

# Uncursing winner's curse: on-line monitoring of directed evolution convergence

Takahiro Nemoto,<sup>1,2,3,\*</sup> Tommaso Ocari,<sup>1</sup> Arthur Planul,<sup>1</sup> Muge Tekinsoy,<sup>1</sup> Emilia A. Zin,<sup>1</sup> Deniz Dalkara,<sup>1,†</sup> and Ulisse Ferrari<sup>1,‡</sup>

<sup>1</sup>*Institut de la Vision, Sorbonne Université, INSERM, CNRS, 17 rue Moreau, 75012, Paris, France*

<sup>2</sup>*Graduate School of Informatics, Kyoto University,*

*Yoshida Hon-machi, Sakyo-ku, Kyoto, 606-8501, Japan*

<sup>3</sup>*Premium Research Institute for Human Metaverse Medicine (WPI-PRIME),*

*Osaka University, Suita, Osaka 565-0871, Japan*

(Dated: This manuscript was compiled on April 6, 2023)

Directed evolution (DE) is a versatile protein-engineering strategy, successfully applied to a range of proteins, including enzymes, antibodies, and viral vectors. However, DE can be time-consuming and costly, as it typically requires many rounds of selection to identify desired mutants. Next-generation sequencing allows monitoring of millions of variants during DE and can be leveraged to reduce the number of selection rounds. Unfortunately the noisy nature of the sequencing data impedes the estimation of the performance of individual variants. Here, we propose ACIDES that combines statistical inference and in-silico simulations to improve performance estimation in DE by providing accurate statistical scores. We tested ACIDES first on a novel random-peptide-insertion experiment and then on several public datasets from DE of viral vectors and phage-display. ACIDES allows experimentalists to reliably estimate variant performance *on the fly* and can aid protein engineering pipelines in a range of applications, including gene therapy.

Keywords: directed evolution | phage display | deep mutational scanning | protein engineering | next generation sequencing

## INTRODUCTION

Directed evolution (DE) [1–3] is a versatile protein engineering strategy to conceive and optimize proteins like enzymes [4–6], antibodies [7, 8] or viral vectors for gene therapy [9–15], culminating in the Nobel Prize in Chemistry 2018 [16]. DE starts from a massive library of random mutants, screens it against a given task over multiple rounds and searches for the variants with the highest performance. As the iteration continues, the best performing variants get enriched and emerge from the bulk, while ineffective ones are instead weeded out. Nowadays, we can rely on next generation sequencing (NGS) [17, 18] to sample millions of variants within the library and monitor their concentrations over multiple rounds or time-points. In this approach, the enrichment of the screened variants is measured to rank the variants depending on their performance. In a similar flavor, Deep mutational scanning (DMS) experiments [19–21] combine extensive mutagenesis with NGS to study the properties of proteins [22–26], promoters [27, 28], small nucleolar RNA [29], or other amino-acid chains. It uses similar techniques to DE and requires similar analysis. The approach presented here can be applied to both DE and DMS experiments, and focus on their common issues and needs.

The analysis of NGS data of multiple selection rounds presents several difficulties. First, variants need to be robustly scored based on their enrichment rates, so-called

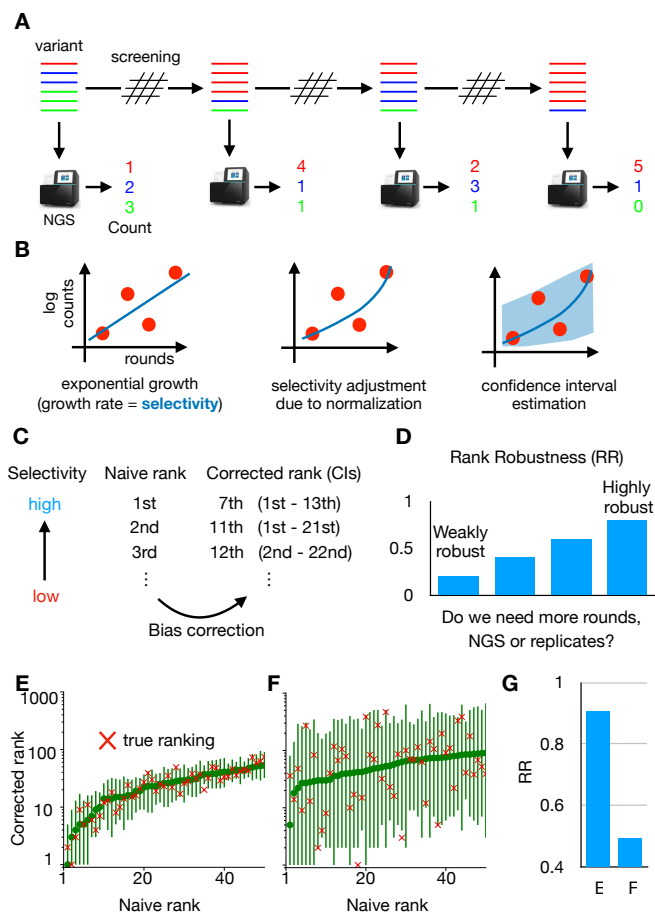
selectivities [30, 31]. This task is complicated by the large noise in the NGS counts introduced by, for example, polymerase chain reaction (PCR) amplification or bacterial cloning, during amplicon preparations [32–34]. This noise needs to be taken into account in the analysis. Second, in order to rank the variants and to identify the best performing ones, the score should come with a precise estimation of its statistical error. As a consequence of the noise in the counts, some irrelevant variants might appear to be highly enriched (winner's curse). This would be anticipated if properly estimated credibility scores are available. Third, when running DE over multiple rounds, it is hard to know when to end the experiment: performing too few rounds could lead to selection of weak variants, not representative of their true ranking. On the other hand, performing too many rounds is costly, time-consuming and even ethically questionable when working with *in-vivo* selections [14, 35]. Similarly, it would be useful to understand the best NGS depth for a given experiment, as deepening the NGS by increasing reads results in better data, but adds an extra expense to the experiment.

In order to account for these issues and needs, we present ACIDES, Accurate Confidence Intervals for Directed Evolution Scores, a computational method to empower the analysis of DE and DMS experiments. We focus on screening experiments on highly diverse libraries where massive NGS data are collected over multiple rounds or multiple time-points (Fig. 1A). Our goal is to develop a method to extract maximal information from noisy NGS data, and allows for scoring and ranking variants with accurate statistical confidence scores. Our approach can be applied to different kinds of experiments, such as *in-vivo* DE [13, 14, 36], and DMS of phage-display

\* nemoto.takahiro.prime@osaka-u.ac.jp

† deniz.dalkara@inserm.fr

‡ ulisse.ferrari@inserm.fr



**FIG. 1. ACIDES framework.** (A) We consider directed evolution (DE) experiments, where protein variants are screened over multiple rounds, and massive NGS datasets are collected. (B) From the obtained count data, we estimate a score (selectivity) for each variant. The higher the score, the better the variant for the task. Each score is estimated together with its 95%-confidence interval (CI). (C) Sorting the scores of all variants in descending order, we obtain a variant rank (naive rank). Due to statistical errors in the scores, the obtained rank is biased in general. To correct for this, using *in-silico* simulations based on the CIs of the scores, we re-estimate the rank with 95%-CI (corrected rank). (D) From the obtained corrected rank, we compute Rank Robustness (RR). RR represents the percentage of the top 50 variants identified in the naive rank that also appear in the top 50 of the corrected rank. (E,F) Examples of rank graphs for two synthetic datasets with different depths of NGS (per round) and numbers of unique variants (respectively, E:  $10^7$ ,  $5 \times 10^4$ ; F:  $10^6$  and  $10^6$ ). The true rank is shown as red crosses. In both cases, most red crosses are within the 95%-CI of the corrected rank. (G) RR for the two synthetic datasets. Note that RR multiplied by 50 (E:  $\sim 45.3$ ; F:  $\sim 24.6$ ) roughly provides the number of the correct top-50 sequences, which are 46 and 23, respectively. (See Fig.s S3 and S4 for more systematic comparison).

[23, 30, 37], yeast two-hybrid [23] and small nucleolar RNA [29] experiments. It is possible to apply ACIDES either *a posteriori* over data collected previously, or along

the course of the experiment as soon as the NGS data become available. The latter strategy allows for monitoring the selection convergence *on the fly*, and to understand when the experiment can be ended. In this way, ACIDES can be integrated into protein engineering pipelines as well as studies of protein function using mutagenesis. The tutorial for using ACIDES, along with an executable code in Python, will be available in GitHub upon publication of this manuscript.

## RESULTS

The first step of ACIDES estimates the selectivity of each individual variant present in the dataset (Fig. 1B) and its 95% confidence interval (95%-CI). In this study the term selectivity means the rate at which each variant increases its concentration with respect to the others. More precisely, we assume an exponential growth as  $\rho_i^{t+\Delta t} \sim \rho_i^t \exp(a^i \Delta t)$ , where  $\rho_i^t$  is the concentration of variant  $i$  at time  $t$ , and  $a^i$  is its selectivity. Compared with previous methods [19–24, 27–31, 38], our approach combines a robust inference framework (maximum likelihood estimation) with a better quantification of the NGS sampling noise [32–34]. For this scope, our approach benefits from a negative binomial distribution [39–42] (Fig. S1) in which the variance of the noise is overdispersed and grows as  $\lambda + \lambda^{2-\alpha}/\beta$ . Here  $\lambda$  is the expected mean count, and  $\alpha, \beta$  are parameters to be inferred (Materials and Methods). Using novel data from a plasmid library, we observed that our negative binomial model realizes a 50- to 70-fold improvement over the Poisson model in the predictive ability of the NGS sampling noise (Fig. S1). The second step of ACIDES uses the estimation of the selectivities and their statistical errors to rank the variants. The rank obtained by sorting the selectivities in descending order (naive rank) is biased due to statistical fluctuations of the selectivities. We correct this bias using *in-silico* simulations (Fig. 1C). The third and last step of ACIDES uses simulations to quantify a Rank Robustness (RR), a measure of the quality of the selection convergence (Fig. 1D). Specifically, RR is the ratio at which the top-50 variants in the naive rank are correctly identified (Materials and Methods). RR ranges from 0 to 1: a low value points out that the variants have not been selected enough, and therefore calls for the necessity to perform more rounds, deeper NGS sampling or possibly more replicates. Conversely, a large value confirms that the selection has properly converged, and suggests that the experiment can be ended without performing additional experimental steps.

Before focusing on experimental data, we apply ACIDES to two synthetic datasets (Materials and Methods) describing two opposite scenarios (See Fig.s S3 and S4 for more systematic comparison): data-rich case (more NGS reads with fewer unique variants) and data-poor case (less NGS reads with more considered variants). In the data-rich case, we first verify that our method

131 reaches high performance in recovering the ground-truth  
 132 values of the selectivities ( $R^2 \simeq 0.92$ , Fig.S3) in a teacher-  
 133 student setting. In this first case, selection convergence is  
 134 reached and the different variants can be robustly ranked  
 135 (Fig. 1E). In the data-poor case, instead, CI-bars are  
 136 large and the ranking is uncertain (Fig. 1F). Consistently,  
 137 the estimated RRs are high and low for, respectively, the  
 138 data-rich and -poor examples (Fig. 1G). Note that, once  
 139 multiplied by 50, RR roughly provides the number of the  
 140 correct top-50 variants in both cases (caption of Fig. 1G).  
 141 Furthermore, we observe that most true rank values (red  
 142 crosses) fall within the 95%-CI in both examples. These  
 143 observations show that our approach can quantify statisti-  
 144 cal errors even in the data-poor regime (See Fig. S4 for  
 145 more systematic comparison).

### 146 Analysis of directed evolution and deep mutational 147 scanning experiments

148 In order to showcase ACIDES, we apply it to sev-  
 149 eral screening datasets, where various proteins (and one  
 150 RNA molecule) are screened using different experimen-  
 151 tal techniques (Table I). Specifically, we consider three  
 152 phage-display screening experiments targeting different  
 153 proteins, such as the breast cancer type 1 susceptibil-  
 154 ity protein (BRCA1) for Data-A, human yes-associated  
 155 protein 65 (hYAP65) for Data-F and immunoglobulin  
 156 heavy chain (IgH) for Data-G, two *in-vivo* DEs of adeno-  
 157 associated virus type 2 (AAV2) vectors targeting canine  
 158 eyes for Data-C and murine lungs for Data-D, a mul-  
 159 tiplexed yeast two-hybrid assay targeting BRCA1 for  
 160 Data-B and a yeast competitive growth screen measur-  
 161 ing the fitness of mutant U3 gene for Data-E. For each  
 162 of these experiments, we rank variants (naive rank) and  
 163 compute the confidence interval of their ranks (corrected  
 164 rank in Fig. 2A-G). The degree of convergence of the  
 165 selection is quantified by RR (2H). When technical repli-  
 166 cates are available (Data-A and Data-B), we compute RR  
 167 over all of them and obtained consistent results (shown  
 168 by the small error-bars in Fig. 2H).

169 We classify the observed RRs into three groups de-  
 170 pending on the quality of the selection convergence: high  
 171 (Data-A and Data-B), intermediate (Data-F and Data-  
 172 G), and low (Data-C, Data-D and Data-E) convergence  
 173 groups. The high group seems to behave similarly to the  
 174 data-rich synthetic data in Fig. 1E. Consistently, RR,  
 175 NGS depth and the number of unique variants are indeed  
 176 of the same order (Table I). In these cases, the obtained  
 177 naive rank is robust, as indicated by the value of RR  
 178 ( $RR > 0.8$ ). In the intermediate group, the value of RR  
 179 ranges between 0.6 and 0.8. The experimental techniques  
 180 used in these datasets are similar to those in the high  
 181 group, but the NGS depths (or the numbers of unique  
 182 variants) are smaller (or larger), which could be the rea-  
 183 son why they result in lower RRs. The low group suffers  
 184 from the noise in the data. In Data-C and Data-D, the  
 185 numbers of unique variants are lower than those of Data-

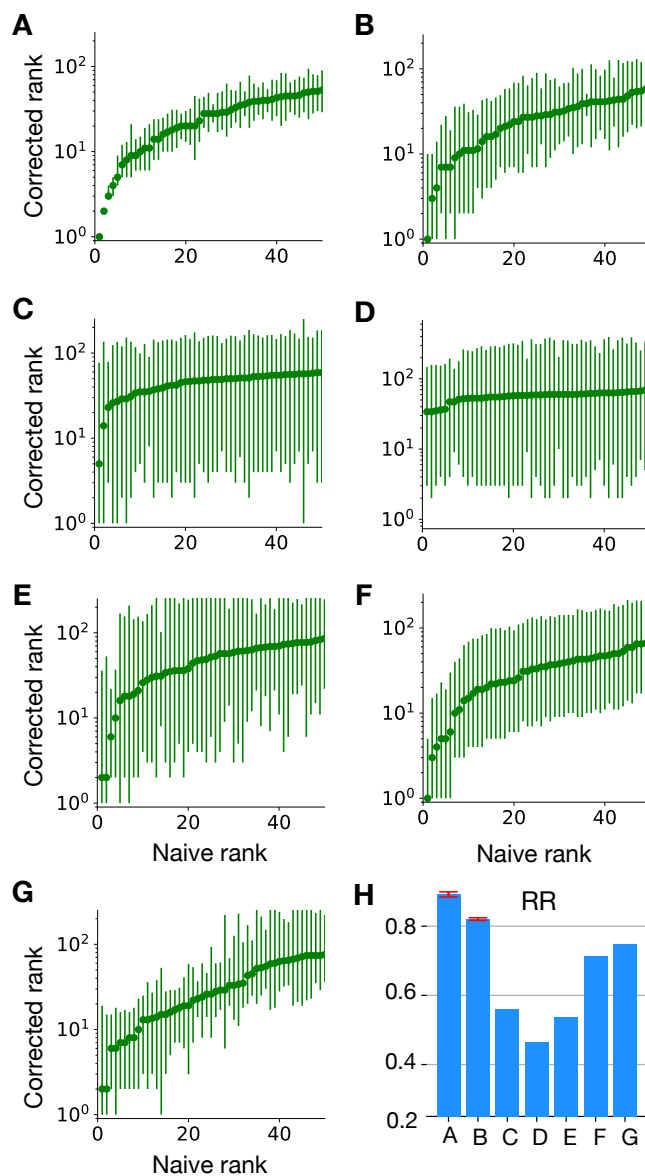


FIG. 2. Rank graph for various experimental datasets. The panel labels A-G correspond to the experiments listed in table I. (H) Rank robustness (RR) for each experiment. When technical replicates are available (Data-A and -B), the mean and standard deviation are shown.

186 A and Data-B, which would normally help these datasets  
 187 with having higher RR, given the same NGS depth. As  
 188 this is not the case, we see that some experiments are *in-*  
 189 *trinsically* more difficult than the others, *i.e.*, *in-vivo* DE  
 190 (Data-C and Data-D) and RNA based screening (Data-  
 191 E) will result in lower RRs than the other experiments if  
 192 the NGS depth and number of variants are similar.

193 In datasets with low RRs, some variants seem to per-  
 194 form better than the others, but the difference between  
 195 their scores is marginal compared with their statistical  
 196 errors. This means that we cannot distinguish if the ob-  
 197 tained variants are selected because of their ability to per-

TABLE I. Next generation sequencing datasets of directed evolution experiments

Label	Experiment	Target	Time-points	Reads/round	# of variants	Replicates	Ref
A	Phage display	BRCA1	T0 → T5	13.6 M	35 k	2 x 3	Starita 2015 [23]
B	Yeast two-hybrid	BRCA1	T0 → T3	13.5 M	27 k	2 x 3	Starita 2015 [23]
C	<i>in-vivo</i> DE (dog eye)	AAV2-7mer	T1 → T5	17 M	5 k	1 x 1	Byrne 2018 [36]
D	<i>in-vivo</i> DE (murine lung)	AAV2-7mer	T0 → T5	6.2 M	0.5 k	1 x 1	Korbelin 2016 [13]
E	Yeast competitive growth	U3 snoRNA	T0 → T4	8 M	24 k	2 x 1	Puchta 2016 [29]
F	Phage display	hYAP65 WW	T0 → T3	5 M	470 k	2 x 1	Araya 2012 [30]
G	Phage display	Ab IgH	T1 → T3	0.1 M	29 k	1 x 1	Boyer 2016 [37]

List and properties of experiments considered in this study. First column introduces dataset label and corresponds to the panels of Fig. 2. Reads/round corresponds to the average NGS counts per time points. # of variants is the number of unique variants that is detected in the NGS at least once during whole experiments. In Replicates,  $x \times y$  means that there are  $x$  replicates that do not share the same initial library, each of which has  $y$  technical replicates (that shares the same initial library).

198 form the task (fitness) or just there due to noise. In these  
 199 cases, experimentalists have two possibilities: (i) based  
 200 on the noisy identified variants, perform further tests in  
 201 addition to DE [13, 14], as for example, study infective  
 202 ability of viral vectors using single-cell RNA-seq [43]. Or  
 203 (ii) increase the quality of the datasets, by performing  
 204 further selection rounds, increasing NGS depths, or repli-  
 205 cating the experiments under the same conditions. This  
 206 second possibility is explored in the next section. Over-  
 207 all our rank-analysis of the different experiments shows  
 208 how our approach can provide an overview of the selec-  
 209 tion convergence, informing about the state of the exper-  
 210 iment and eventually pointing out the necessity of more  
 211 experimental efforts.

### 212 Integration into the experimental pipeline

213 Noise in experimental data can be reduced by perform-  
 214 ing additional selection rounds involving experiments,  
 215 but in general these are expensive, time-consuming and,  
 216 in case of experiments involving animal use, ethically  
 217 problematic [35]. For these reasons, it is important to  
 218 choose accurately the number of rounds and the NGS  
 219 depth. For this scope, ACIDES can be integrated into  
 220 experimental pipelines to obtain an overview on how RR  
 221 depends on these factors. This is to help experimentalists  
 222 with making informed decisions about additional exper-  
 223 imental efforts.

224 ACIDES can estimate RR after each selection round  
 225 (or any time new data become available). This allows  
 226 us to examine the data’s behavior and to quantify the  
 227 degree of convergence in terms of the selection rounds.  
 228 Similarly, for each round, ACIDES can be run on down-  
 229 sampled NGS data to compute RR with smaller NGS  
 230 depth (Materials and Methods). Using these two tech-  
 231 niques, we monitor the need for more selection rounds or  
 232 deeper NGS: a slow increase of RR (or no change in RR)  
 233 upon improving data-quality implies that convergence is  
 234 reached and suggests that the experiment can be ended.  
 235 If, on the other hand, RR increases rapidly when improv-  
 236 ing the rounds and/or NGS depth, it is probably worth  
 237 making further experimental efforts.

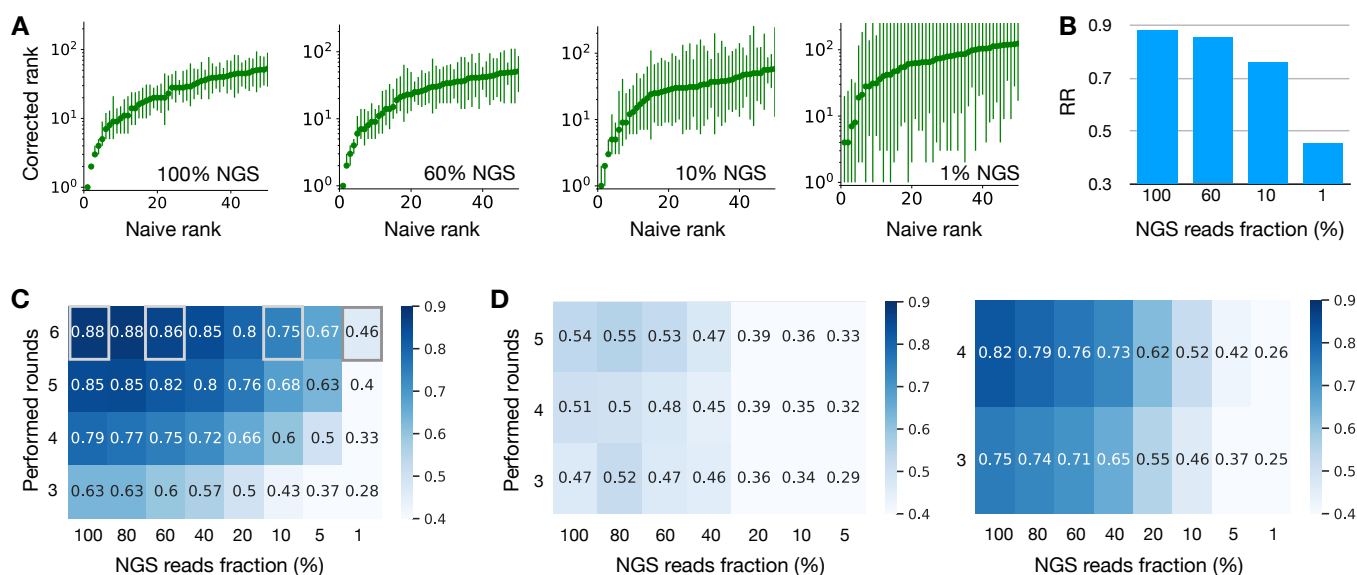
238 In order to showcase our approach, we study how RR

239 depends on the number of screening rounds and NGS  
 240 depth in previous experiments. We start by measuring  
 241 RR in Data-A for different NGS depths. 95%-CI on  
 242 corrected ranks gets larger as the NGS depth becomes  
 243 smaller (Fig.3A). At 1% NGS depth, the variant ordering  
 244 seems largely unreliable: RR is smaller than 0.5 (Fig.3B).  
 245 Importantly, RR does not decrease smoothly as the NGS  
 246 depth decreases, but it remains roughly constant at the  
 247 beginning, and falls only at a very small NGS depth.  
 248 This result suggests that the actual NGS depth of this  
 249 experiment largely exceeds what was necessary (10% of  
 250 the depth would have been sufficient). Next, we quan-  
 251 tify how RR depends on both the number of performed  
 252 rounds and NGS depth (Fig.3C). RR grows from 0.28 (3  
 253 performed rounds with 1% NGS depth) to 0.88 (6 per-  
 254 formed rounds with 100% NGS depth). Saturation of RR  
 255 seems to be observed for  $RR > 0.7$ , which corresponds to  
 256 5 performed rounds with the NGS depth larger than 20%,  
 257 or 4 performed rounds with the NGS depth larger than  
 258 40%. This again indicates that the experiment could have  
 259 been stopped earlier (less rounds and/or lower sequence  
 260 coverage) without much affecting the outcome. Note that  
 261 different datasets show different behaviors. For Data-E  
 262 more selection rounds with a higher number of NGS reads  
 263 is expected to improve RR, while for Data-B they seem  
 264 to have just reached the saturation point (Fig.3D).

265 Overall these results show how our approach can be  
 266 implemented along experimental pipelines. By estimat-  
 267 ing RR while collecting new data, we can understand if  
 268 we should continue/stop adding more rounds or increas-  
 269 ing NGS depth. This could avoid unnecessary, costly  
 270 and time-consuming experimental efforts. Similar analy-  
 271 ses can be done on the number of replicate experiments  
 272 (Fig.S6).

### 273 Comparison with previous work

274 We compare the performance of ACIDES with *En-*  
 275 *rich2*, the state of the art for estimating variant scores  
 276 (selectivities) [31]. *Enrich2* is based on a weighted lin-  
 277 ear fitting of the log-count change along rounds, and the  
 278 first step of ACIDES should be seen as an upgrade for  
 279 this fitting. In order to compare these two approaches



**FIG. 3. How the rank robustness depends on the experimental protocol.** (A) Rank graphs for different NGS depths in Data-A (Table I). Different NGS-depth data are generated using downsampling (Materials and Methods).  $x\%$  means the dataset where the number of NGS reads per round is reduced to  $x\%$  (100% is the original dataset). (B) RR for the rank graphs in the panel A. Note that RR is higher than 0.7 even with the 10% NGS-depth. (C) The heat map showing RR for various NGS depths and performed rounds in Data-A. RR is larger than 0.7 for the data with (i) the 4 performed rounds with the NGS depth larger than or equal to 40% or with (ii) the 5 performed rounds with the NGS depth larger than or equal to 20%. This indicates that the data quality was already high with less experimental efforts. The four grey squares correspond to the four rank graphs in the panel A, respectively. (D) The same graphs as the panel C, but for different datasets. Data-E is used in the left panel, where RR is low and more NGS and/or screening rounds would be useful. Data-B is used in the right panel, where RR takes high values and seems to saturate in NGS depths. Further experimental efforts would probably not be necessary in this dataset.

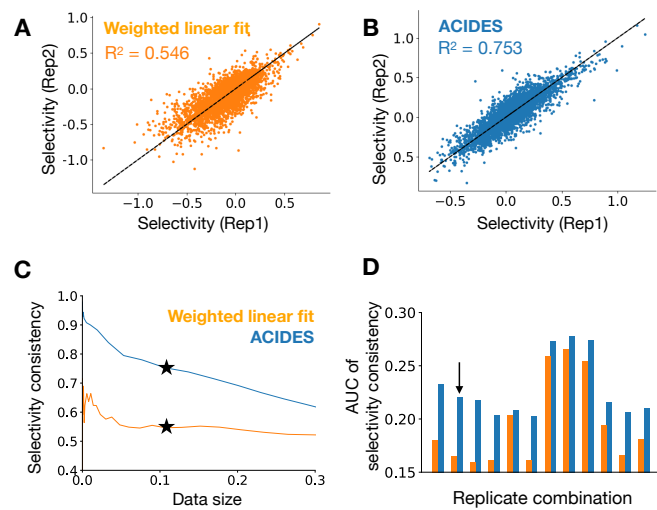
280 we take advantage of replicate datasets. We first investigate if the scores in each method are consistent over  
 281 replicates. For this, we plot the scores obtained from  
 282 one replicate against the scores obtained from the other  
 283 (Fig. 4A, B). The correlation between replicates is estimated  
 284 using the coefficient of determination ( $R^2$ ). The  
 285 correlation quantifies the quality of the method, as higher  
 286 (or lower) correlations imply that the estimated scores  
 287 are more (or less) robust and fewer (or more) replicates  
 288 are needed to obtain reliable results. The figure shows  
 289 that ACIDES outperforms *Enrich2*. Next, we test how  
 290 the comparison depends on the data quality. To this  
 291 goal, we systematically select a set of variants based on  
 292 the magnitude of predicted statistical score-errors (Materials  
 293 and Methods). (Smaller/larger sets include variants  
 294 with smaller/larger predicted statistical errors.) For  
 295 each set, we measure the correlation between two replicates  
 296 as in (Fig. 4A, B), and plot it as a function of the set  
 297 size (Fig. 4C). We observe ACIDES's correlation becomes  
 298 more dominant as the set size decreases, suggesting the  
 299 better quality of both the estimated scores and statistical  
 300 errors. In order to generalize these results, we perform  
 301 the comparison for all possible 12 pairs of technical replicates  
 302 in Data-A and Data-B (Table I). In all cases our  
 303 approach outperforms the competitor (Fig. 4D). We also  
 304 perform an additional test to quantify the consistency of

306 the predicted statistical errors (Supp. Fig. S7).

## DISCUSSION

308 In this work we have presented ACIDES, a method to  
 309 quantify DE and DMS selectivities (fitness), rank variants  
 310 with accurate credibility scores and measure the degree  
 311 of experimental convergence. ACIDES can be used  
 312 *on the fly* to offer an overview of the progress of selection  
 313 experiments, which would help experimentalists  
 314 with making informed decisions on whether new experimental  
 315 efforts are needed. In this way, ACIDES can save significant  
 316 experimental time and resources. We have applied ACIDES  
 317 to several DE and DMS datasets where a number of different  
 318 target proteins and one set of target RNA molecules have  
 319 been screened using different experimental protocols. The  
 320 heterogeneity of these datasets shows that ACIDES is a  
 321 method of general use, applicable to many different experiments.

322 The first step of ACIDES estimates the score (selectivity)  
 323 of each observed variant. This is a necessary step, and  
 324 several alternative methods have been proposed in the past.  
 325 In many applications, such scores are computed as the  
 326 variant enrichment that is defined as the logarithmic ratio  
 327 between the variant frequencies in the last and second to  
 328 last round [13] or between the last  
 329



**FIG. 4. Comparison of our approach with the state of the art.** Using technical replicates in Data-A, we compare ACIDES with a weighted linear least squares method (*Enrich2*) [31]. For both methods (*Enrich2* (A) and ACIDES (B)), the inferred selectivities from one replicate are plotted against the selectivities in the other replicate. The coefficient of determination ( $R^2$ ), which quantifies the consistency between two replicates, is also shown. (C) We next examine how the comparison in the panels A and B depend on data quality. We consider a set of variants in which the estimated statistical errors (Materials and Methods) are smaller than a given threshold. Varying this threshold, sets of variants are systematically selected, where larger/smaller sets include variants with larger/smaller estimated statistical errors. For each set, we estimate  $R^2$  between two replicates, and plot it as a function of the set size. The panel A and B correspond to the stars ★ in C (data size 0.11). (D) In order to test both methods more systematically, we perform the same analysis (as those in the panels A-C) for all possible 12 combinations of technical replicates in Data-A and Data-B. We define the area under curve of  $R^2$  (in the panel C) and plot it for these combinations (D). Our method systematically outperforms the weighted linear fitting method. The replicate combination used for the panels A-C is indicated by the arrow in the panel D.

and first round [14, 19, 20, 22, 27, 38]. These approaches thus make use of data from only two rounds and disregard all the others. For this reason, this strategy is suboptimal and may lead to noisy score estimations. A more sophisticated approach that uses all the data consists in inferring the slope of a linear line fitted to the log-frequencies of variants over all the screening rounds/time points [23, 24, 28, 30]. This method gives the same importance to log-frequencies in all the rounds. Yet as variant counts in the first rounds are typically small and noisy, assuming the same weight on them could result in an overfitting. To fix this effect, *Enrich2* [31] uses the variance of the count data - estimated via a Poisson distribution assumption - as the weights in a linear least squares fitting. ACIDES' first step comes with a three-fold improvement over this last approach. First, in-

stead of relying on the linear least squares fitting, we estimate the score by log-likelihood maximization. A major improvement happens for variants whose log-frequencies do not grow linearly with the rounds, and a simple linear weighted fit may struggle in identifying the correct slope. This is particularly visible in the bulk variants with intermediate scores (Fig.4 A, B). Secondly, instead of a simple exponential growth of the counts, we included a softmax non-linear function (Materials and Methods), where the denominator is inferred from data [44]. This change improves the score estimation when the wildtype (if any) and/or few variants have a large fraction of the total counts and bend the exponential growth of the log-frequencies. Lastly, ACIDES uses a negative binomial distribution to model the count variability [39–42]. This distribution accounts for the large dispersion of next generation sequencing data [32–34] far better than the Poisson distribution (Fig. S1). Additionally, the negative binomial loss in the likelihood maximization allows us to better estimate statistical errors for the inferred scores. Thanks to all these improvements, our approach realizes a more robust and accurate estimation of the variant scores and outperforms the previous method (Fig. 4).

In case of noisy data, the estimated scores of variants come with statistical errors. This means that the rank obtained from the scores (naive rank in our figures) is in general biased: top ranked variants are overvalued, and vice-versa. This simple statistical effect was not taken into account in previous analyses related to DE and DMS experiments. The second step of ACIDES uses a bootstrap method to account for the bias and recover both the corrected rank and its 95%-CI. The deviation between this 95%-CI and the naive rank shows us how much we can trust the naive rank. To quantify it, as a third step of ACIDES, we introduce RR that describes how many of the top-50 variants in the naive rank are correct. RR measures how stable and robust are the ranks of the variant selectivities. As such, it quantifies the degree of convergence of the experimental selection, providing an insightful overview of the state of the experiment.

Although ACIDES demonstrates advantages over the other methods, it has several limitations that may be addressed in the future. First of all, ACIDES does not account for changes in the selection pressure over rounds. This can potentially be included, but has not been done here, as the selection pressure is constant in most datasets we analyzed in this article. Second, ACIDES uses a negative binomial model to describe the dispersion of count data by assuming that the count variance depends only on the frequency of the variant. Although this assumption proves useful to describe NGS count errors (Fig. S1) and is used elsewhere [42], it is possible that dispersions induced by a sequence-dependent procedure, such as error-prone PCR [14, 36, 45], may not be taken into account by our method (Note that Data-C includes an error-prone PCR after the third round of selections, indicating that the estimated results for Data-C may contain biases). We would need to analyze more

404 data from DE experiments using error-prone libraries to  
 405 address this question. Third, statistical errors due to the  
 406 replicates that do not share the same initial library can-  
 407 not be described by ACIDES, provided that the model  
 408 is only trained on a single series of screening rounds. To  
 409 account for this, we would need a framework that gener-  
 410 alizes ACIDES for different sources of variability.

411 Finally, using machine learning techniques, several  
 412 studies have aimed at estimating selectivities from the  
 413 amino-acid sequences of variants. Most of these meth-  
 414 ods rely on supervised algorithms, which are trained to  
 415 predict the selectivity (output) from the sequence of a  
 416 variant (input) [45–53]. Because the performance of these  
 417 methods depends on how the selectivity is estimated from  
 418 data, ACIDES can potentially be incorporated in their  
 419 pipelines to improve the overall performance. We leave  
 420 such analysis for future developments. Other methods  
 421 use instead unsupervised approaches to predict selectiv-  
 422 ities from the sequences of variants [44, 54–57]. Even if  
 423 these methods do not use any sequence scores for their  
 424 training, they often use it to validate and/or test the  
 425 model. Our approach would therefore be useful also in  
 426 these cases.

## 427 METHOD

### 428 Library preparation for Fig.S1

429 To demonstrate that our negative binomial likelihood  
 430 approach outperforms the Poisson counterpart, we con-  
 431 ducted the following experiment: We inserted random 21  
 432 nucleotide oligomers into a RepCap plasmid containing  
 433 adenoassociated virus 2 (AAV2) cap gene using previ-  
 434 ously described methods [58]. The plasmid library ob-  
 435 tained was deep sequenced following generation of ampli-  
 436 cons corresponding to the 7mer insertion region. Since  
 437 the 21 nucleotides are randomly and independently gen-  
 438 erated, we can use a position weight matrix model to pre-  
 439 dict the frequency of each variant in the sample. Based  
 440 on this property, the performance of the two models are  
 441 examined as shown in Fig.S1.

### 442 Model

443 We propose ACIDES for analyzing selection data in  
 444 DE and DMS. Here the mathematical model is described  
 445 in detail. For a given series of samples over screening  
 446 rounds, we perform NGS and denote by  $n_t^i$  the obtained  
 447 count of the  $i$ -th variant ( $i = 1, 2, \dots, M$ ) at round (time-  
 448 point)  $t \in T$ . We denote by  $N_t$  the total count  $N_t =$   
 449  $\sum_i n_t^i$  at  $t$ . For each sample, we define  $\rho_t^i$  as the *expected*  
 450 value of frequency of the  $i$ -th variant at  $t$ . (Note that  
 451 “expected” means that  $\rho_t^i$  itself does not fluctuate due to  
 452 the noise in the experiment.) For each variant, an initial  
 453 frequency  $\rho_0^i$  and a growth rate  $a^i$  are assigned, by which

454 the expected frequency is computed as

$$\rho_{t+\Delta t}^i = C_t \rho_t^i \exp(a^i \Delta t), \quad (1)$$

455 where  $\Delta t$  is the round- (or time-) difference between two  
 456 consecutive NGSs.  $C_t$  is a normalization constant, de-  
 457 fined as  $C_t = 1/\sum_i [\rho_t^i \exp(a^i \Delta t)]$ . We call this model  
 458 ((1)) an exponential model.

459 We use a negative binomial distribution  $\text{NB}(n_t^i|\lambda, r)$   
 460 with two parameters  $\lambda$  and  $r$  to model the noise distri-  
 461 bution of counts  $n_t^i$ . Here  $\lambda$  is the expected value of count  
 462  $n_t^i$  given as  $N_t \rho_t^i$ , while  $r$  is the dispersion parameter that  
 463 describes the deviation of the negative binomial distribu-  
 464 tion from the Poisson distribution. (The negative bino-  
 465 mial distribution is a generalization of the Poisson distri-  
 466 bution with a variance equal to  $\lambda(1 + \lambda/r)$ : the Poisson  
 467 distribution is recovered in the large  $r$  limit.) Here, based  
 468 on Fig. S1 and [42], we assume  $r$  is a power-law function  
 469 of  $\lambda$ :  $r(\lambda) = \beta \lambda^\alpha$  (with  $\alpha, \beta > 0$ ), where  $\alpha$  and  $\beta$  are  
 470 parameters that are common for all the variants in the  
 471 experiment. (The variance is thus  $\lambda + \lambda^{2-\alpha}/\beta$ .) Model  
 472 parameters  $\alpha, \beta$  as well as  $\rho_0^i, a^i$  ( $i = 1, 2, \dots, M$ ) are in-  
 473 ferred from the count data  $n_t^i$  ( $i = 1, 2, \dots, M, t \in T$ ) by  
 474 maximizing the following likelihood function:

$$L(\alpha, \beta, (\rho_0^i)_{i=1}^M, (a^i)_{i=1}^M) = \prod_{i,t} \text{NB} \left( n_t^i | \rho_t^i N_t, \beta (\rho_t^i N_t)^\alpha \right). \quad (2)$$

475 The 95%-CIs of the estimated parameters are computed  
 476 from the curvature of the log-likelihood function at the  
 477 maximum.

### 478 Synthetic data

479 Synthetic count data  $n_t^i$  ( $i = 1, 2, \dots, M, t \in T$ ) are  
 480 generated from the model ((2)) for a given parame-  
 481 ter set  $\alpha, \beta, \rho_0^i, a^i$  ( $i = 1, 2, \dots, M$ ). For Fig.1, we use  
 482  $\alpha, \beta = 0.69, 0.8$  with  $(a^i, \log \rho_0^i)$  generated from the nor-  
 483 mal distribution with the expected values  $(-1, 1)$  and the  
 484 standard deviations  $(0.25, 1)$ .  $(M, N_t)$  are  $(5 \times 10^4, 10^7)$   
 485 for the data-rich case (Fig.1E) and  $(10^6, 10^6)$  for the data-  
 486 poor one (Fig.1F).

### 487 Model inference

488 To maximize the likelihood function, we develop a two-  
 489 step algorithm. The first step infers  $(\rho_0^i, a^i)$ , while the  
 490 second  $(\alpha, \beta)$  and then we iterate the two steps until  
 491 convergence is reached. All inferences are done with a  
 492 gradient descent algorithm, and to reach convergence 10-  
 493 30 iterations are usually sufficient. The first step is itself  
 494 iterative, and loops between the inference of  $(\rho_0^i, a^i)$  and  
 495  $C_t$  by treating  $C_t$  as a parameter. Here we also introduce  
 496 a gauge choice because of the redundancy between  $\rho_0^i, a^i$   
 497 and  $C_t$  (the caption of Fig. S2 for more details). In the  
 498 second step, the inference of  $(\alpha, \beta)$  with a straightforward

549 gradient method produces a bias (Fig. S2E). In order  
500 to correct this, at each iteration the algorithm adopts  
501 a teacher-student framework, runs a simulation of the  
502 count data with the current parameters to obtain an es-  
503 timation of the bias, which is then used to correct the  
504 real inference and update the parameters.

505 In order to reduce computational time and to increase  
506 the stability of the algorithm, we first run the inference  
507 algorithm on a subset of variants to estimate  $\alpha, \beta$ . We  
508 then compute  $(\rho_0^i, a^i)$  of the excluded variants using the  
509 estimated  $\alpha, \beta$ . For this subset, we use the variants that  
510 satisfy the following two criterions: (i) their counts are  
511 larger than 0 more than twice in the selection rounds  
512 and (ii) whose total NGS count (as summed over all the  
513 rounds) is above a threshold. We set this threshold to 100  
514 for all the datasets except for Data-C -D, where 10000  
515 is used. This is because the noise in these experiments  
516 is larger than the others. Results are stable by changing  
517 the threshold value (Fig. S2F).

### 518 Simulated rank and rank robustness (RR)

519 Using the standard deviations  $\delta a^i$  ( $i = 1, \dots, M$ ) of es-  
520 timated scores  $a^i$ , we discard the variants with higher  
521 estimated errors. We keep 5000 variants for further anal-  
522 ysis and denote by  $A$  their indices. We then rearrange  
523 the variant index in  $A$  in descending order of  $a^i$  to de-  
524 fine a *naive rank* (the  $x$ -axis of Fig 2A-G). To obtain a  
525 *corrected rank* (the  $y$ -axis of Fig 2A-G), we first gener-  
526 ate synthetic scores using the normal distribution with  
527 the expected value  $(a^i)_{i \in A}$  and the standard deviation  
528  $(\delta a^i)_{i \in A}$ . Based on the generated scores, we rearrange  
529 the variant index in descending order and define a syn-  
530 thetic naive rank. Repeating this estimation 3000 times,  
531 we then compute the median and 95%-CI of the obtained  
532 synthetic naive ranks. This 95%-CI is defined as the cor-  
533 rected rank.

534 To estimate RR, we compare the top-50 variants in  
535 the naive rank and each synthetic naive rank. We count  
536 the number of overlaps between them and average it over  
537 the 3000 estimations. RR is computed by dividing the  
538 obtained overlap by 50.

### 539 NGS-Downsampling for Fig.3

540 To obtain downsampled count data  $\tilde{n}_t^i$  ( $i = 1, 2, \dots, M$ ,  
541  $t \in T$ ) by a factor  $\epsilon$ , we sample synthetic data from the  
542 likelihood function ((2)) with a reduced number of the to-  
543 tal counts  $\epsilon N_t$  ( $t \in T$ ) and with the estimated parameters  
544  $\rho_0^i, a^i, \alpha, \beta$  ( $i = 1, 2, \dots, M$ ). To obtain a downsampled  
545 RR in Fig. 3, we first re-estimate  $a^i$  ( $i = 1, 2, \dots, M$ ) from

546  $\tilde{n}_t^i$  ( $i = 1, 2, \dots, M, t \in T$ ) using the values of  $(\alpha, \beta)$  that  
547 are already known, and then perform the estimation of  
548 RR described above. Using the synthetic data, we show  
549 that this downsampling method captures well the RR of  
550 actual NGS-read-reduced data (Fig. S5).

### 551 Pre-processing of Data-C and Data-D

552 In their original datasets, Data-C and -D contain a  
553 large number of variants whose total counts are very low  
554 (but not zero). In order to speed up the analysis and  
555 make the analysis more robust we removed the variants  
556 whose total counts are smaller than 1000 (Data-C) and  
557 than 100 (Data-D). The NGS depth and the number of  
558 unique variants shown in Table I are after this prepro-  
559 cessing.

### 560 DATA AVAILABILITY

561 All data analyzed in this article (Table I) are publicly  
562 available except for the random-peptide inserted library  
563 used for Fig. S1. This library will be deposited in a public  
564 database upon publication of this article.

### 565 CODE AVAILABILITY

566 A Python implementation of ACIDES will be available  
567 on GitHub upon publication of this article.

### 568 ACKNOWLEDGMENTS

569 The authors would like to thank A. Rubin and D.  
570 Fowler for kindly providing them with datasets (Data-A  
571 and Data-B) and a working code for enrich2. The authors  
572 also would like to thank L. C. Byrne and T. Mora for  
573 useful comments and discussions, M. Desrosiers and C.  
574 Robert for their technical assistance with the production  
575 of plasmids and viral vectors, and O. Marre for facilitat-  
576 ing and initiating the collaboration and helpful discus-  
577 sions. This work was supported by ERC Starting Grant  
578 (REGENETHER 639888), European Research Council  
579 (ERC) Horizon 2020 Framework Programme Project:  
580 863214 – NEUROPA, UNADEV, the Institut National  
581 de la Santé et de la Recherche Médicale (INSERM), Sor-  
582 bonne Université, The Foundation Fighting Blindness,  
583 Agence National de Recherche (ANR) RHU Light4Deaf,  
584 LabEx LIFESENSES (ANR-10-LABX-65), IHU FORE-  
585 SIGHT (ANR-18-IAHU-01), and JSPS KAKENHI Grant  
586 Number 22K17994.

587 [1] Frances H Arnold. Design by directed evolution. *Ac-*  
588 *counts of chemical research* **31**, 125 (1998).

589 [2] Philip A Romero and Frances H Arnold. Exploring pro-  
590 tein fitness landscapes by directed evolution. *Nature re-*



- views *Molecular cell biology* **10**, 866 (2009).
- [3] Michael S Packer and David R Liu. Methods for the directed evolution of proteins. *Nature Reviews Genetics* **16**, 379 (2015).
- [4] Keqin Chen and Frances H Arnold. Tuning the activity of an enzyme for unusual environments: sequential random mutagenesis of subtilisin E for catalysis in dimethylformamide. *Proceedings of the National Academy of Sciences* **90**, 5618 (1993).
- [5] Nicholas J Turner. Directed evolution drives the next generation of biocatalysts. *Nature chemical biology* **5**, 567 (2009).
- [6] Olga Khersonsky Tawfik and Dan S. Enzyme Promiscuity: A Mechanistic and Evolutionary Perspective. *Annual Review of Biochemistry* **79**, 471 (2010). PMID: 20235827.
- [7] Robert E. Hawkins, Stephen J. Russell, and Greg Winter. Selection of phage antibodies by binding affinity: Mimicking affinity maturation. *Journal of Molecular Biology* **226**, 889 (1992).
- [8] Eric T Boder, Katarina S Midelfort, and K Dane Wittrup. Directed evolution of antibody fragments with monovalent femtomolar antigen-binding affinity. *Proceedings of the National Academy of Sciences* **97**, 10701 (2000).
- [9] Luca Perabo, Hildegard Büning, David M Kofler, Martin U Ried, Anne Girod, Clemens M Wendtner, Jörg Enssle, and Michael Hallek. In vitro selection of viral vectors with modified tropism: the adeno-associated virus display. *Molecular Therapy* **8**, 151 (2003).
- [10] Narendra Maheshri, James T Koerber, Brian K Kaspar, and David V Schaffer. Directed evolution of adeno-associated virus yields enhanced gene delivery vectors. *Nature biotechnology* **24**, 198 (2006).
- [11] Stefan Michelfelder and Martin Trepel. Adeno-associated viral vectors and their redirection to cell-type specific receptors. *Advances in genetics* **67**, 29 (2009).
- [12] Deniz Dalkara, Leah C. Byrne, Ryan R. Klimczak, Meike Visel, Lu Yin, William H. Merigan, John G. Flannery, and David V. Schaffer. In Vivo-Directed Evolution of a New Adeno-Associated Virus for Therapeutic Outer Retinal Gene Delivery from the Vitreous. *Science Translational Medicine* **5**, 189ra76 (2013).
- [13] Jakob Körbelin, Timo Sieber, Stefan Michelfelder, Lars Lunding, Elmar Spies, Agnes Hunger, Malik Alawi, Kleopatra Rapti, Daniela Indenbirken, Oliver J Müller, Renata Pasqualini, Wadiah Arap, Jürgen A Kleinschmidt, and Martin Trepel. Pulmonary Targeting of Adeno-associated Viral Vectors by Next-generation Sequencing-guided Screening of Random Capsid Displayed Peptide Libraries. *Molecular Therapy* **24**, 1050 (2016).
- [14] Leah C Byrne, Timothy P Day, Meike Visel, Jennifer A Strazzeri, Cécile Fortuny, Deniz Dalkara, William H Merigan, David V Schaffer, and John G Flannery. In vivo-directed evolution of adeno-associated virus in the primate retina. *JCI insight* **5** (2020).
- [15] Mohammadsharif Tabeordbar, Kim A. Lagerborg, Alexandra Stanton, Emily M. King, Simon Ye, Liana Tellez, Allison Krunnusz, Sahar Tavakoli, Jeffrey J. Widrick, Kathleen A. Messemer, Emily C. Troiano, Behzad Moghadaszadeh, Bryan L. Peacker, Krystynne A. Leacock, Naftali Horwitz, Alan H. Beggs, Amy J. Wagers, and Pardis C. Sabeti. Directed evolution of a family of AAV capsid variants enabling potent muscle-directed gene delivery across species. *Cell* **184**, 4919 (2021).
- [16] <https://www.nobelprize.org/prizes/chemistry/2018/summary/>.
- [17] Sam Behjati and Patrick S Tarpey. What is next generation sequencing? *Archives of Disease in Childhood - Education and Practice* **98**, 236 (2013).
- [18] Shawn E. Levy and Richard M. Myers. Advancements in Next-Generation Sequencing. *Annual Review of Genomics and Human Genetics* **17**, 95 (2016). PMID: 27362342.
- [19] Douglas M Fowler, Carlos L Araya, Sarel J Fleishman, Elizabeth H Kellogg, Jason J Stephany, David Baker, and Stanley Fields. High-resolution mapping of protein sequence-function relationships. *Nature methods* **7**, 741 (2010).
- [20] Ryan T. Hietpas, Jeffrey D. Jensen, and Daniel N. A. Bolon. Experimental illumination of a fitness landscape. *Proceedings of the National Academy of Sciences* **108**, 7896 (2011).
- [21] Douglas M Fowler and Stanley Fields. Deep mutational scanning: a new style of protein science. *Nature methods* **11**, 801 (2014).
- [22] Alexandre Melnikov, Peter Rogov, Li Wang, Andreas Gnirke, and Tarjei S Mikkelsen. Comprehensive mutational scanning of a kinase in vivo reveals substrate-dependent fitness landscapes. *Nucleic acids research* **42**, e112 (2014).
- [23] Lea M Starita, David L Young, Muhtadi Islam, Jacob O Kitzman, Justin Gullingsrud, Ronald J Hause, Douglas M Fowler, Jeffrey D Parvin, Jay Shendure, and Stanley Fields. Massively Parallel Functional Analysis of BRCA1 RING Domain Variants. *Genetics* **200**, 413 (2015).
- [24] Sebastian Matuszewski, Marcel E Hildebrandt, Ana-Hermina Ghenu, Jeffrey D Jensen, and Claudia Bank. A Statistical Guide to the Design of Deep Mutational Scanning Experiments. *Genetics* **204**, 77 (2016).
- [25] Nathan J Rollins, Kelly P Brock, Frank J Poelwijk, Michael A Stiffler, Nicholas P Gauthier, Chris Sander, and Debora S Marks. Inferring protein 3D structure from deep mutation scans. *Nature genetics* **51**, 1170 (2019).
- [26] Jörn M Schmiedel and Ben Lehner. Determining protein structures using deep mutagenesis. *Nature genetics* **51**, 1177 (2019).
- [27] Rupali P Patwardhan, Choli Lee, Oren Litvin, David L Young, Dana Pe'er, and Jay Shendure. High-resolution analysis of DNA regulatory elements by synthetic saturation mutagenesis. *Nature biotechnology* **27**, 1173 (2009).
- [28] Matthew S Rich, Celia Payen, Alan F Rubin, Giang T Ong, Monica R Sanchez, Nozomu Yachie, Maitreya J Dunham, and Stanley Fields. Comprehensive Analysis of the SUL1 Promoter of *Saccharomyces cerevisiae*. *Genetics* **203**, 191 (2016).
- [29] Olga Puchta, Botond Cseke, Hubert Czaja, David Tollervey, Guido Sanguinetti, and Grzegorz Kudla. Network of epistatic interactions within a yeast snoRNA. *Science* **352**, 840 (2016).
- [30] Carlos L. Araya, Douglas M. Fowler, Wentao Chen, Ike Muniez, Jeffery W. Kelly, and Stanley Fields. A fundamental protein property, thermodynamic stability, revealed solely from large-scale measurements of protein function. *Proceedings of the National Academy of Sciences* **109**, 16858 (2012).
- [31] Alan F Rubin, Hannah Gelman, Nathan Lucas, Sandra M Bajjalieh, Anthony T Papenfuss, Terence P Speed, and Douglas M Fowler. A statistical framework for an-

- alyzing deep mutational scanning data. *Genome biology* **18**, 1 (2017).
- [32] Justus M. Kebschull and Anthony M. Zador. Sources of PCR-induced distortions in high-throughput sequencing data sets. *Nucleic Acids Research* **43**, e143 (2015).
- [33] Katharine Best, Theres Oakes, James M Heather, John Shawe-Taylor, and Benny Chain. Computational analysis of stochastic heterogeneity in PCR amplification efficiency revealed by single molecule barcoding. *Scientific reports* **5**, 1 (2015).
- [34] Vladimir Potapov and Jennifer L Ong. Examining sources of error in PCR by single-molecule sequencing. *PloS one* **12**, e0169774 (2017).
- [35] Simon Festing and Robin Wilkinson. The ethics of animal research. *EMBO reports* **8**, 526 (2007).
- [36] Leah Byrne, Timothy Day, Meike Visel, Deniz Dalkara, Valerie Dufour, Felipe Pompeo Marinho, William Merigan, Gustavo Aguirre, William Beltran, David Schaffer, and John Flannery. Directed Evolution of AAV for Efficient Gene Delivery to Canine and Primate Retina - Raw counts of variants from deep sequencing. *Dryad, Dataset* (2018). <https://doi.org/10.6078/D1895R>.
- [37] Sébastien Boyer, Dipanwita Biswas, Ananda Kumar Soshee, Natale Scaramozzino, Clément Nizak, and Olivier Rivoire. Hierarchy and extremes in selections from pools of randomized proteins. *Proceedings of the National Academy of Sciences* **113**, 3482 (2016).
- [38] Douglas M. Fowler, Carlos L. Araya, Wayne Gerard, and Stanley Fields. Enrich: software for analysis of protein function by enrichment and depletion of variants. *Bioinformatics* **27**, 3430 (2011).
- [39] Simon Anders and Wolfgang Huber. Differential expression analysis for sequence count data. *Nature Precedings* pages 1 (2010).
- [40] Davis J. McCarthy, Yunshun Chen, and Gordon K. Smyth. Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Research* **40**, 4288 (2012).
- [41] Michael I Love, Wolfgang Huber, and Simon Anders. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome biology* **15**, 1 (2014).
- [42] Maximilian Puelma Touzel, Aleksandra M Walczak, and Thierry Mora. Inferring the immune response from repertoire sequencing. *PLOS Computational Biology* **16**, e1007873 (2020).
- [43] Bilge E Öztürk, Molly E Johnson, Michael Kleyman, Serhan Turunç, Jing He, Sara Jabalameli, Zhouhuan Xi, Meike Visel, Valérie L Dufour, Simone Iwabe, Luis Felipe L Pompeo Marinho, Gustavo D Aguirre, José-Alain Sahel, David V Schaffer, Andreas R Pfenning, John G Flannery, William A Beltran, William R Stauffer, and Leah C Byrne. scAAVengr, a transcriptome-based pipeline for quantitative ranking of engineered AAVs with single-cell resolution. *eLife* **10**, e64175 (2021).
- [44] Jorge Fernandez-de Cossio-Diaz, Guido Uguzzoni, and Andrea Pagnani. Unsupervised Inference of Protein Fitness Landscape from Deep Mutational Scan. *Molecular Biology and Evolution* **38**, 318 (2020).
- [45] Zachary Wu, S. B. Jennifer Kan, Russell D. Lewis, Bruce J. Wittmann, and Frances H. Arnold. Machine learning-assisted directed protein evolution with combinatorial libraries. *Proceedings of the National Academy of Sciences* **116**, 8852 (2019).
- [46] Richard J Fox, S Christopher Davis, Emily C Mundorff, Lisa M Newman, Vesna Gavrilovic, Steven K Ma, Loleta M Chung, Charlene Ching, Sarena Tam, Sheela Muley, et al. Improving catalytic function by ProSAR-driven enzyme evolution. *Nature biotechnology* **25**, 338 (2007).
- [47] Philip A. Romero, Andreas Krause, and Frances H. Arnold. Navigating the protein fitness landscape with Gaussian processes. *Proceedings of the National Academy of Sciences* **110**, E193 (2013).
- [48] Jakub Otwinowski, David M. McCandlish, and Joshua B. Plotkin. Inferring the shape of global epistasis. *Proceedings of the National Academy of Sciences* **115**, E7550 (2018).
- [49] Frédéric Cadet, Nicolas Fontaine, Guangyue Li, Joaquin Sanchis, Matthieu Ng Fuk Chong, Rudy Pandjaitan, Iyann Vetrivel, Bernard Offmann, and Manfred T Reetz. A machine learning approach for reliable prediction of amino acid interactions and its application in the directed evolution of enantioselective enzymes. *Scientific reports* **8**, 1 (2018).
- [50] Claire N Bedbrook, Kevin K Yang, J Elliott Robinson, Elisha D Mackey, Viviana Gradinaru, and Frances H Arnold. Machine learning-guided channelrhodopsin engineering enables minimally invasive optogenetics. *Nature methods* **16**, 1176 (2019).
- [51] Kevin K Yang, Zachary Wu, and Frances H Arnold. Machine-learning-guided directed evolution for protein engineering. *Nature methods* **16**, 687 (2019).
- [52] Yuting Xu, Deeptak Verma, Robert P Sheridan, Andy Liaw, Junshui Ma, Nicholas M Marshall, John McIntosh, Edward C Sherer, Vladimir Svetnik, and Jennifer M Johnston. Deep dive into machine learning models for protein engineering. *Journal of chemical information and modeling* **60**, 2773 (2020).
- [53] Drew H Bryant, Ali Bashir, Sam Sinai, Nina K Jain, Pierce J Ogden, Patrick F Riley, George M Church, Lucy J Colwell, and Eric D Kelsic. Deep diversification of an AAV capsid protein by machine learning. *Nature Biotechnology* **39**, 691 (2021).
- [54] Claudia Bank, Ryan T Hietpas, Alex Wong, Daniel N Bolon, and Jeffrey D Jensen. A Bayesian MCMC approach to assess the complete distribution of fitness effects of new mutations: uncovering the potential for adaptive walks in challenging environments. *Genetics* **196**, 841 (2014).
- [55] Jakub Otwinowski. Biophysical Inference of Epistasis and the Effects of Mutations on Protein Stability and Function. *Molecular Biology and Evolution* **35**, 2345 (2018).
- [56] Luca Sesta, Guido Uguzzoni, Jorge Fernandez-de Cossio-Diaz, and Andrea Pagnani. AMaLa: Analysis of Directed Evolution Experiments via Annealed Mutational Approximated Landscape. *International Journal of Molecular Sciences* **22** (2021).
- [57] Andrea Di Gioacchino, Jonah Procyk, Marco Molari, John S. Schreck, Yu Zhou, Yan Liu, Rémi Monasson, Simona Cocco, and Petr Šulc. Generative and interpretable machine learning for aptamer design and analysis of in vitro sequence selection. *PLOS Computational Biology* **18**, 1 (2022).
- [58] James T Koerber, Narendra Maheshri, Brian K Kaspar, and David V Schaffer. Construction of diverse adeno-associated viral libraries for directed evolution of en-

846 hanced gene delivery vehicles. *Nature protocols* **1**, 701  
847 (2006).

## Supplementary information for Uncursing winner's curse: on-line monitoring of directed evolution convergence

Takahiro Nemoto,<sup>1,2,3,\*</sup> Tommaso Ocari,<sup>1</sup> Arthur Planul,<sup>1</sup> Muge Tekinsoy,<sup>1</sup> Emilia A. Zin,<sup>1</sup> Deniz Dalkara,<sup>1,†</sup> and Ulisse Ferrari<sup>1,‡</sup>

<sup>1</sup>*Institut de la Vision, Sorbonne Université, INSERM, CNRS, 17 rue Moreau, 75012, Paris, France*

<sup>2</sup>*Graduate School of Informatics, Kyoto University,  
Yoshida Hon-machi, Sakyo-ku, Kyoto, 606-8501, Japan*

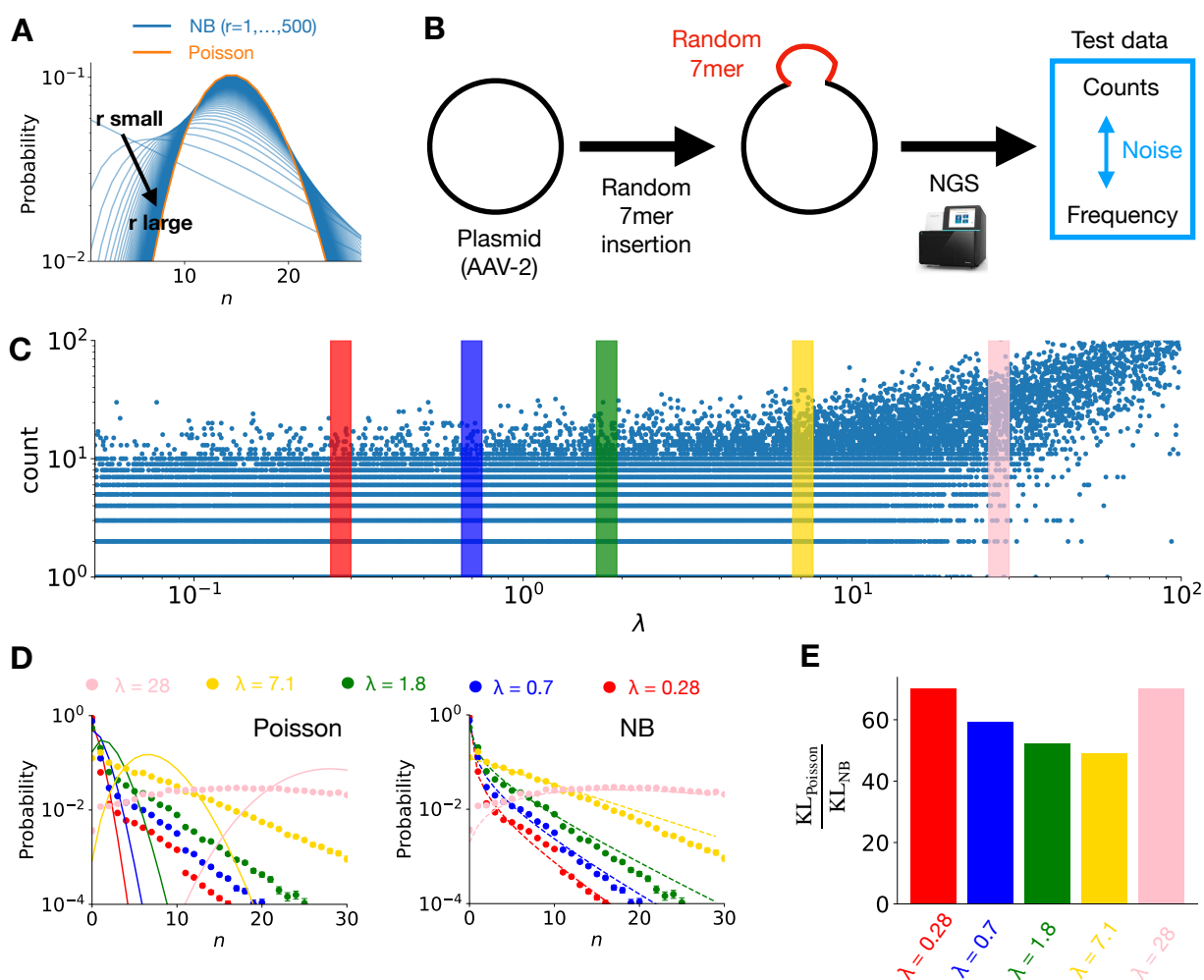
<sup>3</sup>*Premium Research Institute for Human Metaverse Medicine (WPI-PRIME),  
Osaka University, Suita, Osaka 565-0871, Japan*

---

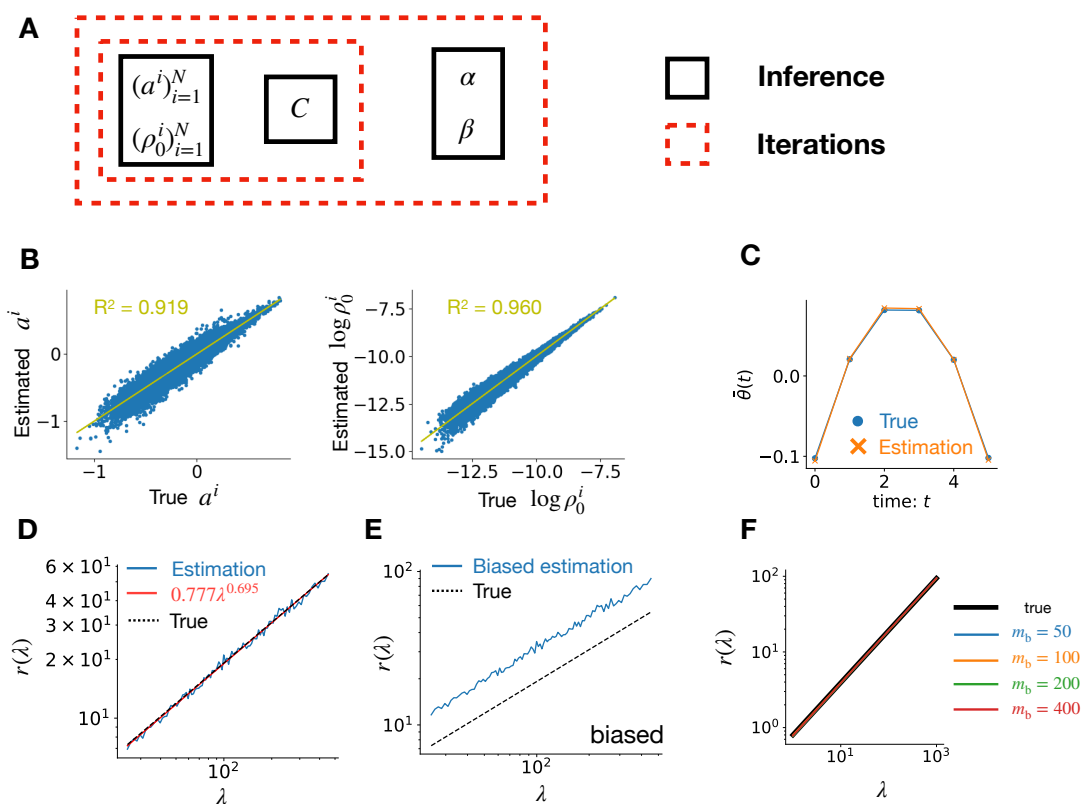
\* nemoto.takahiro.prime@osaka-u.ac.jp

† deniz.dalkara@inserm.fr

‡ ulisse.ferrari@inserm.fr



**FIG. S1. Negative binomial model accounts for NGS count noise better than Poisson model.** (A) The poisson distribution (orange) and the negative binomial distribution (tableau blue) with the expected value  $\lambda = 15$ . The dispersion parameter  $r$  for the negative binomial distribution is set to 1, 2, ..., 500. The negative binomial distribution generalizes the Poisson distribution, allowing for large variances by decreasing  $r$ . It converges to the Poisson distribution in the large  $r$  limit. (B) In order to test the predictive ability of the negative binomial distribution, we performed the following experiment. Using a random peptide (of size 21 corresponding to a 7mer) as a barcode, we first barcoded a plasmid extracted from adeno-associated virus 2 (AAV2) wild type. The obtained 7-mer inserted library was then sent to NGS facility and the corresponding barcoded region was sequenced. Since these 21 nucleotides of barcode are randomly and independently generated, we can use a position weight matrix model to predict the frequency of each variant in the sample. Comparing the predicted frequency with the actual NGS reads, we investigate the noise distribution of NGS counts. (C) The graph showing the obtained counts ( $n$ ) against the predicted frequencies multiplied by the total NGS reads ( $\lambda$ ), where each point corresponds to a variant. We observe that the counts are largely dispersed. (D) The probability distribution of counts for a fixed value of  $\lambda$  together with the model predictions by Poisson distribution (left) and by the negative binomial distribution (right). The probability distribution is estimated in the following way: (i) picking up all the variants within 5 different colored rectangles in the panel (C). (ii) Using the variants corresponding to each rectangle, we then make a histogram of counts, which is plotted in the panel (D) as dots. For the Poisson prediction, we simply use the Poisson distribution with the mean  $\lambda = 0.28, 0.7, 1.8, 7.1, 28$ . For the negative binomial prediction, for each value of  $\lambda$ , we infer the dispersion parameter  $r$  via a maximum likelihood inference and fitted a power-law function  $r = \beta\lambda^\alpha$  to the obtained estimations. The result is  $r = 0.21\lambda^{0.744}$ . Using this relation, the negative binomial distribution is then plotted for each  $\lambda = 0.28, 0.7, 1.8, 7.1, 28$ . (E) Comparison of predictive abilities between the poisson model and the negative binomial model. To quantify the predictive ability of each model, we use Kullback-Leibler divergence (KL) defined as  $KL = \sum_n P_{\text{data}}(n) \log(P_{\text{data}}(n)/P_{\text{model}}(n))$ . The ratio between KL for the poisson model and KL for the negative binomial model is plotted for each value of  $\lambda$ .  $KL_{\text{Poisson}}$  itself is 0.37, 0.59, 0.87, 1.59, 3.29 for  $\lambda = 0.28, 0.7, 1.8, 7.1, 28$ , while  $KL_{\text{NB}}$  is 0.0053, 0.0099, 0.017, 0.032, 0.047.



**FIG. S2. Inference algorithm and synthetic teacher-student examples.** (A) graphical illustration of the inference algorithm with its double loop structure. The internal loop accounts for the parameters of the exponential model (see Materials and Methods) and iterates between the inference of  $(a^i, \rho_0^i)_{i=1}^M$  and  $C$ , while the external loop iterates between this internal loop and the inference of the negative binomial parameters  $(\alpha, \beta)$ . For the internal loop, we first infer  $(a^i, \rho_0^i)_{i=1}^M$  for a fixed  $C, \alpha, \beta$  by maximizing the likelihood function (2).  $C^*$  is then calculated from the obtained  $(a^i, \rho_0^i)_{i=1}^M$  via  $C^* = 1 / \sum_i \rho_0^i \exp(a^i t)$ . The linearly increasing part of  $\log C^*$  is next subtracted as  $\log C^* - (x^* t + y^*)$ , where  $x^*, y^* = \operatorname{argmax}_{x,y} \sum_t [\log C - xt - y]^2$  (fixing a gauge). The obtained quantity is  $\log C$  for the next iteration. For the external loop, to obtain  $(\alpha, \beta)$ , we first infer the dispersion parameter  $r$  for a list of different values of  $\lambda$ . To do so, for a fixed value of  $\lambda$ , we select a set of index  $i$  and the time  $t$  by the condition  $\lambda < \rho_0^i N_t < \lambda + \epsilon$  (with a small parameter  $\epsilon$ ). Only using these  $i$  and  $t$ , we then maximize the likelihood function (2) and determine the value of  $r(\lambda)$ . After obtaining the function  $r(\lambda)$  for several values of  $\lambda$ , we fit a linear function  $r(\lambda) = \alpha \lambda^\beta$  to it and determine  $\alpha$  and  $\beta$ . (B,C,D) Ground truth comparisons for the synthetic data-rich dataset (Fig. 1E) after 30th iterations demonstrates that the algorithm can recover the generating parameters. In the panel B, we plot the estimated parameters  $(a^i)_{i=1}^M$  (left) and  $(\log \rho_0^i)_{i=1}^M$  (right) against their ground truths. The coefficient of determination  $R^2$  is also shown. In the panel C, the normalization coefficient  $(\bar{\theta}(t) \equiv \log C)$  is plotted together with its ground truth. In the panel D, the estimated  $r(\lambda)$  with a fitted line  $\beta \lambda^\alpha$  and its ground truth are shown. (E) While estimating  $r(\lambda)$  for a fixed  $\lambda$ , maximizing the likelihood function (2) results in a biased estimation as shown in the panel E. For fixing this, we generate a synthetic data probe using the current estimation of  $\alpha_0$  and  $\beta_0$  with  $(\rho_0^i)_{i,t}$  and use it to unbiased the  $r$  estimation. More precisely, denoting the biased estimation by  $r_{\text{bias}}(\lambda)$  (panel E) and the estimation of the  $r$  in the probe by  $r_1(\lambda)$ , the unbiased estimator plotted in the panel D is obtained as  $r(\lambda) = r_{\text{bias}}(\lambda) \beta_0 \lambda^{\alpha_0} / r_1(\lambda)$ . (F) To determine  $\alpha$  and  $\beta$ , we use only representative variants, as described in Materials and Methods. For the representative variants, we select the variants whose total counts are larger than a threshold value  $m_b$ . In (F), by using the synthetic data, we show that the inference results ( $\alpha$  and  $\beta$ ) are robust against the change of this parameter  $m_b$ .

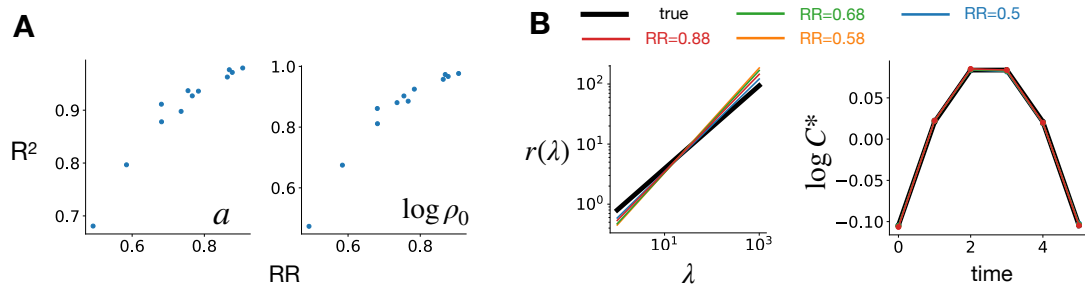


FIG. S3. **RR offers a proxy for the accuracy of the estimated scores.** Here we use synthetic datasets, ranging from data-poor to data-rich regimes, to show that the empirical quantity RR correlates with our ability to recover the true values of scores  $a^i$  and initial frequencies  $\rho_0^i$ . For this, we generate synthetic datasets with different values of total NGS reads  $N$  and number of variants  $M$ :  $(N, M) = (10^7, 5 \times 10^5); (8 \times 10^6, 5 \times 10^5); (6 \times 10^6, 5 \times 10^5); (4 \times 10^6, 5 \times 10^5); (2 \times 10^6, 5 \times 10^5); (10^6, 5 \times 10^5); (10^7, 10^6); (8 \times 10^6, 10^6); (4 \times 10^6, 10^6); (2 \times 10^6, 10^6); (10^6, 10^6)$ . In these datasets, the parameters  $(\alpha, \beta)$  are the same as those used in Fig.1E, F. (A) Coefficients of determination  $R^2$  between inferred and true values for  $a^i$  (left) and  $\rho_0^i$  (right) are plotted against RR. This demonstrates the correlation between  $R^2$  and RR. (B) The inference of  $r(\lambda)$  and  $C^*$  is robust across all synthetic datasets.

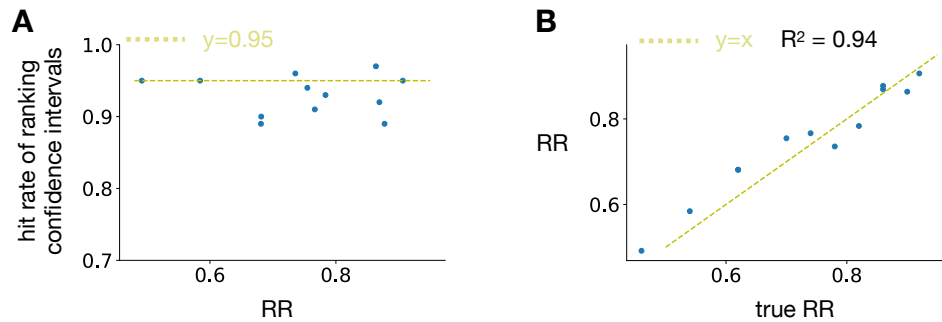


FIG. S4. **Ranking CI and rank robustness are accurately estimated also in the data-poor regime.** A) For all the synthetic datasets introduced in Fig. S3, hit rates of the confidence interval of ranking graphs are plotted against their RR, where the hit rate is defined as the number of true ranking (red crosses in Fig.1E, F for example) that are within the 95%-confidence intervals (green lines in Fig.1E, F), divided by 50. The hit rates fluctuate around 0.95 as expected, demonstrating the quality of our estimation of ranking CI. B) Rank robustness (RR) estimated from inferred parameters is plotted against the true value (ground truth) for the synthetic datasets. The obtained high coefficient of determination shows that ACIDES can estimate RR also in the data-poor regime.

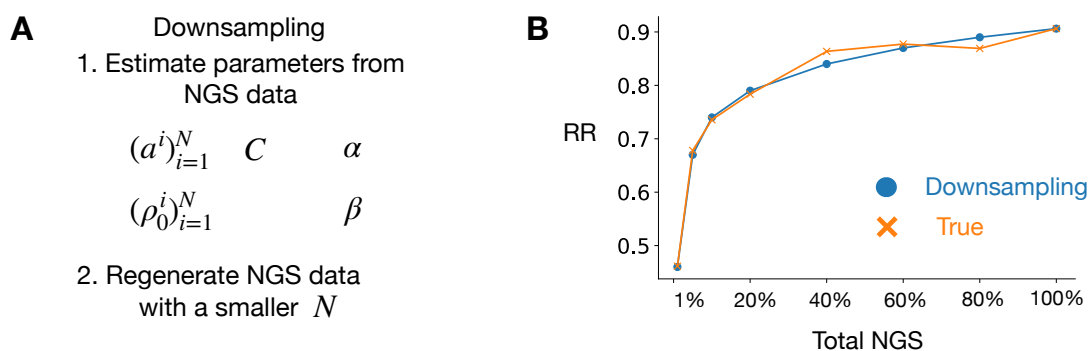


FIG. S5. **Downsampling NGS counts.** (A) For a given dataset of a DE experiment, we downsample the data, *i.e.*, create a synthetic dataset that corresponds to the dataset of the same DE experiments but with smaller values of the total NGS reads. For this, we first estimate the parameters of ACIDES from the dataset. We then generate synthetic NGS data using the likelihood function (2) with smaller numbers of the total NGS reads. For example, if we downsample the data to 40%, we set  $N_t$  to be  $0.4N_t$ . We then estimate RR using ACIDES for this downsampled dataset without reinferring  $\alpha, \beta$ . (B) We show the validity of this down sampling method on the synthetic data. The parameters for the synthetic data are  $(N_t, M) = (10^7, 5 \times 10^5)$ ,  $(\alpha, \beta) = (0.69, 0.8)$  and  $(a^i, \log \rho_0^i)_{i=1}^M$  generated from the normal distribution with the expected values  $(-1, 1)$  and the standard deviation  $(0.25, 1)$ . We plot RRs obtained from this downsampling method (blue circles) and from a standard sampling method (orange crosses) as a function of the total number of NGS (where 100% means the original data). Here the standard sampling method means using ACIDES directly on the dataset with the total number of NGS  $0.01xN_t$ , where  $x$  is the percentage of the total NGS reads ( $x$ -axis in the panel B). We observe that our downsampling method estimates well the RR.

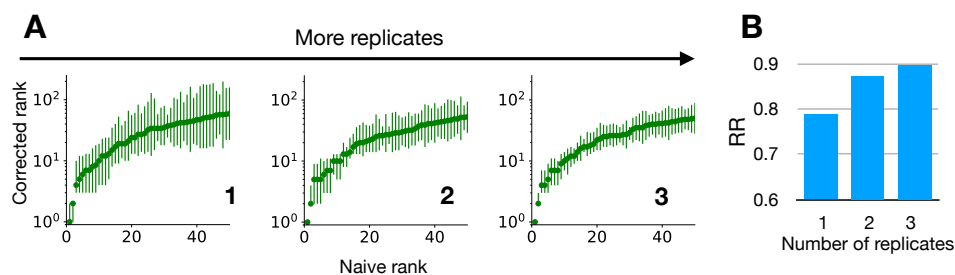


FIG. S6. **Multiple replicates can be combined to increase RR.** We use the first 4 rounds of Data-A (so that RR is relatively small for a single replicate), and perform ACIDES for each replicate. (A) Ranking graphs for one (left), two (middle), three (right) replicates. To combine variant scores of two replicates (denoted by  $a_1, a_2$  with standard deviation  $\delta a_1, \delta a_2$ ), we use  $(a_1 \delta a_2^2 + a_2 \delta a_1^2) / (\delta a_1^2 + \delta a_2^2)$  for the combined score and  $\sqrt{\delta a_1^2 \delta a_2^2} / (\delta a_1^2 + \delta a_2^2)$  for the combined standard deviation. (B) RR for the three cases.



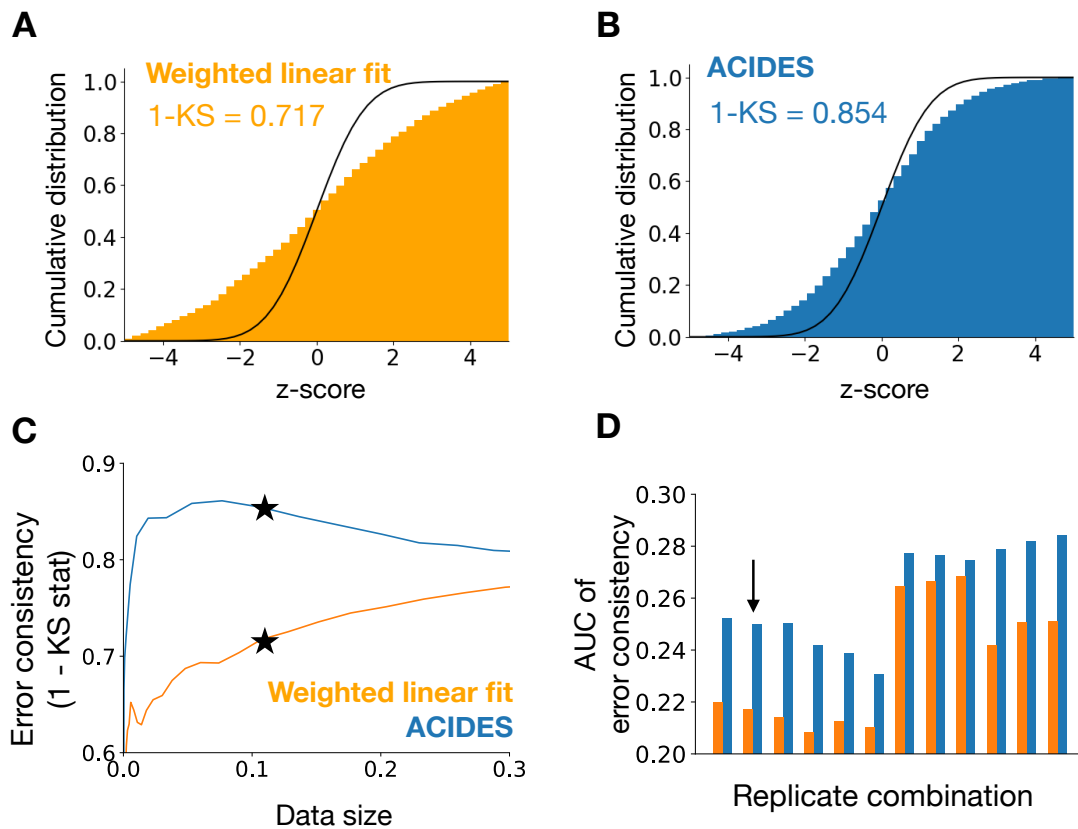


FIG. S7. **ACIDES outperforms previous methods in the estimation of score's statistical errors.** Using replicates in Data-A and Data-B (Table 1), we study the consistency of error bars in ACIDES and in *Enrich2*. Denoting by  $a_1, a_2$  the scores of a variant in replicate 1 and 2 (similarly by  $\delta a_1, \delta a_2$  the standard deviations of the scores), we study the following quantity  $z = (a_1 - a_2) / \sqrt{\delta a_1^2 + \delta a_2^2}$  and compute the histogram of this quantity over different variants. Under the assumption that both scores are distributed following the normal distribution, this obtained histogram is approximated by the standard normal distribution. (A,B) Cumulative distributions of the histograms for *Enrich2* (A) and ACIDES (B) together with the cumulative standard normal distribution. We use 1 - Kolmogorov-Smirnov (KS) statistics (the maximum distance between two distributions measured in the y-direction) to quantify the distance between the histogram and the normal distribution. (C) To study 1 - KS more systematically, we introduce a threshold for the score statistical errors (Materials and Methods) by which we reduce the amount of data. For each fraction of the data, we estimate 1 - KS and plot them in the panel C. The stars in the panel correspond to panels A and B. (D) Finally, using all possible combinations of technical replicates in Data-A and Data-B, we compare ACIDES and *Enrich2*. We compute the area under curve (AUC) of 1 - KS (the panel C) for all combinations. ACIDES always shows better performance than *Enrich2*. The arrow in the panel D indicates the replicate combination used in the panels A-C.