

SCIENTIFIC REPORTS



OPEN

Economical analysis of saturation mutagenesis experiments

Carlos G. Acevedo-Rocha^{1,2,4,5}, Manfred T. Reetz^{1,2} & Yuval Nov³

Received: 16 February 2015

Accepted: 20 April 2015

Published: 20 July 2015

Saturation mutagenesis is a powerful technique for engineering proteins, metabolic pathways and genomes. In spite of its numerous applications, creating high-quality saturation mutagenesis libraries remains a challenge, as various experimental parameters influence in a complex manner the resulting diversity. We explore from the economical perspective various aspects of saturation mutagenesis library preparation: We introduce a cheaper and faster control for assessing library quality based on liquid media; analyze the role of primer purity and supplier in libraries with and without redundancy; compare library quality, yield, randomization efficiency, and annealing bias using traditional and emergent randomization schemes based on mixtures of mutagenic primers; and establish a methodology for choosing the most cost-effective randomization scheme given the screening costs and other experimental parameters. We show that by carefully considering these parameters, laboratory expenses can be significantly reduced.

Following the seminal work of Michael Smith concerning site-directed mutagenesis¹, saturation mutagenesis (SM) has emerged as an indispensable technique in molecular biology for introducing targeted sequence variations at virtually any DNA region. Areas of application include laboratory evolution of enzymes with enhanced activity, stability, and stereoselectivity^{2–6}, manipulation of binding properties of antibodies⁷ and transcription factors⁸. More recently, SM has been used to engineer promoters⁹, transcriptional enhancers¹⁰, ribosome binding sites¹¹, *cis*-regulatory elements and *trans*-acting factors¹², and protein-coding genes using emergent genome engineering tools¹³. Consequently, SM, especially when applied iteratively^{3–5,14,15}, has become an effective tool in protein engineering^{2,4,5}, and holds great potential in the directed evolution of metabolic pathways^{16–18} and genomes^{19–21}.

However, creating high-quality SM libraries is still an experimental challenge⁶: Factors such as the target DNA sequence, G+C content, melting-temperature of DNA duplex, randomization scheme, primer quality, length, and annealing, all play a role in the quality of the resulting library^{22–24}. While some of these factors are inherent to the DNA target sequence, and thus cannot be easily modified, other factors, such as the randomization scheme and primer quality, are decided upon by the experimenter. Yet these decisions carry their costs: more or higher-quality primers, for example, are obviously more expensive. To our knowledge, no study has investigated systematically the decision-making involved in SM experiments from the economical perspective.

In many directed evolution settings, genetic selection of the randomized variants is not an option, so that expensive screening must be undertaken instead²⁵. In these cases, the screening effort increases exponentially as a function of the number of randomized positions. To alleviate this burden, it is advantageous to minimize both redundancy and the frequency of premature stop codons. One step in this direction is to use NNK or NNS degenerate codons (where N = A/C/G/T, K = G/T and S = C/G) instead of NNN (but see²⁶); this way, when randomizing four positions, the theoretical ratio between the most

¹Department of Organic Synthesis, Max-Planck-Institut für Kohlenforschung, Mulheim, 45470, Germany.

²Department of Chemistry, Philipps-Universität Marburg, 35032, Germany. ³Department of Statistics, University of Haifa, Haifa, 31905, Israel. ⁴Prokaryotic Small RNA Biology Group, Max-Planck-Institut für terrestrische Mikrobiologie, Marburg, 35043, Germany. ⁵Landes-Offensive zur Entwicklung Wissenschaftlich-ökonomischer Exzellenz (LOEWE) Centre for Synthetic Microbiology (SYNMIKRO), Philipps-Universität Marburg, 35032, Germany. Correspondence and requests for materials should be addressed to M.T.R. (email: manfred-t.reetz@mpi-mail.mpg.de) or Y.N. (email: yuval@stat.haifa.ac.il)

represented variants (with Arg/Leu/Ser at each randomized position) and the least represented ones (with Met/Trp) drops from 1296:1 to 81:1, and the frequency of prematurely truncated variants drops from 17.5% to 11.9%. However, it is not possible to reduce further the redundancy using a single degenerate primer, while still encoding all 20 amino acids²⁷, so other strategies are called for.

Four such strategies have emerged. The first uses specially prepared mono²⁸, di²⁹, or trinucleotides³⁰ phosphoramidite solutions, or combinations thereof³¹, for synthesizing redundancy-free mutagenic primers. A second strategy, dubbed MAX, relies on the synthesis of a “template” and twenty “selection” oligonucleotides and their hybridization, and shows that library redundancy can be completely eliminated³²; this method has been recently extended to contiguous codons by using ProxiMAX³³. A third strategy builds on solid-phase combinatorial gene synthesis, and results in libraries with the highest diversity and essentially no bias²². Yet in spite of their clear advantages, these three strategies are currently too expensive to be used in routine SM experiments.

The fourth strategy, which we explore in this work, uses mixtures of ordinary primers to achieve zero or near-zero redundancy: Tang *et al.*³⁴ showed that a mixture of four primers per randomized position, NDT, VMA, ATG, and TGG (where D = A/G/T, V = A/C/G, M = A/C) at 12:6:1:1 molar ratio, results in a zero probability for premature stop codons and a perfectly uniform theoretical distribution of 1/20 for each of the 20 amino acids; this scheme is referred to hereafter as “Tang”. The “22c-trick” approach by Kille *et al.*³⁵ uses even fewer primers per randomized position – namely, the three primers NDT, VHG, and TGG (where H = A/C/T), at 12:9:1 molar ratio – and results also in a zero probability for premature stop codons, and in an almost uniform amino-acid distribution: 2/22 for Leu and Val, and 1/22 for each of the remaining 18 amino acids. More sophisticated variations of the primer mixing strategy have been published recently^{36–38}.

The last strategy (Tang and 22c-trick) is the most accessible to any molecular biology lab. However, this strategy requires more primers than NNK and NNS, which means that it is worthwhile to use it only when the screening cost is high enough, so that the reduction in screening effort due to the smaller library size more than offsets the increased library generation cost due to the larger number of primers. An even more extreme option is to discard randomization altogether, and to produce individually, via site-directed mutagenesis, all possible variants; the screening cost here would be minimal, as no duplicates are screened, but mutagenesis will be most expensive, because each variant requires a separate PCR reaction with dedicated primers.

Regardless of the mutagenesis approach chosen, it is expedient to assess experimentally, before engaging in costly library screening, whether the desired DNA diversity has indeed been created, to avoid wasting resources. For this reason we introduced some time ago the semi-quantitative Quick Quality Control (QQC) – a single, pooled sequencing reaction of the library’s plasmids, which provides an early indication for whether the randomization was successful^{24,35,39,40}. Recently, Stewart and colleagues have extended this simple and fast control by introducing the Q-value, an attractive and user-friendly quantitative score that averages the QQC results for estimating overall randomization efficiency²³.

The present study is concerned with all the above *economical* aspects associated with SM. Firstly, we present a cheaper and faster manner to perform the QQC. Secondly, we study how primer purity influences library quality. Thirdly, we assess and compare library quality using primers from three different suppliers. Fourthly, we compare library costs and quality using various redundant and non-redundant randomization schemes. Fifthly, we analyze the bias in favor of the wild-type (WT) codon and its implication on the expected library diversity. Finally, we present a methodology for resolving the trade-off between the cost of library generation and the cost of library screening, i.e., for choosing optimally the randomization scheme given the various experimental parameters.

Materials and Methods

Saturation mutagenesis experiments. The gene coding for P450_{BM3} was cloned into the pRSF-Duet-1 vector (Novagen, Merck, Darmstadt, Germany) under the control of the T7 promoter as described earlier⁴¹. The forward 5'-CGCTTTGATAAAAACTTAXXXCAAGCGCTTAAATTTGTACGTG-3' and reverse 5'-CACGTACAAATTTAAGCGCTTGZZZTAAGTTTTTATCAAAGCG-3' primers (where XXX and ZZZ respectively denote NNK or NNS or NDT+VHG+TGG or NDT+VMA+TGG+ATG and MNN or SNN or AHN+CDB+CCA or AHN+TKB+CCA+CAT) in both desalted (gel filtration column) and HPLC-purification grade, were purchased in Metabion (Martinsried, Germany; supplier 3) and delivered in water at a concentration of 100 μM. The same primers for NNK randomization in the same two purity grades and final concentration were acquired from Integrated DNA technologies (Iowa, US, supplier 1) and Invitrogen or Life Technologies (Carlsbad, US; supplier 2). All primers and primer mixtures were diluted to a concentration of 10 μM. The primer ratios for the 22c-Trick and Tang libraries were mixed as follows: 12 eq. NDT + 9 eq. VHG + 1 eq. TGG, and 12 eq. NDT + 6 eq. VHA + 1 eq. TGG + 1 eq. ATG, respectively. Prior to PCR, 1 μL (10 pmoles) of each forward and reverse primer or primer mixture were mixed with 20 ng of plasmid template in 8 μL of water, plus 10 μL of 1:1 KOD Hot Start Master Mix (Novagen, Merck, Darmstadt, Germany). The following PCR conditions were set: 95 °C for 2 min, followed by 25 cycles at 95 °C for 20 sec, 50 °C for 30 sec and 70 °C for 4 min, ending with 70 °C for 5 min and subsequent cooling. The samples were then dialysed against Millipore-Q water using 0.05 μm Millipore MF membrane filters (Millipore, Merck, Darmstadt, Germany) for 30 min. The methylated template was then treated with 1 μL (20 units) *DpnI* (New England Biolabs, Ipswich, US) for

at least 48 h at 37 °C in appropriate buffer and dialysed as before for 30 min prior to electroporation. From these dialysed samples, 5 µL were used to transform 50 µL of *Escherichia coli* BL21-Gold(DE3) cells using a “MicroPulser” electroporator (BioRad, Hercules, US). The cells were resuspended in 945 µL of LB medium without antibiotic and recovered for 1 h at 37 °C and 220 rpm in a 14 mL Falcon tube. Finally, the 1000 µL of culture were either plated on a big (15 cm diameter) agar plate with kanamycin (50 µg/mL) or left in the tube, to which 4000 µL of LB medium with 60 µg/mL kanamycin were added, resulting in 5 mL LB broth with 48 µg/mL kanamycin. The cultures were incubated at 37 °C with (liquid cultures) or without (solid agar plates) shaking at 220 rpm overnight (for about 16 h).

QQC, Q-values, and sequencing. The Quick Quality Control (QQC) based on solid agar plates was performed as previously described²². Briefly, upon growth, colonies were harvested using a Drigalsky spatula upon addition of 2 mL of pure water. The plasmid DNA was then extracted from the collected cells using the QIAprep Miniprep Kit (Qiagen, Hildesheim, Germany) and sequenced using a T7-coding primer. For QQC based on liquid cultures, the cells from the 14 mL Falcon culture tube were centrifuged for 10 min at 4000 rpm, followed by extraction of the pooled plasmids and sequencing.

To estimate the base distribution at each base position of the randomized codon, the peak heights of the four chromophores representing adenine (A), thymine (T), guanine (G) and cytosine (C) were obtained by moving the mouse cursor over the peaks using the freely available program “Chromas Lite” (Technelysium Pty Ltd, South Brisbane, Australia). The values were then converted to pie diagrams using Microsoft Excel (V. 14.3.2, Microsoft Corporation) and compared to the expected values corresponding to the randomization scheme: NNK or NNS expect 25% of A/T/G/C (N) at the first and second base positions, and 50% of G/T (K) or C/G (S) at the third position. 22c-Trick expects 27%(A)+19%(T)+27%(G)+27%(C) at the first position, 32%(A)+32%(T)+18%(G)+18%(C) at the second, and 55%(T)+45%(G) at the third³⁵. Finally, Tang expects 30%(A)+20%(T)+25%(G)+25%(C) at the first position, 35%(A)+25%(T)+20%(G)+20%(C) at the second, and 30%(A)+60%(T)+10%(G) at the third³⁴. The pooled Q-values, Q_{pool} , were calculated from the base distribution percentages, and the experimental Q-values, Q_{exptl} , were calculated from the sequencing results, as described elsewhere²³.

Upon analyzing library quality, the plasmid sent for QQC can be used to transform newly competent cells. In this case, the same strain *E. coli* BL21-DE3(Gold) was transformed as indicated above, using 50 ng of the plasmid isolated from the agar plate, but only 50 µL of the 1000 µL cell suspension were plated on big LB agar plates, followed the next day by single colony picking into 400 µL of LB medium with appropriate antibiotic (kanamycin at 50 µg/mL), using 2.2 mL 96-well plates (Thermo Scientific Fischer, Waltham, US). The plates were incubated at 37 °C and rotated at 220 rpm overnight; the next day an aliquot of 10 µL was transferred to 96-well solid agar plates with proper antibiotics and sent for plasmid extraction and sequencing using the T7-coding primer (GATC, Konstanz, Germany).

Statistical analysis. All statistical analyses were carried out using the R software (<http://www.r-project.org>). Correlation was estimated through Pearson’s product moment correlation coefficient, and Q-values were compared using Wilcoxon signed-rank test. Yields were compared using a two-sample proportion test. Point estimates and confidence intervals for P_{WT} were computed via a binomial test. The goodness-of-fit tests for annealing uniformity are based on Pearson’s chi-squared test statistic $\chi^2 = \sum_i (O_i - E_i)^2 / E_i$. Because of the low number of observations (sequencing results) relative to the number of categories (possible codons), the distribution of the χ^2 statistic cannot be approximated by the asymptotic chi-squared distribution, and hence the p-values were computed using a Monte Carlo test with 5,000 replicates. Results were considered significant when $p < 0.05$. No adjustments were made for multiple testing.

The library sizes required under the various scenarios – L_{NNK} , L_{NNS} , L_{22c} , and L_{Tang} – were computed by TopLib (<http://stat.haifa.ac.il/~yuval/toplib>)⁴².

Results

The model protein used in this study is the mutant F87G of the cytochrome P450_{BM3} from *Bacillus megaterium*, a self-sufficient and highly active monooxygenase that displays an extraordinary catalytic diversity on many natural and non-natural substrates⁴³. In further mutagenesis work, we chose to target residue Ser72 (S72, WT codon: AGT) because it has been recognized to play an important role in the selectivity of P450_{BM3} towards steroids⁴⁴, a challenging synthetic reaction in organic chemistry⁴⁰.

In total, we generated 12 SM libraries in 12 separate experiments. The experiments differed in the randomization scheme used (NNK / NNS / 22c-trick / Tang), primer supplier and primer purity (desalted/HPLC). See Table 1. The QuikChange protocol^{45,46} was used in all cases but employing a different polymerase from the traditional protocols (see Material and Methods).

After transformation and recovery, cells from each of the 12 libraries were grown in solid and liquid media containing the proper antibiotic. The next day, cells from each of the 24 resulting cultures were harvested and lysed, followed by extraction and sequencing of the pooled plasmids to perform the QQC and determine the Q_{pool} values. Next, a single 96-well plate was sequenced for each library using the samples obtained from the solid agar plates; based on the sequence data, the experimental QQC-like charts and Q_{exptl} values were determined.

Library	Randomization Scheme	Primer		QQC charts and Q-values			
		Supplier	Purity	Theoretical Q-value	Solid Q _{pool}	Liquid Q _{pool}	Experimental Q _{exptl}
1	NNK	IDT Technologies	Desalted	 1.000	 0.470	 0.515	 0.798
2	NNK	IDT Technologies	HPLC	 1.000	 0.721	 0.706	 0.837
3	NNK	Life Technologies	Desalted	 1.000	 0.484	 0.438	 0.800
4	NNK	Life Technologies	HPLC	 1.000	 0.486	 0.620	 0.773
5	NNK	Metabion International	Desalted	 1.000	 0.770	 0.657	 0.863
6	NNK	Metabion International	HPLC	 1.000	 0.766	 0.759	 0.888
7	NNS	Metabion International	Desalted	 1.000	 0.643	 0.549	 0.872
8	NNS	Metabion International	HPLC	 1.000	 0.614	 0.592	 0.844
9	22c-trick	Metabion International	Desalted	 1.000	 0.557	 0.575	 0.979
10	22-trick	Metabion International	HPLC	 1.000	 0.590	 0.579	 0.798
11	Tang	Metabion International	Desalted	 1.000	 0.668	 0.543	 0.834
12	Tang	Metabion International	HPLC	 1.000	 0.581	 0.516	 0.742

Table 1. Quick Quality Control and Q-values. Library specifications with resulting QQC charts and Q-values. The three pie charts in each column/row correspond to the three positions in a codon. Black = guanine; green = adenosine; red = threonine; blue = cytosine.

Quick quality control (QQC) and Q-values. The purpose of QQC is to assess whether the desired diversity was introduced at the target codon, by observing the percentage of bases at each of the three base positions of the codon using a single, pooled DNA sequencing electropherogram. These percentages were converted (see Methods) to a Q-value, a number ranging from $Q = 0$, indicating no randomization, to $Q = 1$, indicating perfect randomization²³. The QQC charts and the Q-values are reported in Table 1.

The Q_{pool} values obtained under the two culturing conditions (solid vs. liquid) were highly correlated ($r = 0.74$, $p = 0.006$), with no statistically significant difference in magnitude between conditions ($p = 0.20$). This finding indicates that QQC under the two conditions carries similar information, so that the two procedures are exchangeable. No statistically significant correlation was found between the Q_{pool} values and the Q_{exptl} values, neither for solid nor for liquid media. There was also no statistically

Library	Successfully randomized	Yield (%)	>1 base per position	Non-target mutations	Primer misinsertions	Suboptimal sequencing	Missed amino acids
1	68	72.3	19	6	1	2	Met, Asp
2	65	68.4	24	6	-	1	Lys, Asn, His
3	55	59.8	29	5	3	4	Met, Lys, Asn, Phe
4	54	58.1	37	2	-	3	Ile
5	50	52.1	42	4	-	-	-
6	64	66.7	27	5	-	-	-
7	39	41.9	50	4	-	3	Met, Ile, Gln, Trp
8	64	67.4	24	6	1	1	Phe, Tyr
9	57	59.4	37	1	1	-	-
10	67	69.8	24	4	1	-	Asp, Tyr
11	51	53.7	41	1	2	1	Asp, Tyr
12	65	69.1	23	6	-	2	Lys, Asp, Tyr

Table 2. Sequencing results summary obtained from 96 single colonies formed on agar plates per library.

significant difference in the Q_{pool} values when comparing libraries according to purity grade (desalted vs. HPLC).

When comparing the six NNK libraries, the ones based on primers from supplier 3 exhibited the highest mean Q_{pool} value, compared to suppliers 1 and 2 (0.738 vs. 0.603 and 0.507, respectively). Therefore, we used only primers from supplier 3 to compare the three other randomization schemes: NNS, Tang and 22c-trick.

Analysis of sequencing data. A small number of the samples (≤ 4 per library) yielded low-quality sequencing results, a common finding observed in other studies^{22,23}. Of the successfully sequenced variants, some exhibited imperfect randomization due to various reasons: Additional mutations outside of the target position, misplaced insertion of the primer within the coding sequence, presence of two or more bases at one or more positions of the target codon, etc. See Table 2. We use the term “yield” to denote the percentage of the successfully randomized variants (including those carrying the WT codon) out of the successfully sequenced ones.

The purity of the primers plays a crucial role in determining library yield: The desalted primers exhibited a significantly lower yield than the HPLC ones (mean yield 56.6% vs. 66.7%, respectively; $p < 0.001$), except for supplier 1 with a yield of 72.3%. In fact, within the libraries based on desalted primers, the yield of the library corresponding to supplier 1 was higher than that of the two other suppliers ($p = 0.0013$); no statistically significant differences in yield between suppliers were found within the libraries based on HPLC primers.

Three of the twelve libraries included all 20 amino acids, whereas the remaining nine libraries missed up to 4 amino acids. The Q -value of the liquid media libraries was negatively correlated with the number of amino acids missed ($r = -0.58$, $p = 0.046$), indicating that the Q -value in such libraries possesses a predictive power regarding final amino-acid diversity. No similar statistically significant correlation was found in the solid media libraries, and no correlation was found between the Q -values and yield (in either solid or liquid media).

The prevalent assumption used in mathematical analyses of SM experiments is that the randomized codon distribution is uniform, i.e., each of the possible codons (32 codons in NNK and in NNS, 22 in 22c-trick, and 20 in Tang) is equally likely to anneal to the template and form a variant^{42,47–49}. An alternative assumption, which we now examine, is that the distribution is not uniform, and in particular, that there is a bias in favor of the WT codon. Figure 1 shows the distribution of the randomized codons for each of the 12 libraries. Statistical analysis reveals a statistically significant bias in favor of the WT codon, AGT, in five of the libraries: HPLC NNK from supplier 1 ($p = 0.004$), the two 22c-trick libraries ($p = 0.044$ for desalted and $p < 0.0001$ for HPLC), and the two Tang libraries ($p = 0.004$ for desalted and $p < 0.0001$ for HPLC). In all of these cases, the randomization resulted in AGT with likelihood significantly higher than expected under a uniform distribution. Since AGT is not among the 32 NNS codons (and indeed, no AGT codon was sequenced in any of the two NNS experiments), the two NNS experiments were excluded from this analysis. In the five remaining libraries, no statistically significant evidence for deviation from uniformity was detected, but the confidence intervals for the probability of WT annealing, P_{WT} , across libraries are too wide and overlapping to rule out the possibility of WT bias also in these cases.

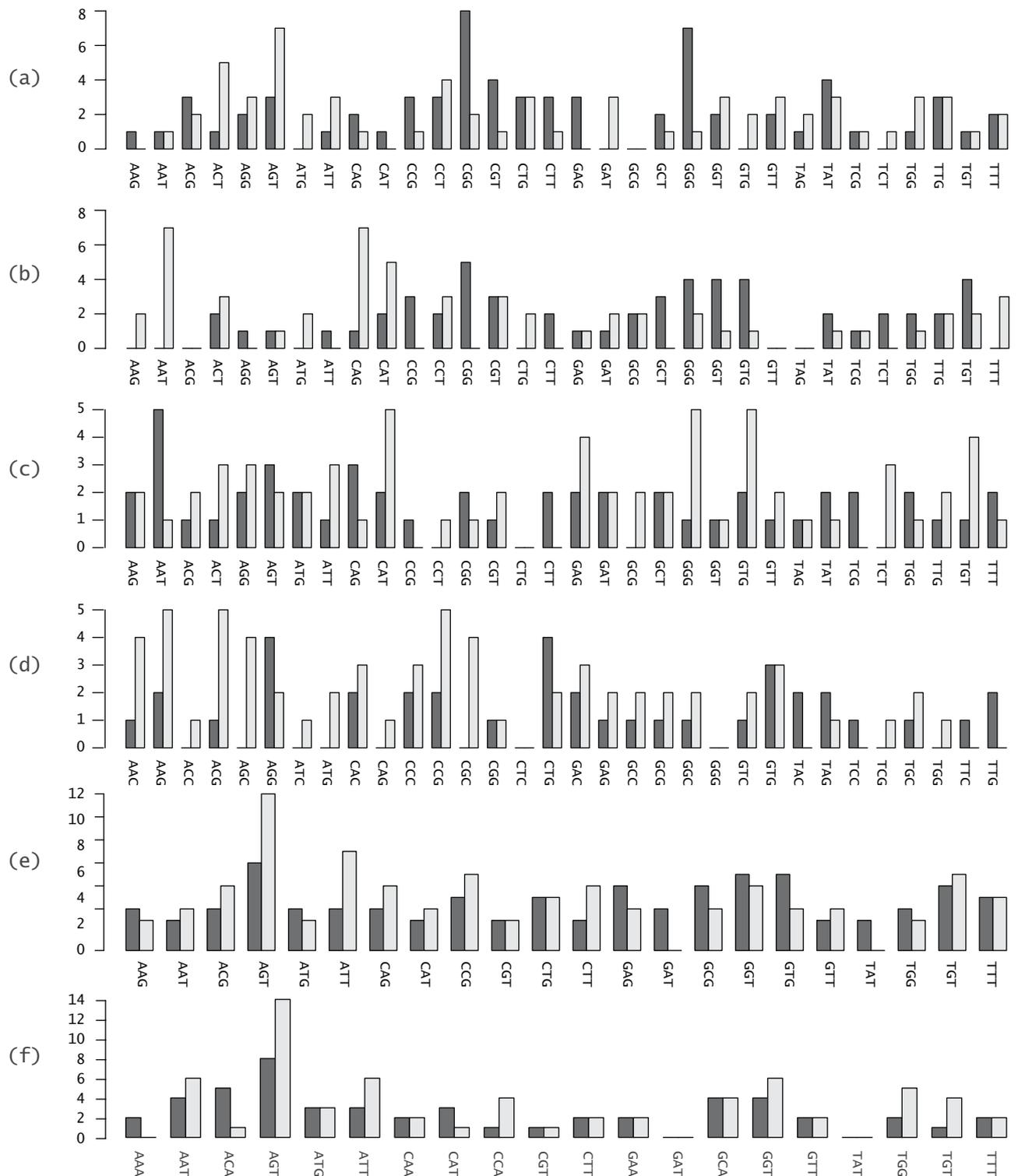


Figure 1. Sequencing results of randomized position S72 (WT codon: AGT) obtained from single colonies formed on agar plates. Vertical axes denote counts (how many codons of each type were found in the sequencing). Black columns denote desalted primers, and grey ones denote HPLC primers. Libraries: (a) 1-2; (b) 3-4; (c) 5-6; (d) 7-8; (e) 9-10; and (f) 11-12.

In all but two of the twelve libraries, the annealing distribution among the remaining, non-WT codons does appear to be uniform. The exceptions that exhibited statistically significant deviation from uniformity were library 1 ($p=0.016$) and library 4 ($p=0.002$).

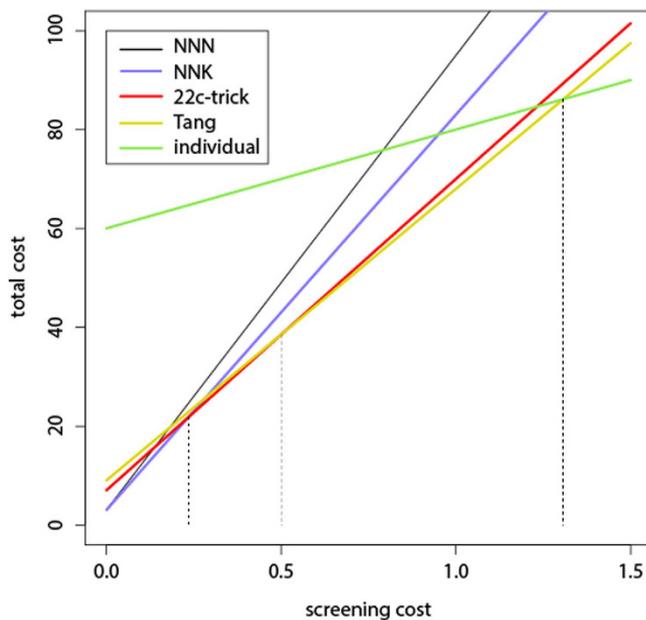


Figure 2. Total cost as a function of screening cost, when randomizing a single position using 5 randomization schemes. Primer cost is $c_{\text{primer}} = 1$.

Cost effectiveness of randomization schemes. Throughout the following analysis, we assume that the QuikChange protocol is used, or a variation thereof employing two partially-overlapping mutagenic primers that may be more suitable for some templates⁵⁰, and that the designer's goal is 95% expected coverage of protein space. This metric, which is sometimes called “95% fractional completeness”, is equivalent to requiring a 0.95 probability for discovering the best variant in a given protein sequence space⁴². The results remain qualitatively unchanged under other library metrics, e.g., requiring 90% expected coverage, or replacing expected coverage with the probability of discovering at least one of the top three variants⁴².

We consider first an idealized SM experiment of a single protein position, in which the yield is 100% and there is no annealing bias. When using NNK (or NNS) randomization, the total cost of the experiment is:

$$c_{\text{NNK}} = 2 \cdot c_{\text{primer}} + c_{\text{fixed}} + L_{\text{NNK}} \cdot c_{\text{screen}}, \quad (1)$$

where c_{primer} is the per-reaction cost of a *single* primer, c_{fixed} is the sum of all other fixed costs associated with a single PCR reaction (labor, dNTPs, polymerase, buffer, template DNA, DpnI, transformation), c_{screen} is the screening cost of a single variant, and $L_{\text{NNK}} = 80$ is the number of variants screened (the library size) required to achieve 95% expected coverage⁴². Similarly, under NNN, 22c-trick, and Tang randomization, the total experimental costs are:

$$c_{\text{NNN}} = 2 \cdot c_{\text{primer}} + c_{\text{fixed}} + L_{\text{NNN}} \cdot c_{\text{screen}} \quad (2)$$

$$c_{\text{22c}} = 6 \cdot c_{\text{primer}} + c_{\text{fixed}} + L_{\text{22c}} \cdot c_{\text{screen}} \quad (3)$$

$$c_{\text{Tang}} = 8 \cdot c_{\text{primer}} + c_{\text{fixed}} + L_{\text{Tang}} \cdot c_{\text{screen}}, \quad (4)$$

where $L_{\text{NNN}} = 92$, $L_{\text{22c}} = 63$, and $L_{\text{Tang}} = 59$. When generating each of the 20 variants individually via site-directed mutagenesis, 20 PCR reactions are required with a pair of primers in each, so the total cost is:

$$c_{\text{indiv}} = 40 \cdot c_{\text{primer}} + 20 \cdot c_{\text{fixed}} + L_{\text{indiv}} \cdot c_{\text{screen}}, \quad (5)$$

where $L_{\text{indiv}} = 20$.

To eliminate the dependence on the currency used, we set $c_{\text{fixed}} = 1$ and measure all costs as multiples of c_{fixed} . This allows focusing on the relative magnitude of the prices, rather than on their actual values, which change rapidly as a result of technological advances.

Figure 2 shows the total cost as a function of c_{screen} , the screening cost, while assuming $c_{\text{primer}} = 1$. The five straight lines correspond to the five randomization schemes. For each value of the screening

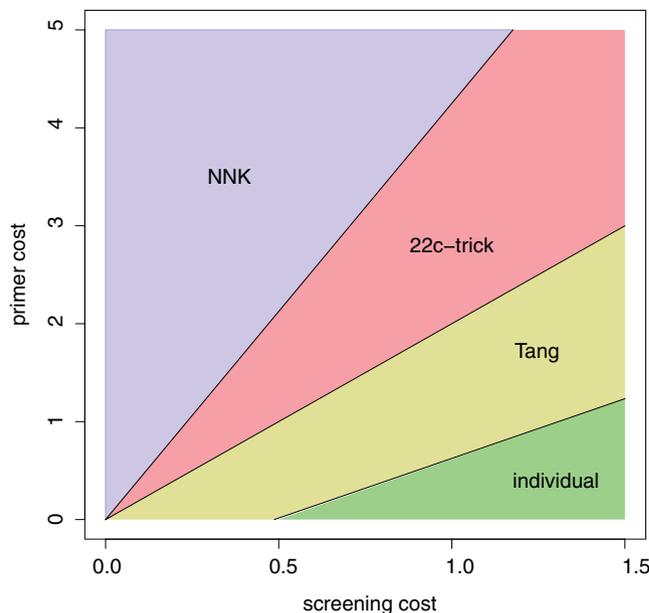


Figure 3. Cost space partitioned into regions according to the optimal randomization scheme (a single randomized position, assuming 100% yield, and no WT bias).

cost (an x value), we are interested in the total cost (y value) that is minimal among the five schemes. These minimal values constitute the so-called *lower envelope* of the five lines. It can be seen that when $0 \leq c_{\text{screen}} \leq 0.24$ the cost is minimized with NNK, when $0.24 \leq c_{\text{screen}} \leq 0.5$ the cost is minimized with 22c-trick, when $0.5 \leq c_{\text{screen}} \leq 1.31$ the cost is minimized with Tang, and when $c_{\text{screen}} \geq 1.31$ it is most economical to generate the 20 variants individually. Note that NNN is never the randomization method of choice.

The differences in cost between 22c-trick and Tang are not very large. On the other hand, the differences between NNK and either 22c-trick or Tang are more substantial, and may exceed 30%. Generating the 20 variants individually becomes the optimal scheme only when screening cost is very high, but when this happens, this method may be considerably cheaper than all other schemes by a large margin.

In Fig. 2 it was assumed that the primer cost equals the fixed cost. Of greater interest is to see how changes in both primer cost and screening cost (relative to the fixed cost of the reaction) determine the optimal randomization scheme. Figure 3 depicts how cost space – which is the two-dimensional plane, whose axes are the primer cost and screening cost – is partitioned into mutually exclusive regions, each corresponding to a different optimal randomization scheme. NNK is optimal when the screening cost is lower than the primer cost by a factor of 4.25 or more (this is the slope of the line separating the NNK and 22c-trick areas); as the screening cost increases, 22c-trick becomes optimal, then Tang, and finally individual generation of the 20 variants.

The partition depicted in Fig. 3 changes under the more realistic assumptions of imperfect yield and some WT bias. The left panel of Fig. 4 shows the partition when the yield is changed from the ideal 100% to 68% (the average yield across the HPLC scenarios), and under a WT bias $P_{\text{WT}} = 0.12$ (the average bias across the scenarios). These changes result in an increased required library sizes: $L_{\text{NNK}} = 121$, $L_{22c} = 99$, and $L_{\text{Tang}} = 93$. The relative position of the four regions corresponding to the four schemes remains the same as in Fig. 3, but the exact location of their boundaries changes. The final total cost under the optimal scheme changes also: for example, when $c_{\text{primer}} = 2.5$ and $c_{\text{screen}} = 0.75$ (the middle point of the region depicted in Figs. 3 and 4 (a), for which 22c-trick happens to be the optimal scheme in both cases), the total cost is 76.0 when assuming 100% yield and no WT bias (as in Fig. 3), and 90.25 when assuming 68% yield and $P_{\text{WT}} = 0.12$ (Fig. 4 (a)).

Under uniform annealing distribution, NNK and NNS randomization are completely identical in terms of the resulting coverage of protein space, as they induce the same distribution over the 20 amino acids and the stop codon. However, when assuming WT bias, differences may occur, which depend on the specific randomized codon. The WT codon in our case is AGT, which is one of the NNK codons but not one of the NNS ones. Thus, under NNS randomization, the WT bias is avoided in this case, and the required library size (assuming 68% yield) is $L_{\text{NNS}} = 118$, slightly lower than the above reported $L_{\text{NNK}} = 121$. Even though NNS randomization is slightly more efficient than NNK in this case, we neglect this difference, and do not include NNS randomization in the analysis.

Another factor that influences the partition of cost space is the number of positions randomized (which was hitherto assumed to be 1). The right panel of Fig. 4 shows the partition of cost space when randomizing two positions, again under 68% yield and $P_{\text{WT}} = 0.12$. It is assumed that the two randomized

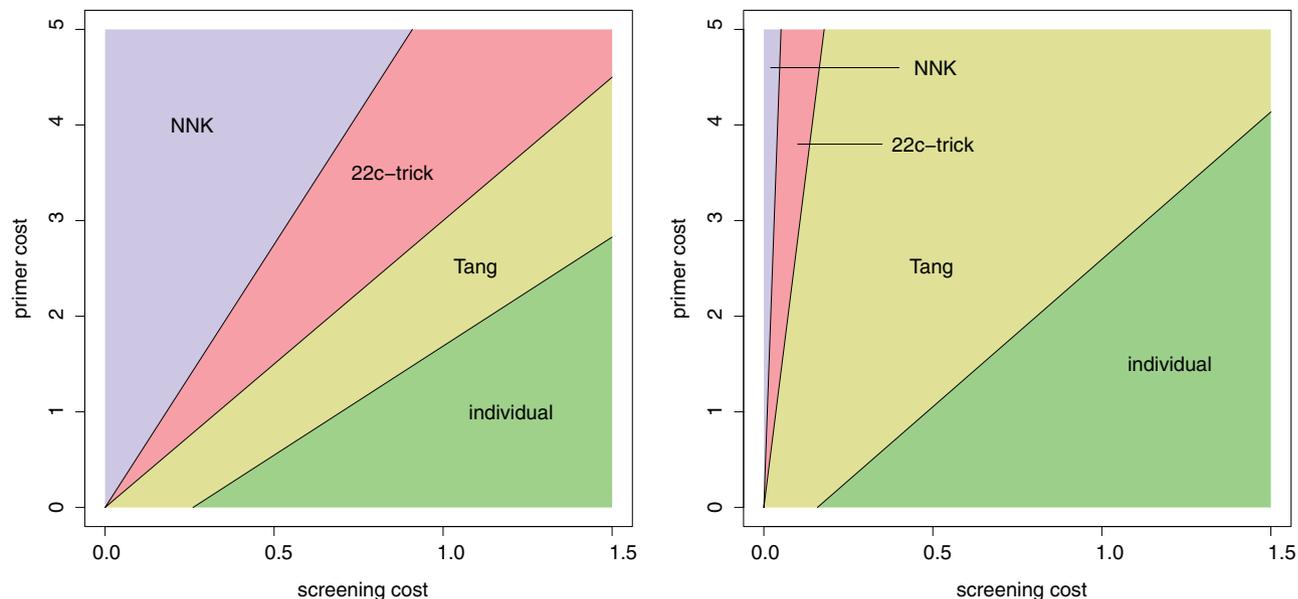


Figure 4. Cost space partition under 68% yield and WT bias $P_{WT} = 0.12$. Left: a single randomized position; right: two randomized positions.

positions are close enough on the primary sequence so that a single primer covers them. The number of NNK primers remains therefore 2, the number of 22c-trick primers is $2 \cdot 3^2 = 18$, the number of Tang primers is $2 \cdot 4^2 = 32$, and the number of primers required when generating individually each of the $20^2 = 400$ variants is $2 \cdot 20^2 = 800$. The required library sizes also increase significantly, to $L_{NNK} = 4881$, $L_{22c} = 3309$, and $L_{Tang} = 2917$.

When compared to Fig. 3, a prominent feature of both panels of Fig. 4 is the smaller area of the NNK region. It can be shown mathematically that this is a general phenomenon, and that this region gets smaller as either the yield decreases or the number of randomized position increases.

Discussion

The common thread in this work is the view of saturation mutagenesis experiments from the economical perspective.

The simple and cheap QQC was developed as a cautionary step to avoid wastage of resources in futile library screening^{24,35,39,40}. When the QQC is performed in solid media, some colonies grow faster than others, creating a “colony bias”. To mitigate this problem, we recommend using liquid cultures, which also represent an economical advantage, since they require as little as a few mL (in our case 4 mL were added to 1 mL of cells in recovery medium) instead of 40 mL in the big solid agar plate format. This way, media costs could be decreased 10-fold, provided that a 14 mL Falcon tube for liquid cultures costs about the same as a large Petri dish. In addition, processing liquid cultures requires less time (only one centrifugation step) than processing solid-agar plates (cell harvesting plus centrifugation step), and the latter approach also requires more incubation space and workload. Thus, the liquid media QQC could be more useful for SM experiments in a high-throughput setting. Stewart and colleagues introduced the Q-value, which summarizes in a single number the QQC results²³. Our statistical analysis shows that the Q-values obtained under the two culture conditions are highly correlated, but that only the liquid Q-value was correlated with the final amino-acid diversity. This finding suggests that the liquid QQC mitigates the “colony bias” problem present in the solid QQC format.

Another economical factor in SM experiments is primer cost. HPLC-purified primers generally cost up to 3 times more than column-purified desalted ones. To our knowledge, only two studies have indicated the role of primer purity in saturation or site-directed mutagenesis. Steffens and Williams reported that desalted primers were successfully used to create 200 single-site SM libraries of a polymerase, but no sequencing data of the resulting library mutants was provided⁵¹. On the other hand, the employment of HPLC-purified primers is recommended for the traditional QuikChange protocol, first by Stratagene⁴⁵ and currently by Agilent⁴⁶, the most widely used method for site-directed mutagenesis and SM⁵². Although no statistically significant difference in Q-value was found between desalted and HPLC libraries, our sequencing data shows that primer purity plays a critical role, with HPLC libraries enjoying significantly higher yield on average than desalted ones. Nevertheless, using desalted primers obtained from supplier 1 gave comparable results in terms of yield to HPLC-purified primers regardless of supplier, indicating that supplier 1 has higher primer quality standards than the two other suppliers regarding non-purified primers.

To further improve the economical efficiency of SM experiments, it is worthwhile to consider alternatives to the traditional NNK and NNS randomization schemes. We³⁵ and others³⁴ recently showed that the respective 22c-trick and Tang strategies can reduce the screening effort of single-site libraries compared to NNK/S, by a factor of about 50%. Other researchers have also demonstrated the practical utility of the 22c-Trick^{53,54} and Tang^{55,56} approaches. However, these “smarter” schemes entail a higher primer budget, and for this reason we established a methodological framework for comparing the cost-effectiveness of the various alternatives. The cost equations we use are perhaps somewhat simplistic, yet they capture the main trade-off between the cost of primers and the cost of screening: by switching to another randomization scheme, one can reduce the former by increasing the latter, or vice versa. The above analysis becomes important especially in large-scale, expensive experiments. This happens either when a large number of a protein's positions are subject separately to SM, or when relatively many (say, ≥ 3) positions are randomized simultaneously. In either case, we recommend that a short, preliminary study be conducted before the main one, to provide estimates for the model parameters (yield, WT bias) required for choosing optimally the randomization scheme and the library sizes under the specific experimental conditions at hand.

We show that under certain experimental settings – namely, when screening cost is high enough – the 22c-trick and Tang schemes are economically superior alternatives to NNK and NNS. Perhaps surprisingly, this is true also when several close-by positions are randomized simultaneously. Indeed, the number of primers required for 22c-trick or Tang randomization grows exponentially as the number of randomized positions increases, but the saving in screening efforts due to the reduction in library size (relative to NNS or NNK) more than offsets the increased cost of primers. When choosing between NNK and NNS, we recommend choosing the scheme that does not include the template's WT codon (as was NNS relative to the WT codon AGT in our case), in order to avoid the WT bias problem and to get a more uniform annealing distribution. The difference in the resulting library size, and hence in screening cost, may be small or even negligible (as was in our case), but when the WT bias is very high, the savings due to choosing the better method might be significant. Another option to reduce WT bias is to modify the base composition ratio during primer synthesis⁵⁷, but this strategy may require multiple PCR optimization steps and may not be thus suitable for high-throughput experiments.

Another factor playing a role in library quality is the primer annealing bias. We modeled this bias only through an increased annealing probability of primers carrying the WT codon, relative to the other primers. However, the annealing bias may manifest itself more finely, according to the number of Watson-Crick pairing mismatches in a codon: the WT codon is an extreme case with zero mismatches, and it is possible that a codon with, say, one mismatch will be more likely to anneal than a codon with two mismatches, etc. Our statistical analysis did not detect such a phenomenon, but also did not rule it out. A much larger data set is required to study the existence and magnitude of this phenomenon.

In summary, we have presented a faster, more economical and reliable method for performing the QQC, based on liquid cultures. Importantly, the QQC should be combined with the Q-values to assess overall library quality. We also demonstrated that primer purity has a significant effect on library yield, but that some suppliers might offer primers of higher quality than others without additional purification steps. In addition, we provided guidelines for choosing optimally a randomization scheme, depending on the screening costs and other experimental parameters. Our guidelines also apply to any PCR-based method for library preparation, including combinatorial gene synthesis⁵⁸, gene assembly⁵⁹ and overlap extension PCR⁶⁰.

References

- Smith, M. Synthetic DNA and Biology (Nobel Lecture). *Angew. Chem. Int. Ed. Engl.* **33**, 1214–1221 (1994).
- Siloto, R. M. P. & Weselake, R. J. Site saturation mutagenesis: Methods and applications in protein engineering. *Biocatal. Agric. Biotechnol.* **1**, 181–189 (2012).
- Valetti, F. & Gilardi, G. Improvement of Biocatalysts for Industrial and Environmental Purposes by Saturation Mutagenesis. *Biomolecules* **3**, 778–811 (2013).
- Reetz, M. T. Biocatalysis in organic chemistry and biotechnology: past, present, and future. *J. Am. Chem. Soc.* **135**, 12480–12496 (2013).
- Reetz, M. T. Laboratory evolution of stereoselective enzymes: a prolific source of catalysts for asymmetric reactions. *Angew. Chem. Int. Ed. Engl.* **50**, 138–174 (2011).
- Gillam, E. M., Copp, J. N. & Ackerley, D. F. Directed Evolution Library Creation. (Springer-Verlag New York (Humana Press), Totowa; 2014).
- Sidhu, S. S. & Kossiakoff, A. A. Exploring and designing protein function with restricted diversity. *Curr. Opin. Chem. Biol.* **11**, 347–354 (2007).
- Pattanaik, S., Werkman, J. R., Kong, Q. & Yuan, L. Site-directed mutagenesis and saturation mutagenesis for the functional study of transcription factors involved in plant secondary metabolite biosynthesis. *Methods Mol. Biol.* **643**, 47–57 (2010).
- Patwardhan, R. P. *et al.* High-resolution analysis of DNA regulatory elements by synthetic saturation mutagenesis. *Nat. Biotechnol.* **27**, 1173–1175 (2009).
- Smith, J. D., McManus, K. F. & Fraser, H. B. A novel test for selection on cis-regulatory elements reveals positive and negative selection acting on mammalian transcriptional enhancers. *Mol. Biol. Evol.* **30**, 2509–2518 (2013).
- Wang, H. H. *et al.* Programming cells by multiplex genome engineering and accelerated evolution. *Nature* **460**, 894–898 (2009).
- Findlay, G. M., Boyle, E. A., Hause, R. J., Klein, J. C. & Shendure, J. Saturation editing of genomic regions by multiplex homology-directed repair. *Nature* **513**, 120–123 (2014).
- Oh, J. H. & van Pijkeren, J. P. CRISPR-Cas9-assisted recombineering in *Lactobacillus reuteri*. *Nucleic Acids Res.* **42**, e131 (2015).
- Acevedo-Rocha, C. G., Hoebenreich, S. & Reetz, M. T. Iterative saturation mutagenesis: a powerful approach to engineer proteins by systematically simulating Darwinian evolution. *Methods Mol. Biol.* **1179**, 103–128 (2014).

15. Reetz, M. T. & Carballeira, J. D. Iterative saturation mutagenesis (ISM) for rapid directed evolution of functional enzymes. *Nat. Protoc.* **2**, 891–903 (2007).
16. Abatemarco, J., Hill, A. & Alper, H. S. Expanding the metabolic engineering toolbox with directed evolution. *Biotechnol. J.* **8**, 1397–1410 (2013).
17. Cobb, R. E., Si, T. & Zhao, H. Directed evolution: an evolving and enabling synthetic biology tool. *Curr. Opin. Chem. Biol.* **16**, 285–291 (2012).
18. Dietrich, J. A., McKee, A. E. & Keasling, J. D. High-throughput metabolic engineering: advances in small-molecule screening and selection. *Annu. Rev. Biochem.* **79**, 563–590 (2010).
19. Esvelt, K. M. & Wang, H. H. Genome-scale engineering for systems and synthetic biology. *Mol. Syst. Biol.* **9**, 641 (2013).
20. Boyle, N. R. & Gill, R. T. Tools for genome-wide strain design and construction. *Curr. Opin. Biotechnol.* **23**, 666–671 (2012).
21. Wang, H. H. & Church, G. M. Multiplexed genome engineering and genotyping methods applications for synthetic biology and metabolic engineering. *Methods Enzymol.* **498**, 409–426 (2011).
22. Hoebenreich, S., Zilly, F. E., Acevedo-Rocha, C. G., Zilly, M. & Reetz, M. T. Speeding up Directed Evolution: Combining the Advantages of Solid-Phase Combinatorial Gene Synthesis with Statistically Guided Reduction of Screening Effort. *ACS Synth. Biol.* **4**, 317–331 (2015).
23. Sullivan, B., Walton, A. Z. & Stewart, J. D. Library construction and evaluation for site saturation mutagenesis. *Enzyme. Microb. Technol.* **53**, 70–77 (2013).
24. Sanchis, J. *et al.* Improved PCR method for the creation of saturation mutagenesis libraries in directed evolution: application to difficult-to-amplify templates. *Appl. Microbiol. Biotechnol.* **81**, 387–397 (2008).
25. Acevedo-Rocha, C. G., Agudo, R. & Reetz, M. T. Directed evolution of stereoselective enzymes based on genetic selection as opposed to screening systems. *J. Biotechnol.* **191**, 3–10 (2014).
26. Baronio, R. *et al.* All-codon scanning identifies p53 cancer rescue mutations. *Nucleic Acids Res.* **38**, 7079–7088 (2010).
27. Cornishbowden, A. Nomenclature for Incompletely Specified Bases in Nucleic-Acid Sequences - Recommendations 1984. *Nucleic Acids Res.* **13**, 3021–3030 (1985).
28. Gaytan, P. & Roldan-Salgado, A. Elimination of redundant and stop codons during the chemical synthesis of degenerate oligonucleotides. Combinatorial testing on the chromophore region of the red fluorescent protein mKate. *ACS Synth. Biol.* **2**, 453–462 (2013).
29. Neuner, P., Cortese, R. & Monaci, P. Codon-based mutagenesis using dimer-phosphoramidites. *Nucleic Acids Res.* **26**, 1223–1227 (1998).
30. Ono, A., Matsuda, A., Zhao, J. & Santi, D. V. The Synthesis of Blocked Triplet-Phosphoramidites and Their Use in Mutagenesis. *Nucleic Acids Res.* **23**, 4677–4682 (1995).
31. Gaytan, P., Contreras-Zambrano, C., Ortiz-Alvarado, M., Morales-Pablos, A. & Yanez, J. TrimerDimer: an oligonucleotide-based saturation mutagenesis approach that removes redundant and stop codons. *Nucleic Acids Res.* **37**, e125 (2009).
32. Hughes, M. D., Nagel, D. A., Santos, A. F., Sutherland, A. J. & Hine, A. V. Removing the redundancy from randomised gene libraries. *J. Mol. Biol.* **331**, 973–979 (2003).
33. Ashraf, M. *et al.* ProxiMAX randomization: a new technology for non-degenerate saturation mutagenesis of contiguous codons. *Biochem. Soc. Trans.* **41**, 1189–1194 (2013).
34. Tang, L. *et al.* Construction of “small-intelligent” focused mutagenesis libraries using well-designed combinatorial degenerate primers. *BioTechniques* **52**, 149–158 (2012).
35. Kille, S. *et al.* Reducing codon redundancy and screening effort of combinatorial protein libraries created by saturation mutagenesis. *ACS Synth. Biol.* **2**, 83–92 (2013).
36. Pines, G. *et al.* Codon Compression Algorithms for Saturation Mutagenesis. *ACS Synth. Biol.* (2014).
37. Tang, L. *et al.* MDC-Analyzer: a novel degenerate primer design tool for the construction of intelligent mutagenesis libraries with contiguous sites. *BioTechniques* **56**, 301–310 (2014).
38. Nov, Y. & Segev, D. Optimal codon randomization via mathematical programming. *J. Theor. Biol.* **335**, 147–152 (2013).
39. Hoguekou, D. J., Kille, S., Taglieber, A. & Reetz, M. T. Directed Evolution of an Enantioselective Enoate-Reductase: Testing the Utility of Iterative Saturation Mutagenesis. *Adv. Synth. Catal.* **351**, 3287–3305 (2009).
40. Kille, S., Zilly, F. E., Acevedo, J. P. & Reetz, M. T. Regio- and stereoselectivity of P450-catalysed hydroxylation of steroids controlled by laboratory evolution. *Nat. Chem.* **3**, 738–743 (2011).
41. Agudo, R., Roiban, G. D. & Reetz, M. T. Achieving regio- and enantioselectivity of P450-catalyzed oxidative CH activation of small functionalized molecules by structure-guided directed evolution. *ChemBioChem* **13**, 1465–1473 (2012).
42. Nov, Y. When second best is good enough: another probabilistic look at saturation mutagenesis. *Appl. Environ. Microbiol.* **78**, 258–262 (2012).
43. Whitehouse, C. J., Bell, S. G. & Wong, L. L. P450(BM3) (CYP102A1): connecting the dots. *Chem. Soc. Rev.* **41**, 1218–1260 (2012).
44. Venkataraman, H. *et al.* A single active site mutation inverts stereoselectivity of 16-hydroxylation of testosterone catalyzed by engineered cytochrome P450 BM3. *ChemBioChem* **13**, 520–523 (2012).
45. Hogrefe, H. H., Cline, J., Youngblood, G. L. & Allen, R. M. Creating randomized amino acid libraries with the QuikChange Multi Site-Directed Mutagenesis Kit. *BioTechniques* **33**, 1158–1160, 1162, 1164–1155 (2002).
46. Hogrefe, H. H. Fine-tuning enzyme activity through saturation mutagenesis. *Methods Mol. Biol.* **634**, 271–283 (2010).
47. Patrick, W. M. & Firth, A. E. Strategies and computational tools for improving randomized protein libraries. *Biomol. Eng.* **22**, 105–112 (2005).
48. Bosley, A. D. & Ostermeier, M. Mathematical expressions useful in the construction, description and evaluation of protein libraries. *Biomol. Eng.* **22**, 57–61 (2005).
49. Denault, M. & Pelletier, J. N. Protein library design and screening: working out the probabilities. *Methods Mol. Biol.* **352**, 127–154 (2007).
50. Xia, Y., Chu, W., Qi, Q. & Xun, L. New insights into the QuikChange™ process guide the use of Phusion DNA polymerase for site-directed mutagenesis. *Nucleic Acids Res.* (2014).
51. Steffens, D. L. & Williams, J. G. Efficient site-directed saturation mutagenesis using degenerate oligonucleotides. *J. Biomol. Tech.* **18**, 147–149 (2007).
52. Tee, K. L. & Wong, T. S. Polishing the craft of genetic diversity creation in directed evolution. *Biotechnol. Adv.* **31**, 1707–1721 (2013).
53. Blomberg, R. *et al.* Precision is essential for efficient catalysis in an evolved Kemp eliminase. *Nature* **503**, 418–421 (2013).
54. McIsaac, R. S. *et al.* Directed evolution of a far-red fluorescent rhodopsin. *Proc. Natl. Acad. Sci. U.S.A.* **111**, 13034–13039 (2014).
55. Lacey, V. K., Louie, G. V., Noel, J. P. & Wang, L. Expanding the library and substrate diversity of the pyrrolysyl-tRNA synthetase to incorporate unnatural amino acids containing conjugated rings. *ChemBioChem* **14**, 2100–2105 (2013).
56. Blikstad, C., Dahlström, K. M., Salminen, T. A. & Widersten, M. Stereoselective Oxidation of Aryl-Substituted Vicinal Diols into Chiral α -Hydroxy Aldehydes by Re-Engineered Propanediol Oxidoreductase. *ACS Catalysis* **3**, 3016–3025 (2013).
57. Airaksinen, A. & Hovi, T. Modified base compositions at degenerate positions of a mutagenic oligonucleotide enhance randomness in site-saturation mutagenesis. *Nucleic Acids Res.* **26**, 576–581 (1998).

58. Currin, A., Swainston, N., Day, P. J. & Kell, D. B. SpeedyGenes: an improved gene synthesis method for the efficient production of error-corrected, synthetic protein libraries for directed evolution. *Protein Eng. Des. Sel.* **27**, 273–280 (2014).
59. Acevedo-Rocha, C. G. & Reetz, M. T. Assembly of Designed Oligonucleotides: a useful tool in synthetic biology for creating high-quality combinatorial DNA libraries. *Methods Mol. Biol.* **1179**, 189–206 (2014).
60. Williams, E. M., Copp, J. N. & Ackerley, D. F. Site-saturation mutagenesis by overlap extension PCR. *Methods Mol. Biol.* **1179**, 83–101 (2014).

Acknowledgements

This work was supported by the Max-Planck-Society, the LOEWE Research Cluster SynChemBio from the state of Hessen, Germany; and the Israeli Science Foundation (grant 286/13). C.G.A.R is holder of a SYNMIKRO postdoctoral fellowship awarded by the LOEWE program of the state of Hessen, Germany.

Author Contributions

All authors devised the experiments, C.G.A.R. conducted the experiments, Y.N. conducted the mathematical and statistical analysis, and all authors wrote the paper.

Additional Information

Competing financial interests: The authors declare no competing financial interests.

How to cite this article: Acevedo-Rocha, C. G. *et al.* Economical analysis of saturation mutagenesis experiments. *Sci. Rep.* **5**, 10654; doi: 10.1038/srep10654 (2015).



This work is licensed under a Creative Commons Attribution 4.0 International License. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in the credit line; if the material is not included under the Creative Commons license, users will need to obtain permission from the license holder to reproduce the material. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>