

Write this

THE GENETIC CODE

by

M. W. Nirenberg

National Heart Institute
National Institutes of Health
Bethesda, Maryland

9/15/67

I. INTRODUCTION

In recent years studies on the genetic code, protein synthesis, and regulation of protein synthesis have expanded to such proportions that investigators in other fields ^{disciplines} and graduate students often find it difficult to assimilate ^{the large body of information that is presented by your discipline} ~~unrelated facts~~ and to formulate general principals from them.

In writing this chapter the objectives have been to ^{it, described how the} ~~first,~~ the approaches ^{which were used to decipher the code,} ~~second,~~ the nature of the code and the ^{apparent logic of the code's design,} ~~third,~~ the apparent logic of the code's design. An

attempt is made to formulate general principals and the ^{apparent} logical design of the code.

This chapter is meant to be a discussion of the code, how it was translated and the experimental data. In writing the chapter an attempt has been made to formulate general ^{principles} principals from the data and to state the apparent logic which underlies the design of the code. Coverage of topics has been selective rather than comprehensive, ^{and} in some ways, the chapter ^{is} more of an essay than a review.

The data which ^{are} ~~is~~ now available on the code, the structure and function of nucleic acids and protein, and the process of protein synthesis

is so extensive that ~~it was~~ ^{it seemed} deemed essential from the outset to be selective ^{selective} rather than comprehensive ~~in discussing the data.~~

An attempt has been made in this chapter to survey the genetic code, to concentrate on fundamental ~~principles~~ ^{principles} ~~principals~~ ^{principals}, especially the apparent ^{logic} design of the code ~~is employed~~ and the logic which is employed. Coverage is selective rather than comprehensive.

B. Base Compositions of Codons

Synthetic RNA preparations containing all possible base combinations have been used as templates for ^{polypeptide} ~~protein~~ synthesis in vitro.

In practice, the major factor ^{limiting} ~~which limits~~ the sensitivity of the assay is the ^{amount} ~~presence~~ of endogenous mRNA in E. coli extracts. As shown in

Fig. _____, the level of endogenous mRNA is greatly reduced by incubating ^{ion} ~~ing~~ of

E. coli extracts in the presence of DNase and all components required for protein synthesis until amino acid incorporation ceases. ~~Extracts then can be dialyzed and stored until needed.~~ Protein synthesis then is almost completely dependent upon the addition of mRNA.

Optimal conditions for in vitro protein synthesis stimulated by synthetic mRNA were determined (), and methods were devised for

rapidly washing radioactive protein precipitates on cellulose nitrate filters (). Most radioactive ^{polypeptide} ~~protein~~ products are washed with 10% tri-

chloroacetic acid. ^T Those rich in proline are washed with 20% trichloroacetic acid () whereas lysine-rich ^{products} ~~proteins~~ are washed with a solution containing

sc ^S ~~m~~ tungstate and trichloroacetic acid ().

The specificity of a randomly-ordered RNA ^{as} ~~template~~ for amino

acid incorporation into ~~protein~~ ^{polypeptide chains} has been studied extensively with

extracts of *E. coli* (). A summary of the minimum ^{base composition of randomly ordered} ~~kinds of bases which~~ ^{polymers}

~~polynucleotides that correspond to amino acids are~~
~~required to code for each amino acid is shown in Table 1.~~

Polynucleotides ^{containing} with one kind of base usually ^{are acting as} is a template ^{for one} for a single

~~amino acid.~~ Little template activity ^{is} was detected with poly G; pre-

sumably G-G interactions inhibit the template activity of RNA (discussed in a later section).

A polynucleotide with two kinds of bases contains eight triplets;

six triplets with two kinds of bases, and two triplets with one ~~kind of~~ ^{kind of}

~~base~~. For example Poly UC contains:

- UUC UUU
- UCU CCC
- CUU
- CCU
- CUC
- UCC

Three preparations, poly UC, poly CG, and poly AG, ^{each serves} ~~as~~ active templates

for four amino acids. ^{per polynucleotide} three other polynucleotide preparations (poly UA,

poly UG, and poly CA) ^{Each one of} serve ^{as} templates for six amino acids ^{when}

~~There are~~ ^{The possible} four polynucleotides ^{containing} with three kinds of bases ^{are}

poly UAG, poly UCG, poly UCA, and poly CAG. A polynucleotide with three

such as kinds of bases, for example, poly UAG, resembles a mixture of seven polynucleotides as follows: poly U, poly ^AU, poly G, poly UA, poly UG, poly AG, poly UAG. ^{Randomly polymerized} Poly UAG contains ^{twenty-seven} ~~three~~ triplets with one kind of base, eighteen triplets with two kinds of bases, and six triplets with three kinds of bases. ~~a total of twenty seven triplets.~~

Each three-base polynucleotide stimulated the incorporation into protein of ^{at least} ten or more amino acids. Four ~~the~~ codons were

found which could not be accounted for by one-base, or two-base codons; poly UAG was a template for methionine and aspartic acid; and poly CAG

was a template for serine and ^{for} aspartic acid. ~~additional~~ additional three-base codons were assigned to amino acids. ~~Such results clearly demonstrate~~

~~the template specificity of codons and the minimum kinds of bases which may be present in codons for each amino acid.~~ ^{specificity of mRNA for amino acids as a function of the base composition of the polynucleotide} In addition the results

^{ed} show that the code is degenerate because amino acids such as leucine, arginine and serine respond to several polynucleotides which differ in base content.

not clear to confirm (10-20)

~~These~~ ^{showed that the template} ~~Such results clearly demonstrate~~ ^{specificity of mRNA for amino acids as a function of the base composition of the polynucleotide}

21
-7-

With additional data it is possible to ~~derive~~ ^{deduce} the relative ~~number~~ ^{number} of

~~proportion~~ ^{per} of bases in ~~many~~ ^{per} codon, as well as the kinds of bases ~~which~~ ^{that}

are present. Base compositions of RNA codons are derived as follows:

The base composition of a polynucleotide ~~can be determined quite accurately,~~ ^{is experimentally} then

~~The~~ expected frequency of each doublet or triplet is calculated easily

once the base ratio of a randomly-ordered polynucleotide is known. By

synthesizing a series of polynucleotides, each containing identical bases,

but ⁱⁿ with different proportions of bases, and ^{subsequently} determining the relative

proportions of amino acids directed into ^{polypeptide material} ~~protein~~ by each polynucleotide,

~~one can estimate~~ the base composition of the codon, as well as the number

of nucleotides per codon ~~can be~~ ^{could be} estimated.

~~It is difficult to compare directly the~~ ^{a direct comparison of} template efficiency of

~~different~~ ^{different} polynucleotide preparation with another, ~~for~~ ^{is difficult since} the efficiency of each

preparation may be influenced strikingly by factors other than base

composition, such as, the conformation of the RNA in solution, the presence

~~of~~ ^{of absence} terminal phosphate, ~~its~~ ^{of} molecular weight, the number of base residues

per molecule and so forth. These factors will be discussed in later

sections. However, ^{by determining} the amount of each amino acid incorporated into protein

due to the addition of ^{particular} a polynucleotide preparation, ~~can be determined~~

incorporated
the relative proportions of amino acids ~~directed into~~
under the influence of
~~protein by~~ different polynucleotide preparations can be compared.

~~the~~ Table ___ *is* shows ³ an example of data obtained with a poly AC preparation; similar data were obtained with four other poly AC preparations, each *different* in base ratio (). The four possible doublet permutations do not contain enough specific information to code for the six amino acids directed into protein by poly AC₃, whereas, the information content of the eight possible triplets is adequate. If every triplet were read, some amino acids would respond to two or more codons. In such cases the sum of the triplet frequencies would then be compared with the corresponding amino acid incorporation data. For example, if (ACA) and (ACC) both corresponded to the same amino acid, the sum of their frequencies would be 24.9 percent, which could not be distinguished from the frequency of the doublet, (CA), which is also 24.9 percent.

The relation between theoretical frequency and the experimentally determined frequency of amino acid incorporation into protein is shown in Fig. _____. *Clearly,* The data demonstrate that histidine, asparagine and *are coded by triplets* glutamine composition of a histidine codon is (CAC)₃ an asparagine codon,

(AAC)^g and a glutamine codon, (CAA). Threonine responds either to two triplets, one of base composition (ACA) ^{and} the other (ACC), or to the doublet (AC).

As shown in Fig. ____, proline responds to two triplets, CCC and (CCA), or to the CC doublet, and ~~lysine~~ lysine responds to the triplet AAA. ~~In this way~~ every triplet base ^{combination} ~~composition~~ in poly AC was assigned to an amino acid as follows:

| | |
|------------|----------------|
| Proline | CCC (CCA) |
| Histidine | (CAC) |
| Threonine | (ACC) (ACA) |
| Glutamine | (CAA) |
| Asparagine | (AAC) |
| Lysine | AAA |

In this way the nucleotide compositions of approximately 50 codons were estimated by Ochoa (), Nirenberg () and their coworkers. A summary ^{of the} is shown in Table _____. Tentative base compositions were estimated for many codons containing three different bases. Most amino acids were found to be coded by multiple words ^(Triplets). Since synonym codons often differ by only one base, the bases which were common to each synonym codon were assumed to occupy the same position within each triplet.

amino acids

Similar results were obtained in both laboratories, although extracts were prepared from E. coli B in the Ochoa laboratory and from E. coli ~~and~~ W3100

(a K12 strain) in the NIH laboratory.

C. Nucleotide Sequence of Codons

~~1. The Effect of Trinucleotides upon AA-tRNA Binding to Ribosomes~~

Each of the 64 trinucleotides have been synthesized, and assayed ^{for} ~~as a~~ template ^{capacity in the} ~~for~~ binding of E. coli AA-tRNA. Since the initial studies showed that AA-tRNA for some amino acids binds to ribosomes in response to trinucleotides at 0.02-0.03 M Mg⁺⁺, but not at 0.01 M Mg⁺⁺, (Nirenberg and Leder, 1964; Leder and Nirenberg, PNAS) a relatively high Mg⁺⁺ concentration (0.03 M) ~~was~~ ^{was} selected for the initial survey of trinucleotide-ribosome-AA-tRNA ^{complexing}. All responses found at 0.03 M Mg⁺⁺ then were reassessed at 0.01 and 0.02 M Mg⁺⁺.

Summaries of responses of unfractionated E. coli AA-tRNA are shown in Tables ___ and ___.

Most trinucleotides have been assayed for template specificity

with 20 AA-tRNA^{purified} preparations from E. coli, each acylated with

one radioactive and 19 unlabeled amino acids (). In surveying

trinucleotide specificity, unfractionated AA-tRNA^{generally} ~~usually~~ ^{is} has been

used^{only} initially, because^{the various} species of tRNA compete with one another

during the formation of AA-tRNA-codon complexes and the specificity of

codon recognition can be altered by changing the concentrations of

two or more species of tRNA.

Insert
discussion
(of tables & data)

Almost all triplets ~~were found to~~ correspond to amino acid.

are

Synonym codons ~~were found to be~~ logically related to one another, and

~~In~~ most cases, ~~synonym codons~~ differ only in the base occupying the third position of the triplet. Only four unique patterns of degeneracy were found, each pattern determined by the kinds of bases which occupy the third positions of synonym triplets.

Patterns of alternate third bases are:

- 1) G
- 2) U = C
- 3) A = G
- 4) U = C = A

*check
Table for
data
C. An*

A fifth pattern, U = C = A = G, ^{occurs} ~~is found~~ which is not necessarily unique, because this pattern would result if two simpler patterns were present, such as [(U = C) + (A = G)] or [(U = C = A) + (G)].

Codons specifying the initiation of protein synthesis may contain alternate bases at the first rather than the third position of the codons. For example, N-formyl-Met-tRNA responds to AUG, GUG, and possibly also UUG (discussed under Punctuation).

~~IV~~

One consequence of logical degeneracy is that many mutations

leading to single base replacements in DNA at sites corresponding to third bases of mRNA codons may not result in amino acid replacement in protein. Hence, many mutations ~~thus~~ are "silent."

The code appears to be arranged so that the effects of base replacements in DNA, or erroneous translation of a base in mRNA, often is minimized. Possible amino acid replacements in protein which would occur as a result of single base changes can be read in

Table ___ by moving horizontally or vertically from the amino acid in question, but not diagonally. ~~The code appears to be arranged~~

~~so that the effects of some errors may be minimized, since amino~~

acids which are structurally or metabolically related often correspond

~~to similar~~ RNA codons. *are coded by similar structure.* For example, Asp-codons, GAU and GAC, are

similar to Glu-codons, GAA and GAG; Ser-codons are related to ~~Thre-~~ codons, and so forth.

*demonstrating
base sequences of codons*

(152-4-10-60)

These results confirm 46 of the 53 codon base compositions which had been assigned on the basis of studies with synthetic, randomly-ordered polynucleotides and the cell-free protein synthesizing systems. (152-4-10-60)

Matthaei and coworkers
determined

28

upon AA-tRNA binding to ribosomes

the effects of sixteen polynucleotide preparations, each with a

one hundred different 5'-terminal doublet followed by approximately ~~100~~ C residues,

such as AU(C)₁₀₀ upon AA-tRNA binding to ribosomes were determined

(Matthaei). The results are shown in Table _____. Each RNA preparation

markedly stimulated binding of Pro-tRNA and Ser-tRNA. Ala-tRNA^{also} responded to most

of the polynucleotides. ~~etc.~~ The ^{marked} high response of Pro-tRNA to every polynucleotide demonstrates a high ^{-off} out-of-phase response to CCC residues.

The very high response of Ser-tRNA to polynucleotides ^{was} ~~is quite~~ un-

expected. Each RNA preparation contains three triplets, depending upon

the phase of codon recognition. For example, ^{Poly}UUCCC(C)₁₀₀ contains the

triplets UUC, UCC, and CCC, according to the reading phase and stimulates

binding of Pro-tRNA, Ser-tRNA, Phe-tRNA, and Ala-tRNA to ribosomes. In some cases,

bases one to three are recognized preferentially; in other cases, bases

two to four are preferred. Therefore, phasing preferences depend upon

the triplet rather than the exact location of the triplet in the 5'-terminal

region. ~~Although such data are of considerable interest in considering~~

~~regard to the preferred modes of phasing codon recognition, interpretation~~

~~of codon base sequence often is complicated by the phasing problem.~~

Nevertheless, Responses of AA-tRNA to approximately half of the sixteen codons tested are sufficiently high ~~so~~ that base sequence assignments can be derived readily. However, ^{ok} interpretation of the remaining data is complicated by the phasing problem. Responses of AA-tRNA to polynucleotides agree well with the responses of AA-tRNA to trinucleotides ^{in ribosome binding} if the assumption is made that almost every possible ^{kind of} triplet in each polynucleotide is recognized to a greater or lesser extent. However, three differences in response of AA-tRNA to polynucleotides and trinucleotides should be noted; Gly-tRNA does not respond to GGC(C)₁₀₀, but does respond to the trinucleotide, GGC; Leu-tRNA does not respond to ^CUC(C)₁₀₀ but does respond to randomly-ordered poly UC, although responses to the trinucleotide, CUC, ^{it may be} ~~are~~ difficult to detect with unfractionated Leu-tRNA (); ^{and finally, the} high response of Ser-tRNA to polynucleotides ^{is} are not observed ^{in the} with trinucleotides ^{= ribosome binding systems,} Khorana and his colleagues () synthesized ^{Polynucleotides with repeating doublet, triplet, or tetramer sequences} ~~were synthesized by Khorana and his colleagues () and were used~~ ^{them} to stimulate amino acid incorporation in E. coli extracts. ^{Table - summarizes} ~~A summary of the results,~~ ~~is shown in Table~~ RNA preparations which do not contain an initiator codon, such as AUG, ^G or GUG, are translated in almost every possible phase during protein synthesis in E. coli extracts. RNA with a repeating doublet

sequence contains two triplets in alternating sequence, and therefore is a template for ^{polypeptide} ~~protein~~ containing two amino acids in alternating sequence.

Most RNA preparations with a repeating triplet sequence are read in three phases, each phase corresponding to a different triplet. For example, poly UUC stimulates the incorporation of radioactive phenylalanine, serine and leucine into protein. Therefore, poly UUC resembles a mixture of poly UUC, poly UCU, and poly CUU; if the reading phase were $(-UUC \cdot UUC-)_n$, the protein product would be polyphenylalanine; if the phase ~~reading~~ were $(-UCU \cdot UCU-)_n$, the product would be polyserine; and if the reading phase were $(-CUU \cdot CUU-)_n$, the expected product would be polyleucine.

A polynucleotide with a repeating tetranucleotide sequence contains four triplets and therefore serves as a template for ^{polypeptides} ~~protein~~ with repeating tetrapeptide sequences. Polymers which stimulate incorporation of less than the expected number of amino acids contain terminator ^{Codons} ~~triplets~~, such as UAA, ~~UAG~~, or ~~the barrier triplet~~ UGA (discussed under Punctuation).

Results obtained with polynucleotides containing repeating ^{base} sequences directly demonstrate nucleotide sequences of terminator codons. Polynucleotides without initiator codons are translated in almost every possible phase and the preferred mode of phasing apparently depends upon the triplet rather

1

2

2

$$8(9) + \frac{1}{2} + \frac{1}{2}$$

3

1 (2)

4

5

4+2

3

than the exact distance of the triplet from the 5'-terminus of the polynucleotide. In most cases, base sequences of codons can be derived only if the reading phase of mRNA is known and is correlated with the phase of amino acid incorporation into protein.

The number of codons per amino acid is shown in Table __. ~~Six~~

The maximum number of codons for one amino acid is six (serine, leucine, and arginine) degenerate codons correspond to serine, five or six to arginine and to leucine and correspond to

~~From four to one codon for each of the remaining amino acids.~~

It should be noted that only one codon corresponds to tryptophan and one

to methionine (*initiator* codons for N-formyl-methionine are discussed under punctuation).