

DETERMINATION OF NUCLEOTIDE SEQUENCES IN DNA

Nobel lecture, 8 December 1980

by

FREDERICK SANGER

Medical Research Council Laboratory of Molecular Biology,
Cambridge, England

INTRODUCTION

In spite of the important role played by DNA sequences in living matter, it is only relatively recently that general methods for their determination have been developed. This is mainly because of the very large size of DNA molecules, the smallest being those of the simple bacteriophages such as ϕ X174 (which contains about 5,000 nucleotides). It was therefore difficult to develop methods with such complicated systems. There are however some relatively small RNA molecules - notably the transfer RNAs of about 75 nucleotides, and these were used for the early studies on nucleic acid sequences (1).

Following my work on amino acid sequences in proteins (2) I turned my attention to RNA and, with G.G. Brownlee and B.G. Barrell, developed a relatively rapid small-scale method for the fractionation of ^{32}P -labelled oligonucleotides (3). This became the basis for most subsequent studies of RNA sequences. The general approach used in these studies, and in those on proteins, depended on the principle of partial degradation. The large molecules were broken down, usually by suitable enzymes, to give smaller products which were then separated from each other and their sequence determined. When sufficient results had been obtained they were fitted together by a process of deduction to give the complete sequence. This approach was necessarily rather slow and tedious, often involving successive digestions and fractionations, and it was not easy to apply it to the larger DNA molecules. When we first studied DNA some significant sequences of about 50 nucleotides in length were obtained with this method (4,5), but it seemed that to be able to sequence genetic material a new approach was desirable and we turned our attention to the use of copying procedures.

Abbreviations The abbreviations C, A and G are used to describe both the ribonucleotides and the deoxyribonucleotides, according to context.

COPYING PROCEDURES

In the RNA field these procedures had been pioneered by C. Weissmann and his colleagues (6) in their studies on the RNA sequence of the bacteriophage $Q\beta$. $Q\beta$ contains a replicase that will synthesize a complementary copy of the single-stranded RNA chain, starting from its 3' end. These workers devised elegant procedures involving pulse-labelling with radioactively labelled nucleotides, from which sequences could be deduced.

For DNA sequences we have used the enzyme DNA polymerase, which copies single-stranded DNA as shown in Fig. 1. The enzyme requires a primer, which is a single-stranded oligonucleotide having a sequence that is complementary to, and therefore able to hybridize with, a region on the DNA being sequenced (the template). Mononucleotide residues are added sequentially to the 3' end of the primer from the corresponding deoxynucleoside triphosphates, making a complementary copy of the template DNA. By using triphosphates containing ^{32}P in the α position the newly synthesized DNA can be labelled. In the early experiments synthetic oligonucleotides were used as primers, but after

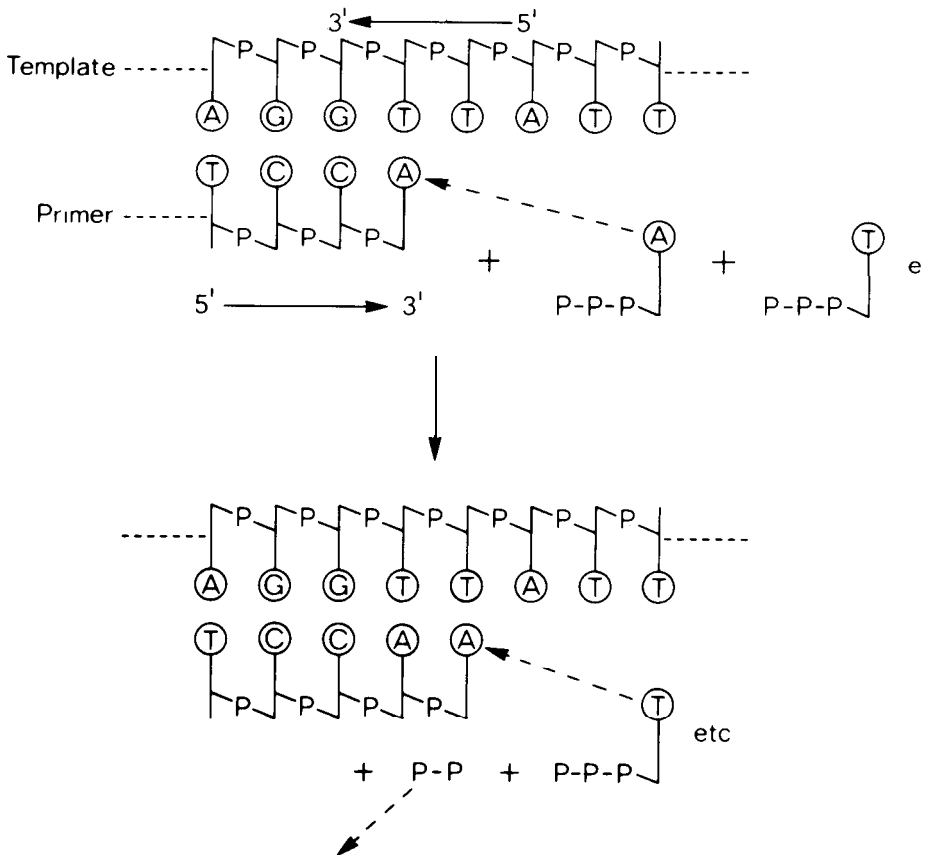


Fig. 1. Specificity requirements for DNA polymerase.

the discovery of restriction enzymes it was more convenient to use fragments resulting from their action as they were much more easily obtained.

The copying procedure was used initially to prepare a short specific region of labelled DNA which could then be subjected to partial digestion procedures. One of the difficulties of sequencing DNA was to find specific methods for breaking it down into small fragments. No suitable enzymes were known that would recognise only one nucleotide. However, Berg, Fancher & Chamberlin (7) had shown earlier that under certain conditions it was possible to incorporate ribonucleotides, in place of the normal deoxyribonucleotides, into DNA chains with DNA polymerase. Thus, for instance, if copying were carried out using riboCTP and the other three deoxynucleoside triphosphates, a chain could be built up in which the C residues were in the ribo form. Bonds involving ribonucleotides could be broken by alkali under conditions where those involving the deoxynucleotides were not, so that a specific splitting at C residues could be obtained. Using this method we were able to extend our sequencing studies to some extent (8). However extensive fractionations and analyses were still required.

THE 'PLUS AND MINUS' METHOD

In the course of these experiments we needed to prepare DNA copies of high specific radioactivity, and in order to do this the highly labelled substrates had to be present in low concentrations. Thus if $\alpha[^{32}\text{P}]\text{-dATP}$ was used for labelling its concentration was much lower than that of the other three triphosphates and frequently when we analysed the newly synthesized DNA chains we found that they terminated at a position immediately before that at which an A should have been incorporated. Consequently a mixture of products was produced all having the same 5' end (the 5' end of the primer) and terminating at the 3' end at the position of the A residues. If these products could be fractionated on a system that separated only on the basis of chain length, the pattern of their distribution on fractionation would be proportional to the distribution of the A's along the DNA chain. And this, together with the distribution of the other three mononucleotides, is the information required for sequence determination. Initial experiments carried out with J.E. Donelson suggested that this approach could be the basis for a more rapid method, and it was found that good fractionations according to size could be obtained by ionophoresis on acrylamide gels.

The method described above met with only limited success but we were able to develop two modified techniques that depended on the same general principle and these provided a much more rapid and simpler method of DNA sequence determination than anything we had used before (9). This, which is known as the "plus and minus" technique, was used to determine the almost complete sequence of the DNA of bacteriophage $\phi\text{X}174$ which contains 5,386 nucleotides (10).

THE 'DIDEOXY' METHOD

More recently we have developed another similar method which uses specific chain-terminating analogues of the normal deoxynucleoside triphosphates (11). This method is both quicker and more accurate than the plus and minus technique. It was used to complete the sequence of ϕ X174 (12), to determine the sequence of a related bacteriophage, G4 (13), and has now been applied to mammalian mitochondrial DNA.

The analogues most widely used are the dideoxynucleoside triphosphates (Fig. 2). They are the same as the normal deoxynucleoside triphosphates but lack the 3' hydroxyl group. They can be incorporated into a growing DNA chain by DNA polymerase but act as terminators because, once they are incorporated, the chain contains no 3' hydroxyl group and so no other nucleotide can be added.

The principle of the method is summarised in Fig. 3. Primer and template are denatured to separate the two strands of the primer, which is usually a restriction enzyme fragment, and then annealed to form the primer-template complex. The mixture is then divided into four samples. One (the T sample) is incubated with DNA polymerase in the presence of a mixture of ddTTP (dideoxy thymidine triphosphate) and a low concentration of TTP, together with the other three deoxynucleoside triphosphates (one of which is labelled

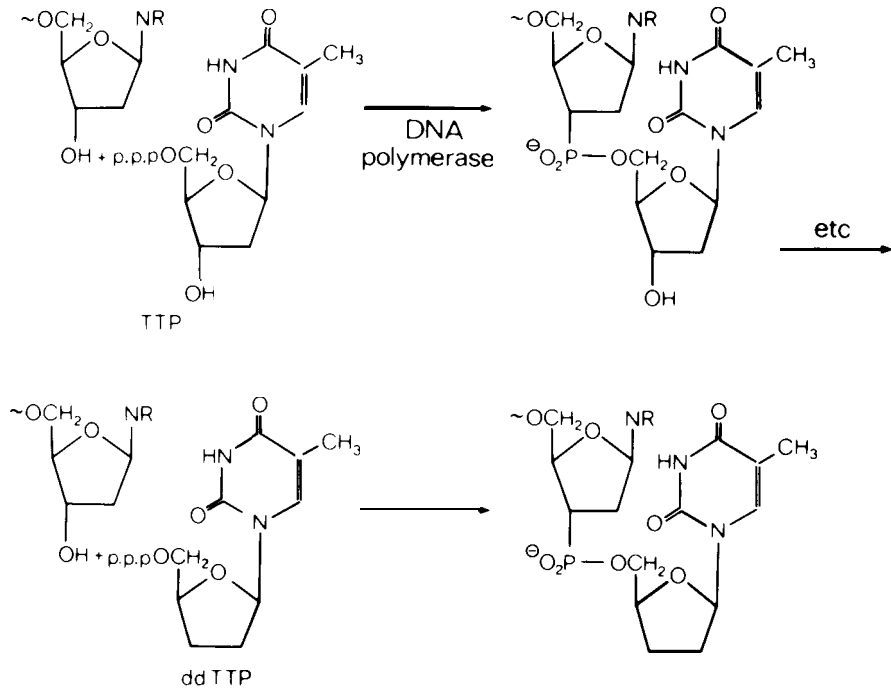


Fig. 2. Diagram showing chain termination with dideoxythymidine triphosphate (ddTTP). The top line shows the DNA polymerase-catalysed reaction of the normal deoxynucleoside triphosphate (TTP) with the 3' terminal nucleotide of the primer: the bottom line the corresponding reaction with ddTTP.

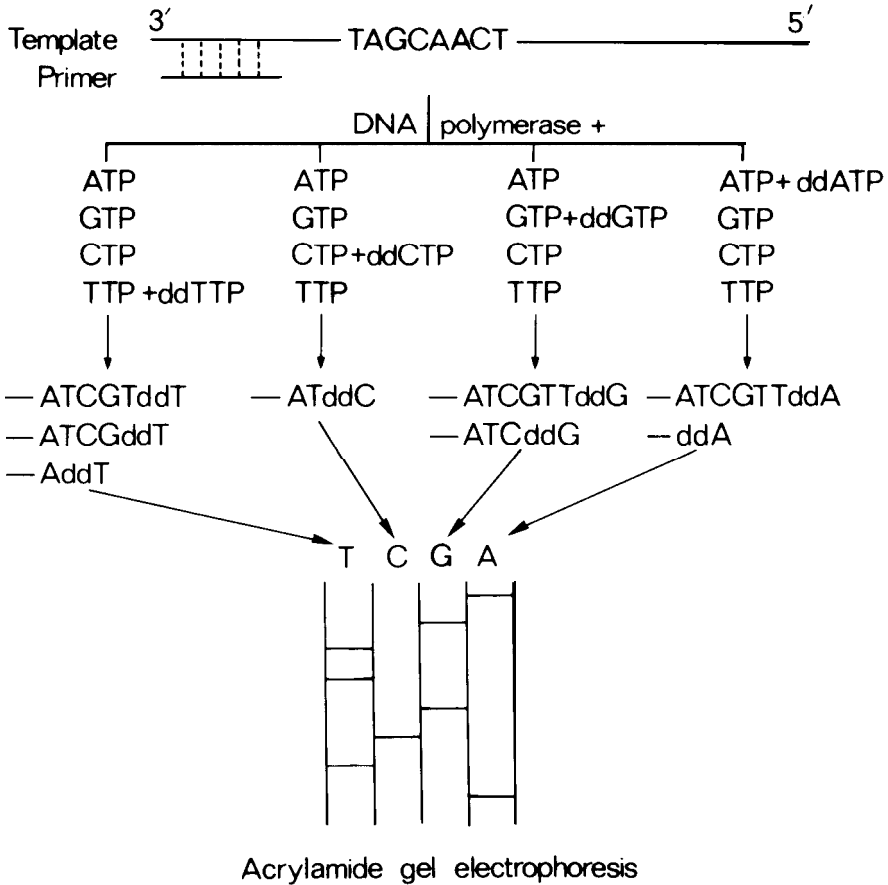


Fig. 3. Principle of the chain-terminating method.

with ^{32}P) at normal concentration. As the DNA chains are built up on the 3' end of the primer the position of the T's will be filled, in most cases by the normal substrate T and extended further, but occasionally by ddT and terminated. Thus at the end of incubation there remains a mixture of chains terminating with T at their 3' end but all having the same 5' end (the 5' end of the primer). Similar incubations are carried out in the presence of each of the other three dideoxy derivatives, giving mixtures terminating at the positions of C, A and G respectively, and the four mixtures are fractionated in parallel by electrophoresis on acrylamide gel under denaturing conditions. This system separates the chains according to size, the small ones moving quickly and the large ones slowly. As all the chains in the T mixture end at T the relative position of the T's in the chain will define the relative sizes of the chains, and therefore their relative positions on the gel after fractionation. The actual sequence can then simply be read off from an autoradiograph of the gel (Fig. 4). The method is comparatively rapid and accurate and sequences of up to about 300 nucleotides from the 3' end of the primer can usually be determined.

2.5 hr
GATC

5hr
GATC

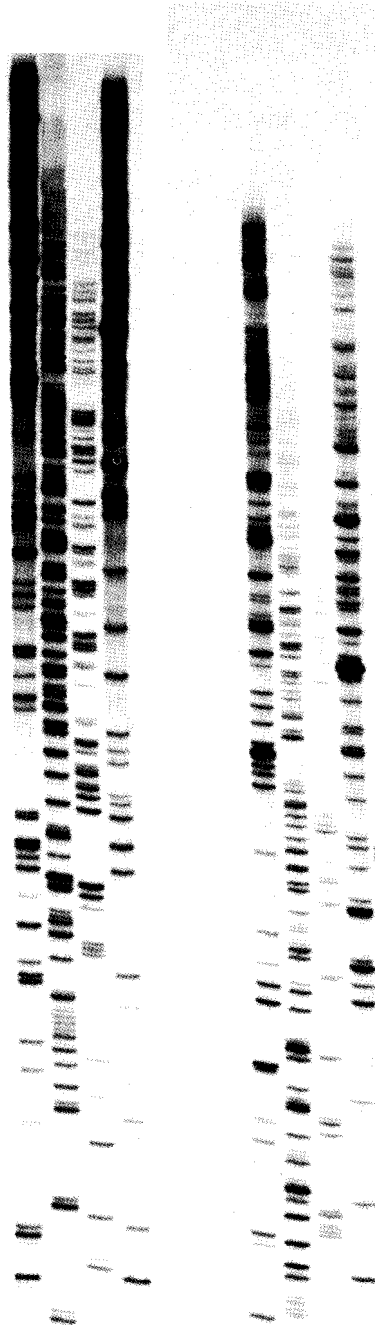


Fig. 4. Autoradiograph of a DNA sequencing gel. The origin is at the top and migration of the DNA chains, according to size, is downwards. The gel on the left has been run for 2.5 hr and that on the right for 5 hr with the same polymerisation mixtures.

Considerably longer sequences have been read off but these are usually less reliable.

One problem with the method is that it requires single-stranded DNA as template. This is the natural form of the DNA in the bacteriophages ϕ X174 and G4, but most DNA is double-stranded and it is frequently difficult to separate the two strands. One way of overcoming this was devised by A.J.H. Smith (14). If the double-stranded linear DNA is treated with exonuclease III (a double-strand specific 3' exonuclease) each chain is degraded from its 3' end, as shown in Fig. 5, giving rise to a structure that is largely single-stranded and can be used as template for DNA polymerase with suitable small primers. This method is particularly suitable for use with fragments cloned in plasmid vectors and has been used extensively in the work on human mitochondrial DNA.

CLONING IN SINGLE-STRANDED BACTERIOPHAGE

Another method of preparing suitable template DNA that is being more widely used is to clone fragments in a single-stranded bacteriophage vector (15-17). This approach is summarised in Fig. 6. Various vectors have been described. We have used a derivative of bacteriophage M 13 developed by Gronenborn & Messing (16) which contains an insert of the β -galactosidase gene with an EcoRI restriction enzyme site in it. The presence of β -galactosidase in a plaque can be readily detected by using a suitable colour-forming substrate (X-gal). The presence of an insert in the EcoRI site destroys the β -galactosidase gene, giving rise to a colourless plaque.

Besides being a simple and general method of preparing single-stranded DNA this approach has other advantages. One is that it is possible to use the same primer on all clones. Heidccker et al. (18) prepared a 96-nucleotide long restriction fragment derived from a position in the M13 vector flanking the EcoRI site (see Fig. 6). This can be used to prime into, and thus determine, a sequence of about 200 nucleotides in the inserted DNA. Smaller synthetic primers have now been prepared (19,20) which allow longer sequences to be determined. The approach that we have used is to prepare clones at random from restriction enzyme digests and determine the sequence with the flanking primer. Computer programmes (21) are then used to store, overlap and arrange the data.

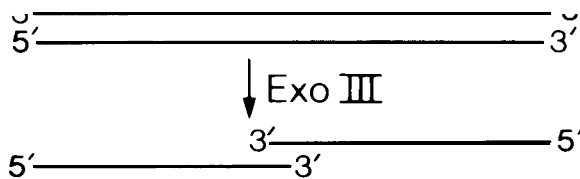


Fig. 5. Degradation of double-stranded DNA with exonuclease III

Another important advantage of the cloning technique is that it is a very efficient and rapid method of fractionating fragments of DNA. In all sequencing techniques, both for proteins and nucleic acids, fractionation has been an important step and major progress has usually been dependent on the development of new fractionation methods. With the new rapid methods for DNA sequencing fractionation is still important and as the sequencing procedure itself is becoming more rapid more of the work has involved fractionation of the restriction enzyme fragments by electrophoresis on acrylamide. This becomes increasingly difficult as larger DNA molecules are studied and may involve several successive fractionations before pure fragments are obtained. In the

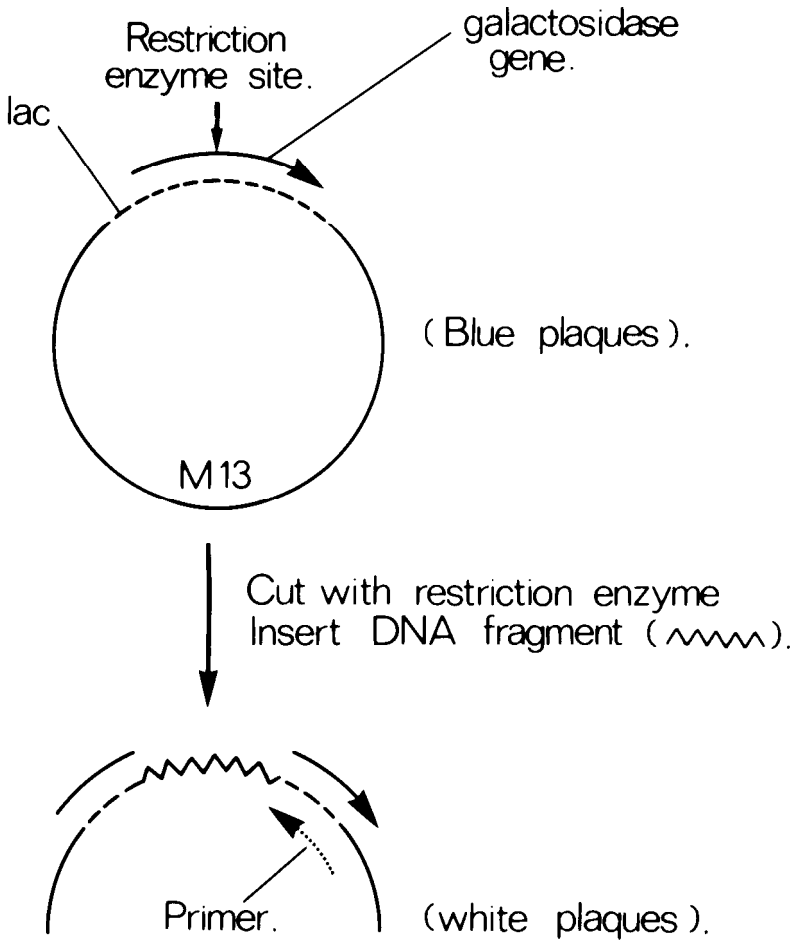


Fig. 6. Diagram illustrating the cloning of DNA fragments in the single-stranded bacteriophage vector M13mp2 (16) and sequencing the insert with a flanking primer.

new method these fractionations are replaced by a cloning procedure. The mixture is spread on an agar plate and grown. Each clone represents the progeny of a single molecule and is therefore pure, irrespective of how complex the original mixture was. It is particularly suitable for studying large DNAs. In fact there is no theoretical limit to the size of DNA that could be sequenced by this method.

We have applied the method to fragments from mitochondrial DNA (22,23) and to bacteriophage λ DNA. Initially new data can be accumulated very quickly (under ideal conditions at about 500-1,000 nucleotides a day). However at later stages much of the data produced will be in regions that have already been sequenced, and progress then appears to be much slower. Nevertheless we find that most new clones studied give some useful data, either for correcting or confirming old sequences. Thus in the work with bacteriophage λ DNA we have about 90% of the molecule identified in sequences and most of the new clones we study contribute some new information. In most studies on DNA one is concerned with identifying the reading frames for protein genes, and to do this the sequence must be correct. Errors can readily occur in such extensive sequences and confirmation is always necessary. We usually consider it necessary to determine the sequence of each region on both strands of the DNA.

Although in theory it would be possible to complete a sequence determination solely by the random approach, it is probably better to use a more specific method to determine the final remaining nucleotides in a sequencing study. Various methods are possible (22,24), but all are slow compared with the random cloning approach.

BACTERIOPHAGE ϕ X174 DNA

The first DNA to be completely sequenced by the copying procedures was from bacteriophage ϕ X174 (10,12) - a single-stranded circular DNA, 5,386 nucleotides long, which codes for ten genes. The most unexpected finding from this work was the presence of 'overlapping' genes. Previous genetic studies had suggested that genes were arranged in a linear order along the DNA chains, each gene being encoded by a unique region of the DNA. The sequencing studies indicated however that there were regions of the ϕ X DNA that were coding for two genes. This is made possible by the nature of the genetic code. Since a sequence of three nucleotides (a codon) codes for one amino acid, each region of DNA can theoretically code for three different amino acid sequences, depending on where translation starts. This is illustrated in Fig. 7. The reading frame or phase in which translation takes place is defined by the position of the initiating ATG codon, following which nucleotides are simply read off three at a time. In ϕ X there is an initiating ATG within the gene coding for the D protein, but in a different reading frame. Consequently this initiates an entirely

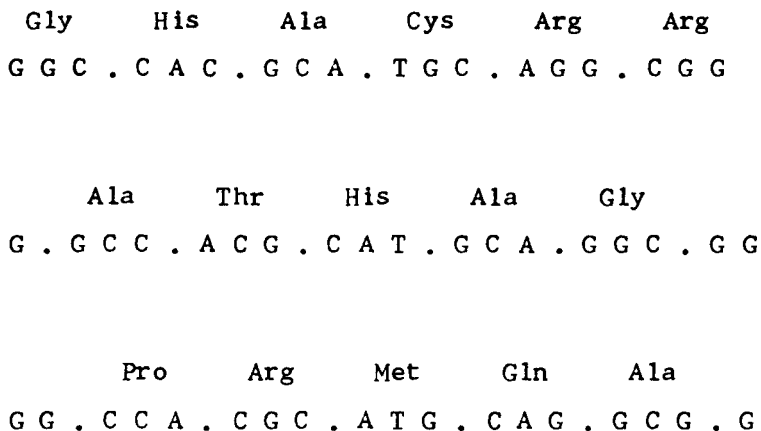


Fig. 7. Diagram illustrating how one DNA sequence can code for three different amino acid sequences. The dots indicate the positions of triplet codons coding for the amino acids.

different sequence of amino acids, which is that of the E protein. Fig. 8 shows the genetic map of ϕX . The E gene is completely contained within the D gene and the B gene within the A.

Further studies (25) on the related bacteriophage, G4, revealed the presence of a previously unidentified gene, which was called K. This overlaps both the A and C genes, and there is a sequence of four nucleotides that codes for part of all three proteins, A, C and K, using all of its three possible reading frames.

It is uncertain whether overlapping genes are a general phenomenon or whether they are confined to viruses, whose survival may depend on their rate of replication and therefore on the size of the DNA: with the overlapping genes more genetic information can be concentrated in a given sized DNA.

Further details of the sequence of bacteriophage $\phi X174$ DNA have been published elsewhere (10,12).

MAMMALIAN MITOCHONDRIAL DNA

Mitochondria contain a small double-stranded DNA (mtDNA) which codes for two ribosomal RNAs (rRNAs), 22-23 transfer RNAs (tRNAs) and about 10-13 proteins which appear to be components of the inner mitochondrial membrane and are somewhat hydrophobic. Other proteins of the mitochondria are encoded by the nucleus of the cell and specifically transported to the mitochondria. Using the above methods we have determined the nucleotide sequence of human mtDNA (23) and almost the complete sequence of bovine mtDNA. The sequence revealed a number of unexpected features which indicated that the transcription and translation machinery of mitochondria is rather different from that of other biological systems.

The genetic code in mitochondria

Hitherto it has been believed that the genetic code was universal. No differences were found in the *E. coli*, yeast or mammalian systems that had been studied. Our initial sequence studies were on human mtDNA. No amino acid sequence of the proteins that were encoded by human mtDNA were known. However Steffans & Buse (26) had determined the sequence of subunit II of cytochrome oxidase (COII) from bovine mitochondria, and Barrell, Bankier & Drouin (27) found that a region of the human mtDNA that they were studying had a sequence that would code for a protein homologous to this amino acid sequence - indicating that it most probably was the gene for the human COII. Surprisingly the DNA sequence contained TGA triplets in the reading frame of the homologous protein. According to the normal genetic code TGA is a termination codon and if it occurs in the reading frame of a protein the polypeptidic chain is terminated at that position. It was noted that in the positions where TGA occurred in the human mtDNA sequence, tryptophan was found in the bovine protein sequence. The only possible conclusion seemed to be that in mitochondria TGA was not a termination codon but was coding for tryptophan. It was similarly concluded that ATA, which normally codes for isoleucine, was coding for methionine. As these studies were based on a

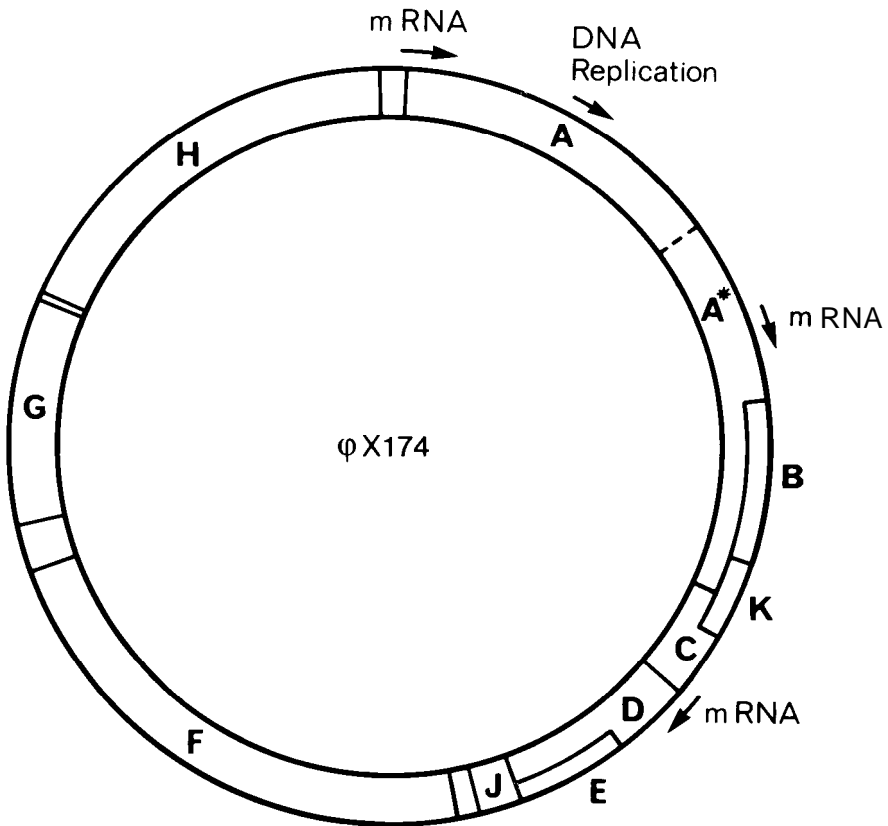


Fig. 8. Gene map of ϕ X174 DNA.

comparison of a human DNA with a bovine protein, the possibility that the differences were due to some species variation, although unlikely, could not be completely excluded. For a conclusive determination of the mitochondrial code it was necessary to compare the DNA sequence of a gene with the amino acid sequence of the protein it was coding for. This was done by Young & Anderson (28) who isolated the bovine mtDNA, determined the sequence of its COII gene and confirmed the above differences.

Fig. 9 shows the human and bovine mitochondrial genetic code and the frequency of use of the different codons in human mitochondria. All codons are used with the exception of UUA and UAG, which are terminators, and AGA and AGG, which normally code for arginine. This suggests that AGA and AGG are probably also termination codons in mitochondria. Further support for this is that no tRNA recognizing the codons has been found (see below) and that some of the unidentified reading frames found in the DNA sequence appear to end with these codons.

In parallel with our studies on mammalian mtDNA, Tzagoloff and his colleagues (29,30) were studying yeast mtDNA. They also found changes in the genetic code, but surprisingly they are not all the same as those found in mammalian mitochondria. These differences are summarised in Table 1.

		SECOND LETTER				
		U	C	A	G	
FIRST LETTER	U	UUU Phe 77	UCU 32	UAU Tyr 46	UGU Cys 5	U
		UUC 140	UCC Ser 99	UAC 89	UGC 17	C
		UUA Leu 73	UCA 83	UAA Ter -	UGA Trp 93	A
		UUG 17	UCG 7	UAG Ter -	UGG 10	G
C	CUU 65	CCU 41	CAU His 18	CGU 7	U	
	CUC Leu 167	CCC Pro 119	CAC 79	CGC Arg 25	C	
	CUA 276	CCA 52	CAA Gln 81	CGA 29	A	
	CUG 45	CCG 7	CAG 9	CGG 2	G	
A	AUU Ile 125	ACU 51	AAU Asn 33	AGU Ser 14	U	
	AUC 196	ACC 155	AAC 130	AGC 39	C	
	AUA Met 166	ACA Thr 133	AAA Lys 85	AGA Ter -	A	
	AUG 40	ACG 10	AAG 10	AGG Ter -	G	
G	GUU 30	GCU 43	GAU Asp 15	GGU 24	U	
	GUC Val 49	GCC Ala 124	GAC 55	GGC Gly 88	C	
	GUA 71	GCA 80	GAA Glu 64	GGA 67	A	
	GUG 18	GCG 8	GAG 24	GGG 34	G	

Fig. 9. The human mitochondrial genetic code, showing the coding properties of the tRNAs (boxed codons) and the total number of codons used in the whole genome shown in Fig. 10. (One methionine tRNA has been detected, but as there is some uncertainty about the number present and their coding properties, these codons are not boxed.)

Table 1. Coding changes in mitochondria

Codon	Amino acid coded		
	Normal	Mammalian mitochondria	Yeast mitochondria
UGA	Term	Trp	Trp
AUA	Ile	Met	Ile
CUN	Leu	Leu	Thr
AGA, AGG	Arg	Term?	Arg
CGN	Arg	Arg	Term?

Transfer RNAs

Transfer RNAs have a characteristic base-pairing structure which can be drawn in the form of the 'cloverleaf' model. By examining the DNA sequence for cloverleaf structures and using a computer programme (31) it was possible to identify genes coding for the mt-tRNAs.

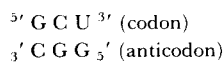
Besides the cloverleaf structure, normal cytoplasmic tRNAs have a number of so-called 'invariable' features which are believed to be important to their biological function. Most of the mammalian mt-tRNAs are anomalous in that some of these invariable features are missing. The most bizarre is one of the serine tRNA in which a complete loop of the cloverleaf structure is missing (32,33). Nevertheless it functions as a tRNA.

In normal cytoplasmic systems at least 32 tRNAs are required to code for all the amino acids. This is related to the 'wobble' effect. Codon-anticodon relationships in the first and second positions of the codons are defined by the normal base-pairing rules, but in the third position G can pair with U. The result of this is that one tRNA can recognise two codons. There are many cases in the genetic code where all four codons starting with the same two nucleotides code for the same amino acid. These are known as 'family boxes'. The situation for the alanine family box is shown in Table 2, indicating that with the normal wobble system two tRNAs are required to code for the four alanine codons.

Table 2. Coding properties of alanine tRNAs

Codon	Anticodon (wobble)	Anticodon (mitochondria)
GCU	GGC	
GCC		UGC
CA GCG	UGC	

Note that the first position of the codon pairs with the third position of the anticodon and vice versa; e. g.



Only 22 tRNA genes could be found in mammalian mtDNA, and for all the family boxes there was only one, which had a T in the position corresponding to the third position of the codon (34). It seems very unlikely that none of the other predicated tRNAs would have been detected and the most feasible explanation is that in mitochondria one tRNA can recognise all four codons in a family box and that a U in the first position of the anticodon can pair with U, C, A or G in the third position of the codon. Clearly in boxes in which two of the codons code for one amino acid and two for a different one, there must be two different tRNAs and the wobble effect still applies. Such tRNAs are found, as expected, in the mitochondrial genes. The coding properties of the mt-tRNAs are shown in Fig. 9. Similar conclusions have been reached by Heckman *et al.* (35) and by Bonitz *et al.* (36), working respectively on neurospora and yeast mitochondria.

Distribution of protein genes

Mitochondrial DNA was known to code for three of the subunits of cytochrome oxidase, probably three subunits of the ATPase complex, cytochrome b, and a number of other unidentified proteins. In order to identify the protein-coding genes, the DNA was searched for reading frames; i.e. long stretches of DNA containing no termination codons in one of the phases and thus being capable of coding for long polypeptide chains. Such reading frames should start with an initiation codon, which in normal systems is nearly always ATG, and end with a termination codon. Fig. 10 summarises the distribution of the reading frames

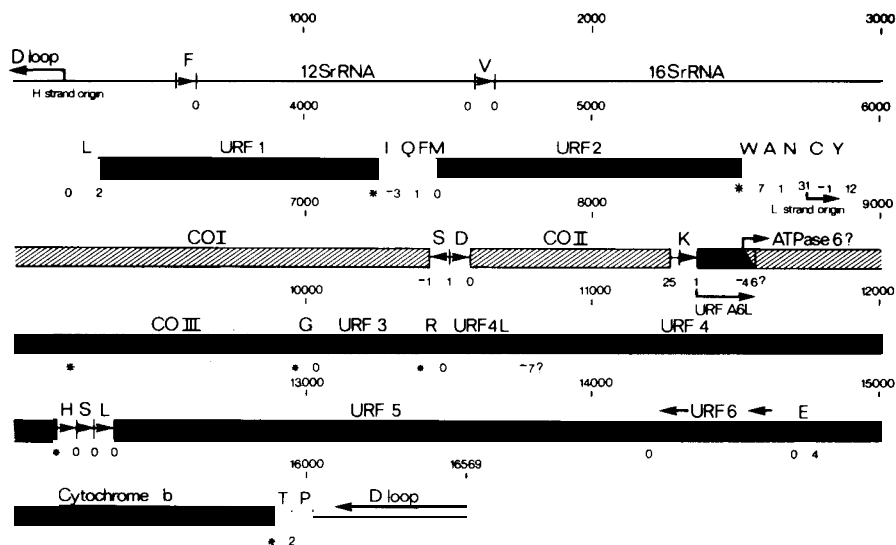


Fig. 10. Gene map of human mtDNA deduced from the DNA sequence. Boxed regions are the predicted reading frames for the proteins. URF = unidentified reading frame. tRNA genes are denoted by the one-letter amino acid code and are either L strand coded (\blacktriangleright) or H strand coded (\blacktriangleleft). Numbers above the genes show the scale in nucleotides and below the predicted number between genes.

* Indicates that termination codons are created by polyadenylation of the mRNA

on the DNA and these are believed to be the genes coding for the proteins. The gene for COII was identified from the amino acid sequence as described above, for subunit I of the cytochrome oxidase from amino acid sequence studies on the bovine protein by J. E. Walker (personal communication), and COIII, cytochrome b and, probably, ATPase 6 were identified by comparison with the DNA sequences of the corresponding genes in yeast mitochondria. Tzagoloff and his colleagues were able to identify these genes in yeast by genetic methods. It has not yet been possible to assign proteins to the other reading frames.

One unusual feature of the mtDNA is that it has a very compact structure. The reading frames coding for the proteins and the rRNA genes appear to be flanked by the tRNA genes with no, or very few, intervening nucleotides. This suggests a relatively simple model for transcription, in which a large RNA is copied from the DNA and the tRNAs are cut out by a processing enzyme, and this same processing leads to the production of the rRNAs and the messenger RNAs (mRNAs), most of which will be monocistronic. Strong support for this model comes from the work of Attardi (37,38) who has identified the RNA sequences at the 5' and 3' ends of the mRNAs, thus locating them on the DNA sequence. One consequence of this arrangement is that the initiation codon is at, or very near, the 5' end of the mRNAs. This suggests that there must be a different mechanism of initiation from that found in other systems. In bacteria there is usually a ribosomal binding site before the initiating ATG codon, whereas in eucaryotic systems the 'cap' structure on the 5' end of the mRNA appears to have a similar function and the first ATG following the cap acts as initiator. It seems that mitochondria may have a more simple, and perhaps more primitive, system with the translation machinery recognising simply the 5' end of the mRNA. Another unique feature of mitochondria is that ATA and possibly ATT can act as initiator codons as well as ATG.

On the basis of the above model, some of the mRNAs will not contain termination codons at their 3' ends after the tRNAs are cut out. However they have T or TA at the 3' end. The mRNAs are generally found polyadenylated at their 3' ends and this process will necessarily give rise to the codon TAA to terminate those protein reading frames.

The compact structure of the mammalian mitochondrial genome is in marked contrast to that of yeast, which is about five times as large and yet codes for only about the same number of proteins and RNAs. The genes are separated by long AT-rich stretches of DNA with no obvious biological function. There are also insertion sequences within some of the genes, whereas these appear to be absent from mammalian mtDNA.

REFERENCES

1. Holley, R. W., Les Prix Nobel, p. 183 (1968)
2. Sanger, F., Les Prix Nobel, p. 134 (1958)
3. Sanger, F., Brownlee, G. G. & Barrell, B. G., *J. Mol. Biol.* 13, 373 (1965)
4. Robertson, H. D., Barrell, B. G., Weith, H. L. & Donelson, J. E., *Nature, New Biol.* 244, 38 (1973)
5. Ziff, E. B., Sedat, J. W. & Galibert, F., *Nature, New Biol.* 241, 34 (1973)
6. Billeter, M. A., Dahlberg, J. E., Goodman, H. M., Hindley, J. & Weissmann, C., *Nature* 224, 1083 (1969)
7. Berg, P., Fancher, H. & Chamberlin, M., Symp. "Informational Macromolecules", pp 467-483, Academic Press: New York & London (1963)
8. Sanger, F., Donelson, J. E., Coulson, A. R., Kössel, H. & Fischer, D., *Proc. Natl. Acad. Sci. USA* 70, 1209 (1973)
9. Sanger, F. & Coulson, A. R. *J. Mol. Biol.* 94, 441 (1975)
10. Sanger, F., Air, G. M., Barrell, B. G., Brown, N. L., Coulson, A. R. Fiddes, J. C., Hutchison, C. A., Slocombe, P. M. & Smith, M., *Nature* 265, 687 (1977)
11. Sanger, F., Nicklen, S. & Coulson, A. R., *Proc. Natl. Acad. Sci. USA* 74, 5463 (1977)
12. Sanger, F., Coulson, A. R., Friedmann, T., Air, G. M., Barrell, B. G., Brown, N. L., Fiddes, J. C., Hutchison, C. A., Slocombe, P. M. & Smith, M., *J. Mol. Biol.* 125, 225 (1978)
13. Godson, G. N., Barrell, B. G., Staden R. & Fiddes, J. C., *Nature* 276, 236 (1978)
14. Smith, A. J. H., *Nucl. Acids Res.* 6, 831 (1979)
15. Barnes, W. M., *Gene* 5, 127 (1979)
16. Gronenborn, B. & Messing, J., *Nature* 272, 375 (1978)
17. Herrmann, R., Neugebauer, K., Schaller, H. & Zentgraf, H., In "The Singlestranded DNA Phages" (Eds. Denhardt, D. T., Dressier, D. & Ray, D. S.) pp 473-476, Cold Spring Harbor Laboratory, New York (1978)
18. Heidecker, G., Messing, J. & Gronenborn, B., *Gene* 10, 69 (1980)
19. Anderson, S., Gait, M. J., Mayol, L. & Young, I. G., *Nucl. Acids. Res.* 8, 1731 (1980)
20. Gait, M. J., Unpublished results (1980)
21. Staden, R., *Nucl. Acids Res.* 8, 3673 (1980)
22. Sanger, F., Coulson, A. R., Barrell, B. G., Smith, A. J. H. & Roe, B. A., *J. Mol. Biol.* 143, 161 (1980)
23. Barrell, B. G., Anderson, S., Bankier, A. T., de Bruijn, M. H. L., Chen, E., Coulson, A. R., Drouin, J., Eperon, I. C., Nierlich, D. P., Roe, B. A., Sanger, F., Schreier, P. H., Smith, A. J. H., Staden, R. & Young, I. G., 31st Mosbach Colloq. "Biological Chemistry of organelle Formation", (Eds. Bücher, T., Sebald, W. & Weiss, H.) pp 11 -25, Springer-Verlag, Berlin (1980)
24. Winter, G. & Fields, S., *Nucl. Acids Res.* 8, 1965 (1980)
25. Shaw, D. C., Walker, J. E., Northrop, F. D., Barrell, B. G., Godson, G. N. & Fiddes, J. C., *Nature* 272, 510 (1978)
26. Steffans, G. J. & Bose, G., *Hoppe-Seyler's Z. Physiol. Chem.* 360, 613 (1979)
27. Barrell, B. G., Bankier, A. T. & Drouin, J., *Nature* 282, 189 (1979)
28. Young, I. G. & Anderson, S., *Gene*, 12, 257 (1980)
29. Macino, G., Coruzzi, G., Nobrega, F. G., Li, M. & Tzagoloff, A., *Proc. Natl. Acad. Sci. USA* 76, 3784 (1979)
30. Macino, G. & Tzagoloff, A., *Proc. Natl. Acad. Sci. USA* 76, 131 (1979)
31. Staden, R., *Nucl. Acids Res.* 8, 817 (1980)
32. de Bruijn, M. H. L., Schreier, P. H., Eperon, I. C., Barrell, B. G., Chen, E. Y., Armstrong, P. W., Wong, J. F. H. & Roe, B. A., *Nucl. Acids Res.* 8, 5213 (1980)
33. Arcari, P. & Brownlee, G. G., *Nucl. Acids. Res.* 8, 5207 (1980)
34. Barrell, B. G., Anderson, S., Bankier, A. T., de Bruijn, M. H. L., Chen, E., Coulson, A. R., Drouin, J., Eperon, I. C., Nierlich, D. P., Roe, B. A., Sanger, F., Schreier, P. H., Smith, A. J. H., Staden, R. & Young, I. G., *Proc. Natl. Acad. Sci. USA* 77, 3164 (1980)

35. Heckman, J. E., Samoff, J., Alzner-De Weerd, B., Yin, S. & RajBhandary, U. L., Proc. Natl. Acad. Sci. USA 77, 3159 (1980)
36. Bonitz, S. G., Berlani, R., Coruzzi, G., Li, M., Macino, G., Nobrega, F. G., Nobrega, M. P., Thalenfeld, B. E. & Tzagoloff, A., Proc. Natl. Acad. Sci. USA, 77, 3167 (1980)
37. Montoya, J., Ojala, D. & Attardi, G., Nature 290,465 (1981)
38. Ojala, D., Montoya, J. & Attardi, G., Nature 290,470 (1981)