for machine learning. *Found. Trends*
981 *Mach. Learn.* **5**, 123–286 (2012).
- 982 49. Nava, E., Mutný, M. & Krause, A. Diversified sampling for batched Bayesian optimization with
983 determinantal point processes. *Proceedings of the 25th International Conference on Artificial*
984 *Intelligence and Statistics* (2022).
- 985 50. Kuiper, B. P., Prins, R. C. & Billerbeck, S. Oligo pools as an affordable source of synthetic DNA
986 for cost-effective library construction in protein- and metabolic pathway engineering.
987 *ChemBioChem* **23**, e202100507 (2022).
- 988 51. Omelina, E. S., Ivankin, A. V., Letiagina, A. E. & Pindyurin, A. V. Optimized PCR conditions
989 minimizing the formation of chimeric DNA molecules from MPRA plasmid libraries. *BMC*
990 *Genomics* **20**, 1–10 (2019).

- 991 52. Thompson, J. R., Marcelino, L. A. & Polz, M. F. Heteroduplexes in mixed-template
992 amplifications: Formation, consequence and elimination by ‘reconditioning PCR’. *Nucleic Acids*
993 *Res.* **30**, 2083–2088 (2002).
- 994 53. Burgener, S., Dačević, B., Zhang, X. & Ward, T. R. Binding interactions and inhibition mechanisms
995 of gold complexes in thiamine diphosphate-dependent enzymes. *Biochemistry* **62**, 3303–3311
996 (2023).
- 997 54. Van Der Maaten, L. & Hinton, G. Visualizing data using t-SNE. *J. Mach. Learn. Res.* **9**, 2579–2605
998 (2008).
- 999 55. Xu, Y. *et al.* Deep dive into machine learning models for protein engineering. *J. Chem. Inf. Model.*
1000 **60**, 2773–2790 (2020).
- 1001 56. Ma, E. J. *et al.* Machine-directed evolution of an imine reductase for activity and
1002 stereoselectivity. *ACS Catal.* **11**, 12433–12445 (2021).
- 1003 57. Carbonell, P., Radivojevic, T. & García Martín, H. Opportunities at the intersection of synthetic
1004 biology, machine learning, and automation. *ACS Synth. Biol.* **8**, 1474–1477 (2019).
- 1005 58. Rapp, J. T., Bremer, B. J. & Romero, P. A. Self-driving laboratories to autonomously navigate the
1006 protein fitness landscape. *Nat. Chem. Eng.* **1**, 97–107 (2024).
- 1007 59. Hsu, C., Nisonoff, H., Fannjiang, C. & Listgarten, J. Learning protein fitness models from
1008 evolutionary and assay-labeled data. *Nat. Biotechnol.* **40**, 1114–1122 (2022).
- 1009 60. Shanehsazzadeh, A., Belanger, D. & Dohan, D. Is transfer learning necessary for protein
1010 landscape prediction? *ArXiv* (2020) doi:10.48550/arXiv.2011.03443.
- 1011 61. Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a language
1012 model. *Science* **379**, 1123–1130 (2023).
- 1013 62. Brandes, N., Goldman, G., Wang, C. H., Ye, C. J. & Ntranos, V. Genome-wide prediction of
1014 disease variant effects with a deep protein language model. *Nat. Genet.* **55**, 1512–1522 (2023).
- 1015 63. Sundar, V., Tu, B., Guan, L. & Esvelt, K. FLIGHTED: Inferring fitness landscapes from noisy high-
1016 throughput experimental data. *NeurIPS* (2023).
- 1017 64. Christoffel, F. *et al.* Design and evolution of chimeric streptavidin for protein-enabled dual gold
1018 catalysis. *Nat. Catal.* **4**, 643–653 (2021).
- 1019 65. Piñero-Fernandez, S., Chimerel, C., Keyser, U. F. & Summers, D. K. Indole transport across
1020 *Escherichia coli* membranes. *J. Bacteriol.* **193**, 1793–1798 (2011).
- 1021 66. Kada, G., Kaiser, K., Falk, H. & Gruber, H. J. Rapid estimation of avidin and streptavidin by
1022 fluorescence quenching or fluorescence polarization. *Biochim. Biophys. Acta* **1427**, 44–8 (1999).
- 1023 67. Pedregosa, F. *et al.* Scikit-Learn: Machine learning in Python. *J. Mach. Learn. Res.* **12**, 2825–
1024 2830 (2011).
- 1025 68. Paszke, A. *et al.* PyTorch: An imperative style, high-performance deep learning library. In
1026 *Advances in Neural Information Processing Systems* 32, 8024–8035 (2019).
- 1027 69. Kunzmann, P. & Hamacher, K. Biotite: a unifying open source computational biology framework
1028 in Python. *BMC Bioinf.* **19**, 1–8 (2018).

- 1029 70. Chaudhury, S., Lyskov, S. & Gray, J. J. PyRosetta: a script-based interface for implementing
1030 molecular modeling algorithms using Rosetta. *Bioinformatics* **26**, 689–691 (2010).
- 1031 71. Virtanen, P. *et al.* SciPy 1.0: Fundamental algorithms for scientific computing in Python. *Nat.*
1032 *Methods* **17**, 261–272 (2020).
- 1033 72. Lotfi, S., Izmailov, P., Benton, G., Goldblum, M. & Wilson, A. G. Bayesian model selection, the
1034 marginal likelihood, and generalization. In *Proceedings of the 39th International Conference on*
1035 *Machine Learning* (2022).
- 1036 73. Neal, R. M. Bayesian Learning for Neural Networks. (Springer Science & Business Media, 2012).
- 1037 74. Nikolov, A. & Singh, M. Maximizing determinants under partition constraints. In *Proceedings of*
1038 *the forty-eighth annual ACM symposium on Theory of Computing*, 192–201 (2016).
- 1039 75. Kathuria, T., Deshpande, A. & Kohli, P. Batched Gaussian process bandit optimization via
1040 determinantal point processes. In *NIPS'16: Proceedings of the 30th International Conference on*
1041 *Neural Information Processing Systems*, 4213–4221 (2016).
- 1042 76. Baker, E. N. & Hubbard, R. E. Hydrogen bonding in globular proteins. *Prog. Biophys. Mol. Biol.*
1043 **44**, 97–179 (1984).
- 1044 77. Park, H. *et al.* Simultaneous optimization of biomolecular energy functions on features from
1045 small molecules and macromolecules. *J. Chem. Theory Comput.* **12**, 6201–6212 (2016).
- 1046