

1 **Deep Directed Evolution of Solid Binding Peptides**
2 **for Quantitative Big-data Generation**

3
4 Deniz T. Yucesoy,^{1,2,&} Siddharth S. Rath,^{1,2,3,&} Jacob L. Rodriguez,^{1,2}
5 Jonathan Francis-Landau,^{1,4} Oliver Nakano-Baker,^{1,2}
6 ⁵ Mehmet Sarikaya^{1,2,4,5,6,*}

7
8 ¹ Genetically Engineered Materials Science and Engineering Center, University of
9 Washington, Seattle, WA, 98195

10 ² Materials Science and Engineering, University of Washington, Seattle, WA, 98195

11 ³ Molecular Engineering and Sciences, University of Washington, Seattle, WA, 98195

12
13 ⁴Mathematics, University of Washington, Seattle, WA, 98195

14
15 ⁵ Chemical Engineering, University of Washington, Seattle, WA, 98195

16
17 ⁶ Oral Health Sciences

18
19 & These authors have contributed equally

20
21 * Corresponding Author: sarikaya@uw.edu

22
23 **Dataset availability:** Sequence Read Archive, NCBI, under Bioproject Accession
24 Number PRJNA598245 and the corresponding BioSample Accession Number
25 SAMN13702265 with unique ID: 3702265. There are 24 SRA accession numbers, one
26 for each Wash and technical replica. The SRA accession numbers start at
27 SRX7483794 and end at SRX7483817, and can be accessed at the following hyperlink:
28 https://www.ncbi.nlm.nih.gov/sra?LinkName=biosample_sra&from_uid=13702265
29 The link contains the raw DNA sequences numbering 288,000,000 identified from the directed
30 evolution experiments through next generation sequencing; see below.

31 **Code Availability at GEMSEC GitHub:** <https://github.com/Sarikaya-Lab-GEMSEC>, which
32 contains all of the software developed, cleaned and curated data relevant to this work.

33
34
35
36
37
38
39

40 **Abstract**

41 Proteins have evolved over millions of years to mediate and carry-out biological
42 processes efficiently. Directed evolution approaches have been used to genetically
43 engineer proteins with desirable functions such as catalysis, mineralization, and target-
44 specific binding. Next-generation sequencing technology offers the capability to discover
45 a massive combinatorial sequence space that is costly to sample experimentally through
46 traditional approaches. Since the permutation space of protein sequence is virtually
47 infinite, and evolution dynamics are poorly understood, experimental verifications have
48 been limited. Recently, machine-learning approaches have been introduced to guide the
49 evolution process that facilitates a deeper and denser search of the sequence-space.
50 Despite these developments, however, frequently used high-fidelity models depend on
51 massive amounts of properly labeled quality data, which so far has been largely lacking
52 in the literature. Here, we provide a preliminary high-throughput peptide-selection protocol
53 with functional scoring to enhance the quality of the data. Solid binding dodecapeptides
54 have been selected against molybdenum disulfide substrate, a two-dimensional
55 atomically thick semiconductor solid. The survival rate of the phage-clones, upon
56 successively stringent washes, quantifies the binding affinity of the peptides onto the solid
57 material. The method suggested here provides a fast generation of preliminary data-pool
58 with ~2 million unique peptides with 12 amino-acids per sequence by avoiding
59 amplification. Our results demonstrate the importance of data-cleaning and proper
60 conditioning of massive datasets in guiding experiments iteratively. The established
61 extensive groundwork here provides unique opportunities to further iterate and modify the
62 technique to suit a wide variety of needs and generate various peptide and protein
63 datasets. Prospective statistical models developed on the datasets to efficiently explore
64 the sequence-function space will guide towards the intelligent design of proteins and
65 peptides through deep directed evolution. Technological applications of the future based
66 on the peptide-single layer solid based bio/nano soft interfaces, such as biosensors,
67 bioelectronics, and logic devices, is expected to benefit from the solid binding peptide
68 dataset alone. Furthermore, protocols described herein will also benefit efforts in medical
69 applications, such as vaccine development, that could significantly accelerate a global
70 response to future pandemics.

71

72 **Introduction**

73 Directed Evolution libraries such as phage display, cell-surface display, mRNA
74 display, and yeast display, are powerful tools for the identification of peptides and
75 proteins, including enzymes and antibodies, with an affinity for a specific target such as
76 antigens, drugs, organic molecules, and inorganic materials.¹⁻⁶ Over the years, the
77 authors and others have successfully applied the tools, particularly M13 phage display
78 and bacterial cell surface display (FLITRX), to study peptide-solid interactions for a myriad
79 of bio-nanotechnological and biomedical applications.² The link between phenotype and
80 genotype of organisms is the common feature in all combinatorial display techniques
81 where the randomized (variant) peptide sequences are displayed as fusion partners with
82 different surface proteins.^{2,6-8} The authors are one of the pioneers of adapting DE
83 techniques to select peptides with affinities to metals, ceramics, semiconductors, and
84 minerals with about 5000+ Solid Binding Peptides (SBP) specific to 25+ different
85 materials.^{2,9-11} SBPs that are successful as fundamental building blocks as molecular

86 linkers, erectors, and assemblers in bio-nanotechnology implementations, are termed
87 Genetically Engineered Peptides for Inorganics (GEPI). Using a unique representation of
88 data and conventional bioinformatics tools, the authors also discovered tiny synthesizers
89 from selected peptides that catalyze solid materials synthesis, e.g., gold, silica, and
90 hydroxyapatite, from ionic precursors.^{11,12} The phage display selection procedure for
91 biomacromolecules is well-established. But the technique can achieve even further
92 improvement in selectivity by integrating simple modifications into the biopanning protocol
93 (e.g., counter-selection step, material specificity testing, etc.) to isolate SBPs with high
94 affinity and specificity to the desired material.^{2,7,13,14} Detailed procedures, therefore, are
95 developed for a particular inorganic material in the powder, thin-film, or in single crystal
96 forms and are demonstrated in numerous publications.^{9,12,15-17}

97 The selected number of peptides has conventionally been small in the literature,
98 up to a few tens of SBPs.^{9,15} The authors ensure that in any SBP-selection, the number
99 of GEPI identified for a given material (either used as a single crystal or a powdered form)
100 are at least 35 (with an average of 50, max 96 peptides). Despite well-established phage
101 display selection procedures and subsequent improvements towards increasing the
102 winner peptides' potential, the methods are still far from leveraging the huge combinatorial
103 potential since the diversity of the libraries is in the order of 10^9 variants.^{2,7,15,18,19} The
104 drawback is primarily due to the limited scalability of clone selection techniques and
105 characterization. Moreover, in low-throughput Sanger sequencing workflows traditionally
106 used, the number of peptides isolated in a reasonable time frame after several rounds of
107 the biopanning process is limited to 10-100 clones. Compared to the vast diversity of the
108 naive library ($\sim 10^9$ unique sequences), the conventional peptide-selection approaches
109 are disappointingly low-throughput (10-100 peptides) and provide a minimal perspective
110 (<0.00001%) of the complete variant (sequence) space.

111 Such a small sample size is, therefore, not only prone to bias from nonspecific,
112 preferentially amplifying false-positive hits but also leads to omitting a large number of
113 promising candidates.^{2,7,15,18} Traditionally, peptides enriched after several rounds of
114 selection are identified by DNA sequencing of the inserts of a limited number (tens to
115 hundreds) of clones. Depending on the sequence diversity remaining in the library after
116 selection, the analysis of a manually chosen limited number of clones does not
117 necessarily result in discovering the most promising candidates. Moreover, phage display
118 screenings are notorious for identifying false-positive hits', e.g., parasitic sequences".^{14,}
119 ^{20,21} Such sequences emerge for two crucial reasons: (a) Binding to materials used during
120 the selection other than the desired substrate (such as plastics or albumin), and (b)
121 Propagation advantages.⁴ A well-known example in the latter category is the greatly
122 accelerated propagation of phages displaying the HAIYPRH peptide in the Ph.D.-7TM
123 library due to a mutation in the Shine-Dalgarno box of the phage protein gllp in this clone.⁵
124 This peptide has been identified in at least 13 independent biopanning experiments.⁴
125 Several web-based tools aid in identifying potential false positives, viz., PepBank can be
126 used to search for peptides already published in other experiments,²² while SAROTUP
127 searches for peptides binding to unintended materials.²³

128 For the accelerated design of soft interfaces, device development, and
129 deployment, leveraging statistical inference and Machine Learning (ML) tools is
130 necessary. Developing ML models is critical not only for a practical exploration of the
131 sequence space but also to make the experimental procedure more efficient and targeted

132 for a variety of bio-nanotechnological applications. To generate large SBP/GEPI data sets
133 to develop an ML algorithm to predictively design peptides, the authors incorporated Next-
134 Generation Sequencing (NGS) tools and screened millions of peptide sequences in one
135 shot. Such a method allows a more comprehensive look at the phage display library's
136 sequence-space, enabling a more practical and higher resolution characterization of the
137 library.²⁴⁻²⁶ However, the lack of a high-throughput fluorescent-microscopy (FM) or
138 spectroscopy-based²⁷ end-point binding characterization method prevents the
139 characterization of binding affinities of SBPs, which is necessary to identify GEPI and for
140 labeling of the massive training data in ML algorithms. In the absence of conventional
141 binding-affinity information, here we report a high-throughput peptide-selection protocol
142 with functional scoring to characterize the peptide affinity to the surface as a function of
143 its count number (survival probability) and ML-based data analysis models to validate the
144 scoring.^{28,29}

145 Machine learning aided directed evolution is more effective at designing SBPs than
146 wet-lab guided DE alone.³⁰⁻³⁵ ML models are paramount in generating and analyzing
147 sequence libraries more targeted towards a specific function, such as the device behavior
148 in bio/nano interfaces. An ML integrated DE, and NGS platform necessitates generating
149 a custom phage library from the predicted sequences that can be screened in a high-
150 throughput manner to validate and benchmark the predictions. The implications of such
151 an integrated experimental and mathematical platform are varied but suffice to mention
152 that it is imperative to meet the demands of modern scientific and application-focused
153 research at the bio/nano interfaces. Identifying and analyzing millions of SBPs in a single
154 shot will decipher underlying mechanisms of evolution towards specific bio/nano
155 interaction. It will also develop superior peptide prediction platforms to design novel
156 peptides with known functions towards various technological and medical
157 implementations (Fig 1).

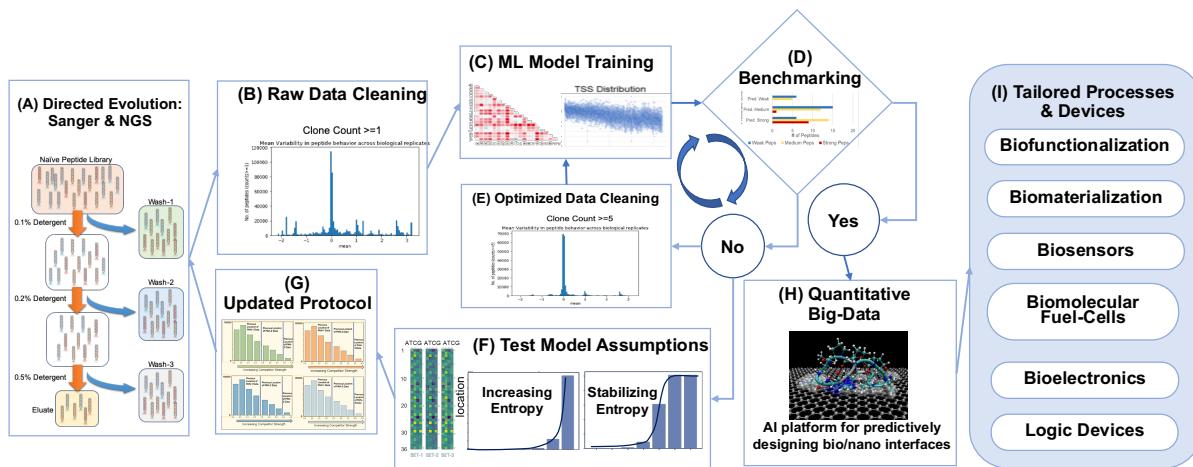


Fig 1. Overall schematic of the iterative deep directed evolution process for smart peptide selection. (A) through (H) enumerate the steps in the procedure. Conventional DE is used, with amplification for creating benchmark datasets while NGS is used for big-data generation with later iterations guided by ML/AI models. (I) Deep directed evolution enables faster design of functional bio/nano materials.

158 **Materials & Methods**

159

160 **Deep Directed Evolution for Big-Data Generation:**

161

162 **Combinatorial Mutagenesis and Biopanning:** The 12-mer Phage display (PhD)
163 library with an estimated diversity of 1.7×10^9 different clones of *M13* bacteriophage is
164 used to select peptide sequences fused to the minor coat protein (pIII) against MoS₂
165 flakes (~300 mesh, 99%, 10129, Alfa Aesar). Before the screening process, MoS₂ flakes
166 are cleaned by sonication sequentially in isopropyl alcohol and de-ionized water, then
167 dried under vacuum. For high-throughput isolation and screening of MoS₂ binding clones,
168 5 mg of MoS₂ flakes are dispersed in 1 mL of potassium phosphate/sodium carbonate
169 buffer (PC, 55 mM KH₂PO₄, 45 mM Na₂CO₃, and 200 mM NaCl, pH 7.4), containing
170 0.02% Tween 20 detergent (Merck, Whitehouse Station, NJ, USA). 10 μ L of the PhD
171 library (1011 pfu, New England Biolabs) is added onto dispersed flakes and incubated for
172 3 hours on a rotator at room temperature and then washed twice before overnight
173 incubation. After incubation, to remove the nonspecifically- or weakly-bound phages, the
174 MoS₂ flakes are washed with PC buffer with increasing detergent concentrations as
175 follows; 0.1% (v/v; Wash Round 1), 0.2% (Wash Round 2), 0.5% (Wash Round 3), at pH:
176 7.4. The remaining 'strong' bound phages are then eluted from the surface in a stepwise
177 manner by applying an elution buffer consisting of 0.2 M Glycine-HCl pH 2.2 (Sigma
178 Aldrich, St. Louis, MO) for 15 min. The washed off, and eluted phages are then transferred
179 to a fresh tube and neutralized. They are labeled as Wash-1, Wash-2, Wash-3, and
180 Eluate. Each phage pool is then purified via PEG/NaCl precipitation and resuspended in
181 de-ionized water. The overall selection procedure is performed in three biological
182 replicates and labeled as Set-1, Set-2, and Set-3 and two technical replicates that are
183 sequenced separately but subsequently combined during data analysis.

184 **DNA Isolation and Next-Generation Sequencing:** DNA amplicons are prepared
185 as previously described with slight modifications.³⁶ Single-stranded DNA (ssDNA) is
186 isolated from wash and eluate phage pools using QIAprep Spin M13 Kit (QIAGEN). The
187 sequencing library is prepared by amplifying the 36 bp peptide-coding variable region. Q5
188 polymerase (New England Biolabs) is used to amplify the target region with forward and
189 reverse primer sequences given below:

190 Forward: CCGCGTGATTACGAGTCGCAAGCTGATAAACCGATACAATTAAAG

191 Reverse: GGGTTAGCAAGTGGCAGCCTACGTTAGTAATGAATTCTGTATGGG.

192 Illumina sequencing adapters containing the p5 and p7 index sequences are
193 attached using a second PCR. The purified PCR products from each phage pool are
194 loaded in duplicates (technical replicates) on the sequencer plate and sequenced on the
195 Illumina NextSeq platform. The Next-Generation Sequencing (NGS) platform yielded
196 288M DNA sequences in total. About 1.3M and 1.1M individual peptides were obtained
197 (replicate-1) in wash-1 and eluate pools, respectively, where ~433124 of the sequences
198 existed in both pools. After combining with the other replicates, upon subsequent
199 translation into amino acid sequences, more than 2 million unique peptides were identified
200 with varying copy numbers that survived on MoS₂ with successively stringent washes.
201 The survival probabilities are then computed as described in later sections to provide a
202 relative label for each unique sequence.

203 A conventional Directed Evolution procedure, outlined in the supplementary
204 document (S1), is used for benchmark data generation. The overall process is visualized
205 in Fig 2. We choose survival probability as a label for peptides obtained from the high
206 throughput deep directed evolution experiment but use spectrophotometry or
207 Fluorescence Microscopy assays to label benchmark peptides. Readers can also find the
208 characterization procedure for benchmark peptides in the supplementary documents (S1
209 and S1 Fig). The discrepancy opens an avenue for transfer learning and metric learning
210 for prospective statistical models. Survival probability with subsequent washes is not yet
211 an established industry standard for binding affinity measurement, and as such,
212 researchers should not label benchmark data with such a metric. However, future ML or
213 AI models could establish survival probability as a reasonable measure for affinity
214 prediction, thereby removing low-throughput characterization bottlenecks and opening up
215 avenues for even more high throughput data collection for a myriad of applications.
216 However, data-scientists must clean massive data appropriately before using it in ML/AI
217 models, so we describe our current data cleaning process in the next section.

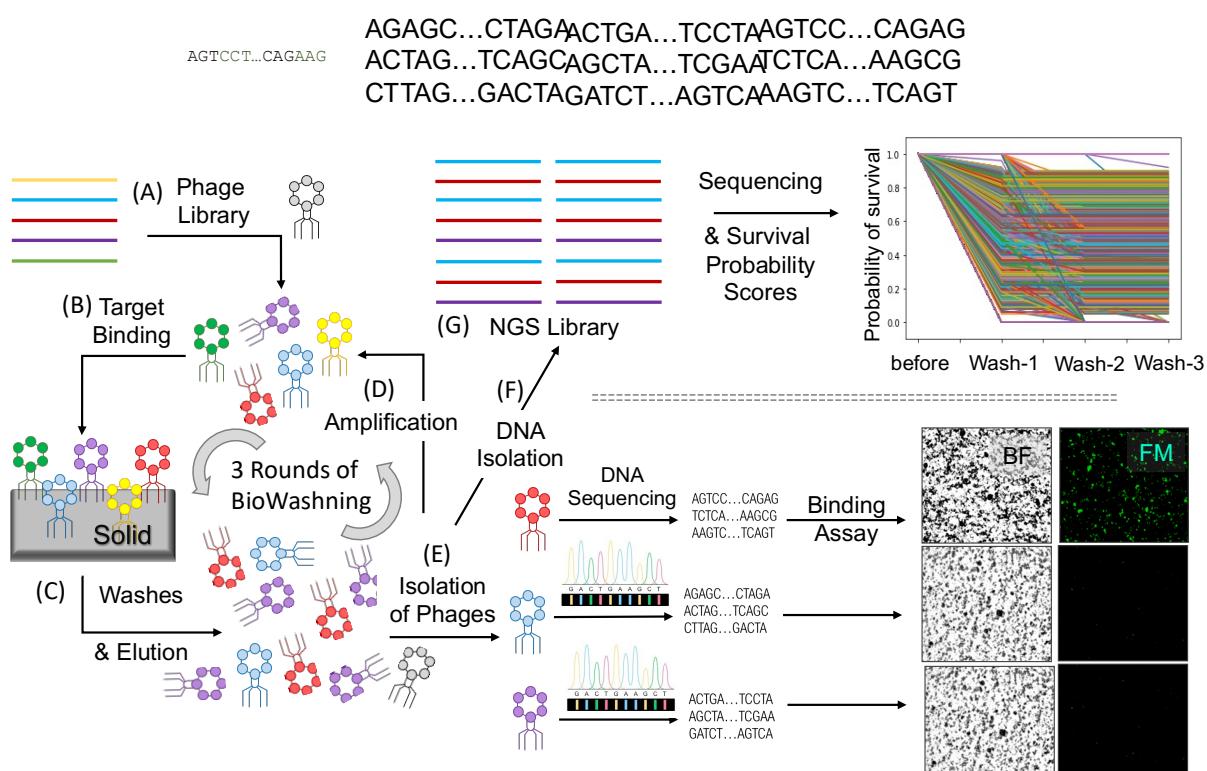


Fig 2. Detail Schematic of the Deep Directed Evolution process. (A) Variant library is created by cloning random DNA fragments into M13 phage and (B) Exposed onto target solid material which is followed by (C) Series of washing and (D) Elution/amplification steps. For benchmark set, this is followed by (E), (F) isolation and sequencing of the phages (before binding affinity characterization) from the colonies while for deep selection, it is followed by (G) next generation sequencing.

218 **Preliminary Data Cleaning:**

219
220 Traditional peptide selection and characterization provide a quantitative analysis
221 of the peptide binding as relative surface coverage or phage-retention, obtained from
222 fluorescence microscopy or spectrophotometric analysis. The NGS data, however, does
223 not provide such information explicitly.^{37,38} It generates a massive amount of sequence
224 data and their copy number in the collected samples, prior and posterior to each round of
225 wash. To quickly characterize the sequences selected through multiple washing steps,
226 one needs to process the raw sequences and their abundances to establish a correlation
227 between the survival probability of each peptide through the successive washing steps
228 and their binding affinity to the surface. Cleaning of the dataset obtained would result in
229 more robustness in the calculation of survival rates for each sequence and impart more
230 clarity as to which mathematical models should be applied to the data. In the following
231 paragraph, we describe the cleaning steps in detail, including consolidation of DNA
232 datasets, translation into peptides, before obtaining and applying thresholds for filtering
233 the dataset from a raw to a cleaned set, ready for ML model application.

234 We convert the raw data obtained as FASTQ files from the NGS experiment to
235 TSV format before the DNA sequences' computational translation into peptide
236 sequences. We then get the raw, but assembled DNA reads where the sequence is
237 labeled with its abundance (termed variously throughout the document as population,
238 count, or copy number) in the sample. We label each sample by its biological replicate
239 (Set-1, 2, or 3), the washing step (wash-1, wash-2, wash-3, or eluate), and its technical
240 replicate (a or b). Overall there are about 288M non-unique sequences with significant
241 overlaps between the replicates.

242 In the next step, we combine the technical replicates and sum the populations of
243 overlapping sequences, followed by concatenating the different datasets for each wash.
244 As a result, we obtain consolidated DNA datasets for each biological replicate. The first
245 column lists all the sequences, and the following four columns list their populations
246 (combined technical replicate populations) in the various washes. We end up with ~55M
247 unique sequences, with individual counts in washes and replicates transformed by a
248 logarithm of base 2. The logarithm transform is crucial because every DNA sequence with
249 n copies ends up with at least $2n$ copies in a usual PCR amplification process. After that,
250 we convert the DNA sequences into Amino Acid sequences. If any two unique DNA
251 sequences encode the same peptide, we sum their populations by wash label (wash-1,
252 2, 3, or Eluate). We impose restrictions on the populations, where, if a DNA sequence
253 occurs three times or less in total, we do not include it in the translated peptide dataset.
254 Therefore, we end up with about 2 million unique peptide sequences after the preliminary
255 cleaning process discussed above.

256
257
258
259
260
261
262
263

264 **Results & Discussions**

265

266 **Data exploration and conditioning:**

267 After the preliminary cleaning and analysis, with >97% overlaps among the
268 biological replicates, we calculate the survival probabilities for each peptide sequence per
269 biological replicate. Per replicate, we sum the copy numbers of each unique peptide over
270 all the washes and eluate to obtain a total 'starting' population. Based on the total
271 population, we calculate the probability of survival as shown in the equations in Fig 3. Fig
272 3 also shows the survival probability trends of a sample of sequences over the washing
273 steps and eluate, for a biological replicate. Finally, we multiply the survival rates together,
274 and stagger them as shown in Fig 3, to obtain a single functional score: 'survival affinity',
275 as a measure of overall phage retention on the surface. We then categorize the peptides
276 as strong, medium, or weak binders based on the distribution of survival affinities. We
277 assume that survival affinity should correlate with the actual phage retention or surface
278 coverage values obtained for the Sanger-based benchmark set (see supplementary S1).

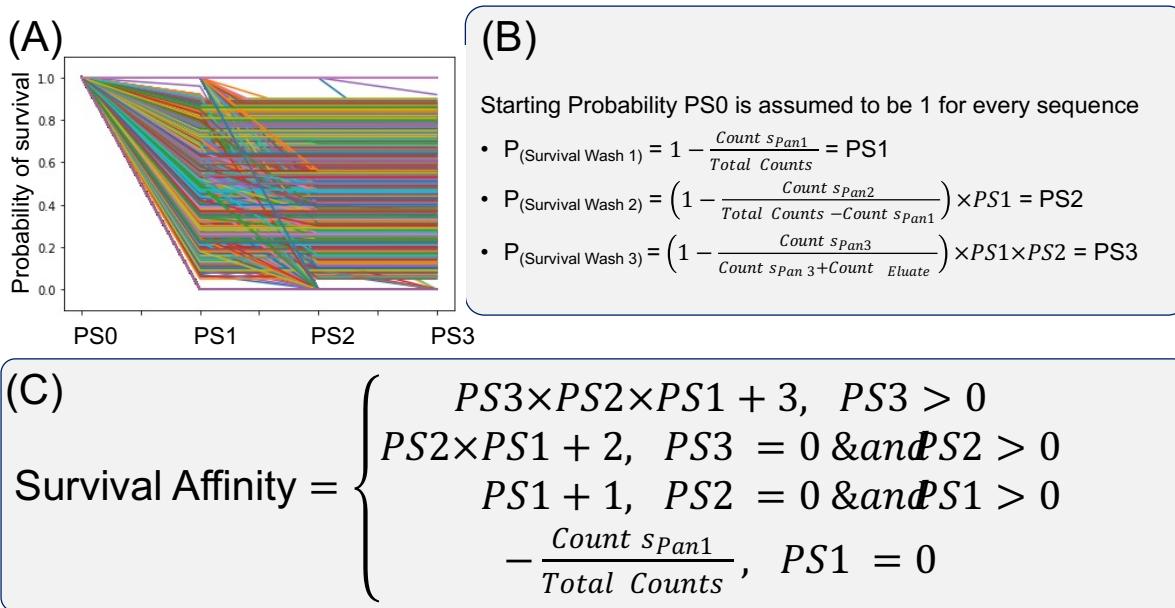


Fig 3. Functional Scoring for labeling massive sequence data. (A) Survival probability trends for a random sample of sequences. PS0, PS1, PS2, and PS3 are the probability of survival (probability of finding a sequence) before washing, after wash-1, after wash-2, and after wash-3 respectively, obtained from equations in (B). (C) Obtaining a single 'survival affinity' metric, the functional score, by multiplying and staggering the survival probabilities as shown.

279 After the cleaning process, we condition the data to make it ready for subsequent
280 data analysis, such as custom statistical algorithms (discussed in later sections) and
281 similarity analysis. We also obtain the 'center of abundance-mass' for each sequence,
282 i.e., the sum of the product of population in the wash and wash index (1 through 4; 4 being
283 the eluate) divided by total population as shown. The new metric is important because
284 many sequences are represented in all the four washes including the eluate. While the
285 survival affinity measure follows a Poisson-like distribution with multiple outliers, the
286 center-of-abundance-mass is Gaussian distributed, also with numerous outliers. Due to
287 the high number of outliers, i.e., significant outlier modes, we impose two other cleaning

288 criteria to condition the dataset properly. The first criterion is that the same peptide's
 289 survival trends in three different biological replicates must be about the same. Therefore,
 290 the mean difference in survival affinity between biological replicates for the same peptide
 291 must be zero. However, the mean difference among sets is nonzero with multiple modes
 292 and noisy distribution (Fig 4). We, therefore, impose a minimum copy number threshold
 293 (Fig 4). After the count restriction, we keep all sequences that behave similarly in at least
 294 two out of three biological replicates (Fig 4) and exclude the rest. As a result, the final
 295 well-cleaned and conditioned dataset contains upwards of a hundred thousand unique
 296 sequences in total, with their count numbers in all the washes and eluate, their survival
 297 affinities, and center-of-abundance-mass values.

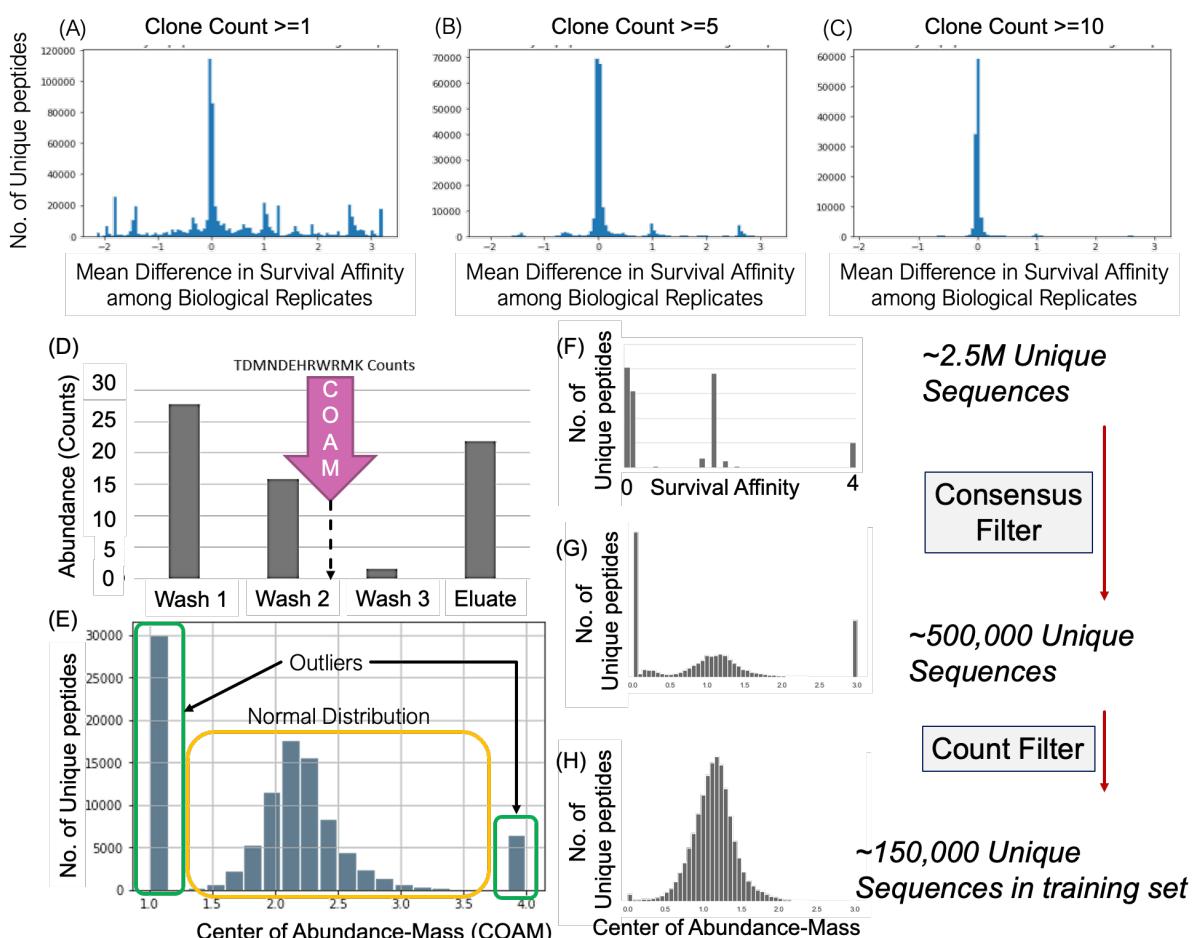


Fig 4. Data Cleaning, Conditioning and the ‘Center of Abundance Mass’ metric. (A)-(C) Upon imposing a minimum count threshold over the sequences, we see a marked decrease in noise and outliers in the mean difference in survival behavior for the same sequence in the three biological replicates. A mean difference of zero implies that the sequences behave in the same manner in all three replicates (Set1, Set2 and Set3). (D),(E) Most sequences appear in all four pools (wash1-3, Eluate). A Center of Abundance Mass (COAM) value for such a sequence is shown to be continuous variable falling anywhere between 1 and 4, with massive outlier modes at the extreme values. (F) We only consider peptides that display similar survival affinities in 2 out of 3 biological replicates. (G) We replace survival-affinity with the COAM value. (H) We impose the clone count threshold, and finally end up with about 150,000 unique sequences in the final dataset.

298 We split the now cleaned and conditioned dataset into two parts. We set aside
299 90% of the sequences to be the training set for predictive modeling, while the remaining
300 10% constitute the test-set for the same. The benchmark dataset is the 96 Sanger-based
301 sequences. In earlier work, authors show that the binding affinity of a peptide of interest
302 is proportional to the sum of the Global Alignment scores to a group of highly functional
303 peptides.^{39,40} We expand upon the concept and apply several custom preliminary
304 bioinformatics and machine learning approaches to the data described in the
305 supplementary information (S2 Fig and S2). Surprisingly, even increasingly complex
306 models failed to reconcile predicted survival-affinity trends satisfactorily with observed
307 trends in the test set and phage-retention trends in the benchmark-set.

308 The various models described in the supplementary (S2) that we apply to the data
309 had an underlying assumption that (a) the eluates are the strong binders and that (b)
310 strong binders have the highest self-similarity in their sequences. We analyze the
311 combinatorial diversity encapsulated in the washes through the means of information
312 entropy.⁴¹ We find that the eluate diversity is higher than all the other washes before;
313 therefore, either (i) the eluate does not contain just the strongest binders, or (ii) strong
314 binders do not have higher self-similarity than weak binders. However, previous studies
315 can rule out the second option as they have conclusively proven that strong binders do
316 indeed have high self-similarity^{42,43}. Therefore, the eluted peptides are not just the
317 strongest binding set. Indeed, peptides from previous washes (wash-1 through 3) also
318 exist in the eluate, although in far fewer numbers. Therefore, the eluate contains the
319 spectrum of strong MoS₂ binding peptides with a mechanism that was attacked by the
320 tween detergent and includes the entire range of binding strengths through various other
321 mechanisms that tween detergent did not attack at all.

322 323 **Analysis of sequence diversity to design future experiments:**

324 Next, we investigate the combinatorial coverage via information entropy analysis,
325 thereby exploring the relative diversity of the DNA sequence space captured for each
326 wash and the eluate. We weight the sequences by their copy numbers to keep the
327 analysis unbiased (i.e., use conditional probabilities instead of regular ones). Shannon's
328 entropy, S_{info} is a concept taken from statistical physics but repurposed to describe
329 information content by Claude Shannon⁴¹ and is mathematically expressed in equation
330 (1):

$$331 S_{info} = - \sum_{i=1}^n P_i \ln (P_i) \quad (1)$$

332 Where P_i is the probability of the i^{th} state of a random variable that can take multiple states
333 from 1 to n with varying probabilities. Applied to the case of DNA sequences, the finite
334 number of states, i, that each location on the sequence can take, are either A, C, T or G.
335 In the current dataset, there exist 36-length DNA sequences that code for 12-amino acid
336 long peptides. Therefore, for any given sample of sequences, a matrix with 36 rows
337 (locations on the sequence) and 4 columns (nucleotides) can be created where the sum
338 is 1 for each row. Using equation (1) on samples in the washes, and conditioning on the
339 same probabilities for samples in the input library (the input library is estimated by pooling
340 the sequences obtained from all the Washes, including eluates together).

341 The total sample information entropy measures the combinatorial coverage of
342 sequences in the sample with n sequences. We estimate maximum entropy for the DNA
343 sequences to be 49.90 (arbitrary units of 'nats'). Applying the above formalism to the
344 washes, the authors observed that the combinatorial coverage of wash-1 is greater than
345 wash-2. Similarly, the combinatorial coverage of wash-2 is greater than wash-3 (Fig 5).

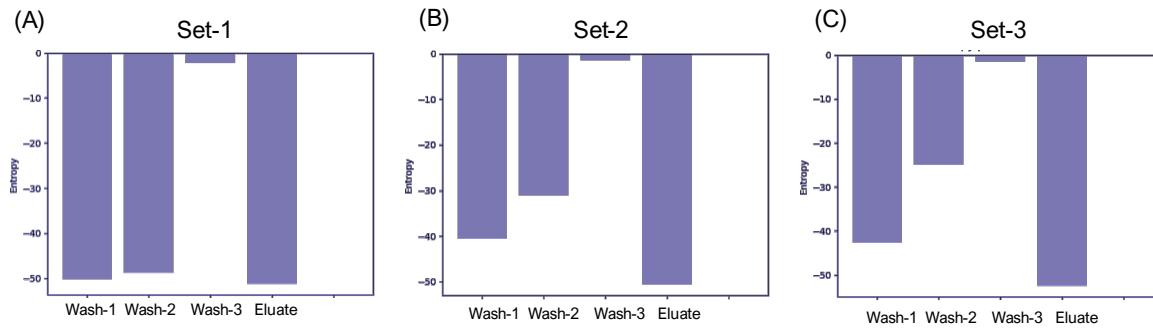


Fig 5. Information Entropy measures combinatorial diversity in the washes and Eluate.

(A)-(C) A greater negative value indicates more Information entropy (arbitrary units), more diversity in the pool (wash1-3, or eluate) across all biological replicates. As can be seen, the eluate has similar or higher diversity than all the washes, implying that the eluted sequences are more dissimilar to each other than those in the washes. We therefore cannot assume that the eluted phages are the strongest binders. Simply that the tween detergent used in the washes did not disrupt the mechanism of interaction between the eluted phages and the substrate. We thus propose a future experiment with washing performed with solutions that target specific mechanisms of substrate binding. A simple denaturation may not be enough to disrupt the interactions at the bio/nano interfaces.

346 Results in Fig 5 tell us that, indeed, sequences considered weak binders, i.e.,
347 sequences washed away in the first wash, are more diverse than the successive washes
348 and have little correlation among themselves, as expected. It comes as no surprise that
349 the sample entropy, a measure of combinatorial coverage and chaos, decreases in further
350 washes, clearly delineating that sequences that are increasingly strong binders to the
351 substrate are indeed more related to each other. The trend should also imply that
352 sequences in the eluate must then be even closely related and that eluate entropy and
353 combinatorial coverage must be less than wash-3. However, surprisingly, the entropy,
354 and therefore combinatorial coverage of the eluate is much higher than the weakest wash
355 (More details in supplementary S3). The supposed strongest binders specific to the same
356 substrate are more diverse and unrelated to each other than the weakest binders to the
357 same substrate. Such an observation seemingly turns our previous assumption that
358 strong binders to a specific substrate must be more alike than weak binders on its head.
359 A word of caution is that the previous assumption does indeed hold, but only if the binding
360 mechanism is the same. Following the previously established phage-display procedure,
361 the commonly used generic tween 20/80 detergent didn't categorically attack specific
362 binding modes and mechanisms. Moreover, many eluted peptides had smaller but
363 significant populations represented in the previous washes due to experimental error. The
364 effect of this 'information leakage' may be mitigated by transforming the survival affinity
365 into the 'center-of-abundance-mass' metric, as shown previously. Models that use the
366 latter metric to regress or classify upon may be more successful at reconciling the trends
367 observed in the training, test, and benchmark datasets.

368 Conclusions and Future steps

369

370 We have developed a universal first step in high-throughput peptide selection
371 approach to obtain solid-binding peptides in high-throughput and identified more than 2
372 Million different MoS₂ binding peptides. Through this paper, we disseminate the first and
373 largest ever solid binding peptide dataset labeled with survival probability scores for
374 binding activity on Molybdenum Disulfide substrate, selected via a high throughput
375 protocol that we adapt into GEPI selection. The dataset is cleaned, conditioned, and
376 analyzed to guide the next stage in high throughput deep directed evolution experiments
377 with functional scoring. Moreover, we also disseminate a carefully selected benchmark
378 GEPI dataset for affinity to MoS₂, characterized and labeled in terms of relative surface
379 coverage under Fluorescent microscopy as well as phage retention from
380 spectrophotometer absorbance. In addition to this, having a large number of
381 experimentally selected peptides enables in-depth exploration of the sequence space. It
382 allows Machine-Learning assisted designing of superior 2nd generation peptides with
383 desirable properties.

384

The efficacy of previously and newly developed machine learning models to
385 characterize the dataset was suboptimal due to the underlying assumption adopted from
386 biology that denaturing the peptides is enough to disrupt the binding events. Just because
387 a sequence is present in greater likelihood in the eluate and less in the washes does not
388 make the sequence a strong binder to the substrate per se, but rather that the generic
389 tween detergents are not disrupting the binding mechanism of that particular sequence.
390 Denaturation is not enough. We conclude that the eluate peptides bind with various
391 mechanisms not targeted by the tween detergent, necessitating a binding-mechanism
392 targeted approach. Explicitly using different washing solutions that target separate
393 substrate-phage interaction mechanisms such as hydrophobicity, aromatic-interaction,
394 amino-acid competition, etc., will be crucial in further resolving the factors that govern
395 binding interactions at the bio/nano interface. Moreover, because we see a significant
396 number of peptides that appeared in washes also turn up in the eluate, refinement and
397 optimization of the eluate via amplification steps coupled with next-generation sequencing
398 protocol are warranted.

399

The conclusion necessitates a multiplexed approach with a greater number of
400 intermediate panning and amplification steps to fully characterize the washes in terms of
401 functional scores that indicate binding affinity. Information entropy, which measures
402 diversity or the combinatorial coverage, is a good measure of the sample's reliability in
403 representing the population's distribution. The entropy analysis outlined above is a simple
404 yet successful method that can guide successive iterations (Fig 1 shows such an iterative
405 process) of the deep directed evolution experiment that we have introduced.

406

Additionally, amplifying the eluted peptides multiple times and repeating the
407 experiment seems like the natural next step to experimentally segregate the sequence
408 space along differences in binding mechanisms. Such a new and improved experimental
409 procedure, combined with the data exploration and quality control methods described in
410 the current work, promises to generate high-quality massive datasets for a thorough
411 exploration of evolutionary dynamics. It will enable highly efficient and rich directed
412 evolution experiments to select novel proteins and peptides for practical implementations
413 in technology and medicine and enrich our understanding of nature's design process.

414 **Acknowledgments**

415
416 The research is financially supported (DTY, SSR, JLR, JFL, ONB, MS) by National
417 Science Foundation (NSF) through the DMREF program (via Materials Genome Initiative)
418 under grant numbers DMREF DMR# 1629071, 1848911, and 1922020. We thank by
419 Jason J. Stephany and Douglas Fowler for technical help in NextGen selection processes
420 and the use of their facilities (UW Genome Sciences), and discussions in machine
421 intelligence and guidance by Kevin Jamieson (Allen School of Computer Science and
422 Engineering, University of Washington).

423
424 **References**

- 425
426 1. Arnold, F. H., Design by directed evolution. *Accounts of chemical research* 1998,
427 31 (3), 125-131.
428 2. Sarikaya, M., Tamerler, C., Jen, A. K.-Y., Schulten, K., Baneyx, F., Molecular
429 biomimetics: nanotechnology through biology. *Nature materials* 2003, 2 (9), 577.
430 3. Tamerler, C., Sarikaya, M., Molecular biomimetics: utilizing nature's molecular
431 ways in practical engineering. *Acta biomaterialia* 2007, 3 (3), 289-299.
432 4. Seeman, N. C., Belcher, A. M., Emulating biology: building nanostructures from
433 the bottom up. *Proceedings of the National Academy of Sciences* 2002, 99 (suppl
434 2), 6451-6455.
435 5. Wilson, D. S., Keefe, A. D., Szostak, J. W., The use of mRNA display to select
436 high-affinity protein-binding peptides. *Proceedings of the National Academy of
437 Sciences* 2001, 98 (7), 3750-3755.
438 6. Brown, S., Metal-recognition by repeating polypeptides. *Nature biotechnology*
439 1997, 15 (3), 269.
440 7. Whaley, S. R., English, D., Hu, E. L., Barbara, P. F., Belcher, A. M., Selection of
441 peptides with semiconductor binding specificity for directed nanocrystal assembly.
442 *Nature* 2000, 405 (6787), 665.
443 8. Di Battista, G., Liotta, G., Whitesides, S. In *The strength of weak proximity*,
444 International Symposium on Graph Drawing, Springer: 1995, pp 178-189.
445 9. Tamerler, C., Oren, E. E., Duman, M., Venkatasubramanian, E., Sarikaya, M.,
446 Adsorption kinetics of an engineered gold binding peptide by surface plasmon
447 resonance spectroscopy and a quartz crystal microbalance. *Langmuir* 2006, 22
448 (18), 7712-7718.
449 10. Tamerler, C., Sarikaya, M., Molecular biomimetics: genetic synthesis, assembly,
450 and formation of materials using peptides. *Mrs Bulletin* 2008, 33 (5), 504-512.
451 11. Gungormus, M., Fong, H., Kim, I. W., Evans, J. S., Tamerler, C., Sarikaya, M.,
452 Regulation of in vitro calcium phosphate mineralization by combinatorially selected
453 hydroxyapatite-binding peptides. *Biomacromolecules* 2008, 9 (3), 966-973.
454 12. Cetinel, S., Dincer, S., Cebeci, A., Oren, E. E., Whitaker, J. D., Schwartz, D. T.,
455 Karaguler, N. G., Sarikaya, M., Tamerler, C., Peptides to bridge biological-platinum
456 materials interface. *Bioinspired, Biomimetic and Nanobiomaterials* 2012, 1 (3),
457 143-153.
458 13. Naik, R. R., Stringer, S. J., Agarwal, G., Jones, S. E., Stone, M. O., Biomimetic
459 synthesis and patterning of silver nanoparticles. *Nature materials* 2002, 1 (3), 169.

- 460 14. Liu, G. W., Livesay, B. R., Kacherovsky, N. A., Cieslewicz, M., Lutz, E., Waalkes,
461 A., Jensen, M. C., Salipante, S. J., Pun, S. H., Efficient identification of murine M2
462 macrophage peptide targeting ligands by phage display and next-generation
463 sequencing. *Bioconjugate chemistry* 2015, 26 (8), 1811-1817.
- 464 15. Yazici, H., Fong, H., Wilson, B., Oren, E., Amos, F., Zhang, H., Evans, J., Snead,
465 M., Sarikaya, M., Tamerler, C., Biological response on a titanium implant-grade
466 surface functionalized with modular peptides. *Acta biomaterialia* 2013, 9 (2), 5341-
467 5352.
- 468 16. Yucesoy, D. T., Khatayevich, D., Tamerler, C., Sarikaya, M., Rationally designed
469 chimeric solid-binding peptides for tailoring solid interfaces. *Medical Devices &*
470 *Sensors* 2020, 3 (3), e10065.
- 471 17. Sedlak, R. H., Hnilova, M., Grosh, C., Fong, H., Baneyx, F., Schwartz, D.,
472 Sarikaya, M., Tamerler, C., Traxler, B., Engineered Escherichia coli silver-binding
473 periplasmic protein that promotes silver tolerance. *Appl. Environ. Microbiol.* 2012,
474 78 (7), 2289-2296.
- 475 18. Naik, R. R., Brott, L. L., Clarson, S. J., Stone, M. O., Silica-precipitating peptides
476 isolated from a combinatorial phage display peptide library. *Journal of nanoscience*
477 *and nanotechnology* 2002, 2 (1), 95-100.
- 478 19. Yucesoy, D. T., Karaca, B. T., Cetinel, S., Caliskan, H. B., Adali, E., Gul-Karaguler,
479 N., Tamerler, C., Direct bioelectrocatalysis at the interfaces by genetically
480 engineered dehydrogenase. *Bioinspired, Biomimetic and Nanobiomaterials* 2015,
481 4 (1), 79-89.
- 482 20. AC't Hoen, P., Jirka, S. M., Bradley, R., Schultes, E. A., Aguilera, B., Pang, K. H.,
483 Heemskerk, H., Aartsma-Rus, A., van Ommen, G. J., den Dunnen, J. T., Phage
484 display screening without repetitious selection rounds. *Analytical biochemistry*
485 2012, 421 (2), 622-631.
- 486 21. Mattock, W. L., Derda, R., Next-generation sequencing of phage-displayed
487 peptide libraries. In *Peptide Libraries*, Springer: 2015, pp 249-266.
- 488 22. Shtatland, T., Guettler, D., Kossodo, M., Pivoarov, M., Weissleder, R., PepBank-
489 a database of peptides based on sequence text mining and public peptide data
490 sources. *BMC bioinformatics* 2007, 8 (1), 280.
- 491 23. Huang, J., Ru, B., Li, S., Lin, H., Guo, F.-B., SAROTUP: scanner and reporter of
492 target-unrelated peptides. *BioMed Research International* 2010, 2010.
- 493 24. Schuster, S. C., Next-generation sequencing transforms today's biology. *Nature*
494 *methods* 2008, 5 (1), 16-18.
- 495 25. Shendure, J., Ji, H., Next-generation DNA sequencing. *Nature biotechnology*
496 2008, 26 (10), 1135.
- 497 26. Ansorge, W. J., Next-generation DNA sequencing techniques. *New biotechnology*
498 2009, 25 (4), 195-203.
- 499 27. So, C. R., Hayamizu, Y., Yazici, H., Gresswell, C., Khatayevich, D., Tamerler, C.,
500 Sarikaya, M., Controlling self-assembly of engineered peptides on graphite by
501 rational mutation. *ACS Nano* 2012, 6 (2), 1648-1656.
- 502 28. Fowler, D. M., Araya, C. L., Fleishman, S. J., Kellogg, E. H., Stephany, J. J., Baker,
503 D., Fields, S., High-resolution mapping of protein sequence-function relationships.
504 *Nature methods* 2010, 7 (9), 741.

- 505 29. Fowler, D. M., Fields, S., Deep mutational scanning: a new style of protein science.
506 *Nature methods* 2014, 11 (8), 801.
- 507 30. Wu Z, Kan SJ, Lewis RD, Wittmann BJ, Arnold FH. Machine learning-assisted
508 directed protein evolution with combinatorial libraries. *Proceedings of the
509 National Academy of Sciences*. 2019 Apr 30;116(18):8852-8.
- 510 31. Yang, K. K., Wu, Z., Arnold, F. H., Machine-learning-guided directed evolution for
511 protein engineering. *Nature methods* 2019, 16 (8), 687-694.
- 512 32. Chen, Y., Xu, D., Understanding protein dispensability through machine-learning
513 analysis of high-throughput data. *Bioinformatics* 2004, 21 (5), 575-581.
- 514 33. Verma, R., Schwaneberg, U., Roccatano, D., Computer-aided protein directed
515 evolution: a review of web servers, databases and other computational tools for
516 protein engineering. *Computational and structural biotechnology journal* 2012, 2
517 (3), e201209008.
- 518 34. Fox, R., Directed molecular evolution by machine learning and the influence of
519 nonlinear interactions. *Journal of theoretical biology* 2005, 234 (2), 187-199.
- 520 35. Saito, Y., Oikawa, M., Nakazawa, H., Niide, T., Kameda, T., Tsuda, K., Umetsu,
521 M., Machine-learning-guided mutagenesis for directed evolution of fluorescent
522 proteins. *ACS synthetic biology* 2018, 7 (9), 2014-2022.
- 523 36. Rentero Rebollo, I., Sabisz, M., Baeriswyl, V., Heinis, C., Identification of target-
524 binding peptide motifs by high-throughput sequencing of phage-selected peptides.
525 *Nucleic acids research* 2014, 42 (22), e169-e169.
- 526 37. Rubin, A. F., Gelman, H., Lucas, N., Bajjaleh, S. M., Papenfuss, A. T., Speed, T.
527 P., Fowler, D. M., A statistical framework for analyzing deep mutational scanning
528 data. *Genome biology* 2017, 18 (1), 150.
- 529 38. Zhao, L., Liu, Z., Levy, S. F., Wu, S., Bartender: a fast and accurate clustering
530 algorithm to count barcode reads. *Bioinformatics* 2017, 34 (5), 739-747.
- 531 39. Oren, E. E., Tamerler, C., Sahin, D., Hnilova, M., Seker, U. O. S., Sarikaya, M.,
532 Samudrala, R., A novel knowledge-based approach to design inorganic-binding
533 peptides. *Bioinformatics* 2007, 23 (21), 2816-2822.
- 534 40. Needleman, S., Needleman-Wunsch algorithm for sequence similarity searches. *J
535 Mol Biol* 1970, 48, 443-453.
- 536 41. Shannon, C. E., Prediction and entropy of printed English. *Bell system technical
537 journal* 1951, 30 (1), 50-64.
- 538 42. Gungormus, M., Oren, E. E., Horst, J. A., Fong, H., Hnilova, M., Somerman, M. J.,
539 Snead, M. L., Samudrala, R., Tamerler, C., Sarikaya, M., Cementomimetics—
540 constructing a cementum-like biominerilized microlayer via amelogenin-derived
541 peptides. *International journal of oral science* 2012, 4 (2), 69.
- 542 43. Ricotta, C., Szeidl, L., Towards a unifying approach to diversity measures:
543 bridging the gap between the Shannon entropy and Rao's quadratic index.
544 *Theoretical population biology* 2006, 70 (3), 237-243.
- 545
- 546
- 547
- 548
- 549

550 **Supplementary Information**

551 **S1: Conventional Directed Evolution for Benchmark data Generation.**

552 The 12-mer Phage display (PhD) library is used to select peptide sequences fused
553 to the minor coat protein (pIII) against MoS₂ flakes as described above. The obtain
554 enough number of clones for subsequent rounds, eluted phages are transferred to an
555 early log phase E. coli ER2738 culture (~OD: 0.4) and amplified for 4 hours. The cell
556 pellet is obtained by centrifugation and purified by polyethylene glycol (PEG) precipitation
557 according to the manufacturer's instructions. Purified phages are obtained in PC buffer
558 volume pH 7.4 with a final volume of 200 µl. For the individual phage isolation, eluate
559 pools are streaked on agar plates and individual colonies were picked. Amplification is
560 performed in the early log phase E. coli ER2738 culture for 4 hours followed by
561 polyethylene glycol (PEG) precipitation.

562 The single-stranded DNA of selected phage plaques are isolated by a QIAprep
563 Spin M13 Kit (Qiagen, Valencia, CA) and amplified via PCR in the presence of dye-
564 labeled terminators (Big dye terminator v3.1, Applied Biosystems, Carlsbad, CA). PCR
565 products are purified by Sephadex G-50 column precipitation. A 96 gIII primer (5'-OH
566 CCC TC TAGTTA GCG TAA CG-3') is used for the amplification of ssDNA. The selected
567 sequences of DNA from clones are analyzed by an Applied Biosystems 310 Avant DNA
568 analyzer. After sequencing, the isolated clones are then tested individually for their
569 relative binding affinities to MoS₂. To measure the relative binding affinity, two different
570 assays are used.

571 Firstly, the individual colonies are mixed with 5 mg of MoS₂ flakes dispersed in 1
572 mL of PC Buffer (pH 7.4), containing 0.02% Tween 20 detergent and incubated for 3
573 hours on a rotator at room temperature. After washing off loosely bound phages from the
574 MoS₂ surface using 1 mL of PC Buffer (pH 7.4), containing 0.1% Tween 20 detergent, the
575 bound clones are fluorescently labeled using anti-M13 antibodies. The samples are
576 characterized by quantitative fluorescent microscopy, employing a Nikon Eclipse TE-
577 2000U fluorescent microscope (Nikon, Melville, NY, USA) equipped with a Hamamatsu
578 ORCA-ER cooled CCD camera (Hamamatsu, Bridgewater, NJ, USA), imaged using a
579 FITC filter (exciter 460–500 nm, dichroic 505 nm, emitter 510–560 nm) and MetaMorph
580 imaging system (Universal Imaging, West Chester, PA, USA). Finally, the binding affinity
581 for each peptide is determined by calculating the ratio of total fluorescent intensity over
582 total surface area of the MoS₂ flakes (n=10). Secondly, for spectrophotometry based
583 analysis, number of bound phages are determined by analyzing the broad optical
584 absorption peak located from 260 nm to 280 nm, with a slight maximum at 269 nm which
585 reflects the nucleotide content of the particular phage, a molar extinction coefficient
586 ($9.006 \times 10^3 \text{ M}^{-1}\text{cm}^{-1}$) where the genome size of M13KE is taken as 7222 base-pairs. Prior
587 to MoS₂ incubation, the initial concentration of each clone is calculated. Next, the relative
588 binding affinities are calculated from the measured absorption intensities of depleted
589 (unbound) phage solutions as a percentage of the original absorption intensity from
590 starting solutions. Each experimental set is performed in triplicate.

591 A total of 144 colonies are isolated and sequenced from each eluate (36 colonies
592 after each biopanning) which yielded a total of 96 unique sequences. As shown in Fig 2,
593 the binding coverage of selected peptides ranges from 20% to 90% demonstrating the
594 presence of strong, moderate, and weak binders. The fluorescent signal is obtained
595 through immunolabeling by fluorescently labeled anti-M13 antibodies and therefore the

596 total intensity is a function of the number of antibodies on the surface. It is important to
597 note that depending on the average number of antibodies on the phage surface, the total
598 intensity can vary significantly. On the other hand, in spectrophotometric method,
599 absorption at maximum at 269 nm is a direct measurement for the total genomic DNA
600 content in a phage solution which provides a direct quantification for phage concentration
601 present on the surface. As shown in Fig 2, the results of each assay have about 70
602 percent overlap, with 4 exceptions (clone 1, 21, 34, 54). One possible explanation for the
603 discrepancy between the two assays, in particular, clone# 1, 21, 34, 54, could be the
604 repeated centrifugation steps where the phages can be trapped between the pelleted
605 flakes and causes false negatives in fluorescent imaging. Moreover, in
606 spectrophotometric characterization, the incomplete sedimentation of tiny MoS₂ particles
607 could cause false positives with increased UV absorption.
608

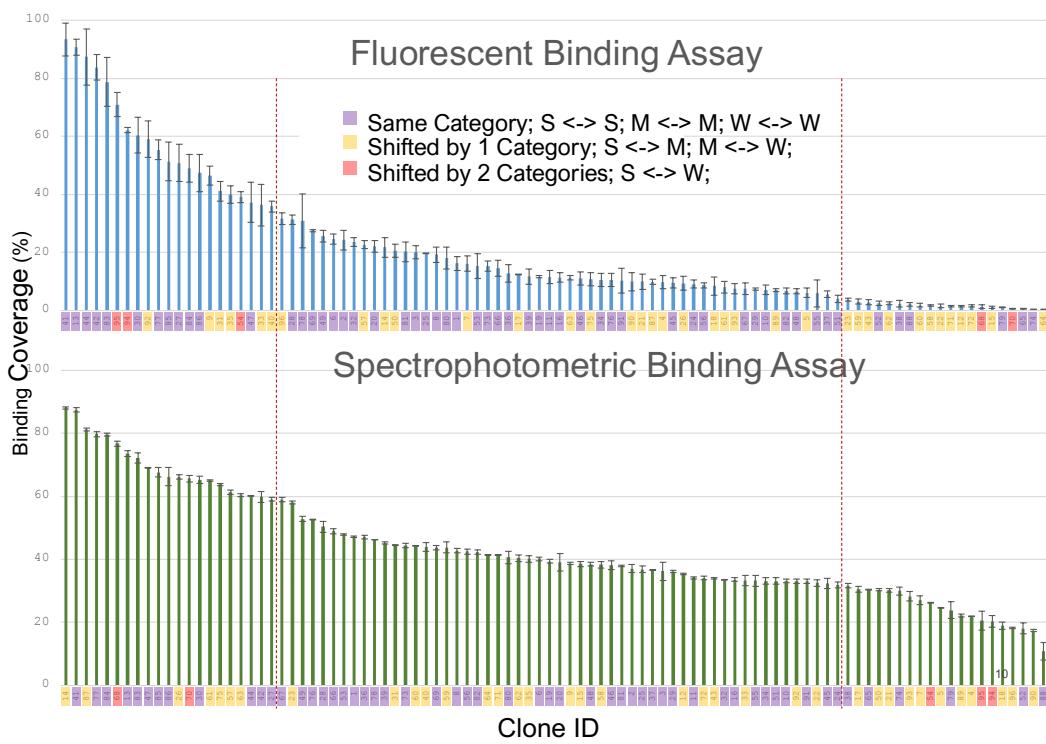


Fig 2: Relative binding affinities of MoS₂ binding peptides based on FM (upper) and Spectrophotometric (lower) selected through traditional sanger based combinatorial mutagenesis screening.

609 **S1 Fig. Fluorescent Microscopy and Spectrophotometry characterization of phage**
610 **retention for sequences obtained through conventional directed evolution and**
611 **sanger sequencing.**

612

613

614

615

616 **S2: Predictive Modelling**

617 In previous work, it was shown that the binding affinity of a peptide of interest is
618 proportional to the sum of the Global Alignment scores to a group of highly functional
619 peptides. Like the novel approach described here, this functional scoring, referred to as
620 the Total Similarity Score (TSS) trend was compared with the experimental binding trend
621 generated by the fluorescence microscopy procedure described above. Traditional
622 models previously generated in GEMSEC were applied in addition to expanded
623 applications of the TSS concept to diversify our analysis.
624

625 **Traditional Similarity Analysis/Iterative Alignment:** Similarity Analysis, or
626 Iterative Alignment, is the first computational scoring method that was developed by
627 GEMSEC, being benchmarked on quartz binding peptides in 2007. Essentially, Similarity
628 Analysis, or Iterative Alignment, optimizes a similarity matrix (distance matrix between
629 amino acids) that ensures the most active (binding, mineralization etc.) peptides are more
630 similar to each other than they are to the least active peptides. The Total Similarity can
631 be explicitly stated as the averaged Needleman-Wunsch alignment score of a group of
632 sequences A to another group B, where the Total Similarity Score of A to B is written
633 $TSS_{A:B}$. The internal similarity of the highly active peptides is captured by the Total
634 Similarity Score Strong-Strong (TSS_{ss}). External similarity of the highly active to the least
635 active is referred to as the Total Similarity Score Strong-Weak (TSS_{sw}). Random changes
636 to the similarity matrix that increase the TSS_{ss} and decrease TSS_{sw} are considered
637 'beneficial changes' and are done until specified by the experimenter. In general, the
638 sequences should become more and more related in sequence as the peptides become
639 more similar in function, leading to a characteristic trend demonstrated by the quartz
640 application (S3 Fig). The seed matrix most commonly used, PAM250 was derived in the
641 seventies by leveraging how mutations affected the function of proteins between closely
642 related species and served as a successful starting point for trained matrices in the
643 original work. In order to generate the characteristic trend of internal similarity to confirm
644 the directed functionality of the dataset, PAM250 matrices were trained using 1000 highly
645 active peptides (highest survival affinity) and 1000 less active peptides (lowest survival
646 affinity) from each biological replicate, including a combined set ensuring consistency of
647 survival affinity across at least 2-3 sets. These matrices were then used to score 5 groups
648 of sequences taken from the three 2.5 million data-points (Strong, Less Strong, Medium,
649 Less Weak and Weak). The 96 Sanger MoS₂ binding peptides were used to benchmark
650 the matrices trained on the larger dataset with affinity for the same material by classifying
651 each sequence by their affinity (Strong, Medium, Weak). Although the trend of internal
652 similarity tended to decrease among all the bar charts generated (Supplementary Section,
653 S3 Fig), the matrices from all datasets except the combined set struggled to place the
654 benchmark sequences in their correct affinity class (Supplementary Section, S3 Fig).
655

656 **Machine-learning on Total Similarity Scores:** The method referred to as ML on
657 TSS (Machine Learning on Total Similarity Scores) is best characterized as a multiple
658 regression on eight independent Total Similarity Scores featurized by 550+
659 physiochemical properties of a peptide set towards prediction of their survival affinity
660 downloaded from a database called the Amino Acid Index, which contains hundreds of
661 curated properties that have been experimentally measured. A schematic is shown in S4

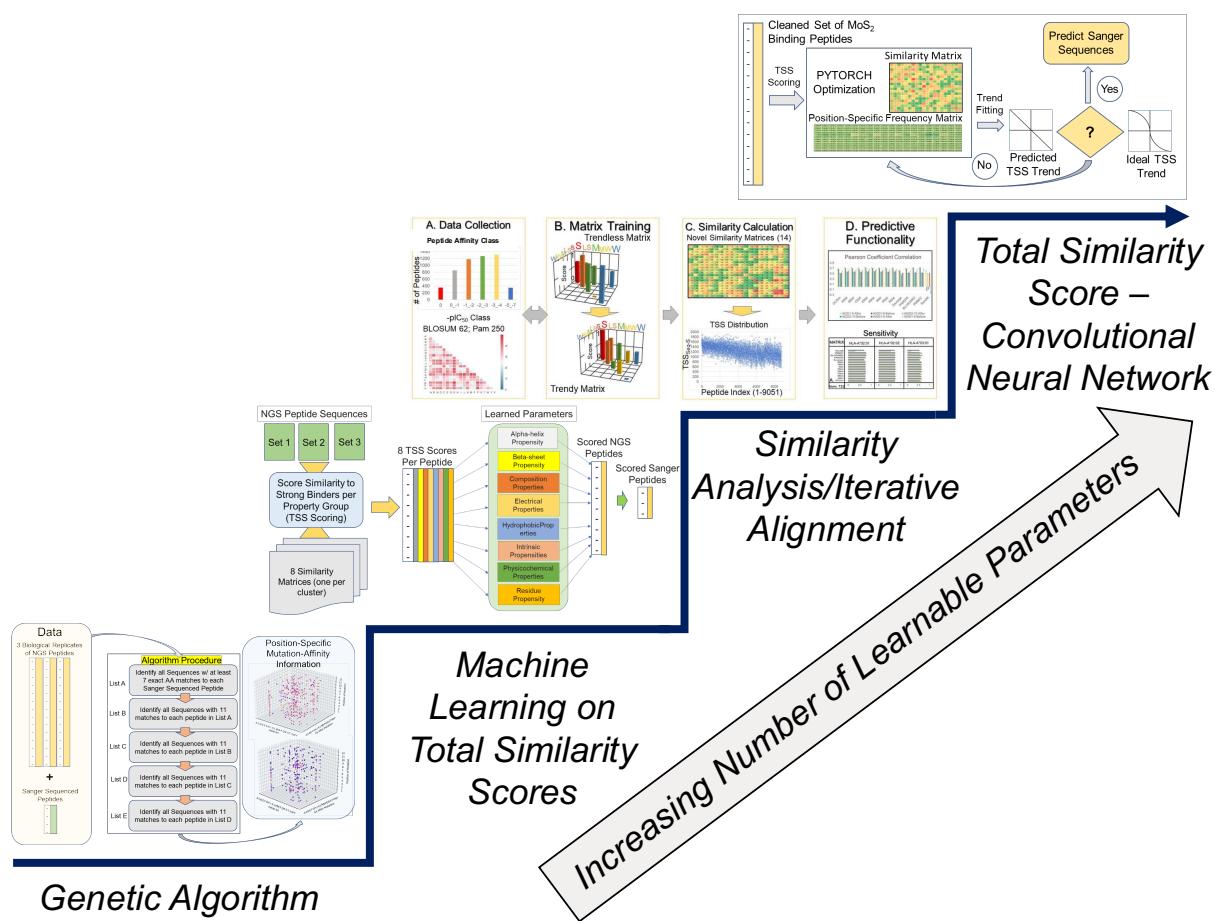
662 Fig. The peptides used as the highly functional group were those sequences with the
663 highest survival affinity values. The distances between AAs were derived by performing
664 Principal Component Analysis (PCA) on each group of properties to generate ~20
665 orthogonal variables by which each amino acid could be described by its Cartesian
666 distance to the others. In the next step, the similarity matrices were populated by the
667 appropriate distances from one AA to another, creating a symmetric matrix. The
668 sequences were scored by calculating their TSS with respect to the most active peptides
669 (the strongest binders with the highest survival affinity) using each similarity matrix,
670 generating 8 scores per peptide. The experiment performed herein was done to explore
671 the effect of the copy number (total number of copies of a particular peptide present during
672 the wash process) on the prediction accuracy. 9 datasets were created that were
673 distinguished by the minimum copy number of peptides that were identified in each
674 biological replicate. The regression was overall able to predict the survival affinities of
675 other sets and themselves with high accuracy (as low of a mean square error as ~7 for
676 the most stringent count discrimination; S4 Fig). As the count discrimination increases,
677 the mean square error decreases, indicating the data containing the most information
678 about the overall trend is present in the sequences having more copies in the selection
679 process.
680

681 **Total Similarity Score – Convolutional Neural Network:** To deepen the learning
682 space of the methods while keeping to the cost function known to describe peptide
683 activity, a constrained convolutional neural network was generated that ensures the
684 parameters necessitated by similarity matrices (symmetric matrix) and position-specific
685 frequency matrices (columns per position must sum to 1). Neural networks have been
686 used to classify and quantify the affinity of smaller datasets of thousands of peptides with
687 reasonable success, indicating a large amount of learning space is required to capture
688 the informatic trend. The TSS-CNN method can be described as a simple constrained
689 neural-network trained to generate a trend of affinity scores that correspond to the $-x^3$
690 curve, the trend always generated by traditional similarity analysis. A schematic of the
691 process is given by S5 Fig. The strength of this method is in its ability to generate the full
692 representative space of how peptides bind to MoS₂ (or the directed function of the
693 dataset) by optimizing the locations and is to identify two matrices of size (12x20) and
694 one of (20x20) are randomly initialized and their constraints applied. Next, the TSS to the
695 frequency matrix representing the strongest binders is calculated per peptide by
696 multiplying the frequency of occurrence of all AAs in the column corresponding to the
697 position of the peptide of interest by the mutability of the AA in that position. The process
698 is repeated across all positions until the entire peptide has been effectively matched with
699 the entire group of high survival affinity peptides (represented by the frequency matrix).
700 The values calculated per position and across the peptide are summed to create a single
701 TSS for the peptide. When applied to a set with directed functionality (i.e. binding to
702 MoS₂), a distribution is generated. Changes to the similarity matrix and frequency
703 matrices were kept dependent on whether they increased the nonlinear correlation of
704 these TSS values with an evenly spaced distribution of values on the $-x^3$ trend. These
705 matrices were expected to be representative of the peptide binding trend defined by the
706 surface. The first iteration of the CNN failed to predict a trend in affinity for the benchmark
707 peptides. Compared to the fluorescence (green line, S5 Fig) and spectrophotometric

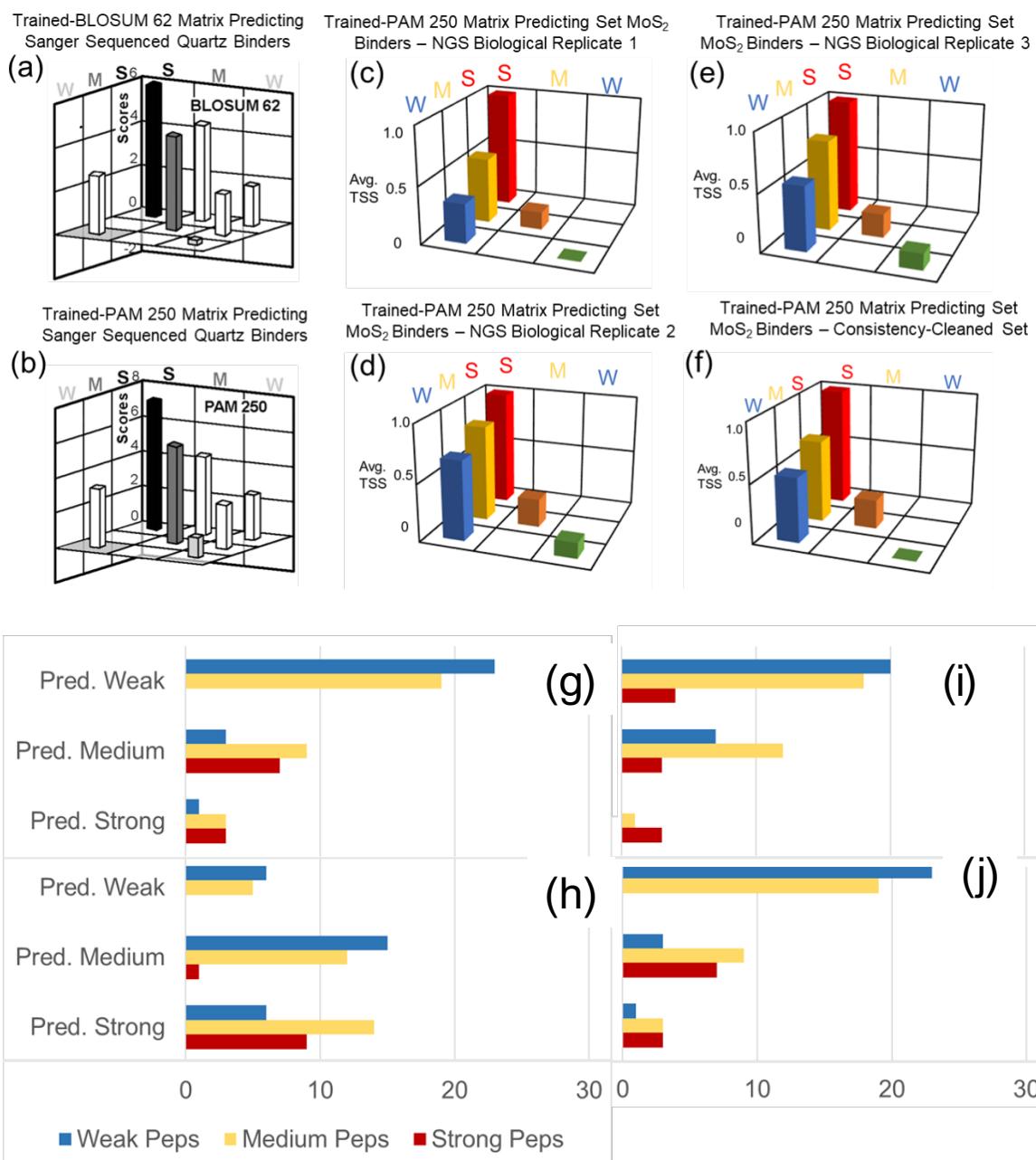
708 (orange line, S5 Fig) trends, the TSS trend (blue lines, S5 Fig) is unable to group binders
709 into semi-quantitative classes of strong, medium, and weak. Although the sum of all the
710 frequencies per position in the Position Specific Frequency Matrix summed to 1, the fifth
711 position blew its values out of proportion with the rest, likely attributable to a wide diversity
712 of amino acids giving a wide range of functions in this position. In the updated method
713 new constraints will be applied such that the values must remain within a certain range.
714

715 **Genetic Algorithm:** The purpose of this method is to diversify the approach
716 beyond machine learning and leverage the size of the dataset to identify changes in
717 sequence that would change the function of the benchmark peptides in a controllable
718 manner. This method focused on identifying mutational changes in peptides similar in
719 sequence to the original 96 MoS₂ binding peptides with the goal of predictively designing
720 the affinity of new MoS₂ binders. A schematic of the full process is provided by S6 Fig.
721 Once the full list of similar peptides had been identified, they were filtered by the
722 consistency of their survival affinity across the 3 biological replicates, to ensure the data
723 was the most accurate. Further, only sequences with at least 8 total copies during the
724 selection process were included in the final lists. After the cleaning process was
725 completed, the differences in binding affinity between sequences only 1 AA apart was
726 calculated and tracked along with the mutation and the position where it occurred. These
727 changes in amino acid sequence were identified by first finding peptides within each of
728 the 3 biological replicates that contained at least 7 amino acid (AA) matches (having the
729 same amino acid in the same position) with the 96 Sanger peptides that were selected
730 for affinity to MoS₂. This process was repeated 5 times to ensure similarity to the original
731 sequences was maintained and to capture a sizeable portion of the dataset. This
732 information was graphed in a 3D matrix (x = Initial AA, y = Final AA, z = Position) where
733 the color indicates the change in affinity associated with the mutation (S6 Fig). The
734 standard deviation of each affinity change was also graphed to ensure mutations tested
735 had a consistent effect across peptides. Overall, the mutations affecting the survival
736 affinity of the peptides were the cysteine to alanine mutations in position 4 (fourth amino
737 acid from the N terminus) which was associated with a dramatic loss of function, while
738 the largest gain-of-function mutations were centered around the mutation of cysteine to
739 leucine and tyrosine to asparagine. It is reasonable to attribute the gain-of-function results
740 to the affinity of asparagine (positively charged) for the negatively charged MoS₂ surface.
741 Further, the loss of a highly negatively charged group (cysteine) to a very hydrophobic
742 chain (leucine) may be indicative of a larger trend of hydrophobicity lending weight to the
743 strength of the interaction.

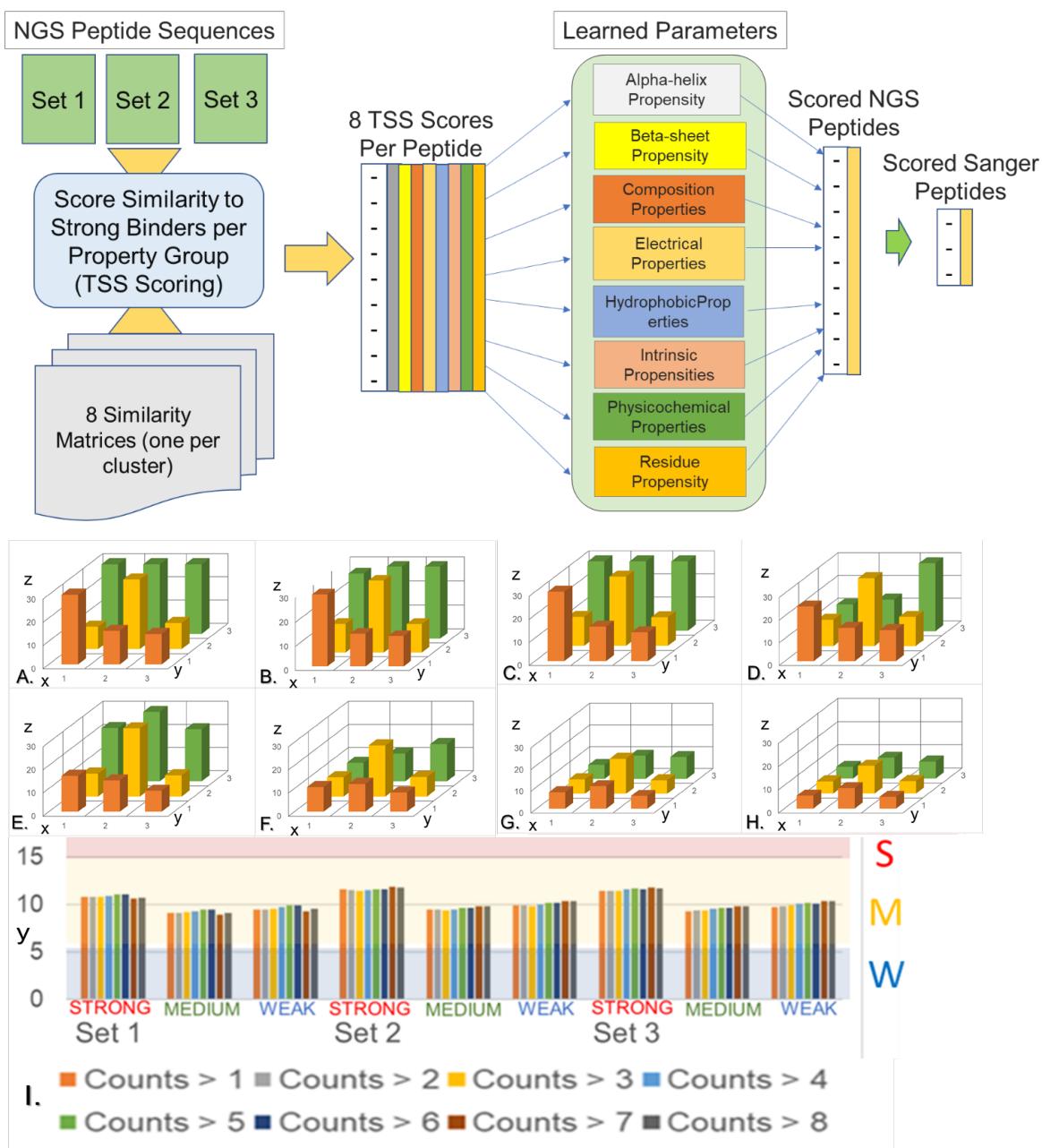
744



745 **S2 Fig. Increasingly complex mathematical models for preliminary predictive**
 746 **analysis on the dataset.**
 747

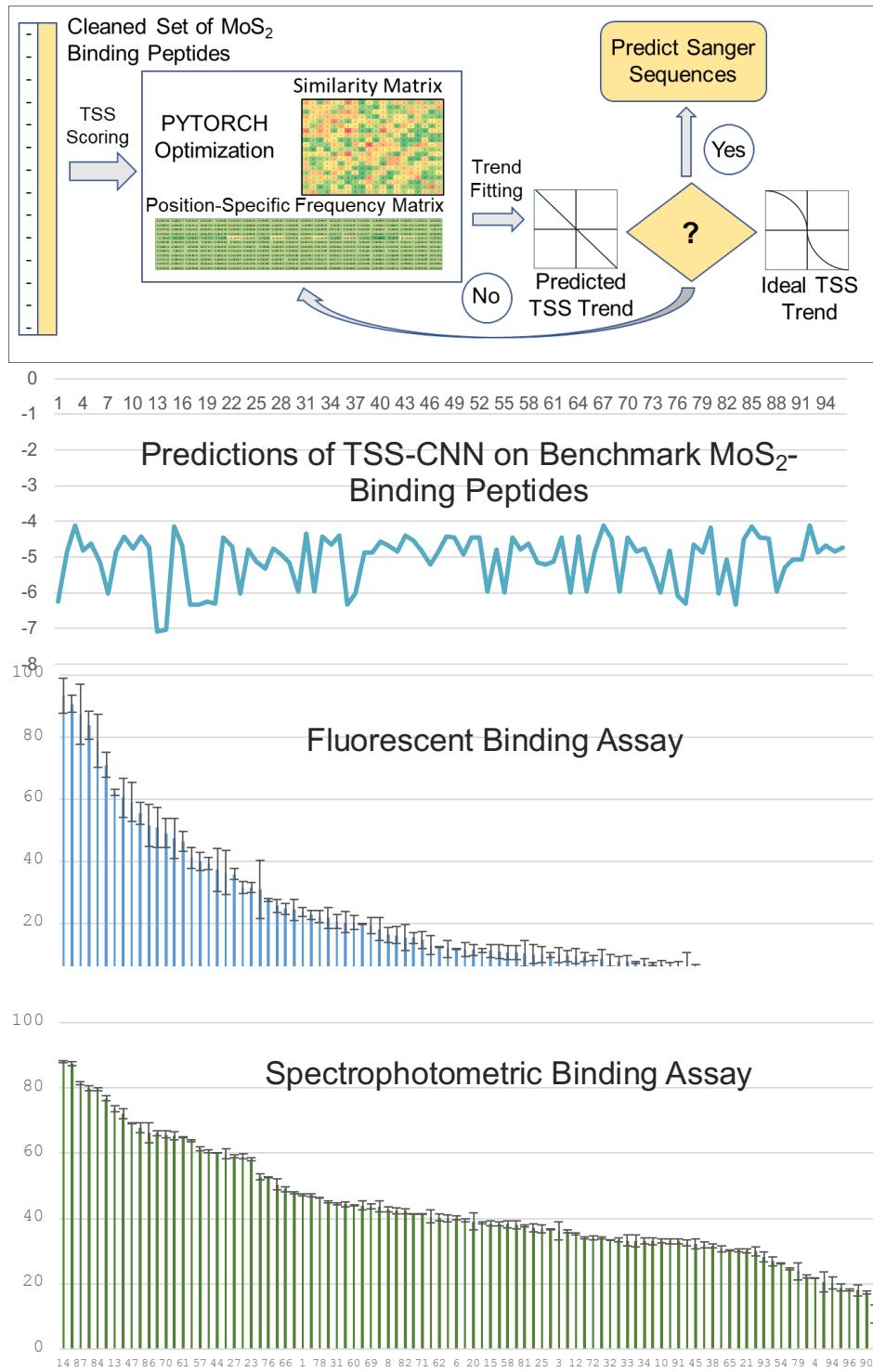


748 **S3 Fig. Traditional Similarity Analysis and Iterative Alignment process**
 749

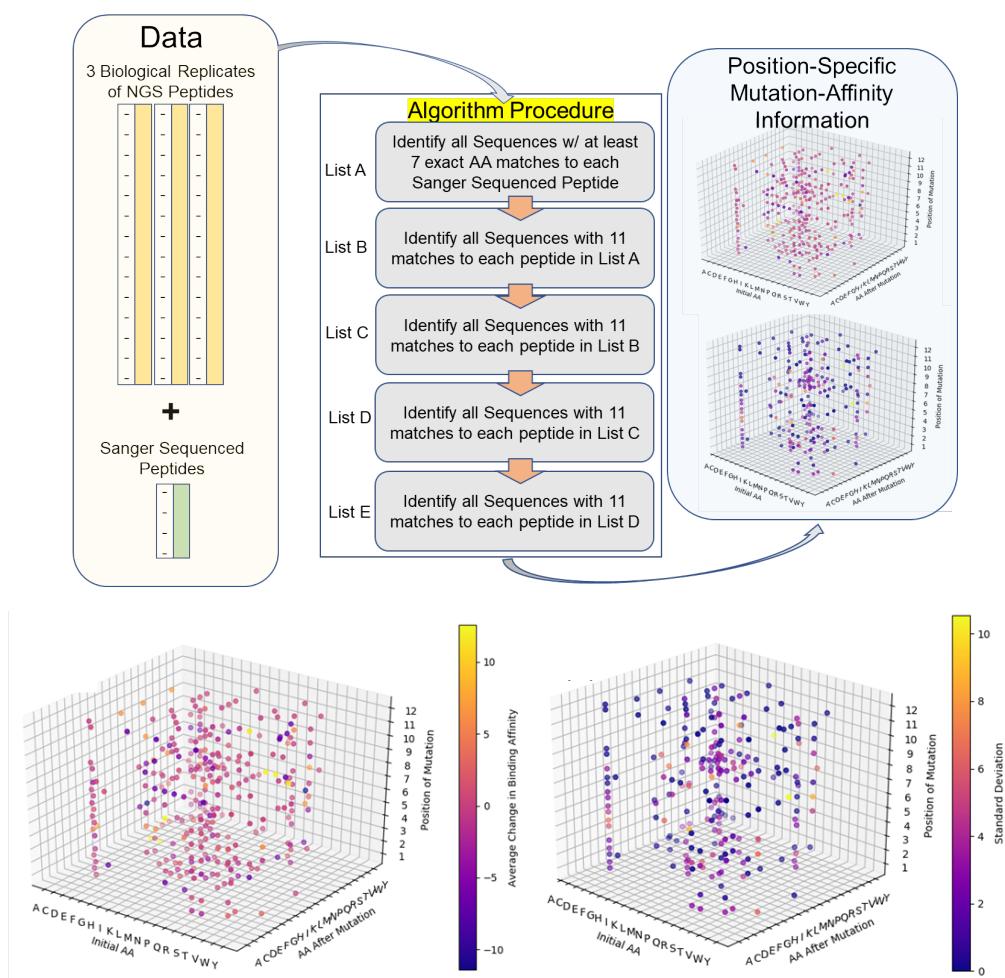


750 **S4 Fig. Machine Learning on TSS Scores.**

751



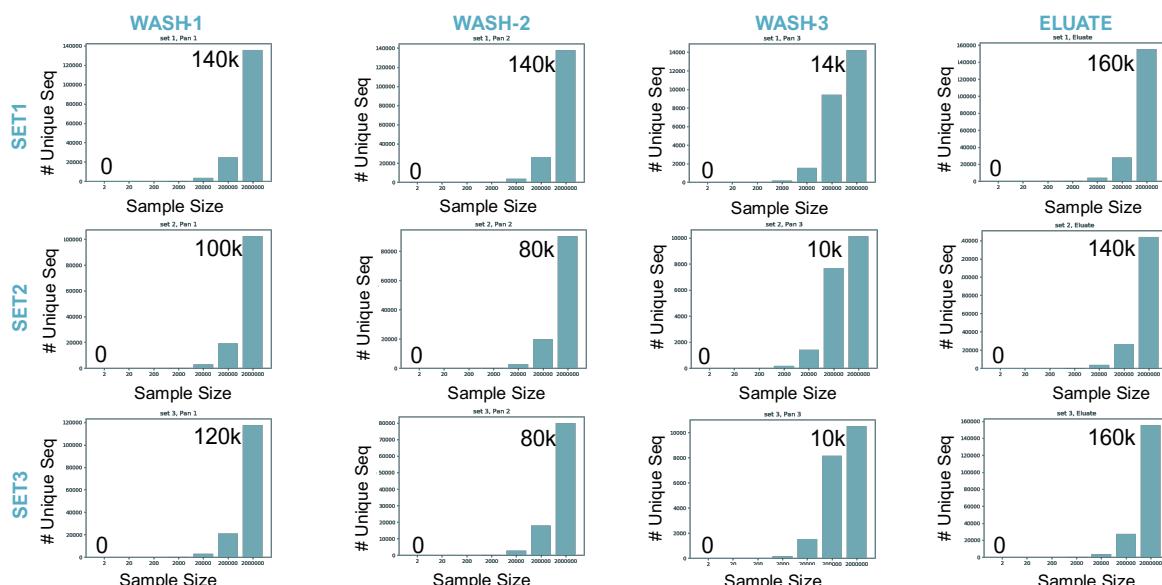
752 **S5 Fig. TSS-CNN: Applying a simple feed forward neural network on the TSS
753 scores.**
754



755 **S6 Fig. Genetic Algorithm Schematic and results of evolutionarily conserved
756 mutations as described in S2.
757**

758 **S3: Assessment of Diversity by Random Sampling of Phage Pools.**

759 To further analyze the diversity of the sequences, we sample 2, then 20, then 200 (with
760 replacement) and continue sampling until we reach 2 million DNA sequences from each
761 wash and eluate phage pool and calculate the number of unique sequences per draw. If
762 the number of unique sequences (by extension, diversity) increases upon continuously
763 sampling more sequences from the same wash/eluate, then the minimum sample size
764 has not been reached that would satisfactorily describe the behavior of the theoretical
765 population in that phage pool. Increasing number of unique sequences upon larger
766 samplings from the same wash/eluate would mean that there are still more combinations
767 of nucleotides available that would display the same behaviors. While a sigmoidal nature
768 of the trend, upon continuously increasing sample size from the same wash/eluate would
769 indicate that we are approaching the sample size that satisfactorily represents all/most
770 possible combinations of nucleotides. In short, increasing trend of entropy indicates that
771 sample size is not large enough to fully describe the combinatorial space of the
772 wash/eluate while a sigmoidal entropy trend indicates the opposite. As seen in the figure
773 below, the uniqueness trends are exponentially increasing for Wash 1 and Wash 2 in all
774 biological replicates of the experiment, while it is stabilizing for Wash 3. The results prove
775 that indeed Washes 1 and 2 do not yet have a representative combinatorial coverage of
776 sequence space to describe weak binders, while Wash 3 is approaching a sample size
777 which would possess a representative combinatorial coverage of sequences to describe
778 the strong binders that bind with a mechanism that is disrupted by the tween detergents.
779 The eluate deviates from this trend and entropy increases as we randomly sampled more
780 and more sequences from the eluate, indicating that the supposed strong binders are
781 indeed binding with multiple mechanisms that necessitate a more targeted approach to
782 directed evolution and NGS protocol than that used in this study. Moreover, because we
783 see a significant representation of peptides that appeared in washes also turn up in the
784 eluate, and since the input library was not explicitly sequenced, a refinement and
785 optimization of the current directed evolution and next generation sequencing protocol is
786 warranted.



S7 Fig. assessment of diversity and sample size requirements in phage pools.