May 14, 1957 7. 6. 51 - 6 31 as muchantes de ware courses low Do Amino Acids Read the Code? 17 permin By Leo Szilard The Enrico Fermi Institute for Nuclear Studies The University of Chicago, Chicago, Illinois

It is generally believed that proteins are synthetized by RNA (ribonucleic acid), and that the sequence of purines and pyrimidines; i.e., adenin, uracil, guanin, and citosin determines the sequence of the amino acids in the polypeptide. It is presumed that these purines and pyrimidines form the basis of a three-letter code. Each of these bases represents one letter of the code, and certain groups of three letters form the word which corresponds to one specific amino acid. For any such code word to which we shall affix the plus sign we can construct the complementary combination (the anti-code word to which we shall affix a minus sign) by substituting uracil for adenin and adenin for uracil, citosin for guanin and guanin for citosin.

It has remained so far a complete mystery in just what way an amino What we he chemical affinity might be that may line up to now what the nature of the chemical affinity might be that may line up the various amino acids -A_i, A_j, A_k, A₁, etc. - in the proper sequence alongside the template that contains the three-letter code.

It is the purpose of the present paper to indicate a conceptually fle mallen which this might be accomplished by the living cell.

We shall assume that the cell contains twenty enzymes (or enzyme systems) which couple each amino acid so a trinucleotide that represents its code word. Thus We shall further assume that one strand of RNA, which may be contained in a double stranded helix is composed of a sequence of anti-code words. We shall Mark Merry Mark Second draft - New Paper

n R to and

Replacement, p. 2 It is the purpose of the present paper to indicate a conceptually simple scheme in which the living cell might solve the above defined problem might be revolved.

We shall assume that the cytoplasm of the cill contains 20 enzymes (or enzyme systems), each of which couples one of the 20 amino He property which which at His print any the eller acids to a trinucleotide (ribose or desory ribose) that we shall designate with (+), and that represents the code which is complementary to the code contained in a RNA (1/2 RNA(-) template. To these trinucleotides we shall affix the (+) sign in contradistinction to the trinucleotides contained in a RNA template that carries the code which is specific

COBRECTION

We shall assume that the cell (presumably the cytoplasm) contains 20 enzymes, E, or enzyme systems which couple each of the 20 enter amino acids to a specific nucleotide that represents the anti-code word for the particular amino acid. To these trinucleotides, we shall attach a (+) sign. To the complementary code words; i.e. to the trinucleotides which are contained in the RNA template alongside of which the proteins are formed, we shall attach the (-) sign. We assume that each amino acid is coupled to the phosphate group on the fifth the leading muclishide affice hi-carbon atom of either the first or the third nucleotide with an acid to and so that we have an acrid and good anhydride bond) so that we have a trinucleotide amino acid (as an alternative, we shall also consider the possibility that we have in place - mono of the trinucleotide amino phosphormino acid, a trinucleotide diphosphoamino acid). This is a high energy bond representing 12,000 calories. According to the notions here presented, a specific protein may

1. man be formed by a specific ribosenucleic acid template(+) in the following manner:

each The trinucleotideswhich carrying the proper amino acid will diffuse in the cytoplasm and each trinucleotide will reversibly attach by means of hydrogen bonds to the complementary ribose trinucleotide contained in the 1/2 RNA(+) template. If the protein to be formed has a molecular weight of 100,000 and contains 1,000 amino acid residues, then there will be 1,000 trinucleotides lined up alongside the template -- which we shall designate for our present purposes as a para-After a certain time the paragene will be completely covered by gene. the proper trinucleotides(+) and somehow a chemical reaction will be In this reaction the acid anhydride bonds split and the triggered. adjacent amino acids will be linked with each other by peptide bonds. Thus a polypeptide with a specific sequence of amino acids, A, A, A, , A, and A1, determined by the paragene will be formed which may be folded in a manner as yet unknown and will yield a protein, for instance, the andrap and enzyme for which this particular paragene is specific.

2-A

a setter

2 2 2 2 2 2

we must now compute at what rate a single paragene can form the corresponding enzyme. We will do so on the assumption that the binding energy between two trinucleotides, which are not complementary to each other, is somewhat smaller than the binding energy between two nucleotides, which are complementary to each other, and that therefore we can assume that we do not have to take into account the fact that the trinucleotides within the template will be occasionally covered for a short period of time by the wrong trinucleotide; i.e. one which is not complementary to it. We shall assume for the sake of argument that the molar concentration for the ribose-trinucleotides, which are present in the cytoplasm of the cell, is the same for each of the 20 different kinds of trinucleotides, and we shall designate the concentration of these trinucleotides(moløs per liter) the rate at which a single trinucleotide within the template is headed by a ribose-trinucleotide(+) in just the right position and with sufficient energy to reversibly combine with it. For this rate we shall write: hit rate

and for the rate at which such trinucleotides will dissociate off from the template, we may write

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

where K is the equilibrium constant for the reaction between one ribosetrinucleotide(+) in the solution or between one ribose-trinucleotide(-) on the template. From this we may compute the maximum rate at which a protein molecule can be formed on the assumption that there are m trinucleotides on the template; i.e. corresponding to m amino acid residues contained in the protein for which the template is specific. As may be seen, if we make the binding energy too small; i.e. K too large, we shall have a low value of

(1)

we must now compute at what rate a single paragen of orm the corresponding enzyme. We will do so on the assumption that the binding energy between two trinucleotides, which are not complementaries and other. somewhat smaller than the binding energy between of nucleotif X thick are complementary to each other, and that therefore we can assume that we do not have to take into account the fact that the trinucleotides within the template will be occasionally covered for a short period of time by the wrong trinucleotide; i.e. one which is not complementary to it. We shall assume for the which are present in the cytoplasm of the cell, is the same for each of the of these trinucleotides (moles per liter) the rate at which a single trinucleotide within the template is headed by a ribose-trinucleotide (+) in just the tide within the template is headed by a ribose-trinucleotide (+) in just the this rate we shall write: hit rate

2m -

and for the rate at which such trinueleotides will dissociate off from the template, we may write $C_{1} = C_{2} = C_{1} = C_{2} = C_$

where K is the equilibrium constant for the reaction between one ribosetrinucleotide(+) in the solution or between one ribose-trinucleotide(-) on the template. From this we may compute the maximum rate at which a protein molecule can be formed on the assumption that there are m trinucleotides on the template; i.e. corresponding to m amino acid residues contained in the protein for which the template is specific. As may be seen, if we make the binding energy too small; i.e. K too large, we shall have a low value of

T = the (1+ lum + m') + -1 -5-A

The maximum too small, we also have a low rate of protein production. /Protein production on the template is obtained if we have $(K_l = K_r)$

$$(\frac{k}{p})^{2} - 2\frac{k}{p} - 1 = 0$$

and/this relationship is fulfilled, then we may write A

Tcomp

For the value of A we can write

if

(5)

(3)

(4)

In this formula, small p is a frantion of the collisions between the ribosetrinucleotide(+) and the target area Σ of a given ribose-trinucleotide on the template and a given complementary ribose-trinucleotide(-) on the template.

Assuming now that the concentration of each ribose Atrinucleotide(+) in the cytoplasm that carries a given amino acid is of the order of magnitude , and assuming for \sum a value of and for p a value of , we find that a single paragene can produce protein at the of rate of 2,000 protein molecules in about 30 minutes. The equilibrium constant, K, we compute from equation (4), which gives a values of This equilibrium constant corresponds to an evaporation rate of nd the corresponding binding energy comes out to be about calories. To each factor of 10 by which the equilibrium constant is increased, there corresponds about 1400 calories by which the binding energy is decreased. As we shall presently see, the explanation of why DNA contains thymin instead of uracil which is contained in RNA might be due to the surmised fact that the binding energy due to the hydrogen bonding between uracil and adenin is

or t= 1 2 1+ lum + m + 1/2 + 2k2

RNA~ DNA

higher than between thymin and adenin. We may presume that the RNA(-) temtheas shrand the perse thought plates which form the proteins are produced by/DNA(+) template, and the of they on a same DNA(+) template may be presumed to produce the DNA(-) template. A number of different schemes may be devised to show how strands of RNA may produce strands of DNA, and how strands of DNA may produce strands of RNA. All of these schemes, however, have one thing in common. The DNA template must produce both/DNA template and the RNA template. This raises the problem of a possible mix-up in the formation of a composite template from the ribose-trinucleotides and the desoxyribose-nucleotides. It is tempting to speculate that the difficulty here mentioned is resolved in the following fashion: DNA may be produced on the RNA template from ribose-trinucleotides which carry on the 5-carbon atom of, say, the leading trinucleotide a hydrogen phosphate, while in contrast to this, the DNA strand is formed on the same DNA template from desoxyribose-trinucleotides which carry a hydrogen phosphate bond in the 5-carbon position, not on the leading but on the trailing nucleotide. The RNA strand may then be synthetized on the DNA template as follows: If the ribose-trinucleotide(-) combines with the trinucleotide, which is in a No. 1 position next to the head of the template, it will attach itself by splitting off one phosphate from the head of the template; if it combines with any other position, even though the trinucleotide with which it conbines is complementary to it, it will evaporate within a very short time. We assume here that the equilibrium constants for the combination of a RNA(+) trinucleotide that contains adenin with the complementary desoxyribose-trinucleotide that contains thymin is perhaps 10 tex times higher than the equilibrium constant for a combination of a trinucleotide and its complementary trinucleotide which contained uracil rather than thymin.

1

If a trinucleotide is attached in the No. 1 position, a trinucleotide that reversibly combines with the trinucleotide in the No. 2 position will have a good chance to combine with the trinucleotide attached to the No. 1 position. So the positions will be occupied one by one in the order of their serial number until m (1000) trinucleotides are all aligned and linked to an RNA(-) strand by phosphate bonds. In each case one phosphorus is split off for each trinucleotide attached to the RNA strand that is being formed. The time required for the formation of such an RNA strand is given by

and the time required for its evaporation is given by

If we now assume that the concentration of the ribose-trinucleotides(-), from which the RNA(-) strand is formed, is about the same as the concentration of the ribose-trinucleotides(+ \Diamond which form the proteins, and if we choose the equilibrium constant, K, to give the shortest possible time, we find that the time it takes for DNA strand(+) to make RNA Strand(-) is about

The optimum K value is higher by a factor of than the optimum and K value for a maximum rate of protein formation **xf** the corresponding binding energy is lower by calories. We have assumed above that the cytoplasm of the cell contains and enzymes which form ribosetrinucleotides, and coupled to each such ribosetrinucleotide through a high energy bond the proper amino acid. We shall further assume that these ribosetrinucleotides, which carry amino acid, are condensed on an RNA(-) template, the paragene, which is contained in the cytoplasm.

* * * 4

We shall now discuss in what manner the RNA(-) template may be formed within the nucleus on a DNA(+) strand. It is tempting to postulate that the RNA(-) strand formed inside the DNA(-) strand formed inside the DNA(-) template is formed through the condensation of ribose-trinucleotides(+)

one of the three trinucleotides again as a high energy phosphate bond either on the 3-carbon position or on the 5-carbon position. We shall call this nucleotide the leading nucleotide of the trinucleotides in contradistinction to the trailing nucleotide of the trinucleotide and the center nucleotide of the trinucleotide. For the sake of easier communication, without we shall specifically assume **that** restricting the general validity of our discussion that the DNA template is a head, and that when the ribose-trinucleotides are lined up alongside the DNA(+) template, then counting the letters of the code, starting at the head of the template, the high energy phosphate bonds of the ribose-trinucleotides may be seen attached to the 5-carbon position of the first letter of the code in every one of the trinucleotides. We shall assume that the RNA(-) template is synthetized alongside the DNA(+) template in the following manner: a ribose-trinucleotide which reversibly combines with the complementary code word on the DNA(-) template will, in general, evaporate, and the rate of evaporation is given by

In this expression K₁ designates the equilibrium constant of the desoxyribosetrinucleotide - ribose-trinucleotide complex. As may be seen later, the value

4.

of K₁ may be assumed to be higher than the value of K₀ for the binding energy of thymin to adenin due to hydrogen bonding being lower than the corresponding binding energy of the triribonucleotide(+)-triribonucleotide(-) complex that is formed in the case of protein synthesis. Thus the ribosetrinucleotide(+) will dissociate off after a short while from the DNA(+) template unless it is permitted to speak of a phosphate, and use the energies here liberated to form a chemical bond either with the head of the template or with an adjacent -- already chemically bound triribonucleotide. Synthesis of the RNA(-) template will accordingly take place as follows: When the first position next to the head of the template is filled by a trimbonucleotide(-), the high energy phosphor group of the hydrogen group will split off, and phosphate and ehemically bind to the head of the template. When position 2, the position, is filled with an RNA(-) trinucleotide, the high energy phosphorus will split off from the hydrogen phosphate group of this nucleotide, and the nucleotide will bind to the nucleotide which occupies the No. 1 position on the DNA(+) template. Thus in succession, position after position adjacent to a triribonucleotide that is already chemically bound to its neighbor will be filled with chemically bound

914 6 8

Correction -- phosphate linked

In this manner m ribotrinucleotides will be linked to form an RNA(-) strand. The time required for this process is given by

The RNA strand will in time detach itself from the DNA template on which it was formed. We presume that the bond tying the first triribonucleotide to the head of the template might be enzymatically broken, and that thereafter in succession the ribonucleotides(-), Nos. 1, 2, 3, etc., may dissociate off from the complementary tridesoxyribonucleotide contained in the

4-A

* 94 5-

The total time needed for formation and dissociation of one strand of RNA is given by

We have to for HOW DO

May 29, 1957

HOW DO AMINO ACIDS READ THE CODE? by Leo Szilard

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

anchion

It is generally believed that proteins are formed alongside of nucleic acid templates. The sequence of purine-pyrimidine bases in the template is supposed somehow to 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, are adenin, uracil, guanin and citosin if the template is an RNA molecule; and if the template is a DNA molecule, thymin takes the place of eitosin.

Because the template which synthetizes protein must carry the same information as the gene but need not necessarily be the gene itself, we shall refer to such a template as a paragene.

It has remained so far a complete mystery in just what conceivable way amino acids could read a code that consists of a sequence of purine-pyrimidine nucleotides. 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 a peptide bond with each other, and thus form a polypeptide.

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 in the living cell. The basic thought underlying this scheme consists in the assumption - serbrages twenty that there are a number of enzymes (or enzyme systems) in the cell, and that each of these catalyzes the formation of a particular trinucleotide which carries a particular amino acid or a particular sequence of three amino acids. The amino acids are tied to the trinucleotide through a high energy bond; either a P or PP bond, and representing an acid anhydride that when split may release an energy of 12,000 calories or 16,000 calories respectively, when they are motive. anothing to the materies have presented anono restols con not read the cade of the prayere at all, But the paraliakides in form togetor punie and pyrind dine loves of the formaliations can attach through the foundation of by drager hand he live proper luctificus on the prove-Jen and this time up the and no weeds in the proper Sequere, from Each Americo acid cares being present in the form of an order autydaythe wones with it the enorgy that nothe relaced in a hended

1-a

... in what manner this might be accomplished in the living cell. The basic thought underlying this scheme consists in the assumption that there are a number of enzymes (or enzyme systems) in the cell and that each of these catalyzes the formation of a particular/trisequence of three amine acids, The amino acids are tied to the trinucleotide through a high energy bond; either a P or PF bend, representing an acld anhydrided that menyspict may felease an energy of 12,000 calories of 16,000 calories respect at of there winthe les tides my Att anti- mole words and to complement easte mort into tag signesses of three underdi the which are unoten my to to the ende winds, The caste wor are certain requerces of three une fills mich ween adams an the mage We con

Sequences of three nucleotides along the paragene represent the code words, and the trinucleotides which carry the amino acids represent the anti-code words. 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 thymin), and where the code word contains uracil (or thymin), the anti-code word contains adenin, and similarly guanin corresponds to citosin and citosin corresponds to guanin. The rationale for this assumption is as follows: in what manner this might be accomplished in the living cell. The basic thought underlying this scheme consists in the assumption that a number of enzymes (or enzyme systems) -- not less than 20 and not more than 64) are contained in the cell and that each one of them catalyzes the formation of a particular trinucleotide which carries a particular minimum call and the cell and that each one of them catalyzes the formation of a particular trinucleotide which carries a particular sequence of three amino acids. These trinucleotides contain the fivecarbon sugar ribose rather than the five-carbon sugar desoxyribose. Each of these five carbon sugars carries a phosphate group in the twocarbon position and an amino acid is attached to each of these phosphate groups through an acid anhydride bond (either P or PP representing an energy of 12,000 calories or 16,000 calories, respectively.)

4 2.

(The possibility that the cell utilizes in fact not trinucleotides but tetranucleotides will be discussed later on in passing. And so will be the possibility that each multinucleotide (trinucleotide or tetranuclectide) might carry one amino acid only, tied with an acid anhydride bond to a phosphate group that in turn hangs either on a threeor a five-carbon atom of either the first or the last nucleotide.) provo Patt tor & erall these trinucleotides (which contain three of the four bases, yh These -adehin, uracil, guanin, or citosin) we shall as I me aroune sume to be complementary" ano thee tellow to the code-words contained in the paragene. To each code word on the template we can construct a complementary code-word by replacing adenin with uracil (or thymin), uracil (or thymin) with adenin, guanin with citosin, and citosin with guanin. The rationale for this assumption is A wante as follows: letter letter

The concept of code-word and complementary code word arese the intermediation originally from the study of the structure of DNA. Your & Watow the We know that in a double stranded DNA structure, adenin pairs with thymin (which presumably plays the same role in DNA as does uracil in RNA) and guanin pairs with citosin. Such pairing is required by the The helical structure of DNA permits such pairing, and hydrogen bonding is possible between adenin and thymin, as well as between guanin and citosin.

We may now tentatively take the view that during protein synthesis the paragene, 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 paragene through the formation of hydrogen bonds between the purine and pyrimidine bases contained in the trinucleotides and the complementary bases on the code of the paragene. When the trinucleotides are lined up in the proper order then, since each trinucleotide carries the proper amino acid, the amino acids are also lined up in the proper order.

We shall consider here now in greater detail one parts forflind up up of alongside the paragene cunascopert amino acids m in the proper order and at the proper distance from each other. This fork' cula solution the following around the

The homelentude cantain 3-A the Sh 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). To each of these phosphate (or diphosphate groups there is attached an amino acid. These acid anhydrides represent an energy-rich P(or PP bond so that each amino acid carries with it the energy which is necessary to form peptide bonds between adjacent amino acids. We assume that during protein synthesis the nucleic acid strand that functions as a template (the paragene) Mury a here around = take up 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 purine and pyrimidine bases of the paragene, and if they are so lined up, then the amino acids carried by the trinucleotides might come to lie at just about the right distance from each other to permit the formation of a peptide bond between the adjacent amino acids. A chemical reaction chain -- starting from the head of a paragene and the askot and move down alongside the paragene, split off one or two phosphates marci from the number 2 carbon positions of the ribose moieties, and neered supply the energy for the formation of the peptide bonds between adjacent amino acids. room If indeed proteins are formed in this manner, then there no madel imposes certain restrictions imposed on the possible amino acid sequences

that paragenes can produce. As we shall presently see, however, this restriction is not a very serious one. If we have four letters to choose from the can form 64 different three-letter words.

energy for the formation of the peptide bonds between adjacent amino acids.

later

If indeed proteins are formed in this manner, then there are certain restrictions imposed on the possible amino acid sequences that paragenes can produce. As we shall presently see, however, this restriction is not a very serious one.

If our code consists in three-letter words, if all 64 possinevelt ble three-letter combinations form a code word, and if the nucleic acid strand assumes at the time of the formation of the polypeptide the helical enfilie configuration discussed above, then it follows that the code on the paragene must be read consecutively from one end) -- say, "the"head "of the danou word : such as This is so because this helical structure does not paragene H4 onward provide for commas between the individual code words, and in a 64-word, three-letter word, code every 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 antialongspole of the paragene conserved T rather than from one code words, were to assemble simultaneously end on consecutively -- alongside of the paragene. If we want simultaneous assembly of the trinucleotides alongside of the, comma-less, paragene template, then we must be satisfied with only 20 code words instead of the 64 code words that are possible if there are commas between the words.

The notion of such a 20-word code which needs no commas, was fur introduced by F.H.C.Crick, J.S.Griffith, and L.E.Orgelfin a memorandum circulated in May, 1956 among workers interested in the subject of protein synthesis They have shown that, if we have four letters at our disposal From such a code we must demand that while the letters 1, 2, 3 on the template form a code word, and the letters 4, 5, 6 also form a code words, 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 and that a code which requires no commas may be constructed with that that a code which requires no commas may be constructed with that that a group of four different three-letter words with the letters drawn from a group of four different letters words with the letters area for a group of four different letters words.

21-0 -7, 29927-, 40

Honing shown His,) from which we form three-letter words, there can be constructed a code consisting of 20 words which requires no commas. They raised the question of whether the number 20 might represent a more than fortuitous coincimall dence in view of the fact that there might be just about 20 essential different amino acids that go into the formation of proteins. On the basis of the notions presented above, this coincidence would have to ache è n be regarded as fortuitous. to reyard this wine dury hat the pute From such a code which requires no commas we must demand that while the letters 1, 2, 3 alongside the template form a code word, and the letters 4, 5, 6 also form a code words, 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. [Grick and his co-workers have shown that their 20-letter code meets this demand.

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

appeny Each of the 20 amino acids may be carried once as the first appalrent letter and once as the last letter of one of the 20/trinucleotide) anticode 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. X there some of the amino acid sequences that may be found experimentally in sequential analysis of proteins and polypeptides will prove that the restrictions imposed by this system of any model polypeptide synthesis on the possible amino acid sequences are too severe , sould this might forme tins ho be untof the model that an remains to be seen. formitted There is, how inherent reason why we should have a pure three-letter code or why four-letter code words should not be utilized also. for instand a lerface number of

testays mathat

5.

If we had a pure four-letter word code and demanded that this code requires no commas, the number of code words available will be greater than 20.)

(Professor Leonard J. Savage of the University of Chicago informs me that in such a pure four-letter code the number of words is 46 or greater.

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 synthetized in the manner described above should be a multiple of three The number of amino acid residues in the proteins and polywhen here here (with antennate and formed formed) peptides so far analyzed are as follows:

(2) Insulin, chain A 21; insulin, chain B 30; corticotropin, B
39; oxytocin, 9; vasopressin, 9; intermedin B 18
(11) of these would
fit a pure three-letter code.

(4) Intermedin A 13; glucagon 29 and pancreas ribonuclease 124; Pthese 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.

obviously in a mixed system of three- and four-letter words, one can -- in the case of sufficiently large numbers -- never conclude that more than two 4-letter words must be included.

DP

unst he

6.

abrement

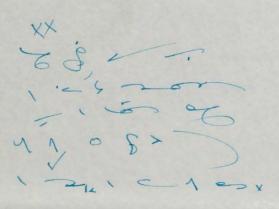
Rate of enzyme synthesis

-In general the rate of enzyme synthesis in bacteria is sound Most enzymes are synthetized in bacteria at a very low rate f which quite low. and according to the notions here presented these rates do not reflect the maximum synthetic capacity of the corresponding paragenes. We know, however, that the rate of production of a given enzyme may be greatly enhanced if the enzyme is induced. The enzyme, B-galactosidase, for instance, which splits lactose, can be induced -- as Jacques Monod and his co-workers have shown -- by certain chemical analogues of lactose; for instance, thiomethyl The rate of the production of these enzymes goes up galactoside (TNG). cit may youp almost instantaneously upon adding the inducer to the medium by a factor land pertraps were 10, of several thousand to about 5 enzyme molecules per second.

Because the increase in the production rate of the enzyme is almost instantaneous upon the addition of the inducer, we are inclined that Freezers we Het Torrestort - Storeth to assume that the inducer does not act by increasing the number of templates that make this particular enzyme but rather by increasing the rate at which one particular template makes the corresponding enzyme. Wit seems, furthermore, reasonable to believe that all templates which make enzymes of about the same molecular weight -- say, of 100,000 are potentially capable of synthetizing their enzyme at the same maximum rate. //Since the paragenes might consist of RNA and since there may be about five times as much RNA in the bacterial cell as there is DNA, it is conceivable that the same overy enzyme is made by perhaps five paragenes rather than just one paragene, and one might then conclude that a paragene must be capable of synflee somether for at at ter thetizing the enzyme for which it is specific at a rate of the that is order of magnitude of one per second, and the maximum rate might in fact be higher than this.

7.

Ahut In these circumstances, we are led to conclude that a lower limit for the order of magnitude for the rate of enzyme synthesis by one paragene to the one per second. The arter of magnitude of the order of magnitude of the order of magnitude for the order of magnitude of



8.

Computed rate of enzyme synthesis

We shall now attempt to compute at what rate a paragene may be able to synthetize the corresponding enzyme on the basis of the for the mo Bar por marked mechanism of protein synthesis that we have postulated. We shall assume that the molecular weight of the enzyme is 100,000 pire. that We have about 1,000 amino acid residues, in the enzyme, and accordingly we would Anyonde the propero have to assemble m = 300 trinucleotides, each of which is "loaded" with three amino acids. //In the approximation which we shall consider, the minimum time, T, needed for the formation of the polypeptide is composed of two terms, T1 and T2.), T= EI+L2 free free C If a polypeptide has been assembled and after all the amino acids, assembled/alongside the template have been joined in p polypeptide and assuming that this polypeptide is at once removed, a certain time, T,, (to be computed later) will elapse until the trinucleotides, which arrive of the anno accels loaded with amino acids but which are now denuded) evaporate from the template and their place is taken by trinucleotides which are loaded with he port amino acids. We shall assume that the concentration of denuded trinucleotide in the cell is very small compared to the concentration of trinucleotides which are loaded with the proper amino acids so that after the denuded trinucleotides evaporate, the loaded trinucleotides do not have to compete with the denuded trinucleotides for legitimate position along the paragene. The time, C1, which isnecessary to permit evaporation of all trinucleotides and to assemble all i.e. m, loaded trinucleotides, we shall in phin place compute here on the assumption that once a loaded trinucleotide has found Becunst its position alongside the template, it will not evaporate again. assumption is, of course, not correct, the the correction that is needed Atort me must is represented by the second term, (2. Knorefor while

9.

melinde one free freety to Hust In order to compute this second term, 2, we must consider the equilibrium between the loaded trinucleotide in the cell and the complex which these trinucleotides can form with the complementary nucleo-1 one requeree of molafiles in the call tide (the anti-code word) on the template. We shall in the following assume that the concentration for in the cell is the same for each kind of loaded tri-We shall further assume that the equilibrium constant, Ko nucleotide. for the equilibrium between, the loaded trinucleotide and one given comrequirel of the muchalides plementary code word is also the same for each kind of loaded trinucleotide. The probability that in equilibrium a given code word on the template is not covered by the proper trinucleotide, which is loaded with the proper amino acids, is given by and the total number of such "holes" alongside the template in equilibrium is given in equilibrium by which curtains 3 m meteophiles is yornen by - m T+ P/12 When) After the amino acids are all lined up, and if the formation of the polypeptide between the adjacent amino acids is somehow triggered an from the head of the template, then these holes will have to be filled

consecutively, and we may write for the time that this will take

 $T_2 = \frac{1}{AS} \frac{m}{1+S14}$ Here $AS \neq \beta$ represents the rate at which a loaded tri-

nucleotide combines with a given code word that is complementary to it.

B= AP

We may now compute the time, C_1 , needed for the evaporation of the naked trinucleotides and the condensation of the loaded trinucleotides in their place. The rate, α , at which a loaded trinucleotide would x cue to produce the template is given by

L= 2AK

We shall, for the sake of simplicity, assume that a denuded trinucleotide evaporates at the same rate \mathcal{A}

evaporates at the same rate of the Appendix As will be shown, we may then write within the approximation here attempted for the time. T.

$$T_{1} = \frac{1}{4g} \left\{ 1 + ln m - ln \frac{G}{\chi} \right\} \frac{1}{24}$$
when $\frac{G}{\chi} >>1$

For the total time, $\overline{L}_0 = \overline{C}_1 + \overline{C}_2$ we thus obtain $\int_{1+\frac{1}{K}} \frac{m}{1+\frac{1}{K}} + \frac{1}{2k} \int_{1+\frac{1}{K}} \frac{m}{2k} - \frac{1}{2k} \int_{1+\frac{1}{K}} \frac{m}{2k} \int_{1+\frac{1}{K}} \frac{m}{2k$

$$\frac{f}{k} = \sqrt{\frac{2m}{1 + lmm - ln\beta}} - 1$$

From this we obtain $f \approx 10$

as to hav

and, therefore, we obtain for ma 300

To ~ 50 AP

the ... For proper May 3/157 A = 61023 × 103 vop $A_{1} = 10$ M = 1000 $P = \frac{10}{4\pi(7)^{2}} \times \frac{1}{5}$ $A_{1} = 10^{7}$ $A_{1} = 10^{7}$ $M = 10^{7}$ Ag = 100/sec. ANA & m=300 To Time for anouly ~ 1 3 24 S 1 + lu m - lu B} $\beta >> d_{x} \qquad \beta = A g \qquad f = \frac{g}{ak}$ To fine to fill yaps . -Tr= Ap 2 m 1+ g 3 $T = \frac{1}{Ag} \left\{ \frac{m}{1+g} + \frac{1}{2} \frac{1}{1+g} + \frac{1}{2} \frac{1}{1+g} \frac{1}{2} \frac{$ S==9 T= 50 m yp=10 Mall T= 50 m yp=10 Mall

In order to compute this second term, Lo, we may consider the equilibrium between the loaded trinucleotide in the cell and the complex which these trinucleotides can form with the complementary nucleotide (the anti-code word) on the template.

The probability that in equilibrium a given code word on the template is not covered by the proper trinucleotide, which is loaded, with the proper amino acids, is given by

and the total number of such holes alongside the template in equilibrium is given in equilibrium by

SAR

We shall in the following assume that the concentration

of each kind of trinucleotide, loaded by the proper anino and we here at an and the required for the acids, is about the same. We may then write for the rate, β , at which 2 louder a free trinucleotide on the template combines with the proper trinucleotide combined with (that is loaded with the proper amino acids) france caste moved that is winp For a molekuter much of M = 1000 we have Vo= 5103 cm/sec it me may write impose

Here the coefficient, A, stands for $A_0 = 6/0^{20} v_0 v_0 p_0$ where v, the molecular velocity $v_0 = \sqrt{\frac{2RT}{TM}}$

and signa is the area of the target that must be hit for hydrogen bonding is to take place between the trinucleotide and the complementary code word and 0~~ 10 am2 We shall assume

M = malehaber mer

and finally p denotes the probability that the trinucleotide, when hitting the code, is in just the right position to permit hydrogen bonding to take

place between the three complementary pairs of bases that are involved, which Taking the molecular weight for a trinucleotide as about 1,000, we obtain

and if we estimate p = 1/300, we obtain for A me mog take as a nely rough estimate pox 300 these istimates give Ao = 10 per see E10 \$ 10 10 300 = 10/persee The simon 300 = 10/persee For a Turne for the partice If the concentration of each kind me flows at wants M landed himdenderders in malle to in the all is aleast Hen p= 10 - 5 mme/l = 10 - 50 Mp = 100 and the minimum Shime meder for and the minimum Shime meder and prep and marken for form and prep fran No To = 50 = 1/2 sec ane thes the maximum me at which the parti paragene can make enspire

is allant 2 per secured in A is af course promute man alle flice place between the three configure pairs of bases their are involved. The notecular weight for a trincicotide as about 1,000, we ob the the molecular weight for a trincicotide as about 1,000, we ob the the second are a trincicotide as about 1,000, we ob Vour rough estimate for 0 p might well be eff by a pretor 10 j Thus op might beten Kimes higher them any the name me queted and P might then he is kines lang i. e & it might he to millite. arming falo the interpranting Knalies nuge from HAMBE - 7 10 10-6 K=10-7 10 10-6 K=10 mand inmermud to him Hing ingly AH of 4 to the former that the the We can compose from K the landing mong AH lubmen ter In millentude and heart the the place proper toto incation of the providing. -for kK = 10 - 13 - Kr 7.36 2 = 10 0 1= 13×2.3 A A= 30 23 = 30 Q = 18,000 Cal or 3000 Cal band

(2) Amini within 10 10 A A preitracint 11

June 7, 1957

INO 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 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. (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 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 therexmight 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) G. Gamow, Nature, Vol. 173, p. 318 (1954).

Because a template which synthetizes 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 <u>paragene</u>.

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 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 --iwhich 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 <u>not</u> 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.**

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

Nature	, Vol	. 173,	р.	318 19	953.		
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. limited We might, however, be **xextrivited** to the use of 20 out of the 64 words that limitation 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 paragene assumes at the time of the formation of the polypeptide a helical configuration similar to the hælical 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, 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.

3**

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-wood 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 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 <u>proper distance</u> 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 <u>three</u> 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 represent an energy-rich P, or PP, bond which, when split, may release 12,000 or 16,000 calories, respectively.

4***

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 EXN 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 -- starting perhaps 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 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° . The fundamental repeating distance in a fully extended polypeptide chain is about 7° . 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 paragene at a 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.

5.**

Observed rate of enzyme synthesis

According to the notions here adopted most enzymes are synthetized in growing bacteria at a rather low rate which does not represent the maximum synthetizing 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 β -galactosidase which splits lactose. Jacques Monod and his coworkers 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 10¹² 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 pacterium. We thus find for the rate, at which this enzyme is produced in fully induced wild type bacteria growing in minimal medium, about 2 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. If the paragene is ENAL 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 synthetized somewhat faster than larger enzyme molecules.

In the basis of the figure given above, we are thus led to believe 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 <u>order of</u> magnitude of one per second per paragene.

Computation of T_2

First we shall now compute this second term, T_2 . This computation will be based on the fact that (because of reevaporation of the loaded trinucleotides, which reversibly combined with the anti-code words of the paragene) there will be - no matter how long we wait - always a certain number of code words (sequences of three nucleotides) on the paragene which are not "covered". We shall assume that the equilibrium constant, K₀, for the complexing of one code word by the proper loaded trinucleotide, is the same for the different kinds of trinucleotides, and that each code word complexes only with the proper trinucleotide. Similarly we shall assume that the different kinds of loaded trinucleotides are all present at the same concentration, f_0 , in the cell.

Willy to be replaced

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

(1)

Accordingly, in equilibrium, the total number of such gaps along the paragene which contains 3 m nucleotides is given by

(2)

"number of gaps" = $\frac{m}{1 + \frac{P_0}{K_0}}$

 $f = \frac{1}{1 + \frac{1}{1$

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

$$(3) \qquad \qquad \frac{30}{K_0} > > 1$$

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, T_2 , 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 have to be filled consecutively, and the average time, $\frac{1}{1/5}$, that it takes to fill one given gap.)

 \leq Thus, for \mathcal{T}_2 we may write

 $(4) \qquad \overline{C_2} = \frac{1}{p_0} - \frac{m_1}{1 + \frac{p_0}{k_D}}$

In this expression $1/\beta$ is the average time that it takes for a given gap to be filled by the proper trinucleotide, and we may write for β

 $\beta = A_o \int_o^{\infty} f_{o} = \int_o^{\infty} \int_\partial^{\infty} \int_\partial^$

Computation of I and to What share the partity

We may now compute the average time, \mathcal{T}_1 , needed for the evaporation of the naked trinucleotides and the assembling of the loaded trinucleotides in their place. The rate, α , at which a loaded trinucleotide evaporates from the template is given (see appendix) by

(5)
$$\alpha = 2 A_0 K_0 = 2 \frac{K_0}{\beta_0} \beta$$

(6) and for
$$\frac{f_{2}}{k} >> |$$
 we have $\beta >> \infty$

We shall, for the sake of simplicity, assume that a denuded trinucleotide evaporates at the same rate, α .

mation here ased for the time, 71,

For the total time, $\widehat{l_o} = \widetilde{l_1} + \widetilde{l_2}$ we thus obtain

(8)
$$T_0 \approx \frac{1}{R} \int \frac{m}{1+\frac{2}{K_0}} + \frac{1}{2} \frac{s_0}{k_0} \ln m_f$$

If we wish to make this time as small as possible, we have to choose K_{so} as to have for $\int_{K_{so}}^{O}$

(9)
$$\frac{f_0}{K_0} \approx \sqrt{\frac{2m}{m}}$$

Substituting this value into (b) and writing for
$$A$$
, so obtain
(10) $\tilde{c}_{0} \approx \frac{\sqrt{2}}{A \circ f \circ} \sqrt{m} \ln m$

For a polypeptide containing 1,000 amino acid residues $\frac{i.e.}{m}$ a paragene containing about 300 code words, we may write m = 300, and thus we obtain from (8) and (10)

(12)

Estimate for the value of A and

If \int_{0}^{0} is expressed in mol/liter, then the coefficient A is given by 20

0

To ~ 50 Aogo

v is the molecular velocity

(14)
$$v_{1} = \sqrt{\frac{2RT}{\pi M}}$$

and for a molecular weight of M \approx 1000, we have v $\approx 5 \times 10^3$ cm/sec.

 σ is the target area that must be hit if hydrogen bonding is to κ take place between three adjacent nucleotides on the paragene and the com- κ plementary loaded trinucleotides that move about freely within the cell. We assume for σ_{o} the value of $\sigma_{o} = 10^{-15} \text{ cm}^{2}$.

po 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 polas a very rough estimate

Thus we obtain $\sigma_0 \rho_0 = 1/3 \ 10^{-13}$.

With the above quoted values we obtain $A_{p} = 10^7/\text{sec.}$

If the concentration \mathcal{G}_{p} of each kind of loaded trinucleotide in the cell is $\mathcal{G}_{p} = 10^{-5}$ mol/liter, then we have

A So = 100 hits/sec. (17)

Substituting this value for A_{μ} into (12) gives for the average time, C_{ρ} , that it takes for one paragene to form one polypeptide : $C_{\rho} = 0.53$ e.

This means that the paragene may form about 2 polypeptide Jmolecules per second.

The value of \hat{l}_o depends only on the product of A_{p} and f_o . It might well be that the value which we estimated for $\mathcal{O}_o \rho_o$ is too low by a

factor of 10 and consequently the value which we obtained for A is also low by the same factor. If this were the case, then the value we chose for \int_{0}^{0} is ten times too high and the correct value would rather be $\int_{0}^{0} = 10^{6}$ mol/liter.

It is not possible for the present to estimate these values any closer. In these circumstances our assumption for the equilibrium constant, $K_{*}(f_{*} \approx 10)$ gives only a rough estimate

From the value of K we may estimate the binding energy, ΔH , for the combination of the trinucleotide with the paragene by writing from

(19)
$$2A_{0}K_{0} = 10^{13}e^{-\frac{4M}{RT}}$$

This gives for $K = 10^{-7}$ mol/liter, $\Delta H \approx 18,000$ calories, or since six hydrogen bonds are involved, about 3,000 calories per hydrogen bond.

If K is ten times larger, then ΔH is about 1400 calories lower.

Conclusion

These considerations show that the theory which we postulated would be able to explain the high rate of enzyme synthesis which one observes in bacteria when the rate of formation of an enzyme is enhanced by the use of an inducer. The basic thought of the theory here given consists in the assumption that trinucleotides read the code of the paragene and that these trinucleotides carry amino acids. We worked out the details in the case of a particular model for protein synthesis which assumed that each trinucleotide carries a sequence of three amino acids. Naturally I chose for the detailed theory the model which appeared to me to be the most likely to be correct. This does not mean, however, that other models need not be considered also.

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 first or the third 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.

If we make this assumption, however, we can not assume at the same time that during protein synthesis the paragene takes up the helical configuration which we postulated above. In the case of such a helical configuration the trinucleotides lined up along the paragene would not place the amino acids at the right distance from each other to permit adjacent amino acids to be linked to a polypeptide.

It is conceivable that the paragene might assume during protein synthesis some entirely different configuration which meets the requirement of bringing adjacent amino acids at the right distance from each other when the trinucleotides which carry one amino acid each line up alongside the paragene. But unless it is possible to indicate a plausible, non-helical, and yet regular configuration that meets this requirement, it does not appear useful to carry the discuss of such an alternate model for protein synthesis any further.

which hind to gabacof Such these trinucleotides (+) may take up positions face to face with the tempolate. corresponding trinucleotides (-) contained in the 1/2 RNA(-) strand. We assume that the amino acids in linked to the trinucleotide through a hydrogen bond, presumably an acid anhydride bond, between the phosphate group of the trinucleotide and the carboxyl group of the amino acid (estimated) 12,000 calories). One of several things might now happent the amount

(1) The adjacent acids from the peptide bonds join with each other so that there arises a polypeptide.

The adjacent ribonucleotides may link up with each other

and form a strand of RNA which we may designate as 1/2 RNA(+).

and form a protein, and the 1/2 RNA(+) strand can detach itself from the 1/2 RNA(-) template, or it can form a two stranded RNA molecule designated by RNA(+). by the case of the 2 RN H (-) henro take . -

To which polynucleotide one may assign an arror pointing from left to right or from right to left, which indicates the manner in which the phosphoric acid group is linked to the OH groups of the 5-carbon sugar. If we assign to the trinucleotides that carry the 20 different amino acids an arrow pointing from left to right, then the 1/2 RNA(+) strand which would arise in the process described above would also carry an arrow pointing from left to right.

each

Using this terminology, the 20 enzymes that catalyzed the formation of the postulated amino acid trinucleotides may be designated as . The basic thought here presented permits devising a number of different schemes for the synthesis of RNA and DNA. What schemes one might prefer depend to some extent on whether one believes that it might be possible to line up in the groove left free by oneand two-stranded helix formed by 1/2 RNA(+) and 1/2 RNA(-) strand of RNA. One may expect to have lined up How can we account for the existence of the 1/2 RNA(-) templates in the cytoplasm? We assume that this template synthetized in the nucleus. There are within the nucleus, so we shall assume, of anzymes which form twenty is ribosetrinucleotides(-) which represent the anticode word for the 20 amino acids. These enzymes couple each of these ribose-trinucleotides(-) to certain amino acids, perhaps predominantly arginine, histidine and lysin. We shall designate these enzymes with: E(AA, mith horizon unclean inter-)

to

There is, so we assume, within the nucleus a strand of DNA, which we shall designate as 1/2 DNA(+), which contains desoxyribose-trinucleotides (the code words) in the sequence which corresponds to the amino acid sequence, A, A, A, of the specific protein. //In analogy to what happened in the cytoplasm, the ribose-trinucleotides (-) will fuche come face to face with the corresponding desoxyribose-trinucleotides(+) contained in the 1/2 DNA(+) template , and with the help of the energy supplied by the acid anhydride bond of the ribose-trinucleotides will tink up to form a strand of RNA - 1/2 RNA(-). The molecules adjacent thus formed in the nucleus will diffuse out into the cytoplasm where they will serve as a template for protein formation as described above. The molecules remaining within the nucleus serve as a template for the synthesis of a strand of DNA -- 1/2 DNA(+). This strand of DNA is synthetized from desoxyribose-trinucleotides(+) which are formed and coupled to certain amino acids, perhaps predominantly arginine, histidine and lysin, by enzymes which are present within the nucleus. These enzymes may be designated by

The acid anhydride bonds of the ribose-trinucleotides that are lined up inside the template supply the energy for forming peptide bonds between the adjacent amino acids that are carried by the ribosetrinucleotides and to link the adjacent ribose-trinucleotides with each other through ester bonds. Thus the polypeptide and the RNA strand, 1/2 RNA(-), are formed but in contrast what happens in the formation of protein and polypeptide and the polynucleotide need not separate so that we have a ribonucleic protein which we may designate with AA - 1/2 RNA(-). Some of these will diffuse into the cytoplasm and may serve there as a template for the formation of proteins as discussed above. Insert - New paper

merifile The ribose trinucleotides carrying the various amino acids will freely diffuse around in the cytoplasm of the cell. Assuming that the concentration which an amino acid which is thus coupled to the corresponding trinucleotide is of the order of a milligram per liter, we may expect that the right trinucleotide will come to lie face to face with the corresponding code word of the 1/2 RNA(-) template. Assuming an activation energy of 0, the code word will complex with the anti-code word at the rate of about 100 times per second. If the complex remains undissociated for a sufficiently long period of time and if the template synthetizes an enzyme which contains 1000 amino acid residues, then the time which it takes on an average to have all the code words of the 1/2 RNA(-) template thus complexed by the anti-code word will be proportionate to the natural logarithm of M, and hus in the case of the synthesis of an enzyme which is between 100 and 100,000 amino acid residues, it may take about 1/10th of a second to assemble all the anti-code words on the template. If the binding energy of the word-anti-word complex is about 18,000 calories, then the anti-code words stick to the template long enough to permit the completion of the whole amino acid sequence. It is presumed that the six high energy bonds that can be formed between the two complementary trinucleotides will supply binding energy in excess of this value. We may assume that the free phosphate group/trinucleotides is linked to the amino acid which the trinucleotide carries by an oxide-anhydride bond representing about 12,000 calories. #After all the ribosd trinucleotides(+) are lined up on the 1/2 RNA(-) template, something might trigger a chemical reaction in which a phospha 18 is poples antyon & leaved) split off in which the amino acids carryed by the adjacent trinucleoann no and ne tides form a peptide linkage and the adjacent trinucleotides are linked by a phosphate ester bond and at the same time one phosphate is split off. In this way we would obtain simultaneously a poly-

4 +B(2)

peptide and 1/2 RNA(+) strand. The polypeptide has an amino acid sequence which is determined by the 1/2 RNA(-) template and the RNA strand, 1/2 RNA(+), that is newly formed is complementary to the 1/2 RNA(-) template.

\$N# OH mining no second pro.

Since it is presumed that protein synthesis can occur without accompanying net RNA synthesis, we have to consider several possibilities:

a) That when the amino acids that are lined up are linked to a polypeptide, contrary to what we said above, the ribose trinucleotides remain unlinked and zere returned to the free ribose trinucleotide pool.

b) That the 1/2 RNA(+) strand is formed but that such naked RNA strands are hydrolized to mononucleotides.

c) That the 1/2 RNA(-) template was part of a double stranded structure where the other strand is the 1/2 RNA(+) strand. These two complementary strands of RNA might form a helix and in the groove left free there may conceivably fit in the ribose trinucleotides(+) which carry the amino acids. When the polypeptide is formed simulaXtaneously a new 1/2 RNA(+) strand is formed. One of the two 1/2 RNA(+) strands, either the original strand or the newly made strand, may remain united with the 1/2 RNA(-) template, and the other strand which is now naked might be hydrolized. If protein synthesis by the RNA template in the cytoplasm follows this pattern, we would then expect a considerable turn-over of RNA to accompany protein synthesis.

5.

1/2 RNA(-) template, we presume, is formed at 1/2 DNA(+) strand in the nucleus. It is tempting to speculate that the RNA and DNA strands are fommed in the nucleus on the basis of a principle very similar to the one described above in connection with the synthesis of specific proteins, with the following difference: The RNA and DNA strands must be synthetized from ribose trinucleotides which have a free pyrophosphate

from trinucleotide phosphates from trinucleotiphosphate amino acids In the former case the synthesis may result in ribose or desoxyribonucleic acid strands, whereas in the latter case ribose or desoxy ribonucleX&oprotein may be formed.

This means that there would have to be in the nucleus enzymes which catalyze the formation of ribotrinucleotide diphosphates and desoxyribonucleotide diphosphates, or in the case of the second alterrative mentioned ribose trinucleotide amino acids, diphospho-amino acids and desoxyribose and trinucleotide amino acids. Among the amino acids thus coupled to trinucleotides, we presume that lysin, arginine, and histidine were predominant.