

Total Synthesis of the Structural Gene for the Precursor of a Tyrosine Suppressor Transfer RNA from *Escherichia coli*

1. GENERAL INTRODUCTION*

(Received for publication, March 19, 1975)

H. GOBIND KHORANA, KAN L. AGARWAL,[‡] PETER BESMER,[§] HENRY BÜCHI,[¶] MARVIN H. CARUTHERS,^{||} PETER J. CASHION,^a MATI FRIDKIN,^b ERNEST JAY,^c KJELL KLEPPE,^d RUTH KLEPPE,^d ASHOK KUMAR,^e PETER C. LOEWEN,^f ROBERT C. MILLER,^g KATSUMARO MINAMOTO,^h AMOS PANET,ⁱ UTTAM L. RAJBHANDARY, BELAGAJE RAMAMOORTHY, TAKAO SEKIYA, TATSUO TAKEYA, AND J. HANS VAN DE SANDE^j

From the Departments of Biology and Chemistry, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139 and, the Institute for Enzyme Research of the University of Wisconsin, Madison, Wisconsin 53706

With the ultimate objective of the total synthesis of a tRNA gene including its transcriptional signals, an *Escherichia coli* tyrosine suppressor tRNA gene was chosen. The arguments in favor of this choice are presented. A plan for the total synthesis of the 126-nucleotide-long DNA duplex corresponding to a precursor (Altman S., and Smith, J. D. (1971) *Nature New Biol.* **233**, 35) to the above tRNA is formulated. The plan involves: (a) the chemical synthesis of 26 deoxyriboooligonucleotide segments, (b) polynucleotide ligase-catalyzed joining of several segments at a time to form a total of four DNA duplexes with appropriate complementary single-stranded ends, and (c) the joining of the duplexes to form the entire DNA duplex. Ten accompanying papers describe the experimental realization of this objective.

Methods have been developed in recent years for the synthesis of bihelical DNA of defined nucleotide sequences. These involve: (a) the chemical synthesis of short deoxyriboooligo-

*This work has been supported by grants from the National Cancer Institute of the National Institutes of Health, United States Public Health Service (CA05178 and CA11981), the National Science Foundation, Washington, D.C. (GB-7434X, GB-21053X, GB-36881X, and BMS73-06757), the American Cancer Society (NP-140) and by funds made available to the Massachusetts Institute of Technology by the Sloan Foundation. This is Paper CXXXI in the series "Studies on Polynucleotides." The preceding paper is Ref. 1.

[‡] Present address, Department of Biochemistry, University of Chicago, Chicago, Illinois 60637.

[§] Present address, Department of Biology, Center for Cancer Research, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139.

[¶] Present address, Neue Wangerstrasse, Haus Aurora, 7320 Sargans, Switzerland.

^{||} Present address, Department of Chemistry, University of Colorado, Boulder, Colorado 80302.

^a Present address, Biology Department, University of New Brunswick, Fredericton, New Brunswick, Canada.

^b Present address, Department of Organic Chemistry, The Weizmann Institute of Science, Rehovot, Israel.

^c Present address, Department of Biochemistry and Molecular Biology, Cornell University, Ithaca, New York 14850.

^d Present address, Department of Biochemistry, University of Bergen, 5000 Bergen, Norway.

^e Present address, Department of Biochemistry, All India Institute of Medical Sciences, New Delhi, India.

^f Present address, Department of Microbiology, University of Manitoba, Winnipeg, Manitoba, Canada.

nucleotide segments corresponding to the entire two strands of the intended DNA, (b) phosphorylation of the 5'-hydroxyl end groups in the synthetic oligonucleotides using polynucleotide kinase, and (c) the head to tail joining of the appropriate segments when they are aligned to form bihelical complexes using the T₄-polynucleotide ligase. This methodology has been successfully applied to the total synthesis of the 77-nucleotide-long DNA corresponding to the major yeast alanine tRNA (2). While the accomplishment of this synthesis established confidence in the general methodology for DNA synthesis, and the availability of several relatively short DNA duplexes of defined nucleotide sequences made it possible to study aspects of transcription (3, 4) and of DNA enzymology (5-7), the synthetic DNA corresponding to the yeast alanine tRNA proved, at least for some time, unsuitable for studies of certain problems of central biochemical interest. For example, it had been hoped that the availability of synthetic DNAs would permit further studies of the following two problems: (a) the mechanism of initiation and termination of transcription and (b) precise structure-function relationship

^g Present address, Department of Microbiology, University of British Columbia, Vancouver, British Columbia, Canada.

^h Present address, Department of Applied Organic Chemistry, Nagoya University, Nagoya, Japan.

ⁱ Present address, Department of Virology, Hebrew University, The Hadassah Medical School, Jerusalem, Israel.

^j Present address, Division of Medical Biochemistry, Faculty of Medicine, University of Calgary, Calgary, Alberta, Canada.

in tRNA. With the continued hope of being able to apply the synthetic approach to these and related problems, the total synthesis of the DNA corresponding to an *Escherichia coli* transfer RNA gene was undertaken. We now wish to report the total synthesis of a DNA corresponding to the entire length (126 nucleotides) of the precursor to an *E. coli* tyrosine suppressor tRNA. The present paper gives the main arguments for the choice of this RNA and introduces the synthetic plan, while ten accompanying papers document the experimental realization of the objective (8-17). Brief reports on portions of this work have appeared during the last 4 years (18-21).

The first requirement for undertaking synthesis of a DNA is the specification of its sequence. For RNAs whose sequences are known, the sequences of the genes can be deduced directly. Further, among RNAs, the choice was made in favor of tRNA genes because of a variety of reasons. Of the various classes of gene products, the tRNAs are easily the most intriguing in regard to structure and function. These molecules have to be recognized by a rather large number of components of the protein-synthesizing machinery, such as by the aminoacyl-tRNA synthetases, by the nucleotidyltransferase which repairs the C-C-A end, by the ribosomes and by several proteins involved in protein chain initiation, elongation, and termination, and finally by messenger RNA. Also, tRNA molecules abound in modified bases and the nascent tRNA molecules have to be recognized by several modifying enzymes. Indeed, the tRNAs are a unique class of molecules, which evidently possess common secondary structure characteristic of nucleic acids but they also undergo folding to adopt tertiary structures. This has been amply demonstrated by the establishment of tertiary structures in a number of cases and the elucidation of the structure by x-ray diffraction methods (22, 23). Despite this recent progress, understanding of the structure-function relationships is largely lacking. It is hoped that chemical synthesis could, in principle, offer a definitive approach of wide

scope. Different parts of the tRNA structure could be systematically modified at the gene level. The modifications could involve additions, deletions, or substitutions of single or a few bases, or could be more extensive, such as the replacements of loops and stems by those present in different tRNAs.

CHOICE OF ESCHERICHIA COLI TYROSINE tRNA SUPPRESSOR GENE

The first major consideration in favor of an *E. coli* tRNA gene was the fact that biochemical work in the tRNA field is much more advanced with *E. coli* than with other organisms. Thus, the cell-free protein-synthesizing system, its characterization, the biochemistry of the ribosomes, and the understanding of the various factors required for initiation, elongation, and termination of polypeptide chains are all much better understood than with other systems. Specifically in the case of tyrosine tRNA, the aminoacyl-tRNA synthetase had been purified and characterized by Calendar and Berg (24).

A second consideration was the accuracy of the nucleotide sequence of the tRNA chosen. An assurance on this account at the start of the synthetic work was obviously desirable. The sequence of the *E. coli* tyrosine tRNAs (including the amber suppressor tRNA) was first determined by Goodman *et al.* (25). Fortunately, the same sequence was determined independently for this tRNA by RajBhandary, Nishimura, and their co-workers (26) by using a separate set of methods.

From the standpoint of studies on structure-function relationships, among many other general considerations, two specific lines of reasoning in favor of tyrosine tRNA were as follows. Firstly, a comparison of the *E. coli* tRNA^{Tyr} and *E. coli* tRNA^{Met} sequences (Fig. 1) (27) showed remarkable similarities in parts of the cloverleaf structures but a striking difference was in the size of the loop III in the two tRNAs. It seemed reasonable to investigate the minimal changes in tRNA^{Tyr} which would be required to elicit an initiator function in protein synthesis (see also Kleppe *et al.* (17)). Secondly, the

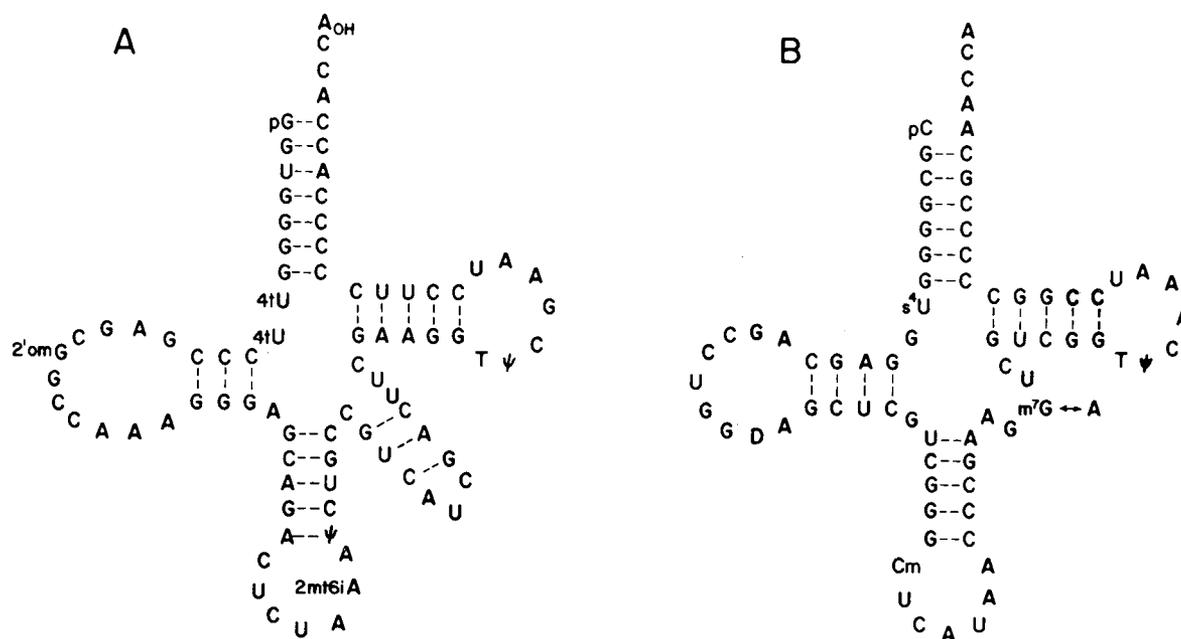


Fig. 1. Cloverleaf models of the primary nucleotide sequences of an *Escherichia coli* tyrosine suppressor tRNA (A) and the *E. coli* tRNA^{Met} (B).

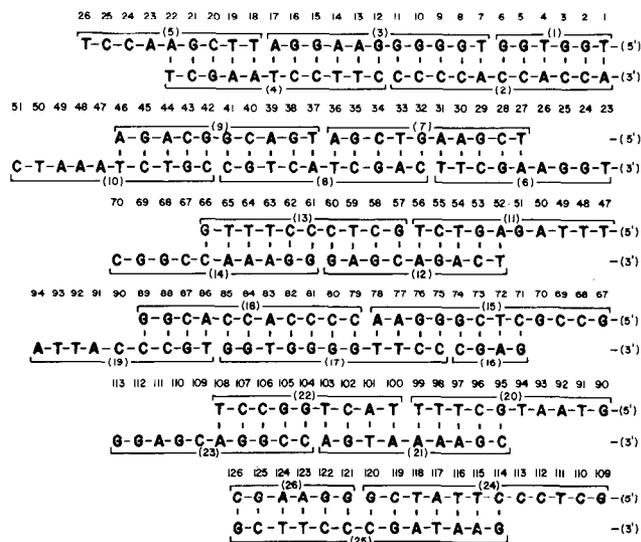


Fig. 3. The plan showing the 26 segments to be synthesized chemically in the total synthesis of the bihelical DNA corresponding to the precursor for the *Escherichia coli* tyrosine suppressor tRNA.

T-C-C-C-T-T-A), also derived from the above heptanucleotide and belonging to segment 8 in the previous work, now was used in segment 24. Further, a systematic search for sequences occurring more than once in the present DNA² showed that the nonanucleotide sequence d(C-C-C-C-A-C-C-A-C) occurs in segment 2 as well as in segment 18. Similarly, the hexanucleotide sequence, d(T-T-C-G-A-A), occurs twice (nucleotides 18–23 and 26–31) in the same strand and because it is self-complementary, it is also present twice in the complementary strand. Shorter sequences can be found to repeat with increasing frequency, but their use may not always be allowed by other considerations.

As mentioned above, the hexanucleotide sequence d(T-T-C-G-A-A) occurs four times. The best use of this common sequence would be to place it at the 5'-end of four segments and the same protected hexanucleotide intermediate could then be rapidly elongated to give the four segments. However, it is not possible to do so because of the previous experience with a situation of this type (35). Thus, a duplex such as the one shown in Fig. 4A would rather undergo dimerization to give B in Fig. 4 than to add the required additional segment. The plan adopted (Fig. 3) (segments 4 to 7) circumvented this side reaction by actually using the pentanucleotide sequence d(T-C-G-A-A) in the synthesis of two segments (segments 4 and 7) and, similarly, the symmetrical hexanucleotide sequence d(T-T-C-G-A-A) in the synthesis of segments 5 and 6. The joining reactions were to be so arranged that the two segments containing the symmetrical hexanucleotide sequence would not be present in the same reaction mixture. The plan consisting of the 26 segments (Fig. 3) accommodates the considerations described above.

Attention may be drawn to the situation around segments 16 and 17 in the plan shown in Fig. 3, the former being only a tetranucleotide and the latter being especially guanine-rich. Indeed, this part of the gene has remarkably high G:C content. The plan in this region underwent revision twice, the main

² We are grateful to Professor Walter Fitch of the Department of Physiological Chemistry, University of Wisconsin for doing this search with his computer program for detecting common sequences in proteins and nucleic acids.

reason being that initially the aim was to synthesize only the 85-nucleotide unit-long DNA corresponding to the tyrosine tRNA, the synthetic plan in this region was reconsidered. The plan proposed aimed at the chemical synthesis of the pentadecanucleotide containing the undecanucleotide of segment 17 and the tetranucleotide shown as segment 16. Unfortunately, the synthesis of segment 17 alone proved to be rather overwhelming and an extension of this synthesis to include segment 16 was not practical. It was therefore hoped that conditions might be found for the subsequent enzymatic reactions such that the tetranucleotide (segment 16) would join to the neighboring segments 14 and 17. The enzymatic joinings were indeed carried out successfully although the yields left a great deal to be desired (15). An improvement in the enzymatic joinings in this part is still under investigation by undertaking the synthesis of the tetradecanucleotide which combines the present segment 16 with segment 14. The results of this study will be reported upon at a later date.

Finally, it may be noted that in the plan shown in Fig. 3, single-stranded runs (hexanucleotide sequences) are available at the ends of the double-stranded DNA for extension to the regions which would, presumably, form the promoter and terminator regions for the transcription of the gene.¹

ENZYMATIC JOINING OF CHEMICALLY SYNTHESIZED SEGMENTS

As mentioned above, the results of joining experiments are frequently unpredictable. In experimental systems containing three or four segments, the yields vary very widely. Therefore, a large amount of empirical work is necessary to determine the combination of segments which would give optimal yields in the overall joining reactions. Following extensive experimentation, the 26 chemically synthesized segments (Fig. 3) were divided into four groups shown in Fig. 5. While the detailed arguments are presented in the individual papers dealing with different sections (13–16), it may simply be mentioned here that duplex [I] could only go as far as segment 5. In duplex [II], as many as eight segments could be used in a one-step joining reaction without any ambiguity. Duplex [III] required particularly detailed investigation for reasons mentioned above. The yield was rather low. Duplex [IV] consisted of six to seven segments (segments 19 to 25). A great surprise was the failure to join segment 26 to the remainder of duplex [IV].

Having prepared the four duplexes shown in Fig. 5, the next step was to quantitatively phosphorylate the terminal 5'-OH groups in the duplexes in preparation for the ligase-catalyzed joining to complete the synthesis of the total duplex. Work with a number of defined duplexes³ showed, however, that the rates and extent of phosphorylation of 5'-OH groups at the termini of DNA duplexes by the polynucleotide kinase are influenced very much by the duplex structures around the 5'-OH groups. To ensure facile and complete phosphorylation, it seemed clearly desirable to have the terminal 5'-OH groups at the protruding single-stranded ends of the duplexes. The grouping was therefore amended in regard to duplexes [II], [III], and [IV], as shown in Fig. 6. The modified grouping largely met the above requirement and no difficulty was experienced in the phosphorylation reactions with the duplexes and, therefore, in the completion of the total synthesis (17).

Five accompanying papers (8–12) describe in a condensed

³ Present work and unpublished work from the laboratory of Dr. K. Kleppe, Bergen, Norway.

gene. Towards these objectives, the sequence of 23 nucleotides in the region adjoining the C-C-A end has been determined (32) and the corresponding DNA duplex has already been synthesized (36).

Similarly, the sequence of 29 nucleotides immediately adjacent to the initiation point of transcription of the precursor to the tRNA has also been determined (33). Studies are continuing on the sequence work, as well as on synthesis as the sequence becomes known. Concurrently, work is in progress on the mechanism of action of the DNA-dependent RNA polymerase. In addition, a number of other promoters which are recognized by the *E. coli* polymerase are under intensive study in a number of laboratories and progress in sequence determination is rapid. It seems very likely that insights into the mode of binding of the enzyme, selection of the initiation site and related aspects of the mechanism of transcription will be gained in the near future. Synthetic work could further aid in more precisely defining the chemistry of the various steps in the overall process. Consequently, controlled transcription of the synthetic gene for the precursor to the tyrosine tRNA should be possible.

With the recent dramatic progress in methodology for DNA sequencing and its successful application in the determination of the nucleotide sequences in the control regions in a variety of genetic systems, it seems certain that synthesis as exemplified in the present series of papers will play an important role in understanding the mechanisms of the DNA-protein interactions and the expression of genetic information in general.

REFERENCES

- Kleid, D. G., Agarwal, K. L., and Khorana, H. G. (1975) *J. Biol. Chem.* **250**, 5574-5582
- Khorana, H. G., Agarwal, K. L., Büchi, H., Caruthers, M. H., Gupta, N. K., Kleppe, K., Kumar, A., Ohtsuka, E., Raj-Bhandary, U. L., van de Sande, J. H., Sgaramella, V., Terao, T., Weber, H., and Yamada, T. (1972) *J. Mol. Biol.* **72**, 209, and accompanying papers
- Kleppe, R., and Khorana, H. G. (1972) *J. Biol. Chem.* **247**, 6149
- Terao, T., Dahlberg, J. E., and Khorana, H. G. (1972) *J. Biol. Chem.* **247**, 6157
- Gupta, N. K., and Khorana, H. G. (1968) *Proc. Natl. Acad. Sci. U. S. A.* **61**, 215
- Kleppe, K., Ohtsuka, E., Kleppe, R., Molineux, I. J., and Khorana, H. G. (1971) *J. Mol. Biol.* **56**, 341
- van de Sande, J. H., Loewen, P. C., and Khorana, H. G. (1972) *J. Biol. Chem.* **247**, 6140
- van de Sande, J. H., Caruthers, M. H., Kumar, A., and Khorana, H. G. (1976) *J. Biol. Chem.* **251**, 571-586
- Minamoto, K., Caruthers, M. H., Ramamoorthy, B., van de Sande, J. H., Sidorova, N., and Khorana, H. G. (1976) *J. Biol. Chem.* **251**, 587-598
- Agarwal, K. L., Caruthers, M. H., Fridkin, M., Kumar, A., van de Sande, J. H., and Khorana, H. G. (1976) *J. Biol. Chem.* **251**, 599-608
- Jay, E., Cashion, P. J., Fridkin, M., Ramamoorthy, B., Agarwal, K. L., Caruthers, M. H., and Khorana, H. G. (1976) *J. Biol. Chem.* **251**, 609-623
- Agarwal, K. L., Caruthers, M. H., Büchi, H., van de Sande, J. H., and Khorana, H. G. (1976) *J. Biol. Chem.* **251**, 624-633
- Sekiya, T., Besmer, P., Takeya, T., and Khorana, H. G. (1976) *J. Biol. Chem.* **251**, 634-641
- Loewen, P. C., Miller, R. C., Panet, A., Sekiya, T., and Khorana, H. G. (1976) *J. Biol. Chem.* **251**, 642-650
- Panet, A., Kleppe, R., Kleppe, K., and Khorana, H. G. (1976) *J. Biol. Chem.* **251**, 651-657
- Caruthers, M. H., Kleppe, R., Kleppe, K., and Khorana, H. G. (1976) *J. Biol. Chem.* **251**, 658-666
- Kleppe, R., Sekiya, T., Loewen, P. C., Kleppe, K., Agarwal, K. L., Büchi, H., Besmer, P., Caruthers, M. H., Cashion, P. J., Fridkin, M., Jay, E., Kumar, A., Miller, R. C., Minamoto, K., Panet, A., RajBhandary, U. L., Ramamoorthy, B., Sidorova, N., Takeya, T., van de Sande, J. H., and Khorana, H. G. (1976) *J. Biol. Chem.* **251**, 667-694
- Besmer, P., Agarwal, K., Caruthers, M. H., Cashion, P. J., Fridkin, M., Jay, E., Kumar, A., Loewen, P. C., Ohtsuka, E., van de Sande, J. H., Sidorova, N., RajBhandary, U. L. (1971) *Fed. Proc.* **30**, 1314
- Khorana, H. G. (1973) *Naturwiss. Rundsch.* **26**, 137
- Khorana, H. G. (1974) *Pure and Applied Chemistry* **2**, 19-43
- Khorana, H. G. (1975) in *Proceedings of the International Symposium on Macromolecules, Rio de Janeiro* (Mano, E. B., ed) pp. 371-395, Elsevier Scientific Publishing, Amsterdam
- Kim, S. H., Suddath, F. L., Quigley, G. J., McPherson, A., Sussman, J. L., Wang, A. H. J., Seeman, N. C., and Rich, A. (1974) *Science* **185**, 435-440 and related papers
- Klug, A., Robertus, J. D., Lardner, J. E., Brown, R. S., and Finch, J. T. (1974) *Proc. Natl. Acad. Sci. U. S. A.* **71**, 3711
- Calendar, R., and Berg, P. (1966) *Biochemistry* **5**, 1681, 1690
- Goodman, H. M., Abelson, J., Landy, A., Brenner, S., and Smith, J. D. (1968) *Nature* **217**, 1019; Goodman, H. M., Abelson, J. N., Landy, A., Zadrzil, S., and Smith, J. D. (1970) *Eur. J. Biochem.* **13**, 461
- Harada, F., Gross, H. J., Kimura, F., Chang, S. H., Nishimura, S., and RajBhandary, U. L. (1968) *Biochem. Biophys. Res. Commun.* **33**, 299; RajBhandary, U. L., Chang, S. H., Gross, H. J., Harada, F., Kimura, F., and Nishimura, S. (1969) *Fed. Proc.* **28**, 409
- Dube, S. K., Marcker, K. A., Clark, B. F. C., and Cory, S. (1969) *Eur. J. Biochem.* **8**, 244; Dube, S. K., and Marcker, K. A. (1969) *Eur. J. Biochem.* **8**, 256
- Smith, J. D., Barnett, L., Brenner, S., and Russell, R. L. (1970) *J. Mol. Biol.* **54**, 1; Smith, J. D., Anderson, K., Cashmore, A., Hooper, M. L., and Russell, R. L. (1970) *Cold Spring Harbor Symp. Quant. Biol.* **35**, 21; Altman, S., Brenner, S., and Smith, J. D. (1971) *J. Mol. Biol.* **56**, 195; and others cited in *Handbook of Nucleic Acid Sequences* (1974) (Barrell, B. G., and Clark, B. F. C., eds) pp. 51-52, Joynson-Bruvvers, Ltd., Oxford
- Altman, S., and Smith, J. D. (1971) *Nature New Biol.* **233**, 35
- Miller, R. C., Jr., Besmer, P., Khorana, H. G., Fiantdt, M., and Szybalski, W. (1971) *J. Mol. Biol.* **56**, 363
- Besmer, P., Miller, R. C., Jr., Caruthers, M. H., Kumar, A., Minamoto, K., van de Sande, J. H., Sidorova, N., and Khorana, H. G. (1972) *J. Mol. Biol.* **72**, 503
- Loewen, P. C., Sekiya, T., and Khorana, H. G. (1974) *J. Biol. Chem.* **249**, 217
- Sekiya, T., and Khorana, H. G. (1974) *Proc. Natl. Acad. Sci. U. S. A.* **71**, 2978; Sekiya, T., van Ormondt, H., and Khorana, H. G. (1975) *J. Biol. Chem.* **250**, 1087
- Söll, D. (1971) *Science* **173**, 293; Schäfer, K. P., and Söll, D. (1974) *Biochimie* **56**, 795
- Sgaramella, V., and Khorana, H. G. (1972) *J. Mol. Biol.* **72**, 427
- Ramamoorthy, B., Lees, R. G., Kleid, D. G., and Khorana, H. G. (1976) *J. Biol. Chem.* **251**, 676-694