

kon

June 7, 1957

HOW DO AMINO ACIDS READ THE CODE?

by Leo Szilard

(Submitted by Joseph E. Mayer)

The Enrico Fermi Institute for Nuclear Studies  
The University of Chicago, Chicago, Illinois

It is generally believed that proteins are formed alongside of nucleic acid templates. The sequence of purine-pyrimidine bases in the template is supposed to represent a code that may somehow determine the sequence of the amino acids in the particular polypeptide (protein) that a given template will form. The purine and pyrimidine bases of the template, the letters of the code, are adenine, uracil, guanine and cytosine if the template be an RNA molecule; and if the template be a DNA molecule, thymine takes the place of uracil.

Because a template which synthesizes protein must carry the same information as the corresponding gene but need not necessarily be the gene itself, we shall here refer to such a template, for the sake of brevity, as a paragene.

It has remained so far a complete mystery in just what conceivable way amino acids could read such a code. In what manner can chemical forces of the kind we know to exist -- line up amino acids alongside such a template in the proper sequence and at the proper distance from each other, so that there might be initiated a chemical reaction chain through which adjacent amino acids might ~~form~~ be linked through peptide bonds with each other?

It is the purpose of the present paper to indicate a conceptually simple scheme that will -- at least by way of an example -- illustrate in what manner this might ~~be accomplished~~ take place in the living cell.

The basic thought underlying this scheme consists in the assumption that there are a number of enzymes (or enzyme systems) -- perhaps twenty altogether -- in the cell, and that each of these catalyzes the formation of a particular trinucleotide which carries either one particular amino acid or, more likely perhaps, a particular sequence of three amino acids. Each amino acid may be tied to a nucleotide through a high energy bond, either a P or PP bond, and ~~it is assumed that~~ such acid anhydrides may release an energy of 12,000 calories or 16,000 calories respectively, when they split themselves

According to the notions here presented, amino acids can not read/ the code of the parogene. ~~strat~~. But the purine and pyrimidine bases of the trinucleotides, which carry the proper amino acids, may attach through the formation of 6 hydrogen bonds to the proper sequence of three purine or pyrimidine bases on the parogene, and thus the amino acids may be lined up in the proper sequence along the parogene. Each amino acid being present in the form of an acid anhydride carries with it the energy that may be released in a chemical reaction chain that links adjacent amino acids through peptide bonds to each other.

Accordingly sequences of three nucleotides along the parogene represent the code words, and the trinucleotides which carry the amino acids would represent the anti-code words. We assume that these anti-code words are complementary to the code words in the sense that where the code word contains adenine the anti-code word contains uracil (or thymine), where the code word contains uracil (or thymine) the anti-code word contains adenine and similarly guanine corresponds to cytosine and cytosine corresponds to guanine. The rationale for this assumption is as follows:

The concept of code-letter and complementary code-letter arose originally from the interpretation of the structure of DNA given by J. D.

Watson and F. H.C. Crick.<sup>(1)</sup> They showed that in a double stranded DNA structure, adenine pairs with thymine (which presumably plays the same role in DNA as does uracil in RNA) and guanine pairs with cytosine. The helical structure of DNA permits just such pairing, and hydrogen bonding is possible between adenine and thymine as well as between guanine and cytosine.

We may now tentatively adopt the view that during protein synthesis the parogene, whether it be a single DNA or a single RNA strand, assumes a somewhat similar helical configuration. The amino acids carried by the proper trinucleotides (the anti-code words) may then be lined up in the proper sequence along the parogene through the formation of hydrogen bonds between the purine and pyrimidine bases of the trinucleotides and the complementary bases on the parogene. When the trinucleotides are lined up in the proper order, then -- since each trinucleotide carries the proper amino acids -- the amino acids are also lined up in the proper order.

We shall now single out for more detailed examination one conceivable model for protein synthesis which might account for the lining up of the amino acids alongside the parogene, both in the proper order and at the proper distance from each other. This particular model is based on the following assumptions:

The trinucleotides which form the anti-code words contain the sugar ribose rather than the sugar desoxyribose. Each particular ribose trinucleotide (the anti-code word) carries a particular sequence of three amino acids. A phosphate (or diphosphate) group is attached through an oxygen atom to the (2) carbon atom of the ribose moiety of each nucleotide (ester linkage) and one amino acid is attached to each of these phosphate (or diphosphate) groups. The amino acid anhydrides represent an energy-rich P (or PP)

~~bond, which is broken when the amino acid is released from the trinucleotide.~~

(1) Proc. Roy. Soc., Vol. 223, p. 80, 1954.

~~XX~~

During protein synthesis the nucleic acid strand that functions as a template (the paragene) may take up -- so we here assume -- a helical configuration resembling the helical configuration of a DNA strand in the double strand DNA helix. The ribose trinucleotides may then line up alongside the helical paragene with their purine and pyrimidine bases paired with the complementary bases of the paragene, and if they are so lined up, then the amino acids carried by the trinucleotides may come to lie at just about the right distance from each other to permit the formation of a peptide bond between adjacent amino acids. A chemical reaction chain -- starting from the head of a paragene -- may then move down along the paragene, split the acid anhydrides, and thus free the amino acids as well as ~~supply~~ <sup>made available</sup> the energy needed for the formation of peptide bonds between adjacent amino acids.

Clearly, this model imposes certain restrictions on the possible amino acid sequences that paragenes can produce. What are these restrictions?

If we have four letters to choose from, then -- in a code which utilizes three-letter words -- we can construct 64 different words. We might, however, be restricted to the use of 20 out of the 64 words that are available. The reasons for this restriction would be as follows:

If all 64 possible three-letter combinations form in fact a code word, and if the nucleic acid strand assumes at the time of the formation of the polypeptide the helical configuration mentioned above, then it follows that the code on the paragene must be read consecutively from one end -- say, from the "head" of the paragene downward. This is so because such a helical structure does not provide for commas between the individual code words, and in a 64-word, three-letter-word, code any three consecutive letters form a word. The letters 1, 2, 3 form a code word which was meant

to be conveyed and so do the letters 4, 5, 6, but sequences of three letters which encroach on two adjacent words (such as 2, 3, 4 or 3, 4, 5, for example) form code words which are not meant to be conveyed.

In these circumstances, the code would be misread if the trinucleotides, which represent the anti-code words, assemble alongside the paragenes simultaneously, rather than -- from one end/<sup>on</sup>-- consecutively. If we want to have simultaneous assembly of the trinucleotides alongside the comma-less paragenes, then we are restricted to 20 code words.

The notion of such a 20-word code, which needs no commas, was introduced by F.H.C. Crick, J.S. Griffith, and L.E. Orgel of the Medical Research Council Unit at the Cavendish Laboratory, Cambridge, in a memorandum circulated in May, 1956 among workers interested in the subject of protein synthesis. From such a code we must demand that ~~which~~ the letters 1, 2, 3 on the template form a code word, and the letters 4, 5, 6 also form<sup>a</sup> code words, but /sequences of three letters, which encroach on two adjacent words (such as 2, 3, 4 or 3, 4, 5, for example) form no code word. Crick and his co-workers have shown that this demand can be met, that a code which requires no commas may be constructed and that it can accommodate 20 three-letter code words.

Applying the concept of a 20-letter code that requires no commas to our particular model of protein synthesis, we may now say the following:

Each of the 20 amino acids may appear once attached to the first letter and once attached to the last letter of the 20 (trinucleotide) anti-code words. Therefore, among the polypeptides that can be formed, each amino acid may precede any other amino acid, and each amino acid may follow any other amino acid. This does not, of course, mean that any amino acid sequence is possible.

Some of the amino acid sequences that may be found experimentally in sequential analysis of proteins and polypeptides might show that the restrictions imposed by our model on the possible amino acid sequence are too severe, and that the model has to be modified to accommodate established facts. As will be presently seen some modification of our model may be indicated for other reasons also. There is certainly no inherent reason why we should have a pure three-letter code -- as we have assumed above -- and why, for instance, a certain number of four-letter code words should not be utilized also.

If we had a pure three-letter code, we would have to demand that the number of amino acid residues of all polypeptides or proteins synthesized in the manner described above should be a multiple of three.

So far the number of amino acid residues in polypeptides and proteins have been determined only within one rather special class; all of them represent substances which are secreted by mammalian tissues. The number of amino acid residues found in such polypeptides and proteins, which have been analyzed with adequate accuracy, are as follows:

a) Insulin chain A: 21; insulin chain B: 30; corticotropin B: 39; oxytocin: 9; vasopressin: 9; Intermedin B: 18.

All of these would fit a pure three-letter code.

b) Intermedin A: 13; glucagon: 29 and pancreas ribonuclease: 124.

These do not fit a pure three-letter code. Intermedin and ribonuclease would have to include at least one 4-letter word and glucagon at least two 4-letter words. In a mixed system of three- and four-letter words, one can obviously never draw the conclusion that more than two four-letter words have been included.

Observed rate of enzyme synthesis

According to the notions here adopted most enzymes are synthesized in growing bacteria at a rather low rate which does not represent the maximum synthesizing capacity of the corresponding paragenes. The rate of production of a given enzyme may, however, be greatly enhanced if the enzyme is induced, and what we are interested to learn is the ~~maximum~~ <sup>maximal</sup> rate at which a paragene may be able to form the corresponding enzyme.

One of the most studied cases of enzyme induction is the induction of the enzyme  $\beta$ -galactosidase which splits lactose. Jacques Monod and his co-workers have shown that the production rate of this enzyme in bacteria can be greatly enhanced by certain chemical analogues of lactose, and that the rate of production of the enzyme goes up almost instantaneously upon adding the inducer to the medium. We are thus led to believe that the inducer may act by increasing the rate at which one template produces the enzyme rather than by increasing the number of templates that produce the enzyme at an unchanged rate.

The highest rate at which this enzyme is produced in fully induced wild type bacteria growing in minimal medium is about  $2 \cdot 10^{-18}$  grams per cell per second.

Assuming a molecular weight of 100,000 for this enzyme, we would obtain a rate of about 15 enzyme molecules per second per bacterium. Or if we assume a molecular weight of a million (Jacques Monod and ~~Max~~ <sup>Melvin</sup> Cohen estimate the molecular weight of this enzyme at about 800,000), we obtain a rate of 1.5 enzyme molecules per cell per second. The number of paragenes per cell is not known. If the paragene is a single strand of RNA, there might be a few paragenes present per cell rather than just one, and the number of paragenes might ~~in fact~~ be of the order of magnitude of 10.

In these circumstances we are led to believe, on the basis of the figures given above, that for enzyme ~~molecules~~ of high molecular weight - when the enzyme is fully induced and enzyme synthesis proceeds at its maximal rate -- the rate of formation of the enzyme is of the order of magnitude of one per second per paragene.

We shall now attempt to compute at what rate a parogene may be able to synthesize the corresponding enzyme on the basis of the model that we have postulated. For the purposes of this computation, we shall assume that the molecular weight of the enzyme is about 100,000. We then have about 1,000 amino acid residues in the enzyme, and accordingly we would have to assemble alongside the parogene  $m = 300$  trinucleotides, each of which is "loaded" with three amino acids.

In the approximation which we shall ~~use~~ <sup>use</sup> ~~consider~~, the minimum time,  $\tau_0$ , needed for the formation of the polypeptide is composed of two terms,  $\tau_1$  and  $\tau_2$ .

$$\tau_0 = \tau_1 + \tau_2$$

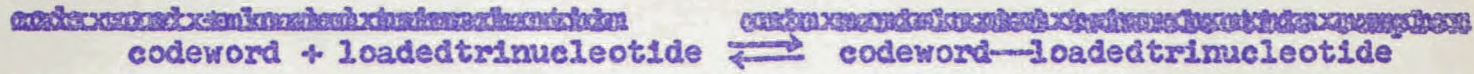
After all the amino acids, assembled alongside the template, have been linked into one polypeptide, and ~~assuming~~ <sup>if</sup> this polypeptide is at once removed, a certain time,  $\tau_1$ , will elapse until the trinucleotides, which are now denuded of amino acids, evaporate from the template and their place is taken by trinucleotides which are loaded with the proper amino acids. We ~~still~~ assume that the concentration of denuded trinucleotides in the cell is very small so that after the denuded trinucleotides evaporate, the loaded trinucleotides do not have to compete with denuded trinucleotides for their legitimate position along the parogene. The time,  $\tau_1$ , which is necessary to permit evaporation of  $m$  denuded trinucleotides and to assemble  $m$  loaded trinucleotides in their place we shall compute here on the assumption that once a loaded trinucleotide has found its position alongside the template, it will not evaporate again. Because this assumption is, ~~probably~~ not valid, we must make a correction which is represented by the second term,  $\tau_2$ .

Computation of  $\tau_2$

First we shall now compute this second term,  $\tau_2$ . This computation ~~will~~ <sup>is</sup> ~~be~~ based on the fact that (because of reevaporation of the loaded trinucleotides, which ~~are~~ reversibly combined with the anti-code words of the parogene) there will be - no matter how long we wait - always a certain number of code words (sequences of three nucleotides) on the parogene which are not "covered".



We have to deal with a reversible reaction which may be written as follows:



The rate at which the reaction proceeds from left to right; i.e. the number of successful hits per unit time is given by

$$\text{hit rate} = A \rho$$

$\rho$  denotes the concentration of the particular kind of loaded trinucleotide in mol/cc. We shall, for the sake of simplicity, assume that the concentration in the cell of each kind of loaded trinucleotide is the same.

A stands for

$$A = 6 \times 10^{23} v \sigma p$$

*Handwritten notes:*  $5 \times 10^3$  (pointing to  $6 \times 10^{23}$ ),  $\frac{1}{3} 10^{-15}$  (pointing to  $\sigma$ ),  $\$ 10 \frac{24}{10} \frac{3}{10} \frac{-17}{10}$  (pointing to  $p$ )

v is the molecular velocity,  $v = \sqrt{\frac{2RT}{\pi M}}$  where M is the molecular weight of the loaded trinucleotide. If  $M \approx 1000$ , we have  $v \approx 5 \times 10^3$  cm/sec.

$\sigma$  is the target area that must be hit by the loaded trinucleotide if Hydrogen bonding is to take place with the three adjacent nucleotides on the paragene. We may assume for  $\sigma$  a value of  $\sigma = 10^{-15}$  cm<sup>2</sup>.

p denotes the probability that the loaded trinucleotide, when hitting the code word, is in just the right geometrical position to permit hydrogen bonding to take place between the three complementary pairs of bases that are involved. We may take for p, as a very rough estimate  $p = 1/300$ ; This would give for  $\sigma p$  the value  $\sigma p = 1/3 \times 10^{-17}$ .

The rate at which the reaction proceeds from left to right is given by the rate  $\alpha$  at which the ~~codeword-loaded trinucleotide~~ complex dissociates or, as we may also say, the rate  $\alpha$  at which the loaded trinucleotide evaporates from the paragene. For this rate we may write

$$(2) \quad \alpha = \text{rate of evaporation} \approx 10^{13} e^{-\frac{\Delta H}{RT}}$$

where  $\Delta H$  is the binding energy for the loaded trinucleotide.

In equilibrium the hit rate and evaporation rate must be equal. The equilibrium constant,  $K$ , of this reversible reaction is defined as the concentration of the loaded trinucleotides at which the code word is covered half of the time. Thus we may write

$$(3) \quad \begin{aligned} AK &= \frac{1}{2} 10^{13} e^{-\frac{\Delta H}{RT}} = \frac{1}{2} \alpha \\ \alpha &= 2AK \end{aligned}$$

$$AK = \frac{1}{2} \alpha$$

In equilibrium the probability,  $f$ , for a given code word on the parogene not being covered by the proper loaded trinucleotide is given by

$$(4) \quad f = \frac{1}{1 + \frac{p}{K}}$$

Accordingly, in equilibrium, the total number of such gaps along the parogene which contains  $m$  nucleotides is given by

$$(5) \quad \text{"number of gaps"} = \frac{m}{1 + \frac{p}{K}}$$

We shall assume that most code words are "covered" in equilibrium and this means that

$$(6) \quad \frac{p}{K} \gg 1$$

We presume that after such equilibrium is established a chemical reaction chain is somehow triggered, and, moving down along the parogene, links adjacent amino acids into a polypeptide. The average time,  $\tau$ , needed for the formation of the polypeptide from the amino acids assembled along the parogene is given by the product of the "number of gaps", that

Dr Müller hat zwei Entwürfe  
er besmerkt aber die wesentliche der Erscheinung

Uns mein am Leben Wahrheit und Obedy

Man  
wenn man sich auf seine Nase  
nicht verlassen kann ja warum  
dann nicht man sich den man  
Verlassen?

er wollen wir hoffen

Sch. habe jedoch nicht die Absicht  
mich wegen der Zeichnung in Ernst zu  
dem ohne Erklärung gibt es keine  
Falschheit es gibt es keine Wahrheit

Wolfsche Stammes  
schädelstein Kristalline Form; ich habe es  
Was für eine Form haben  
ist nicht nicht die die Wahrheit in Shakespeare  
nascendi.

have to be filled consecutively, and the average time,  $\frac{1}{A\beta}$ , that it takes to fill one given gap. Thus, for  $\tau_2$  we may write

$$(7) \quad \tau_2 = \frac{1}{A\beta} \frac{m}{1 + \frac{\beta}{K}}$$

Computation of  $\tau_1$  and  $\tau_0 = \tau_1 + \tau_2$

When the polypeptide is formed and leaves the paragens, the code words are covered with the denuded trinucleotides. We may now compute the average time, , needed for the evaporation of all the denuded trinucleotides and the assembling of all the loaded trinucleotides in their place. We shall, for the sake of simplicity, assume that a denuded trinucleotide evaporates at the same rate,  $\alpha$ , as a loaded trinucleotide. The rate,  $\alpha$ , at which a loaded trinucleotide evaporates from the template is given by (3) and we may also write this in the form

$$(8) \quad \alpha = A\beta \frac{K}{P}$$

Since we have assumed  $\frac{P}{K} \gg 1$ , we have

$$(9) \quad A\beta \gg \alpha$$

As may be shown, ~~minimum time for~~ for very large values of  $m$ , ( $\ln m \gg 1$ ) we may write for

$$(10) \quad \tau_1 \approx \frac{1}{A\beta} \frac{P}{K} \ln m = \frac{1}{AK} \ln m$$

$\sim \frac{1}{\sqrt{2} A\beta} \frac{m \ln m}{\sqrt{2} m} + \frac{1}{\sqrt{2} A\beta} \frac{m \ln m}{\sqrt{2} m}$   
two terms equal for minimum  $\tau_0$

For the total time,  $\tau_0 = \tau_1 + \tau_2$  we thus obtain

$$(11) \quad \tau_0 = \frac{1}{A\beta} \left\{ \frac{m}{1 + \frac{\beta}{K}} + \frac{P}{K} \ln m \right\} \approx \frac{1}{A\beta} \frac{m}{1 + \frac{\beta}{K}} + \frac{1}{2AK} \ln m$$

If we wish to make this time as small as possible, we have to choose  $K$  so as to have for  $\frac{P}{K}$

$$(12) \quad \frac{P}{K} \approx \sqrt{\frac{2m}{\ln m}}$$

$\frac{d}{dK} \left( \frac{m}{1 + \frac{\beta}{K}} + \frac{1}{2AK} \ln m \right) = 0$

$m \approx \frac{2(1 + \frac{\beta}{K})^2}{2AK} = 1.7 \times 10^3$

$\tau_0 = \frac{2}{2AK} \frac{P}{K} \left\{ \frac{P}{K} = 12, K = \frac{P}{20}, \frac{P}{K} = 1.7 \times 10^3 \right\}$

Substituting this value into (11) gives

(13)  $\tau_0 \approx \frac{2 \times \cancel{A^2}}{AP} \sqrt{m \ln m}$        $\tau_0 = \frac{2\sqrt{3}}{2AK\frac{P}{K}} \sqrt{m \ln m}$

*see back of M.S.*

For a polypeptide containing 1,000 amino acid residues, i.e. a paragene containing about 300 code words, we may write  $m = 300$ , and thus we obtain from (12) and (13)

(14)  $\frac{P}{K} \approx 10$

and

(15)  $\tau_0 \approx \frac{50}{AP}$

*compute  $\tau_0$  for  $m = 1000$*

$$\tau_0 = \frac{\sqrt{2}}{AP} \sqrt{10^3 \times 7} = \frac{120 \times 1.4}{AP}$$

$$= \frac{170}{AP}$$

Estimates for the values of A and  $\sigma_p$

If we estimate  $\sigma_p = 1/3 \cdot 10^{-17}$ , then we have  $A = 10^{10}$ . For  $\rho = 5 \cdot 10^{-9}$  mol/cc ( $\rho = 5 \cdot 10^{-6}$  mol/liter) we then obtain  $\tau_0 \approx 1$  sec., which means that one polypeptide is formed per paragene per second.

It might well be that  $\sigma_p$  is ten times higher:  $\sigma_p = 1/3 \cdot 10^{-16}$  which gives  $A = 10^{11}$ . In this case  $\tau_0 \approx 1$  sec is obtained for  $\rho = 5 \cdot 10^{-10}$  mol/cc (or  $\rho = 5 \cdot 10^{-7}$  mol/liter).

If one were to attempt to find experimentally trinucleotides in bacterial cells, one would have to look for concentrations of this order of magnitude.

Since we have assumed  $\frac{P}{K} \approx 10$ , we must expect to find for K:

(16)  $5 \cdot 10^{-8} < K \text{ (in mol/liter)} < 5 \cdot 10^{-7}$

or

$5 \cdot 10^{-11} < K \text{ (in mol/cc)} < 5 \cdot 10^{-10}$

$\frac{170}{AP} = 1 \text{ sec}$

*for  $A \rho \approx 270$*

$A = 10^{11}$        $\rho = 2 \cdot 10^{-4} / \mu = 2 \cdot 10^{-6} / \text{liter}$

$A = 10^{10}$        $\rho = 2 \cdot 10^{-5} / \text{liter}$

From (3) we obtain for a given pair of values A and K the binding energy  $\Delta H$  between the loaded trinucleotide and the code word on the parogene. For a value of  $A = 10^{-10}$  and  $K = 10^{-7}$  mol/liter, or  $K = 10^{-10}$  mol/cc, we obtain  $\Delta H \approx 18,000$  calories. Since 6 hydrogen bonds are involved, this would mean about 3,000 calories per hydrogen bond.

For the same value of A and a value of K which is ten times larger  $\Delta H$  would decrease by about 1400 calories.

### Conclusion

These considerations show that the theory which we postulated ~~was~~ should be able to explain the high rate of enzyme synthesis which one observes in bacteria when the rate of formation of an enzyme is ~~is~~ <sup>maximally</sup> enhanced by the use of an inducer. The basic thought of ~~the~~ <sup>this</sup> theory ~~has~~ consists in the assumption that trinucleotides read the code of the parogene and that these trinucleotides carry amino acids. ~~Maximal~~ ~~rate~~ ~~of~~ ~~enzyme~~ ~~synthesis~~ ~~is~~ ~~observed~~ ~~when~~ ~~the~~ ~~rate~~ ~~of~~ ~~formation~~ ~~of~~ ~~an~~ ~~enzyme~~ ~~is~~ ~~enhanced~~ ~~by~~ ~~the~~ ~~use~~ ~~of~~ ~~an~~ ~~inducer~~. ~~One~~ ~~particular~~ ~~model~~ ~~for~~ ~~protein~~ ~~synthesis~~ ~~which~~ ~~assumed~~ ~~that~~ ~~each~~ ~~trinucleotide~~ ~~carries~~ ~~a~~ ~~sequence~~ ~~of~~ ~~three~~ ~~amino~~ ~~acids~~, ~~was~~ ~~singled~~ ~~out~~ ~~for~~ ~~detailed~~ ~~discussion~~ ~~because~~ ~~this~~ ~~model~~ ~~which~~ ~~appeared~~ ~~to~~ ~~be~~ ~~the~~ ~~most~~ ~~plausible~~. This does not mean, however, that other models need not be considered. ~~It~~ ~~is~~ ~~not~~ ~~clear~~ ~~that~~ ~~the~~ ~~model~~ ~~which~~ ~~we~~ ~~postulated~~ ~~is~~ ~~the~~ ~~most~~ ~~plausible~~.

Rather than to assume that each kind of trinucleotide carries a particular sequence of three amino acids, it would be in some ways more appealing to assume that each trinucleotide carries only one amino acid. In this case the amino acid might be carried by a phosphate group linked by an oxygen atom (ester linkage), either to the third or the fifth of the 5 carbon sugar of either the leading or the trailing nucleotide. Assuming twenty different trinucleotides, each carrying one particular amino acid, we could have a code that requires no commas, with no restrictions imposed on the possible amino acid sequences of the proteins formed by the paragenes.

The particular model for protein synthesis here considered cannot be modified simply by saying that each trinucleotide carries one amino acid instead of carrying three amino acids, for adjacent amino acids would not then be at the right distance from each other to be linked into a polypeptide.

Such an alternate model for protein synthesis would, therefore, require additional ideas concerning the configuration -- not necessarily a helical one -- which the parogene might adopt during protein synthesis.

I am grateful for the discussion of a variety of problems arising out of the considerations here presented which I had at the University of Chicago with Prof. Herbert S. Anker, Dr. Nandor L. Balazs, Mr. Hirono Kuki, Prof. Joseph E. Mayer, and Prof. Leonard J. Savage.

$$(x-1)$$

$$\frac{m}{k} + \ln m$$

Diagram showing a circle with  $1 + \frac{1}{x}$  inside, and  $x$  written below it. There are some scribbles below the circle.

$$\frac{m}{x^2} = \ln m$$

$$x^2 = \frac{\ln m}{\frac{m}{x^2}}$$

$$I \sqrt{\ln m}$$

$$x = \sqrt{\frac{\ln m}{m}}$$

$$II (x-1) \ln m$$

$$T_0 = 2 \sqrt{\ln m}$$



kan.

June 7, 1957

HOW MAY AMINO ACIDS READ THE NUCLEOTIDE CODE?

by Leo Szilard

~~(Submitted by Joseph E. Mayer)~~

6000 - 4  
1400

The Enrico Fermi Institute for Nuclear Studies

The University of Chicago, Chicago, Illinois

*(communicated by Joseph E. Mayer June 10 1957.)*

It is now generally believed that proteins are formed along ~~side~~ <sup>in a pattern determined by</sup> nucleic acid templates. The sequence of purine and pyrimidine bases in the template is supposed to represent a code that may somehow determine the sequence of the amino acids in the particular polypeptide (protein) that a given template will form. (1) The purine and pyrimidine bases of the template, the letters of the code, are adenine, uracil, guanine and cytosine if the template be an RNA molecule; and if the template be a DNA molecule, thymine takes the place of uracil.

It has remained so far a ~~complete~~ <sup>summary of a</sup> mystery in just what ~~conceivable~~ way amino acids could read such a code. In what manner can chemical forces of the kind we know to exist -- line up amino acids alongside such a template in the proper sequence and at the proper distance from each other, so that ~~straight~~ a chemical reaction chain may link adjacent amino acids through peptide bonds with each other?

It is the purpose of the ~~present~~ <sup>this</sup> paper to ~~indicate~~ <sup>present</sup> a conceptually simple scheme that will -- at least by way of an example -- illustrate in what manner this might ~~take~~ <sup>functionally</sup> place in the living cell.

~~A. L. Dawance, Enzymologia 15, 251, 1952~~  
(1) G. Gamow, Nature, Vol. 173, p. 318 (1954).

(2) ~~Added in proof:~~ <sup>by F. H. Crick, J. S. Griffith</sup> and L. E. Orgel which appeared in the May issue Proc. Nat. Acad. Sci. Vol 43 p 416 1957. containing a mysterious word on an oral communication by S. Brenner is

Because a template which synthesizes protein need not necessarily be the gene itself but must carry the same information as the corresponding gene, we shall here refer to such a template for the sake of brevity as a paragene.

The basic thought underlying the scheme here presented consists in the assumption that there are a number of enzymes (or enzyme systems) -- ~~perhaps twenty altogether~~ -- in the cell, and that each of these catalyzes the formation of a particular trinucleotide, which carries, ~~either one parti-~~ <sup>(2)</sup> ~~cular amino acid or, more likely perhaps,~~ a particular sequence of three amino acids. If the amino acid is carried by the nucleotide on a phosphate or pyrophosphate group as an acid anhydride -- which is a high energy compound -- then the energy needed for the formation of the peptide bonds will become free when the amino acid is split off. In this sense one can say that each amino acid may carry the energy needed for forming its peptide bond.

According to the notions here presented, amino acids can not/read ~~themselves~~ the code of the paragene. But the trinucleotides, which carry the proper amino acids, may attach with their three bases through the formation of 6 hydrogen bonds to the proper sequence of three bases on the paragene, and thus the amino acids may be lined up in the proper sequence along the paragene.

Accordingly, sequences of three nucleotides along the paragene represent the code words, and the trinucleotides which carry the amino acids represent the anti-code words. We assume that these anti-code words are complementary to the code words in the sense that, where the code word contains adenine the anti-code word contains uracil (or thymine) where the code word contains uracil (or thymine) the anti-code word contains adenine; and similarly guanine corresponds to cytosine and cytosine corresponds to guanine. The rationale for this assumption is as follows:

(2) ~~added on~~

(2) Note added June 17th. The May issue of the Proc.Nat.Acad.Sci., Vol. 43, p. 416 (1957) which was belatedly received here, contains an article by F.H.Crick, J.S.Griffith, and L.E.Orgel which is in part identical with their previously circulated memorandum referred to in the text. In addition, however, the authors discuss -- in conjunction with a detailed oral communication of S. Brenner to them -- how amino acids might read the nucleotide code. The relevant passage of their paper is quoted in full in Appendix C.

The concept of code-letter and complementary code-letter arose originally from the interpretation of the structure of DNA given by J.D. Watson and F.H.C.Crick. (2) They showed that in a double stranded DNA

---

(3) Nature, Vol. 173, p. 518 1953.  
Proc. Roy. Soc., Vol. 223, p. 80, 1954.

---

X <sup>helix</sup>~~structure~~, adenine of one strand pairs with thymine of the other strand (which presumably plays the same role in DNA as does uracil in RNA) and similarly guanine pairs with cytosine. The helical structure of DNA permits just such pairing, and hydrogen bonding is possible between adenine and thymine, as well as between guanine and cytosine.

If the sequence of bases along one strand of DNA represents a coded message which consists in three letter-words then, because we have four letters to choose from, such a message could utilize 64 different words. We might, however, be <sup>limited</sup>~~restricted~~ to the use of 20 out of the 64 words that are available. <sup>limitation</sup>~~restriction~~ The reasons for this ~~restriction~~ would be as follows:

If all 64 possible three-letter combinations form ~~infinite~~ a code word, and if the parogene assumes at the time of the formation of the polypeptide a helical configuration similar to the helical configuration of DNA, then it follows that the code on the parogene must be read consecutively from one end -- say, from the "head" of the parogene downward. This is so because such a helical structure does not provide for commas between the individual code words, and in a code containing 64 words any three consecutive letters form a word. If we number the letters along the parogene, say from the head of the parogene downward, then the first three letters, the letters 1, 2, 3, form a word which was meant to be conveyed and so do the next three letters, the letters 4, 5, 6. But sequences of three letters which encroach on two adjacent words (such as 2, 3, 4 or 3, 4, 5, for example) form code words which are not meant to be conveyed.

In these circumstances, the code would be misread if the trinucleotides, which represent the anti-code words, assemble alongside the parogene simultaneously, rather than -- from one end on -- consecutively. If we want to have simultaneous assembly of the trinucleotides alongside the comma-less parogene, then we are restricted to 20 code words.

The notion of such a 20-word code, which needs no commas, was introduced by F.H.C.Crick, J.S.Griffith, and L.E.Orgel of the Medical Research Council Unit at the Cavendish Laboratory, Cambridge, in a memorandum circulated in May, 1956 among workers interested in the subject of protein synthesis. (2) From such a code we must demand that the letters 1, 2, 3 on the template form a code word, and the letters 4, 5, 6 also form a code word, but sequences of three letters, which encroach on two adjacent words (such as 2, 3, 3 or 3, 4, 5, for example) form no code word. Crick and his co-workers have shown that this demand can be met, that a code which requires no commas may be constructed and that it can accommodate 20 three-letter code words.

We shall now single out for more detailed examination one conceivable model for protein synthesis which might provide for the lining up of the amino acids alongside the parogene, both in the proper order and at the proper distance from each other. This particular model is based on the following assumptions:

The trinucleotides which form the anti-code words contain the sugar ribose rather than the sugar desoxyribose. Each particular ribose trinucleotide (the anti-code word) carries a particular sequence of three amino acids. A phosphate (or diphosphate) group is attached to the (2) carbon atom of the ribose moiety of each nucleotide and an amino acid is attached to each of these phosphate (or diphosphate) groups. The amino acid anhydrides <sup>containing</sup> represent an energy-rich P, (or PP,) bond which, when split, may release 12,000 or 16,000 calories, respectively.

During protein synthesis the nucleic acid strand that functions as a template (the parogene) may take up -- so we here assume -- a helical configuration resembling the helical configuration of a DNA strand in the double stranded DNA helix. The trinucleotides may then line up alongside the helical parogene with their purine and pyrimidine bases paired with the complementary bases of the parogene, and if they are so lined up, then the amino acids carried by the trinucleotides ~~may~~ come to lie at just about the right distance from each other to permit the formation of a peptide bond between adjacent amino acids. A chemical reaction chain -- starting perhaps from the head of a parogene -- may then move down along the parogene, split the acid anhydrides, and thus free the amino acids as well as make available the energy needed for the formation of peptide bonds between adjacent amino acids.

Adjacent amino acids can be linked only if the distance from each other is smaller or equal, but not appreciably larger, than the fundamental repeating distance in a polypeptide chain, which is  $7.27\text{\AA}$ . The fundamental repeating distance in a fully extended polypeptide chain is about  $7\text{\AA}$ . Since before they are linked into a polypeptide, the amino acids can rotate around the chemical bond which ties them to the phosphate group, they might well be assembled along the parogene at a smaller distance from each other than the fundamental repeating distance of the polypeptide chain.

*(at the (2) carbon atom of the ribose)*  
somewhat

Applying the concept of a 20-letter code, that requires no commas, to our particular model of protein synthesis, we may now say the following:

Each of the 20 amino acids ~~may~~ appear once attached to the leading letter and once attached to the trailing letter of the 20 (trinucleotide) anti-code words. Therefore, among the polypeptides that can be formed, each amino acid ~~may~~ precede any other amino acid, and each amino acid may follow any other amino acid. This does not, of course, mean that any amino acid sequence is possible.

Some of the amino acid sequences that may be found experimentally in sequential analysis of proteins and polypeptides might show that the restrictions imposed by our model on the possible amino acid sequence are too severe, and that the model has to be modified to accommodate established facts. As will be presently seen some modification of our model may be indicated for other reasons also. There is certainly no inherent reason why we should have a pure three-letter code -- as we have assumed above -- and why, for instance, a certain number of four-letter code words should not be utilized also. (4)

If we had a pure three-letter code, we would have to demand that the number of amino acid residues of all polypeptides or proteins synthesized in the manner described above should be a multiple of three.

X

So far the number of amino acid residues in <sup>non-cyclic</sup> polypeptides and proteins have been determined only within one rather special class; all of them represent substances which are secreted by mammalian tissues. The number of amino acid residues found in such polypeptides and proteins, which have been analyzed with adequate accuracy, are as follows:

⚡

- a) Insulin chain A: 21; insulin chain B: 30; corticotropin B: 39; oxytocin: 9; vasopressin: 9; Intermedin B: 18.

All of these would fit a pure three-letter code.

- b) Intermedin A: 13; glucagon: 29 and pancreas ribonuclease: 124.

These do not fit a pure three-letter code. Intermedin and ribonuclease would have to include at least one 4-letter word and glucagon at least two 4-letter words. ~~In a mixed system of three- and four-letter words, one can obviously never draw the conclusion that more than two four-letter words have been included.~~

a multiple of 3.

(4) See Appendix A

(All numbers of longer chains can be represented as sum of 4 or 8 and ...)

### Observed rate of enzyme synthesis

According to the notions here adopted most enzymes are synthesized in growing bacteria at a rather low rate which does not represent the maximum synthesizing capacity of the corresponding paragenes. The rate of production of a given enzyme may, however, be greatly enhanced if the enzyme is induced, and what we are interested to learn is the maximal rate at which a paragene may be able to form the corresponding enzyme.

One of the most studied cases of enzyme induction is the induction of the enzyme  $\beta$ -galactosidase which splits lactose. Jacques Monod and his co-workers have shown that the production rate of this enzyme in bacteria can be greatly enhanced by certain chemical analogues of lactose, which act as inducers, and that the rate of production of the enzyme goes up almost instantaneously upon adding such an inducer to the medium. We are thus led to believe that the inducer may act by increasing the rate at which one template produces the enzyme rather than by increasing the number of templates that produce the enzyme at an unchanged rate.

In fully induced wild type bacteria growing in minimal medium this enzyme is contained in the amount of about 8  $\mu$ g., per  $10^{12}$  bacteria and thus amounts to about 8% of the total proteins. We obtain the rate, at which this enzyme is produced in minimal medium per bacterium, by dividing the amount contained in one bacterium by 1.44 times the doubling time (40 minutes) of the bacterium. We thus find for the rate, at which this enzyme is produced in fully induced wild type bacteria growing in minimal medium, about  $2 \cdot 10^{-18}$  grams per cell per second.

If we assume a molecular weight of a million (Jacques Monod and Melvin Cohn estimate the molecular weight of this enzyme at about 800,000), we obtain a rate of 1.5 enzyme molecules per cell per second. The number of paragenes per cell is not known. ~~It is not known if there are many paragenes~~ There might be a few paragenes present per cell rather than just one, and the number of paragenes might be of the order of magnitude of 10. On the other hand, smaller enzyme molecules might be synthesized somewhat faster than larger enzyme molecules.

On the basis of the figure given above, we are thus led to believe ~~that~~ <sup>in bacteria --</sup> when an enzyme is fully induced and enzyme synthesis proceeds at its maximal rate -- the rate of formation of the enzyme may be of the order of magnitude of one per second per paragene.

## Computed rate of enzyme synthesis

We shall now attempt to compute at what rate a parogene may be able to synthesize the corresponding enzyme on the basis of the model that we have postulated. For the purposes of this computation, we shall assume that the molecular weight of the enzyme is about 100,000. We then have about 1,000 amino acid residues in the enzyme, and accordingly we would have to assemble alongside the parogene  $m = 300$  trinucleotides, each of which is "loaded" with three amino acids.

use

In the approximation which we shall ~~consider~~, the minimum time,  $\tau$ , needed for the formation of the polypeptide is composed of two terms,  $\tau_1$  and  $\tau_2$ .

$$\tau = \tau_1 + \tau_2$$

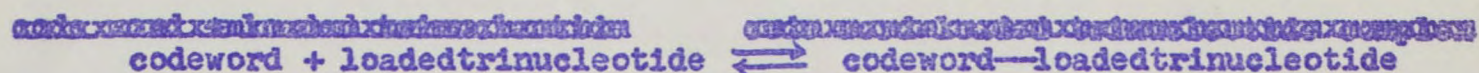
After all the amino acids, assembled alongside the template, have been linked into one polypeptide, and <sup>if</sup> ~~assuming that~~ this polypeptide is at once removed, a certain time,  $\tau_1$ , will elapse until the trinucleotides, which are now denuded of amino acids, evaporate from the template and their place is taken by trinucleotides which are loaded with the proper amino acids. We ~~shall~~ assume that the concentration of denuded trinucleotides in the cell is very small so that after the denuded trinucleotides evaporate, the loaded trinucleotides do not have to compete with denuded trinucleotides for their legitimate positions along the parogene. The time,  $\tau_1$ , which is necessary to permit evaporation of  $m$  denuded trinucleotides and to assemble  $m$  loaded trinucleotides in their place we shall compute here on the assumption that once a loaded trinucleotide has found its position alongside the template, it will not evaporate again. Because this assumption is ~~probably~~ not valid, we must make a correction which is represented by the second term,  $\tau_2$ .

### Computation of $\tau_2$

First we shall now compute this second term,  $\tau_2$ . This computation <sup>is</sup> ~~will~~ be based on the fact that (because of reevaporation of the loaded trinucleotides which ~~are~~ reversibly combined with the anti-code words of the parogene) there will be - no matter how long we wait - always a certain number of code words (sequences of three nucleotides) on the parogene which are not "covered".



We have to deal with a reversible reaction which may be written as follows:



The rate at which the reaction proceeds from left to right; i.e. the number of successful hits per unit time is given by

$$\text{hit rate} = A \rho$$

$\rho$  denotes the concentration of the particular kind of loaded trinucleotide, in mol/cc. We shall, for the sake of simplicity, assume that the concentration in the cell of each kind of loaded trinucleotide is the same.

A stands for

$$A = 6 \cdot 10^{23} v \sigma p$$

$v$  is the molecular velocity,  $v = \sqrt{\frac{2RT}{\pi M}}$  where  $M$  is the molecular weight of the loaded trinucleotide. If  $M \approx 1000$ , we have  $v \approx 5 \times 10^3$  cm/sec.

$\sigma$  is the target area that must be hit by the loaded trinucleotide if hydrogen bonding is to take place with the three adjacent nucleotides on the parogene. We may assume for  $\sigma$  a value of  $\sigma = 10^{-15}$  cm<sup>2</sup>.

$p$  denotes the probability that the loaded trinucleotide, when hitting the code word, is in just the right geometrical position to permit hydrogen bonding to take place between the three complementary pairs of bases that are involved. *assuming zero activation energy* We may take for  $p$ , as a very rough estimate  $p = 1/300$ ; This would give for  $\sigma p$  the value  $\sigma p = 1/3 \cdot 10^{-17}$ .

The rate at which the reaction proceeds from left to right is given by the rate  $\alpha$  at which the ~~codeword—loaded trinucleotide~~ complex dissociates or, as we may also say, the rate  $\alpha$  at which the loaded trinucleotide evaporates from the parogene. For this rate we may write

$$(2) \quad \alpha = \text{rate of evaporation} \approx 10^{13} e^{-\frac{\Delta H}{RT}}$$

where  $\Delta H$  is the binding energy for the loaded trinucleotides.

In equilibrium the hit rate and evaporation rate must be equal.

The equilibrium constant,  $K$ , of this reversible reaction is defined as the concentration of the loaded trinucleotides at which the code word is covered half of the time. Thus we may write

$$(2) \quad AK \approx \frac{1}{2} 10^{13} e^{-\frac{\Delta H}{RT}} = \frac{1}{2} \alpha$$

$$\alpha = 2AK$$

In equilibrium the probability,  $f$ , for a given code word on the parogene not being covered by the proper loaded trinucleotide is given by

$$(4) \quad f = \frac{1}{1 + \frac{P}{K}}$$

Accordingly, in equilibrium, the total number of such gaps along the parogene which contains  $m$  nucleotides is given by

$$(5) \quad \text{"number of gaps"} = \frac{m}{1 + \frac{P}{K}}$$

We shall assume that most code words are "covered" in equilibrium and this means that

$$(6) \quad \frac{P}{K} \gg 1$$

We presume that after such equilibrium is established a chemical reaction chain is somehow triggered, and, moving down along the parogene, links adjacent amino acids into a polypeptide. The average time,  $\tau$ , needed for the formation of the polypeptide from the amino acids assembled along the parogene is given by the product of the "number of gaps", that

\* One may verify that the reaction is not diffusion <sup>and that the</sup> ~~limited~~ <sup>11.\*</sup> ~~by~~ <sup>inclusion of the</sup> ~~the~~ <sup>triplets</sup>

have to be filled consecutively, and the average time,  $\frac{1}{AP}$ , that it takes to fill one given gap. Thus, for  $\tau_2$  we may write

(7) 
$$\tau_2 = \frac{1}{AP} \frac{m}{1 + \frac{P}{K}} \approx \frac{1}{50} \cdot 30 = \frac{1}{2}$$

in the vicinity of the paragenesis is not appreciably lower than?

Computation of  $\tau_1$  and  $\tau_0 = \tau_1 + \tau_2$

When the polypeptide is formed and leaves the paragenesis, the code words are covered with the ~~denuded~~ <sup>stripped</sup> trinucleotides. We may now compute the average time,  $\tau_1$ , needed for the evaporation of all the ~~denuded~~ <sup>stripped</sup> trinucleotides and the assembling of all the loaded trinucleotides in their place. We shall, for the sake of simplicity, assume that a ~~denuded~~ <sup>stripped</sup> trinucleotide evaporates at the same rate,  $\alpha$ , as a loaded trinucleotide. The rate,  $\alpha$ , at which a loaded trinucleotide evaporates from the template is given by (3) and we may also write this in the form

(8) 
$$\alpha = AP \frac{2K}{P} \approx 10$$

Since we have assumed  $\frac{P}{K} \gg 1$ , we have

(9) 
$$AP \gg \alpha$$

As may be shown, ~~known as the~~ <sup>(5)</sup> ~~for very large~~ <sup>values of m,  ( $\ln m \gg 1$ ), we may write for  $\tau_1$</sup>

(10) 
$$\tau_1 \approx \frac{1}{AP} \frac{P}{2K} \ln m = \frac{1}{10} 5.6$$

For the total time,  $\tau_0 = \tau_1 + \tau_2$  we thus obtain

(11) 
$$\tau_0 \approx \frac{1}{AP} \left\{ \frac{m}{1 + \frac{P}{K}} + \frac{1}{2} \frac{P}{K} \ln m \right\}$$

If we wish to make this time as small as possible, we have to choose K so as to have for  $\frac{P}{K}$

(12) 
$$\frac{P}{K} \approx \sqrt{\frac{2m}{\ln m}}$$

(5) See Appendix B

Substituting this value into (11) gives

$$(13) \quad \tau_0 \approx \frac{\sqrt{2}}{A\beta} \sqrt{m \ln m}$$

For a polypeptide containing 1,000 amino acid residues, i.e. a paragene containing about 300 code words, we may write  $m = 300$ , and thus we obtain from (12) and (13)

$$(14) \quad \frac{\beta}{K} \approx 10$$

and

$$(15) \quad \tau_0 \approx \frac{50}{A\beta} \text{ sec}$$

XXX

Estimates for the values of  $\sigma p$ ,  $\Lambda$ ,  $\alpha$  and  $K$

See.

XXX

As we may see from (15) we obtain  $\tau_0 = 1$  if  $A\beta = 50$ . This means that for this particular value of  $A\beta$  one enzyme molecule is produced per paragene per second. As we have seen before, this is the order of magnitude of the rate at which highly induced bacteria produce the enzyme  $\beta$ -galactosidase per paragene.

We shall, therefore, in the following assume  $A\beta = 50$ , and compute from it  $\beta$ , the concentration at which the different kinds of trinucleotides may be present in the cell.

If we use for  $\sigma p$  the value of  $\sigma p = \frac{1}{3} 10^{-17} \text{ cm}^2$ , then we obtain from (1)  $\Lambda$ ;  $\Lambda = 10^{10}$  and accordingly we have  $\beta = 5 \cdot 10^{-9} \text{ mol/cc}$  ( $\beta = 5 \cdot 10^{-6} \text{ mol/liter}$ ).

It might well be, however, that  $\sigma p$  is ten times higher so that we have  $\sigma p = \frac{1}{3} 10^{-16} \text{ cm}^2$ , and then we obtain from (1) for  $\Lambda$ ;  $\Lambda = 10^{11}$ , so that we have  $\beta = 5 \cdot 10^{-10} \text{ mol/cc}$  or ( $\beta = 5 \cdot 10^{-7} \text{ mol/liter}$ ).

Thus the concentration,  $\beta$ , of the different kinds of trinucleotides in the cell is likely to be between  $5 \times 10^{-6}$  and  $5 \times 10^{-7} \text{ mol/liter}$ , and one would have to look for concentrations of this order of magnitude in order to obtain experimental confirmation of their presence.

Since we have assumed  $\frac{p}{K} \approx 10$ , and since we believe that we have  $A_j = 50$ , it follows that we must have for the rate of evaporation of the trinucleotides from the parogene,  $\alpha$

$$(16) \quad \alpha = 2AK = 10/\text{sec.} \quad \text{or} \quad 1/\alpha = 1/10 \text{ sec.}$$

For values of A between  $10^{10}$  and  $10^{11}$ , we have

$$(17) \quad 5 \cdot 10^{-11} < K \text{ (in mol/cc)} < 5 \cdot 10^{-10}$$

$$5 \cdot 10^{-8} < K \text{ (in mol/liter)} < 5 \cdot 10^{-7} \quad \text{or}$$

From (3) we obtain for a given pair of values A and K the binding energy  $\Delta H$  between the loaded trinucleotide and the code word on the parogene. For a value of  $A = 10^{10}$  and  $K = 10^{-10}$  mol/cc ( $K = 10^{-7}$  mol/liter), we obtain  $\Delta H$  18,000 calories. Since 6 hydrogen bonds are involved, this would mean about 3,000 calories per hydrogen bond.

For the same value of A and a value of K which is ten times larger,  $\Delta H$  would decrease by about 1400 calories.

### Conclusion

These considerations show that the theory which we postulated should be able to explain the high rate of enzyme synthesis which one observes in bacteria when the rate of formation of an enzyme is maximally enhanced by the use of an inducer. The basic thought of this theory consists in the assumption that trinucleotides and possibly also tetranucleotides read the code of the parogene, and that these oligonucleotides carry amino acids. One particular model for protein synthesis, which assumed that each trinucleotide or tetranucleotide carries a sequence of three or four amino acids respectively, was singled out for detailed discussion because this model appeared to be the most plausible. This does not mean, however, that other models can not be considered.

For instance, rather than to assume that each kind of trinucleotide carries a particular sequence of three amino acids, one might wish to explore the possibility that each trinucleotide might carry only one amino acid. In this case the amino acid might be carried by a phosphate group linked by an oxygen atom (ester linkage), either to the third or the fifth carbon of the 5 carbon sugar of either the leading or the trailing nucleotide. Assuming twenty different trinucleotides, each carrying one particular amino acid, we could have a code that requires no commas, with no restrictions imposed on the possible amino acid sequences of the proteins formed by the paragenes.

However, the particular model for protein synthesis here considered cannot be modified by simply saying that each trinucleotide carries one amino acid instead of carrying three amino acids, for adjacent amino acids would not then be at the right distance from each other to be linked into a polypeptide. Such an alternate model for protein synthesis would, therefore, require additional ideas concerning the formation of the polypeptide.

I am grateful for the discussion of a variety of problems arising out of the considerations here presented which I had at the University of Chicago with Prof. Herbert S. Anker, Dr. Nandor L. Balazs, Mr. Hircendo Kulci, Prof. Joseph L. Moyer, and Prof. Leonard J. Savage.

## APPENDIX

(Added June 17, 1957 )

A.) If, in addition to ribose trinucleotides, carrying three amino acids, we have in the role of anti-code words also ribose tetranucleotides, carrying four amino acids, then we must postulate that only three of the bases of the tetranucleotide pair with complementary bases on the paragene. The fourth base, if it is a purine, must pair with the wrong pyrimidine, and if it is a pyrimidine, it must pair with the wrong purine. The reason for this is as follows:

As formula (17) shows,  $\alpha$ , the rate of evaporation of a trinucleotide from the paragene can be estimated to be about 10/sec., corresponding to a binding energy,  $\Delta H$ , of about 18,000 calories or 3,000 calories per hydrogen bond. If all four bases of the tetranucleotide were paired with the complementary bases on the paragene, we would then have two more hydrogen bonds and presumably an additional binding energy of about 6,000 calories. The equilibrium constant,  $K$ , of the tetranucleotide would therefore be lower than the value computed for  $K$  of the trinucleotides by a factor of about  $10^4$ , and the rate of evaporation,  $\alpha$ , would be lower by the same factor. This would then make,  $\tau_0$ , the minimal time it takes for a paragene to form a polypeptide too long to be compatible with the observed rate of production of  $\beta$ -galactosidase in highly induced bacteria.

The number of words that may be constructed in a mixed three-letter word and four-letter word code of the kind described above, <sup>where</sup> if we demand that the code require no commas, has so far not been determined.

B.) The problem of computing  $\tau_1$  amounts to the solving of the following problem:

There are  $n$  boxes, each of which can hold one white ball or one black ball. Initially, at time,  $\lambda = 0$ , all of these boxes contain one

white ball -- a stripped trinucleotide. These white balls evaporate at the rate,  $\alpha$ , so that at time  $\lambda$  the probability,  $W(\lambda)$  of having no white ball in the box is given by

$$(18) \quad W(\lambda) = 1 - e^{-\alpha\lambda}$$

If the rate at which black balls fall into an empty box is designated by  $\beta$ , then the probability,  $y$ , that a given box does not contain a black ball at time,  $t$ , is given by

$$(19) \quad y(t) = e^{-\alpha t} + \int_{\lambda=0}^{\lambda=t} e^{-\beta(t-\lambda)} \frac{dW}{d\lambda} d\lambda$$

we may write

$$(20) \quad \int_{\lambda=0}^{\lambda=t} e^{-\beta(t-\lambda)} \frac{dW}{d\lambda} d\lambda = \frac{\alpha}{\beta-\alpha} \left[ e^{-\alpha t} - e^{-\beta t} \right]$$

so that we have

$$(21) \quad y(t) = \frac{1}{\beta-\alpha} (\beta e^{-\alpha t} - \alpha e^{-\beta t})$$

Thus we may write for the probability,  $x$ , that a given box contains a black ball

$$(22) \quad x(t) = 1 - y(t) = 1 - \frac{\beta}{\beta-\alpha} e^{-\alpha t} + \frac{\alpha}{\beta-\alpha} e^{-\beta t}$$



From this we obtain for  $P(t)$ , the probability that all  $m$  boxes contain one black ball

$$(23) \quad P(t) = x(t)^m = \left( 1 - \frac{\beta}{\beta - \alpha} e^{-\alpha t} + \frac{\alpha}{\beta - \alpha} e^{-\beta t} \right)^m$$

The average time,  $\tau_1$ , needed for the evaporation of the white balls from all  $m$  boxes and the filling of all  $m$  boxes with black balls is given by

$$(24) \quad \tau_1 = \int_0^{\infty} t \frac{dP}{dt} dt$$

If  $\beta \gg \alpha$ , we may write

$$(25) \quad \tau_1 = \int_0^{\infty} t \frac{d}{dt} (1 - e^{-\alpha t})^m dt$$

The expression  $\frac{d}{dt} (1 - e^{-\alpha t})^m$  has a maximum at some value of  $t$ ;  $t = \tau_0$  and -- for large values of  $m$  -- it becomes small very rapidly both below and above  $t = \tau_0$ . Therefore, if  $m$  is large we may write

$$(26) \quad \tau_1 \approx \tau_0$$

We obtain  $\tau_0$  by writing

$$(27) \quad \left\{ \frac{d^2}{dt^2} (1 - e^{-\alpha t})^m \right\}_{t=\tau_0} = 0$$

and from this we obtain

$$(28) \quad \tau_0 = \frac{1}{\alpha} \ln m$$

Thus for  $\beta \gg \alpha$  and  $m \gg 1$  we may write (26)

$$(29) \quad \tau_1 \approx \frac{1}{\alpha} \ln m$$

It can be shown from (23) and (24) that in the next higher approximation we have

$$(30) \quad \tau_1 \approx \frac{1}{\alpha} \left\{ \ln(m+1) + C \right\} + \Delta ; \text{ where } \begin{cases} C = 0.577 \text{ (Euler's const.)} \\ \frac{1}{\beta} < \Delta < \frac{1}{\beta - \alpha} . \end{cases}$$

C.) F. H. Crick, J. S. Griffith and L. E. Orgel write in the May issue of the Proc. Nat. Acad. Sci. (Vol. 43, pp. 419 and 420, 1957):

|| To fix ideas, we shall describe a simple model to illustrate the advantages of such a code. Imagine that a single chain of RNA, held in a regular configuration, is the template. Let the intermediates in protein synthesis be 20 distinct molecules, each consisting of a trinucleotide chemically attached to one amino acid. The bases of each trinucleotide are chosen according to the code given above. Let these intermediate molecules combine, by hydrogen bonding between bases, with the RNA template and there await polymerization. Now imagine that such an amino acid-trinucleotide were to diffuse into an incorrect place on the template, such that two of its bases were hydrogen-bonded, though not the third. We postulate that this incomplete attachment will only retain the intermediate for a very brief time (for example, less than 1 millisecond) before the latter breaks loose and diffuses elsewhere. However, when it eventually diffuses to the correct place, it will be held by hydrogen bonds to all three bases and will thus be retained, on the average, for a much longer time (say, seconds or minutes). Now the code we have described insures that this more lengthy attachment can occur only at the points where the intermediate is needed. If one of the 20 intermediates could

stay for a long time on one of the false positions, it would effectively block the two positions it was straddling and hold up the polymerization process. Our code makes this impossible. This scheme, therefore, allows the intermediates to accumulate at the correct positions on the template without ever blocking the process by settling, except momentarily, in the wrong place. It is the feature which gives it an advantage over schemes in which the intermediates are compelled to combine with the template one after the other in the correct order.

The example given here is only for illustration, but it brings out the physical idea behind the concept of a comma-less code.

In passing, it should be mentioned that while the idea of making three nonoverlapping nucleotides code for one amino acid at first sight entails certain stereochemical difficulties, these are not insuperable if it is assumed that the polypeptide chain, when polymerized, does not remain attached to the template. A detailed scheme along these lines has been described to us by Dr. S. Brenner (personal communication).<sup>4</sup>

June 7, 1957

X30  
10 m/12 #8a  
Use Webster's  
dictionary

10/12 caps ital #8c

HOW MANY AMINO ACIDS READ THE NUCLEOTIDE CODE?

Paper

10/12 caps #8a By Leo Szilard

Communicated by Joseph E. Mayer

June 10, 1957

8/10 ital #8c

Enrico Fermi Institute for Nuclear Studies,

University of Chicago, Chicago, Illinois

transpose

It is now generally believed that proteins are formed alongside

nucleic acid templates. The sequence of purine and pyrimidine bases in the template is supposed to represent a code that may somehow determine the sequence of the amino acids in the particular polypeptide (protein) that a given template will form. The purine and pyrimidine bases of the template, the letters of the code, are adenine, uracil, guanine, and cytosine if the template <sup>is</sup> an RNA molecule; and if the template <sup>is</sup> a DNA molecule, thymine takes the place of uracil.

Letters ital only where underlined

It has remained so far <sup>somewhat of a</sup> ~~complete~~ mystery in just what conceivable way amino acids could read such a code. In what manner can chemical forces of the kind we know to exist <sup>can</sup> line up amino acids alongside such a template in the proper sequence and at the proper distance from each other, so that ~~there might~~ a chemical reaction chain may link adjacent amino acids through peptide bonds with each other?

It is the purpose of the present paper to indicate a conceptually simple scheme that will <sup>can</sup> at least by way of an example <sup>can</sup> illustrate in what manner this might take place in the living cell.

transpose to p. 19

~~A. L. Dounce, Enzymologia 15, p. 251, 1952~~  
~~G. Gamow, Nature, Vol. 173, p. 318 (1954)~~

④ 1/A. L. Dounce, Enzymologia, 15, 251, 1952; G. Gamow, Nature, 173, 318, 1954.

8/10 #8a #8c +25

5/2

Because a template which synthesizes protein need not necessarily be the gene itself but must carry the same information as the corresponding gene, we shall here refer to such a template for the sake of brevity as a paragene.

The basic thought underlying the scheme here presented consists in the assumption that there are a number of enzymes (or enzyme systems) <sup>em</sup> perhaps twenty <sup>R#</sup> altogether <sup>em</sup> in the cell, and that each of these catalyzes the formation of a particular trinucleotide which carries either one particular amino acid <sup>(2)</sup> or, more likely, perhaps, a particular sequence of three amino acids. If the amino acid is carried by the nucleotide on a phosphate or pyrophosphate group as an acid anhydride <sup>em</sup> -- which is a high-energy compound <sup>em</sup> -- then the energy needed for the formation of the peptide bonds will become free when the amino acid is split off. In this sense one can say that each amino acid may carry the energy needed for forming its peptide bond.

According to the notions here presented, amino acids can <sup>themselves</sup> not read ~~the~~ the code of the paragene. But the trinucleotides, which carry the proper amino acids, may attach with their three bases through the formation of <sup>Six</sup> hydrogen bonds to the proper sequence of three bases on the paragene, and thus the amino acids may be lined up in the proper sequence along the paragene.

Accordingly, sequences of three nucleotides along the paragene represent the code words, and the trinucleotides which carry the amino acids represent the anti-code words. We assume that these anti-code words are complementary to the code words in the sense that, where the code word contains adenine, the anti-code word contains uracil (or thymine); where the code word contains uracil (or thymine), the anti-code word contains adenine; and, similarly, guanine corresponds to cytosine and cytosine corresponds to guanine. The rationale for this assumption is as follows:

Paragene

transfer to p. 19

Proc. Nat. Acad. Sci. 1957

Note added June 17th. The May issue of the Proc. Nat. Acad. Sci. ~~43~~ 43, 416, 1957, which was belatedly received here, contains an article by F.H.Crick, J.S.Griffith, and L.E.Orgel which is in part identical with their previously circulated memorandum referred to in the text. In addition, however, the authors discuss <sup>em</sup> in conjunction with a detailed oral communication of S. Brenner to them <sup>em</sup> how amino acids might read the nucleotide code. The relevant passage of their paper is quoted in full in Appendix C.

8/10  
#62  
+FC  
+25

*transpose to p. 19*

The concept of code<sup>#</sup>-letter and complementary code<sup>#</sup>-letter arose originally from the interpretation of the structure of DNA given by J.D. Watson and F.H.C. Crick. <sup>3</sup> They showed that in a double-stranded DNA

Nature, ~~173~~ 173, 318, 1953; Proc. Roy. Soc. London, A, 223, 80, 1954.  
~~Proc. Roy. Soc., Vol. 223, p. 80, 1954.~~

810  
#89  
+60  
+255

structure, adenine of one strand pairs with thymine of the other strand (which presumably plays the same role in DNA as does uracil in RNA), and similarly guanine pairs with cytosine. The helical structure of DNA permits just such pairing, and hydrogen bonding is possible between adenine and thymine, as well as between guanine and cytosine.

If the sequence of bases along one strand of DNA represents a coded message which consists <sup>of</sup> ~~of~~ three-letter<sup>#</sup> words, then, because we have four letters to choose from, such a message could utilize 64 different words. We might, however, be <sup>limited</sup> ~~restricted~~ to the use of 20 out of the 64 words that are available. The reasons for this <sup>limitation</sup> ~~restriction~~ would be as follows:

If all 64 possible three-letter combinations form ~~in fact~~ a code word, and if the parogene assumes at the time of the formation of the polypeptide a helical configuration similar to the helical configuration of DNA, then it follows that the code on the parogene must be read consecutively from one end <sup>we</sup> say, from the "head" of the parogene downward. This is so because such a helical structure does not provide for commas between the individual code words, and in a code containing 64 words any three consecutive letters form a word. If we number the letters along the paraf<sup>f</sup> gene, ~~we~~ from the head of the parogene downward, then the first three letters, the letters 1, 2, 3, form a word which was meant to be conveyed, and so do the next three letters, the letters 4, 5, 6. But sequences of three letters which encroach on two adjacent words (such as 2, 3, 4 or 3, 4, 5, for example) form code words which are not meant to be conveyed.

In these circumstances, the code would be misread if the trinucleotides, which represent the anti-code words, assembled alongside the parogene simultaneously, rather than <sup>aw</sup> from one end on <sup>aw</sup> consecutively. If we want to have simultaneous assembly of the trinucleotides alongside the comma-less parogene, then we are restricted to 20 code words.

The notion of such a 20-word code, which needs no commas, was introduced by F.H.C. Crick, J.S. Griffith, and L.E. Orgel, of the Medical Research Council Unit at the Cavendish Laboratory, Cambridge, in a memorandum circulated in May, 1956, among workers interested in the subject of protein synthesis. <sup>(2)</sup> From such a code we must demand that the letters 1, 2, 3 on the template form a code word, and the letters 4, 5, 6 also form a code word, but sequences of three letters, which encroach on two adjacent words (such as 2, 3, 3 or 3, 4, 5, for example) form no code word. Crick and his co-workers have shown that this demand can be met, that a code which requires no commas may be constructed, and that it can accommodate 20 three-letter code words.

We shall now single out for more detailed examination one conceivable model for protein synthesis which might provide for the lining-up of the amino acids alongside the parogene, both in the proper order and at the proper distance from each other. This particular model is based on the following assumptions:

The trinucleotides which form the anti-code words contain the sugar ribose rather than the sugar desoxyribose. Each particular ribose trinucleotide (the anti-code word) carries a particular sequence of three amino acids. A phosphate (or diphosphate) group is attached to the (2) carbon atom of the ribose moiety of each nucleotide, and an amino acid is attached to each of these phosphate (or diphosphate) groups. The amino acid anhydrides <sup>contain</sup> ~~represent~~ an energy-rich P, or PP, bond which, when split, may release 12,000 or 16,000 calories, respectively.

see duplicate reference

During protein synthesis the nucleic acid strand that functions as a template (the paragene) may take up <sup>an</sup> -- so we here assume <sup>an</sup> -- a helical configuration resembling the helical configuration of a ~~///~~ DNA strand in the double-stranded DNA helix. The trinucleotides may then line up alongside the helical paragene with their purine and pyrimidine bases paired with the complementary bases of the paragene, and, if they are so lined up, then the amino acids carried by the trinucleotides may come to lie at just about the right distance from each other to permit the formation of a peptide bond between adjacent amino acids. A chemical reaction chain <sup>an</sup> -- starting perhaps from the head of a paragene <sup>an</sup> -- may then move down along the paragene, split the acid anhydrides, and thus free the amino acids as well as make available the energy needed for the formation of peptide bonds between adjacent amino acids.

OK?

Adjacent amino acids can be linked only if the <sup>in</sup> distance <sup>is</sup> from each other <sup>are</sup> ~~is~~ smaller <sup>than</sup> or equal <sup>to,</sup> but not appreciably larger <sup>than</sup> the fundamental repeating distance in a polypeptide chain, which is 7.27Å. The fundamental repeating distance in a fully extended polypeptide chain is about 7Å. Since before they are linked into a polypeptide, <sup>(at the 2 carbon atom of the ribose)</sup> the amino acids can rotate around the chemical bond which ties them to the phosphate group, they might well be assembled along the paragene at a <sup>somewhat</sup> smaller distance from each other than the fundamental repeating distance of the polypeptide chain.

Applying the concept of a 20-letter code, that requires no commas, to our particular model of protein synthesis, we may now say the following:

Each of the 20 amino acids may appear once attached to the leading letter and once attached to the trailing letter of the 20 (trinucleotide) <sup>words</sup> anti-code words. Therefore, among the polypeptides that can be formed, each amino acid may precede any other amino acid, and each amino acid may follow any other amino acid. This does not, of course, mean that any amino acid sequence is possible.



Some of the amino acid sequences that may be found experimentally in sequential analysis of proteins and polypeptides might show that the restrictions imposed by our model on the possible amino acid sequence are too severe, and that the model has to be modified to accommodate established facts. As will be presently seen, some modification of our model may be indicated for other reasons also. There is certainly no inherent reason why we should have a pure three-letter code <sup>am</sup>-- as we have assumed above <sup>am</sup>-- and why, for instance, a certain number of four-letter code words should not be utilized also. ~~4/1~~

If we had a pure three-letter code, we would have to demand that the number of amino acid residues of all polypeptides or proteins synthesized in the manner described above should be a multiple of three.

So far the number<sup>s</sup> of amino acid residues in <sup>nontcyclic</sup> polypeptides and proteins ~~have~~ <sup>have</sup> been determined only within one rather special class; all of them represent substances which are secreted by mammalian tissues. The number<sup>s</sup> of amino acid residues found in such polypeptides and proteins, which have been analyzed with adequate accuracy, are as follows:

(a) Insulin chain A, 21; insulin chain B, 30; corticotropin B, 39; oxytocin, 9; vasopressin, 9; Intermedin B, 18.

All of these would fit a pure three-letter code.

(b) Intermedin A, 13; glucagon, 29; and pancreas ribonuclease, 124.

These do not fit a pure three-letter code. Intermedin and ribonuclease would have to include at least one <sup>four-</sup>letter word and glucagon at least two <sup>four-</sup>letter words. ~~In a mixed system of three- and four-letter words, one can obviously never draw the conclusion that more than two four-letter words have been included.~~

4/1 Lee Appendix A.

transfer to p. 19

# A Observed Rate of Enzyme Synthesis.

According to the notions here adopted, most enzymes are synthesized in growing bacteria at a rather low rate which does not represent the maximum synthesizing capacity of the corresponding paragenes. The rate of production of a given enzyme may, however, be greatly enhanced if the enzyme is induced, and what we are interested <sup>in</sup> ~~to~~ <sup>learn</sup> ~~is~~ <sup>ing</sup> is the maximal rate at which a paragene may be able to form the corresponding enzyme.

9b # 26 One of the most studied cases of enzyme induction is the induction of the enzyme  $\beta$ -galactosidase which splits lactose. Jacques Monod and his co-workers have shown that the production rate of this enzyme in bacteria can be greatly enhanced by certain chemical analogues of lactose, which act as inducers, and that the rate of production of the enzyme goes up almost instantaneously upon adding such an inducer to the medium. We are thus led to believe that the inducer may act by increasing the rate at which one template produces the enzyme rather than by increasing the number of templates that produce the enzyme at an unchanged rate.

In fully induced wild-type bacteria growing in minimal medium this enzyme is contained in the amount of about 8  $\mu$ g. per  $10^{12}$  bacteria and thus amounts to about 8% <sup>per cent</sup> of the total proteins. We obtain the rate at which this enzyme is produced in minimal medium per bacterium, by dividing the amount contained in one bacterium by 1.44 times the doubling time (40 minutes) of the bacterium. We thus find, for the rate at which this enzyme is produced in fully induced wild-type bacteria growing in minimal medium, about  $2 \times 10^{-18}$  <sup>gm.</sup> ~~grams~~ per cell per second. times

If we assume a molecular weight of a million (Jacques Monod and Malvin Cohn estimate the molecular weight of this enzyme at about 800,000), we obtain a rate of 1.5 enzyme molecules per cell per second. The number of paragenes per cell is not known. ~~It is not known whether~~ There might be a few paragenes present per cell rather than just one, and the number of paragenes might be of the order of magnitude of 10. On the other hand, smaller enzyme molecules might be synthesized somewhat faster than larger enzyme molecules.

5/? On the basis of the figure given above, we are thus led to believe <sup>in bacteria</sup> ~~that~~ <sup>← EM dash</sup> when an enzyme is fully induced and enzyme synthesis proceeds at its maximal rate <sup>the</sup> ~~the~~ rate of formation of the enzyme may be of the order of magnitude of one per second per paragene.

# 4 Computed Rate of Enzyme Synthesis. em

5/7

We shall now attempt to compute at what rate a paragene may be able to synthesize the corresponding enzyme on the basis of the model that we have postulated. For the purposes of this computation, we shall assume that the molecular weight of the enzyme is about 100,000. We then have about 1,000 amino acid residues in the enzyme, and accordingly we would have to assemble alongside the paragene  $m = 300$  trinucleotides, each of which is "loaded" with three amino acids.

gr #43

In the approximation which we shall ~~use,~~ <sup>use,</sup> ~~assume,~~ the minimum time,  $\tau_0$ , needed for the formation of the polypeptide is composed of two terms,  $\tau_1$  and  $\tau_2$ .

$\tau_0 = \tau_1 + \tau_2$ .

zero

After all the amino acids, assembled alongside the template, have been linked into one polypeptide, ~~and~~ <sup>if</sup> ~~assuming that~~ this polypeptide is at once removed, a certain time,  $\tau_1$ , will elapse until the trinucleotides, which are now denuded of amino acids, evaporate from the template and their place is taken by trinucleotides which are loaded with the proper amino acids. We ~~shall~~ assume that the concentration of denuded trinucleotides in the cell is very small, so that after the denuded trinucleotides evaporate, the loaded trinucleotides do not have to compete with denuded trinucleotides for their legitimate positions along the paragene. The time,  $\tau_1$ , which is necessary to permit evaporation of  $m$  denuded trinucleotides and to assemble  $m$  loaded trinucleotides in their place we shall compute here on the assumption that once a loaded trinucleotide has found its position alongside the template, it will not evaporate again. Because this assumption is ~~probably~~ not valid, we must make a correction which is represented by the second term,  $\tau_2$ .

## 4 Computation of $\tau_2$ em

First we shall ~~now~~ compute this second term,  $\tau_2$ . This computation <sup>is</sup> ~~is~~ based on the fact that (because of reevaporation of the loaded trinucleotides, which ~~are~~ <sup>always</sup> reversibly combined with the anti-code words of the paragene) there will be ~~no~~ <sup>em</sup> matter how long we wait ~~always~~ a certain number of code words (sequences of three nucleotides) on the paragene which are not "covered".

We have to deal with a reversible reaction which may be written as follows:

10/12  
1/2  
x32  
Quack

minus?

6

~~code word + loaded trinucleotide~~  $\rightleftharpoons$  ~~code word - loaded trinucleotide~~  
The rate at which the reaction proceeds from left to right, i.e., the number of successful hits per unit time, is given by

qk #41

Hit rate =  $A f$

qk #41

where  $f$  denotes the concentration of the particular kind of loaded trinucleotide in mol/cm<sup>3</sup>. We shall, for the sake of simplicity, assume that the concentration in the cell of each kind of loaded trinucleotide is the same

$A$  stands for is defined by

$A = 6 \times 10^{23} v \sigma p$

$6 \times 10^{23} v \sigma p$

run eq. not flush right (1)

where  $v$  is the molecular velocity,  $v = \sqrt{\frac{2RT}{\pi M}}$ , where  $M$  is the molecular weight of the loaded trinucleotide. If  $M \approx 1000$ , we have  $v \approx 5 \times 10^3$  cm/sec.

$\sqrt{\frac{2RT}{\pi M}}$

$\sigma$  is the target area that must be hit by the loaded trinucleotide if hydrogen bonding is to take place with the three adjacent nucleotides on the paragene. We may assume for  $\sigma$  a value of  $10^{-15}$  cm<sup>2</sup>.

$p$  denotes the probability that the loaded trinucleotide, when hitting the code word, is in just the right geometrical position to permit hydrogen bonding to take place between the three complementary pairs of bases that are involved. We may take for  $p$ , as a very rough estimate,  $1/300$ ; This would give for  $\sigma p$  the value  $1/3 \times 10^{-13}$ .

3-piece mech. built

The rate at which the reaction proceeds from left to right is given by the rate  $\alpha$  at which the ~~code word - loaded trinucleotide~~ complex dissociates or, as we may also say, the rate at which the loaded trinucleotide evaporates from the paragene. For this rate we may write

minus?

(2)  $\alpha =$  (rate of evaporation)  $\approx 10^{13} e^{-\frac{\Delta H}{RT}}$ , where  $\Delta H$  is the binding energy for the loaded trinucleotide.

(4) In equilibrium the hit rate and evaporation rate must be equal.

The equilibrium constant,  $K$ , of this reversible reaction is defined as the concentration of the loaded trinucleotides at which the code word is covered half of the time. Thus we may write

(3)  $AK \approx \frac{1}{2} \times 10^{13} e^{-\frac{\Delta H}{RT}} = \frac{1}{2} \alpha$   
 or  $\alpha = 2AK$

In equilibrium the probability,  $f$ , for a given code word on the paragene not being covered by the proper loaded trinucleotide is given by

(4)  $f = \frac{1}{1 + \frac{P}{K}}$

Accordingly, in equilibrium, the total number of such gaps along the paragene which contains  $m$  nucleotides is given by

(5) "Number of gaps"  $= \frac{m}{1 + \frac{P}{K}}$

We shall assume that most code words are "covered" in equilibrium, and this means that

(6)  $\frac{P}{K} \gg 1$

We presume that after such equilibrium is established, a chemical reaction chain is somehow triggered, and, moving down along the paragene, links adjacent amino acids into a polypeptide. The average time,  $\tau_1$ , needed for the formation of the polypeptide from the amino acids assembled along the paragene is given by the product of the "number of gaps" that

Operator: in centered formulas from here on, letters incl unless marked rom  
 or  
 Flush  
 on separate line (lines up with (3))

Operator: cap K ital throughout

act #49

act #4

act #4

have to be filled consecutively, and the average time,  $\frac{1}{AP}$ , that it takes to fill one given gap. Thus, for  $\tau_2$ , we may write

(7)  $\tau_2 = \frac{1}{AP} \frac{m}{1 + \frac{P}{K}}$

Computation of  $\tau_1$  and  $\tau_0 = \tau_1 + \tau_2$

When the polypeptide is formed and leaves the paragne, the code words are covered with the denuded trinucleotides. We may now compute the average time,  $\tau_1$ , needed for the evaporation of all the denuded trinucleotides and the assembling of all the loaded trinucleotides in their place. We shall, for the sake of simplicity, assume that a denuded trinucleotide evaporates at the same rate,  $\alpha$ , as a loaded trinucleotide. The rate,  $\alpha$ , at which a loaded trinucleotide evaporates from the template is given by

equation (3), and we may also write this in the form

(8)  $\alpha = AP \frac{2K}{P}$

Since we have assumed  $\frac{P}{K} \gg 1$ , we have

(9)  $AP \gg \alpha$

As may be shown, for very large values of m,  $(\ln \frac{m}{2} \gg 1)$ , we may write for  $\tau_1$

(10)  $\tau_1 \approx \frac{1}{AP} \frac{P}{2K} \ln m$

For the total time,  $\tau_0 = \tau_1 + \tau_2$ , we thus obtain

(11)  $\tau_0 = \frac{1}{AP} \left\{ \frac{m}{1 + \frac{P}{K}} + \frac{1}{2} \frac{P}{K} \ln m \right\}$

If we wish to make this time as small as possible, we have to choose K so as to have, for  $\frac{P}{K}$ ,

(12)  $\frac{P}{K} \approx \sqrt{\frac{2m}{\ln m}}$

(5) See Appendix B. Transpose to p. 19

Substituting this value into <sup>equation</sup> (11) gives

(13) 
$$\tau_0 \approx \frac{\sqrt{2}}{A \rho} \sqrt{m \ln m}$$

*300* *from  $\sqrt{m \ln m}$*

For a polypeptide containing 1,000 amino acid residues, i.e., a paragene containing about 300 code words, we may write  $m = 300$ , and thus we obtain from <sup>relations</sup> (12) and (13)

(14) 
$$\frac{\rho}{K} \approx 10$$

and

(15) 
$$\tau_0 \approx \frac{50}{A \rho}$$

*300*

Estimates for the Values of  $\sigma_p$ ,  $A$ ,  $\alpha$ , and  $K$ . *cm*

As we may see from <sup>relation</sup> (15), we obtain  $\tau_0 = 1$  if  $A \rho = 50$ . This means that for this particular value of  $A \rho$  one enzyme molecule is produced per paragene per second. As we have seen before, this is the order of magnitude of the rate at which highly induced bacteria produce the enzyme  $\beta$ -galactosidase per paragene.

We shall, therefore, in the following, assume  $A \rho = 50$ , and compute from it  $\rho$ , the concentration at which the different kinds of trinucleotides may be present in the cell.

If we use for  $\sigma_p$  the value ~~of  $\frac{1}{3} \times 10^{-17} \text{ cm}^2$~~   $\frac{1}{3} \times 10^{-17} \text{ cm}^2$ , then we obtain from <sup>equation</sup> (1)  $A = 10^{10}$ , and, accordingly, we have  $\rho = 5 \times 10^{-9} \text{ mol/cm}^3$  ( $\rho = 5 \times 10^{-6} \text{ mol/liter}$ ).

*3-piece each. built* *times* *cm<sup>2</sup>* *at no period*

It might well be, however, that  $\sigma_p$  is ten times higher, so that we have  $\sigma_p = \frac{1}{3} \times 10^{-16} \text{ cm}^2$ , and then we obtain from <sup>equation</sup> (1)  $A = 10^{11}$ , so that we have  $\rho = 5 \times 10^{-10} \text{ mol/cm}^3$  or  $5 \times 10^{-7} \text{ mol/liter}$ .

*3-piece each. built* *times* *10<sup>-16</sup> cm<sup>2</sup>* *times* *times*

Thus the concentration,  $\rho$ , of the different kinds of trinucleotides in the cell is likely to be between  $5 \times 10^{-6}$  and  $5 \times 10^{-7} \text{ mol/liter}$ , and one would have to look for concentrations of this order of magnitude in order to obtain experimental confirmation of their presence.

*times*

Since we have assumed  $\frac{\rho}{K} \approx 10$ , and since we believe that we have  $\Delta f = 50$ , it follows that we must have, for the rate of evaporation of the trinucleotides from the parogene, \*

$$(16) \quad \alpha = 2AK = 10 \frac{\#}{\text{sec}} \quad \text{or} \quad \frac{1}{\alpha} = 1/10 \text{ sec.}$$

10 sec<sup>-1</sup>      #      -1      or      #

For values of  $A$  between  $10^{10}$  and  $10^{11}$ , we have

$$(17) \quad 5 \times 10^{-11} < K \text{ (in mol/cm}^3\text{)} < 5 \times 10^{-10}$$

times      cm<sup>3</sup>      times

$$5 \times 10^{-8} < K \text{ (in mol/liter)} < 5 \times 10^{-7}$$

times      or

From (5) we obtain for a given pair of values  $A$  and  $K$  the binding energy  $\Delta H$  between the loaded trinucleotide and the code word on the parogene. For a value of  $A = 10^{10}$  and  $K = 10^{-10}$  mol/cm<sup>3</sup> ( $K = 10^{-7}$  mol/liter), we obtain  $\Delta H$  18,000 calories. Since <sup>six</sup> hydrogen bonds are involved, this would mean about 3,000 calories per hydrogen bond.

For the same value of  $A$  and a value of  $K$  which is ten times larger,  $\Delta H$  would decrease by about 1,400 calories.

### Conclusion. <sup>cm</sup>

These considerations show that the theory which we postulated should be able to explain the high rate of enzyme synthesis which one observes in bacteria when the rate of formation of an enzyme is maximally enhanced by the use of an inducer. The basic thought of this theory consists in the assumption that trinucleotides and possibly also tetranucleotides read the code of the parogene, and that these oligonucleotides carry amino acids. One particular model for protein synthesis, which assumed that each trinucleotide or tetranucleotide carries a sequence of three or four amino acids, respectively, was singled out for detailed discussion because this model appeared to be the most plausible. This does not mean, however, that other models can not be considered.



For instance, rather than ~~to~~ assume that each kind of trinucleotide carries a particular sequence of three amino acids, one might wish to explore the possibility that each trinucleotide might carry only one amino acid. <sup>(2)</sup> In this case the amino acid might be carried by a phosphate group linked by an oxygen atom (ester linkage), either to the third or the fifth carbon of the 5 carbon sugar of either the leading or the trailing nucleotide. Assuming <sup>20</sup> ~~twenty~~ different trinucleotides, each carrying one particular amino acid, we could have a code that requires no commas, with no restrictions imposed on the possible amino acid sequences of the proteins formed by the paragenes.

However, the particular model for protein synthesis here considered cannot be modified by simply saying that each trinucleotide carries one amino acid instead of carrying three amino acids, for adjacent amino acids would not then be at the right distance from each other to be linked into a polypeptide. Such an alternate model for protein synthesis would, therefore, require additional ideas concerning the formation of the polypeptide.

12 pt #

I am grateful for the discussion of a variety of problems arising out of the considerations here presented which I had at the University of Chicago with Prof. <sup>essor</sup> Herbert S. Anker, Dr. Nandor L. Balazs, Mr. Hirono Kulci, Prof. <sup>essor</sup> Joseph E. Mayer, and Prof. <sup>essor</sup> Leonard J. Savage.

Appendix follows

a/?

set APPENDIX in 9 pt. type

MDL  
6/14/57

WPCo.:

The formulas on  
pp. 16-18  
could be set up in  
10 pt (the same  
as regular formulas  
in text) if it would  
make assembly or  
composition easier  
(just the formulas -  
text should still  
be 9/11)

Jan X30

#8a  
9/11 caps ch X30 6

APPENDIXES

9/11 calc ital ch X30  
(Added June 17, 1957) 6

9/11 circ #8 APPENDIX A 4

135

If, in addition to ribose trinucleotides, carrying three amino acids, we have in the role of anti-code words also ribose tetranucleotides, carrying four amino acids, then we must postulate that only three of the bases of the tetranucleotide pair with complementary bases on the paragene. The fourth base, if it is a purine, must pair with the wrong pyrimidine, and if it is a pyrimidine, it must pair with the wrong purine. The reason for this is as follows:

Qk #49

10 sec<sup>-1</sup>

As formula (17) shows,  $\alpha$ , the rate of evaporation of a trinucleotide from the paragene, can be estimated to be about 10<sup>#</sup> sec<sup>-1</sup>, corresponding to a binding energy,  $\Delta H$ , of about 18,000 calories, or 3,000 calories per hydrogen bond. If all four bases of the tetranucleotide were paired with the complementary bases on the paragene, we would then have two more hydrogen bonds and presumably an additional binding energy of about 6,000 calories. The equilibrium constant,  $K$ , of the tetranucleotide would therefore be lower than the value computed for  $K$  of the trinucleotides by a factor of about 10<sup>4</sup>, and the rate of evaporation,  $\alpha$ , would be lower by the same factor. This would then make  $\tau_0$ , the minimal time it takes for a paragene to form a polypeptide, too long to be compatible with the observed rate of production of  $\beta$ -galactosidase in highly induced bacteria.

The number of words that may be constructed in a mixed three-letter-word and four-letter-word code of the kind described above, ~~we~~ <sup>where</sup> we demand that the code require no commas, has so far not been determined.

Qk #30 Qk #8a APPENDIX B 6

The problem of computing  $\tau_1$  amounts to the solving of the following problem:

There are  $n$  boxes, each of which can hold one white ball or one black ball. Initially, at time  $\lambda = 0$ , all these boxes contain one

Qk #35 zero throughout

white ball -- a stripped trinucleotide. These white balls evaporate at the rate  $\alpha$ , so that at time  $\lambda$  the probability,  $W(\lambda)$ , of having no white ball in the box is given by

(18)  $W(\lambda) = 1 - e^{-\alpha\lambda}$ .  
Operator: letters from here on unless marked otherwise. (gh#35) (gh#49)

If the rate at which black balls fall into an empty box is designated by  $\beta$ , then the probability,  $y$ , that a given box does not contain a black ball at time  $t$ , is given by

(19)  $y(t) = e^{-\beta t} + \int_{\lambda=0}^{\lambda=t} e^{-\beta(t-\lambda)} \frac{dW}{d\lambda} d\lambda$ .  
letters ital (24#1) (gh#26)

We may write

(20)  $\int_{\lambda=0}^{\lambda=t} e^{-\beta(t-\lambda)} \frac{dW}{d\lambda} d\lambda = \frac{\alpha}{\beta-\alpha} [e^{-\alpha t} - e^{-\beta t}]$ ,  
(24#1)

so that we have

(21)  $y(t) = \frac{1}{\beta-\alpha} (\beta e^{-\alpha t} - \alpha e^{-\beta t})$ .  
start commas

Thus we may write, for the probability,  $x$ , that a given box contains a black ball,

(22)  $x(t) = 1 - y(t) = 1 - \frac{\beta}{\beta-\alpha} e^{-\alpha t} + \frac{\alpha}{\beta-\alpha} e^{-\beta t}$ .

9/11/77

dt x30

dt x30

flush

dt x30

dt x30

dt x30

eg. nos flush right

From this we obtain for  $P(t)$ , the probability that all  $m$  boxes contain one black ball,

(23) 
$$P(t) = x(t)^m = \left( 1 - \frac{\beta}{\beta - \alpha} e^{-\alpha t} + \frac{\alpha}{\beta - \alpha} e^{-\beta t} \right)^m$$

The average time,  $\tau_1$ , needed for the evaporation of the white balls from all  $m$  boxes and the filling of all  $m$  boxes with black balls is given by

(24) 
$$\tau_1 = \int_0^{\infty} t \frac{dP}{dt} dt$$

If  $\beta \gg \alpha$ , we may write

(25) 
$$\tau_1 = \int_0^{\infty} t \frac{d}{dt} (1 - e^{-\alpha t})^m dt$$

The expression  $\frac{d}{dt} (1 - e^{-\alpha t})^m$  has a maximum at some value of  $t = \tau_0$ , and for large values of  $m$  it becomes small very rapidly both below and above  $t = \tau_0$ . Therefore, if  $m$  is large, we may write

(26) 
$$\tau_1 \approx \tau_0$$

We obtain  $\tau_0$  by writing

(27) 
$$\left\{ \frac{d^2}{dt^2} (1 - e^{-\alpha t})^m \right\}_{t=\tau_0} = 0$$

and from this we obtain

(28) 
$$\tau_0 = \frac{1}{\alpha} \ln m$$

from  
ln

Thus, for  $\beta \gg \alpha$  and  $m \gg 1$ , we may write (26) as

(29)  $\tau_1 \approx \frac{1}{\alpha} \ln m$   
*from equations*

It can be shown from (23) and (24) that in the next higher approximation we have

(30)  $\tau_1 \approx \frac{1}{\alpha} \{ \ln(m+1) + C \} + \Delta$   
*ln from cap* *where*  $\left\{ \begin{array}{l} C = 0.577 \dots \text{(Euler's const.)} \\ \frac{1}{\beta} < \Delta < \frac{1}{\beta - \alpha} \end{array} \right.$

APPENDIX C

F. H. Crick, J. S. Griffith, and L. E. Orgel write in the May issue of the *Proceedings of the Nat. Acad. Sci.* (vol. 43, pp. 419 and 420, 1957):

To fix ideas, we shall describe a simple model to illustrate the advantages of such a code. Imagine that a single chain of RNA, held in a regular configuration, is the template. Let the intermediates in protein synthesis be 20 distinct molecules, each consisting of a trinucleotide chemically attached to one amino acid. The bases of each trinucleotide are chosen according to the code given above. Let these intermediate molecules combine, by hydrogen bonding between bases, with the RNA template and there await polymerization. Now imagine that such an amino acid-trinucleotide were to diffuse into an incorrect place on the template, such that two of its bases were hydrogen-bonded, though not the third. We postulate that this incomplete attachment will only retain the intermediate for a very brief time (for example, less than 1 millisecond) before the latter breaks loose and diffuses elsewhere. However, when it eventually diffuses to the correct place, it will be held by hydrogen bonds to all three bases and will thus be retained, on the average, for a much longer time (say, seconds or minutes). Now the code we have described insures that this more lengthy attachment can occur only at the points where the intermediate is needed. If one of the 20 intermediates could

9/11 #8a  
 (29)  
 (30)  
 7/11 #8a #8a #8a  
 opening  
 quote

equation  
 18

from  
 10  
 4

9/11  
#8a

stay for a long time on one of the false positions, it would effectively block the two positions it was straddling and hold up the polymerization process. Our code makes this impossible. This scheme, therefore, allows the intermediates to accumulate at the correct positions on the template without ever blocking the process by settling, except momentarily, in the wrong place. It is the <sup>is</sup> feature which gives it an advantage over schemes in which the intermediates are compelled to combine with the template one after the other in the correct order.

The example given here is only for illustration, but it brings out the physical idea behind the concept of a comma-less code.

In passing, it should be mentioned that while the idea of making three nonoverlapping nucleotides code for one amino acid at first sight entails certain stereochemical difficulties, these are not insuperable if it is assumed that the polypeptide chain, when polymerized, does not remain attached to the template. A detailed scheme along these lines has been described to us by Dr. S. Brenner (personal communication).

closing quote

6 pt #

insert notes from pp. 1, 2, 3, 6, 11

#

June 7, 1957

*Complete*  
~~with some parts~~

HOW MAY AMINO ACIDS READ THE NUCLEOTIDE CODE?

by Leo Szilard

(Submitted by Joseph E. Mayer)

The Enrico Fermi Institute for Nuclear Studies  
The University of Chicago, Chicago, Illinois

It is now generally believed that proteins are formed along ~~side~~ nucleic acid templates. The sequence of purine and pyrimidine bases in the template is supposed to represent a code that may somehow determine the sequence of the amino acids in the particular polypeptide (protein) that a given template will form. <sup>(1)</sup> The purine and pyrimidine bases of the template, the letters of the code, are adenine, uracil, guanine and cytosine if the template be an RNA molecule; and if the template be a DNA molecule, thymine takes the place of uracil.

It has remained so far ~~an unsolved~~ <sup>somewhat of a</sup> mystery in just what ~~way~~ way amino acids could read such a code. In what manner can chemical forces of the kind we know to exist -- line up amino acids alongside such a template in the proper sequence and at the proper distance from each other, so that ~~there might~~ a chemical reaction chain may link adjacent amino acids through peptide bonds with each other?

It is the purpose of the present paper to indicate a conceptually simple scheme that will -- at least by way of an example -- illustrate in what manner this might take place in the living cell.

---

(1) A. L. Dounce, *Enzymologia*, 15, 251 (1952)  
G. Gamow, *Nature*, Vol. 173, p. 318 (1954).



Because a template which synthesizes protein need not necessarily be the gene itself but must carry the same information as the corresponding gene, we shall here refer to such a template for the sake of brevity as a paragene.

The basic thought underlying the scheme here presented consists in the assumption that there are a number of enzymes (or enzyme systems) -- perhaps twenty altogether -- in the cell, and that each of these catalyzes the formation of a particular trinucleotide which carries either one particular amino acid <sup>(2)</sup> or, more likely perhaps, a particular sequence of three amino acids. If the amino acid is carried by the nucleotide on a phosphate or pyrophosphate group as an acid anhydride -- which is a high energy compound -- then the energy needed for the formation of the peptide bonds will become free when the amino acid is split off. In this sense one can say that each amino acid may carry the energy needed for forming its peptide bond.

According to the notions here presented, amino acids can not read <sup>themselves</sup> ~~themselves~~ the code of the paragene. But the trinucleotides, which carry the proper amino acids, may attach with their three bases through the formation of 6 hydrogen bonds to the proper sequence of three bases on the paragene, and thus the amino acids may be lined up in the proper sequence along the paragene.

Accordingly, sequences of three nucleotides along the paragene represent the code words, and the trinucleotides which carry the amino acids represent the anti-code words. We assume that these anti-code words are complementary to the code words in the sense that, where the code word contains adenine the anti-code word contains uracil (or thymine), where the code word contains uracil (or thymine) the anti-code word contains adenine; and similarly guanine corresponds to cytosine and cytosine corresponds to guanine. The rationale for this assumption is as follows:

(2) Note added June 17th. The May issue of the Proc. Nat. Acad. Sci., Vol. 43, p. 416 (1957) which was belatedly received here, contains an article by F.H. Crick, J.S. Griffith, and L.E. Orgel which is in part identical with their previously circulated memorandum referred to in the text. In addition, however, the authors discuss -- in conjunction with a detailed oral communication of S. Brenner to them -- how amino acids might read the nucleotide code. The relevant passage of their paper is quoted in full in Appendix C.

The concept of code-letter and complementary code-letter arose originally from the interpretation of the structure of DNA given by J.D. Watson and F.H.C. Crick. <sup>(2)</sup> They showed that in a double stranded DNA

---

<sup>(2)</sup> Nature, Vol. 173, p. 318 1955.  
Proc. Roy. Soc., Vol. 223, p. 80, 1954.

---

structure, adenine of one strand pairs with thymine of the other strand (which presumably plays the same role in DNA as does uracil in RNA) and similarly guanine pairs with cytosine. The helical structure of DNA permits just such pairing, and hydrogen bonding is possible between adenine and thymine, as well as between guanine and cytosine.

If the sequence of bases along one strand of DNA represents a coded message which consists in three letter-words then, because we have four letters to choose from such a message could utilize 64 different words. We might, however, be ~~restricted~~ <sup>limited</sup> to the use of 20 out of the 64 words that are available. The reasons for this ~~restriction~~ <sup>limitation</sup> would be as follows:

If all 64 possible three-letter combinations form ~~in fact~~ a code word, and if the parogene assumes at the time of the formation of the polypeptide a helical configuration similar to the helical configuration of DNA, then it follows that the code on the parogene must be read consecutively from one end -- say, from the "head" of the parogene downward. This is so because such a helical structure does not provide for commas between the individual code words, and in a code containing 64 words any three consecutive letters form a word. If we number the letters along the parogene, ~~say~~ from the head of the parogene downward, then the first three letters, the letters 1, 2, 3, form a word which was meant to be conveyed and so do the next three letters, the letters 4, 5, 6. But sequences of three letters which encroach on two adjacent words (such as 2, 3, 4 or 3, 4, 5, for example) form code words which are not meant to be conveyed.

In these circumstances, the code would be misread if the trinucleotides, which represent the anti-code words, assemble alongside the paragene simultaneously, rather than -- from one end on -- consecutively. If we want to have simultaneous assembly of the trinucleotides alongside the comma-less paragene, then we are restricted to 20 code words.

The notion of such a 20-word code, which needs no commas, was introduced by F.H.C.Crick, J.S.Griffith, and L.E.Orgel of the Medical Research Council Unit at the Cavendish Laboratory, Cambridge, in a memorandum circulated in May, 1956 among workers interested in the subject of protein synthesis.<sup>(?)</sup> From such a code we must demand that the letters 1, 2, 3 on the template form a code word, and the letters 4, 5, 6 also form a code word, but sequences of three letters, which encroach on two adjacent words (such as 2, 3, 3 or 3, 4, 5, for example) form no code word. Crick and his co-workers have shown that this demand can be met, that a code which requires no commas may be constructed and that it can accommodate 20 three-letter code words.

We shall now single out for more detailed examination one conceivable model for protein synthesis which might provide for the lining up of the amino acids alongside the paragene, both in the proper order and at the proper distance from each other. This particular model is based on the following assumptions:

The trinucleotides which form the anti-code words contain the sugar ribose rather than the sugar desoxyribose. Each particular ribose trinucleotide (the anti-code word) carries a particular sequence of three amino acids. A phosphate (or diphosphate) group is attached to the (2) carbon atom of the ribose moiety of each nucleotide and an amino acid is attached to each of these phosphate (or diphosphate) groups. The amino acid anhydrides ~~contain~~<sup>contain</sup> an energy-rich P, or PP, bond which, when split, may release 12,000 or 16,000 calories, respectively.

During protein synthesis the nucleic acid strand that functions as a template (the parogene) may take up -- so we here assume -- a helical configuration resembling the helical configuration of a B-DNA strand in the double stranded DNA helix. The trinucleotides may then line up alongside the helical parogene with their purine and pyrimidine bases paired with the complementary bases of the parogene, and if they are so lined up, then the amino acids carried by the trinucleotides may come to lie at just about the right distance from each other to permit the formation of a peptide bond between adjacent amino acids. A chemical reaction chain -- starting perhaps from the head of a parogene -- may then move down along the parogene, split the acid anhydrides, and thus free the amino acids as well as make available the energy needed for the formation of peptide bonds between adjacent amino acids.

Adjacent amino acids can be linked only if the distance from each other is smaller or equal but not appreciably larger than the fundamental repeating distance in a polypeptide chain which is  $7.27\text{\AA}$ . The fundamental repeating distance in a fully extended polypeptide chain is about  $7\text{\AA}$ . Since before they are linked into a polypeptide, the amino acids can rotate around the chemical bond which ties them to the phosphate group, <sup>(at the (2) carbon atom of the ribose)</sup> they might well be assembled along the parogene at a <sup>somewhat</sup> smaller distance from each other than the fundamental repeating distance of the polypeptide chain.

Applying the concept of a 20-letter code that requires no commas to our particular model of protein synthesis, we may now say the following:

Each of the 20 amino acids may appear once attached to the leading letter and once attached to the trailing letter of the 20 (trinucleotide) anti-code words. Therefore, among the polypeptides that can be formed, each amino acid may precede any other amino acid, and each amino acid may follow any other amino acid. This does not, of course, mean that any amino acid sequence is possible.

Some of the amino acid sequences that may be found experimentally in sequential analysis of proteins and polypeptides might show that the restrictions imposed by our model on the possible amino acid sequence are too severe, and that the model has to be modified to accommodate established facts. As will be presently seen some modification of our model may be indicated for other reasons also. There is certainly no inherent reason why we should have a pure three-letter code -- as we have assumed above -- and why, for instance, a certain number of four-letter code words should not be utilized also.<sup>(4)</sup>

If we had a pure three-letter code, we would have to demand that the number of amino acid residues of all polypeptides or proteins synthesized in the manner described above should be a multiple of three.

So far the number of amino acid residues in <sup>non-cyclic</sup> polypeptides and proteins have been determined only within one rather special class; all of them represent substances which are secreted by mammalian tissues. The number of amino acid residues found in such polypeptides and proteins, which have been analyzed with adequate accuracy, are as follows:

a) Insulin chain A: 21; insulin chain B: 30; corticotropin B: 39; oxytocin: 9; vasopressin: 9; Intermedin B: 18.

All of these would fit a pure three-letter code.

b) Intermedin A: 13; glucagon: 29 and pancreas ribonuclease: 124.

These do not fit a pure three-letter code. Intermedin and ribonuclease would have to include at least one 4-letter word and glucagon at least two 4-letter words. ~~Therefore, we can conclude that the conclusion that the code is a pure three-letter code is not valid.~~

(4) See Appendix A

## Observed rate of enzyme synthesis

According to the notions here adopted most enzymes are synthesized in growing bacteria at a rather low rate which does not represent the maximum synthesizing capacity of the corresponding paragenes. The rate of production of a given enzyme may, however, be greatly enhanced if the enzyme is induced, and what we are interested to learn is the maximal rate at which a paragene may be able to form the corresponding enzyme.

One of the most studied cases of enzyme induction is the induction of the enzyme  $\beta$ -galactosidase which splits lactose. Jacques Monod and his co-workers have shown that the production rate of this enzyme in bacteria can be greatly enhanced by certain chemical analogues of lactose, which act as inducers, and that the rate of production of the enzyme goes up almost instantaneously upon adding such an inducer to the medium. We are thus led to believe that the inducer may act by increasing the rate at which one template produces the enzyme rather than by increasing the number of templates that produce the enzyme at an unchanged rate.

In fully induced wild type bacteria growing in minimal medium this enzyme is contained in the amount of about 8 mgm. per  $10^{12}$  bacteria and thus amounts to about 8% of the total proteins. We obtain the rate at which this enzyme is produced in minimal medium per bacterium by dividing the amount contained in one bacterium by 1.44 times the doubling time (40 minutes) of the bacterium. We thus find for the rate, at which this enzyme is produced in fully induced wild type bacteria growing in minimal medium, about  $2 \cdot 10^{18}$  grams per cell per second.

If we assume a molecular weight of a millicion (Jacques Monod and Melvin Cohn estimate the molecular weight of this enzyme at about 800,000), we obtain a rate of 1.5 enzyme molecules per cell per second. The number of paragenes per cell is not known. ~~There might be a few paragenes~~ There might be a few paragenes present per cell rather than just one, and the number of paragenes might be of the order of magnitude of 10. On the other hand, smaller enzyme molecules might be synthesized somewhat faster than larger enzyme molecules.

On the basis of the figure given above, we are thus led to believe ~~in bacteria --~~ that when an enzyme is fully induced and enzyme synthesis proceeds at its maximal rate -- the rate of formation of the enzyme may be of the order of magnitude of one per second per paragene.

### Computed rate of enzyme synthesis

We shall now attempt to compute at what rate a parogene may be able to synthesize the corresponding enzyme on the basis of the model that we have postulated. For the purposes of this computation, we shall assume that the molecular weight of the enzyme is about 100,000. We then have about 1,000 amino acid residues in the enzyme, and accordingly we would have to assemble alongside the parogene  $m = 300$  trinucleotides, each of which is "loaded" with three amino acids.

In the approximation which we shall ~~consider~~<sup>use</sup>, the minimum time,  $\tau_0$ , needed for the formation of the polypeptide is composed of two terms,  $\tau_1$  and  $\tau_2$ ;

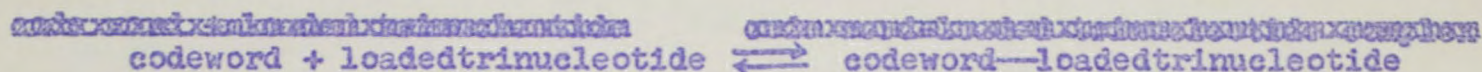
$$\tau_0 = \tau_1 + \tau_2$$

After all the amino acids, assembled alongside the template, have been linked into one polypeptide, and ~~assuming that~~<sup>if</sup> this polypeptide is at once removed, a certain time,  $\tau_1$ , will elapse until the trinucleotides, which are now denuded of amino acids, evaporate from the template and their place is taken by trinucleotides which are loaded with the proper amino acids. We ~~shall~~ assume that the concentration of denuded trinucleotides in the cell is very small so that after the denuded trinucleotides evaporate, the loaded trinucleotides do not have to compete with denuded trinucleotides for their legitimate positions along the parogene. The time,  $\tau_1$ , which is necessary to permit evaporation of  $m$  denuded trinucleotides and to assemble  $m$  loaded trinucleotides in their place we shall compute here on the assumption that once a loaded trinucleotide has found its position alongside the template, it will not evaporate again. Because this assumption is ~~not~~ not valid, we must make a correction which is represented by the second term,  $\tau_2$ .

### Computation of $\tau_1$

First we shall now compute this second term,  $\tau_2$ . This computation ~~will~~<sup>is</sup> be based on the fact that (because of reevaporation of the loaded trinucleotides, which ~~are~~ reversibly combined with the anti-code words of the parogene) there will be - no matter how long we wait - always a certain number of code words (sequences of three nucleotides) on the parogene which are not "covered".

We have to deal with a reversible reaction which may be written as follows:



The rate at which the reaction proceeds from left to right; i.e. the number of successful hits per unit time is given by

$$\text{hit rate} = A \rho$$

$\rho$  denotes the concentration of the particular kind of loaded trinucleotide in mol/cc. We shall, for the sake of simplicity, assume that the concentration in the cell of each kind of loaded trinucleotide is the same.

A stands for

$$(1) \quad A = 6 \cdot 10^{23} v \sigma p$$

$v$  is the molecular velocity,  $v = \sqrt{\frac{2RT}{\pi M}}$  where  $M$  is the molecular weight of the loaded trinucleotide. If  $M \approx 1000$ , we have  $v \approx 5 \times 10^3$  cm/sec.

$\sigma$  is the target area that must be hit by the loaded trinucleotide if hydrogen bonding is to take place with the three adjacent nucleotides on the paragene. We may assume for  $\sigma$  a value of  $\sigma = 10^{-15}$  cm<sup>2</sup>.

$p$  denotes the probability that the loaded trinucleotide, when hitting the code word, is in just the right geometrical position to permit hydrogen bonding to take place between the three complementary pairs of bases that are involved. We may take for  $p$ , as a very rough estimate  $p = 1/300$ ; This would give for  $\sigma p$  the value  $\sigma p = 1/3 \cdot 10^{-13}$ .

The rate at which the reaction proceeds from left to right is given by the rate  $\alpha$  at which the ~~codeword-loaded trinucleotide~~ complex dissociates or, as we may also say, the rate  $\alpha$  at which the loaded trinucleotide evaporates from the paragene. For this rate we may write

$$A = 6 \cdot 10^{23} \cdot 5 \cdot 10^3 \cdot 10^{-15} \cdot \frac{1}{300} \cdot \frac{1}{300}$$

$$10^{22} \cdot 10^3 \cdot 10^{-15} = 10^{+10}$$



$$(2) \quad \alpha = \text{rate of evaporation} \approx 10^{13} e^{-\frac{\Delta H}{RT}}$$

where  $\Delta H$  is the binding energy for the loaded trinucleotide.

In equilibrium the hit rate and evaporation rate must be equal.

The equilibrium constant,  $K$ , of this reversible reaction is defined as the concentration of the loaded trinucleotides at which the code word is covered half of the time. Thus we may write

$$(3) \quad AK \approx \frac{1}{2} 10^{13} e^{-\frac{\Delta H}{RT}} = \frac{1}{2} \alpha$$

$$\alpha = 2AK$$

Handwritten notes and calculations:

$\alpha = 2AK$

$\alpha = 2 \cdot 10^{-4}$

$\alpha = 2 \cdot \frac{A \cdot P \cdot K}{P} = 10^{13} e^{-\frac{\Delta H}{RT}}$

$\frac{P}{K} = 10^4$

$\frac{1 \text{ sec}}{10^{-4}} = 10^4$

$\frac{N_{tr}}{N_{total}}$

In equilibrium the probability,  $f$ , for a given code word on the parogene not being covered by the proper loaded trinucleotide is given by

$$(4) \quad f = \frac{1}{1 + \frac{P}{K}}$$

Accordingly, in equilibrium, the total number of such gaps along the parogene which contains  $m$  nucleotides is given by

$$(5) \quad \text{"number of gaps"} = \frac{m}{1 + \frac{P}{K}}$$

We shall assume that most code words are "covered" in equilibrium and this means that

$$(6) \quad \frac{P}{K} \gg 1$$

We presume that after such equilibrium is established a chemical reaction chain is somehow triggered, and, moving down along the parogene, links adjacent amino acids into a polypeptide. The average time,  $\tau_1$ , needed for the formation of the polypeptide from the amino acids assembled along the parogene is given by the product of the "number of gaps", that

have to be filled consecutively, and the average time,  $\frac{1}{Ap}$ , that it takes to fill one given gap. Thus, for  $\tau_2$  we may write

$$(7) \quad \tau_2 = \frac{1}{Ap} \frac{m}{1 + \frac{P}{K}}$$

Computation of  $\tau_1$  and  $\tau_0 = \tau_1 + \tau_2$

When the polypeptide is formed and leaves the paragne, the code words are covered with the denuded trinucleotides. We may now compute the average time,  $\tau_1$ , needed for the evaporation of all the denuded trinucleotides and the assembling of all the loaded trinucleotides in their place. We shall, for the sake of simplicity, assume that a denuded trinucleotide evaporates at the same rate,  $\alpha$ , as a loaded trinucleotide. The rate,  $\alpha$ , at which a loaded trinucleotide evaporates from the template is given by (3) and we may also write this in the form

$$(8) \quad \alpha = Ap \frac{2K}{P}$$

Since we have assumed  $\frac{P}{K} \gg 1$ , we have

$$(9) \quad Ap \gg \alpha$$

As may be shown, ~~xx~~ for very large values of  $m$ , ~~xxxxxxxxxxxxxxxx~~ (  $\ln m \gg 1$  ) we may write for  $\tau_1$

$$(10) \quad \tau_1 \approx \frac{1}{Ap} \frac{P}{2K} \ln m$$

For the total time,  $\tau_0 = \tau_1 + \tau_2$  we thus obtain

$$(11) \quad \tau_0 = \frac{1}{Ap} \left\{ \frac{m}{1 + \frac{P}{K}} + \frac{1}{2} \frac{P}{K} \ln m \right\}$$

If we wish to make this time as small as possible, we have to choose  $K$  so as to have for  $\frac{P}{K}$

$$(12) \quad \frac{P}{K} \approx \sqrt{\frac{2m}{\ln m}}$$

(5) See Appendix B

$$\frac{1}{2} \gg 1$$

$$\frac{1}{2} \sqrt{\frac{2m}{\ln m}} \approx \frac{1}{2} \sqrt{\frac{2m}{\ln m}}$$

$$\frac{1}{2} \sqrt{\frac{2m}{\ln m}} \approx \frac{1}{2} \sqrt{\frac{2m}{\ln m}}$$

April

Substituting this value into (11) gives

$$(13) \quad \tau_0 \approx \frac{\sqrt{2}}{A\phi} \sqrt{m \ln m}$$

For a polypeptide containing 1,000 amino acid residues, i.e. a paragene containing about 300 code words, we may write  $m = 300$ , and thus we obtain from (12) and (13)

$$(14) \quad \frac{\phi}{K} \approx 10$$

and

$$(15) \quad \tau_0 \approx \frac{50}{A\phi}$$

#### Estimates for the values of $\sigma\phi$ , $A$ , $\alpha$ and $K$

As we may see from (15) we obtain  $\tau_0 = 1$  if  $A\phi = 50$ . This means that for this particular value of  $A\phi$  one enzyme molecule is produced per paragene per second. As we have seen before, this is the order of magnitude of the rate at which highly induced bacteria produce the enzyme  $\beta$ -galactosidase per paragene.

We shall, therefore, in the following assume  $A\phi = 50$ , and compute from it  $\phi$ , the concentration at which the different kinds of trinucleotides may be present in the cell.

If we use for  $\sigma\phi$  the value of  $\sigma\phi = \frac{1}{3} 10^{-11} \text{ cm}^2$ , then we obtain from (1)  $A$ ;  $A = 10^{10}$  and accordingly we have  $\phi = 5 \cdot 10^{-9} \text{ mol/cc}$  ( $\phi = 5 \cdot 10^{-6} \text{ mol/liter}$ ).

It might well be, however, that  $\sigma\phi$  is ten times higher so that we have  $\sigma\phi = \frac{1}{3} 10^{-10} \text{ cm}^2$ , and then we obtain from (1) for  $A$ ;  $A = 10^{11}$ , so that we have  $\phi = 5 \cdot 10^{-10} \text{ mol/cc}$  or ( $\phi = 5 \cdot 10^{-7} \text{ mol/liter}$ ).

Thus the concentration,  $\phi$ , of the different kinds of trinucleotides in the cell is likely to be between  $5 \times 10^{-6}$  and  $5 \times 10^{-7} \text{ mol/liter}$ , and one would have to look for concentrations of this order of magnitude in order to obtain experimental confirmation of their presence.

$$\alpha = 10^{13} e^{-\frac{x}{kT}} \approx 10^{13} \quad 15****$$

$$e^{13 \times 2.3} \quad \frac{x}{kT} = \frac{173 \times 2.3}{600 \times 2.3 \times 10^3}$$

and since we believe

Since we have assumed  $\frac{p}{k} \approx 10$  that we have  $A_j = 50$ , it follows that we must have for the rate of evaporation of the trinucleotides from the parogene,  $\alpha$

$$(16) \quad \alpha = 2AK = 10/\text{sec.} \quad \text{or} \quad 1/\alpha = 1/10 \text{ sec.}$$

For values of A between  $10^{10}$  and  $10^{11}$ , we have

$$(17) \quad 5 \cdot 10^{-11} < K \text{ (in mol/cc)} < 5 \cdot 10^{-10}$$

or

$$5 \cdot 10^{-8} < K \text{ (in mol/liter)} < 5 \cdot 10^{-7}$$

~~47~~ =

From (3) we obtain for a given pair of values A and K the binding energy  $\Delta H$  between the loaded trinucleotide and the code word on the parogene. For a value of  $A = 10^{10}$  and  $K = 10^{-10}$  mol/cc ( $K = 10^{-7}$  mol/liter), we obtain  $\Delta H$  18,000 calories. Since 6 hydrogen bonds are involved, this would mean about 3,000 calories per hydrogen bond. ( $\Delta F = RT \times 7 \times 2.3$ )

For the same value of A and a value of K which is ten times larger,  $\Delta H$  would decrease by about 1400 calories.

### Conclusion

These considerations show that the theory which we postulated should be able to explain the high rate of enzyme synthesis which one observes in bacteria when the rate of formation of an enzyme is maximally enhanced by the use of an inducer. The basic thought of this theory consists in the assumption that trinucleotides and possibly also tetranucleotides read the code of the parogene, and that these oligonucleotides carry amino acids. One particular model for protein synthesis, which assumed that each trinucleotide or tetranucleotide carries a sequence of three or four amino acids respectively, was singled out for detailed discussion because this model appeared to be the most plausible. This does not mean, however, that other models can not be considered.

For instance, rather than to assume that each kind of trinucleotide carries a particular sequence of three amino acids, one might wish to explore the possibility that each trinucleotide might carry only one amino acid. In this case the amino acid might be carried by a phosphate group linked by an oxygen atom (ester linkage), either to the third or the fifth carbon of the 5 carbon sugar of either the leading or the trailing nucleotide. Assuming twenty different trinucleotides, each carrying one particular amino acid, we could have a code that requires no commas, with no restrictions imposed on the possible amino acid sequences of the proteins formed by the paragenes.

However, the particular model for protein synthesis here considered cannot be modified by simply saying that each trinucleotide carries one amino acid instead of carrying three amino acids, for adjacent amino acids would not then be at the right distance from each other to be linked into a polypeptide. Such an alternate model for protein synthesis would, therefore, require additional ideas concerning the formation of the polypeptide.

I am grateful for the discussion of a variety of problems arising out of the considerations here presented which I had at the University of Chicago with Prof. Herbert S. Anker, Dr. Nandor L. Balazs, Mr. Hirono Kulci, Prof. Joseph E. Mayer, and Prof. Leonard J. Savage.

## APPENDIX

(Added June 17/57)

A.) If, in addition to ribose trinucleotides, carrying three amino acids, we have in the role of anti-code words also ribose tetranucleotides, carrying four amino acids, then we must postulate that only three of the bases of the tetranucleotide pair with complementary bases on the paragene. The fourth base, if it is a purine, must pair with the wrong pyrimidine, and if it is a pyrimidine, it must pair with the wrong purine. The reason for this is as follows:

As formula (17) shows,  $\alpha$ , the rate of evaporation of a trinucleotide from the paragene can be estimated to be about 10/sec., corresponding to a binding energy,  $\Delta H$ , of about 18,000 calories or 3,000 calories per hydrogen bond. If all four bases of the tetranucleotide were paired with the complementary bases on the paragene, we would then have two more hydrogen bonds and presumably an additional binding energy of about 6,000 calories. The equilibrium constant,  $K$ , of the tetranucleotide would therefore be lower than the value computed for  $K$  of the trinucleotides by a factor of about  $10^4$ , and the rate of evaporation,  $\alpha$ , would be lower by the same factor. This would then make,  $\tau_0$ , the minimal time it takes for a paragene to form a polypeptide too long to be compatible with the observed rate of production of  $\beta$ -galactosidase in highly induced bacteria.

The number of words that may be constructed in a mixed three-letter word and four-letter word code of the kind described above, <sup>WHERE</sup> ~~IF~~ we demand that the code require no commas, has so far not been determined.

B.) The problem of computing  $\tau$ , amounts to the solving of the following problem:

There are  $m$  boxes, each of which can hold one white ball or one black ball. Initially, at time,  $\lambda = 0$ , all of these boxes contain one

white ball -- a stripped trinucleotide. These white balls evaporate at the rate,  $\alpha$ , so that at time  $\lambda$  the probability,  $W(\lambda)$  of having no white ball in the box is given by

$$(18) \quad W(\lambda) = 1 - e^{-\alpha\lambda}$$

If the rate at which black balls fall into an empty box is designated by  $\beta$ , then the probability,  $y$ , that a given box does not contain a black ball at time,  $t$ , is given by

$$(19) \quad y(t) = e^{-\alpha t} + \int_{\lambda=0}^{\lambda=t} e^{-\beta(t-\lambda)} \frac{dW}{d\lambda} d\lambda$$

we may write

$$(20) \quad \int_{\lambda=0}^{\lambda=t} e^{-\beta(t-\lambda)} \frac{dW}{d\lambda} d\lambda = \frac{\alpha}{\beta-\alpha} [e^{-\alpha t} - e^{-\beta t}]$$

so that we have

$$(21) \quad y(t) = \frac{1}{\beta-\alpha} (\beta e^{-\alpha t} - \alpha e^{-\beta t})$$

Thus we may write for the probability,  $x$ , that a given box contains a black ball

$$(22) \quad x(t) = 1 - y(t) = 1 - \frac{\beta}{\beta-\alpha} e^{-\alpha t} + \frac{\alpha}{\beta-\alpha} e^{-\beta t}$$

From this we obtain for  $P(t)$ , the probability that all  $m$  boxes contain one black ball

$$(23) \quad P(t) = x(t)^m = \left( 1 - \frac{\beta}{\beta - \alpha} e^{-\alpha t} + \frac{\alpha}{\beta - \alpha} e^{-\beta t} \right)^m$$

The average time,  $\tau_1$ , needed for the evaporation of the white balls from all  $m$  boxes and the filling of all  $m$  boxes with black balls is given by

$$(24) \quad \tau_1 = \int_0^{\infty} t \frac{dP}{dt} dt$$

If  $\beta \gg \alpha$ , we may write

$$(25) \quad \tau_1 = \int_0^{\infty} t \frac{d}{dt} (1 - e^{-\alpha t})^m dt$$

The expression  $\frac{d}{dt} (1 - e^{-\alpha t})^m$  has a maximum at some value of  $t$ ;  $t = \tau_0$  and -- for large values of  $m$  -- it becomes small very rapidly both below and above  $t = \tau_0$ . Therefore, if  $m$  is large we may write

$$(26) \quad \tau_1 \approx \tau_0$$

We obtain  $\tau_0$  by writing

$$(27) \quad \left\{ \frac{d^2}{dt^2} (1 - e^{-\alpha t})^m \right\}_{t=\tau_0} = 0$$

and from this we obtain

$$(28) \quad \tau_0 = \frac{1}{\alpha} \ln m$$



Thus for  $\beta \gg \alpha$  and  $m \gg 1$  we may write (26)

$$(29) \quad \tau_1 \approx \frac{1}{\alpha} \ln m$$

It can be shown from (23) and (24) that in the next higher approximation we have

$$(30) \quad \tau_1 \approx \frac{1}{\alpha} \{ \log(m+1) + C \} + \Delta ; \text{ where } \begin{cases} C = 0.577 \text{ (Euler's const.)} \\ \frac{1}{\beta} < \Delta < \frac{1}{\beta - \alpha} . \end{cases}$$

C.) F. H. Crick, J. S. Griffith and L. E. Orgel write in the May issue of the Proc. Nat. Acad. Sci. (Vol. 43, pp. 419 and 420, 1957):

"To fix ideas, we shall describe a simple model to illustrate the advantages of such a code. Imagine that a single chain of RNA, held in a regular configuration, is the template. Let the intermediates in protein synthesis be 20 distinct molecules, each consisting of a trinucleotide chemically attached to one amino acid. The bases of each trinucleotide are chosen according to the code given above. Let these intermediate molecules combine, by hydrogen bonding between bases, with the RNA template and there await polymerization. Now imagine that such an amino acid-trinucleotide were to diffuse into an incorrect place on the template, such that two of its bases were hydrogen-bonded, though not the third. We postulate that this incomplete attachment will only retain the intermediate for a very brief time (for example, less than 1 millisecond) before the latter breaks loose and diffuses elsewhere. However, when it eventually diffuses to the correct place, it will be held by hydrogen bonds to all three bases and will thus be retained, on the average, for a much longer time (say, seconds or minutes). Now the code we have described insures that this more lengthy attachment can occur only at the points where the intermediate is needed. If one of the 20 intermediates could

stay for a long time on one of the false positions, it would effectively block the two positions it was straddling and hold up the polymerization process. Our code makes this impossible. This scheme, therefore, allows the intermediates to accumulate at the correct positions on the template without ever blocking the process by settling, except momentarily, in the wrong place. It is the feature which gives it an advantage over schemes in which the intermediates are compelled to combine with the template one after the other in the correct order.

The example given here is only for illustration, but it brings out the physical idea behind the concept of a comma-less code.

In passing, it should be mentioned that while the idea of making three nonoverlapping nucleotides code for one amino acid at first sight entails certain stereochemical difficulties, these are not insuperable if it is assumed that the polypeptide chain, when polymerized, does not remain attached to the template. A detailed scheme along these lines has been described to us by Dr. S. Brenner (personal communication)."

June 7, 1957

*Second printing*

HOW MAY AMINO ACIDS READ THE NUCLEOTIDE CODE?

*Final proof of the manuscript*

by Leo Szilard

(Submitted by Joseph E. Mayer)

The Enrico Fermi Institute for Nuclear Studies  
The University of Chicago, Chicago, Illinois

It is now generally believed that proteins are formed along ~~the~~ nucleic acid templates. The sequence of purine and pyrimidine bases in the template is supposed to represent a code that may somehow determine the sequence of the amino acids in the particular polypeptide (protein) that a given template will form. <sup>(1)</sup> The purine and pyrimidine bases of the template, the letters of the code, are adenine, uracil, guanine and cytosine if the template be an RNA molecule; and if the template be a DNA molecule, thymine takes the place of uracil.

It has remained so far ~~rather a~~ <sup>somewhat of a</sup> mystery in just what ~~xxxxxxxxxxxx~~ way amino acids could read such a code. In what manner can chemical forces -- of the kind we know to exist -- line up amino acids alongside such a template in the proper sequence and at the proper distance from each other, so that ~~thus might~~ a chemical reaction chain may link adjacent amino acids through peptide bonds with each other?

It is the purpose of the present paper to ~~indicate~~ <sup>discuss in detail</sup> a conceptually simple scheme that will -- at least by way of an example -- illustrate in what manner this might take place in the living cell.

(1) A. L. Dounce, Enzymologia, 15, 251 (1952)  
G. Gamow, Nature, Vol. 173, p. 318 (1954).

Because a template which synthesizes protein need not necessarily be the gene itself but must carry the same information as the corresponding gene, we shall here refer to such a template for the sake of brevity as a paragene.

The basic thought underlying the ~~scheme~~ <sup>considerations</sup> here presented consists in the assumption that there are a number of enzymes (or enzyme systems) -- perhaps twenty altogether -- in the cell, and that each of these catalyzes the formation of a particular trinucleotide <sup>trinucleotides</sup> which carries ~~either~~ <sup>one</sup> particular amino acid ~~or, more likely perhaps, a particular sequence of three amino acids.~~ <sup>or, more likely perhaps, a particular sequence of three amino acids.</sup> If the amino acid is carried by the nucleotide on a phosphate or pyrophosphate group as an acid anhydride -- which is a high energy compound -- then the energy needed for the formation of the peptide bonds will become free when the amino acid is split off. In this sense one can say that each amino acid may carry the energy needed for forming its peptide bond. *Add one sentence here*

According to the notions here ~~presented~~ <sup>used</sup>, amino acids can ~~not~~ <sup>themselves</sup> read ~~the~~ <sup>the</sup> code of the paragene. But ~~the~~ <sup>the</sup> trinucleotides, which carry the proper amino acids, may attach with their three bases through the formation of 6 hydrogen bonds to the proper sequence of three bases on the paragene, and thus the amino acids may be lined up in the proper sequence along the paragene.

Accordingly, sequences of three nucleotides along the paragene represent the code words, and the trinucleotides which carry the amino acids represent the anti-code words. We assume that these anti-code words are complementary to the code words in the sense that, where the code word contains adenine the anti-code word contains uracil (or thymine), where the code word contains uracil (or thymine) the anti-code word contains adenine; and similarly guanine corresponds to cytosine and cytosine corresponds to guanine. The rationale for this assumption is as follows:

- (2) Note added June 17th. The May issue of the Proc. Nat. Acad. Sci., Vol. 43, p. 416 (1957) which was belatedly received here, contains an article by F.H. Crick, J.S. Griffith, and L.E. Orgel which is in part identical with their previously circulated memorandum referred to in the text. In addition, however, the authors discuss -- in conjunction with a detailed oral communication of S. Brenner to them -- how amino acids might read the nucleotide code. The relevant passage of their paper is quoted in full in Appendix *A*.

The concept of code-letter and complementary code-letter arose originally from the interpretation of the structure of DNA given by J.D. Watson and F.H.C. Crick. <sup>(2)</sup> They showed that in a double stranded DNA

---

(2) Nature, Vol. 173, p. 318, 1955.  
Proc. Roy. Soc., Vol. 223, p. 80, 1954.

---

structure, adenine of one strand pairs with thymine of the other strand (which presumably plays the same role in DNA as does uracil in RNA) and similarly guanine pairs with cytosine. The helical structure of DNA permits just such pairing, and hydrogen bonding is possible between adenine and thymine, as well as between guanine and cytosine.

If the sequence of bases along one strand of DNA represents a coded message which consists in three letter-words then, because we have four letters to choose from such a message could utilize 64 different words. We might, however, be ~~restricted~~ <sup>limited</sup> to the use of 20 out of the 64 words that are available. The reasons for this ~~restriction~~ <sup>limitation</sup> would be as follows:

If all 64 possible three-letter combinations form ~~in fact~~ a code word, and if the paragene assumes at the time of the formation of the polypeptide a helical configuration similar to the helical configuration of DNA, then it follows that the code on the paragene must be read consecutively from one end -- say, from the "head" of the paragene downward. This is so because such a helical structure does not provide for commas between the individual code words, and in a code containing 64 words any three consecutive letters form a word. If we number the letters along the paragene, say from the head of the paragene downward, then the first three letters, the letters 1, 2, 3, form a word which was meant to be conveyed and so do the next three letters, the letters 4, 5, 6. But sequences of three letters which encroach on two adjacent words (such as 2, 3, 4 or 3, 4, 5, for example) form code words which are not meant to be conveyed.

In these circumstances, the code would be misread if the trinucleotides, which represent the anti-code words, assemble alongside the paragene simultaneously, rather than -- from one end on -- consecutively. If we want to have simultaneous assembly of the trinucleotides alongside the comma-less paragene, then we are restricted to 20 code words.

The notion of such a 20-word code, which needs no commas, was introduced by F.H.C.Crick, J.S.Griffith, and L.E.Orgel of the Medical Research Council Unit at the Cavendish Laboratory, Cambridge, in a memorandum circulated in May, 1956 among workers interested in the subject of protein synthesis. (2) From such a code we must demand that the letters 1, 2, 3 on the template form a code word, and the letters 4, 5, 6 also form a code word, but sequences of three letters, which encroach on two adjacent words (such as 2, 3, 3 or 3, 4, 5, for example) form no code word. Crick and his co-workers have shown that this demand can be met, that a code which requires no commas may be constructed and that it can accommodate 20 three-letter code words.

We shall now <sup>concentrate</sup> ~~single out~~ for more detailed examination <sup>on</sup> one conceivable model for protein synthesis which might provide for the lining up of the amino acids alongside the paragene, both in the proper order and at the proper distance from each other. This particular model is based on the following assumptions:

The ~~trinucleotides~~ <sup>oligonucleotides</sup> which form the anti-code words contain the sugar ribose rather than the sugar desoxyribose. Each particular ribose trinucleotide (the anti-code word) carries a particular sequence of three amino acids. A phosphate (or diphosphate) group is attached to the (2) carbon atom of the ribose moiety of each nucleotide and an amino acid is attached to each of these phosphate (or diphosphate) groups. The amino acid anhydrides <sup>contain</sup> ~~possess~~ an energy-rich P, or PP, bond which, when split, may release 12,000 or 16,000 calories, respectively.

During protein synthesis the nucleic acid strand that functions as a template (the parogene) may take up -- so we here assume -- a helical configuration resembling the helical configuration of a DNA strand in the double stranded DNA helix. The trinucleotides may then line up alongside the helical parogene with their purine and pyrimidine bases paired with the complementary bases of the parogene, and if they are so lined up, then the amino acids carried by the trinucleotides may come to lie at just about the right distance from each other to permit the formation of a peptide bond between adjacent amino acids. A chemical reaction chain -- starting perhaps from the head of a parogene -- may then move down along the parogene, split the acid anhydrides, and thus free the amino acids as well as make available the energy needed for the formation of peptide bonds between adjacent amino acids.

Adjacent amino acids can be linked only if the distance from each other is smaller or equal but not appreciably larger than the fundamental repeating distance in a polypeptide chain which is  $7.27\text{\AA}$ . The fundamental repeating distance in a fully extended polypeptide chain is about  $7\text{\AA}$ .

Since before they are linked into a polypeptide, the amino acids can rotate around the chemical bond <sup>(at the (2) carbon atom of the ribose)</sup> which ties them to the phosphate group, they might well be assembled along the parogene at a <sup>somewhat</sup> smaller distance from each other than the fundamental repeating distance of the polypeptide chain.

*If we had a 20 letter word code and if*  
~~Applying the concept of a 20 letter code that requires no commas to our particular model of protein synthesis, we may now say the following:~~

Each of the 20 amino acids may appear once attached to the leading letter and once attached to the trailing letter of the 20 ~~(trinucleotides)~~ anti-code words. <sup>then</sup> ~~therefore~~ among the polypeptides that can be formed, each amino acid may precede any other amino acid, and each amino acid may follow any other amino acid. This does not, of course, mean that any amino acid sequence is possible.

~~Some of the amino acid sequences that may be found experimentally in sequential analysis of proteins and polypeptides might show that the restrictions imposed by our model on the possible amino acid sequence are too severe, and that the model has to be modified to accommodate established facts. As will be presently seen some modification of our model may be indicated for other reasons also. <sup>in the case of our model</sup> There is certainly no inherent reason why we should have a pure three-letter <sup>code</sup> -- as we have assumed above -- and why, for instance, a certain number of four-letter <sup>code words</sup> should not be utilized also. <sup>(4)</sup> Nevertheless for the sake of simplicity~~

~~If we had a pure three-letter code, we would have to demand that the number of amino acid residues of all polypeptides or proteins synthesized in the manner described above should be a multiple of three. <sup>This is not the case.</sup> non-cyclic~~

So far the number of amino acid residues in polypeptides and proteins have been determined only within one rather special class; all of them represent substances which are secreted by mammalian tissues. The number of amino acid residues found in such polypeptides and proteins, which have been analyzed with adequate accuracy, are as follows:

- a) Insulin chain A: 21; insulin chain B: 30; corticotropin B: 39; oxytocin: 9; vasopressin: 9; Intermedin B: 18.

All of these would fit a pure three-letter code.

- b) Intermedin A: 13; glucagon: 29 and pancreas ribonuclease: 124.

These do not fit a pure three-letter code. Intermedin and ribonuclease would have to include at least one 4-letter word and glucagon at least two 4-letter words. ~~Therefore, we can only conclude that the conclusion that only three-letter words have been included.~~

(4) See Appendix A



## Observed rate of enzyme synthesis

According to the notions here adopted most enzymes are synthesized in growing bacteria at a rather low rate which does not represent the maximum synthesizing capacity of the corresponding paragenes. The rate of production of a given enzyme may, however, be greatly enhanced if the enzyme is induced, and what we are interested to learn is the maximal rate at which a paragene may be able to form the corresponding enzyme.

One of the most studied cases of enzyme induction is the induction of the enzyme  $\beta$ -galactosidase which splits lactose. Jacques Monod and his co-workers have shown that the production rate of this enzyme in bacteria can be greatly enhanced by certain chemical analogues of lactose, which act as inducers, and that the rate of production of the enzyme goes up almost instantaneously upon adding such an inducer to the medium. We are thus led to believe that the inducer may act by increasing the rate at which one template produces the enzyme rather than by increasing the number of templates that produce the enzyme at an unchanged rate.

In fully induced wild type bacteria growing in minimal medium this enzyme is contained in the amount of about 8 mgm. per  $10^{12}$  bacteria and thus amounts to about 8% of the total proteins. We obtain the rate at which this enzyme is produced in minimal medium per bacterium by dividing the amount contained in one bacterium by 1.44 times the doubling time (40 minutes) of the bacterium. We thus find for the rate, at which this enzyme is produced in fully induced wild type bacteria growing in minimal medium, about  $2 \cdot 10^{18}$  grams per cell per second.

If we assume a molecular weight of a million (Jacques Monod and Melvin Cohn estimate the molecular weight of this enzyme at about 800,000), we obtain a rate of 1.5 enzyme molecules per cell per second. The number of paragenes per cell is not known. ~~There might be a few paragenes~~ There might be a few paragenes present per cell rather than just one, and the number of paragenes might be of the order of magnitude of 10. On the other hand, smaller enzyme molecules might be synthesized somewhat faster than larger enzyme molecules.

On the basis of the figure given above, we are thus led to believe ~~that~~ in bacteria -- that when an enzyme is fully induced and enzyme synthesis proceeds at its maximal rate -- the rate of formation of the enzyme may be of the order of magnitude of one per second per paragene.

Computed rate of enzyme synthesis

We shall now attempt to compute at what rate a parogene may be able to synthesize the corresponding enzyme on the basis of the model that we have postulated. For the purposes of this computation, we shall assume that the molecular weight of the enzyme is about 100,000. We then have about 1,000 amino acid residues in the enzyme, and accordingly we would have to assemble alongside the parogene  $m = 300$  trinucleotides, each of which is "loaded" with three amino acids. ~~[assumed that amino acids are attached to the trinucleotides]~~

In the approximation which we shall ~~use~~ use ~~the minimum time,  $\tau_0$ ,~~ the minimum time,  $\tau_0$ , needed for the formation of the polypeptide is composed of two terms,  $\tau_1$  and  $\tau_2$ ;

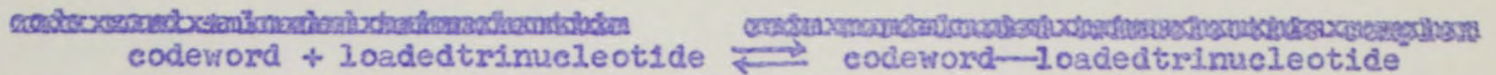
$$\tau_0 = \tau_1 + \tau_2$$

After all the amino acids, assembled alongside the template, have been linked into one polypeptide, and <sup>if</sup> ~~assuming that~~ this polypeptide is at once removed, a certain time,  $\tau_1$ , will elapse until the trinucleotides, which are now denuded of amino acids, evaporate from the template and their place is taken by trinucleotides which are loaded with the proper amino acids. We ~~shall~~ assume that the concentration of denuded trinucleotides in the cell is very small so that after the denuded trinucleotides evaporate, the loaded trinucleotides do not have to compete with denuded trinucleotides for their legitimate positions along the parogene. The time,  $\tau_1$ , which is necessary to permit evaporation of  $m$  denuded trinucleotides and to assemble  $m$  loaded trinucleotides in their place we shall compute here on the assumption that once a loaded trinucleotide has found its position alongside the template, it will not evaporate again. Because this assumption ~~is not~~ is not valid, we must make a correction which is represented by the second term,  $\tau_2$ .

Computation of  $\tau_2$ 

First we shall now compute this second term,  $\tau_2$ . This computation ~~will~~ <sup>is</sup> based on the fact that (because of reevaporation of the loaded trinucleotides, which ~~are~~ reversibly combined with the anti-code words of the parogene) there will be - no matter how long we wait - always a certain number of code words (sequences of three nucleotides) on the parogene which are not "covered".

We have to deal with a reversible reaction which may be written as follows:



The rate at which the reaction proceeds from left to right; i.e. the number of successful hits per unit time is given by

$$\text{hit rate} = A \rho$$

$\rho$  denotes the concentration of the particular kind of loaded trinucleotide in mol/cc. We shall, for the sake of simplicity, assume that the concentration in the cell of each kind of loaded trinucleotide is the same.

A stands for

$$(1) \quad A = 6 \times 10^{23} v \sigma p$$

$v$  is the molecular velocity,  $v = \sqrt{\frac{2RT}{\pi M}}$  where  $M$  is the molecular weight of the loaded trinucleotide. If  $M \approx 1000$ , we have  $v \approx 5 \times 10^3$  cm/sec.

$\sigma$  is the target area that must be hit by the loaded trinucleotide if Hydrogen bonding is to take place with the three adjacent nucleotides on the parogene. We may assume for  $\sigma$  a value of  $\sigma = 10^{-15}$  cm<sup>2</sup>.

$p$  denotes the probability that the loaded trinucleotide, when hitting the code word, is in just the right geometrical position to permit hydrogen bonding to take place between the three complementary pairs of bases that are involved. We may take for  $p$ , as a very rough estimate  $p = 1/300$ ; This would give for  $\sigma p$  the value  $\sigma p = 1/3 \times 10^{-13}$ .

The rate at which the reaction proceeds from left to right is given by the rate  $\alpha$  at which the ~~codeword-loaded trinucleotide~~ complex dissociates or, as we may also say, the rate  $\alpha$  at which the loaded trinucleotide evaporates from the parogene. For this rate we may write

$$(2) \quad \alpha = \text{rate of evaporation} \approx 10^{13} e^{-\frac{\Delta H}{RT}}$$

where  $\Delta H$  is the binding energy for the loaded trinucleotide.

In equilibrium the hit rate and evaporation rate must be equal.

The equilibrium constant,  $K$ , of this reversible reaction is defined as the concentration of the loaded trinucleotides at which the code word is covered half of the time. Thus we may write

$$(2) \quad AK \approx \frac{1}{2} 10^{13} e^{-\frac{\Delta H}{RT}} = \frac{1}{2} \alpha$$

$$\alpha = 2AK$$

In equilibrium the probability,  $f$ , for a given code word on the parogene not being covered by the proper loaded trinucleotide is given by

$$(4) \quad f = \frac{1}{1 + \frac{P}{K}}$$

Accordingly, in equilibrium, the total number of such gaps along the parogene which contains  $m$  nucleotides is given by

$$(5) \quad \text{"number of gaps"} = \frac{m}{1 + \frac{P}{K}}$$

We shall assume that most code words are "covered" in equilibrium and this means that

$$(6) \quad \frac{P}{K} \gg 1$$

We presume that after such equilibrium is established a chemical reaction chain is somehow triggered, and, moving down along the parogene, links adjacent amino acids into a polypeptide. The average time,  $\tau_1$ , needed for the formation of the polypeptide from the amino acids assembled along the parogene is given by the product of the "number of gaps", that

have to be filled consecutively, and the average time,  $\frac{1}{A\rho}$ , that it takes to fill one given gap. Thus, for  $\tau_2$  we may write

$$(7) \quad \tau_2 = \frac{1}{A\rho} \frac{m}{1 + \frac{P}{K}}$$

Computation of  $\tau_1$  and  $\tau_0 = \tau_1 + \tau_2$

When the polypeptide is formed and leaves the paragne, the code words are covered with the denuded trinucleotides. We may now compute the average time,  $\tau_1$ , needed for the evaporation of all the denuded trinucleotides and the assembling of all the loaded trinucleotides in their place. We shall, for the sake of simplicity, assume that a denuded trinucleotide evaporates at the same rate,  $\alpha$ , as a loaded trinucleotide. The rate,  $\alpha$ , at which a loaded trinucleotide evaporates from the template is given by (3) and we may also write this in the form

$$(8) \quad \alpha = A\rho \frac{2K}{P}$$

Since we have assumed  $\frac{P}{K} \gg 1$ , we have

$$(9) \quad A\rho \gg \alpha$$

As may be shown, <sup>(5)</sup> ~~we may write the following expression for very large values of m,  $\frac{P}{K} \gg 1$~~  (  $\ln m \gg 1$  ) we may write for  $\tau_1$

$$(10) \quad \tau_1 \approx \frac{1}{A\rho} \frac{P}{2K} \ln m$$

For the total time,  $\tau_0 = \tau_1 + \tau_2$  we thus obtain

$$(11) \quad \tau_0 \approx \frac{1}{A\rho} \left\{ \frac{m}{1 + \frac{P}{K}} + \frac{1}{2} \frac{P}{K} \ln m \right\}$$

If we wish to make this time as small as possible, we have to choose K so as to have for  $\frac{P}{K}$

$$(12) \quad \frac{P}{K} \approx \sqrt{\frac{2m}{\ln m}}$$

$$\sqrt{\frac{2000}{3 \times 2.3}} \approx 17$$

(5) See Appendix B

Substituting this value into (11) gives

$$(13) \quad \tau_0 \approx \frac{\sqrt{2}}{A\phi} \sqrt{m \ln m}$$

For a polypeptide containing 1,000 amino acid residues, i.e. a paragene containing about ~~500~~<sup>1000</sup> code words, we may write  $m = \frac{1000}{500}$ , and thus we obtain from (12) and (13)

$$(14) \quad \frac{\sigma\phi}{K} \approx 10 \quad \text{✓}$$

and

$$(15) \quad \tau_0 \approx \frac{50}{A\phi} \sqrt{12000} \approx \frac{100}{A\phi}$$

#### Estimates for the values of $\sigma\phi$ , $A$ , $\alpha$ and $K$

As we may see from (15) we obtain  $\tau_0 = 1$  if  $A\phi = 100$ . This means that for this particular value of  $A\phi$  one enzyme molecule is produced per paragene per second. As we have seen before, this is the order of magnitude of the rate at which highly induced bacteria produce the enzyme  $\beta$ -galactosidase per paragene.

We shall, therefore, in the following assume  $A\phi = 100$  and compute from it  $\phi$ , the concentration at which the different kinds of trinucleotides may be present in the cell.

If we use for  $\sigma\phi$  the value of  $\sigma\phi = \frac{1}{3} 10^{-17} \text{ cm}^2$ , then we obtain from (1)  $A$ ;  $A = 10^{10}$  and accordingly we have  $\phi = \frac{5 \cdot 10^9}{10^{-8}} \text{ mol/cc}$  ( $\phi = \frac{5 \cdot 10^{-6}}{10^{-5}} \text{ mol/liter}$ ).

It might well be, however, that  $\sigma\phi$  is ten times higher so that we have  $\sigma\phi = \frac{1}{3} 10^{-16} \text{ cm}^2$ , and then we obtain from (1) for  $A$ ;  $A = 10^{11}$ , so that we have  $\phi = \frac{5 \cdot 10^{10}}{10^{-9}} \text{ mol/cc}$  or ( $\phi = \frac{5 \cdot 10^7}{10^{-6}} \text{ mol/liter}$ ).

Thus the concentration,  $\phi$ , of the different kinds of trinucleotides in the cell is likely to be between  $5 \cdot 10^{-6}$  and  $5 \cdot 10^{-7} \text{ mol/liter}$ , and one would have to look for concentrations of this order of magnitude in order to obtain experimental confirmation of their presence.

Since we have assumed  $\frac{p}{k} \approx 17$ , and since we believe that we have  $A_p = 50$ , it follows that we must have for the rate of evaporation of the trinucleotides from the parogene,  $\alpha$

$$(16) \quad \alpha = 2AK \approx 10/\text{sec.} \quad \text{or} \quad 1/\alpha \approx 1/10 \text{ sec.}$$

For values of A between  $10^{10}$  and  $10^{11}$ , we have

$$(17) \quad 5 \cdot 10^{-11} < K \text{ (in mol/cc)} < 5 \cdot 10^{-10}$$

or

$$5 \cdot 10^{-8} < K \text{ (in mol/liter)} < 5 \cdot 10^{-7}$$

From (5) we obtain for a given pair of values A and K the binding energy  $\Delta H$  between the loaded trinucleotide and the code word on the parogene. For a value of  $A = 10^{10}$  and  $K = 10^{-10}$  mol/cc ( $K = 10^{-7}$  mol/liter), we obtain  $\Delta H$  18,000 calories. Since 6 hydrogen bonds are involved, this would mean about 3,000 calories per hydrogen bond.

For the same value of A and a value of K which is ten times larger,  $\Delta H$  would decrease by about 1400 calories.

### Conclusion

These considerations show that the theory which we postulated should be able to explain the high rate of enzyme synthesis which one observes in bacteria when the rate of formation of an enzyme is maximally enhanced by the use of an inducer. The basic thought of this theory consists in the assumption that trinucleotides and possibly also tetranucleotides read the code of the parogene, and that these oligonucleotides carry amino acids. *and perhaps higher oligonucleotides* ~~One particular model for protein synthesis, which assumed that each trinucleotide or tetranucleotide carries a sequence of three or four amino acids respectively, was singled out for detailed discussion because this model appeared to be the most plausible. This does not mean, however, that other models can not be considered.~~ *and shows it*

For instance, rather than to assume that each kind of trinucleotide carries a particular sequence of three amino acids, one might wish to explore the possibility that each trinucleotide ~~might carry~~ <sup>carries</sup> only one amino acid. In this case the amino acid might be carried by a phosphate group linked by an oxygen atom (ester linkage), either to the third or the fifth carbon of the 5 carbon sugar of either the leading or the trailing nucleotide. Assuming twenty different trinucleotides, each carrying one particular amino acid, we could have a code that requires no commas, with no restrictions imposed on the possible amino acid sequences of the proteins formed by the paragenes. *Amber, Crick et al for commaless code*

However, the particular model for protein synthesis here considered cannot be modified by simply saying that each trinucleotide carries one amino acid instead of carrying three amino acids, for adjacent amino acids would not then be at the right distance from each other to be linked into a polypeptide. Such an alternate model for protein synthesis would, therefore, require additional ideas concerning the formation of the polypeptide. <sup>(2)</sup>

I am grateful for the discussion of a variety of problems arising out of the considerations here presented which I had at the University of Chicago with Prof. Herbert S. Anker, Dr. Nandor L. Balazs, Mr. Hirono Kulci, Prof. Joseph E. Mayer, and Prof. Leonard J. Savage.



APPENDIX

(Added June 17/57)

A.) If, in addition to ribose trinucleotides, carrying three amino acids, we have in the role of anti-code words also ribose tetranucleotides, carrying four amino acids, then we must postulate that only three of the bases of the tetranucleotide pair with complementary bases on the paragene. The fourth base, if it is a purine, must pair with the wrong pyrimidine, and if it is a pyrimidine, it must pair with the wrong purine. The reason for this is as follows:

As formula (17) shows,  $\alpha$ , the rate of evaporation of a trinucleotide from the paragene can be estimated to be about 10/sec., corresponding to a binding energy,  $\Delta H$ , of about 18,000 calories or 3,000 calories per hydrogen bond. If all four bases of the tetranucleotide were paired with the complementary bases on the paragene, we would then have two more hydrogen bonds and presumably an additional binding energy of about 6,000 calories. The equilibrium constant,  $K$ , of the tetranucleotide would therefore be lower than the value computed for  $K$  of the trinucleotides by a factor of about  $10^4$ , and the rate of evaporation,  $\alpha$ , would be lower by the same factor. This would then make,  $\tau_0$ , the minimal time it takes for a paragene to form a polypeptide too long to be compatible with the observed rate of production of  $\beta$ -galactosidase in highly induced bacteria.

The number of words that may be constructed in a mixed three-letter word and four-letter word code of the kind described above, <sup>WHERE</sup> we demand that the code require no commas, has so far not been determined.

B.) The problem of computing  $\tau_1$  amounts to the solving of the following problem:

There are  $m$  boxes, each of which can hold one white ball or one black ball. Initially, at time,  $\lambda = 0$ , all of these boxes contain one

white ball -- a stripped trinucleotide. These white balls evaporate at the rate,  $\alpha$ , so that at time  $\lambda$  the probability,  $W(\lambda)$  of having no white ball in the box is given by

$$(18) \quad W(\lambda) = 1 - e^{-\alpha\lambda}$$

If the rate at which black balls fall into an empty box is designated by  $\beta$ , then the probability,  $y$ , that a given box does not contain a black ball at time,  $t$ , is given by

$$(19) \quad y(t) = e^{-\alpha t} + \int_{\lambda=0}^{\lambda=t} e^{-\beta(t-\lambda)} \frac{dW}{d\lambda} d\lambda$$

we may write

$$(20) \quad \int_{\lambda=0}^{\lambda=t} e^{-\beta(t-\lambda)} \frac{dW}{d\lambda} d\lambda = \frac{\alpha}{\beta-\alpha} \left[ e^{-\alpha t} - e^{-\beta t} \right]$$

so that we have

$$(21) \quad y(t) = \frac{1}{\beta-\alpha} (\beta e^{-\alpha t} - \alpha e^{-\beta t})$$

Thus we may write for the probability,  $x$ , that a given box contains a black ball

$$(22) \quad x(t) = 1 - y(t) = 1 - \frac{\beta}{\beta-\alpha} e^{-\alpha t} + \frac{\alpha}{\beta-\alpha} e^{-\beta t}$$

From this we obtain for  $P(t)$ , the probability that all  $m$  boxes contain one black ball

$$(23) \quad P(t) = x(t)^m = \left( 1 - \frac{\beta}{\beta - \alpha} e^{-\alpha t} + \frac{\alpha}{\beta - \alpha} e^{-\beta t} \right)^m$$

The average time,  $\tau_1$ , needed for the evaporation of the white balls from all  $m$  boxes and the filling of all  $m$  boxes with black balls is given by

$$(24) \quad \tau_1 = \int_0^{\infty} t \frac{dP}{dt} dt$$

If  $\beta \gg \alpha$ , we may write

$$(25) \quad \tau_1 = \int_0^{\infty} t \frac{d}{dt} (1 - e^{-\alpha t})^m dt$$

The expression  $\frac{d}{dt} (1 - e^{-\alpha t})^m$  has a maximum at some value of  $t$ ;  $t = \tau_0$  and -- for large values of  $m$  -- it becomes small very rapidly both below and above  $t = \tau_0$ . Therefore, if  $m$  is large we may write

$$(26) \quad \tau_1 \approx \tau_0$$

We obtain  $\tau_0$  by writing

$$(27) \quad \left\{ \frac{d^2}{dt^2} (1 - e^{-\alpha t})^m \right\}_{t=\tau_0} = 0$$

and from this we obtain

$$(28) \quad \tau_0 = \frac{1}{\alpha} \ln m$$

Thus for  $\beta \gg \alpha$  and  $m \gg 1$  we may write (26)

$$(29) \quad \tau_1 \approx \frac{1}{\alpha} \ln m$$

It can be shown from (23) and (24) that in the next higher approximation we have

$$(30) \quad \tau_1 \approx \frac{1}{\alpha} \{ \log(m+1) + C \} + \Delta ; \text{ where } \begin{cases} C = 0.577 \text{ (Euler's const.)} \\ \frac{1}{\beta} < \Delta < \frac{1}{\beta - \alpha} . \end{cases}$$

(\*) F. H. Crick, J. S. Griffith and L. E. Orgel write in the May issue of the Proc. Nat. Acad. Sci. (Vol. 43, pp. 419 and 420, 1957):

"To fix ideas, we shall describe a simple model to illustrate the advantages of such a code. Imagine that a single chain of RNA, held in a regular configuration, is the template. Let the intermediates in protein synthesis be 20 distinct molecules, each consisting of a trinucleotide chemically attached to one amino acid. The bases of each trinucleotide are chosen according to the code given above. Let these intermediate molecules combine, by hydrogen bonding between bases, with the RNA template and there await polymerization. Now imagine that such an amino acid-trinucleotide were to diffuse into an incorrect place on the template, such that two of its bases were hydrogen-bonded, though not the third. We postulate that this incomplete attachment will only retain the intermediate for a very brief time (for example, less than 1 millisecond) before the latter breaks loose and diffuses elsewhere. However, when it eventually diffuses to the correct place, it will be held by hydrogen bonds to all three bases and will thus be retained, on the average, for a much longer time (say, seconds or minutes). Now the code we have described insures that this more lengthy attachment can occur only at the points where the intermediate is needed. If one of the 20 intermediates could

stay for a long time on one of the false positions, it would effectively block the two positions it was straddling and hold up the polymerization process. Our code makes this impossible. This scheme, therefore, allows the intermediates to accumulate at the correct positions on the template without ever blocking the process by settling, except momentarily, in the wrong place. It is the feature which gives it an advantage over schemes in which the intermediates are compelled to combine with the template one after the other in the correct order.

The example given here is only for illustration, but it brings out the physical idea behind the concept of a comma-less code.

In passing, it should be mentioned that while the idea of making three nonoverlapping nucleotides code for one amino acid at first sight entails certain stereochemical difficulties, these are not insuperable if it is assumed that the polypeptide chain, when polymerized, does not remain attached to the template. A detailed scheme along these lines has been described to us by Dr. S. Brenner (personal communication)."