

Compensated pathogenic deviations

Anja Barešić and Andrew C.R. Martin*

Institute of Structural and Molecular Biology, Division of Biosciences, University College London, Darwin Building, Gower Street, London WC1E 6BT, UK

*Corresponding author
e-mail: andrew@bioinf.org.uk

Abstract

Deleterious or ‘disease-associated’ mutations are mutations that lead to disease with high phenotype penetrance: they are inherited in a simple Mendelian manner, or, in the case of cancer, accumulate in somatic cells leading directly to disease. However, in some cases, the amino acid that is substituted resulting in disease is the wild-type native residue in the functionally equivalent protein in another species. Such examples are known as ‘compensated pathogenic deviations’ (CPDs) because, somewhere in the second species, there must be compensatory mutations that allow the protein to function normally despite having a residue which would cause disease in the first species. Depending on the nature of the mutations, compensation can occur in the same protein, or in a different protein with which it interacts. In principle, compensation can be achieved by a single mutation (most probably structurally close to the CPD), or by the cumulative effect of several mutations. Although it is clear that these effects occur in proteins, compensatory mutations are also important in RNA potentially having an impact on disease. As a much simpler molecule, RNA provides an interesting model for understanding mechanisms of compensatory effects, both by looking at naturally occurring RNA molecules and as a means of computational simulation. This review surveys the rather limited literature that has explored these effects. Understanding the nature of CPDs is important in understanding traversal along fitness landscape valleys in evolution. It could also have applications in treating diseases that result from such mutations.

Keywords: co-adaptation; co-evolution; deleterious mutations; disease-associated mutations; epistasis; single nucleotide polymorphisms.

Introduction

It has frequently been observed that, when deleterious single amino acid mutations are surveyed, mutated amino acid types with detrimental effects in one species are found as the native wild-type residue in homologous proteins of other species, with neutral effect on the fitness of the latter species. The most likely scenario explaining such observations is that the two homologous proteins provide slightly different struc-

tural environments for the same residue, thus compensating for the deleterious effect of the residue in the first protein. Generally researchers have looked at cases of human disease-causing ‘deleterious’ or ‘disease-associated’ mutations (DAMs) and observed that the mutant (disease-causing) amino acid is the native (wild-type) amino acid in another species. Such cases are known as ‘compensated pathogenic deviations’ (CPDs).

Figure 1 shows an example of two DAMs in human antithrombin III (ANT3), one of which is compensated and the other uncompensated. In the human protein, the mutations Ala416 → Pro and Ala416 → Ser both cause susceptibility to thrombophilia as a result of antithrombin III deficiency. Details of these mutations can be seen at Online Mendelian Inheritance in Man (OMIM) Entries 107300.0007 and 107300.0027 (<http://www.ncbi.nlm.nih.gov/omim/107300>). Although OMIM states that the mutation occurs at residue 384, this equates to residue 416 in the UniProtKB/SwissProt sequence (UniProtKB/SwissProt accession P01008). Our online resource at <http://www.bioinf.org.uk/omim/> provides a validated mapping of residue numbers in OMIM to UniProtKB/SwissProt residue numbers. As the alignment shows, this residue is a conserved alanine in all the sequences examined – neither proline nor serine is seen in any other species and the two disease-causing mutations seen in humans are therefore classified as ‘pathogenic deviations’ (PDs, see below). However, a mutation of Ala419 → Val (as described in OMIM Entry 107300.0042, OMIM residue number 387), which also leads to antithrombin III deficiency, occurs at a residue which is not conserved in the alignment. In fact, sheep and cows have a valine at this position in the native sequence and thus the Ala419 → Val mutation in humans is classified as a CPD.

The question, therefore, is how do the sheep and bovine proteins function properly with a valine at position 419? Presumably, during the evolution of human, sheep and bovine ANT3 proteins from a common ancestor, some other amino acid difference(s) have occurred in the sheep and bovine proteins compared with the human protein that somehow compensate for what, in the human protein, is the negative effect of having a valine at position 419. How the compensation is achieved in this example is not clear.

Compensation of mutations is also important at the RNA level. Stable Watson-Crick base pairing in RNA can bring together remote parts of the molecule to form stable three-dimensional structures of functional importance. Thus, mutations in the RNA must undergo compensatory events to maintain the necessary base pairing requiring the crossing of valleys on the fitness landscape. Not only has this been studied using real RNA sequences (1), but RNA has also been used in computational models designed to understand compensatory mutation (2, 3).

```

                                                                    416 419
                                                                    |  |
P01008|ANT3_HUMAN  GFSLKEQLQDMGLVDLFSPEKSKLPGIVAEGRDDLYVSDAFHKAFLEVNEEGSEAAASTAVVI
Q5R5A3|ANT3_PONPY  GFSLKEQLQDMGLVDLFSPEKSKLPGIVAEGRDDLYVSDAFHKAFLEVNEEGSEAAASTAVVI
P32261|ANT3_MOUSE  GFSLKEQLQDMGLIDLFSPEKSQLPGIVAGGRDDLYVSDAFHKAFLEVNEEGSEAAASTSVVI
P41361|ANT3_BOVIN  SFSVKEQLQDMGLEDLFSPEKSRPLPGIVAEGRSDLYVSDAFHKAFLEVNEEGSEAAASTVVISI
P32262|ANT3_SHEEP  SFSVKEQLQDMGLEDLFSPEKSRPLPGIVAEGRNDLYVSDAFHKAFLEVNEEGSEAAASTVVISI
. **:***** * :***** :***** ** .***** :***** : *

```

Figure 1 Examples of two disease-associated mutations (DAMs) reproduced from our structural analysis (10). The Figure shows the alignment of the human antithrombin III (ANT3) protein sequence with non-human functionally equivalent homologous proteins. Highlighted are columns 416 and 419 which represent an uncompensated pathogenic deviation (PD) and a compensated pathogenic deviation (CPD), respectively.

Body of review

The term ‘compensated mutations’ was introduced by Kimura (4), who demonstrated that two mutually compensatory mutations could become fixed in a population as a result of random genetic drift. Kimura defined ‘compensatory neutral mutations’ as linked deleterious mutations; in other words, two mutations each of which, by itself, has a deleterious effect, but together have a neutral (or potentially even a beneficial) effect on overall fitness. The ability of one mutation to compensate for the pathogenic effects of another newly introduced mutation is an important mechanism in evolution. Using the same analogy used by Wright (5) and used extensively by Dawkins (6), the fitness landscape can be viewed as mountains of high fitness separated by valleys of low fitness. Thus, compensation of mutations allows bridging the valleys of low fitness.

Terminology

Because the analysis and understanding of CPDs crosses the boundaries of structural and evolutionary biology, it is useful to define several terms that are used in the field before we go into any more discussion.

‘Single nucleotide polymorphisms’ (SNPs) are single DNA base changes. Strictly the term is applied only to instances where the mutation is observed in at least 1% of a ‘normal’ population. In other words, they will either have a completely neutral phenotype or a low-penetrance phenotype where there is no clear Mendelian inheritance. Such SNPs can be involved in more complex conditions such as heart disease or simply give a propensity towards disease through interaction with external factors. However, it should be noted that many researchers use the term SNP to refer to *any* single base change, even when no frequency data are available. In our previous study looking at the effects of mutation on protein structure (7), we tried to use the term SNP in the correct way (with the assumption that they do not lead to high-penetrance Mendelian inherited disease) and contrasted these with mutations that do lead to disease. However, even dbSNP (8), the primary repository for SNP data at the National Center for Biotechnology Information (NCBI), includes data on lower frequency mutations.

SNPs can occur in coding or non-coding regions of DNA. Both coding (cSNPs) and non-coding SNPs (ncSNPs) can

have effects on gene expression or mRNA splicing; cSNPs can (i) be synonymous in terms of the resultant amino acid (sSNP), (ii) lead to a premature stop or ‘nonsense’ codon (nSNP), or (iii) be non-synonymous (an nsSNP) resulting in a single amino acid change (Figure 2).

‘Single amino acid polymorphisms’ (SAAPs) are single amino acid mutations resulting from nsSNPs. We use the term as defined by Hurst et al. (7) to apply both to mutations resulting from strictly defined nsSNPs (i.e., those that occur in at least 1% of a normal population) and to deleterious mutations (DAMs) as defined below (Figure 2).

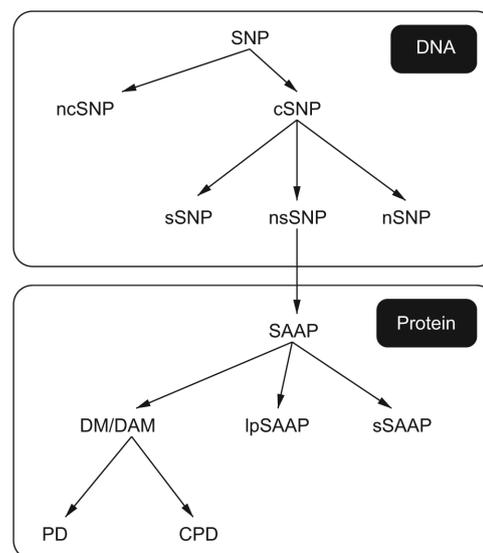


Figure 2 Hierarchy of SNPs, mutations, and their effects. SNPs (defined in the general sense to mean any single base DNA mutation) can be non-coding (ncSNPs) or coding (cSNPs). cSNPs can be synonymous (sSNPs), nonsense (nSNPs), or non-synonymous (nsSNPs). nsSNPs result in a single amino acid polymorphism (SAAP) at the protein level. These can be phenotypically silent (sSAAP), low-penetrance (lpSAAP), or high-penetrance deleterious mutations (DMs) also known as disease-associated mutations (DAMs). A DAM can be compensated in another species (a compensated pathogenic deviation, CPD) or uncompensated (a pathogenic deviation, PD). Note that all forms of SNPs can have effects on expression as they can affect regulatory regions or splice sites. Note also that lpSAAPs form a continuum between phenotypically silent and high-penetrance disease-associated mutations.

‘Deleterious mutations’ also referred to as DAMs (9) are SAAPs that result in high-penetrance disease phenotypes. In this review, we use the term to encompass both PDs and CPDs as defined below (Figure 2).

‘Pathogenic deviation’ (PD) is often used as a synonym for DAM, but in the discussion of CPDs (see below), we generally refer to PDs as disease-causing mutations that are *not* observed to be compensated in any other species and that is the definition we use throughout this review. As discussed by Barešić et al. (10), this definition of PDs is not completely reliable because it is based on a negative observation. Mutations are classified as PDs rather than CPDs simply because the residue is not observed as the native residue in any other species, but until we have the sequence of every species, we cannot conclusively know that it is not compensated in at least one other species. See Figure 2 and column 416 in Figure 1 for an example of a PD.

‘Compensated pathogenic deviations’ (CPDs) have also been referred to as ‘potential compensated mutations’ (9). Their existence was first discussed by Kimura (4), who termed them ‘compensatory neutral mutations’, whereas the term CPD was first defined by Kondrashov et al. (11). A CPD is a SAAP (as defined above) associated with a disease phenotype (i.e., a DAM), usually in a human protein, where the mutated amino acid type is found as the native (phenotypically neutral) residue at the same position in an ortholog of another species. See Figure 2 and column 419 in Figure 1 for an example of a CPD.

‘Functionally equivalent proteins’ (FEPs) are orthologs which have maintained the same function during evolution, as discussed by McMillan and Martin (12). Homologous genes (or proteins) have descended from a common ancestor, while orthologs are the subset of homologs that arise from speciation events (13). However, if two species have diverged sufficiently, the function of one of the pair of orthologous proteins could diverge. For example, Shibata et al. (14) showed that although the general function of exportin-5 proteins (nuclear export of miRNAs and tRNAs) is conserved across different species, substrate specificity varies.

‘Co-evolution’: at the molecular level, evolution of each protein molecule is affected by (potentially numerous) interaction partners and environmental factors. When a similar evolutionary pattern is detected for the two molecules, they are said to be co-evolving. This shared evolutionary history can be a consequence of their co-adaptation, shared cellular pathway or localization, or a shared expression pattern (15). In examining CPDs, we are only interested in the first of these – the co-adaptation of two amino acids which affect each other’s evolutionary paths.

‘Epistasis’ is defined as the effects of one gene being modified by one (or several) other genes (sometimes termed ‘modifier genes’). Typically the phenotype of one gene (the ‘epistatic’ gene) is expressed, whereas the other (the ‘hypostatic’ gene) is altered or suppressed. This interaction of *different* genetic loci contrasts with normal Mendelian effects, where one allele is ‘dominant’ over another ‘recessive’ allele at the *same* locus. In a more general way, epistasis is defined as an interdependence between two gene loci as discussed

by Cordell (16). In the context of population genetics, ‘epistasis’ refers to the interaction between alleles at different loci in such a way that the effect on the individual cannot be predicted from merely adding up effects of interacting loci. In the case of CPDs we are interested in the change of fitness of a protein caused by a change of a single amino acid. Fitness can be modified by differences (i.e., amino acid changes) at other locations. Although the term ‘epistasis’ should strictly be applied only to changes in other proteins, when discussing CPDs it is further generalized to refer to changes at other locations within the *same* protein.

‘Sign epistasis’ also known as fitness reversal, refers to the situation in which there is a deleterious mutation which co-evolves with a mutation having an epistatic effect that more than compensates for the deleterious effects of the other mutation. Thus, the overall fitness change becomes positive (or at least neutral) rather than negative. Sign epistasis facilitates sampling protein space for novel amino acid combinations and provides a mechanism of escape from local fitness minima (2). In some cases ‘fitness reversal’ can be used as a more general term (perhaps influenced by epigenetic effects), whereas sign epistasis specifically refers to the effect of compensatory mutations. In this review, we use the terms interchangeably.

RNA as a model of compensation

Although it is clear that protein and RNA are very different molecules, the simple nature of RNA models has, in general, been widely applied to study evolution. Understanding the importance of compensatory events during evolution is no exception. RNA consists of just four nucleotides: adenine (A), guanine (G), cytosine (C) and uracil (U). Just as in DNA, stable Watson-Crick pairing can occur between A and U, and between G and C. This can bring relatively remote parts of the molecule together to form stable three-dimensional structures composed of features such as ‘stems’ (helical base-paired regions) and unpaired regions which form ‘loops’ (at the end of a stem) or ‘bulges’ (in the middle of a stem). One can view the RNA sequence as being a ‘genotype’, whereas the manifestation of a stable folded structure is the ‘phenotype’. The simple nature of RNA folding means that it can be simulated in a computer with a high degree of accuracy using freely available software [see, e.g., Zuker’s MFold software (17, 18) and Schuster’s ViennaRNA (19)]. More recent software can even predict the shapes of RNA molecules during interactions with other molecules [e.g., the work of the Hofacker group (20) and of the Mathews group (21, 22)].

Computational models of RNA evolution typically simulate a large population of RNA molecules and apply the standard strategy of random mutation followed by natural selection. On the basis that most functional RNA molecules have shapes that are extremely conserved throughout evolution, because shape has a dominant role in determining function (23), the fitness of an RNA molecule is determined by predicting its shape and then applying a fitness function based on similarity to some predetermined ideal target shape. Having evaluated the fitness, molecules are allowed to rep-

licate in proportion to their fitness and, during the replication, random mutations are allowed to occur.

The application of RNA models to understanding evolution is reviewed by Cowperthwaite and Meyers (3) and, in an earlier paper, Cowperthwaite et al. (2) used these models to examine fitness reversal. They observed that RNA mutations that can be regarded as ‘pathogenic’ in the model system accumulate more rapidly than expected based on their effect on overall population fitness. Furthermore, they observed that the drop in fitness was not as severe as would be expected based on the accumulation of deleterious variations. Because deleterious effects were not additive, compensatory events were clearly occurring. Indeed, mutations that initially were deleterious accumulated at nearly the same rate as mutations that were immediately beneficial and fixations of more than half of the initially deleterious mutations led to fitness reversals. The fixation of initially deleterious mutations led to a substantial positive effect on the total fitness of the genome. When other mutational events such as ‘hitchhiking’ and random drift were considered, their model showed that some 80% of PDs were fixed through fitness reversal or co-adaptation with a compensatory mutation.

In a related study, but using real sequences rather than computer simulations, Meer et al. (1) attempted to address the question of whether valleys on the fitness landscape (corresponding to low-fitness genotypes) can be crossed to reach isolated fitness peaks. In particular, they examined the switch between AU and GC Watson-Crick nucleotide pairs at equivalent sites in the mitochondrial tRNA stem regions in 83 mammalian species. Clearly, to switch from an AU pair to a GC pair either needs A → G and U → C mutations to occur simultaneously (thus jumping from one fitness peak to another – an unlikely event), or requires one mutation to occur before the other thus passing through a valley of low fitness where there will be a Watson-Crick mismatch. Because of the need to traverse low-fitness valleys, they found that these ‘Watson-Crick switches’ occurred 30–40 times more slowly than did pairs of neutral substitutions (where base pairing was not a factor). However, they found that substitutions leading to a Watson-Crick switch were strongly correlated. They were able to estimate the depths of the fitness valleys and showed that AC intermediates are slightly more deleterious than GU intermediates. Nevertheless, the compensatory evolutionary events that do occur must proceed via rare disfavoured intermediate variants that never become fixed in the population.

Analysis of compensatory events in proteins

As discussed above, computer simulations in RNA and studies of RNA molecules have shown that compensatory events do indeed allow traversal of valleys in the fitness landscape. RNA, having only four nucleotides is clearly a much simpler system than proteins composed of 20 amino acids, but we know that compensatory events must also occur in proteins. It is difficult to say whether the fact that there are 20 amino acids with a wide variety of chemical and physical properties makes it harder or easier to compensate in proteins than in RNA. On the one hand, the subtlety and complexity of inter-

actions made by amino acids could mean that compensatory events are difficult; on the other hand, a change that is damaging might be somewhat small in nature and therefore only need a small compensatory event, perhaps by a conservative substitution in a nearby amino acid. The compensating event might (if it happens first) not have a particularly negative effect.

Over the past decade, several groups have started to look at CPDs in proteins, but although the definition of a CPD is the same, different approaches have been taken to gathering CPD data.

CPDs are identified by (i) identifying missense mutations that lead to disease (generally in humans), (ii) identifying a set of homologous proteins, (iii) performing a multiple alignment of the human sequences with the homologous sequences, and (iv) identifying cases where the pathogenic mutation is observed as the native residue in at least one other species. Thus, not surprisingly, datasets of CPDs are highly dependent on (i) the alignment building method, (ii) the thresholds used to detect homologous proteins, and (iii) the choice of species to be tested for homologs. Several methods are summarized in Table 1 showing a variety of species, cut-off values for identifying homologies to be included in the dataset, and multiple sequence alignment methods. In particular, Poon et al. approach (24) was rather different from the others. They analyzed deleterious missense mutations from a range of proteins in different species. Rather than use a sequence-comparison approach as used in the other datasets, they analyzed data from publications identified using relevant keyword searches. Thus, their data show a very high fraction of deleterious mutations that are compensated because their analysis focused only on these mutations. In addition, they considered mutations introduced with mutagenesis-inducing agents as well as evolutionary events.

Once the data have been collected, some authors performed various analyses to compare and contrast compensated mutations with the rest of the dataset to try to understand whether the nature of mutations that are seen to be compensated (CPDs) is different from those that are not seen to be compensated (PDs). As described in the definitions above, although we use the term PDs strictly to refer to uncompensated mutations, the identification of a clear uncompensated set is not completely rigorous as it is based on a negative observation. Thus, the fact that no compensatory event has been observed could simply be because a species that has a compensatory event has not yet been sequenced. Similarly, sequence quality is always a concern (25) and it is possible that apparent CPDs are actually a result of sequencing errors.

Excluding the Poon dataset which is deliberately biased towards compensated mutations, Table 1 shows that the fraction of disease-causing mutations that are seen to be compensated varies from 0.14% in the Zhang dataset (62/44348) to 19.5% in the Barešić dataset (453/2328): our contribution to this field. This clearly shows that the number of compensatory events is correlated with the evolutionary distance between the species considered. In the Zhang dataset, only humans, chimpanzees, and neanderthals were examined,

Table 1 Datasets of compensated pathogenic deviations described in the literature.

Dataset	Species	Identity cut-off value	Alignment method	Human proteins	# DAMs	# CPDs
Kondrashov et al. (11)	Any mammals ^a	>50%	CLUSTALW	32 3	4880 ^b	608 20
Kulathinal et al. (26)	Diptera			475 ^c	1527	6
Ferrer-Costa et al. (33) ^d	Any mammals	>10% (>60%)	Pfam	287 (24) 184	9334	1658 (52) 847
Barešić et al. (10)	Any	None ^e	MUSCLE	245	2328	453
Zhang et al. (missense) (9)	Human, neanderthal, chimpanzee		ANFO	2628	44348	62
Poon et al. (24) Set A ^f	Any			43 ^g	115	88
Poon et al. (24) Set B ^f	Any			17 ^g	59	49

^aKondrashov et al. tested all found orthologs (with no sequence identity threshold) for CPDs and then switched to mammalian-only orthologs to identify compensatory mutations.

^bPrecise numbers are somewhat unclear. They report 608 CPDs and that this is approximately 10% of DAMs. In Table 1 of their paper, there are 4272 ‘known missense’ mutations which we believe to be PDs because the last row of the table has more CPDs than ‘known missense’ mutations. This makes a total of 4880 (4272+608) DAMs.

^cIn the Kulathinal group dataset, the reference species is *Drosophila melanogaster* instead of human.

^dNumbers in parentheses refer to the CPDs with structural data available, used for structural analysis.

^eFunctional equivalence among homologs used instead of a sequence identity threshold.

^fThe Poon group Set A includes mutations brought about by mutagenic agents, whereas Set B does not.

^gThere is no reference species in the Poon et al. study.

whereas the high fraction in the Barešić dataset results from the fact that no limit was applied to the divergence of the homologous sequences. As sequences diverge more, the environment around any given residue is likely to be more different and therefore a residue change is more likely to be tolerated, or indeed, required. Kondrashov et al. found that, when using a dataset containing only homologs with at least 50% identity to the reference sequence, on average around 1 in 10 disease-causing mutations is seen to be compensated in other species (11). By contrast, alignments of recently diverged sequences [e.g., three Dipteran genomes (26) or chimpanzee, neanderthal, and modern human (9)] show far fewer CPDs (0.4% in the Kulathinal dataset and 0.14% in the Zhang dataset).

The motivation for not using any sequence identity threshold in our study (10) was that we wished to compare the local structural effects of mutations that could be compensated with those that could not. Therefore, having a set of CPDs that was as broad as possible meant that our uncompensated PD dataset was likely to be more accurate. The dataset was built using only FEPs [as defined by McMillan and Martin (12)]. Thus, whereas other groups identify homologs using a BLAST (27) search with default parameters (11), or manually curated alignments from Pfam (28), we selected all orthologs where function has been conserved as defined by annotations in UniprotKB/SwissProt. These data are available in our FOSTA database (12).

Properties of compensated mutations and mechanisms of compensation

Maintaining a functional protein requires a delicate balance between the residues present to obtain proteins having a narrow range of thermodynamic stability, a range of ΔG from

-3 to -10 kcal/mol. If the stability is any lower, then the protein will start to unfold, becoming a target for degradation; higher stability means that the protein cannot be turned over effectively and therefore often becomes unresponsive to cell regulation or can lose its activity (29). In addition, mutated proteins of both lower and higher stability than optimal often show increased propensity for aggregation, although aggregation potential is not solely dependent on protein stability. Amino acid substitutions result in an average $\Delta\Delta G$ of 0.5–5 kcal/mol (29), so it is clear that most SAAPs will have a significant effect on protein stability and consequently protein function and the individual’s fitness.

From a structural perspective, compensated mutations have been shown to have less damaging effects than uncompensated mutations. Henikoff and Henikoff (30) created the BLOSUM amino acid substitution matrices from around 2000 blocks of aligned sequence segments from more than 500 groups of related proteins to show how frequently one amino acid can substitute for another in homologous proteins. These matrices were designed for use in protein sequence alignment and are familiar to most biologists as the default similarity matrix for use with the BLAST sequence searching tool (27). Ferrer-Costa et al. (31) showed that CPDs show significantly larger BLOSUM62 scores than PDs – in other words the amino acid replacements observed in CPDs are more frequently observed to occur in general in homologous proteins, whereas the replacements seen in PDs are less commonly observed between homologous proteins. They also found that CPDs are characterized by less extreme changes in amino acid volume and hydrophobicity when compared with uncompensated PDs.

In a previous study (10), we examined 14 different local structural effects covering stability and folding of the protein,

as well as binding effects and functional annotations. We found that CPDs are less likely to display any of these effects, especially if the structural effect is likely to require several consecutive compensatory mutations for full fitness reversal rather than it being possible to compensate using a single substitution. For example, a buried mutation, where a small residue is replaced by a larger residue, could cause a clash. However, although it is theoretically possible that a single mutation could do so, compensation of a clash is most likely to be achieved by making several smaller changes. Both Ferrer-Costa et al. (31) and Barešić et al. (10) found that CPDs have a higher average solvent accessibility. In other words, they are much more likely to be found on, or near, the protein surface.

Compensatory mutations in evolution

In the context of evolution, compensated mutations become fixed in a population through ‘co-adaptation’ or, more precisely, through ‘sign epistasis’ as defined above. At the protein level, depending on the context and role of the deleterious mutation (*D*), the compensatory mutation (*C*) can be on the same protein, or on an interacting partner protein. The compensatory mutation, *C*, could have no effect on fitness or could itself be somewhat deleterious, but at such a level that it can exist in the population. However, the main feature of *C* is that, when it co-occurs together with the deleterious mutation, *D*, it reverses the negative fitness effect of *D* to a neutral or positive one and, if *C* by itself has any deleterious effect, the combination of *C* and *D* will have a neutral or positive effect. Thus during evolution, when fitness landscapes are explored, compensation provides a path through the valleys of lower fitness, allowing individuals to travel from one peak to another (5).

As discussed above, numerous cases of compensation have been identified and documented in proteins (9–11). Although a classic compensatory event to achieve fitness reversal would result from *C* being a single amino acid change, in proteins it is also perfectly possible – and indeed more likely – for *C* to consist of a complete change in environment from multiple amino acid changes.

Poon et al. (24) set out to study how many different compensatory mutations act on a given deleterious mutation. They performed a maximum-likelihood analysis of experimental data collected from the literature on suppressor mutations (which are equivalent to compensatory mutations) to determine the shape of the statistical distribution for the number of compensatory mutations per deleterious mutation. They found that the data were best explained by an L-shaped gamma distribution which predicted an average of 11.8 compensatory mutations per deleterious mutation to achieve full sign epistasis and compensate for the deleterious effect of a DAM (24). Interestingly, they also found that, when they partitioned the data into viruses, prokaryotes and eukaryotes, there was a significant improvement in the fit to the model: on average, there were fewer compensatory mutations in viruses than in prokaryotes or eukaryotes. They suggested that the differences in genome size and gene length in viruses compared with prokaryotes and eukaryotes means that the

number of possible interactions within and between gene products is constrained.

In our more recent structural study (10), we showed that CPDs are surrounded by significantly more diverged residues than PDs. As described above, we created sequence alignments of functionally equivalent homologous proteins for each instance in which a human deleterious mutation (DAM) is known [typically identified from OMIM (32, 33), but also from a number of locus-specific mutation databases (7)]. The DAMs were then assigned as CPDs or PDs depending on whether the damaging mutant residue was observed as the native in another species. Where a structure was available for the human protein, we identified amino acids within an 8 Å sphere around the DAM. Having identified these structural neighbors in the human protein which form the environment surrounding the DAM, we mapped their positions back onto the sequence alignment. We were then able to calculate the fraction of these structurally neighboring residues that were mutated in each of the sequences when compared with the human sequence. For CPDs this was done just with the sequences in which compensation was observed, whereas for PDs it was done for each sequence in the alignment. We then plotted this local fraction of mutated residues against the overall (whole protein) sequence identity between the human and non-human sequence.

We found that this environmental ‘sphere’ compensation appeared on average to occur as a result of random drift in the sequence. We fitted a straight line to the data imposing the biologically obvious constraint that the line had to pass through the 100% identity, zero mutations point – if the sequences are 100% identical then there can be no mutations within the local environment. Allowing for the fact that sequence identity ranges from 0% to 100%, whereas our fraction of mutations scale runs from 0 to 1 (and that one scale is scoring conservation, whereas the other is scoring mutations), this fitting revealed a slope of 1.007 for CPDs and 0.9 for PDs. The slope of approximately 1 for CPDs implies that the environment around a CPD is mutated at the same rate as the sequence overall such that compensation occurs as a result of random drift in the sequence. By contrast, the environment around PDs is more conserved than the sequence as a whole.

Although this ‘sphere compensation’ is probably the most common compensatory mechanism in proteins, the alternative classical ‘one-on-one compensation’ can also occur where one deleterious SAAP is compensated by another single mutation in the structural vicinity. This type of compensation is easier to detect, especially in analyses where only recently diverged homologs are considered (9, 11, 26). Two examples of one-on-one compensations are shown in the case study presented below.

The Poon et al. study (24) also investigated whether compensatory mutations are intragenic (i.e., occur within the same gene and hence the same protein chain as the deleterious mutation) or intergenic (i.e., occur within a different gene and protein chain from the DAM). Overall, from their dataset of 129 CPDs, they found that the majority (78%) of compensatory mutations were intragenic suggesting that the

complexity of interactions between proteins is likely to be less important than the complexity of the protein itself. However, when they studied different taxa separately, they found that compensation is much less likely to be intragenic in viruses (69%) than in prokaryotes (92%) and eukaryotes (90%). They proposed that this is probably a result of the fact that viral genes tend to be shorter, thus limiting the number of internal interactions.

Research performed by Povolotskaya and Kondrashov (34) suggests that compensated pathogenic deviations are unidirectional drivers of evolution; once compensation occurs, it is unlikely that sequences will revert to the original wild-type state. Their investigation of divergence of proteins in sequence space showed that, at any given point in time, only 2% of all possible missense mutations are allowed in order to avoid non-functional protein products. If we assume that only one missense mutation at a time can be introduced into the sequence, then we can consider how this observation affects a protein chain consisting of 100 residues. For every residue there are 19 possible substitutions, so at any one time 1900 (100×19) mutational events could occur. Given that 2% of these are 'allowed', 38 missense mutations will result in a functional protein. Let us assume that an allowed mutation of residue X to residue Y occurs at position n (i.e., $X_n \rightarrow Y$). At the next step, there will again be 38 allowed missense mutations, one of which would be the reversal of the mutation that occurred in the previous step (i.e., $Y_n \rightarrow X$). Thus, there is a 1 in 38 (2.6%) chance that this will occur, but a 97.4% chance that another mutation will occur. From this statistic, we do not know how the 38 allowed mutations will be distributed across the 100 amino acid positions of the protein. Thus, the second mutation could result in $Y_n \rightarrow Z$, but in general it is much more likely that the mutation will occur at a location m that is different from n . Thus, we are much more likely to obtain a double mutant after the second step than we are to obtain a reversion to the original sequence or to introduce a different amino acid at position n . Consequently, subsequent mutational events will lead to a drift away from the original sequence and it is intuitive that compensation will be observed significantly more often than reversal to the original sequence.

The question remains as to the timeline of compensatory events. As discussed in our previous study (10), DePristo et al. (29) proposed two hypotheses of CPD evolution based on models of biophysical properties. In the first scenario, a compensatory mutation C is phenotypically neutral and stable, thus fixing itself quickly in the population. The deleterious mutation D is unstable and can only become fixed in the population if it occurs *after* the compensatory mutation C . Thus, D will exist as a CPD because of the compensatory effect of C . In the second model, both D and C are individually deleterious, but either can exist in the population at low levels; it is known that small frequencies of low-fitness mutations exist in large populations. Consequently, if D is present in the population at low levels, then C can occur later and fix the D - C genotype in the population because of the epistatic effect of the mutant pair. Cowperthwaite et al. (2), in their *in silico* RNA models discussed earlier, con-

firmed that the deleterious mutation, D , can occur first. Equally, the compensatory mutation C can be present in the population at low levels and D can occur later leading to fixation of the D - C genotype in the same way. A less likely, but possible, scenario is that both C and D occur simultaneously. Provided the mutation rate is sufficiently high, epistatic selection with compensatory mutations is the most prevalent mechanism for fixation of otherwise deleterious mutations.

Artificial compensatory events

With recent advances in sequencing technology (35), sequencing large amounts of genomic data is becoming cheaper, faster, and more accessible, providing new opportunities in biomedical research. Genome-wide association studies (GWAS) are becoming more and more widespread, associating mutations with both high- and low-penetrance disease phenotypes. An important area of interest is the ability to predict whether a given mutation – particularly a SAAP – will lead to disease. Numerous tools have been developed both to analyze the local structural effects of mutations and to predict whether mutations will be damaging, many of these working mostly at the sequence level. Among these are SAAPdb (7), SNPs3D (36), stSNP (37), ModSNP (38), MutDB (39), LS-SNP (40), TopoSNP (41), SIFT (42), SNPeff (43), PolyPhen (44, 45), subPSEC (46), and nsSNPAnalyzer (47).

Recently, Critical Assessment of Genome Interpretation (CAGI) (48), a community experiment to assess computational methods for predicting the phenotypic impacts of genomic variation objectively, organized by Steve Brenner, John Moulton and Susanna Repo, was run for the first time. Participants were provided with genetic variant data and asked to make predictions of the resulting molecular, cellular, or organismal phenotype. Results from over 100 prediction submissions from eight countries were evaluated against experimental data by independent assessors and discussed at a workshop in December 2010 (see <http://genomeinterpretation.org/>).

One of the prediction datasets was particularly interesting in the context of compensated mutations. A dataset of p53 mutations (see <http://genomeinterpretation.org/content/p53/>) was designed to test prediction of 'cancer rescue mutations'. p53 is a tumor suppressor protein which plays a central role in detecting DNA damage, slowing the cell cycle to allow DNA repair enzymes to do their work (49), or if DNA damage is too severe, triggering programmed cell death (apoptosis) (50, 51). If p53 is rendered non-functional as a result of mutation, this central checkpoint is lost, allowing other mutations to accumulate in the DNA eventually leading to cancer. Unusually for tumor suppressor genes, in which most mutations tend to be frameshifts or nonsense codons, the majority of mutations in p53 are single DNA base changes resulting in a SAAP. In some cases, mutations at second intragenic sites are known to rescue the function reactivating otherwise inactive p53 (52–54) and are therefore acting as compensatory mutations. The Lathrop group at the University of California, Irvine, has been performing a complete

functional census of these cancer rescue mutations (55). In this case, the aim of the CAGI prediction experiment was to predict whether a given mutation is able to rescue the function of p53 and thus act as a compensatory mutation (56, 57). Although the results of the CAGI prediction experiment have not been published at the time of writing this review, we suspect the field of compensation prediction will progress significantly in the near future. If a disease-associated deleterious mutation is amenable to compensation (i.e., it is seen to be a CPD), it is likely that other (non-mutational) mechanisms of compensation could also be applied. For example, the Fersht group in Cambridge has shown that some p53 mutations can be compensated by binding small peptides that stabilize the p53 core domain (58, 59). More recently, small molecules which are more likely to be usable drug leads have been used successfully in the same way (60–63).

Case study: two compensated mutations and their environment

There are many examples of compensation which include the p53 rescue mutations described above where, in some cases, crystal structures have been solved to study the mechanism of compensation (53). Here, we will discuss two examples of compensated mutations. First, a CPD in human GTP cyclohydrolase (GTPCH) is presented, with an obvious destabilizing effect on the protein structure, while a compensating mutation has a stabilizing effect restoring enzyme activity. The second example, in ornithine transcarbamylase (OTC), is less obvious at the structural level, despite being confirmed by *in vitro* enzyme activity experiments.

GTPCH, encoded by the gene GCH1, plays a role in the folate and biopterin biosynthesis pathways and hydrolyses guanosine triphosphate (GTP) to form 7,8-dihydroneopterin-3'-triphosphate. This is the first step in the biosynthesis of tetrahydrobiopterin, an essential cofactor required by aromatic amino acid hydroxylase (AAAH) and nitric oxide synthase (NOS). These, in turn, are involved in the biosynthesis of monoamine neurotransmitters such as serotonin, melatonin, dopamine, noradrenaline, adrenaline, and nitric oxide. Mutations are associated with phenylketonuria (PKU) and

hyperphenylalaninemia (HPA), as well as levodopa-responsive dystonia.

Figure 3A shows the whole wild-type GTP cyclohydrolase I which consists of five identical chains, with mutually parallel C helices stabilizing the pentameric structure (64). The images, rendered with PyMol (<http://www.pymol.org/>), are based on coordinates obtained from Protein Databank entry 1FB1 accessible online at <http://www.pdb.org/> (65). Figure 3B shows details of the wild-type residues that are mutated (residues 249 and 250). The wild-type Arg249 in one chain and Ser250 in the next chain form a tight ring-like structure.

An Arg249→Ser mutation is associated with disease, causing a severe decrease in enzyme activity and resulting in recessive levodopa-responsive dystonia (OMIM: 600225.0016). Figure 3C shows the effect of introducing an Arg249→Ser mutation in all five chains modeled using the minimum perturbation protocol (66) implemented in the program Mutmodel (67). The non-covalent interactions between residues 249 and 250 are reduced, presumably destabilizing the complex and leading to disease. However, the functionally equivalent protein in *Rickettsia bellii* has a serine at 249 in the wild-type enzyme, but also has a compensatory lysine at 250, which is also modeled into the structure in Figure 3D restoring and, indeed, enhancing the ring-like set of interactions.

A less clear example of a compensatory mutation is seen in ornithine transcarbamylase (OTC) which catalyzes the reaction between carbamoyl phosphate and ornithine to form citrulline and phosphate. In prokaryotes and plants, it is involved in arginine biosynthesis, whereas in mammals it is a key enzyme of the urea cycle. Figure 4A shows one monomer of the enzyme which exists as a trimer. The structure for OTC in the Protein Databank shows only a monomer (PDB ID: 10TH), but the assembly of the biologically relevant trimer can be obtained from PISA (68) available online at http://www.ebi.ac.uk/msd-srv/prot_int/pistart.html. OTC deficiency, although a rare condition occurring in around 1 in 80 000 births, is the most common disorder of the urea cycle which removes ammonia from the body. Mutations in OTC lead to an accumulation of toxic ammonia which can lead to developmental delay and mental retardation, progres-

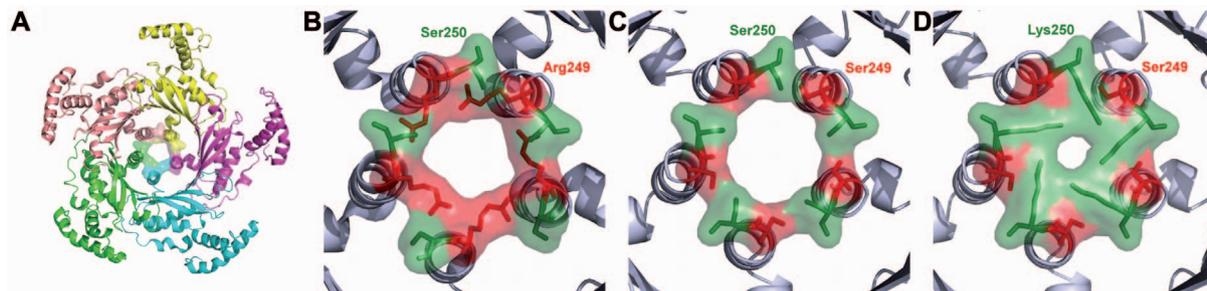


Figure 3 Compensated mutation in human GTP cyclohydrolase I. Residues 249 and 250 are shown with a surface in all five chains. (A) Structure of the wild-type homopentamer with each chain shown in a different color. (B) Zoomed view of residues 249 and 250 from all five chains with the residues shown in green and red, respectively. (C) The disease-causing Arg249→Ser mutation modeled into all five chains. (D) The compensatory Ser250→Lys mutation modeled into all five chains as well as the Arg249→Ser mutation.

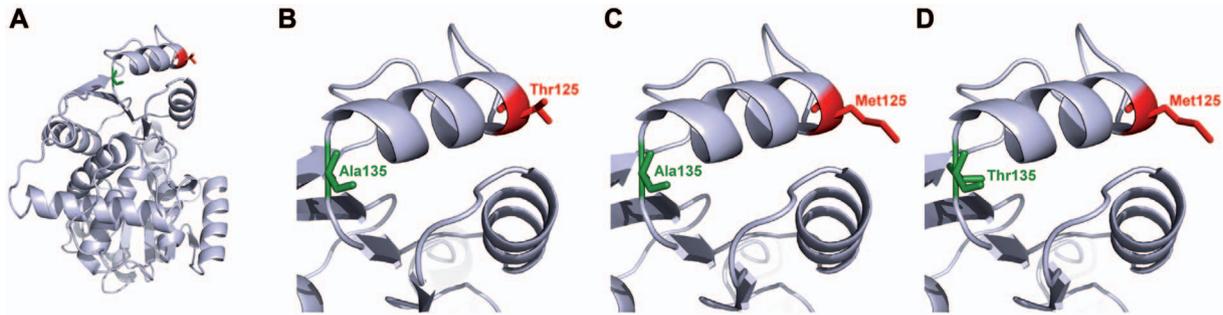


Figure 4 Compensated mutation in human ornithine transcarbamylase. (A) Structure of wild-type human OTC. (B) View on helix 3, with residues 125 and 135 shown in red and green, respectively. (C) The disease-causing Thr125→Met mutation, modeled structure. (D) A compensatory Ala135→Thr in addition to the Thr125→Met mutation.

sive liver damage, skin lesions, poorly controlled breathing, seizures, coma, and death.

Figure 4B shows helix 3 from PDB entry 10TH and highlights residues 125 and 135 in red and green, respectively (69). Thr125→Met is a known disease causing mutation in humans resulting in lethal neonatal congenital hyperammonemia (OMIM: 311250). Suriano et al. (70) showed that the human enzyme with the Thr125→Met mutation has a negligible rate of enzyme activity in *in vitro* constructs. However, this mutation is a CPD as Met is the native residue in chimpanzees. The only other residue which differs between human and chimp OTC is residue 135 where there is an Ala→Thr mutation which must compensate for the deleterious effect of the Thr125→Met mutation. However, the mechanism of compensation is unclear as Figure 4 shows that residues 125 and 135 are not in direct contact and this is also the case in the trimer. However, as previously suggested by Azevedo et al. (71), the presence of Thr125 might be crucial at the end of helix 3 because this helix is involved in trimerization of human OTC, and the chimpanzee compensates for its loss by having a threonine introduced at the other end of helix 3 (at position 135), restoring enzyme activity to rates similar to human wild-type. Interestingly, Suriano et al. (70) also suggested that the ancestral genotype could have had threonines at both positions 125 and 135 and had a higher enzyme activity than either the human or chimpanzee enzymes. If this is the case, then this mutation is an example of two species starting to explore fitness ridges, in search of another local optimum.

Expert opinion

In conclusion, although there is also the possibility that epigenetic effects can also be compensatory (i.e., some difference in the non-protein environment), compensation of deleterious mutations through epistatic protein mutations is a very common effect. The frequency of these compensatory mutations depends on the time elapsed from the common ancestor and the data in Table 1 show that there is a correlation between the frequency of CPDs and the diversity of the homologs used to detect CPDs. For example, our dataset (10) (where we apply no constraint on the sequence identity

between functionally equivalent homologs) shows a higher ratio of CPDs compared with the dipteran-only (26) or mammalian-only (9, 11, 31) datasets.

Study of the evolution of RNA molecules and *in silico* models of RNA evolution show clear examples of one-on-one compensation (2). Although compensation in proteins is often more complex, involving multiple compensatory events changing the environment in which a residue exists, there are also several examples of one-on-one compensation including ‘cancer rescue’ mutations in p53. As shown by DePristo et al. (29), any mutation has an average effect on protein stability ($\Delta\Delta G$) of around 0.5–5 kcal/mol. Restoring protein stability and hence regulated activity will often need compensatory mutations that restore stability to the acceptable range of free energies. From a structural analysis perspective, compensated mutations are preferentially on the protein surface (10, 31). As shown by Poon et al. (24), compensatory events are most commonly intragenic, so the surface location is likely to be a result of it being easier to accumulate compensatory events (probably before the CPD mutation occurs) rather than it being anything to do with interactions with other proteins. In addition, CPDs have ‘milder’ effects on the protein structure than uncompensated mutations (10, 31) and tend to be more conservative in nature (31).

Outlook

CPDs will continue to be an interesting area of research in understanding evolution and traversal of the fitness landscape. As more species are sequenced, the identification of true PDs will become more accurate. This will allow us to compare CPDs and PDs in a more accurate manner and therefore understand more completely which mutations can be easily compensated and which cannot. The CAGI experiment described above has led the way with the challenge of predicting which mutations will be ‘cancer rescue’ mutations in p53 and this will become a more significant area of research. The fact that certain mutations can be rescued or compensated by an amino acid change will also allow us to identify types of mutations that, in general, can be more easily rescued leading us towards the possibility of drugs that

can rescue protein function. Consequently, studying CPDs is not only of interest in understanding evolution, but is also important in developing future drugs.

Highlights

- Compensation of deleterious mutations through epistatic protein mutations is a very common effect.
- The frequency of compensated mutations depends on the time elapsed from the common ancestor – more diverged sequences are more likely to show compensatory events.
- Study of RNA molecules and *in silico* models of RNA evolution clearly show one-on-one compensation.
- Compensation in proteins is more likely to involve multiple compensatory events, but there are also several examples of one-on-one compensation.
- ‘Cancer rescue mutations’ in p53 are an example of one-on-one compensation.
- CPDs are more likely to occur on the protein surface, be more conservative in nature, and be less damaging in structural terms than PDs.
- Prediction of compensable mutations could allow design of drugs that are able to compensate for the effects of a damaging mutation.

Acknowledgments

A.B. was supported by a UK Overseas Research Scholarship and by a UCL Graduate School Research Scholarship.

References

1. Meer MV, Kondrashov AS, Artzy-Randrup Y, Kondrashov FA. Compensatory evolution in mitochondrial tRNAs navigates valleys of low fitness. *Nature* (London) 2010; 464: 279–82.
2. Cowperthwaite MC, Bull JJ, Meyers LA. From bad to good: fitness reversals and the ascent of deleterious mutations. *PLoS Comput Biol* 2006; 2: e141.
3. Cowperthwaite MC, Meyers LA. How mutational networks shape evolution: lessons from RNA models. *Annu Rev Ecol Evol Syst* 2007; 38: 203–30.
4. Kimura M. The role of compensatory neutral mutations in molecular evolution. *J Genet* 1985; 64: 7–19.
5. Wright S. The roles of mutation, inbreeding, crossbreeding and selection in evolution. *Proc VI Intl Cong Genet* 1932; 1: 356–66.
6. Dawkins R. *Climbing mount improbable*. W.W. Norton, New York, 1996.
7. Hurst JM, McMillan LEM, Porter CT, Allen J, Fakorede A, Martin ACR. The SAAPdb web resource: a large-scale structural analysis of mutant proteins. *Hum Mut* 2009; 30: 616–24.
8. Sherry ST, Ward MH, Kholodov M, Baker J, Phan L, Smigielski EM, Sirotkin K. dbSNP: the NCBI database of genetic variation. *Nucleic Acids Res* 2001; 29: 308–11.
9. Zhang G, Pei Z, Krawczak M, Ball EV, Mort M, Kehrer-Sawatzki H, Cooper DN. Triangulation of the human, chimpanzee, and neanderthal genome sequences identifies potentially compensated mutations. *Hum Mut* 2010; 31: 1286–93.
10. Barešić A, Hopcroft LEM, Rogers HH, Hurst JM, Martin ACR. Compensated pathogenic deviations: analysis of structural effects. *J Mol Biol* 2010; 396: 19–30.
11. Kondrashov AS, Sunyaev S, Kondrashov FA. Dobzhansky-Muller incompatibilities in protein evolution. *Proc Natl Acad Sci USA* 2002; 99: 14878–83.
12. McMillan LEM, Martin ACR. Automatically extracting functionally equivalent proteins from SwissProt. *BMC Bioinf* 2008; 9: 418.
13. Fitch WM. Homology a personal view on some of the problems. *Trends Genet* 2000; 16: 227–31.
14. Shibata S, Sasaki M, Miki T, Shimamoto A, Furuichi Y, Katakira J, Yoneda Y. Exportin-5 orthologues are functionally divergent among species. *Nucleic Acids Res* 2006; 34: 4711–21.
15. Pazos F, Valencia A. Protein co-evolution, co-adaptation and interactions. *EMBO J* 2008; 27: 2648–55.
16. Cordell HJ. Epistasis: what it means, what it doesn't mean, and statistical methods to detect it in humans. *Hum Mol Genet* 2002; 11: 2463–8.
17. Zuker M, Stiegler P. Optimal computer folding of large RNA sequences using thermodynamics and auxiliary information. *Nucleic Acids Res* 1981; 9: 133–48.
18. Zuker M. On finding all suboptimal foldings of an RNA molecule. *Science* 1989; 244: 48–52.
19. Hofacker I, Fontana W, Stadler P, Bonhoeffer L, Tacker M, Schuster P. Fast folding and comparison of RNA secondary structures. *Monatsh Chem* 1994; 125: 167–88.
20. Bernhart SH, Tafer H, Mückstein U, Flamm C, Stadler PF, Hofacker IL. Partition function and base pairing probabilities of RNA heterodimers. *Algorithms Mol Biol* 2006; 1: 3.
21. Mathews D. Revolutions in RNA secondary structure prediction. *J Mol Biol* 2006; 359: 526–32.
22. Mathews D, Turner D. Prediction of RNA secondary structure by free energy minimization. *Curr Opin Struct Biol* 2006; 16: 270–8.
23. Doudna JA. Structural genomics of RNA. *Nat Struct Biol* 2000; 7 (Suppl): 954–6.
24. Poon A, Davis BH, Chao L. The coupon collector and the suppressor mutation: estimating the number of compensatory mutations by maximum likelihood. *Genetics* 2005; 170: 1323–32.
25. Dawson K, Thorpe RS, Malhotra A. Estimating genetic variability in non-model taxa: a general procedure for discriminating sequence errors from actual variation. *PLoS One* 2010; 5: e15204.
26. Kulathinal RJ, Bettencourt BR, Hartl DL. Compensated deleterious mutations in insect genomes. *Science* 2004; 306: 1553–4.
27. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J Mol Biol* 1990; 215: 403–10.
28. Finn RD, Mistry J, Tate J, Cogill P, Heger A, Pollington JE, Gavin OL, Gunasekaran P, Ceric G, Forslund K, Holm L, Sonnhammer ELL, Eddy SR, Bateman A. The Pfam protein families database. *Nucleic Acids Res* 2010; 38: D211–22.
29. DePristo MA, Weinreich DM, Hartl DL. Missense meanderings in sequence space: a biophysical view of protein evolution. *Nat Rev Genet* 2005; 6: 678–87.
30. Henikoff S, Henikoff JG. Amino acid substitution matrices from protein blocks. *Proc Natl Acad Sci USA* 1992; 89: 10915–9.
31. Ferrer-Costa C, Orozco M, de la Cruz X. Characterization of compensated mutations in terms of structural and physicochemical properties. *J Mol Biol* 2007; 365: 249–56.
32. Amberger J, Bocchini CA, Scott AF, Hamosh A. McKusick's Online Mendelian Inheritance in Man (OMIM). *Nucleic Acids Res* 2009; 37: D793–6.

33. McKusick VA. Online Mendelian Inheritance in Man (OMIM) (TM). McKusick-Nathans Institute for Genetic Medicine, Johns Hopkins University (Baltimore, MD) and National Center for Biotechnology Information, National Library of Medicine (Bethesda, MD), 2011. Available at: <http://www.ncbi.nlm.nih.gov/omim/>.
34. Povolotskaya IS, Kondrashov FA. Sequence space and the ongoing expansion of the protein universe. *Nature* 2010; 465: 922–6.
35. Bentley DR. Whole-genome re-sequencing. *Curr Opin Genet Dev* 2006; 16: 545–52.
36. Yue P, Melamud E, Moulton J. SNPs3D: candidate gene and SNP selection for association studies. *BMC Bioinf* 2006; 7: 166.
37. Uzun A, Leslin CM, Abyzov A, Ilyin V. Structure SNP (St-SNP): a web server for mapping and modeling nsSNPs on protein structures with linkage to metabolic pathways. *Nucleic Acids Res* 2007; 35: W384–92.
38. Yip YL, Scheib H, Diemand AV, Gattiker A, Famiglietti LM, Gasteiger E, Bairoch A. The Swiss-Prot variant page and the ModSNP database: a resource for sequence and structure information on human protein variants. *Hum Mut* 2004; 23: 464–70.
39. Dantzer J, Moad C, Heiland R, Mooney S. MutDB services: interactive structural analysis of mutation data. *Nucleic Acids Res* 2005; 33: W311–4.
40. Karchin R, Diekhans M, Kelly L, Thomas DJ, Pieper U, Eswar N, Haussler D, Sali A. LS-SNP: large-scale annotation of coding non-synonymous SNPs based on multiple information sources. *Bioinformatics* 2005; 21: 2814–20.
41. Stitzel NO, Binkowski TA, Tseng YY, Kasif S, Liang J. topoSNP: a topographic database of non-synonymous single nucleotide polymorphisms with and without known disease association. *Nucleic Acids Res* 2004; 32: D520–2.
42. Ng PC, Henikoff S. SIFT: predicting amino acid changes that affect protein function. *Nucleic Acids Res* 2003; 31: 3812–4.
43. Reumers J, Maurer-Stroh S, Schymkowitz J, Rousseau F. SNPeffect v2.0: a new step in investigating the molecular phenotypic effects of human non-synonymous SNPs. *Bioinformatics* 2006; 22: 2183–5.
44. Ramensky V, Bork P, Sunyaev S. Human non-synonymous SNPs: server and survey. *Nucleic Acids Res* 2002; 30: 3894–900.
45. Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P, Kondrashov AS, Sunyaev SR. A method and server for predicting damaging missense mutations. *Nat Methods* 2010; 7: 248–9.
46. Thomas PD, Campbell MJ, Kejariwal A, Mi H, Karlak B, Daverman R, Diemer K, Muruganujan A, Narechania A. PANTHER: a library of protein families and subfamilies indexed by function. *Genome Res* 2003; 13: 2129–41.
47. Bao L, Zhou M, Cui Y. nsSNPAnalyzer: identifying disease-associated nonsynonymous single nucleotide polymorphisms. *Nucleic Acids Res* 2005; 33: W480–2.
48. Callaway E. Mutation-prediction software rewarded, 2010. *Nat News* DOI: 10.1038/news.2010.679.
49. Vogelstein B, Lane D, Levine AJ. Surfing the p53 network. *Nature (London)* 2000; 408: 307–10.
50. Lakin ND, Jackson SP. Regulation of p53 in response to DNA damage. *Oncogene* 1999; 18: 7644–55.
51. Chao C, Saito S, Kang J, Anderson CW, Appella E, Xu Y. p53 transcriptional activity is essential for p53-dependent apoptosis following DNA damage. *EMBO J* 2000; 19: 4967–75.
52. Nikolova PV, Wong KB, DeDecker B, Henckel J, Fersht AR. Mechanism of rescue of common p53 cancer mutations by second-site suppressor mutations. *EMBO J* 2000; 19: 370–8.
53. Joerger AC, Ang HC, Veprintsev DB, Blair CM, Fersht AR. Structures of p53 cancer mutants and mechanism of rescue by second-site suppressor mutations. *J Mol Biol* 2005; 280: 16030–7.
54. Danziger SA, Swamidass SJ, Zeng J, Dearth LR, Lu Q, Chen JH, Cheng J, Hoang VP, Saigo H, Luo R, Baldi P, Brachmann RK, Lathrop RH. Functional census of mutation sequence spaces: the example of p53 cancer rescue mutants. *IEEE/ACM Trans Comput Biol Bioinform* 2006; 3: 114–25.
55. Baronio R, Danziger SA, Hall LV, Salmon K, Hatfield GW, Lathrop RH, Kaiser P. All-codon scanning identifies p53 cancer rescue mutations. *Nucleic Acids Res* 2010; 38: 7079–88.
56. Danziger SA, Zeng J, Wang Y, Brachmann RK, Lathrop RH. Choosing where to look next in a mutation sequence space: active learning of informative p53 cancer rescue mutants. *Bioinformatics* 2007; 23: i104–14.
57. Danziger SA, Baronio R, Ho L, Hall L, Salmon K, Hatfield GW, Kaiser P, Lathrop RH. Predicting positive p53 cancer rescue regions using Most Informative Positive (MIP) active learning. *PLoS Comput Biol* 2009; 5: e1000498.
58. Friedler A, Hansson LO, Veprintsev DB, Freund SMV, Rippl TM, Nikolova PV, Proctor MR, Rudiger S, Fersht AR. A peptide that binds and stabilizes p53 core domain: chaperone strategy for rescue of oncogenic mutants. *Proc Natl Acad Sci USA* 2002; 99: 937–42.
59. Friedler A, DeDecker BS, Freund SMV, Blair C, Rudiger S, Fersht AR. Structural distortion of p53 by the mutation R249S and its rescue by a designed peptide: implications for ‘mutant conformation’. *J Mol Biol* 2004; 336: 187–96.
60. Boeckler FM, Joerger AC, Jaggi G, Rutherford TJ, Veprintsev DB, Fersht AR. Targeted rescue of a destabilized mutant of p53 by an in silico screened drug. *Proc Natl Acad Sci USA* 2008; 105: 10360–5.
61. Girardini JE, Del Sal G. Improving pharmacological rescue of p53 function: RITA targets mutant p53. *Cell Cycle* 2010; 9: 2059–62.
62. Zhao CY, Grinkevich VV, Nikulenkov F, Bao W, Selivanova G. Rescue of the apoptotic-inducing function of mutant p53 by small molecule RITA. *Cell Cycle* 2010; 9: 1847–55.
63. Zhao CY, Szekely L, Bao W, Selivanova G. Rescue of p53 function by small-molecule RITA in cervical carcinoma by blocking E6-mediated degradation. *Cancer Res* 2010; 70: 3372–81.
64. Auerbach G, Herrmann A, Bracher A, Bader G, Gutlich M, Fischer M, Neukamm M, Garrido-Franco M, Richardson J, Nar H, Huber R, Bacher A. Zinc plays a key role in human and bacterial GTP cyclohydrolase I. *Proc Natl Acad Sci USA* 2000; 97: 13567–72.
65. Berman HM, Battistuz T, Bhat TN, Bluhm WF, Bourne PE, Burkhardt K, Feng Z, Gilliland GL, Iype L, Jain S, Fagan P, Marvin J, Padilla D, Ravichandran V, Schneider B, Thanki N, Weissig H, Westbrook JD, Zardecki C. The Protein Data Bank. *Acta Crystallogr D Biol Crystallogr* 2002; 58: 899–907.
66. Shih HL, Brady J, Karplus M. Structure of proteins with single-site mutations: a minimum perturbation approach. *Proc Natl Acad Sci USA* 1985; 82: 1697–700.
67. Martin ACR, Facchiano AM, Cuff AL, Hernandez-Boussard T, Olivier M, Hainaut P, Thornton JM. Integrating mutation data and structural analysis of the tp53 tumor-suppressor protein. *Hum Mut* 2002; 19: 149–64.
68. Krissinel E, Henrick K. Inference of macromolecular assemblies from crystalline state. *J Mol Biol* 2007; 372: 774–97.
69. Shi D, Morizono H, Ha Y, Aoyagi M, Tuchman M, Allewell NM. 1.85-Å resolution crystal structure of human ornithine

- transcarbamoylase complexed with N-phosphonacetyl-L-ornithine. Catalytic mechanism and correlation with inherited deficiency. *J Mol Biol* 1998; 273: 34247–54.
70. Suriano G, Azevedo L, Novais M, Boscolo B, Seruca R, Amorim A, Ghibaudi EM. In vitro demonstration of intra-locus compensation using the ornithine transcarbamylase protein as model. *Hum Mol Genet* 2007; 16: 2209–14.
71. Azevedo L, Suriano G, van Asch B, Harding RM, Amorim A. Epistatic interactions: how strong in disease and evolution? *Trends Genet* 2006; 22: 581–5.