# Experimental Support for Regarding Functional Classes of Proteins to be Highly Isolated from Each Other

**Michael J. Behe**

*Department of Biological Science*
*Lehigh University*

## Preliminary Remarks

In writing on the topic of naturalism and evolution the problem arises of what to call the contending camps. The difficulty comes from the fact that, although the term "evolutionist" is often used to refer to persons who demand the unrelenting application of physical laws to all phenomena in the universe, many other persons who are opposed to this view are perfectly willing to concede that a limited number of phenomena can be explained by Darwinistic principles. Similarly, although a term like "creationist" brings to mind champions of a young-earth theory, it is often applied to persons who do not defend that thesis but do contend that natural laws have at some points been superseded by a supernatural agency. Since the focus of this symposium is the sufficiency of natural law, and in order to avoid the confusing terminology discussed above, in this essay I will use the term "believer" for those who believe in the universal application of natural law and the term "skeptic" for those who doubt it. This has the advantage of using terms for each side that the opposite side generally regards positively. Perhaps this will go a little way toward promoting the good will that this conference strives for.

## Introduction

Several years ago the fossilized remains of an extinct species of whale were unearthed in the Zeuglodon valley of Egypt. The particular aspect of the fossil which excited archeologists and science writers was the fact that the whale apparently had functional legs and feet. From the condition of the fossilized leg bones it could be discerned by trained eyes that the legs were well-muscled and thus must have been actively used during the life of the whale. Now, the Washington Post ran a story on the discovery. Along with the article was a drawing of a modern whale and an ancient whale, showing the differences in their shapes but similarities in their lengths. Also included in the illustration, down in the lower right hand corner, was a drawing of an animal that looked for all the world like a scruffy dog. Underneath the dog was the caption "Mesonychid, the ancestor of the whales". In the story it was explained that

> Most researchers agree the earliest whales descended from a line of large carnivorous beasts the size of wolves and bears. These furry land mammals, known as mesonychids, ran around on four legs. But for unknown reasons, some mesonychids evolved into forms that returned to the sea, from which all life originally arose. The legs found on primitive whales are remnants from their time on land.

> Washington Post, July 13,1990.

Even allowing for the enthusiasms of the popular press the story reflects the way in which a theory, here evolution, is allowed to supply "facts" which the evidence in no way justifies. I discussed this article with my students in a course I teach for Freshmen, entitled "Popular Arguments on Evolution".

The course is intended to develop critical reasoning skills, using popular books that have opposing viewpoints on evolution as the vehicle. This past semester we read, side by side, Richard Dawkins' The Blind Watchmaker and Michael Denton's Evolution: A Theory in Crisis. This forced the students to argue over the meaning of observations, without the automatic social support that usually goes to proponents of evolution in academic settings. The students themselves, after reading the Post's article, pointed out that there is no reason to suppose that the ancient whale appeared on earth before the modern whale since modern whales have vestigial legs which could have developed into the functional legs of the Zeuglodon whale. For the same reason, the students noted, the discovery does not represent the development of a new trait or even the loss of an old one. Finally, and most glaringly obvious, if random evolution is true there must have been a large number of transitional forms between the mesonychid and the ancient whale: Where are they? It seems like quite a coincidence that of all the intermediate species that must have existed between the mesonychid and whale, only species that are very similar to the end species have been found. The students concluded that the fossil whale, although a fascinating discovery for natural history, was no evidence for the Post's evolutionary scenario.

I have started my contribution to this symposium with a discussion of the Zeuglodon whale because it is a paradigmatic example of evolutionary argumentation: a small change in a preexisting structure is used to argue to massive changes involving completely new structures or functions. It is like arguing that because a man can jump over a fissure five feet wide, then given enough time he could jump over the Grand Canyon. Now, a believer in the unabating rule of natural law would argue that the man could jump over the Grand Canyon if there were ledges and buttes for him to use as stepping stones. The skeptic would ask to be shown the stepping stones. This essay will examine how the search is going for stepping stones in one area of biochemistry: that of protein structure. We will see that, without a prior committment to naturalism, there is little reason to suppose that stepping stones exist in the canyon separating functional classes of proteins.

## Protein Structure

I ask for the patience of those who already have a working knowledge of protein structure, but in order to make sure that everyone reading this essay has the necessary background I will spend a little time discussing some fundamentals.

Although most people think of proteins as something you eat, one of the major food groups, when they reside in the body of an uneaten animal or plant proteins serve a different purpose. Proteins are the machinery of living tissue that builds the structures and carries out the chemical reactions necessary for life. For example, the conversion of foodstuffs to biologically- usable forms of energy is carried out, step by step, by part of a group of proteins called enzymes. Skin is made in large measure of a protein called collagen. When light impinges on your retina it interacts first with a protein called rhodopsin. As can be seen even by this limited number of examples proteins carry out amazingly diverse functions. However, in general a given protein can perform only one or a few functions: rhodopsin can not form skin and collagen can not interact usefully with light. Therefore a typical cell contains thousands and thousands of different types of proteins to perform the many tasks necessary for life, much like a carpenter's workshop might contain many different kinds of tools for various carpentry tasks.

What do these versatile tools look like? The basic structure of proteins is quite simple: they are formed by hooking together in a chain discrete subunits called amino acids. Now, although the protein chain can consist of anywhere from about 50 to about 1,000 amino acid links, each position can only contain one of twenty different amino acids. In this they are much like words: words can come in various lengths but they are made up from a discrete set of 26 letters. As a matter of fact, biochemists often refer to each amino acid by a single letter abbreviation - G for glycine, S for serine, H for histidine,

and so forth. Each different kind of amino acid has a different shape and different chemical properties; for example, W is large but A is small, R carries a positive charge but E carries a negative charge, S prefers to be dissolved in water but I prefers oil, etc. Now, a protein in a cell does not float around like a floppy chain; rather, it folds up into a very precise structure which can be quite different for different types of proteins. This is done automatically through interactions such as a positively charged amino acid trying to get near a negatively charged one, oil-preferring amino acids trying to huddle together to exclude water, large amino acids being excluded from small spaces, etc. When all is said and done two different amino acid sequences - two different proteins - can be folded to structures as specific and different from each other as a three-eighths inch wrench and a jigsaw. And like the household tools, if the shape of the proteins is significantly warped then they fail to do their jobs.

## Proteins and Language

Because amino acid residues are often abbreviated by letters, because there are a similar number of letters and amino acids (26 vs. 20, respectively), and because a small protein consists of about 100 amino acids, many commentators have likened a functional protein (i.e. - one which has the correct shape to be able to do a particular job) to a functional sentence (i.e. - one which obeys the rules of English grammar) of about 100 letters. My students in "Popular Arguments on Evolution" found it particularly interesting that both believers and skeptics used this kind of analogy in their writings, but that their reasonings brought them to opposite conclusions. The skeptic typically argues that a monkey banging away at a typewriter (monkeys and typewriters are very popular) would be very unlikely to produce an intelligible, grammatically-correct sentence like "Drop the anchor in one hour." in a reasonable length of time. Near misses don't count for the skeptic since the change of even one letter would break a spelling or grammar rule, or change the sense of the sentence. Needless to say the hour would most likely pass, and the anchor remain undropped, before the monkey produced the correct sentence.

Believers in the universal application of physical law take a different approach with their monkey and typewriter. Their argument generally goes something like this. Suppose in his first try the monkey typed "bsqm dshcbbbk,RR .nsurlei aknex". Admittedly this is poor grammar, but it's the only sentence we've got. Since living systems reproduce, and since there is Darwinian competition, the bad sentence will be reproduced until a better one comes along. Now suppose in his second try the monkey typed a 'p' in the fourth position and a 'u' in the penultimate position. Well, since these are closer to the target sentence we will throw out the original sentence and keep "bsqp dshcbbbk,RR .nsurlei aknux". After a few more rounds perhaps the monkey has got a few more letters correct, say a 'd' in the first position and a 'ch' in the 13 and 14 positions. Now we have "dsqp dshcbbbchRR .nsurlei aknux". Since this has more matches with the target sentence we'll keep it and throw out the last sentence. After perhaps 50 rounds we get to "dsop dhe abchRR in uneei hnur." Breed from this. In another 50 rounds or so we arrive triumphantly at our target "Drop the anchor in one hour."

The above argument in its pure form can only be convincing to persons already convinced. It asserts a functional difference between two nonsensical strings of letters. No person, or machine for that matter, looking for a sentence would notice a difference between "bsqm dshcbbbk,RR .nsurlei aknex" and "bsqp dshcbbbk,RR .nsurlei aknux." It is only because the believer has a distant goal in mind that he chooses one nonsense character string over the other. In the believers' argument the analogy of proteins to language is implicitly abandoned in the first rounds of the monkey's typing, since the character string does not have to obey any rules of spelling or grammar. The analogy to language is used simply to try to impress the unwary with the apparent production of sense from nonsense. My students in "Popular Arguments on Evolution" were uneasy with this argument when they read it in Dawkins' book, but they could not refute it. It is not easy for the casual reader to see that the illusion of steady, gradual evolution to a functional sentence is produced by an intellect, either the believer's

directly or in some cases a computer program written by him, guiding the result to a distant goal. This of course is the antithesis of Darwinian evolution.

But perhaps there is a middle ground between the skeptic's insistence on absolute grammatical correctness and the believer's abandonment of grammatical rules. Suppose we allowed the vowels in the sentence to vary to produce something like "Drep tha enchir on une hoir". Such a sentence could probably still be recognized by someone, perhaps a sailor, even though all the words are misspelled. Or, alternatively, suppose we vary some consonants: "Trof tte ankhow im ode hous". Clearly some misspelled words would be easier to recognize than others and some letter substitutions ('t' for 'd', 'k' for 'c') would be easier to follow than others ('r' for 't', 'l' for 'g'). The ability of a sentence like that to function would depend a lot on the reader and the context.

To put this back into a protein context, it might be possible for a protein to tolerate a lot of amino acid substitutions and remain functional. (Again, when talking about proteins `functional' means folded to a discrete, stable structure.) And in fact it has been known for a long time that this is true. Analogous proteins from different species, for example human hemoglobin and horse hemoglobin, have differences between their amino acid sequences, yet fold to discrete and closely similar structures. But what is the limit to tolerance for amino acid changes? Are proteins significantly more tolerant to changes in 'spelling' than words are? Is there a point at which, like our sentences above, further changes will render a protein nonfunctional? What then is the probability of finding some member of a particular class in a reasonable time in a nondirected search? These are empirical questions and, although they can be speculated upon in the absence of relevant data, such speculations must be radically curtailed when data are available. A direct approach to the question, How isolated are functional protein sequences? would have been experimentally impossible twenty years ago, before the molecular biological revolution. But since the development of powerful tools to probe the molecules of life an answer to that question appears to be within reach. Progress in this area is the topic of the following sections.

## How Rare Are Functional Proteins?

The observation that analogous proteins from different species could differ from each other, often by quite a bit, and yet retain the same compact shape led workers in the field to speculate that perhaps the exact identity of an amino acid at a particular position in a protein was not as important as its overall chemical properties. So, for example, if one finds an I at position 10 of hedgehog hemoglobin and an L in position 10 of the analogous protein from skunk, then perhaps the important feature is that both I and L prefer an oily environment, and maybe any other amino acid, such as W, F, or V, that prefers a similar environment would also be suitable at that position. This is something like saying that in a language perhaps all of the vowels are interchangeable. Taking the idea further, perhaps amino acids, such as S, A, H, and T, that prefer a watery environment could form an interchangeable group, and perhaps charged amino acids (E, D, R, and K) another group.

Fifteen years ago a man named Hubert Yockey published an article in the *Journal of Theoretical Biology* (1) showing that these considerations could enormously reduce the odds against finding a functional protein by trial and error. If we do not insist on the perfect diction of the typical skeptic, but allow some slurred speech in proteins, then the probability of finding a small, functional protein of 100 amino acids in length is reduced from 1 in 10 to the 130 power to 1 in 10 to the 65 power - a reduction of 65 orders of magnitude! Yockey went on to show in the article that his calculation of 1 in 10&%, which he obtained from theoretical considerations, fit very closely with the number that could be calculated from considerations of the known sequence variability of the protein cytochrome c among many different species.

Now, the problem with Yockey's calculation for a believer in the sufficiency of natural law is that, although 10&% is enormously smaller than 10!# , it still is quite a big number. It has been calculated that there are about 10&% atoms in a galaxy. Thus, if Yockey was correct, the odds of finding a functional protein are about the same as finding one particular atom in the Milky Way. Not too likely. Well, if you were a believer how might you answer this challenge? One way is through obfuscation, like the production of sentences from nonsense character strings, as was discussed above. A second way is by claiming that Yockey's calculation is inaccurate and that the known sequences of cytochrome c that he used to buttress his work do not reflect all the possible sequences that could produce a folded protein. The best way, though, in the absence of relevant data, is to produce your own calculation, starting from a separate set of independent principles, and show that the odds are not quite so long as Yockey thought. This is what has been done in an elegant series of calculations from the laboratory of Ken Dill (2,3) at the University of California at San Francisco.

Dill's laboratory asked a question which can be paraphrased as follows: Given a ten-by-ten square matrix (like a big checkerboard) and a string of pearls containing both black beads and white beads, in how many ways can a string of 100 pearls be laid on the checkerboard so that each square contains one and only one pearl, and most of the black pearls are in the middle spaces of the board? This analogy is intended to represent a folding protein comprised of two types of amino acids - ones that prefer watery surroundings and ones that do not. After feeding this scenario into a computer the surprising result Dill's group obtained was that it wasn't that hard to fit the pearl necklace on the checkerboard in the right way. They then mathematically extrapolated from the two dimensional checkerboard to three dimensional space, and finally arrived at the conclusion that about 1 in 10! amino acid sequences would yield a folded protein. This is a much smaller number than Yockey's (the federal government spends 10! dollars -ten billion dollars- every three days) and brings the spontaneous generation of functional proteins into the realm of the credible.

Now the problem for a skeptic is how to refute Dill's calculation. It isn't easy since few people are as mathematically talented as he and since it's hard to disprove the simplifying assumptions his model contains. Skeptics are free to criticize the assumptions, but there is enough uncertainty in such things to allow believers to credibly tout Dill's calculation over Yockey's. To resolve this dilemma, to gain firm ground to stand on, hard experimental results are required. Fortunately in the past several years such results have been forthcoming from the laboratory of Robert Sauer (4-6) in the Department of Biology at the Massachusetts Institute of Technology. We now turn to those crucial experiments.

## Very Rare

In the past twenty years the science of molecular biology has made enormous strides. It is now literally possible, in laboratories with such expertise, to cut up a gene, rearrange it to suit yourself, and place it back in a functioning biological system. Since genes code for proteins, one can also produce proteins made-to-order in this manner. Sauer's laboratory, in order to answer questions about protein structure that interested them, took the genes for several viral proteins, systematically took out small pieces of them (corresponding to instructions for three amino acids at a time) and inserted altered pieces back in the genes. They did this, three amino acids 'codons' at a time, for the whole length of the gene. By clever manipulation of the altered pieces they were able to screen codons for all twenty amino acids at each position of the protein. This is like trying all 26 letters of the alphabet in turn at each position of a word. The altered genes were then placed in bacteria, which read the DNA code and produced chains of amino acids from them. It turns out that bacteria quickly destroy proteins that are not folded, so Sauer's group looked for the altered proteins that were not destroyed. By determining their sequences they could tell which amino acids in a given position were compatible with producing a folded, functional protein. And what did they see?

In some positions of the protein Sauer's group saw that a great deal of amino acid diversity could be tolerated. Up to 15 of the twenty amino acids could occur at some positions and still yield a functional, folded protein. However, at other positions in the amino acid sequence very little diversity could be tolerated. Many positions could accomodate only 3 or 4 different amino acids. Other positions had an absolute requirement for a particular amino acid; this means that if, say, a P does not appear at position 78 of a given protein the protein will not fold regardless of the proximity of the rest of the sequence to the natural protein. In terms of our sentence analogy, this is like saying that, yes, all vowels are interchangeable, but that if the last `r' is changed to any other letter, such as 's' ("Drop the anchor in one hous"), the protein sentence is no longer understandable.

Sauer's results can be used to calculate the probability of finding a given protein structure (6). We proceed in the following manner. If any of ten amino acids can appear in the first position of a given functional protein sequence then the odds are 1 in 2 that a nondirected search will place one of the allowed group there. If any of four amino acids can appear in the second position then the odds are 1 in 5 of finding one of that group, and the odds of finding the correct amino acids next to each other in the first two positions are one-half times one-fifth, which is one-tenth. Suppose in the third position there is an absolute requirement for G. Then the odds of getting a G at that position are one in twenty and the odds of getting the first three amino acids right are now up to one in two hundred. In this aspect it is like winning a trifecta in horse racing. Over the course of 100 amino acids in our small protein the odds quickly reach astronomical numbers.

From the actual experimental results of Sauer's group it can easily be calculated that the odds of finding a folded protein are about 1 in 10 to the 65 power (6). To put this fantastic number in perspective imagine that someone hid a grain of sand, marked with a tiny 'X', somewhere in the Sahara Desert. After wandering blindfolded for several years in the desert you reach down, pick up a grain of sand, take off your blindfold, and find it has a tiny 'X'. Suspicious, you give the grain of sand to someone to hide again, again you wander blindfolded into the desert, bend down, and the grain you pick up again has an 'X'. A third time you repeat this action and a third time you find the marked grain. The odds of finding that marked grain of sand in the Sahara Desert three times in a row are about the same as finding one new functional protein structure. Rather than accept the result as a lucky coincidence, most people would be certain that the game had been fixed.

The number of 1 in 10&%, arrived at by Sauer's experimental route, is virtually identical to the results obtained by Yockey's theoretical calculation and his deduction from natural cytochrome c sequences! It therefore strongly reinforces our confidence that a correct result has been obtained. Sauer's group obtained closely similar results for two different proteins: arc repressor (4) and lamda repressor (5,6). This means that all proteins that have been examined to date, either experimentally or by comparison of analogous sequences from different species, have been seen to be surrounded by an almost infinitely wide chasm of unfolded, nonfunctional, useless protein sequences. There are no ledges, no buttes, no stepping stones to cross the chasm. The conclusion that a reasonable person draws from this is that the laws of nature are insufficient to produce functional proteins and, therefore, functional proteins have not been produced through a nondirected search.

## Implications of Protein Sequence Isolation

The numerical concreteness of Sauer's and Yockey's results is breathtaking. When a skeptic sees a drawing of Mesonychid next to the Zeuglodon whale he intuitively realizes that the transformation is highly improbable. But how improbable? There is no way to put a quantitative measure on the difference between a dog-like animal and a whale, and believers in the relentless application of physical law take advantage of this by verbally minimizing the differences. The situation is otherwise with proteins. Because there is a discrete set of amino acids and a finite number of positions in a given protein, the odds of attaining a folded, functional protein can be calculated quite closely, but only if

the tolerance of proteins to amino acid substitution is known. Thanks to Sauer and Yockey we now have such quantitative data.

It is important to realize that Sauer's and Yockey's results hold whether or not the system can replicate and is subject to Darwinian selection. The odds against finding a new functional protein structure remain astronomical in either case. This is because Darwinian selection can only discriminate based on function and, with the exception of those found in living organisms, virtually all protein sequences are functionless. An amino acid sequence can be replicated and mutated in living organisms till the cows come home and the odds are still 1 in 10&% that a new functional protein class will be produced.

The problem of the isolation of functional protein sequences is a vivid illustration of the truth of the symposium thesis,

> Darwinism and neo-Darwinism as generally held and taught in our society carry with them an a priori commitment to metaphysical naturalism, which is essential to make a convincing case in their behalf.

The skeptic can accept Sauer's and Yockey's results with equanimity because his world is not necessarily limited to those phenomena that can be explained by naturalism. Furthermore, the skeptic can happily concede that many biological phenomena are explained by natural laws. He can agree that beak shape and wing color can change under selective pressure, or that different proteins in the same structural class, such as the alpha and beta chains of hemoglobin, may have arisen through Darwinistic mechanisms. But the believer in the universal application of physical law is stuck. He must maintain, against the evidence, that different protein classes, like cytochromes and immunoglobulins, found their way by raw luck through the vast, dark sea of nonfunctional sequences to the tiny islands of function we observe experimentally. He must maintain, without any evidence, that Mesonychid gave birth over time to the whale. And why, we ask, must he maintain these positions against impossible odds and without supporting evidence? Because, he replies, I can only measure material phenomena, and therefore nothing else exists.

In closing I would like to paraphrase Hubert Yockey (7), who in his career repeatedly pointed out facts that are not supposed to be mentioned in polite scientific company: "Since science has not the vaguest idea how (proteins) originated, it would only be honest to admit this to students, the agencies funding research, and the public."

# References

1. Yockey, H. P. (1978) "A Calculation of the Probability of Spontaneous Biogenesis by Information Theory", *Journal of Theoretical Biology* 67, 377-398.

2. Lau, K. F., & Dill, K. A. (1989) "A Lattice Statistical Mechanics Model of the Conformational and Sequence Spaces of Proteins", *Macromolecules* 22, 3986-3994.

3. Chan, H. S., & Dill, K. A. (1990) "Origins of Structure in Globular Proteins", *Proceedings of the National Academy of Sciences USA* 87, 6388-6392.

4. Bowie, J. U., & Sauer, R. T. (1989) "Identifying Determinants of Folding and Activity for a Protein of Unknown Structure", *Proceedings of the National Academy of Sciences USA* 86, 2152-2156.

5. Bowie, J. U., Reidhaar-Olson, J. F., Lim, W. A., & Sauer, R. T. (1990) "Deciphering the Message in Protein Sequences: Tolerance to Amino Acid Substitution", *Science* 247, 1306-1310.

6. Reidhaar-Olson, J. F., & Sauer, R. T. (1990) "Functionally Acceptable Substitutions in Two -Helical Regions of Repressor", *Proteins: Structure, Function, and Genetics* 7, 306-316.

7. Yockey, H. P. (1981) "Self Organization Origin of Life Scenarios and Information Theory", *Journal of Theoretical Biology* 91, 13-31.