# 10    Directed evolution using multiple amino acid substitutions

© 2023 Romas Kazlauskas

**Summary.**  Large changes in protein properties require multiple amino acid substitutions. If the substitutions act independently, then they can be discovered as beneficial single substitutions, then combined for a larger effect. Nearby substitutions or substitutions that cause shifts in conformation may act cooperatively. These cooperatively acting substitutions can be discovered by recombination of beneficial single substitutions, by accumulation of beneficial single substitutions in multiple rounds of mutagenesis and screening, or directly by screening large libraries containing multiple simultaneous substitutions. The fitness landscape metaphor compares climbing a mountain to improving protein properties by amino acid substitutions. Steps on this landscape correspond to single amino acid substitutions, while jumps on this landscape correspond to multiple simultaneous amino acid substitutions.

### Key learning goals

- Many beneficial single substitutions act independently, combining all of them yields a combined benefit. Other beneficial single substitutions interact with each other to change the size of the effect or even to become detrimental. Recombination of these substitutions finds the best subset by eliminating detrimental combinations.

- Some beneficial substitutions are hidden because they become beneficial only by cooperation with other substitutions. Accumulation of beneficial substitutions in rounds of mutagenesis and screening can find some of these hidden substitutions. A common type of hidden beneficial substitution is a destabilizing, but otherwise beneficial substitution. Pre-stabilization of the starting protein eliminates the destabilizing effect and reveals these hidden beneficial substitutions.

- Screening multiple simultaneous substitution libraries can find complex cooperative effects as well as find independently acting beneficial single substitutions. Two limitations of multiple simultaneous substitution libraries are their vast size and the masking of beneficial effects of substitution by a strongly detrimental substitution. DNA shuffling or smart libraries, which limit replacement amino acids to those that appear in homologs, minimize the fraction of strongly detrimental substitutions.

Making large changes in protein properties or changes in multiple properties of a protein requires correspondingly large changes in the structure of the protein, which means making multiple amino acid substitutions. For example, switching the coenzyme specificity of isopropyl malate dehydrogenase from NAD to NADP required five amino acid replacements, each of which contributed to the change in specificity (Lunzer et al., 2005), Table 10.1. The 100-fold preference for NAD changed to a 200-fold preference for NADP, which corresponds to a 20,000-fold change in selectivity. The engineered enzyme was as active with the NADP as the original enzyme was with NAD.

**Table 10.1.** Selected examples of multiple amino acid substitutions to improve an enzyme.

| enzyme | improvement | amino acid replacements | reference |
|---|---|---|---|
| isopropyl malate dehydrogenase | 20,000-fold shift in coenzyme selectivity from NAD to NADP | five - in coenzyme binding site | Lunzer et al., 2005 |
| *p*-nitrobenzyl esterase | 100-fold increase in activity in organic solvent-water mixtures | seven - outside the active site | Spiller et al., 1999 |
| phosphotriesterase | 135-fold faster hydrolysis of cyclosarin, a nerve toxin | five - in substrate binding site | Gupta et al., 2011 |
| transaminase | >40,000-fold increase in volumetric productivity | 27 - throughout protein | Savile et al., 2010 |
| halohydrin dehalogenase | ~4000-fold increase in volumetric productivity | 35 - throughout protein | Fox et al., 2007 |

Similarly, improvements in several different properties will also usually require multiple amino acid substitutions. For example, a common engineering goal for enzymes used in manufacturing is higher volumetric productivity, which is the ability to produce more product for a given amount of enzyme, time and reaction volume. This ability depends on the catalytic efficiency of the enzyme, its stability to reaction conditions, inhibition by substrates and products, as well as other factors. Multiple substitutions are required to improve these various properties. Multiple mutations - 27 of the 330 amino acids of the transaminase (8% of the total) and 35 of the 254 amino acids in halohydrin dehalogenase (14% of the total) - dramatically improved the volumetric productivity of industrial enzymes, Table 10.1.

This chapter surveys the different approaches protein engineers use to find these multiple substitutions. Some approaches focus on discovering individual beneficial substitutions, then combining them. Other approaches make multiple substitutions simultaneously to find cooperatively-acting combinations.

# 10.1 Independently acting beneficial substitutions

In many cases beneficial single substitutions continue to be equally beneficial when combined with other substitutions. This independent behavior occurs when the two amino acids do not interact with each other, which is most likely when the amino acid substitutions are well-separated from each other and make minimal changes in structure (Wells, 1990). Combinations of independently acting substitutions can be discovered by first discovering beneficial single substitutions and then combining all of them to combine their effects. The best combination contains all of the beneficial single substitutions.

When amino acid substitutions act independently, their free energy effects, $\Delta G$, are additive. If mutations X and Y independently improve a protein by free energy amounts $\Delta\Delta G_X$ and $\Delta\Delta G_Y$, then combination of X and Y improves the protein by the sum of these free energies. For example, if two independently acting mutations X and Y each decrease the transitions state energy 1.4 kcal/mol, then the combination XY decreases the transition state by 2.8 kcal/mol.

$$\Delta\Delta G_{XY} = \Delta\Delta G_X + \Delta\Delta G_Y \qquad (10.1)$$

This additive behavior of free energies leads to multiplicative behavior of the properties. Recall that free energy is related to the logarithm of the size of the effect, eq. 10.2. Subtracting two terms with logarithms requires dividing the expressions within the logarithm. Adding two free energy terms corresponds to multiplying the terms within the logarithm, eq. 10.3. Thus, the effects of each mutation multiply upon combination. In the example above, the decrease in transition state energy by 1.4 kcal/mol corresponds to increase catalytic activity by a factor of 10 at 25 °C and the decrease of 2.8 kcal/mol corresponds to a 100-fold increase in catalytic activity.

$$\Delta\Delta G_X = \Delta G_X - \Delta G_{wt} = -RT\ln(x) - (-RT\ln(wt)) = -RT\ln(x/wt)$$

$$(10.2)$$

$$\Delta\Delta G_X + \Delta\Delta G_Y = -RT\ln(x/wt) + [- RT\ln(y/wt)] = -RT\ln(x/wt \cdot y/wt)$$

$$(10.3)$$

Protein-stabilizing substitutions often act independently because they involve local interactions. For example, substitution with proline decreases the flexibility of the unfolded state at the site of the substitution. For this reason stabilizing substitutions are often discovered as single-stabilizing substitutions, then combined to create a stabilized protein. For example, Rath and Davidson (2000) first identified three single amino acid substitutions that stabilized a small SH3 domain protein: $\Delta\Delta G_{unfold}$: Glu7Leu 1.23 kcal/mol, Val21Lys 0.38 kcal/mol, Asn23Gly 1.41 kcal/mol. Combining these three substitutions into a triple substitution variant increased the $\Delta\Delta G_{unfold}$ by 3.35 kcal/mol relative to wild type, which is nearly identical to the sum of the three single substitutions: 3.02 kcal/mol.

Additive behavior of substitutions simplifies protein engineering. Improvements can be discovered independently and then combined to make larger improvements. An example of additive behavior is the engineering of the product distribution of the γ-humulene synthase (Yoshikuni et al., 2006). This terpene cyclase catalyzes the cyclization of farnesyl diphosphate to γ-humulene in 45% yield, but forms 51 other sesquiterpenes in smaller amounts. This mixture of products forms because the substrate can fold in different ways and the carbocation intermediate can rearrange in different ways. The researchers independently altered the 19 amino acid residues that shape the active site by saturation mutagenesis, measured the changes in product distribution, and, using an additive model, predicted combinations of substitutions to form predominantly one or another product. The important locations were separated from each other within the active site, so the assumption of additive behavior was reasonable. For example, the starting enzyme formed only 23 mol% sibirene, but a triple substitution increased the fraction to 78%, while the additive model predicted 81%. Protein engineering starts with the assumption that the effects of amino acid substitutions will be additive when the amino acids don't interact directly, but recognizes that the strategies may need to be adjusted if cooperativity is detected.

In most cases, the magnitude of the combined effect does not exactly match the sum of the individual effects. Some combinations yield more than expected improvement; others yield less

than expected. This difference measure the cooperativity of the interactions. In the SH3 domain stabilization example above the triple substitution was 0.33 kcal/mol more stable than expected from the sum of the effects of the three single substitutions (3.35 vs. 3.02 kcal/mol). In this case the cooperativity creates a gain of 0.33 kcal/mol in stability. This not-quite additive behavior is the first type of cooperativity, magnitude cooperativity, which will be discussed in the next section. In this case, the cooperativity is positive since the benefit of the combination is larger than the sum of the individual substitutions. Negative cooperativity yields a benefit smaller than the sum. A possible source of the magnitude cooperativity in the SH3 domain is an interaction between the substitutions at positions 21 and 23, which both lie in the same β-turn.

## 10.2 Cooperatively-acting substitutions

Cooperativity means that the changes do not act independently; instead, their effects are non-additive when combined. Cooperativity is also sometimes called epistasis, which originally referred to cooperative effects between genes, but now also refers to cooperative effects within a protein. Cooperativity is more common when the changes are nearby. Direct contact or large structure changes cause each substitution to affect the environment of the other and alter its effect. Non-additive behavior is also possible between distant substitutions. They can interact with each other by electrostatic interactions, which act at distances up to 10 Å, or via cooperative behavior like a reaction mechanism or conformational changes including protein folding or unfolding.

One way to classify cooperativity between single substitutions considers whether the contribution of the single substitution reverses when combined (Poelwijk et al., 2007), Table 10.2. The first type of cooperativity is magnitude cooperativity, which shows no reversal in the beneficial effect of the single substitutions. Both substitutions remain beneficial, but the combined beneficial effect is larger or smaller than the sum of the single mutations. That is, the effects are only approximately additive. The previous section mentioned an example of positive magnitude cooperativity. Even though each substitution is beneficial, the best combination may be a subset of the beneficial substitutions. This subset maximizes combinations that yield more than the sum of the independent effects and minimizes combinations that yield less than the sum of the independent effects.

**Table 10.2.** Additive behavior and four types of cooperativity between single amino acid substitutions.

| Type of cooperativity | Effect of individual changes | Effect of combined changes |
| --- | --- | --- |
| none = additive (+) + (+) yields 2(+) | A and B are both beneficial. | AB is beneficial to the same extent as A & B individually. |
| 1. Magnitude differs (+) + (+) yields >2(+) or <2(+) | A and B are both beneficial. | AB is beneficial to a different extent as A & B individually. |
| 2. One sign changes (+) + (+) yields (– or 0) | A and B are both beneficial. | AB is detrimental or neutral. |
| 3. One sign changes (+) + (– or 0) yields >(+) | A is beneficial, but B is detrimental or neutral. | AB is more beneficial that A alone; B becomes beneficial when A is present. |
| 4. Both signs change (– or 0) + (– or 0) yields (+) | A and B are both detrimental or neutral. | A and B become beneficial only when the other is present. |

The second and third types of cooperativity are one-sign change cooperativity where the effect of one of the substitutions reverses upon combination. In the second type a beneficial single substitution becomes detrimental upon combination. These single substitutions are analogous to false positives since upon combination they become detrimental. Engineering requires eliminating these false positives from the combinations of substitutions.
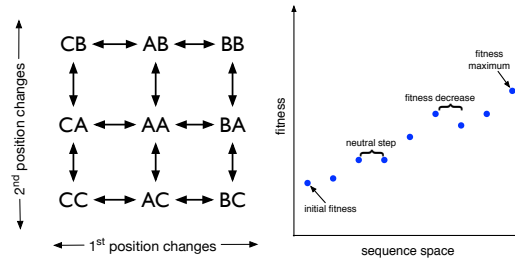
In the third type a neutral or detrimental substitution becomes beneficial upon combination. These are analogous to false negatives since the single substitution does not appear beneficial, but will become beneficial when combined with other substitutions. Since a substitution must show a benefit to be discovered, these false negative substitutions cannot be discovered until the other substitutions are present. An example of this third type of cooperativity is a destabilizing substitution that also benefits catalysis or selectivity. By itself this substitution is neutral or detrimental due to the destabilizing effect, but when combined with a stabilizing substitution, its beneficial effect on catalysis or selectivity becomes evident.

The fourth type of cooperativity is two-sign-change cooperativity. Both substitutions reverse their effect from neutral or detrimental to beneficial. Neither substitution can be discovered independently; both must be present for the benefit. For example, the combination of adding a hydrogen bond donor and shifting the loop containing it can create a new favorable interaction with the substrate, but neither the hydrogen bond donor nor the loop shift are individually beneficial. (See Fig 10.3 later in this text.)

Measuring pairwise cooperativity experimentally is straightforward. One compares the measured effects (free energy contributions) of two individual substitutions on the wild-type protein to the measured effect of the combined substitutions on the wild-type protein. An exact match indicates additive behavior, while any difference indicates cooperativity. The amount of the difference is the cooperative interaction energy between the two substitutions.

**The fitness landscape metaphor**

Protein sequence space is an imaginary multidimensional space that represents all possible amino acid sequences for a protein. Adjacent points differ by a single amino acid substitution, while more distant points differ by multiple amino acid substitutions. For a dipeptide where there are three amino acid choices (A, B, C) one can sketch the sequence space showing the nine possibilities (3^2), Figure 10.1a. For typical proteins, sequence space is impossible to imagine since there are multiple dimensions representing substitutions in different positions and substitutions with different amino acids.



**Figure 10.1.** a) Protein sequence space is an imaginary multidimensional space where each point represents a possible sequence. The sequence space for a dipeptide consisting of three possible amino acids (A, B, C) contain nine points. Adjacent sequences differ by a single substitution while more distant sequences differ by two substitutions. b) A two-dimensional fitness landscape combines fitness along the y-axis with a simplified sequence space along the x-axis. The points show a possible path to increasing fitness. Most steps increase fitness, but one is neutral and another decreases fitness before increasing it in subsequent steps. Strong selection pressure requires each step to increase fitness, while weak selection pressure tolerates neutral and slightly deleterious steps.
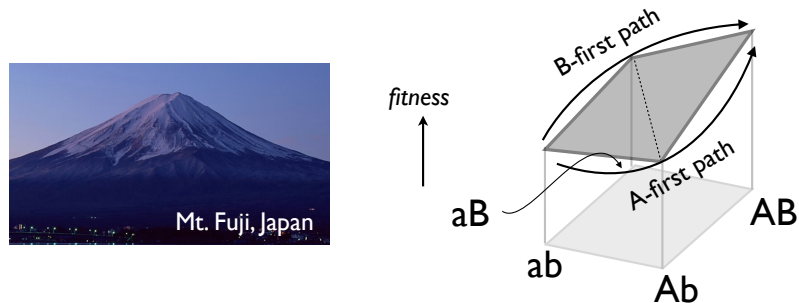
A fitness landscape adds the dimension of fitness to protein sequence space, Figure 10.1b. Fitness of an organism is its ability to survive and reproduce, while fitness in protein engineering refers to progress toward the protein engineering goal. To draw a two-dimensional fitness landscape, one abstracts multidimensional sequence space to a single axis. Movements along this axis represent single substitutions, but now they may occur anywhere within the protein sequence and to any amino acid. Since evolution selects for fitness, the sequence of substitutions during evolution must be an upward path. A substitution must increase fitness (or at least be neutral) for the substitution to be accessible.

A mountainous landscape is a useful metaphor for a three dimensional fitness landscape (Maynard-Smith, 1970; Wright, 1932). Height corresponds to fitness, while movement along the surface represents changes in amino acid sequence. Single steps represent single amino acid substitutions, while jumps to a new location represent simultaneous changes in several amino acids. To reach the peak with single amino acid substitutions, there must be at least one continuously upward path to the top. Most protein sequences yield unfolded or otherwise inactive proteins so most of the landscape is a flat plain of low fitness. The rare mountains represent local optima of fitness consisting of functional, closely related proteins.

Directed evolution is a series of local searches on the fitness landscape. Upward steps or jumps are a possible route to the peak. One must measure improvements at each step to know that the path is still upward. Protein engineering differs from mountaineering because the landscape is invisible during protein engineering. One cannot see ahead to know if a detrimental substitution will eventually become useful. One must measure an improvement to know that one is on an upward path to the peak. The ruggedness of the mountains represents cooperativity

between amino acids. A smooth mountain represents no cooperativity between amino acid substitutions, while increasingly rugged slopes represent increasing cooperative behavior.
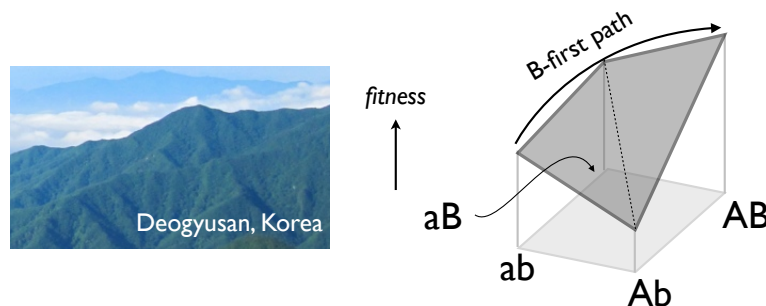
The independently-acting substitutions case corresponds the smooth landscape of a volcanic mountain like Mt. Fuji in Japan, Figure 10.2. There is one peak and the slope toward the peak is always upward. At any point on the slope, there are many possible directions upward and all paths lead to the peak. This landscape represents independently-acting substitutions. Each upward step (beneficial substitution) gets you closer to the peak. Many continuously upward paths to the top exist, which correspond to adding beneficial substitution in any order.



**Figure 10.2.** A smooth volcanic mountain like Mt. Fuji in Japan represents additive (independent) behavior of individual substitutions. It has many continuously upward paths leading to a single maximum. Moving from ab to AB can follow either the path that adds A first or the one that adds B first. Both paths lead upward at each step.
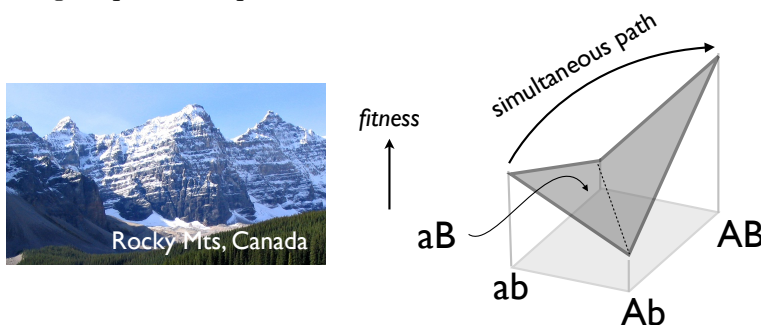
The first type of cooperativity, magnitude cooperativity, correspond to a slightly rougher landscape. The slope is still continuously upward to a single peak and has many paths upward, but some paths are more bumpy than others.

The second and third types of cooperativity corresponds to increased ruggedness like the old mountains in the Appalachians in the USA or the Baekdudaegan in Korea, Figure 10.3. The number of continuously upward paths to the peak is limited. Small peaks within this landscape represent false positives, where the individual steps appear beneficial, but do not lead to the highest peak. Valleys within this represent false negatives where the individual step leads downward, but can also eventually lead upward on the other side of the valley. An alternate route around the valley exists which avoids the downward steps. This alternate route represents a different order of adding the mutations. The rugged nature of this landscape limits the number of continuously upward paths to the peak. Some upward paths do not lead to the peak. Paths through a valley are inaccessible by single steps because one must continuously move upward.

**Figure 10.3.** An old, slightly rugged mountain like Deogyusan in South Korea represents a fitness landscape with second and third type of cooperativity. There are a limited number of continuously upward paths to the peak. Local peaks not leading to the maximum represent false positives, while valleys on the way to the peak represent false negatives. At least one path exists that avoids local peaks and valleys such at that each step is upward. For example, the B-first path is a possible route from ab to AB. Fitness increases upon adding B and further increases upon adding A. The A-first path is not a continuously upward path.

The most rugged landscape of newly formed mountains like the Rocky Mountains in Canada represent cooperative behavior where substitutions are beneficial only when both are present, Figure 10.4. Individually, they are detrimental. There is no path to the highest peak where each step is upward. One must descend into valleys or jump over them in order to reach the highest peak. This fitness landscape is the most difficult to explore since one cannot see where higher regions of the fitness landscape exist and one cannot find these regions by single substitutions. To explore this fitness landscape one must make multiple simultaneous substitutions, which correspond to leaps within the fitness landscape. By leaping around, one can cross a valley, then continue upward in single upward steps.



**Figure 10.4.** A highly rugged landscape like the Rocky Mountains in Canada represent cooperative behavior where several substitutions are required simultaneously to create a beneficial effect. There is no path that leads to the highest peak where each single step is upward. Reaching the highest peak by single steps requires descents into valleys or leaps across the valley, which represent several mutations added simultaneously. Substitutions A and B are individually detrimental, but when combined, they are beneficial. Only simultaneous addition of both A and B can discover this beneficial effect.

Real fitness landscapes for proteins are likely mixtures of different cases. Some amino acids pairs act cooperatively, while others independently. Like real mountains, fitness landscapes may contain both rugged and smooth regions. In the next sections, we consider ways to find beneficial combinations of cooperatively acting substitutions. Since these are more difficult than those for independently acting substitutions, researchers often use mixed approaches. One may first find beneficial independently-acting substitutions, then search for cooperative cases, then return to optimize with additional independently-acting substitutions.

## 10.3 Discovering and combining beneficial single substitutions

Different types of cooperativity require different approaches to find beneficial combinations, Table 10.3. In the simplest cases, independently acting substitutions or magnitude cooperativity one simply finds all the beneficial single substitutions from one library of single substitutions and

*combines* all of them. A single round of mutagenesis and screening suffices. If cooperativity creates false positives, then one *recombines* the beneficial single substitutions to find the best subset. The recombination requires, after identifying the beneficial single substitutions, the creation of a second library that recombines the beneficial single substitutions. Screening this library identifies the best subset. recombinants and screens. If cooperativity creates false negatives, then one must *accumulate* beneficial substitutions through multiple cycles of mutagenesis and screening. In addition, this accumulation may use parallel lineages of different sets of beneficial substitutions. This accumulation approach finds substitutions that act cooperatively with substitutions added in earlier rounds. Finally, screening libraries with *simultaneous multiple substitutions* finds sets that are beneficial only when combined. These library are vastly larger than libraries of single substitutions.

**Table 10.3.** Strategies to add single substitutions for different cooperativity cases.

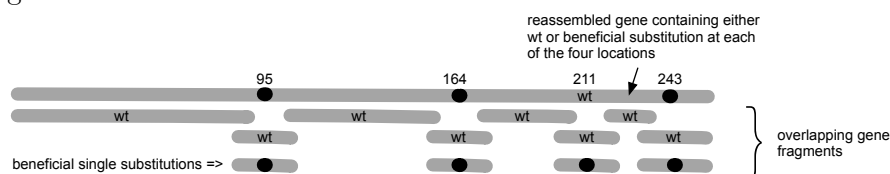| Expected fitness landscape | Effect of individual versus combined substitutions | Strategy to find beneficial multiple substitutions |
| --- | --- | --- |
| smooth, single peak | Beneficial effects of A & B add (or approximately add) in AB | Combine all beneficial single substitutions |
| rough with local peaks | Beneficial effects of A & B reverse to become detrimental in AB | Recombine beneficial single substitutions to find the best subset |
| rough with valleys | Beneficial effect of B only realized once A has been added | Accumulate single beneficial substitutions stepwise in multiple parallel lineages |
| highly rugged, multiple peaks & valleys | Beneficial effects of A & B only realized when both are present | Discover multiple substitutions simultaneously |

Each directed evolution experiment may use multiple approaches. First, one does not know if cooperativity is important or the type of cooperativity before starting the directed evolution experiment. One must chose an approach based on a best guess and be willing to change the approach if one does not find sufficient improvement. Second, the landscape will change as the fitness increases. There may be groups of substitution that act cooperatively along with other groups that act additively. Choosing the approach balances the effort involved with the need to find cooperative interactions.

**One-sign reversal: recombination eliminates false positives**

In this second type of cooperativity, the cooperativity causes the combination of two beneficial changes to be detrimental. One of the beneficial substitutions changes has changed sign to become detrimental. It is not possible to determine which beneficial substitution has changed its effect; in principle, both may have changed. This type of cooperativity is analogous to false positives in other experiments. Individually beneficial substitutions appear to contribute to an improved multi-substitution variant (a positive result), but when combined, they do not and were therefore false positives. The solution is to eliminate the substitutions that turn detrimental using recombination of beneficial single substitutions. Recombination tests different subsets of beneficial substitutions and eliminates unfavorable combinations. Recombination eliminates

'false positive' substitutions to find the subset of beneficial single substitutions that yields the highest benefit when combined.

To completely search all possible subsets, one creates a library that recombines all the beneficial single substitutions. Gene assembly PCR, also known as Gibson assembly, is a convenient way to create this library (Gibson et al., 2009), Figure 10.4. A PCR reaction connects overlapping DNA fragments to create the full length gene. For DNA fragments spanning the location of a beneficial mutation, one uses a mixture of two DNA fragments. One encodes the wild-type sequence, while the other encodes the beneficial substitution. The number of ways to recombine the single substitutions is $2^N$ where N is the number of individual beneficial changes. The 2 indicates that the combination can either contain or not contain the change. For example, if researchers identified ten beneficial single substitutions in a protein, then the number of possible variants is $2^{10}$ or 1,024. Screening this library will identify the best subsets of the ten beneficial single substitutions.



**Figure 10.5.** Gene assembly mutagenesis to recombine four beneficial single substitutions. The overlapping DNA fragments encode either the wild-type amino acid or the beneficial substitution identified in earlier experiments. In this hypothetical example the beneficial substitutions occurred at positions 95, 164, 211, and 243 (amino acid numbering). Reassembly of the gene incorporates either the wild type codon or the codon for the beneficial substitution yielding $2^4 = 16$ recombinant genes. The example reassembled gene shows one of these 16 with a wt amino acid at position 211 and beneficial single substitutions at positions 95, 164, and 243.

In the engineering of an enzyme for an enantioselective reaction, van der Meer et al. (2016) identified four beneficial single substitutions that favored the opposite enantiomer. The wild-type enzyme favored the 2R3S enantiomer by a factor of two, while the following single substitution variants favored the opposite 2S3R enantiomer: H6I (by a factor of 11.5 or a 23-folds shift from wild type), M45Y (by at least a factor of five; >10-fold shift from wild type), M45H (by a factor of five; 10-fold shift from wild type), F50A (by a factor of 13; 27-fold shift from wild type). Note that two of the beneficial substitutions are different amino acids at position 45. To combine the benefits of the substitutions, they recombined these substitutions to create the 5 possible double substitutions and the two triple substitutions. Three combinations showed low or no activity indicating unfavorable cooperative interactions. The other four favored the 2S3R enantiomer by a factor of 24, which is higher than the best single substitution (13), indicating at least partial additivity of the substitutions. The best variant was the double substitution of M45Y/F50A, which favored the 2S3R enantiomer by a factor of 99. The two residues lie nearby each other and the structure of this variant showed that the substitutions created a new substrate binding pocket.

Another example of this recombination approach is increasing the stability of a xylanase (Dumon et al., 2008). Xylanase cleaves the hemicellulose portion of lignocellulose, and is used in paper manufacture to help release cellulose fibers. First, saturation mutagenesis at all sites one at a time identified 69 single substitution variants with improved stability. Next, if the effects were purely additive, the researchers could simply combine the best ones. Instead, the researchers assumed that the effects would be non additive such that some combinations would decrease stability, so that a subset of beneficial single substitutions would be best. To find these best combinations they recombined the fifteen best single substitutions (twelve positions), which had melting temperatures 1-8 °C higher that wild type. The library contained 12,288 possible unique variants ($2^{10}×3×4$). There was a single beneficial substitution identified at ten positions, two beneficial substitutions at an eleventh position and three beneficial substitutions at a twelfth position. Screening the entire library identified a variant with seven amino acid substitutions and melted ~25 °C higher that wild type. An x-ray structure explained that the T13F (4.7 °C increase) substitution increased hydrophobic interactions and the S9P (7.5 °C increase) substitution decreased flexibility of the unfolded form, but did not explain how the other five substitutions increased stability.

In a third example, seven amino acid substitutions improved the promiscuous glycoaldehyde synthase activity of benzaldehyde lyase 70-fold, which corresponds to an average contribution of 1.8-fold increase for each substitution (Lu et al., 2019). Two of the substitutions were from a previous protein design of a homologous enzyme for a similar reaction (Siegel et al., 2015). The remaining five were discovered by first identifying fourteen beneficial single substitutions (saturation mutagenesis at 25 locations), then recombining these fourteen to find the best subset, which was a set of five substitutions.

As the number of beneficial single substitutions increases, the number of recombinants increases exponentially making it impractical to test all of them. For example, recombining 40 beneficial single substitutions requires testing $2^{40}$ ~ $10^{12}$ variants. To reduce the number of recombinants, one can assume that only a few substitutions are false positives and most will remain beneficial when combined. In the three examples above, the best recombinant contained about half of the beneficial single substitution: 2 of 4, 7 of 12, and 5 of 14. Thus, one can assume that the best recombinant will contain more than half of the substitutions. One can ignore those with fewer substitutions. For example, Hamamatsu and coworkers (2005) identified 14 single substitutions that increase the thermostability of prolyl endopeptidase. Recombining all 14 would create a library of $2^{14}$ or 16,384 variants. Instead of making a library containing all variants (single substitutions, two substitutions, etc.), they biased the library to contain an average of twelve substitutions by increasing the proportion of the DNA fragments encoding the substitution over the DNA fragments containing the wild-type amino acid. This smaller recombination library required screening only 2000 colonies to find a good one with twelve substitutions.

Another way to reduce the size of the recombination library is to use several cycles to accumulate substitutions. For example, one could choose ten of the 40 beneficial substitutions

and find the best combination by screening 1,024 variants, then add ten more substitutions and find the best combination again. Just four cycles of recombining 10 substitutions finds the optimum or nearly optimum combination for 40 substitutions by making and testing only 4000 variants instead of $10^{12}$ variants. This sparse search approach yields an improved variant with at least 98% of the best possible fitness gain (Fox & Huisman, 2008) because it accumulates beneficial substitutions.

The most drastic accumulation approach is stepwise addition of the beneficial substitutions starting from the most beneficial. Including the substitutions with the most benefit is best if possible. Each addition of a beneficial substitution should improve stability; if it does not, then it has become detrimental due to interaction with other substitutions (it is a false positive) and is removed. For example, Arabnejad et al. (2017) identified 23 locations where single substitutions stabilized halohydrin dehalogenase as measured by an increase in the apparent unfolding temperature of at least 1 °C. They combined 13 of the most stabilizing substitutions by adding them stepwise starting from the most stabilizing single substitution. This stepwise addition tested whether each new addition contributed to the the stability of the combined variant. Twelve of the thirteen increased the apparent unfolding temperature as they were accumulated in the combined variant, but one did not. This substitution (T134V, the fifth-most stabilizing) decreased the stability when added to the variant already containing four other stabilizing substitutions. This destabilizing substitution was removed yielding a twelve substitution variant, which had an apparent unfolding temperature 25.5 °C higher that the wild type protein. This approach required making thirteen successive site-directed mutations and testing only these thirteen. Since only thirteen recombinants out of a possible $2^{13} = 8,192$ were tested, the recombinant identified is unlikely the best possible one, but it may be good enough to reach the engineering goal.

**One-sign reversal: accumulation minimizes false negatives**

In the third type of cooperativity, a neutral or detrimental substitution (A) cooperates to become beneficial upon interactions with another beneficial substitution (B). Cooperativity reverses the effect of B from neutral or detrimental to beneficial. Because of this reversal, the order of discovery of the beneficial substitutions is critical. Substitution B cannot be discovered before substitution A. Substitution A must be present for B to be beneficial. This third type of cooperativity is analogous to false negatives. Individually, substitution B falsely appears non-beneficial. A common case of this third type of cooperativity is the combination of a stabilizing substitution with a destabilizing, but otherwise beneficial substitution. Once the stabilizing substitution compensates for the destabilizing effect, the beneficial effect can be observed.

In these false-negative-like cases, all beneficial substitutions cannot be discovered from the first library of variants. Instead, beneficial substitutions must accumulate stepwise using multiple rounds of mutagenesis and screening. Instead of finding all single beneficial substitutions independently, one accumulates single substitutions. One makes a first library of single substitutions variants and screens to find beneficial substitutions. Next, one makes additional libraries (parallel libraries) from some of these beneficial variants and screens again. Through

multiple rounds of library generation and screening one accumulate beneficial substitutions, some of which may act cooperatively with each other. The randomization approach can be error prone PCR, saturation mutagenesis, or recombination. The critical part is repeated cycles of mutagenesis and screening to accumulate mutations.

Within the fitness landscape metaphor, this third type of cooperativity limits the number of paths to the peak. The path to the combination must add B before A, so all paths that do not add B first are not accessible because A cannot be identified as beneficial.
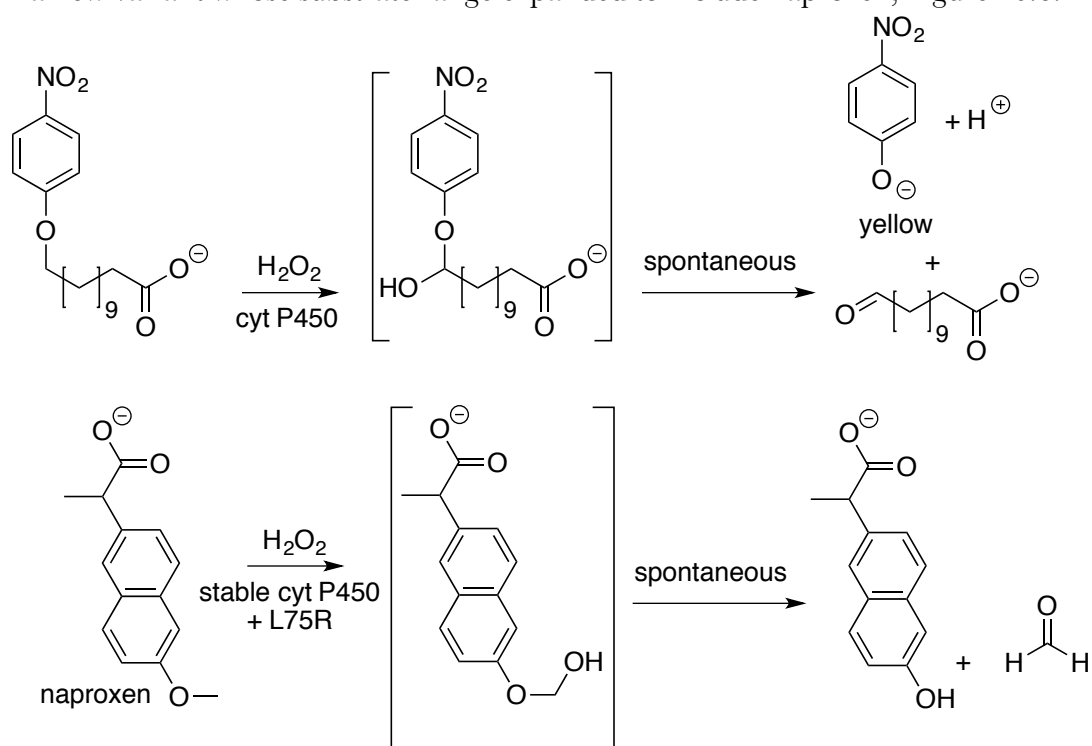
A natural example of the limited number of paths to a fitness peak is the accumulation of five mutations in a β-lactamase gene (Weinreich et al., 2005). Bacteria expressing the wild-type β-lactamase are sensitive to a third generation β-lactam antibiotic, cefotaxime, because the wild-type enzyme cannot catalyze hydrolysis of this β-lactam to inactivate it. The five mutations increase the resistance of the bacteria ~100,000-fold. All five of these mutations contribute to increasing the resistance. Four mutations caused amino acid substitution in the enzyme, while one mutation encodes the same amino acid, but alters the RNA sequence to increase the amount of protein produced.

There are 5! or 120 ways to introduce five mutations. There are five choices for the first mutation, four for the second and so on. For natural selection of a mutation as beneficial, it must increase the resistance of the bacteria to cefotaxime. An accessible path is one that increases the resistance at each step. Only 18 of the 120 possible paths were accessible. Most paths were inaccessible because at least one mutation step decreased resistance to cefotaxime. For example, mutation B did not affect the rate, but stabilized protein. The wild-type β-lactamase is sufficiently stable so mutation B alone does not increase resistance. Mutation B cannot be identified as beneficial by adding it to the wild-type enzyme so the path that starts with mutation B is inaccessible. Mutation A increased the reaction rate, but destabilized the protein. The overall effect was slightly beneficial and it could be selected. Once mutation A was present, mutation B became beneficial because it compensated for the destabilizing effect of mutation A. In another example, mutation C increased protein expression, but that increased expression led to aggregation and overall lower performance. Thus, by itself mutation C was deleterious. However, if the stabilizing mutation B were present, it reduced aggregation. The combination of mutations C and B gives more protein and avoids aggregation, so is beneficial.

Consider the word puzzle to convert the word 'WORD' to the word 'GENE' by changing one letter at a time, where each step must be an English word (Maynard Smith, 1970). All four letters must be changed, so four changes are needed. There are 4! or 24 possible paths to make this switch: four changes for the first step, three for the second, two for the third, and one for the last. For the first step, changing each of the four letters to the letter in GENE yields: GORD, WERD, WOND, WORE. Only the last one is an English word, so the first three paths are not possible. Similar reasoning shows that only one choice is possible at the second and third change so only one of the 24 paths yields a solution: WORD –> WORE –> GORE –> GONE –> GENE.

**Pre-stabilizing proteins to remove cooperativity due to stabilization.** A common case of false-negative-like cooperativity is the combination of a beneficial, but destabilizing, substitution with a stabilizing substitution. One way to eliminate this cooperativity is to start with an especially stable protein. The extra stability allows the protein to tolerate destabilizing substitutions, but still fold correctly. The extra stability removes the cooperative interaction; the beneficial substitution can now be discovered independently.

For example, Bloom and coworkers (2005) sought to expand the substrate range of a cytochrome P450 from fatty acid derivatives to naproxen, Figure 10.2. They screened random mutagenesis libraries generated from either a marginally stable cytochrome P450 enzyme or a pre-stabilized variant of the same enzyme. They hypothesized that the marginally stable P450 would require at least two substitutions: one to expand the substrate range and the second to stabilize the enzyme. Two beneficial substitution would like be rare in the library. In contrast, the thermostable starting point would not require an additional stabilizing substitution, only the substrate-expanding substitution is needed, so solutions would be more common. In agreement with their hypothesis, only the library derived from the thermostable variant contained a solution – a new variant whose substrate range expanded to include naproxen, Figure 10.6.



**Figure 10.6.** Cytochrome P450 enzyme catalyzes the hydroxylation of fatty acid derivatives, which leads to C–O bond cleavage in these ethers. Upon hydroxylation, the assay substrate (top reaction) releases the yellow *p*-nitrophenoxide. The L75R variant of cytochrome P450 also catalyzes the hydroxylation of naproxen (bottom reaction). The L75R substitution destabilizes the enzyme, so only a pre-stabilized variant of cytochrome P450 tolerates this substitution. The carboxyl groups of naproxen and the fatty acid derivative likely bind in the different places within the enzyme. The use of a pre-stabilized enzyme enabled discovery of an otherwise deleterious (destabilizing), but function-enabling substitution.

This cytochrome P450 enzyme catalyzes the hydroxylation of fatty acids near their hydrophobic end as they bind to the hydrophobic active site. The Leu75Arg substitution expands the substrate range of this cytochrome P450 to include naproxen. Hydroxylation of naproxen, a 2-arylpropionic acid (a non-steroidal anti-inflammatory; trade name Aleve) requires binding the negatively charged carboxylate in the hydrophobic pocket. The Leu75Arg substitution provides the compensating positive charge for this binding, but placing a charged amino acid within a hydrophobic pocket destabilized the enzyme. (The temperature that inactivates half of the protein within 10 min dropped by 8 °C with the Leu75Arg substitution.) The thermostable starting variant tolerated this substitution yielding a catalytically active variant with expanded substrate range. In contrast, the marginally stable cytochrome P450 variant did not tolerate the Leu75Arg substitution. Site-directed mutagenesis to introduce this substitution yielded only inactive, unfolded protein. Not surprisingly then, the error prone PCR library derived from the marginally stable enzyme did not contain any variants whose substrate range had expanded to include naproxen. In principle, a combination of Leu75Arg and a stabilizing substitution would have yielded a solution, but this combination is rare and much less likely than the Leu75Arg substitution alone.

Another example of directed evolution starting from a pre-stabilized enzyme is the engineering of phosphotriesterase to increase its reaction rate toward a nerve toxin, cyclosarin (Gupta et al., 2011). Five amino acid replacements within the substrate-binding site of the pre-stabilized phosphotriesterase increased the rate of hydrolysis 135-fold. Similarly, affinity-enhancing substitutions destabilized an antibody, so engineering increased affinity of this antibody required the addition of stabilizing substitutions (Julian et al., 2017). An alternative to pre-stabilization of proteins to minimize the effect of destabilizing substitutions is to start with a stable protein from a thermophile.

## 10.4 Multiple simultaneous substitutions

Multiple simultaneous substitution corresponds to jumps on the fitness landscape. These jumps can escape a local fitness maximum by crossing a fitness valley to reach a region of higher fitness. High mutation rates yield libraries with multiple simultaneous substitutions. For example. Boder et al. (2000) created libraries of variants of single chain antibody fragments using a error-prone PCR and recombination of DNA fragments. Screening for tighter binding of fluorescein-biotin identified variants containing an average of 1.7-4.5 amino acid substitutions. The best variant after four rounds of mutagenesis and screening contained twelve substitutions and bound the target 14,000-fold more tightly. The researchers achieved their goal of tighter binding and did not attempt to identify the roles of the individual substitutions.

Two advantages of making multiple substitution simultaneously are faster discovery of beneficial substitutions and the potential to discover more complex types of cooperativity. One disadvantage is the vast library sizes. In the example above, researchers screened the libraries of $10^8$ yeast cells using fluorescence-activated cell sorting. Another disadvantage is that strongly

detrimental mutations, such as those that prevent protein folding, mask the beneficial effects of other mutation within the same protein.

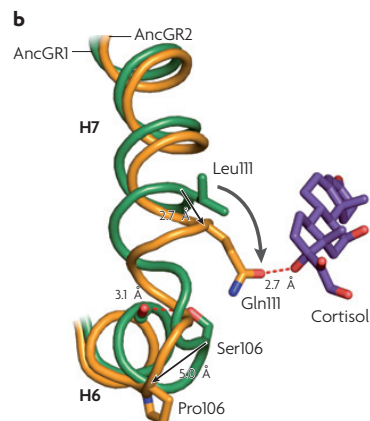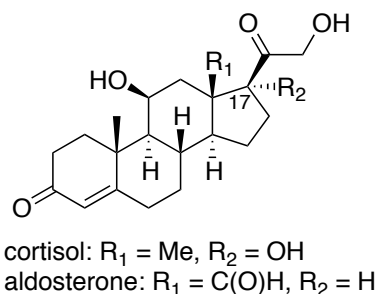**More complex cooperativity: two-sign reversal**

The most difficult type of cooperativity to discover is where detrimental or neutral changes combine to create a benefit. This fourth type of cooperativity is called two-sign reversal, Table 10.2 above. These two-sign reversals are impossible to discover stepwise because the neither substitution is beneficial individually; only the combination is beneficial. Discovering these beneficial combinations requires multiple simultaneous substitutions, which is difficult because the number of possible multiple simultaneous substitutions is vast. Two simultaneous substitutions in a 300-aa protein create ~16 million possibilities.

In principle, another type of two-sign reversal is possible where two beneficial substitutions cooperate to become detrimental. This another example of a false positive, which can be removed by recombination of beneficial substitutions as described in the previous section. This section focuses on the difficult to discover case: a two-sign reversal where two neutral or detrimental substitution cooperate to create a beneficial effect.

A simple example of a two-sign reversal is an improved lock & key. Changing either the lock or the key causes a misfit, but changing both can improve it. Many interactions within proteins require a similar complementarity. Binding a protein interface to another requires both to match, folding a protein requires matching of shape and interaction between secondary structures elements like helices. Reshaping a binding site can require both more space on one side and less space on the other side. Thus, one expects two-sign reversal cooperativity to be common in proteins.

An example of two-sign reversal of the effect of substitutions is the change in ligand preference of the steroid hormone receptor, Figure 10.6 (Ortlund et al., 2007). The original ligand is aldosterone, which has a hydrogen at $R_2$, while the new ligand is cortisol, which has a hydrogen bond donor (hydroxyl) at $R_2$. Leu11Gln and Ser106Pro act cooperatively to switch the ligand preference.

**Figure 10.6.** The change in ligand preference of steroid receptor from aldosterone to cortisol requires two amino acid replacements. a) Cortisol has a hydrogen-bond donor (hydroxyl) at R2, while aldosterone does not. b) An overlay of X-ray crystal structures show the two amino acid substitutions needed to increase the binding of cortisol. The Leu111Gln substitution adds a hydrogen bond acceptor, while the Ser106Pro substitution distorts the helix to position it for hydrogen bonding. Panel b is from Ortlund et al., 2007.

Leu111Gln has little effect on the receptor–ligand interactions so is a neutral change. Ser106Pro destabilizes the interactions of all ligands so is a deleterious change. The combination, however, changes the ligand preference. The Leu111Gln adds a hydrogen bond acceptor in the side chain, while Ser106Pro introduces a kink into the protein main-chain, repositioning a helix that borders the ligand-binding pocket. The change moves the amino acid at position 111 bringing it closer to the hydroxyl group at $R_2$ of cortisol. The cooperation between the proline substitution to shift helix and the glutamine substitution to form the H-bond changes the ligand preference.

Table 10.2 above defines cooperativity between pairs of amino acids, but more complex cooperativity is possible between three or more amino acids. The simultaneous multiple substitutions can also find these more complex interactions. A word puzzle analogy of a four-sign reversal is converting the word ALSO to the word GENE via other English language words. No single substitution of the letters in GENE into the ALSO yields an English word: GLSO, AESO, ALNO, ALSE. Similarly, double substitutions (GESO, GLNO, GLSE, AENO, AESE, ALNE) and even triple substitutions (GENO, GLNE, GESE, AENE) do not yield English words. Only four simultaneous substitutions of letters can convert the word ALSO to the word GENE.

Searching for cooperative behavior like two-sign reversals requires simultaneous multiple substitutions.  The two problems with high mutation rates in libraries are 1) the masking of beneficial effects by strongly deleterious substitutions and 2) the astronomically large numbers of possible variants. The next two sections describe strategies to limit these problems.

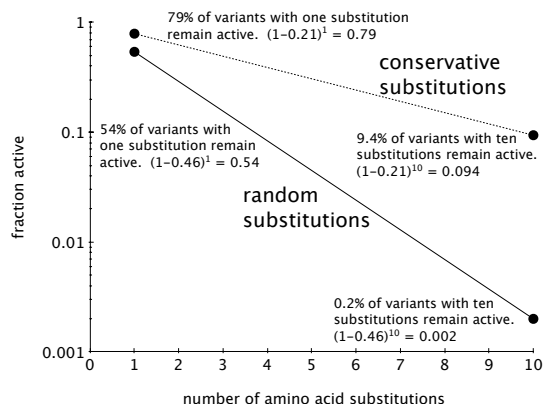### 10.4.1 Minimize strongly deleterious substitutions

Strongly deleterious substitutions are those that disrupt protein function by preventing protein folding or other essential steps like binding or catalysis. If the variant contains one or more strongly deleterious substitution, then any beneficial effects of the other substitution are hidden by this disruption of function. Avoiding strongly deleterious substitutions reveals the effect of other substitutions to allow discovery of beneficial substitutions.

Typically, 30–50% of single amino acid substitutions are strongly deleterious (Guo et al., 2004; Drummond et al., 2005). This fraction varies for each enzyme and with the precise definition of strongly deleterious. For the purposes of protein engineering, a drop to <10% of wild-type activity is usually considered strongly deleterious.

As the number of substitutions in a variant increases, the likelihood that at least one of those substitutions is strongly deleterious increases exponentially, eq. 10.4. In other words, the fraction of catalytically active variants decreases exponentially as the number of substitutions increases. This equation assumes that the strongly deleterious substitutions act independently so that the fraction remains constant as the substitutions accumulate.

$$\text{fraction active} = (1\text{-}F_{\text{det}})^m \tag{10.4}$$

Here $F_{\text{det}}$ is fraction of random substitutions that are strongly deleterious and $m$ is the number of random substitutions. For a single substitution, the chance of a detrimental substitution is $F_{\text{det}}$ and the chance of not having a detrimental substitution is $(1\text{-}F_{\text{det}})$. For the second substitution, the chance of not picking a detrimental substitutions remains the same, $(1\text{-}F_{\text{det}})$, so the change of not picking a detrimental substitution over both the first and second is $(1\text{-}F_{\text{det}})^2$. In the general case, the chance of not picking a detrimental substitution over $m$ picks in $(1\text{-}F_{\text{det}})^m$. Since $1\text{-}F_{\text{det}}$ is less than one, the fraction folded decreases as m increases. For β-lactamase TEM-1, Drummond and coworkers (2005) measured that 46% of random substitution were strongly deleterious ($F_{\text{det}} = 0.46$). Equation 10.4 predicts that only 4.6% of variants containing five random substitutions would retain catalytic activity, only 0.2% of variants containing ten random substitutions and 0.0004% of those containing twenty substitutions.

**Figure 10.7.** The fraction of active proteins decreases exponentially (logarithmic y-axis scale) as the number of substitutions increases (linear x-axis scale). For β-lactamase TEM-1, 46% of random single substitutions were strongly deleterious, so only 54% of the variants with one substitution retained activity. Among variants containing ten substitutions, only 0.2% retained activity according to eq. 10.4 because the others contained at least one strongly detrimental substitution. Conservative substitutions (amino acids that occur in homologs) are less likely to be deleterious. For β-lactamase TEM-1, only 21% of single amino acid exchanges with a homolog were strongly deleterious. Among variants containing ten substitutions, a 47-fold higher fraction of variants with ten substitutions (9.4%) retain activity according to eq. 10.4 when the substitutions are conservative as compared to random.
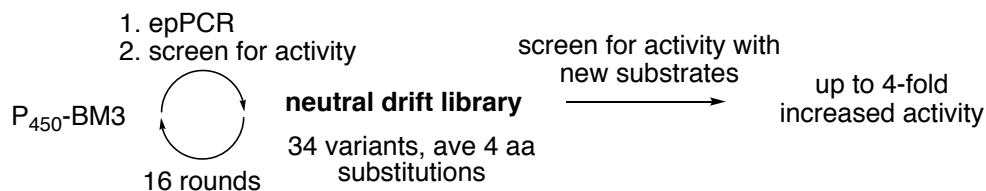
Several different approaches can reduce the probability that variants in a library contain a strongly detrimental substitution, Table 10.4. These approaches are discussed in the following sections.

**Table 10.4.** Limiting strongly deleterious substitutions when adding multiple substitutions.

| approach | substitutions | rationale & assumptions |
|---|---|---|
| screen & remove inactive variants (neutral drift library) | no restrictions, typically error-prone PCR (low # of substitutions, ~5) | - eliminates deleterious substitutions by requiring original activity be retained |
| use only conservative substitutions (DNA shuffling) | only aa found in homologs (high # of substitutions, 10-100 | - substitutions found in homologues are less likely to be deleterious than random substitutions |
| gene assembly, statistical analysis | typically 10-20 | - statistical analysis identifies beneficial & deleterious substitutions |

**Neutral drift libraries.** A direct, but tedious, way to eliminate strongly deleterious substitutions is to screen the variants and discard the inactive ones. Such a sifted library is called a neutral drift library, which refers to the natural process where mutations give rise to random substitutions while natural selection eliminates the strongly deleterious substitutions. Thus, natural populations of genes contain many variations that retain the activity needed for the fitness of the organism.

Neutral drift libraries in laboratory evolution are collections of variants that maintain a minimum activity defined by the researcher. For example, Bloom and coworkers (2007) created a neutral drift library of P450 monooxygenases by random mutagenesis using epPCR (an average of 1.4 nucleotide substitutions) followed by screening for a defined minimum activity with a convenient test substrate, Figure 10.8. The test substrate was a *p*-nitrophenoxy derivative whose reaction was conveniently monitored by a color change. Screening the neutral drift library for activity toward new substrates identified variants with up to four-fold increased activity for the new substrates.

**Figure 10.8**. Sixteen rounds of random mutagenesis (average of 1.4 nucleotide substitutions) of a $P_{450}$-BM3 peroxygenase followed by screening for activity toward a convenient test substrate (*p*-nitrophenoxy derivative in top equation of Figure 10.2) created a neutral drift library of variants. Approximately half of the variants maintained the defined minimum activity toward the test substrate at each round. Next, the neutral drift library was screened for activity toward new substrates, which involved a more complex assay.

Protein engineering rarely uses neutral drift libraries because generation of the libraries is tedious. It requires screening for retained activity as well as sequencing to eliminate wild-type and duplicates. Unlike most libraries in laboratory evolution, neutral drift libraries are usually small, only 34 variants in the example above. If the goal is activity toward many new substrates, then prescreening for activity by creating a neutral drift library may save time. However, if the goal is activity toward one new substrate, then it may be simpler to directly screen for the desired activity.

**Limit substitutions to conservative ones (DNA shuffling).** Conservative replacements are those amino acids that occur at the corresponding locations in homologs. Exchanging substitutions between homologs is less likely to be strongly deleterious than random substitutions. Substitutions are strongly deleterious due to incompatibles with the existing amino acids in the protein. Since homologs have amino acids in common, exchanging a substitution between them has fewer opportunities for incompatibilities. For example, a random substitution in a 300-amino acid protein creates 299 new pairwise interactions between the new amino acid and the 299 remaining amino acids. Some of these interactions may be strongly deleterious. In contrast, exchanging a substitution between two homologous 300-amino acid proteins that differ by 100 amino acids creates only 99 new pairwise interactions. There are fewer differences between the proteins so fewer new pairwise interactions. There are fewer opportunities for incompatibilities when exchanging amino acids between homologs than when making random substitutions. (Not all pairwise interactions are equally important. Nearby amino acids are more likely to create incompatibilities than distant amino acids, but if the differences are randomly distributed throughout the protein, then the average conclusion should hold.)

For example, the two β-lactamases PSE-4 and TEM-1 share 40% identical amino acids and their structures are similar. While 46% of random substitutions were strongly deleterious in these proteins, only 21% of the exchanges between the two homologs were strongly deleterious (Drummond et al., 2005). Because of the exponential relationship in equation 10.4, this approximately two-fold shift in the fraction of strongly deleterious substitutions creates large differences in the fraction folded with multiple substitutions, see Figure 10.7 above. For five exchanges between homologs, equation 10.4 predicts that 31% would avoid a strongly deleterious substitution, which is 6.7-fold larger than the 4.6% when making random substitutions. Similarly, with ten exchanges between homologs, 9.4% remain folded, which is 47-fold more than with random substitutions and with twenty exchanges between homologs 0.9% remain folded, which is 2000-fold more than with random substitutions.

Two ways to limit the substitutions to conservative substitutions are first, to directly exchange fragments between several homologs by recombination of homologous genes (DNA shuffling) or,

second and more commonly, to assembly new genes from synthetic DNA fragments where a multiple sequence alignment identified the conservative replacements at each location (synthetic DNA shuffling).

Homologous recombination is the exchange of fragments between two similar DNA strands. This exchange creates substitutions at the locations that differ between the DNA strands. Sexual reproduction involves homologous recombination of parental DNA. The progeny inherit DNA consisting of recombined fragments of the parents. Homologous recombination in protein engineering, usually called DNA shuffling, also involves the exchange of DNA fragments, but it occurs in vitro, may have multiple starting strands (parents) and the genes may come from different species, Figure 10.9.



**Figure 10.9** Recombination of two homologous genes ($h = 2$) at two crossover points (three segments, $n = 3$) yields six recombinant genes; $2^3 - 2 = 6$.

The recombined genes from different species and the corresponding proteins are called chimeras. In Greek mythology, the chimera was a hybrid creature. It was a fire-breathing she-monster usually represented as a composite of a lion, goat, and serpent. Using this word for biological materials emphasizes that the genes or fragments of genes come from different species. This association of recombinations with a monster has contributed to public resistance to genetic engineering technology. Viewing these recombined genes as a set of conservative substitutions would create a more benign image of the same thing.

The number of possible recombined genes or chimeras depends on the number of parent sequences, $h$, and the number of segments, $n$, eq. 10.5.

$$\text{number of chimeras} = h^n - h \tag{10.5}$$

For example, two parents ($h = 2$), fragmented into thirds ($n = 3$), yields six unique chimeras: $2^3$-2 = 6, Figure 10.9 above.

Besides the advantage of avoiding strongly deleterious substitutions by making conservative substitutions, recombination can also remove detrimental substitutions and combine beneficial ones (McDonald et al., 2016). Consider a gene with one beneficial and one detrimental mutation. Random substitutions are unlikely to improve this gene since beneficial mutations are rarer than detrimental ones. However, recombination can improve this gene by separating the beneficial

substitution from the detrimental one. Recombination with another gene yields chimeras with neither substitution (25% chance), both substitutions (25%), only the beneficial one (25%) or only the detrimental one (25%). This recombination has a good chance (25%) of creating an improved variant by separating the beneficial mutation from the detrimental one. Next, consider two genes containing different beneficial substitutions. Random mutations are unlikely to add the second beneficial mutation since the beneficial mutations are rare. Recombination of the two genes, however, has a 25% chance to create chimeras with both beneficial mutations.

The first example of DNA shuffling (Stemmer, 1994) used a single gene, but later examples used multiple genes as described above (Ness et al., 1999). The crossover location may be chosen at specific sites (Heinzelman et al., 2009) or it may be random. The random DNA shuffling involves digestion of the parental DNA with an non-specific endonuclease followed by reassembly of the gene using PCR. The switch from one gene to another occur at homologous regions where a fragment from a different gene serves as the template for DNA synthesis.

**Smart libraries.** Smart libraries assume that substitutions that occur in one sequence will be tolerated in related sequences. Multiple sequence alignments yield a set of possible substitutions at each position. These substitutions are fewer than all possible substitutions, but are more likely to yield correctly folded, stable and catalytically active enzymes because these substitutions occur in related proteins. To make only these substitutions, one can use degenerate codons. For example, the degenerate codon NDT mentioned above encodes for the amino acids: Phe, Leu, Ile, Val, Tyr, His, Asn, Asp, Cys, Arg, Ser, and Gly. This approach increases the quality of the library. The hypothesis is that substitutions that occur rarely in nature are likely to be deleterious. AA-Calculator is a web tool to identify degenerate codons: http://guinevere.otago.ac.nz/stats.html (Firth & Patrick, 2008).

Directed evolution of prolyl endopeptidase to resist cleavage by pepsin used the smart library approach. Digestive proteases cleave gluten to proline-rich peptides. If people develop allergic reactions to these peptides, they become gluten intolerant and must avoid gluten. A potential treatment is dietary supplement with prolyl endopeptidases that can fragment the proline-rich peptides and prevent an allergic reaction. These prolyl endopeptidases must avoid cleavage by pepsin in the stomach. Khosla and co-workers (Ehren et al., 2008) aligned the amino acid sequences of 100 peptidase homologs and identified 30 specific, potentially beneficial substitutions. DNA synthesis yielded the genes for 47 variants that contained different combinations of these 30 amino acid substitutions. Testing these 47 variants, followed by statistical analysis identified 22 of the substitutions as beneficial. A second round of variants and testing yielded a five-amino-acid-substitution variant with 200-fold increased resistance to pepsin.

While smart libraries come from analysis of multiple sequence alignments, similar libraries can also come from computational design. Guntas and colleagues (2010) evolved a protein to bind its non-natural partner. The interface between the two proteins consisted of 13 amino acid positions. The theoretical number of possible sequences is $\sim 10^{17}$. They constructed this naïve library, but could screen only a tiny fraction, $\sim 10^{7}$. This approach found only three improved

variants, with only a three-fold increase in binding, likely because most of the variants screened contained at least one strongly deleterious substitution. Next, they use computational methods to identify positions more tolerant to mutation, without any attempt to predict specific favorable interactions with the binding partner. The role of this computation was similar to that of multiple sequence alignments in the design of a smart library: to reduce the probability of a strongly deleterious mutation. The theoretical number of variants in the designed library was much smaller, ~$10^8$. Screening ~$10^7$ variants identified a number of improved variants with up to 10,000-fold improved affinity. Enriching the library with well-folded sequences was sufficient to identify tightly-binding variants without predictions of specific binding interaction with the partner.

### 10.4.2. Limit the search for cooperative interactions

Current in-vitro directed evolution technology cannot find multiple simultaneous substitutions because the libraries are too large. Three simultaneous substitutions in a 300-aa protein require 30 billion variants for a complete library, while four simultaneous substitutions require 40 trillion. These are impossible numbers for any technique requiring plasmid transfer into bacteria. Although molecular biology methods can create the required numbers DNA molecules, plasmid transfer is limited to approximately a million transformants. In practice this number is often lower.

**Limit locations to pairs within the active site.** A second approach to limit the site of libraries is to assume that cooperativity is most likely between nearby residues within the active site. Rather than testing all possible pairs of substitutions, one tests only adjacent pairs within the active site. For example, the combinatorial active site saturation test (CASTing) assumes that amino acids adjacent to each other in space and nearby in sequence are most likely to to act cooperatively. For this reason, pairs of nearby amino acids are randomized simultaneously. More distant amino acids, for example on the other side of the active site, are assumed to act additively, so improved pairs can be added stepwise around the binding pocket (Boccola et al., 2005). For example, one pair of substitutions acted cooperatively to increase the activity of a dehydrogenase (LeADH) toward a pharmaceutical intermediate, while this same pair acted additively with another substitution (Jiao et al., 2016). Substitutions Asn235His and Pro236His in LeADH had little effect separately (2.2 and 2.0 U/mg, respectively, as compared to 2.1 U/mg for wild type), but when combined showed a cooperative benefit of a 2.6 fold increase to 5.5 U/mg. A more distant substitution, Ile87Phe, increased the activity 4.5 fold to 9.6 U/mg. Combining all three substitutions yielded a 12-fold improvement over wild type to 27 U/mg showing that the Asn235His and Pro236His pair acted additively with the Ile87Phe substitution (4.5 * 2.6 = 12). The cooperatively acting substitutions (235, 236) were adjacent to each other, while the additively acting substitution (87) was on the other side of the active site. In another example, researchers expanded substrate binding site of a monoamine oxidase to fit a larger substrate (Rowle et al., 2012). They simultaneously varied the amino acids at Phe210 & Leu213 and found a beneficial double substitution that increased reactivity 100-fold. The researchers searched for

substitutions at Phe210 and Leu213 only in combination, not separately, because they expected these nearby amino acids to act cooperatively as a set. They treated this pair of beneficial substitutions as a single set. In another experiment, they simultaneously varied another pair of nearby amino acids – Met242 and Met246 – and found another beneficial double substitution that improved reactivity 330-fold.

Phe210Leu/Leu213Thr = 100-fold improvement

Met242Gln/Met246Thr = 330-fold improvement

Phe210Leu/Leu213Thr/Met242Gln/Met246Thr = 990-fold improvement

The combination of these two sets of two substitutions made a four substitution variant. If the beneficial effects acted independently, the four substitution variant would be 33,000-fold better, but it was only 990-fold faster than wild type. Thus, negative cooperativity between the two sets of beneficial double substitutions reduced the benefit of the combination.

**Multiple locations within the active site, but limited replacement possibilities.** Cooperativity can involve more than pairs of amino acids. For example, the catalytic triad of serine proteases requires all three members of the triad for efficient catalysis. To find cooperativity between multiple sites, one must test multiple sites simultaneously, but to reduce the number of possibilities one can limit the number of replacement amino acids. For example, Sandström et al. (2012) increased the size of a binding site tunnel by simultaneously mutating nine sites. To keep the library size manageable, the replacement amino acids were limited to one or a few amino acids. In a similar example, Sun and coworkers (2016) further assumed that only some of the sites act cooperatively so that the cooperating groups could be optimized separately. To change the enantioselectivity of an epoxide hydrolase, the researchers focused on the ten residues forming the substrate binding site and allowed only Val, Phe, or Tyr as replacement amino acids. These assumptions limit the number of variants to $10^4$, which is still a large number. They divided the ten residues into three groups and optimized each group separately. This division drops the number of variants to $3^4 + 3^4 + 4^4 = 418$. This approach could find cooperative effects within the groups of three or four positions, but not between groups.

# 10.5 Nature's evolution

In most ways, laboratory evolution mimics nature's evolution. Both create variants and select the better ones. But Nature's evolution differs in subtle ways from laboratory evolution. When little is known about potential solutions, Nature's approach may better because it efficiently explores variants and accumulates beneficial mutations. When an approximate solution is known, then laboratory evolution may be better because the experimenter can focus on mutations most likely to be beneficial. In addition, laboratory evolution can select for an artificial definition of fitness, while nature only selects for survival and reproduction.

Replacing 8-14% of the amino acids in a protein during protein engineering can be viewed as the equivalent of converting mouse proteins into human proteins since the amino-acid sequences of similar proteins in mice and human typically differ by 13% (Mouse Genome Sequencing Consortium, 2002). This protein engineering is equivalent to compressing the 75,000,000-yr evolution of an early mammal into modern-day mice and humans into several months of laboratory work.

### Nature's evolution acts on populations

One difference between laboratory and Nature's evolution is that laboratory evolution starts with a few variants, while Nature's evolution starts with large populations of variants. One liter of bacterial culture contains ~10 trillion individual bacteria. In nature, random changes and selection continue indefinitely. Most changes are neutral, that is, they do not affect function, so a population accumulates many variations of a protein. Natural populations are a neutral drift library. Genes in natural populations of microbes vary extensively. DNA sequence comparison of independently isolated *E. coli* strains revealed extensive genetic variation (Dixit et al., 2015). Pairs of aligned 1000-bp segments differed by 4-25 nucleotide substitutions. These differences are due to both copying errors when bacteria divide (similar to errors in error-prone PCR) and due to recombinations between bacteria (similar to gene shuffling).

These different starting points create vastly larger libraries in nature. For example, creating an average of one substitution in a 300-aa protein starting from a single protein yields 5700 variants. If the starting point is not a single protein, but a million variants, then the same random substitution creates a million-fold more variants. The small starting population in directed evolution, known as a bottleneck in evolutionary biology, limits the possible solutions. Each subsequent round of laboratory evolution creates another population bottleneck since experimenters pass only one or a few variants to the next round of evolution. These multiple population bottlenecks limit the genetic variation in the experiment.

These larger libraries make it more likely to find cooperative interactions. For example, substitution A may require substitution B to be effective. B might be a particular stabilizing substitution. If you make substitution A in large population, then it is likely that one of the variants already contains substitution B and you will find this cooperative, beneficial effect. The likelihood that a single starting point has substitution B is lower. This ability to find cooperative effects allows Nature's evolution to avoid being stuck at a local optimum in sequence space and more likely to find new optima to explore. Nature's proteins are more robust (able to withstand changes) than laboratory proteins because Nature starts with a large pool of variants.

In most cases, this smaller genetic variation in directed evolution is undesirable because it excludes possible solutions. However, the population bottleneck can be an advantage if the starting variants are close to a suitable solution. For example, starting with a variant with good activity to the desired substrate is more likely to yield an efficient enzyme than a large population with little or no activity. Similarly, if the starting variant is especially stable, then most new substitutions, even the destabilizing ones, will yield folded proteins.

### Nature selects for organism fitness at each step

Nature's evolution is limited by historical constraints because selection favors fitness. Evolution adapts existing features like the protein fold to new situations. In fitness landscape terms, natural evolution remains on and optimizes a single fitness peak.

In contrast, during laboratory evolution, the experimenter can chose less fit individuals as starting points for the next round of evolution. For example, an experimenter could select unfolded, inactive enzymes during evolution with the goal of changing the protein fold. Thus, laboratory evolution can yield larger changes thereby shifting to a new fitness peak. (Need example of laboratory evolution that passes through less fit variants).

Nature's evolution also can't select for non-natural properties because they do not provide a fitness advantage. For example, nature cannot select for increased reactivity or enantioselectivity toward a synthetic intermediate, stability or activity in an organic solvent or other unnatural conditions.

### Nature's evolution has a low mutation rate

Low mutation rates change only one nucleotide within a codon, which yields conservative amino acid substitutions. Higher mutation rates can change several nucleotides in the same codon, which yields non-conservative substitutions. In addition, laboratory evolution can focus changes to particular regions, often the active site, using methods like saturation mutagenesis.

For example, researchers identified a pair of residues (Lys211 and Arg212) in the protease subtilisin where mutations significantly increased thermostability. The most stable variants identified by saturation mutagenesis and screening were nonconservative replacements with hydrophobic residues (Pro/Ala, Pro/Val, Leu/Val, and Trp/Ser). These substitutions required multiple (two to three) nucleotide substitutions in a single codon, which are extremely rare in a point mutation library.

The limitation of natural evolution is that it requires a long time. For example, nature evolved bacteria that can degrade the herbicide atrazine (Seffernick & Wackett, 2001), but this evolution required ~30 years. Initially, no bacteria could degrade atrazine. After 30 y, several groups identified bacteria with nine amino acid substitutions in melamine deaminase enzyme, which now catalyzed the dechlorination of atrazine. The exact evolutionary path of the conversion is not known. The low natural mutation rate required 30 y to create enough variants to include the beneficial combination.

evolution = generation of variants combined with selection for fitness

directed evolution mimics natural evolution by the creation of variants and selection for fitness.

In natural evolution, variant generation and selection occur simultaneously

In directed evolution, they occur stepwise.

the consequence of this difference is

# References

H. Arabnejad, M. Dal Lago, P. A. Jekel, R. J. Floor, A.-M. W. H. Thunnissen, A. C. Terwissscha van Scheltinga, H. J. Wijma, D. B. Janssen (2017) A robust cosolvent-compatible halohydrin dehalogenase by computational library design. *Protein Eng. Des. Sel.* **30**, 173–87. https://doi.org/10.1093/protein/gzw068

K. A. Bava, M. M. Gromiha, H. Uedaira, K. Kitajimi, A. Sarai (2004) ProTherm, version 4.0: thermodynamic database for proteins and mutants. *Nucleic Acids Res.* **32**, D120-1.

J. D. Bloom, P. A. Romero, Z. Lu, F. H. Arnold (2007) Neutral genetic drift can alter promiscuous protein functions, potentially aiding functional evolution. *Biol. Direct* **2**, 17; https://doi.org/10.1186/1745-6150-2-17

J. D. Bloom, J. J. Silberg, C. O. Wilke, D. A. Drummond, C. Adami, F. A. Arnold (2005) Thermodynamic prediction of protein neutrality. *Proc. Natl. Acad. Sci. USA* **102**, 606–11.

E. T. Boder, K. S. Midelfort, K. D. Wittrup (2000) Directed evolution of antibody fragments with monovalent femtomolar antigen-binding affinity. *Proc. Natl. Acad. Sci. U. S. A.* **97**, 10701–5. https://doi.org/10.1073/pnas.170297297

E. Dellus-Gur, Á. Tóth-Petróczy, M. Elias, D. S. Tawfik (2013) What makes a protein fold amenable to functional innovation? Fold polarity and stability trade-offs. *J. Mol. Biol.* **425**, 2609–21. https://doi.org/10.1016/j.jmb.2013.03.033

P. D. Dixit, T. Y. Pang, F. W. Studier, S. Maslov (2015) Recombinant transfer in the basic genome of *Escherichia coli*. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 9070-5. https://doi.org/10.1073/pnas.1510839112

D. A. Drummond, J. J. Silberg, M. M. Meyer, C. O. Wilke, F. H. Arnold (2005) On the conservative nature of intragenic recombination. *Proc. Natl. Acad. Sci. USA* **102**, 5380–5.

C. Dumon, A. Varvak, M. A. Wall, J. E. Flint, R. J. Lewis, J. H. Lakey, C. Morland, P. Luginbühl, S. Healey, T. Todaro, G. DeSantis, M. Sun, L. Parra-Gessert, X. Tan, D. P. Weiner, H. J. Gilbert (2008) Engineering hyperthermostability into a GH11 xylanase is mediated by subtle changes to protein structure, *J. Biol. Chem.* **283**, 22557–64.

J. Ehren, S. Govindarajan, B. Moron, J. Minshull, C. Khosla (2008) Protein engineering of improved prolyl endopeptidases for celiac sprue therapy. *Protein Eng. Des. Sel.* **21**, 699-707.

A. E. Firth, W. M. Patrick (2008) GLUE-IT and PEDEL-AA: new programmes for analyzing protein diversity in randomized libraries. *Nucl. Acids Res.* **36**, W281–5; https://doi.org/10.1093/nar/gkn226

R. J. Fox (2005) Directed molecular evolution by machine learning and the influence of nonlinear interactions. *J. Theoretical Biol.* **234**(2), 187–99. https://doi.org/10.1016/j.jtbi.2004.11.031

R. J. Fox, S. C. Davis, E. C. Mundorff, L. M Newman, V. Gavrilovic, S. K. Ma, L. M. Chung, C. Ching, S. Tam, S. Muley, J. Grate, J. Gruber, J. C. Whitman, R. A. Sheldon, G. W. Huisman (2007) Improving catalytic function by ProSAR-driven enzyme evolution. *Nat. Biotechnol.* **25**, 338–44.

R. J. Fox, G. W. Huisman (2008) Enzyme optimization: moving from blind evolution to statistical exploration of sequence–function space. *Trends Biotechnol.* **26**, 132–8. https://doi.org/10.1016/j.tibtech.2007.12.001

D. G. Gibson, L. Young, R.-Y. Chuang, J. C. Venter, C. A. Hutchison, H. O. Smith (2009) Enzymatic assembly of DNA molecules up to several hundred kilobases. *Nat. Meth.* **6**, 343–5. https://doi.org/10.1038/nmeth.1318

G. Guntas, C. Purbeck, B. Kuhlman (2010) Engineering a protein-protein interface using a computationally designed library. *Proc. Natl. Acad. Scl. U. S. A.* **107**, 19296–301; http://doi.org/10.1073/pnas.1006528107

H. H. Guo, J. Choe, L. A. Loeb (2004) Protein tolerance to random amino acid change. *Proc. Natl. Acad. Sci. U. S. A.* **101**, 9205–10; https://doi.org/10.1073/pnas.0403255101

R. D. Gupta, D. S. Tawfik (2008) Directed enzyme evolution via small and effective neutral drift libraries. *Nat. Meth.* **5**, 939–42.

R. D. Gupta, M. Goldsmith, Y. Ashani, Y.Simo, G. Mullokandov, H. Bar, M. Ben-David, H. Leader, R. Margalit, I. Silman, J. L. Sussman, D. S. Tawfik (2011) Directed evolution of hydrolases for prevention of G-type nerve agent intoxication. *Nat. Chem. Biol.* **7**, 120–5.

N. Hamamatsu, T. Aita, Y. Nomiya, H. Uchiyama, M. Nakajima, Y. Husimi, Y. Shibanaka (2005) Biased mutation-assembling: an efficient method for rapid directed evolution through simultaneous mutation accumulation. *Protein Eng. Des. Sel.* **18**, 265–71; https://doi.org/10.1093/protein/gzi028

P. Heinzelman, C. D. Snow, I. Wu, C. Nguyen, A. Villalobos, S. Govindarajan, J. Minshull, F. H. Arnold (2009) A family of thermostable fungal cellulases created by structure-guided recombination. Proc Natl Acad Sci USA 106, 5610–5.

C. C. Hsu, Z. Hong, M. Wada, D. Franke, C. H. Wong (2005) Directed evolution of D-sialic acid aldolase to L-3-deoxy-manno-2-octulosonic acid (L-KDO) aldolase. *Proc. Natl Acad. Sci. USA* **102,** 9122–6.

X.-C. Jiao, Y.-J. Zhang, Q. Chen, J. Pan, J.-H. Xu (2016) A green-by-design system for efficient bio-oxidation of an unnatural hexapyranose into chiral lactone for building statin side-chains. *Catal. Sci. Technol.*, **6**, 7094–100. https://doi.org/10.1039/C6CY01085GM

C. Julian, L. Li, S. Garde, R. Wilen, P. M. Tessier (2017) Efficient affinity maturation of antibody variable domains requires co-selection of compensatory mutations to maintain thermodynamic stability. *Sci. Rep.* **7**, 45259; https://doi.org/ 10.1038/srep45259

X. Lu, Y. Liu, Y. Yang, S. Wang, Q. Wang, X. Wang, Z. Yan, J. Cheng, C. Liu, X. Yang,, H. Luo, S. Yang, J. Gou, L. Ye, L. Lu, Z. Zhang, Y. Guo, Y. Nie, J. Lin, S. Li, C. Tian, T. Cai, B. Zhuo, H. Ma, W. Wang, Y. Ma, Y. Liu, Y. Li, H. Jiang (2019) Constructing a synthetic pathway for acetyl-coenzyme A from one-carbon through enzyme design. *Nat. Commun.* **10**, 1378; https://doi.org/10.1038/s41467-019-09095-z

M. Lunzer, G. B. Golding, A. M. Dean (2010) Pervasive cryptic epistasis in molecular evolution. *PLoS Genet.* **6**, e1001162; https://doi.org/10.1371/journal.pgen.1001162

M. Lunzer, S. P. Miller, R. Felsheim, A. M. Dean (2005) The biochemical architecture of an ancient adaptive landscape. *Science* **310,** 499–501.

J. Maynard Smith (1970) Natural selection and the concept of a protein space. *Nature* **225**, 563–4; https://doi.org/ 10.1038/225563a0

M. J. McDonald, D. P. Rice, M. M. Desai (2016) Sex speeds adaptation by altering the dynamics of molecular evolution. *Nature* **531**, 233–6; https://doi.org/10.1038/nature17143

Mouse Genome Sequencing Consortium (2002) Initial sequencing and comparative analysis of the mouse genome. *Nature*, **420**, 520–62.

E. A. Ortlund, J. T. Bridgham, M. R. Redinbo, J. W. Thornton (2007) Crystal structure of an ancient protein: evolution of a new function by conformational epistasis. *Science*, **317**, 1544–8. https://doi.org/101126/ science.1142819.

S. Oue, A. Okamoto, T. Yano, H. Kagamiyama (1999) Redesigning the substrate specificity of an enzyme by cumulative effects of the mutations of non-active site residues. *J. Biol. Chem.* **274,** 2344–9.

F. J. Poelwijk, D. J. Kiviet, D. M. Weinreich, S. J. Tans (2007) Empirical fitness landscapes reveal accessible evolutionary paths. *Nature* **445**, 383–6. https://doi.org/10.1038/nature05451

A. Rath, A. R. Davidson (2000) The design of a hyperstable mutant of the Abp1p SH3 domain by sequence alignment analysis. *Protein Sci.* **9**, 2457–69.

M. T. Reetz, M. Bocola, J. D. Carballeira, D. Zha, A. Vogel (2005) Expanding the range of substrate acceptance of enzymes: combinatorial active-site saturation test. *Angew. Chem. Int. Ed.* **44**, 4192–6; https://doi.org/10.1002/ anie.200500767

M. T. Reetz, D. Kahakeaw, R. Lohmer (2008) Addressing the numbers problem in directed evolution. *ChemBioChem* **9**, 1797–804; https://doi.org/10.1002/cbic.200800298

P. A. Romero, F. H. Arnold (2009) Exploring protein fitness landscapes by directed evolution. *Nat. Rev. Mol. Cell Biol.* **10**, 866–76.

S. C. Rothman, M. Voorhies, J. F. Kirsch (2004) Directed evolution relieves product inhibition and confers *in vivo* function to a rationally designed tyrosine aminotransferase. *Protein Sci.* **13,** 763–772.

I. Rowles, Kirk J. Malone, L. L. Etchells, S. C. Willies, N. J. Turner (2012) Directed evolution of the enzyme monoamine oxidase (MAO-N): highly efficient chemo-enzymatic deracemisation of the alkaloid (±)-crispine A, *ChemCatChem*, **4**, 1259–61.

A. G. Sandström, Y. Wikmark, K. Engström, J. Nyhlén, J.-E. Bäckvall (2012) Combinatorial reshaping of the *Candida antarctica* lipase A substrate pocket for enantioselectivity using an extremely condensed library. *Proc. Natl. Acad. Sci. U. S. A. 109*, 78–83. https://doi.org/10.1073/pnas.1111537108

C. K. Savile, J. M. Janey, E. C. Mundorff, J. C. Moore, S. Tam, W. R. Jarvis, J. C. Colbeck, A. Krebber, F. J. Fleitz, J. Brands, P. N. Devine, G. W. Huisman, G. J. Hughes (2010) Biocatalytic asymmetric synthesis of chiral amines from ketones applied to sitagliptin manufacture. *Science* **329**, 305–9.

J. L. Seffernick, L. P. Wackett (2001) Rapid evolution of bacterial catabolic enzymes: a case study with atrazine chlorohydrolase. *Biochemistry* **40**, 12747–53. https://doi.org/10.1021/bi011293r

J. B. Siegel, A. L. Smith, S. Poust, A. J. Wargacki, A. Bar-Even, C. Louw, B. W. Shen, C. B. Eiben, H. M. Tran, E. Noor, J. L. Gallaher, J. Bale, Y. Yoshikuni, M. H. Gelb, J. D. Keasling, B. L. Stoddard, M. E. Lidstrom, D. Baker (2015) Computational protein design enables a novel one-carbon assimilation pathway. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 3704–9; https://doi.org/10.1073/pnas.1500545112

B. Spiller, A. Gershenson, F. H. Arnold, R. C. Stevens (1999) A structural view of evolutionary divergence. *Proc. Natl Acad. Sci. U. S. A.* **96,** 12305–10.

Z. Sun, R. Lonsdale, L. Wu, G. Li, A. Li, J. Wang, J. Zhou, M. T. Reetz (2016) Structure-guided triple-code saturation mutagenesis: efficient tuning of the stereoselectivity of an epoxide hydrolase. *ACS Catal.* **6**, 1590–1597. https://doi.org/10.1021/acscatal.5b02751

N. Tokuriki, D. S. Tawfik (2009) Protein dynamism and evolvability. *Science* **324**, 203–7. https://doi.org/10.1126/science.1169375

J.-Y. van der Meer, H. Poddar, B.-J. Baas, Y. Miao, M. Rahimi, A. Kunzendorf, R. van Merkerk, P. G. Tepper, E. M. Geertsema, A.-M. W. H. Thunnissen, W. J. Quax, G. P. Poelarends (2016) Using mutability landscapes of a promiscuous tautomerase to guide the engineering of enantioselective Michaelases. *Nat. Commun.* **7**, 10911. https://doi.org/10.1038/ncomms10911

D. M. Weinreich, N, F. Delaney, M. A. DePristo, D. L. Hart (2006) Darwinian Evolution can follow only very few mutational paths to fitter proteins, *Science*, **312**, 111–4.

J. A. Wells (1990) Additivity of mutational effects in proteins. *Biochemistry* **29**, 8509–17.

Y. Yoshikuni, T. E. Ferrin, J. D. Keasling (2006) Designed divergent evolution of enzyme function. *Nature* **440**, 107882.