

The Function of the Hereditary Materials: Biological Catalyses Reflect the Cell's Evolutionary History^{1,2}

BRUCE M. ALBERTS

Department of Biochemistry and Biophysics, University of California, San Francisco, San Francisco, California 94143-0448

SYNOPSIS. The recent discovery of specialized RNA molecules that function like enzymes suggests that cells evolved before there were proteins. Such RNA-based cells would have contained large numbers of mutually supportive RNA molecules, each with a different catalytic function. Protein synthesis probably evolved later and was catalyzed by some of these RNA molecules. Because DNA must have been a relatively late addition to the cell, it is reasonable to assume that all DNA functions evolved in the presence of powerful protein catalysts.

The above evolutionary perspective helps to explain why two different classes of catalytic mechanisms are used in present-day cells. The ancient processes of protein synthesis and pre-mRNA splicing are catalyzed by ribonucleoprotein particles, in which RNA catalysis still seems to play an important role. In contrast, late-evolving functions like DNA replication are catalyzed by efficient protein machines. By analogy, protein machines are also likely to mediate the processes that control the transcription of eucaryotic genes.

INTRODUCTION

The topic that John Moore assigned me for this article, "the function of the hereditary materials," is not only a broad one, but one about which we have an enormous amount of knowledge. Fortunately, however, I have been asked to try to present a fresh view here, rather than to summarize facts that are already well explained in textbooks (see, for example, Stent, 1971; Watson, 1976; Alberts *et al.*, 1983; Lewin, 1985). I have therefore decided to discuss the hereditary materials from a particular evolutionary perspective that I have found useful in thinking about the cell. As for all evolutionary analyses, the discussion can only be speculative, but I believe that the major ideas presented generate enough predictions about how cells function to be testable as we learn more details about their mechanisms.

In its most general terms, a cell is nothing more (or less) than a highly elaborate

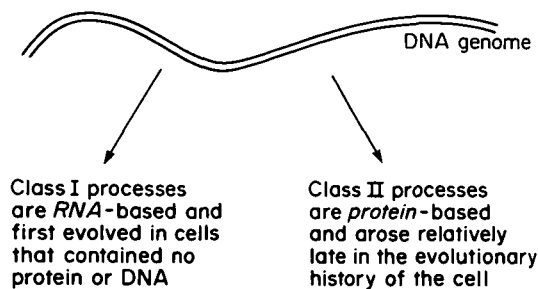
set of self-reproducing chemical reactions, and it therefore relies upon a large set of highly specific catalysts for its existence. In all of the cells we know, these catalysts are encoded by a DNA genome. As an overview of the argument to be pursued in this article, this DNA can be thought of as specifying two different families of catalysts, arbitrarily designated here as class I and class II, that derive from two very different epochs in the cell's evolutionary history (Fig. 1). The class II catalysts are the major type, being the enzymes that fill biochemistry textbooks; each is a protein molecule with an elaborately-folded catalytic surface. But at one time early in cell evolution there were no proteins as we know them, and it seems likely that the cell had to survive with RNA molecules as its major catalysts. The remnants of these early RNA-based reactions that survive constitute the class I reactions. Class I reactions can be recognized because they are catalyzed by ribonucleoprotein complexes, in which catalysis by RNA molecules still plays a significant role. These reactions are relatively complicated and unwieldy compared to class II reactions, which are often catalyzed by highly-efficient protein machines.

THE CENTRAL DOGMA IS MISLEADING

In any discussion of the hereditary materials, it is customary to start with the central dogma: DNA → RNA → protein.

¹ From the Symposium on *Science as a Way of Knowing—Genetics* presented at the Annual Meeting of the American Society of Zoologists, 27-30 December 1985, at Baltimore, Maryland.

² I would like to dedicate this article to John A. Moore, whose inexhaustible idealism and clear conceptual expositions of wide areas of biology have inspired and set standards for two generations of students, teachers, and authors; his textbook *Principles of Zoology* introduced me to "science as a way of knowing" nearly 30 years ago.



Class I processes are RNA-based and first evolved in cells that contained no protein or DNA

Class II processes are protein-based and arose relatively late in the evolutionary history of the cell

FIG. 1. As indicated, the DNA genome in today's cells specifies two classes of catalytic mechanisms, here designated as class I and class II. Many of the catalysts in each class evolved before cells contained DNA (see Fig. 8, below).

Today this dogma has become such an integral part of our thinking that many of us find it difficult to imagine any kind of life without DNA. Proteins are also crucial, being the end product of the pathway. They form the enzyme catalysts that determine what reactions occur in cells, as well as ensuring the necessary coordination between these reactions. RNA in this scheme is left with a subsidiary role, acting mainly as a "messenger boy" to carry out the instructions for protein synthesis that emanate from the DNA. This is the picture that emerges in all current textbooks, and it is what all of us teach to our students.

How far have the mighty fallen! For there is every reason to believe that, as far as life is concerned, RNA used to be at center stage. For example, it now seems certain that DNA was only a relatively late addition in the evolutionary history of the cell. Before there was DNA, RNA played the informational role in cells, specifying proteins by itself. And what about earlier, even before proteins? Many of us used to think that there could have been no life before proteins. This view has been shattered in the last few years by the startling (but in retrospect quite reasonable) discovery that RNA molecules can fold in complex ways that give them the same type of highly-specific and sophisticated catalytic activity that biochemists previously associated only with enzymes (Cech, 1985; Cech and Bass, 1986). It is the widespread implications of this revolutionary finding that I wish to discuss here, emphasizing the important

contributions this discovery makes to our understanding of cells today.

RNA MOLECULES CAN FUNCTION AS HIGHLY-SELECTIVE CATALYSTS

Efficient enzyme-like catalysis by a protein-free RNA molecule was first discovered in a rather obscure organism—the ciliated protozoon *Tetrahymena*. A *Tetrahymena* ribosomal RNA molecule had been shown to be produced from a larger RNA precursor molecule by RNA splicing, an event that seems to be ubiquitous in eucaryotic cells. (In RNA splicing, two non-contiguous stretches of RNA sequence in an RNA molecule are joined together with the concomitant removal of the nucleotide sequence between them; the latter sequence is called an *intron* sequence and it is normally discarded by the cell.) The surprise came when the scientists attempted to reproduce the ribosomal RNA splicing reaction in an *in vitro* system, so as to be able to study its mechanism. Although it was assumed that the reaction was catalyzed by enzymes, and thus would require a protein extract of lysed *Tetrahymena* cells, the control reactions in which the protein-free RNA was incubated without enzymes also underwent the splicing reaction (Cech *et al.*, 1981; Kruger *et al.*, 1982). It was subsequently shown that the intron sequence itself has an enzyme-like catalytic activity which carries out the reaction in two steps (Fig. 2).

More recently, the 400 nucleotide intron sequence has been synthesized in a test tube and studied in isolation from the rest of the ribosomal RNA transcript (Cech *et al.*, 1983; Bass and Cech, 1984). This sequence folds up to form a complex surface that functions like an enzyme in several reactions. For example, it can bind two specific substrates tightly ($K_m \approx 10^{-5} M$)—a guanine nucleotide and an RNA chain, catalyzing their covalent addition and severing the RNA chain at a unique sequence (Fig. 3).

In this model reaction, which mimics the first step in Figure 2, the same intron sequence can act over and over to cut many different RNA substrate chains. Although autocatalyzed RNA splicing is relatively rare, self-splicing RNAs with intron

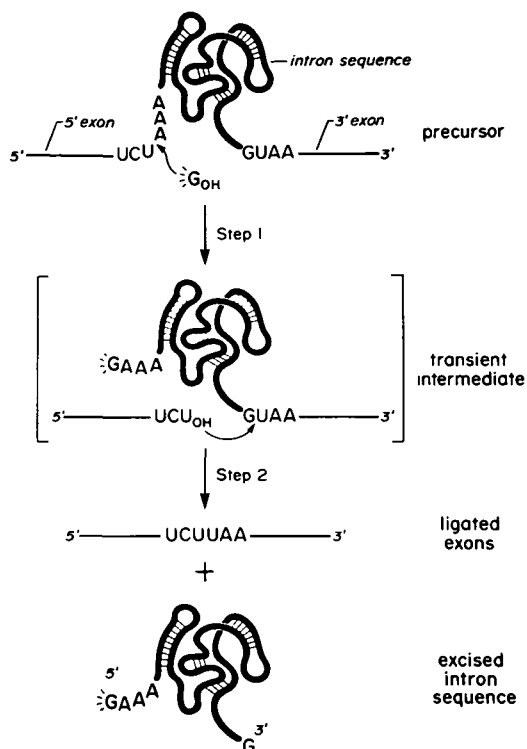


FIG. 2. Diagram of the self-splicing reaction in which an intron sequence catalyzes its own excision from a *Tetrahymena* ribosomal RNA molecule. As shown, the reaction is initiated when a G nucleotide adds to the intron sequence and cleaves the RNA chain; the new end of an RNA chain created then adds to the other side of the intron to complete the process (Zaug *et al.*, 1983).

nucleotide sequences that are related to the one found in *Tetrahymena* have been discovered in other types of cells, including fungi and bacteria (Garriga and Lambowitz, 1984; Belfort *et al.*, 1985). It has therefore been suggested that this RNA sequence may derive from a very ancient one, predating the time that the eucaryotic and procaryotic lineages branched off from each other about 1.5 billion years ago. Other families of catalytic RNAs also exist; for example, an RNA-protein complex that recognizes tRNA precursors and cleaves them at specific sites has RNA and not protein as its major catalyst (Guerrier-Takada *et al.*, 1983). Last but not least, the ribosome itself—lying at the center of cellular biochemistry as the mediator of protein synthesis—is now thought to be an RNA-

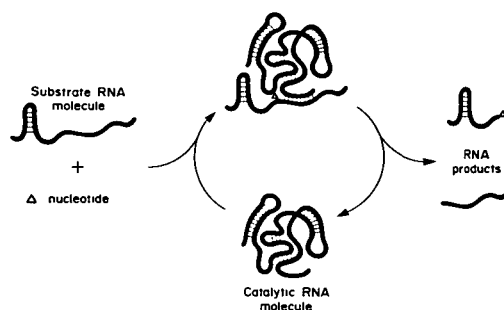


FIG. 3. Schematic diagram of an enzyme-like reaction catalyzed by the purified *Tetrahymena* intron sequence. In this reaction, which models step 1 in Figure 2, both a specific substrate RNA molecule and a G nucleotide become tightly bound to the surface of the catalytic RNA. The nucleotide is then added to the substrate RNA molecule, cleaving it at a specific site. The release of the two RNA product chains frees the intron sequence for further cycles of reaction. This figure is based on unpublished experiments that have been independently carried out by T. Cech and J. Szostack (personal communications).

based catalyst. Most of the ribosomal proteins are suspected to play only a supporting role to the ribosomal RNAs, which make up more than half of the mass of the ribosome and are postulated to have direct catalytic roles in protein synthesis (Woese, 1980; Noller, 1984).

How is it possible for an RNA molecule to act like an enzyme? It is misleading to think of RNA as a linear polymer containing occasional Watson-Crick base-pair interactions: as do proteins, RNA molecules often fold up in highly specific ways. Most of what we know about RNA folding has been derived from structural studies of a family of unusually short RNA molecules, the transfer RNAs. These molecules, which are only about 70 to 90 nucleotides long, have the three-dimensional conformation outlined in Figure 4B, as determined by X-ray crystallographic analyses (Rich and Kim, 1978). The highly folded molecule is held together by a substantial number of tertiary bonding interactions, some of which are indicated on the simple "cloverleaf" representation of the same tRNA molecule in Figure 4A. Thus, it is incorrect to think of RNA as being capable of only simple Watson-Crick base pairing. Once one admits the possibility of tertiary bonds, it is easy to see how larger

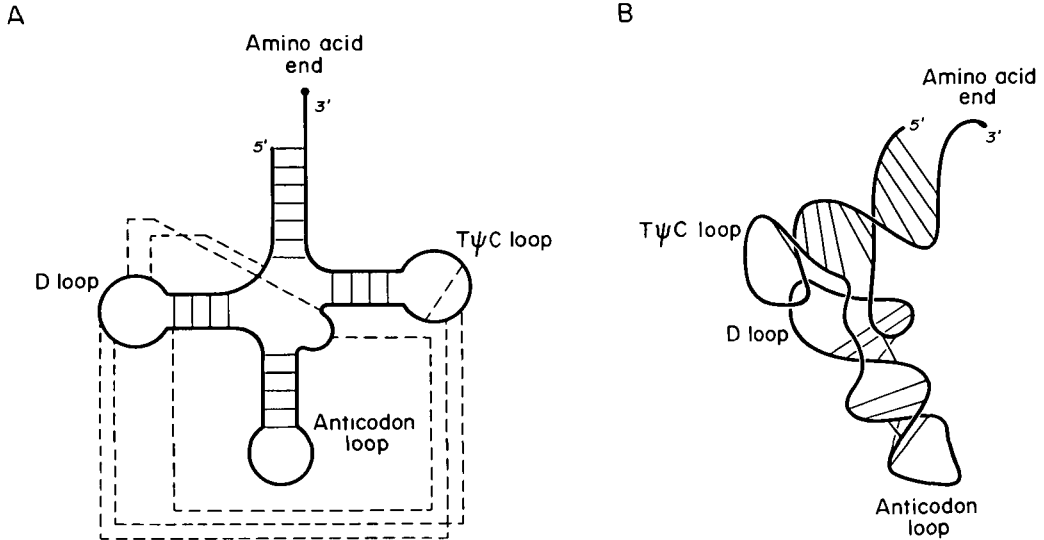


FIG. 4. An outline of the structure of a transfer RNA (tRNA) molecule. (A) A two dimensional view of the structure, shown in its open "cloverleaf form." The thin solid lines represent base pairs in short double-helical regions. The dotted lines connect some of the nucleotides that form tertiary bonding interactions with each other in the folded form that is illustrated in (B) (for further details, see Rich and Kim, 1978).

RNAs, such as the 400 nucleotide intron sequence just discussed, can fold up to form very complex and sophisticated surfaces with a powerful catalytic activity.

IN THE FIRST CELLS, RNA MOLECULES
MAY HAVE FUNCTIONED BOTH AS THE
SOURCE OF GENETIC INFORMATION
AND AS THE MAIN CATALYSTS

The catalytic potential of folded RNA molecules makes it much easier to imagine how the first cells arose on the earth. One suspects that a crucial early event was the evolution of an RNA molecule that could catalyze its own replication, thereby reproducing itself by autocatalysis (Fig. 5, stage 1). Eventually, evolution would select for a mutually-supportive collection of catalytic RNA molecules, some catalyzing the replication of the others (Fig. 5, stage 2). For example, one of these catalytic RNAs might have helped in the production of nucleotide precursors for RNA synthesis. Another might have catalyzed the accumulation of lipid-like molecules to form primitive membranes that isolated each self-replicating RNA family from its neighbors. Once such primitive "cells" were formed, very efficient cycles of mutation

and natural selection would occur: different collections of mutually-supportive RNA molecules could now be selected for their increasing fitness as self-reproducing units (Eigen *et al.*, 1981).

Under the pressure of evolution, the RNA molecules in these primitive RNA-based cells would be expected to acquire many of the same properties that enzymes have in cells today. For example, some of these RNAs presumably bound small molecule "coenzymes" to their active surface, which allowed them to increase the chemical versatility of their catalyses. Moreover, to permit homeostasis, feedback regulation could have evolved; such regulation would be mediated through allosteric changes in the structure of RNA catalysts caused by the binding of specific metabolites. Finally, RNA molecules could have harnessed chemical energy to do useful work through organized allosteric changes in their shape; as is found for proteins, the energetically-favorable hydrolysis of ligands bound to the RNA surface could induce these shape changes. Now that we realize that RNA molecules can be such powerful catalysts, it seems reasonable to postulate that RNA-based cells of this type

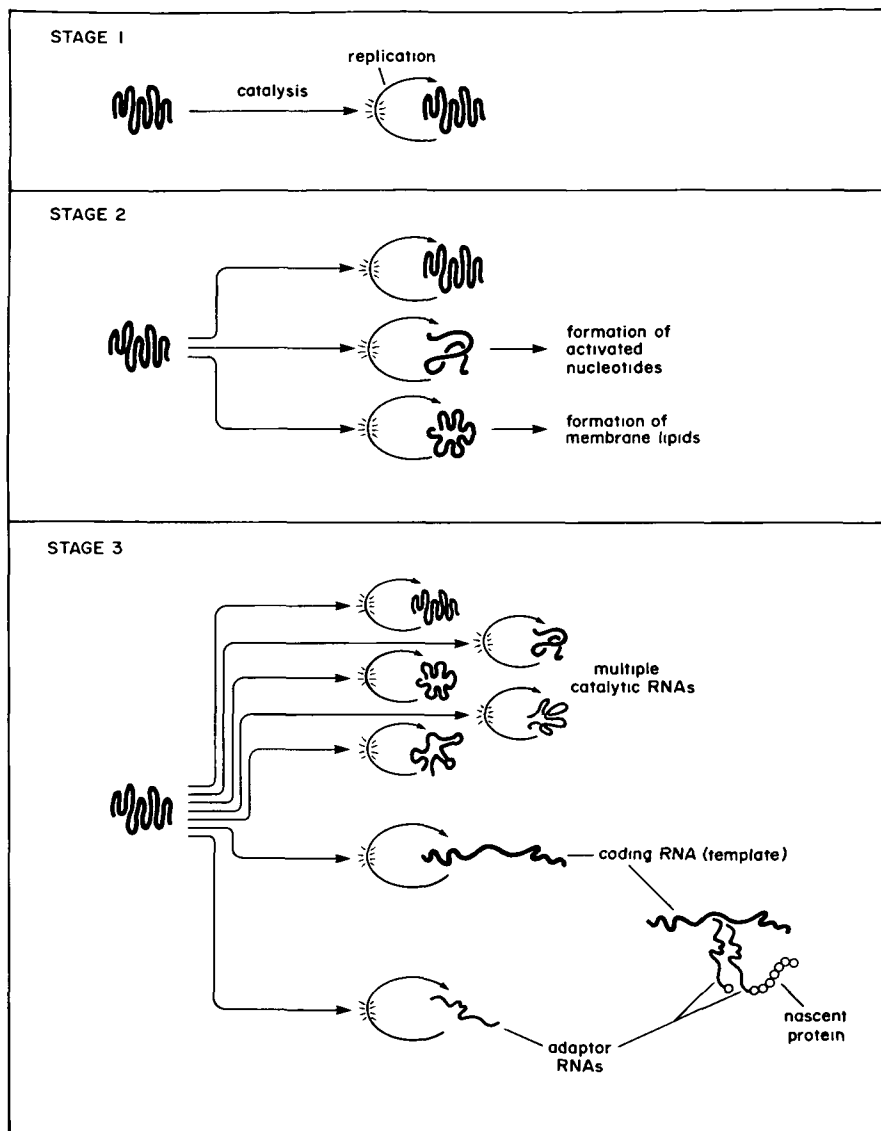


FIG. 5. Schematic illustration of some possible early stages in the evolution of cells. In stage 1, an RNA molecule that catalyzes its own synthesis is illustrated. In reality, such autocatalysis of replication may have required a set of several RNA molecules. The first cells are suggested to have formed in stage 2, when membranes enclosed a set of mutually-supportive catalytic RNA molecules. In stage 3, protein synthesis evolved in these RNA-based cells.

became quite sophisticated chemically. This hypothesis also makes it much easier to understand how the complex process of protein synthesis eventually evolved.

The chemistry involved in protein synthesis presumably developed over a long period of time. Initially, various catalytic RNA molecules would have experimented

with joining amino acids together without a template; in this way they could produce short peptides with useful chemical reactivities. In its first version, template-directed protein synthesis probably required only a coding RNA molecule and a set of "adaptor RNAs," as illustrated in Figure 5 (stage 3). The early adaptors, the

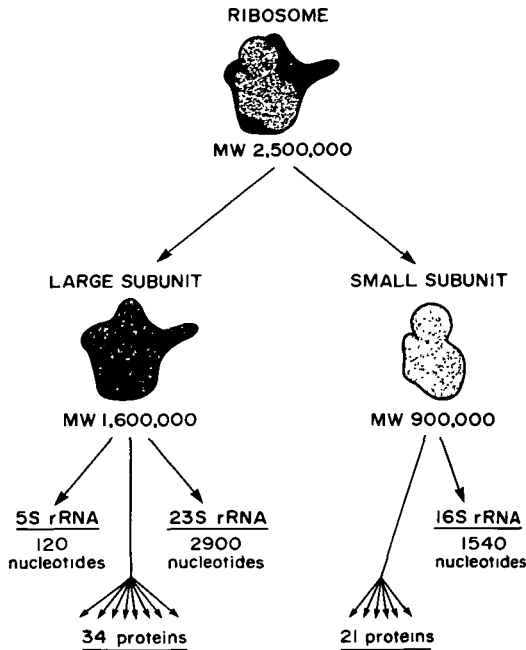


FIG. 6. The structure of the bacterial ribosome. More than 60 percent of the mass of the ribosome is RNA, the remainder being protein (Lake, 1985).

precursors of modern transfer RNAs, are likely to have bound specific amino acids directly, activating them for the subsequent synthesis of polypeptides without requiring other catalysts (Crick, 1968). Obviously, there must have been a relatively simple set of adaptor RNAs and only a limited repertoire of amino acids in such early cells (Crick, 1968; Orgel, 1968; Crick *et al.*, 1976; Hopfield, 1978; Woese, 1980; Crothers, 1982).

Another important early development in protein synthesis was presumably the evolution of new RNA-catalysts that promoted adaptor RNA binding to the coding RNA strand and the subsequent polymerization of the amino acids attached to these adaptors. These catalysts would have been the precursors of modern ribosomal RNA molecules.

While the general evolutionary pathway sketched in Figure 5 seems reasonable, the details are not crucial for the subsequent arguments. The important point is that, whatever the evolutionary pathway, pro-

tein synthesis must have evolved in a cell that lacked proteins as we now know them. It is therefore clear that specific RNA molecules were the major catalysts that made the evolution of protein synthesis possible.

NOTHING ABOUT PROTEIN SYNTHESIS MAKES SENSE EXCEPT IN THE LIGHT OF EVOLUTION

Today, protein synthesis is catalyzed by the ribosome, an enormously complex machine whose structure is outlined in Figure 6. The ribosome is composed of a large and a small subunit, which in bacteria contain a total of three ribosomal RNAs and more than 50 different proteins (Lake, 1985). More than 60 percent of its mass is RNA, which was originally thought to play a structural role, helping to position the ribosomal proteins. However, attempts to find specific proteins that catalyze peptide bond formation failed, and it is now widely believed that the ribosomal RNA is itself the major catalyst (Noller, 1984; Moore, 1986). This view is supported by the evolutionary conservation of the RNA components of the ribosome: the structures of the ribosomal RNA molecules appear to be highly conserved, being similar in organisms as diverse as bacteria and humans (Noller, 1984; Gutell *et al.*, 1985).

An outline of the structure of the ribosomal RNA molecule in the small subunit of bacterial ribosomes is presented in Figure 7. This representation is only two-dimensional, being analogous to the view of the tRNA molecule shown earlier in Figure 4A. We do not yet know how this molecule folds up into its compact three-dimensional form, but with a size twenty times larger than a tRNA molecule, the possibilities for creation of a complex and interesting surface are certainly impressive.

I would like to stress two facts about the ribosome. First, its catalysis of protein synthesis appears to be an RNA-based process, as expected from the pathway by which the process of protein synthesis evolved (Fig. 5). Second, the mechanism of protein synthesis seems complex and awkward compared to other biological processes that

evolved later and were therefore based on protein catalysts (this point will be explored later, when we discuss DNA replication as a class II reaction). Some important, although tentative, conclusions can be derived from these two observations. It seems, first of all, that RNA-based catalyses are considerably less powerful than protein-based catalyses. As a consequence, it takes much more molecular mass to carry out a reaction catalyzed by RNA molecules than to carry out the same reaction catalyzed by proteins. In terms of a familiar analogy, the early cells that used only RNA catalysis were like a computer based on vacuum tube technology: very slow for their size. This is presumably why those cells that developed protein synthesis proliferated at the expense of their neighbors, and came to dominate the earth to such an extent that no cells lacking proteins have survived.

If they are less efficient than protein catalysts, why do any RNA catalysts still exist in cells? The suggestion is that cells, unlike those of us who have recently purchased computers, have been unable to escape the past. Thus, while a "microchip solution" to the synthesis of proteins would presumably be more efficient for the cell, the old mechanism clearly works well enough in its present patched-together form (in which ribosomal proteins have been added on as appendages to help the ribosomal RNAs) to be retained. In other words, cells—unlike computers—are not optimally designed. Instead what they are today is in large part a reflection of their past history (Jacob, 1977). The ribosome is a notable example. As a machine for making proteins, the ribosome seems so awkward as to be a bore both for teachers to teach and for students to learn. Its many pieces seem to make no conceptual sense at all, especially when compared to the elegantly-designed pieces of a DNA replication machine (see below). Only when viewed as a historical relic does the ribosome come alive. Now it suddenly turns into a fascinating object that can help us to understand the pathway by which protein synthesis evolved, and even how early cells

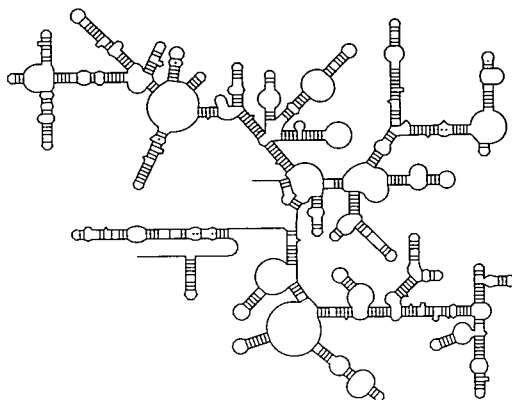


FIG. 7. The structure of 16S rRNA, the RNA molecule in the small subunit of bacterial ribosomes. Only a two-dimensional view of this molecule, which contains 1,542 nucleotides, is shown; how the molecule is folded in three dimensions is unknown (courtesy of H. F. Noller; from Gutell *et al.*, 1985).

might have worked before there were proteins.

ALL DNA FUNCTIONS EVOLVED IN AN ENVIRONMENT RICH IN PROTEIN CATALYSTS

If history still shows in a cell as claimed, then knowing the path by which the cell evolved should be useful for understanding many cellular processes. In present-day cells, a great deal of action centers on DNA and the manner in which gene expression is controlled. Yet DNA itself is generally considered to have arisen only at a relatively late step in the evolution of the cell. As outlined in Figure 8, the postulated RNA-based cells were succeeded by cells in which proteins carried out more and more of cellular catalysis. As these cells became more complex and sophisticated, there would have been pressure to evolve specialized RNA molecules that stored the cell's genetic information in an RNA double helix (Strickberger, 1986). In this form, each nucleotide sequence would be stored in duplicate, and RNA repair mechanisms (analogous to present-day DNA repair mechanisms) could operate to stabilize the genetic information against the inevitable random damage that is inflicted by chemical decay. Only in this way could the many

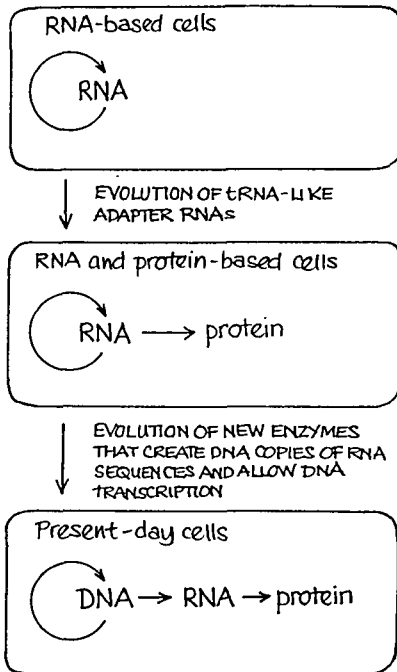


FIG. 8. Three postulated steps in the evolution of cells, culminating in the central dogma. Because DNA is a relatively late addition to the cell, it is likely that many enzymes evolved to an advanced state in cells that lacked DNA.

sequences of nucleotides required to specify complex cells be stably maintained.

RNA differs from DNA in having an extra 2'OH group on each of its sugars. The extra hydroxyl groups on RNA molecules are no doubt very important for imparting specific catalytic activities to its various folded forms. This difference is sufficient to explain why RNA and not DNA was the nucleic acid that formed the basis for the first cells. But the same chemical reactivity that is useful for catalysis is harmful in a molecule designed to store genetic information because it increases the rate of spontaneous chemical decay processes. It is presumably for this reason that double-stranded DNA evolved to take the place of RNA as the nucleic acid that stores the genetic information in present-day cells. The replacement process must have occurred gradually, with the evolution of intermediate cells containing both RNA and DNA information stores. During this period, a large entourage of new enzymes

had to be developed to handle DNA—including RNA polymerase and a variety of gene regulatory proteins.

The pressure to switch to a more stable genome should have arisen after cells became reasonably complex, and carried so much nucleotide sequence information that the increased chemical stability of DNA was important. It therefore seems safe to assume that DNA became the main information store in cells only relatively late in cell evolution. Because such cells would have contained a large number of very efficient protein catalysts, the various DNA functions must have evolved in an intracellular environment where RNA catalysis had become largely obsolete. Accordingly, by our arguments one would expect all of the processes occurring on DNA to be carried out by mechanisms that are much more efficient than those found for protein synthesis, with protein catalysis being used exclusively (*i.e.*, no RNA catalysis).

DNA REPLICATION DEMONSTRATES THE POWER OF PROTEIN CATALYSIS

DNA replication is an example of a process that occurs on DNA by a mechanism that is well understood (Kornberg, 1980). For this reason, it serves as a good model of a late-evolving catalytic mechanism and provides a useful comparison to protein synthesis. The action takes place at a structure called a "replication fork"; here the parental DNA double helix is opened into its two separate strands, so that each old strand can serve as a template for the formation of a new strand. As a result, two daughter DNA double helices are formed, each with one old strand and one new strand (the so-called semiconservative mode of DNA synthesis that was predicted by Watson and Crick [1953]).

A number of proteins with discrete functions are involved in moving a replication fork, and these cooperate to form a multienzyme "protein machine" that synthesizes DNA (Alberts, 1985). The replication fork with its bound proteins is displayed in a two-dimensional representation in Figure 9. Because the two strands of the DNA double helix are oriented in opposite direc-

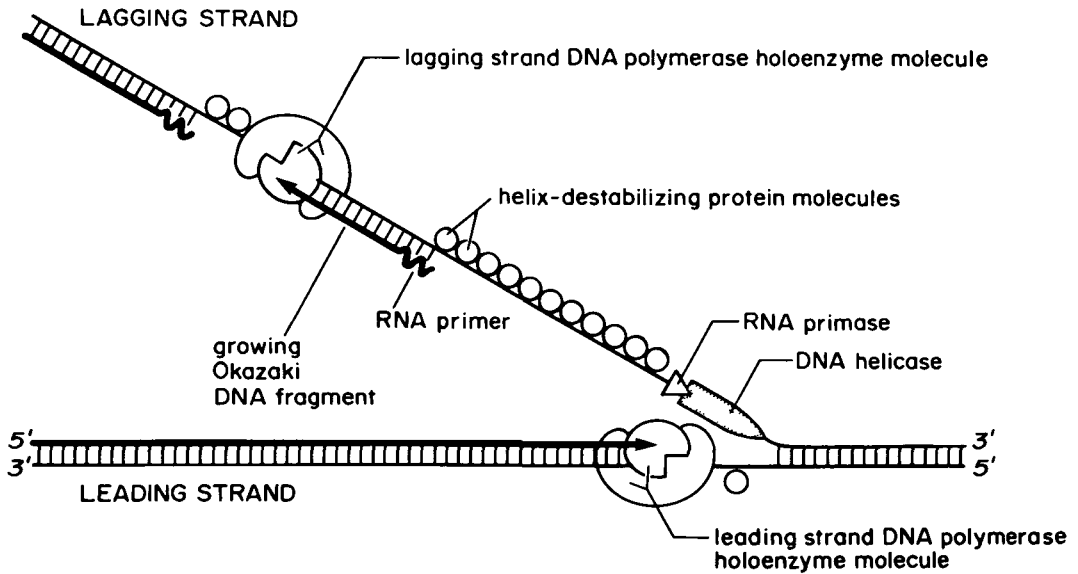


FIG. 9. A two-dimensional view of the DNA replication fork, showing the major proteins present (see text for details).

tions, while the polymerization of nucleotides occurs only in one direction, the replication fork is asymmetric—with “leading” and “lagging” strands. The actual synthesis of DNA is catalyzed by a complex of several proteins called a DNA polymerase “holoenzyme,” with a separate polymerase molecule on each strand. The DNA polymerase on the leading strand moves continuously, whereas the other DNA polymerase molecule works discontinuously, being forced by the orientation of its template strand to synthesize the lagging strand as a series of short Okazaki fragments. Behind the replication fork, these fragments are stitched together by a DNA repair process, creating a continuous daughter DNA strand.

The cell was faced with several technical problems in the design of the replication mechanism, and these were solved by the evolution of a variety of cooperating replication proteins. First of all, the DNA helix ahead of the replication fork had to be opened at a rapid rate, exposing the DNA bases on the template in a single-stranded form. This problem is solved by a DNA helicase molecule that uses ATP hydrolysis energy to propel itself rapidly along a DNA

single strand (Fig. 10). This enzyme runs along the lagging strand at the fork, pushing open the helix ahead of it as it goes. Helix-destabilizing protein molecules help the helicase by binding in clusters to the newly-opened DNA strands; these protein molecules manage to bind tightly to a DNA single strand while leaving the DNA bases on the strand freely available for base pairing (see Fig. 9, above).

A second problem is that all DNA polymerases require a 3'OH end on which to polymerize nucleotides; this pre-existing “primer chain” must be base-paired to the template strand to be copied (Kornberg, 1980). Thus, the synthesis of every Okazaki fragment must be started by a separate oligonucleotide primer. The primer used for this purpose is a short RNA molecule, which is synthesized by a separate enzyme called an *RNA primase*; this primase enzyme is kept at the proper position on the lagging strand by virtue of its attachment to the moving DNA helicase (see Fig. 9, above).

The replication fork in three-dimensions is even more impressive. As shown in Figure 11, the DNA on the lagging strand of the fork is apparently folded back on

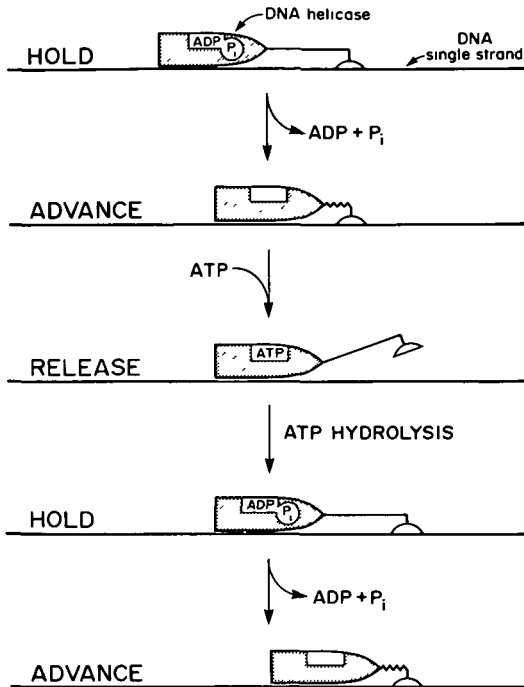


FIG. 10. Schematic diagram showing how a DNA helicase is moved along a DNA strand by allosteric changes in its conformation caused by the hydrolysis and release of ATP molecules. The fact that one of the conformational transitions is directly coupled to ATP hydrolysis makes this cycle of three shape changes unidirectional. As a result, the helicase moves consistently in a single direction along the DNA, as shown. The rate of helicase movement in bacterial cells is about 500 nucleotides per second.

itself to allow the lagging strand DNA polymerase molecule to form a complex with the leading strand DNA polymerase molecule. This linkage makes it possible for the DNA polymerase molecule on the lagging strand to be used over and over for successive rounds of Okazaki fragment synthesis, as schematically illustrated in Figure 12.

The entire set of replication proteins acts like a tiny sewing machine, powered by nucleoside triphosphate hydrolyses that move the individual protein parts relative to each other without disassembly of the complex. The mechanism is incredibly efficient: in bacteria, the replication fork moves at a rate of about 500 nucleotides per second, and the fidelity of templating is such that less than one nucleotide incorporation

error is made per every 10^8 base pairs replicated. Compare this performance with that of a ribosome, where proteins are synthesized at a rate of only 20 amino acids per second with about one error in every 10^4 amino acids polymerized. Yet the ribosome has a total mass more than three times greater than that of the replication apparatus. In my opinion, this state of affairs does not make sense unless we view the ribosome as a historical relic that evolved during an age where cells were capable only of "vacuum tube technology."

THE MECHANISMS USED FOR NUCLEAR PRE-mRNA SPLICING SUGGEST THAT THIS REACTION IS OF ANCIENT ORIGIN

Eucaryotic cells evolved from procarcotic cells, and yet only eucaryotes make extensive use of an externally-catalyzed form of RNA splicing that removes internal sequences (intron sequences) from their primary RNA transcripts (Chambon, 1981). Therefore, either bacteria have lost an old mechanism (Gilbert, 1978), or nuclear pre-mRNA splicing evolved relatively late in cell evolution. Recent biochemical studies probing the mechanism of the eucaryotic type of RNA splicing have revealed that the reaction is carried out by a large complex of ribonucleoprotein particles, whose total size approaches the size of a ribosome (Brody and Abelson, 1985; Grabowski *et al.*, 1985). According to our arguments, this reaction should therefore be an ancient one that first evolved in RNA-based cells and was present in the ancestors of all bacteria.

EUCARYOTIC GENE EXPRESSION, LIKE DNA REPLICATION, IS LIKELY TO BE MEDIATED BY PROTEIN MACHINES

The control of eucaryotic gene expression is currently an area of intense research. These controls program the cells in a multicellular organism to become different according to their position during embryonic development, as required to produce a complex organism. Whether a particular gene is expressed or not can probably be regulated at any one of the many different steps required to translate a DNA sequence into a protein sequence. However, the pre-

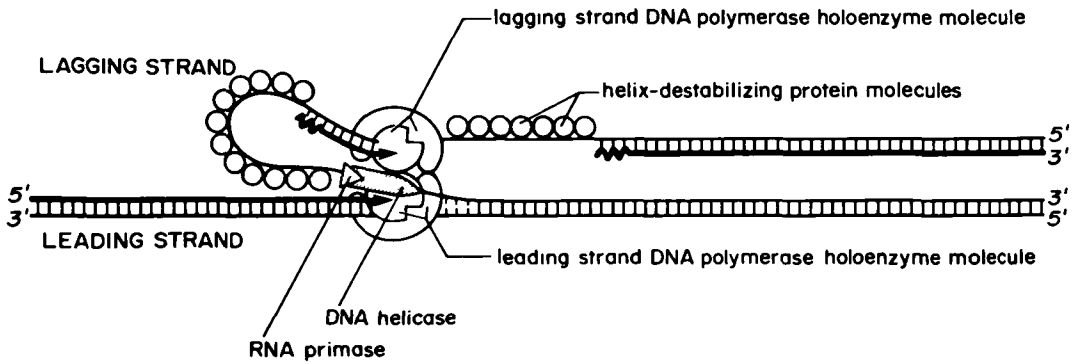


FIG. 11. The proteins of a DNA replication machine as they are thought to exist in an actual replication fork. The two-dimensional replication fork in Figure 9 has been converted into the structure shown by folding the DNA on the lagging strand in such a way as to bring the DNA polymerase on the lagging strand close to the DNA polymerase molecule on the leading strand. The lagging strand DNA polymerase molecule is thereby held to the rest of the replication proteins, allowing it to be retained for many successive cycles of Okazaki fragment synthesis, as shown in Figure 12.

dominant level of control is exerted at the first step of gene expression, when the decision is made to transcribe a given region of the DNA into RNA (Darnell, 1985). This level of control also predominates in bacteria, where many of the mechanisms involved have been worked out in great detail (Lewin, 1985). Here proteins bind to specific DNA sequences about 10 to 20 nucleotides long to control the expression of genes. Such proteins are called *gene regulatory proteins*, and they can either activate or repress the process of transcribing the adjacent region of DNA into an RNA sequence. Repression is an especially simple process. RNA polymerase, the enzyme that synthesizes RNA from a DNA template, binds to a specific DNA sequence called a *promoter* to initiate the process of DNA transcription. Gene regulatory proteins that work as repressors bind to a DNA sequence that overlaps the promoter sequence, preventing RNA polymerase from binding at that site and thereby blocking the transcription of the adjacent gene. Gene activation is only slightly more complicated. In this case, a gene is normally turned off because its promoter sequence is an altered one that RNA polymerase by itself is unable to use efficiently. However, this defective promoter becomes a good promoter for the RNA polymerase when a gene regulatory protein binds to a specific DNA sequence just upstream. Dur-

ing such gene activation, the gene regulatory protein is believed to touch the RNA polymerase at the adjacent promoter site in a way that helps this enzyme to begin its RNA synthesis (Ptashne *et al.*, 1980).

In bacteria, both gene activation and gene repression usually involve the binding of gene regulatory proteins at or very near to the promoter, which in turn contains the start site for RNA synthesis. There are some theoretical reasons why one might expect a difference between bacterial and eucaryotic gene control mechanisms in this regard. A gene in a complex multicellular organism appears to turn on or off in response to the sum of many different inputs, and we now know that the controlling mechanism interprets the cumulative effect of many gene regulatory proteins acting simultaneously on each gene (Yamamoto, 1985). This type of *combinatorial* control is advantageous for the cell, because it allows a large number of genes to be regulated by relatively few gene regulatory proteins (Gierer, 1974). It is not clear how such multifactorial combinatorial control would be accomplished by the simple type of mechanism just described, where only a small region on the DNA is allotted for regulating a gene's activity.

Fortunately for the theorists, when investigators started dissecting the mechanisms of eucaryotic gene regulation through the use of recombinant DNA methods, a some-

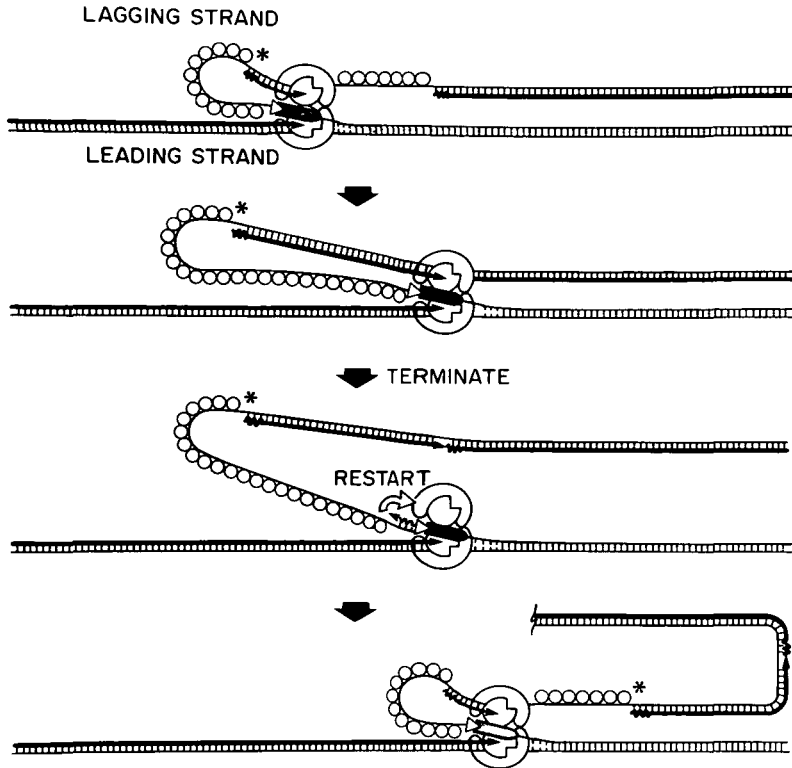


FIG. 12. A model for the movement of the replication fork shown in Figure 11. The crucial step in the cycle shown occurs when the lagging strand DNA polymerase molecule finishes the synthesis of each Okazaki fragment; in this "termination" step, the DNA on the lagging strand is released, freeing the polymerase to start the next Okazaki fragment. As indicated, the termination event also appears to trigger the synthesis of an RNA primer by the adjacent molecule of RNA primase (from Alberts, 1985).

what different form of gene regulation was discovered. Although gene regulatory proteins that bind to specific DNA sequences are again involved, one class of these can control the transcription of eucaryotic genes by binding to a site that is hundreds to thousands of DNA base pairs away from the promoter sequence at which RNA polymerase starts. Moreover, the sites that bind these proteins can be experimentally relocated at many different positions relative to the promoter without losing their effect. While most such sites are located upstream from the transcription start site, sometimes they are found in the middle of a transcribed DNA sequence or even at the far end of a gene. The first such gene control sites found were called "enhancers" (Serfling *et al.*, 1985), because when they bound a gene regulatory protein the tran-

scription of a nearby gene increased (Fig. 13A). More recently, sites that seem to act in the opposite way—turning off genes from a distance—have been discovered (Brand *et al.*, 1985; Johnson and Herskowitz, 1985; Struhl, 1985); these sites have tentatively been called "silencers" (Fig. 13A).

As more and more information has been obtained concerning the regulatory regions near higher eucaryotic genes, it has become apparent that many eucaryotic genes have the general structure that is schematically illustrated in Figure 13B. Whether a gene is on or off depends on the sum of multiple inputs from a number of different gene regulatory proteins, some tending to turn the gene on and others tending to turn the gene off. These regulatory proteins bind to specific sites that can be scattered over

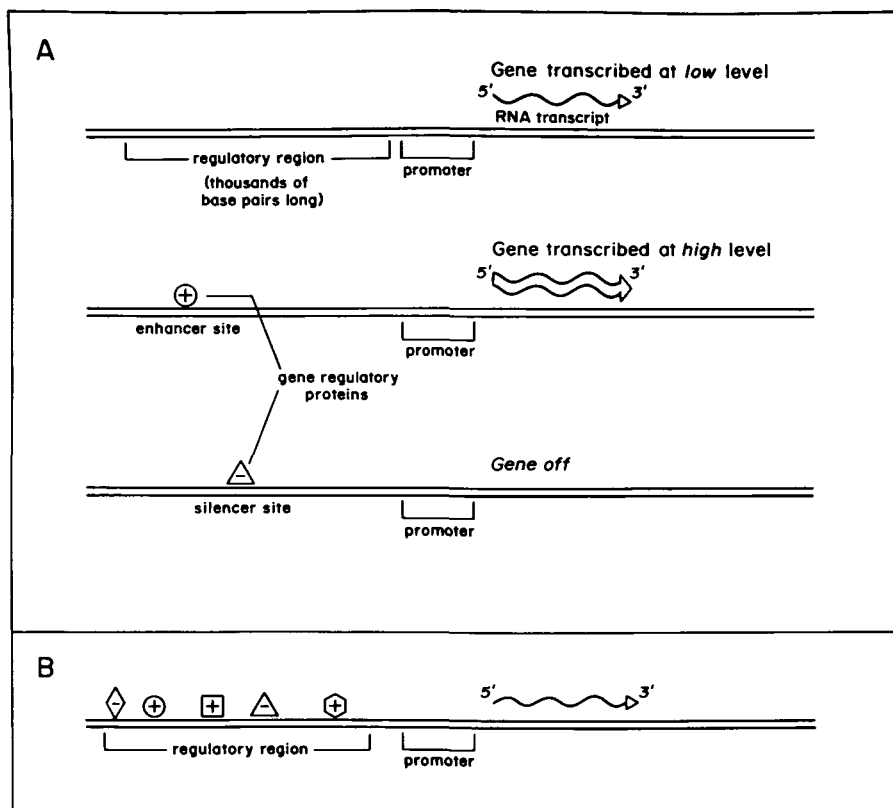


FIG. 13. Diagram illustrating the control of eucaryotic gene expression by gene regulatory proteins that bind to enhancer and silencer sites. The same gene is shown in three different states of activity in (A). In (B), a eucaryotic gene is shown that is regulated by the combined action of many different gene regulatory proteins. Such combinatorial gene regulation is common in eucaryotes.

a region of the DNA that is 5,000 or more base pairs long. By unknown means, the signals from all of these protein binding sites are integrated to control DNA transcription at the promoter.

Our knowledge of the effect of different enhancer and silencer sites on the patterns of gene expression in an organism is increasing especially rapidly in *Drosophila*, where mutants in a series of different gene regulatory proteins have been identified by geneticists, each of which has a major effect on patterns of gene expression in the early embryo (Lewis, 1978; Nusslein-Volhard and Wieschaus, 1980; Mahowald and Hardy, 1985). Several such regulatory genes have now been characterized and shown to be expressed in highly-specific spatial patterns across the early 6,000-cell blastoderm embryo (for example, see

Hafen *et al.*, 1984; Kornberg *et al.*, 1985). These spatial patterns of gene expression change when other putative gene regulatory proteins are mutated in the embryo, and it is thought that many of the gene regulatory proteins defined by these mutants affect each others' syntheses (Desplan *et al.*, 1985; Struhl and White, 1985). Specific alterations in patterns of gene expression are also observed when different short sections of DNA, each apparently containing one or more enhancer or silencer sites, are removed from the regulatory region next to a gene (Hiromi *et al.*, 1985).

As yet, biochemical studies of the mechanism by which enhancers and silencers act to control eucaryotic gene expression have lagged far behind the physiological studies that have defined these elements. Thus,

how the gene regulatory proteins act from a distance and how effects from different sites are combined and interpreted remain unknown. Another important unknown is how "cell memory" works—*i.e.*, how cells in eucaryotes can inherit a gene that remains either on or off in a cell clone (Brown, 1984).

In the context of this article, it is worth remembering that all these mechanisms of eucaryotic gene expression must have evolved quite late in the evolutionary history of the cell. Thus, while the mechanisms are not known, we can expect them to resemble DNA replication in being catalyzed by elegant and efficient protein machines. This means, first of all, that one expects proteins and not ribonucleoprotein particles to be involved. Second, one should not be surprised to find systems of interacting proteins that move relative to each other without dissociating, some of which hydrolyze nucleoside triphosphates to create ordered conformational changes in either the DNA or themselves. It should be a fascinating story, with many of the details appearing in the next ten years. Be sure to keep tuned!

CONCLUSIONS AND CAVEATS

In this report I have suggested that cells contain two broad classes of catalytic mechanisms. Class I mechanisms are processes that are carried out by large ribonucleoprotein complexes and appear to involve RNA catalysis. Two examples are protein synthesis and RNA splicing. These mechanisms seem complicated and unwieldy, and are best explained as historical relics of processes that arose early in the evolution of the cell, when only relatively inefficient catalysts were available. Class II processes do not involve RNA catalysis; instead they are likely to be carried out by multienzyme "protein machines." An example of such a mechanism is DNA replication. These mechanisms seem elegant and efficient when compared to class I mechanisms. They are likely to have evolved later, during a period when cells contained a large repertoire of very efficient protein catalysts. This binary characterization of catalytic reactions is no doubt oversimplified,

but it seems useful in two respects. First of all, as witnessed by our discussion of nuclear pre-mRNA splicing, this view can allow the origin of a biological process to be positioned with respect to the evolutionary history of the cell, once its mechanism is known. Second, it predicts that the mechanisms used for all DNA-mediated processes will resemble those found for DNA replication much more than those found for protein synthesis. Such a prediction is of practical value, since most of these mechanisms are not yet well-understood, and the optimal experimental approach to deciphering them can depend on the type of catalysis involved (for example, see Alberts, 1984).

The above analysis incorporates several unproven assumptions, some of which may not be obvious to most readers. First of all, previous discussions of the origin of protein synthesis have usually assumed that this process originated in a relatively primitive, pre-cellular "soup" of macromolecules, and therefore that there were no cells before there were proteins (*e.g.*, see Eigen *et al.*, 1981). In my opinion, the recent discovery of efficient catalysis by RNA molecules makes it unreasonable to insist that cells could not exist without proteins. In turn, the possibility of RNA-based cells allows protein synthesis to evolve in the presence of a sophisticated family of RNA catalysts, making the spontaneous origin of genetically-specified proteins on the earth seem much more feasible from a chemical standpoint.

There is a second important assumption that needs to be exposed. The two processes that I have termed class I reactions—protein synthesis and nuclear pre-mRNA splicing—both involve the accurate recognition of nucleotide sequences in single-stranded RNA molecules. One could therefore explain the observed role of ribonucleoprotein complexes in these reactions simply by postulating that proteins by themselves have a great deal of difficulty in recognizing specific sequences in a folded nucleic acid chain, and that such recognitions are best accomplished by a second polynucleotide sequence (Cech and Bass, 1986). In such a

view, RNA catalysts are present in cells not because they are historical relics, but because they can accomplish certain types of catalyses better than protein molecules alone. It is hard for me to accept this hypothesis, because specific nucleotide sequences in single-stranded DNA chains are known to be recognized very rapidly and efficiently by proteins during DNA replication (*e.g.*, see Eisenberg *et al.*, 1977; Schlomai and Kornberg, 1980; Tabor and Richardson, 1981; Cha and Alberts, 1986). However, we can never be certain that it is possible to design an RNA-free protein machine that would improve upon protein synthesis—unless of course someone discovers a novel bacterium that can divide every two minutes by making proteins without ribosomes!

In this article, I have consistently stressed the relative inefficiency of RNA catalysts in comparison to protein catalysts. However, it is important to realize that the RNA molecules in the first cells must have been considerably more diverse and sophisticated as catalysts than the very limited set of catalytic RNA molecules thus far known. While most of these early catalysts would have undergone extinction when more efficient protein catalysts evolved to supplant them, one would expect others to have survived. Detailed studies of various catalyses carried out by ribonucleoprotein particles are now underway; hopefully, some of the results will reveal much more about the fascinating RNA catalysts that are suspected to have made life possible before proteins.

ACKNOWLEDGMENTS

I would like to thank Keith Roberts for the design of Figures 3, 5, and 8, Gert Weil for drawing the final figures, Kathleen Rañeses for preparation of the manuscript, and Beth Alberts for many helpful comments on the writing.

REFERENCES

- Alberts, B. M. 1984. The DNA enzymology of protein machines. Cold Spring Harbor Symp. Quant. Biol. 49:1–12.
- Alberts, B. M. 1985. Protein machines mediate the basic genetic processes. Trends in Genetics 1:26–30.
- Alberts, B. M., D. Bray, M. Raff, K. Roberts, J. Lewis, and J. D. Watson. 1983. *Molecular biology of the cell*. Garland Publishing, Inc., New York.
- Bass, B. L. and T. R. Cech. 1984. Specific interaction between the self-splicing RNA of *Tetrahymena* and its guanosine substrate: Implications for biological catalysis by RNA. Nature 308:820–826.
- Belfort, M., J. Pedersen-Lane, D. West, K. Ehrenman, G. Maley, F. Chu, and F. Maley. 1985. Processing of the intron-containing thymidylate synthase (*td*) gene of phage T4 is at the RNA level. Cell 41:375–382.
- Brand, A. H., L. Breeden, J. Abraham, R. Sternglanz, and K. Nasmyth. 1985. Characterization of a “silencer” in yeast: A DNA sequence with properties opposite to those of a transcriptional enhancer. Cell 41:41–48.
- Brody, E. and J. Abelson. 1985. The “spliceosome”: Yeast pre-messenger RNA associates with a 40S complex in a splicing-dependent reaction. Science 228:963–967.
- Brown, D. D. 1984. The role of stable complexes that repress and activate eucaryotic genes. Cell 37:359–365.
- Cech, T. R. 1985. Self-splicing RNA: Implications for evolution. Int. Rev. Cytol. 93:3–22.
- Cech, T. R. and B. L. Bass. 1986. Biological catalysis by RNA. Ann. Rev. Biochem. 55. (In press)
- Cech, T. R., N. K. Tanner, I. Tinoco, Jr., B. R. Weir, M. Zuker, and P. S. Perlman. 1983. Secondary structure of the *Tetrahymena* ribosomal RNA intervening sequence: Structural homology with fungal mitochondrial intervening sequences. Proc. Natl. Acad. Sci. U.S.A. 80:3903–3907.
- Cech, T. R., A. J. Zaig, and P. J. Grabowski. 1981. In vitro splicing of the ribosomal RNA precursor of *Tetrahymena*: Involvement of a guanosine nucleotide in the excision of the intervening sequence. Cell 27:487–496.
- Cha, T.-A. and B. M. Alberts. 1986. Studies of the DNA helicase-RNA primase unit from bacteriophage T4: A trinucleotide sequence on the DNA template starts RNA primer synthesis. J. Biol. Chem. 261:7001–7010.
- Chambon, P. 1981. Split genes. Sci. Amer. 244(5):60–71.
- Crick, F. H. C. 1968. The origin of the genetic code. J. Mol. Biol. 38:367–379.
- Crick, F. H. C., S. Brenner, A. Klug, and G. Pieczenik. 1976. A speculation on the origin of protein synthesis. Origins of Life 7:389–397.
- Crothers, D. M. 1982. Nucleic acid aggregation geometry and the possible evolutionary origin of ribosomes and the genetic code. J. Mol. Biol. 162:379–391.
- Darnell, J. E. 1985. RNA. Sci. Amer. 253(4):68–78.
- Desplan, C., J. Theis, and P. H. O’Farrell. 1985. The *Drosophila* developmental gene, *engrailed*, encodes a sequence-specific DNA binding activity. Nature 318:630–635.
- Eigen, M., W. Gardiner, P. Schuster, and R. Winkler-Oswatitsch. 1981. The origin of genetic information. Sci. Amer. 244(4):88–118.
- Eisenberg, S., J. Griffith, and A. Kornberg. 1977.

